

cloudera[®]

Cloudera Security

Important Notice

© 2010-2021 Cloudera, Inc. All rights reserved.

Cloudera, the Cloudera logo, and any other product or service names or slogans contained in this document are trademarks of Cloudera and its suppliers or licensors, and may not be copied, imitated or used, in whole or in part, without the prior written permission of Cloudera or the applicable trademark holder. If this documentation includes code, including but not limited to, code examples, Cloudera makes this available to you under the terms of the Apache License, Version 2.0, including any required notices. A copy of the Apache License Version 2.0, including any notices, is included herein. A copy of the Apache License Version 2.0 can also be found here: <https://opensource.org/licenses/Apache-2.0>

Hadoop and the Hadoop elephant logo are trademarks of the Apache Software Foundation. All other trademarks, registered trademarks, product names and company names or logos mentioned in this document are the property of their respective owners. Reference to any products, services, processes or other information, by trade name, trademark, manufacturer, supplier or otherwise does not constitute or imply endorsement, sponsorship or recommendation thereof by us.

Complying with all applicable copyright laws is the responsibility of the user. Without limiting the rights under copyright, no part of this document may be reproduced, stored in or introduced into a retrieval system, or transmitted in any form or by any means (electronic, mechanical, photocopying, recording, or otherwise), or for any purpose, without the express written permission of Cloudera.

Cloudera may have patents, patent applications, trademarks, copyrights, or other intellectual property rights covering subject matter in this document. Except as expressly provided in any written license agreement from Cloudera, the furnishing of this document does not give you any license to these patents, trademarks copyrights, or other intellectual property. For information about patents covering Cloudera products, see <http://tiny.cloudera.com/patents>.

The information in this document is subject to change without notice. Cloudera shall not be liable for any damages resulting from technical errors or omissions which may be present in this document, or from use of this document.

Cloudera, Inc.
395 Page Mill Road
Palo Alto, CA 94306
info@cloudera.com
US: 1-888-789-1488
Intl: 1-650-362-0488
www.cloudera.com

Release Information

Version: Cloudera Enterprise 5.6.x
Date: February 3, 2021

Table of Contents

About this Guide	10
-------------------------	-----------

Security Overview for an Enterprise Data Hub	11
-----------------------------------------------------	-----------

Facets of Hadoop Security	11
Levels of Hadoop Security	11
Hadoop Security Architecture	12
Overview of Authentication Mechanisms for an Enterprise Data Hub	13
<i>Basic Kerberos Concepts</i>	14
<i>Types of Kerberos Deployments</i>	15
<i>TLS/SSL Requirements for Secure Distribution of Kerberos Keytabs</i>	21
<i>Configuring Kerberos Authentication on a Cluster</i>	22
<i>Authentication Mechanisms used by Hadoop Projects</i>	22
Overview of Data Protection Mechanisms for an Enterprise Data Hub	23
<i>Protecting Data At-Rest</i>	23
<i>Protecting Data In-Transit</i>	26
<i>Data Protection within Hadoop Projects</i>	28
Overview of Authorization Mechanisms for an Enterprise Data Hub	28
<i>Authorization Mechanisms in Hadoop</i>	29
<i>Integration with Authentication Mechanisms for Identity Management</i>	35
<i>Authorization within Hadoop Projects</i>	35
Overview of Data Management Mechanisms for an Enterprise Data Hub	36
<i>Cloudera Navigator</i>	37
<i>Integration within an EDH</i>	38
<i>Auditing in Hadoop Projects</i>	38

How to Configure TLS Encryption for Cloudera Manager	40
-------------------------------------------------------------	-----------

Generate TLS Certificates	40
<i>Generate the Cloudera Manager Server Certificate</i>	40
<i>Generate the Cloudera Manager Agent Certificates</i>	42
Configuring TLS Encryption for the Cloudera Manager Admin Console	43
<i>Step 1: Enable HTTPS for the Cloudera Manager Admin Console</i>	43
<i>Step 2: Specify SSL Truststore Properties for Cloudera Management Services</i>	44
<i>Step 3: Restart Cloudera Manager and Services</i>	44
Configuring TLS Encryption for Cloudera Manager Agents	45
<i>Step 1: Enable TLS Encryption for Agents in Cloudera Manager</i>	45
<i>Step 2: Enable TLS on Cloudera Manager Agent Hosts</i>	45

- Step 3: Restart Cloudera Manager Server and Agents.....45
- Step 4: Verify that the Cloudera Manager Server and Agents are Communicating.....45
- Enabling Server Certificate Verification on Cloudera Manager Agents.....45
- Configuring Agent Certificate Authentication.....46
- Step 1: Export the Private Key to a File.....46
- Step 2: Create a Password File.....47
- Step 3: Configure the Agent to Use Private Keys and Certificates.....47
- Step 4: Enable Agent Certificate Authentication.....47
- Step 5: Restart Cloudera Manager Server and Agents.....47
- Step 6: Verify that Cloudera Manager Server and Agents are Communicating.....48

Configuring Authentication.....49

- Configuring Authentication in Cloudera Manager.....49
- Cloudera Manager User Accounts.....50
- Configuring External Authentication for Cloudera Manager.....51
- Kerberos Concepts - Principals, Keytabs and Delegation Tokens.....57
- Enabling Kerberos Authentication Using the Wizard.....59
- Enabling Kerberos Authentication for Single User Mode or Non-Default Users.....69
- Configuring a Cluster with Custom Kerberos Principals.....70
- Viewing and Regenerating Kerberos Principals.....72
- Using a Custom Kerberos Keytab Retrieval Script.....72
- Mapping Kerberos Principals to Short Names.....73
- Moving Kerberos Principals to Another OU Within Active Directory.....74
- Using Auth-to-Local Rules to Isolate Cluster Users.....74
- Enabling Kerberos Authentication Without the Wizard.....75
- Configuring Authentication in the Cloudera Navigator Data Management Component.....86
- Configuring External Authentication for the Cloudera Navigator Data Management Component.....86
- Managing Users and Groups for the Cloudera Navigator Data Management Component.....91
- Configuring Authentication in CDH Using the Command Line.....93
- Enabling Kerberos Authentication for Hadoop Using the Command Line.....93
- FUSE Kerberos Configuration.....115
- Using kadmin to Create Kerberos Keytab Files.....115
- Configuring the Mapping from Kerberos Principals to Short Names.....117
- Enabling Debugging Output for the Sun Kerberos Classes.....119
- Flume Authentication.....119
- Configuring Flume's Security Properties.....120
- Configuring Kerberos for Flume Thrift Source and Sink Using Cloudera Manager.....121
- Configuring Kerberos for Flume Thrift Source and Sink Using the Command Line.....122
- Flume Account Requirements.....123
- Testing the Flume HDFS Sink Configuration.....124
- Writing to a Secure HBase cluster.....124
- HBase Authentication.....125
- Configuring Kerberos Authentication for HBase.....125

<i>Configuring Secure HBase Replication.....</i>	<i>131</i>
<i>Configuring the HBase Client TGT Renewal Period.....</i>	<i>131</i>
<i>HCatalog Authentication.....</i>	<i>131</i>
<i> Before You Start.....</i>	<i>132</i>
<i> Step 1: Create the HTTP keytab file.....</i>	<i>132</i>
<i> Step 2: Configure WebHCat to Use Security.....</i>	<i>132</i>
<i> Step 3: Create Proxy Users.....</i>	<i>132</i>
<i> Step 4: Verify the Configuration.....</i>	<i>133</i>
<i>Hive Authentication.....</i>	<i>133</i>
<i> HiveServer2 Security Configuration.....</i>	<i>133</i>
<i> Hive Metastore Server Security Configuration.....</i>	<i>139</i>
<i> Using Hive to Run Queries on a Secure HBase Server.....</i>	<i>140</i>
<i>HttpFS Authentication.....</i>	<i>141</i>
<i> Configuring the HttpFS Server to Support Kerberos Security.....</i>	<i>141</i>
<i> Using curl to access an URL Protected by Kerberos HTTP SPNEGO.....</i>	<i>142</i>
<i>Hue Authentication.....</i>	<i>143</i>
<i> Hue Security Enhancements.....</i>	<i>143</i>
<i> Configuring Kerberos Authentication for Hue.....</i>	<i>144</i>
<i> Integrating Hue with LDAP.....</i>	<i>147</i>
<i> Configuring Hue for SAML.....</i>	<i>151</i>
<i>Impala Authentication.....</i>	<i>154</i>
<i> Enabling Kerberos Authentication for Impala.....</i>	<i>154</i>
<i> Enabling LDAP Authentication for Impala.....</i>	<i>157</i>
<i> Using Multiple Authentication Methods with Impala.....</i>	<i>159</i>
<i> Configuring Impala Delegation for Hue and BI Tools.....</i>	<i>159</i>
<i>Llama Authentication.....</i>	<i>160</i>
<i> Configuring Llama to Support Kerberos Security.....</i>	<i>161</i>
<i>Oozie Authentication.....</i>	<i>161</i>
<i> Configuring Kerberos Authentication for the Oozie Server.....</i>	<i>162</i>
<i> Configuring Oozie HA with Kerberos.....</i>	<i>163</i>
<i>Search Authentication.....</i>	<i>163</i>
<i> Enabling Kerberos Authentication for Search.....</i>	<i>164</i>
<i> Enabling LDAP Authentication for Search.....</i>	<i>165</i>
<i> Using Kerberos with Search.....</i>	<i>167</i>
<i>Spark Authentication.....</i>	<i>170</i>
<i> Configuring Kerberos Authentication for Spark Using the Command Line.....</i>	<i>170</i>
<i> Configuring Spark Authentication With a Shared Secret Using Cloudera Manager.....</i>	<i>171</i>
<i> Configuring Spark on YARN for Long-Running Applications.....</i>	<i>171</i>
<i>Sqoop 2 Authentication.....</i>	<i>171</i>
<i> Create the Sqoop 2 Principal and Keytab File.....</i>	<i>171</i>
<i> Configure Sqoop 2 to Use Kerberos.....</i>	<i>172</i>
<i>ZooKeeper Authentication.....</i>	<i>172</i>
<i> Configuring ZooKeeper Server for Kerberos Authentication.....</i>	<i>172</i>
<i> Configuring the ZooKeeper Client Shell to Support Kerberos Security.....</i>	<i>173</i>

<i>Verifying the Configuration</i>	174
Hadoop Users in Cloudera Manager and CDH.....	174
Configuring a Cluster-dedicated MIT KDC with Cross-Realm Trust.....	180
<i>When to use kadmind.local and kadmind</i>	180
<i>Setting up a Cluster-Dedicated KDC and Default Realm for the Hadoop Cluster</i>	180
Integrating Hadoop Security with Active Directory.....	185
<i>Configuring a Local MIT Kerberos Realm to Trust Active Directory</i>	186
Integrating Hadoop Security with Alternate Authentication.....	187
<i>Configuring the AuthenticationFilter to use Kerberos</i>	188
<i>Creating an AltKerberosAuthenticationHandler Subclass</i>	188
<i>Enabling Your AltKerberosAuthenticationHandler Subclass</i>	188
<i>Example Implementation for Oozie</i>	189
Authenticating Kerberos Principals in Java Code.....	190
Using a Web Browser to Access an URL Protected by Kerberos HTTP SPNEGO.....	190
Troubleshooting Authentication Issues.....	194
<i>Sample Kerberos Configuration files: krb5.conf, kdc.conf, kadmind5.acl</i>	194
<i>Potential Security Problems and Their Solutions</i>	196

Configuring Encryption.....205

TLS/SSL Certificates Overview.....	205
<i>Creating Certificates</i>	206
<i>Creating Java Keystores and Truststores</i>	207
<i>Private Key and Certificate Reuse Across Java Keystores and OpenSSL</i>	210
Configuring TLS Security for Cloudera Manager.....	211
<i>Configuring TLS Encryption Only for Cloudera Manager</i>	212
<i>Level 1: Configuring TLS Encryption for Cloudera Manager Agents</i>	216
<i>Level 2: Configuring TLS Verification of Cloudera Manager Server by the Agents</i>	217
<i>Level 3: Configuring TLS Authentication of Agents to the Cloudera Manager Server</i>	219
<i>HTTPS Communication in Cloudera Manager</i>	224
<i>Troubleshooting TLS/SSL Connectivity</i>	226
<i>Using Self-Signed Certificates for TLS</i>	226
Configuring TLS/SSL for the Cloudera Navigator Data Management Component.....	227
Configuring TLS/SSL for Cloudera Management Service Roles.....	228
Configuring TLS/SSL Encryption for CDH Services.....	228
<i>Prerequisites</i>	229
<i>Hadoop Services as TLS/SSL Servers and Clients</i>	229
<i>Compatible Certificate Formats for Hadoop Components</i>	229
<i>Configuring TLS/SSL for HDFS, YARN and MapReduce</i>	229
<i>Configuring TLS/SSL for HBase</i>	232
<i>Configuring TLS/SSL for Flume Thrift Source and Sink</i>	234
<i>Configuring Encrypted Communication Between HiveServer2 and Client Drivers</i>	235
<i>Configuring TLS/SSL for Hue</i>	237
<i>Configuring TLS/SSL for Impala</i>	241

<i>Configuring TLS/SSL for Oozie</i>	243
<i>Configuring TLS/SSL for Solr</i>	244
<i>Configuring TLS/SSL for HttpFS</i>	248
<i>Encrypted Shuffle and Encrypted Web UIs</i>	250
<i>Deployment Planning for Data at Rest Encryption</i>	256
<i>Data at Rest Encryption Reference Architecture</i>	256
<i>Data at Rest Encryption Requirements</i>	257
<i>Resource Planning for Data at Rest Encryption</i>	261
<i>HDFS Transparent Encryption</i>	261
<i>Use Cases</i>	261
<i>Architecture</i>	262
<i>DistCp Considerations</i>	262
<i>Attack Vectors</i>	263
<i>Optimizing for HDFS Data at Rest Encryption</i>	264
<i>Enabling HDFS Encryption Using the Wizard</i>	266
<i>Managing Encryption Keys and Zones</i>	272
<i>Configuring the Key Management Server (KMS)</i>	274
<i>Securing the Key Management Server (KMS)</i>	278
<i>Migrating Keys from a Java KeyStore to Cloudera Navigator Key Trustee Server</i>	289
<i>Configuring CDH Services for HDFS Encryption</i>	290
<i>Troubleshooting HDFS Encryption</i>	296
<i>Cloudera Navigator Key Trustee Server</i>	297
<i>Backing Up and Restoring Key Trustee Server</i>	298
<i>Initializing Standalone Key Trustee Server</i>	302
<i>Configuring a Mail Transfer Agent for Key Trustee Server</i>	304
<i>Verifying Cloudera Navigator Key Trustee Server Operations</i>	304
<i>Managing Key Trustee Server Organizations</i>	304
<i>Managing Key Trustee Server Certificates</i>	307
<i>Cloudera Navigator Key HSM</i>	309
<i>Initializing Navigator Key HSM</i>	309
<i>HSM-Specific Setup for Cloudera Navigator Key HSM</i>	311
<i>Validating Key HSM Settings</i>	313
<i>Creating a Key Store with CA-Signed Certificate</i>	313
<i>Managing the Navigator Key HSM Service</i>	313
<i>Integrating Key HSM with Key Trustee Server</i>	314
<i>Cloudera Navigator Encrypt</i>	316
<i>Registering Cloudera Navigator Encrypt with Key Trustee Server</i>	316
<i>Preparing for Encryption Using Cloudera Navigator Encrypt</i>	320
<i>Encrypting and Decrypting Data Using Cloudera Navigator Encrypt</i>	325
<i>Migrating eCryptfs-Encrypted Data to dm-crypt</i>	328
<i>Cloudera Navigator Encrypt Access Control List</i>	329
<i>Maintaining Cloudera Navigator Encrypt</i>	334
<i>Configuring Encryption for Data Spills</i>	338
<i>MapReduce v2 (YARN)</i>	338

<i>HBase</i>	338
<i>Impala</i>	338
<i>Hive</i>	339
<i>Flume</i>	339
<i>Using an On-disk Encrypted File Channel</i>	339
Configuring Encrypted HDFS Data Transport.....	341
<i>Using Cloudera Manager</i>	342
<i>Using the Command Line</i>	342
Configuring Encrypted HBase Data Transport.....	342
<i>Using Cloudera Manager</i>	342
<i>Using the Command Line</i>	343

Configuring Authorization.....344

Cloudera Manager User Roles.....	344
<i>User Roles</i>	344
<i>Determining the Role of the Currently Logged in User</i>	346
<i>Removing the Full Administrator User Role</i>	346
Cloudera Navigator Data Management Component User Roles.....	347
<i>User Roles</i>	347
<i>Determining the Roles of the Currently Logged in User</i>	347
HDFS Extended ACLs.....	347
<i>Enabling HDFS Access Control Lists</i>	348
<i>Commands</i>	348
<i>HDFS Extended ACL Example</i>	349
Configuring LDAP Group Mappings.....	350
<i>Using Cloudera Manager</i>	352
<i>Using the Command Line</i>	353
Authorization With Apache Sentry (Incubating).....	354
<i>Architecture Overview</i>	354
<i>Sentry Integration with the Hadoop Ecosystem</i>	356
<i>The Sentry Service</i>	359
<i>Sentry Policy File Authorization</i>	392
<i>Enabling Sentry Authorization for Impala</i>	411
<i>Enabling Sentry Authorization for Search using the Command Line</i>	422
Configuring HBase Authorization.....	432
<i>Understanding HBase Access Levels</i>	432
<i>Enable HBase Authorization</i>	434
<i>Configure Access Control Lists for Authorization</i>	435

Sensitive Data Redaction.....437

Cloudera Manager API Redaction.....	438
Enabling Log and Query Redaction Using Cloudera Manager.....	438

Configuring the Cloudera Navigator Data Management Component to Redact PII.....	439
---------------------------------------------------------------------------------	-----

Overview of Impala Security.....441

Security Guidelines for Impala.....	441
Securing Impala Data and Log Files.....	442
Installation Considerations for Impala Security.....	443
Securing the Hive Metastore Database.....	443
Securing the Impala Web User Interface.....	443

Miscellaneous Topics.....445

Jsvc, Task Controller and Container Executor Programs.....	445
<i>MRv1 and YARN: The jsvc Program.....</i>	<i>445</i>
<i>MRv1 Only: The Linux TaskController Program.....</i>	<i>445</i>
<i>YARN Only: The Linux Container Executor Program.....</i>	<i>445</i>
<i>Task-controller and Container-executor Error Codes.....</i>	<i>446</i>
<i>MRv1 ONLY: Task-controller Error Codes.....</i>	<i>446</i>
<i>YARN ONLY: Container-executor Error Codes.....</i>	<i>448</i>
Sqoop, Pig, and Whirr Security Support Status.....	449
Setting Up a Gateway Node to Restrict Cluster Access.....	450
<i>Installing and Configuring the Firewall and Gateway.....</i>	<i>450</i>
<i>Accessing HDFS.....</i>	<i>450</i>
<i>Submitting and Monitoring Jobs.....</i>	<i>451</i>
Logging a Security Support Case.....	451
<i>Kerberos Issues.....</i>	<i>451</i>
<i>TLS/SSL Issues.....</i>	<i>451</i>
<i>LDAP Issues.....</i>	<i>451</i>
Using Antivirus Software on CDH Hosts.....	452

Appendix: Apache License, Version 2.0.....453

About this Guide

This guide is intended for system administrators who want to secure a cluster using data encryption, user authentication, and authorization techniques. This topic also provides information about Hadoop security programs and shows you how to set up a gateway to restrict access.

Security Overview for an Enterprise Data Hub

As the adoption of Hadoop increases, the volume of data and the types of data handled by Hadoop deployments have also grown. For production deployments, a large amount of this data is generally sensitive, or subject to industry regulations and governance controls. In order to be compliant with such regulations, Hadoop must offer strong capabilities to thwart any attacks on the data it stores, and take measures to ensure proper security at all times. The security landscape for Hadoop is changing rapidly. However, this rate of change is not consistent across all Hadoop components, which is why the degree of security capability might appear uneven across the Hadoop ecosystem. That is, some components might be compatible with stronger security technologies than others.

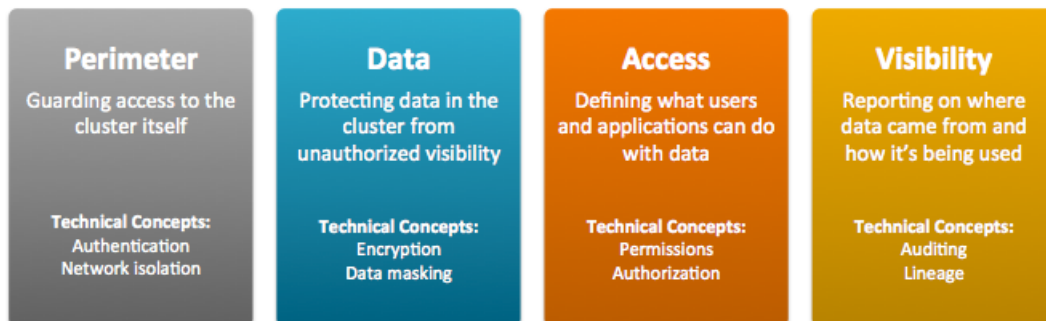
This topic describes the facets of Hadoop security and the different stages in which security can be added to a Hadoop cluster:

The rest of this security overview outlines the various options organizations have for integrating Hadoop security (authentication, authorization, data protection and governance) in enterprise environments. Its aim is to provide you with enough information to make an informed decision as to how you would like to deploy Hadoop security in your organization. In some cases, it also describes the architectural considerations for setting up, managing, and integrating Hadoop security.

Facets of Hadoop Security

Hadoop security can be viewed as a series of business and operational capabilities, including:

- **Perimeter Security**, which focuses on guarding access to the cluster, its data, and its various services. In information security, this translates to [Authentication](#).
- **Data Protection**, which comprises the protection of data from unauthorized access, at rest and in transit. In information security, this translates to [Encryption](#).
- **Entitlement**, which includes the definition and enforcement of what users and applications can do with data. In information security, this translates to [Authorization](#).
- **Transparency**, which consists of the reporting and monitoring on the where, when, and how of data usage. In information security, this translates to [Auditing](#).



Application developers and IT teams can use security capabilities inherent in the security-related components of the Hadoop cluster as well as leverage external tools to encompass all aspects of Hadoop security. The Hadoop ecosystem covers a wide range of applications, datastores, and computing frameworks, and each of these security components manifest these operational capabilities differently.

Levels of Hadoop Security

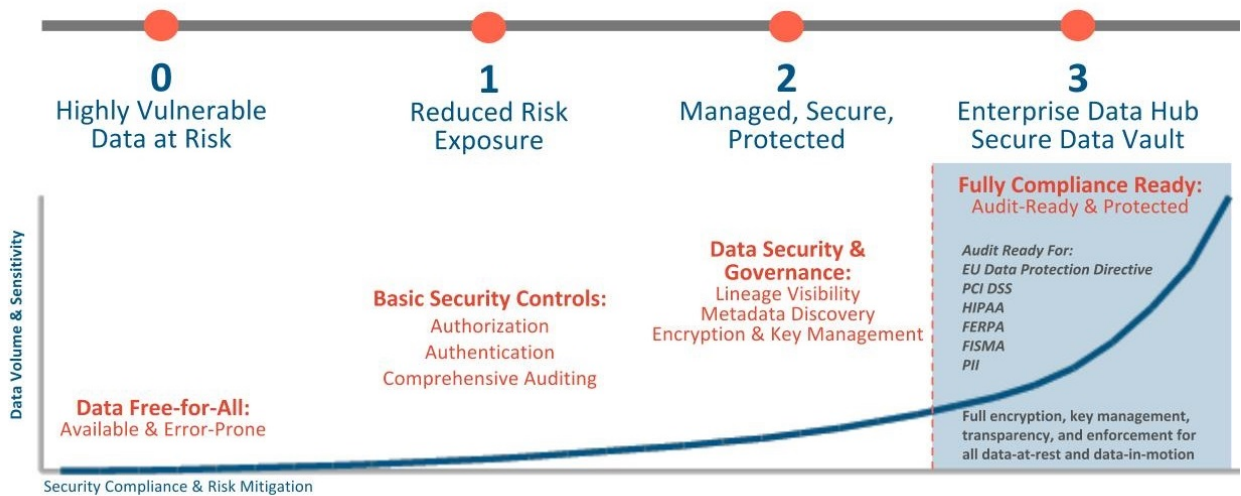
Security can typically be added to a Hadoop cluster in several stages:

Security Overview for an Enterprise Data Hub

- **Level 0:** Cloudera recommends you first start with a fully-functional non-secure cluster before you begin adding security to it. However, a non-secure cluster is very vulnerable to attacks and must never be used in production.
- **Level 1:** Start with the basics. First, set up authentication checks to prove that users/services accessing the cluster are who they claim to be. You can then implement some simple authorization mechanisms that allow you to assign access privileges to users and user groups. Auditing procedures to keep a check on who accesses the cluster and how, can also be added in this step. However, these are still very basic security measures. If you do go to production with only authentication, authorization and auditing enabled, make sure your cluster administrators are well trained and that the security procedures in place have been certified by an expert.
- **Level 2:** For more robust security, cluster data, or at least sensitive data, must be encrypted. There should be key-management systems in place for managing encrypted data and the associated encryption keys. Data governance is an important aspect of security. Governance includes auditing accesses to data residing in metastores, reviewing and updating metadata, and discovering the lineage of data objects.
- **Level 3:** At this level, all data on the cluster, at-rest and in-transit, must be encrypted, and the key management system in use must be fault-tolerant. A completely secure enterprise data hub (EDH) is one that can stand up to the audits required for compliance with PCI, HIPAA and other common industry standards. The regulations associated with these standards do not apply just to the EDH storing this data. Any system that integrates with the EDH in question is subject to scrutiny as well.

Leveraging all four levels of security, Cloudera's EDH platform can pass technology reviews for most common compliance regulations.

Achieve Scale and Cost Effectiveness via a Secure Data Vault

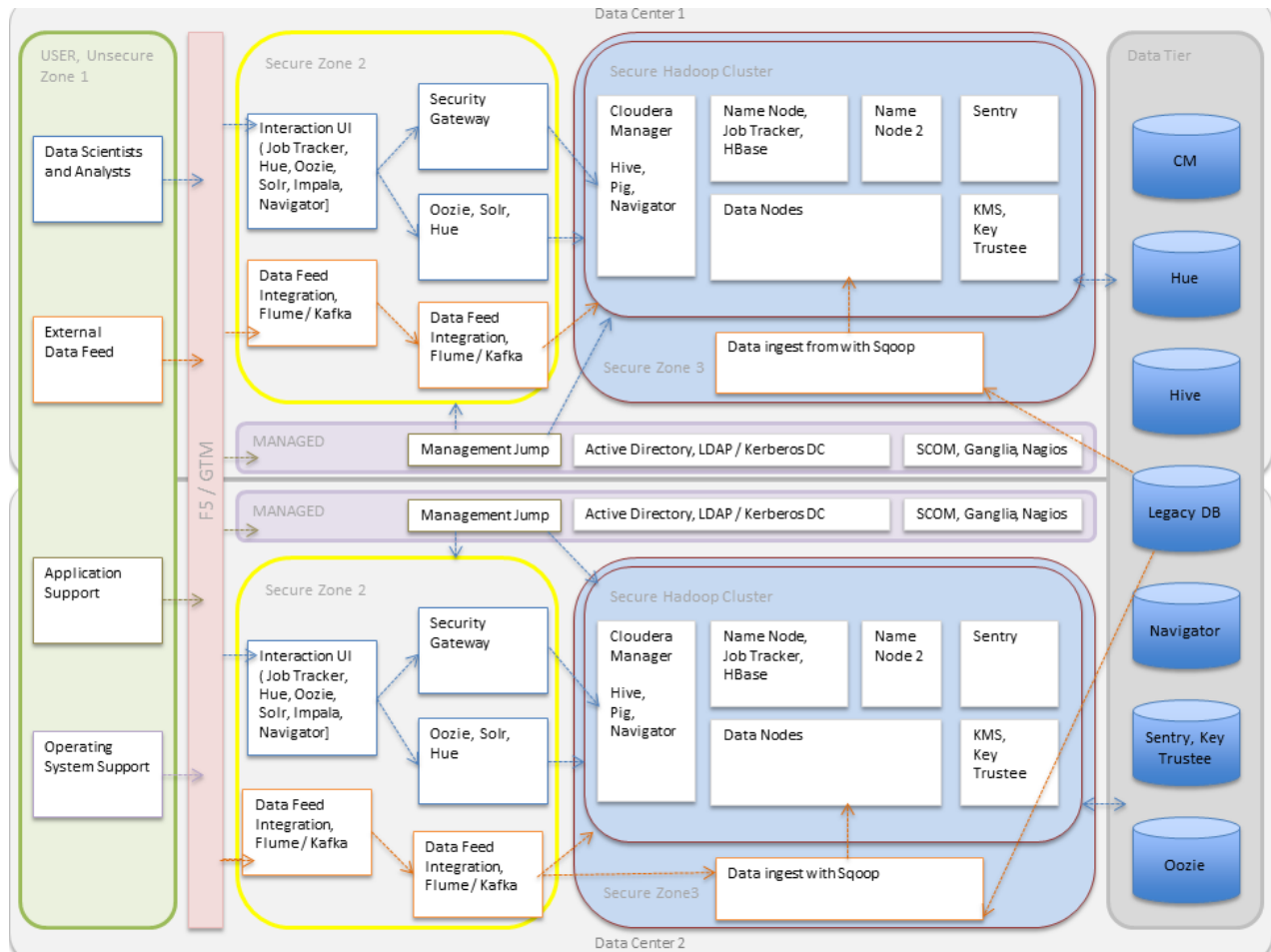


Hadoop Security Architecture

What follows is a detailed depiction of the Hadoop ecosystem in particular as it shows the interactions between different Cloudera Enterprise, security, and user management components. It also shows how a production environment with a couple of datacenters and assorted users and data feeds, both internal and external, will need to deal with receiving and authenticating so many insecure connections.

- As illustrated, external data streams can be authenticated by mechanisms in place for Flume and Kafka. Any data from legacy databases is ingested using Sqoop. Users such as data scientists and analysts can interact directly with the cluster using interfaces such as Hue or Cloudera Manager. Alternatively, they could be using a service like Impala for creating and submitting jobs for data analysis. All of these interactions can be protected by an Active Directory Kerberos deployment.

- Encryption can be applied to data at-rest using transparent HDFS encryption with an enterprise-grade Key Trustee Server. Cloudera also recommends using Navigator Encrypt to protect data on a cluster associated with the Cloudera Manager, Cloudera Navigator, Hive and HBase metastores, and any log files or spills.
- Authorization policies can be enforced using Sentry (for services such as Hive, Impala and Search) as well as HDFS Access Control Lists.
- Auditing capabilities can be provided by using Cloudera Navigator.



Overview of Authentication Mechanisms for an Enterprise Data Hub

The purpose of authentication in Hadoop is simply to prove that users or services are who they claim to be. Typically, authentication for CDH applications is handled using Kerberos, an enterprise-grade authentication protocol. Kerberos provides strong security benefits including capabilities that render intercepted authentication packets unusable by an attacker. It virtually eliminates the threat of impersonation by never sending a user's credentials in cleartext over the network. Several components of the Hadoop ecosystem are converging to use Kerberos authentication with the option to manage and store credentials in Active Directory (AD) or Lightweight Directory Access Protocol (LDAP). The rest of this topic describes in detail some basic Kerberos concepts as they relate to Hadoop, and the different ways in which Kerberos can be deployed on a CDH cluster.

For enterprise components such as Cloudera Manager, Cloudera Navigator, Hue, Hive, and Impala, that face external clients, Cloudera also supports external authentication using services such as AD, LDAP, SAML, and so on. For instructions on how to configure external authentication, refer the [Cloudera Security Guide](#).

Basic Kerberos Concepts



Important: If you are not familiar with how the Kerberos authentication protocol works, make sure you refer the following links before you proceed:

- [MIT Kerberos Home](#)
- [MIT Kerberos Documentation](#)

This section describes how Hadoop uses Kerberos principals and keytabs for user authentication. It also briefly describes how Hadoop uses delegation tokens to authenticate jobs at execution time, to avoid overwhelming the KDC with authentication requests for each job.

Kerberos Principals

A user in Kerberos is called a **principal**, which is made up of three distinct components: the primary, instance, and realm. A **Kerberos principal** is used in a Kerberos-secured system to represent a unique identity. The first component of the principal is called the **primary**, or sometimes the user component. The primary component is an arbitrary string and may be the operating system username of the user or the name of a service. The primary component is followed by an optional section called the **instance**, which is used to create principals that are used by users in special roles or to define the host on which a service runs, for example. An instance, if it exists, is separated from the primary by a slash and then the content is used to disambiguate multiple principals for a single user or service. The final component of the principal is the realm. The **realm** is similar to a domain in DNS in that it logically defines a related group of objects, although rather than hostnames as in DNS, the Kerberos realm defines a group of principals. Each realm can have its own settings including the location of the KDC on the network and supported encryption algorithms. Large organizations commonly create distinct realms to delegate administration of a realm to a group within the enterprise. Realms, by convention, are written in uppercase characters.

Kerberos assigns tickets to Kerberos principals to enable them to access Kerberos-secured Hadoop services. For the Hadoop daemon principals, the principal names should be of the format `username/fully.qualified.domain.name@YOUR-REALM.COM`. In this guide, *username* in the `username/fully.qualified.domain.name@YOUR-REALM.COM` principal refers to the username of an existing Unix account that is used by Hadoop daemons, such as `hdfs` or `mapred`. Human users who want to access the Hadoop cluster also need to have Kerberos principals; in this case, *username* refers to the username of the user's Unix account, such as `joe` or `jane`. Single-component principal names (such as `joe@YOUR-REALM.COM`) are acceptable for client user accounts. Hadoop does not support more than two-component principal names.

Kerberos Keytabs

A **keytab** is a file containing pairs of Kerberos principals and an encrypted copy of that principal's key. A keytab file for a Hadoop daemon is unique to each host since the principal names include the hostname. This file is used to authenticate a principal on a host to Kerberos without human interaction or storing a password in a plain text file. Because having access to the keytab file for a principal allows one to act as that principal, access to the keytab files should be tightly secured. They should be readable by a minimal set of users, should be stored on local disk, and should not be included in host backups, unless access to those backups is as secure as access to the local host.

Delegation Tokens

Users in a Hadoop cluster authenticate themselves to the NameNode using their Kerberos credentials. However, once the user is authenticated, each job subsequently submitted must also be checked to ensure it comes from an authenticated user. Since there could be a time gap between a job being submitted and the job being executed, during which the user could have logged off, user credentials are passed to the NameNode using delegation tokens that can be used for authentication in the future.

Delegation tokens are a secret key shared with the NameNode, that can be used to impersonate a user to get a job executed. While these tokens can be renewed, new tokens can only be obtained by clients authenticating to the NameNode using Kerberos credentials. By default, delegation tokens are only valid for a day. However, since jobs can last longer than a day, each token specifies a JobTracker as a *renewer* which is allowed to renew the delegation token

once a day, until the job completes, or for a maximum period of 7 days. When the job is complete, the JobTracker requests the NameNode to cancel the delegation token.

Token Format

The NameNode uses a random `masterKey` to generate delegation tokens. All active tokens are stored in memory with their expiry date (`maxDate`). Delegation tokens can either expire when the current time exceeds the expiry date, or, they can be canceled by the owner of the token. Expired or canceled tokens are then deleted from memory. The `sequenceNumber` serves as a unique ID for the tokens. The following section describes how the Delegation Token is used for authentication.

```
TokenID = {ownerID, renewerID, issueDate, maxDate, sequenceNumber}
TokenAuthenticator = HMAC-SHA1(masterKey, TokenID)
Delegation Token = {TokenID, TokenAuthenticator}
```

Authentication Process

To begin the authentication process, the client first sends the `TokenID` to the NameNode. The NameNode uses this `TokenID` and the `masterKey` to once again generate the corresponding `TokenAuthenticator`, and consequently, the Delegation Token. If the NameNode finds that the token already exists in memory, and that the current time is less than the expiry date (`maxDate`) of the token, then the token is considered valid. If valid, the client and the NameNode will then authenticate each other by using the `TokenAuthenticator` that they possess as the secret key, and MD5 as the protocol. Since the client and NameNode do not actually exchange `TokenAuthenticators` during the process, even if authentication fails, the tokens are not compromised.

Token Renewal

Delegation tokens must be renewed periodically by the designated renewer (`renewerID`). For example, if a JobTracker is the designated renewer, the JobTracker will first authenticate itself to the NameNode. It will then send the token to be authenticated to the NameNode. The NameNode verifies the following information before renewing the token:

- The JobTracker requesting renewal is the same as the one identified in the token by `renewerID`.
- The `TokenAuthenticator` generated by the NameNode using the `TokenID` and the `masterKey` matches the one previously stored by the NameNode.
- The current time must be less than the time specified by `maxDate`.

If the token renewal request is successful, the NameNode sets the new expiry date to `min(current time+renew period, maxDate)`. If the NameNode was restarted at any time, it will have lost all previous tokens from memory. In this case, the token will be saved to memory once again, this time with a new expiry date. Hence, designated renewers must renew all tokens with the NameNode after a restart, and before relaunching any failed tasks.

A designated renewer can also revive an expired or canceled token as long as the current time does not exceed `maxDate`. The NameNode cannot tell the difference between a token that was canceled, or has expired, and one that was erased from memory due to a restart, since only the `masterKey` persists in memory. The `masterKey` must be updated regularly.

Types of Kerberos Deployments

Kerberos provides strong security benefits including capabilities that render intercepted authentication packets unusable by an attacker. It virtually eliminates the threat of impersonation by never sending a user's credentials in cleartext over the network. Several components of the Hadoop ecosystem are converging to use Kerberos authentication with the option to manage and store credentials in LDAP or AD. Microsoft's Active Directory (AD) is an LDAP directory that also provides Kerberos authentication for added security. Before you configure Kerberos on your cluster, ensure you have a working KDC (MIT KDC or Active Directory), set up. You can then use Cloudera Manager's Kerberos wizard to automate several aspects of configuring Kerberos authentication on your cluster.

Without Kerberos enabled, Hadoop only checks to ensure that a user and their group membership is valid in the context of HDFS. However, it makes no effort to verify that the user is who they say they are.

With Kerberos enabled, users must first authenticate themselves to a Kerberos Key Distribution Centre (KDC) to obtain a valid Ticket-Granting-Ticket (TGT). The TGT is then used by Hadoop services to verify the user's identity. With Kerberos, a user is not only authenticated on the system they are logged into, but they are also authenticated to the network. Any subsequent interactions with other services that have been configured to allow Kerberos authentication for user access, are also secured.

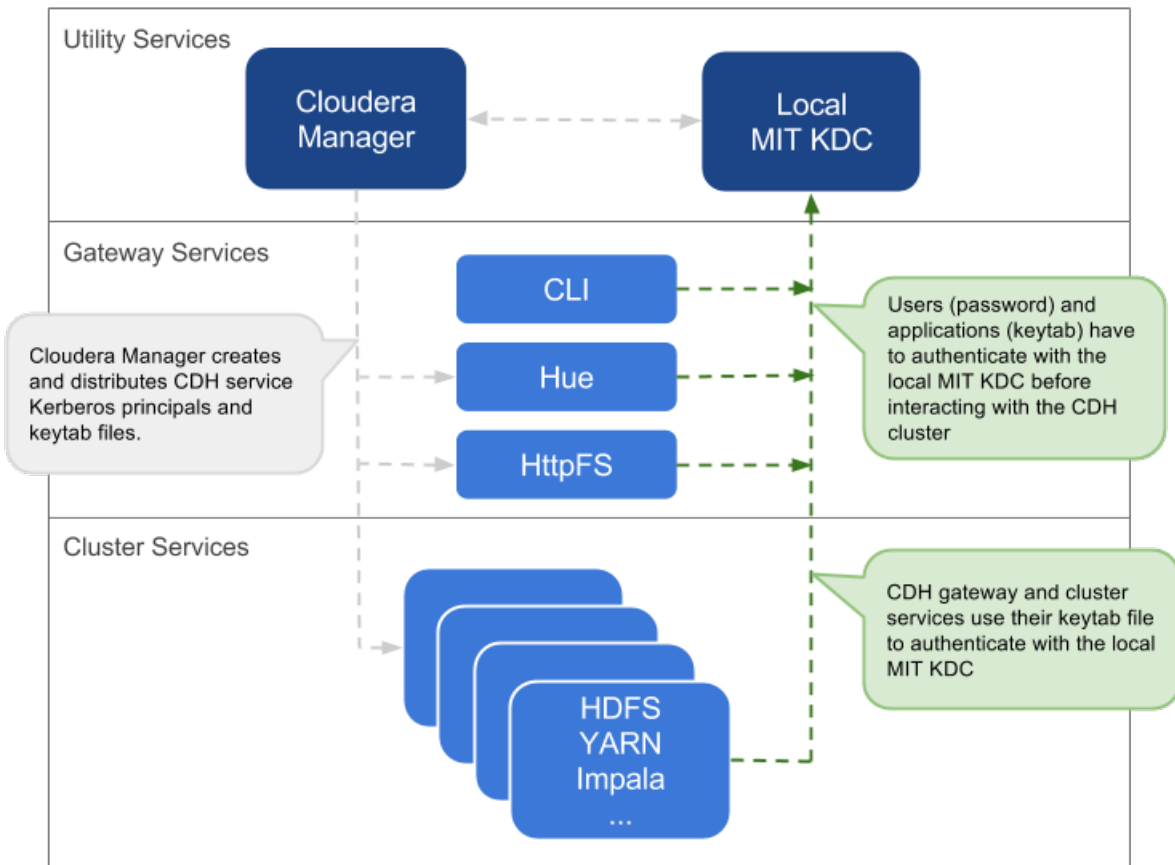
This section describes the architectural options that are available for deploying Hadoop security in enterprise environments. Each option includes a high-level description of the approach along with a list of pros and cons. There are three options currently available:

Local MIT KDC

This approach uses an MIT KDC that is local to the cluster. Users and services will have to authenticate with this local KDC before they can interact with the CDH components on the cluster.

Architecture Summary:

- An MIT KDC and a distinct Kerberos realm is deployed locally to the CDH cluster. The local MIT KDC is typically deployed on a Utility host. Additional replicated MIT KDCs for high-availability are optional.
- All cluster hosts must be configured to use the local MIT Kerberos realm using the `krb5.conf` file.
- All **service and user principals** must be created in the local MIT KDC and Kerberos realm.
- The local MIT KDC will authenticate both the service principals (using keytab files) and user principals (using passwords).
- Cloudera Manager connects to the local MIT KDC to create and manage the principals for the CDH services running on the cluster. To do this Cloudera Manager uses an admin principal and keytab that is created during the security setup. This step has been automated by the Kerberos wizard. Instructions for manually creating the Cloudera Manager admin principal are provided in the Cloudera Manager [security documentation](#).
- Typically, the local MIT KDC administrator is responsible for creating all other user principals. If you use the Kerberos wizard, Cloudera Manager will create these principals and associated keytab files for you.



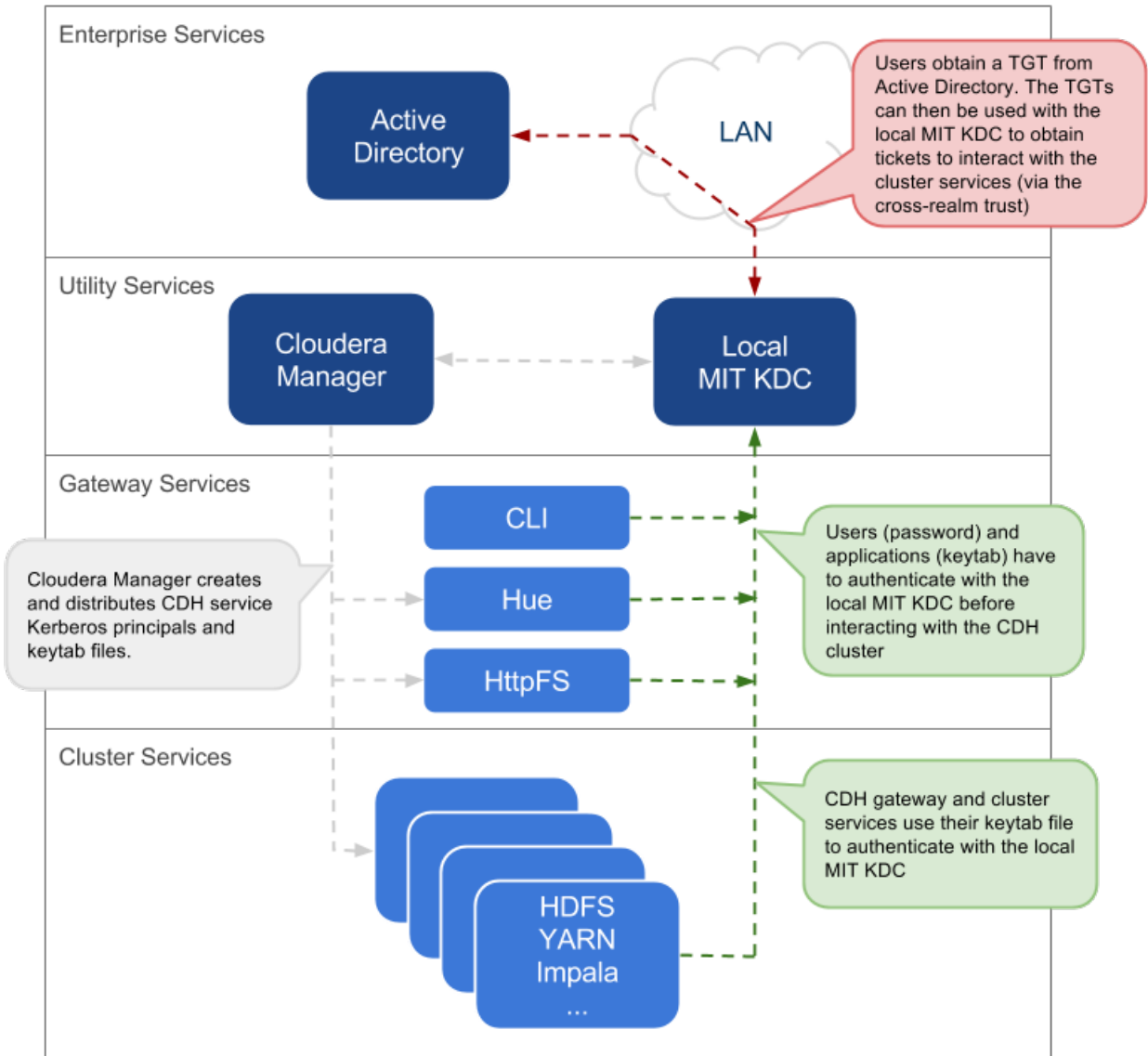
Pros	Cons
The authentication mechanism is isolated from the rest of the enterprise.	This mechanism is not integrated with central authentication system.
This is fairly easy to setup, especially if you use the Cloudera Manager Kerberos wizard that automates creation and distribution of service principals and keytab files.	User and service principals must be created in the local MIT KDC, which can be time-consuming.
	The local MIT KDC can be a single point of failure for the cluster unless replicated KDCs can be configured for high-availability.
	The local MIT KDC is yet another authentication system to manage.

Local MIT KDC with Active Directory Integration

This approach uses an MIT KDC and Kerberos realm that is local to the cluster. However, Active Directory stores the user principals that will access the cluster in a central realm. Users will have to authenticate with this central AD realm to obtain TGTs before they can interact with CDH services on the cluster. Note that CDH service principals reside only in the local KDC realm.

Architecture Summary:

- An MIT KDC and a distinct Kerberos realm is deployed locally to the CDH cluster. The local MIT KDC is typically deployed on a Utility host and additional replicated MIT KDCs for high-availability are optional.
- All cluster hosts are configured with both Kerberos realms (local and central AD) using the `krb5.conf` file. The default realm should be the local MIT Kerberos realm.
- **Service principals** should be created in the local MIT KDC and the local Kerberos realm. Cloudera Manager connects to the local MIT KDC to create and manage the principals for the CDH services running on the cluster. To do this, Cloudera Manager uses an admin principal and keytab that is created during the security setup. This step has been automated by the Kerberos wizard. Instructions for manually creating the Cloudera Manager admin principal are provided in the Cloudera Manager [security documentation](#).
- A one-way, cross-realm trust must be set up from the local Kerberos realm to the central AD realm containing the **user principals** that require access to the CDH cluster. There is no need to create the service principals in the central AD realm and no need to create user principals in the local realm.



Pros	Cons
<p>The local MIT KDC serves as a shield for the central Active Directory from the many hosts and services in a CDH cluster. Service restarts in a large cluster create many simultaneous authentication requests. If Active Directory is unable to handle the spike in load, then the cluster can effectively cause a distributed denial of service (DDOS) attack.</p>	<p>The local MIT KDC can be a single point of failure (SPOF) for the cluster. Replicated KDCs can be configured for high-availability.</p>
<p>This is fairly easy to setup, especially if you use the Cloudera Manager Kerberos wizard that automates creation and distribution of service principals and keytab files.</p> <p>Active Directory administrators will only need to be involved to configure the cross-realm trust during setup.</p>	<p>The local MIT KDC is yet another authentication system to manage.</p>

Pros	Cons
Integration with central Active Directory for user principal authentication results in a more complete authentication solution.	
Allows for incremental configuration. Hadoop security can be configured and verified using local MIT KDC independently of integrating with Active Directory.	

Direct to Active Directory

This approach uses the central Active Directory as the KDC. No local KDC is required. Before you decide upon an AD KDC deployment, make sure you are aware of the following possible ramifications of that decision.

Considerations when using an Active Directory KDC

Performance:

As your cluster grows, so will the volume of Authentication Service (AS) and Ticket Granting Service (TGS) interaction between the services on each cluster server. Consider evaluating the volume of this interaction against the Active Directory domain controllers you have configured for the cluster before rolling this feature out to a production environment. If cluster performance suffers, over time it might become necessary to dedicate a set of AD domain controllers to larger deployments.

The central Active Directory KDC is not shielded from the CDH cluster. Service restarts in a large cluster create many simultaneous authentication requests. If Active Directory is unable to handle the spike in load, then the cluster can effectively cause a distributed denial of service (DDOS) attack.

Network Proximity:

By default, Kerberos uses UDP for client/server communication. Often, AD services are in a different network than project application services such as Hadoop. If the domain controllers supporting a cluster for Kerberos are not in the same subnet, or they're separated by a firewall, consider using the `udp_preference_limit = 1` setting in the `[libdefaults]` section of the `krb5.conf` used by cluster services. Cloudera strongly recommends *against* using AD domain controller (KDC) servers that are separated from the cluster by a WAN connection, as latency in this service will significantly impact cluster performance.

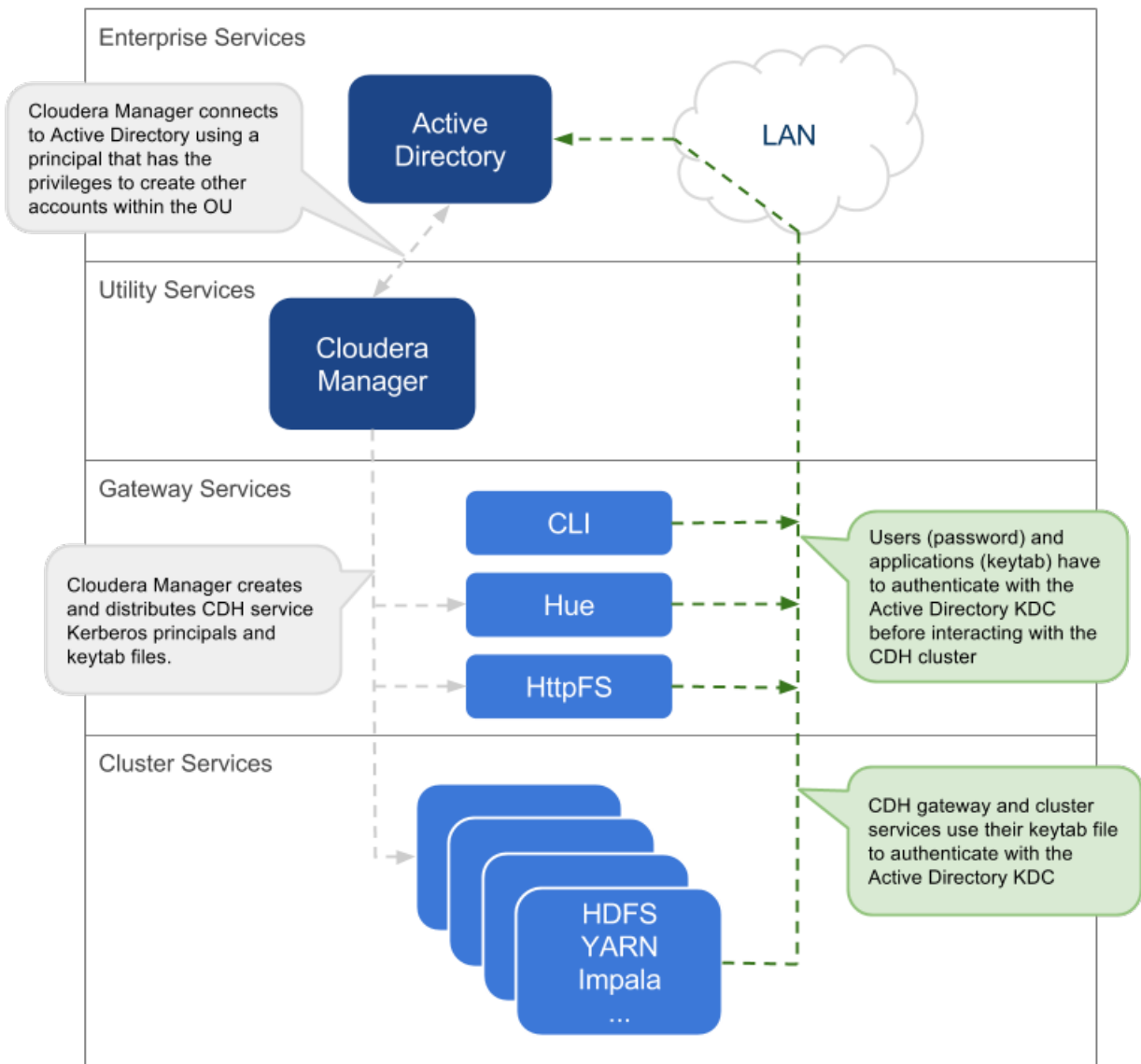
Process:

Troubleshooting the cluster's operations, especially for Kerberos-enabled services, will need to include AD administration resources. Evaluate your organizational processes for engaging the AD administration team, and how to escalate in case a cluster outage occurs due to issues with Kerberos authentication against AD services. In some situations it might be necessary to [enable Kerberos event logging](#) to address desktop and KDC issues within windows environments.

Also note that if you decommission any Cloudera Manager roles or nodes, the related AD accounts will need to be deleted manually. This is required because Cloudera Manager will not delete existing entries in Active Directory.

Architecture Summary

- All **service and user principals** are created in the Active Directory KDC.
- All cluster hosts are configured with the central AD Kerberos realm using `krb5.conf`.
- Cloudera Manager connects to the Active Directory KDC to create and manage the principals for the CDH services running on the cluster. To do this, Cloudera Manager uses a principal that has the privileges to create other accounts within the given Organisational Unit (OU). This step has been automated by the Kerberos wizard. Instructions for manually creating the Cloudera Manager admin principal are provided in the Cloudera Manager [security documentation](#).
- All service and user principals are authenticated by the Active Directory KDC.



Note: If it is not possible to create the Cloudera Manager admin principal with the required privileges in the Active Directory KDC, then the CDH services principals will need to be created manually. The corresponding keytab files should then be stored securely on the Cloudera Manager Server host. Cloudera Manager's [Custom Kerberos Keytab Retrieval script](#) can be used to retrieve the keytab files from the local filesystem.

Identity Integration with Active Directory

A core requirement for enabling Kerberos security in the platform is that users have accounts on all cluster processing nodes. Commercial products such as Centrify or Quest Authentication Services (QAS) provide integration of all cluster hosts for user and group resolution to Active Directory. These tools support automated Kerberos authentication on login by users to a Linux host with AD. For sites not using Active Directory, or sites wanting to use an open source solution, the Site Security Services Daemon (SSSD) can be used with either AD or OpenLDAP compatible directory services and MIT Kerberos for the same needs.

For third-party providers, you may have to purchase licenses from the respective vendors. This procedure requires some planning as it takes time to procure these licenses and deploy these products on a cluster. Care should be taken to ensure that the identity management product does not associate the service principal names (SPNs) with the host

principals when the computers are joined to the AD domain. For example, Centrifyp by default associates the HTTP SPN with the host principal. So the HTTP SPN should be specifically excluded when the hosts are joined to the domain.

You will also need to complete the following setup tasks in AD:

- **Active Directory Organizational Unit (OU) and OU user** - A separate OU in Active Directory should be created along with an account that has privileges to create additional accounts in that OU.
- **Enable SSL for AD** - Cloudera Manager should be able to connect to AD on the LDAPS (TCP 636) port.
- **Principals and Keytabs** - In a direct-to-AD deployment that is set up using the Kerberos wizard, by default, all required principals and keytabs will be created, deployed and managed by Cloudera Manager. However, if for some reason you cannot allow Cloudera Manager to manage your direct-to-AD deployment, then unique accounts should be manually created in AD for each service running on each host and keytab files must be provided for the same. These accounts should have the AD User Principal Name (UPN) set to `service/fqdn@REALM`, and the Service Principal Name (SPN) set to `service/fqdn`. The principal name in the keytab files should be the UPN of the account. The keytab files should follow the naming convention: `servicename_fqdn.keytab`. The following principals and keytab files must be created for each host they run on: [Hadoop Users in Cloudera Manager and CDH](#) on page 174.
- **AD Bind Account** - Create an AD account that will be used for LDAP bindings in Hue, Cloudera Manager and Cloudera Navigator.
- **AD Groups for Privileged Users** - Create AD groups and add members for the authorized users, HDFS admins and HDFS superuser groups.
 - Authorized users – A group consisting of all users that need access to the cluster
 - HDFS admins – Groups of users that will run HDFS administrative commands
 - HDFS super users – Group of users that require superuser privilege, that is, read/wwrite access to all data and directories in HDFS

Putting regular users into the HDFS superuser group is *not* recommended. Instead, an account that administrators escalate issues to, should be part of the HDFS superuser group.
- **AD Groups for Role-Based Access to Cloudera Manager and Cloudera Navigator** - Create AD groups and add members to these groups so you can later configure role-based access to Cloudera Manager and Cloudera Navigator. Cloudera Manager roles and their definitions are available here: [Cloudera Manager User Roles](#) on page 344. Cloudera Navigator roles and their definitions are available here: [Cloudera Navigator Data Management Component User Roles](#) on page 347
- **AD Test Users and Groups** - At least one existing AD user and the group that the user belongs to should be provided to test whether authorization rules work as expected.

TLS/SSL Requirements for Secure Distribution of Kerberos Keytabs

Communication between Cloudera Manager Server and Cloudera Manager Agents must be encrypted so that sensitive information such as Kerberos keytabs are not distributed from the Server to Agents in cleartext.

Cloudera Manager supports three levels of TLS security.

- **Level 1 (Good)** - This level encrypts communication between the browser and Cloudera Manager, and between Agents and the Cloudera Manager Server. See [Configuring TLS Encryption Only for Cloudera Manager](#) on page 212 followed by [Level 1: Configuring TLS Encryption for Cloudera Manager Agents](#) on page 216 for instructions. Level 1 encryption prevents snooping of commands and controls ongoing communication between Agents and Cloudera Manager.
- **Level 2 (Better)** - This level encrypts communication between the Agents and the Server, and provides strong verification of the Cloudera Manager Server certificate by Agents. See [Level 2: Configuring TLS Verification of Cloudera Manager Server by the Agents](#) on page 217. Level 2 provides Agents with additional security by verifying trust for the certificate presented by the Cloudera Manager Server.

- Level 3 (Best) - This includes encrypted communication between the Agents and the Server, strong verification of the Cloudera Manager Server certificate by the Agents, *and* authentication of Agents to the Cloudera Manager Server using self-signed or CA-signed certs. See [Level 3: Configuring TLS Authentication of Agents to the Cloudera Manager Server](#) on page 219. Level 3 TLS prevents cluster Servers from being spoofed by untrusted Agents running on a host. Cloudera recommends that you configure Level 3 TLS encryption for untrusted network environments before enabling Kerberos authentication. This provides secure communication of keytabs between the Cloudera Manager Server and verified Agents across the cluster.

This means, if you want to implement Level 3 TLS, you will need to provide TLS certificates for every host in the cluster. For minimal security, that is, Level 1 TLS, you will at least need to provide a certificate for the Cloudera Manager Server host, and a certificate for each of the gateway nodes to secure the web consoles.

If the CA that signs these certificates is an internal CA, then you will also need to provide the complete certificate chain of the CA that signed these certificates. The same certificates can be used to encrypt Cloudera Manager, Hue, HiveServer2 & Impala JDBC/ODBC interfaces, and for encrypted shuffle. If any external services such as LDAPS or SAML use certificates signed by an internal CA, then the public certificate of the Root CA and any intermediate CA in the chain should be provided.

Configuring Kerberos Authentication on a Cluster

Before you use the following sections to configure Kerberos on your cluster, ensure you have a working KDC (MIT KDC or Active Directory), set up.

You can use *one* of the following ways to set up Kerberos authentication on your cluster using Cloudera Manager:

- Cloudera Manager 5.1 introduced a new wizard to automate the procedure to set up Kerberos on a cluster. Using the KDC information you enter, the wizard will create new principals and keytab files for your CDH services. The wizard can be used to deploy the `krb5.conf` file cluster-wide, and automate other manual tasks such as stopping all services, deploying client configuration and restarting all services on the cluster.

If you want to use the Kerberos wizard, follow the instructions at [Enabling Kerberos Authentication Using the Wizard](#) on page 59.

- If you do not want to use the Kerberos wizard, follow the instructions at [Enabling Kerberos Authentication Without the Wizard](#) on page 75.

Authentication Mechanisms used by Hadoop Projects

Project	Authentication Capabilities
HDFS	Kerberos, SPNEGO (HttpFS)
MapReduce	Kerberos (<i>also see HDFS</i>)
YARN	Kerberos (<i>also see HDFS</i>)
Accumulo	Kerberos (partial)
Flume	Kerberos (starting CDH 5.4)
HBase	Kerberos (<i>HBase Thrift and REST clients must perform their own user authentication</i>)
HiveServer	None
HiveServer2	Kerberos, LDAP, Custom/pluggable authentication
Hive Metastore	Kerberos
Hue	Kerberos, LDAP, SAML, Custom/pluggable authentication
Impala	Kerberos, LDAP, SPNEGO (Impala Web Console)
Oozie	Kerberos, SPNEGO

Project	Authentication Capabilities
Pig	Kerberos
Search	Kerberos, SPNEGO
Sentry	Kerberos
Spark	Kerberos
Sqoop	Kerberos
Sqoop2	Kerberos (starting CDH 5.4)
Zookeeper	Kerberos
Cloudera Manager	Kerberos, LDAP, SAML
Cloudera Navigator	<i>See Cloudera Manager</i>
Backup and Disaster Recovery	<i>See Cloudera Manager</i>

Overview of Data Protection Mechanisms for an Enterprise Data Hub

The goal of data protection is to ensure that only authorized users can view, use, or contribute to a data set. These security controls add another layer of protection against potential threats by end-users, administrators and any other potentially-malicious actors on the network. The means to achieving this goal consists of two parts: protecting data when it is persisted to disk or other storage mediums, commonly called *data-at-rest*, and protecting data while it moves from one process or system to another, that is, *data in transit*.

Several common compliance regulations call for data protection in addition to the other security controls discussed in this paper. Cloudera provides security for data in transit, through TLS and other mechanisms - centrally deployed through Cloudera Manager. Transparent data-at-rest protection is provided through the combination of transparent HDFS encryption, Navigator Encrypt, and Navigator Key Trustee.

Protecting Data At-Rest

For data-at-rest, there are several risk avenues that may lead to unprivileged access, such as exposed data on a hard drive after a mechanical failure, or, data being accessed by previously-compromised accounts. Data-at-rest is also complicated by longevity; where and how does one store the keys to decrypt data throughout an arbitrary period of retention. Key management issues, such as mitigating compromised keys and the re-encryption of data, also complicate the process.

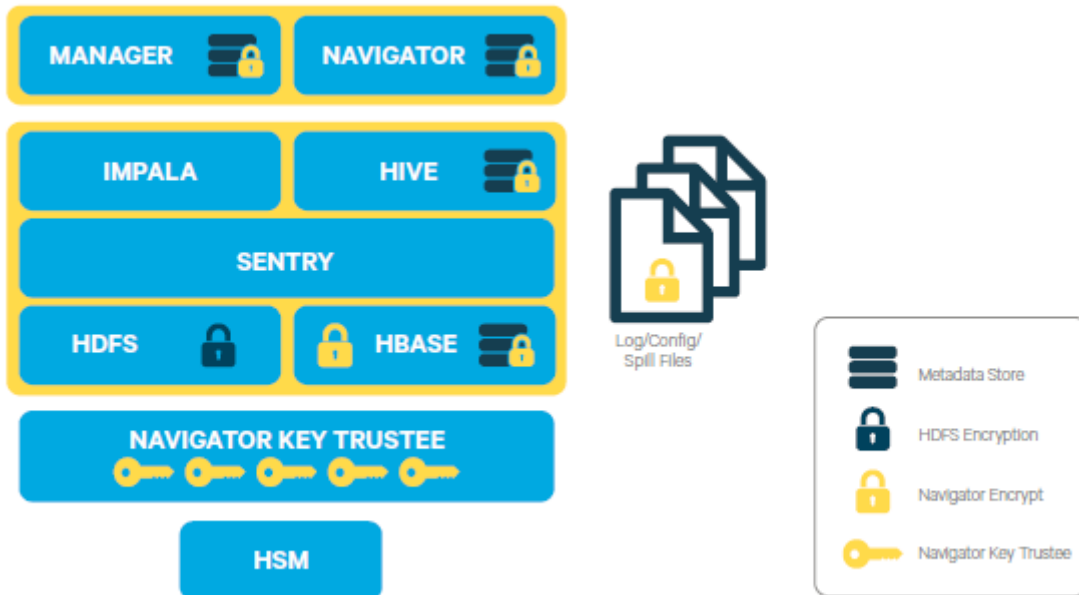
However, in some cases, you might find that relying on encryption to protect your data is not sufficient. With encryption, sensitive data may still be exposed to an administrator who has complete access to the cluster. Even users with appropriate ACLs on the data could have access to logs and queries where sensitive data might have leaked. To prevent such leaks, Cloudera now allows you to mask personally identifiable information (PII) from log files, audit data and SQL queries.

Encryption Options Available with Hadoop

Hadoop, as a central data hub for many kinds of data within an organization, naturally has different needs for data protection depending on the audience and processes while using a given data set. CDH provides transparent HDFS encryption, ensuring that all sensitive data is encrypted before being stored on disk. This capability, together with enterprise-grade encryption key management with Navigator Key Trustee, delivers the necessary protection to meet regulatory compliance for most enterprises. HDFS Encryption together with Navigator Encrypt (available with Cloudera Enterprise) provides transparent encryption for Hadoop, for both data and metadata. These solutions automatically encrypt data while the cluster continues to run as usual, with a very low performance impact. It is massively scalable, allowing encryption to happen in parallel against all the data nodes - as the cluster grows, encryption grows with it.

Additionally, this transparent encryption is optimized for the Intel chipset for high performance. Intel chipsets include AES-NI co-processors, which provide special capabilities that make encryption workloads run extremely fast. Not all AES-NI is equal, and Cloudera is able to take advantage of the latest Intel advances for even faster performance.

As an example, the figure below shows a sample deployment that uses CDH's transparent HDFS encryption to protect the data stored in HDFS, while Navigator Encrypt is being used to protect all other data on the cluster associated with the Cloudera Manager, Cloudera Navigator, Hive and HBase metadata stores, along with any logs or spills. Navigator Key Trustee is being used to provide robust and fault-tolerant key management.



Encryption can be applied at a number of levels within Hadoop:

- **OS Filesystem-level** - Encryption can be applied at the Linux operating system filesystem level to cover all files in a volume. An example of this approach is [Cloudera Navigator Encrypt](#) on page 316 (formerly Gazzang zNcrypt) which is available for Cloudera customers licensed for Cloudera Navigator. Navigator Encrypt operates at the Linux volume level, so it can encrypt cluster data inside and outside HDFS, such as temp/spill files, configuration files and metadata databases (to be used only for data related to a CDH cluster). Navigator Encrypt must be used with [Cloudera Navigator Key Trustee Server](#) on page 297 (formerly Gazzang zTrustee).

CDH components, such as Impala, MapReduce, YARN, or HBase, also have the ability to encrypt data that lives temporarily on the local filesystem outside HDFS. To enable this feature, see [Configuring Encryption for Data Spills](#) on page 338.

- **Network-level** - Encryption can be applied to encrypt data just before it gets sent across a network and to decrypt it just after receipt. In Hadoop, this means coverage for data sent from client user interfaces as well as service-to-service communication like remote procedure calls (RPCs). This protection uses industry-standard protocols such as TLS/SSL.



Note: Cloudera Manager and CDH components support either TLS 1.0, TLS 1.1, or TLS 1.2, but not SSL 3.0. References to SSL only continue because of its widespread use in technical jargon.

- **HDFS-level** - Encryption applied by the HDFS client software. [HDFS Transparent Encryption](#) on page 261 operates at the HDFS folder level, allowing you to encrypt some folders and leave others unencrypted. HDFS transparent encryption cannot encrypt any data outside HDFS. To ensure reliable key storage (so that data is not lost), use Cloudera Navigator Key Trustee Server; the default Java keystore can be used for test purposes. For more information, see [Enabling HDFS Encryption Using Cloudera Navigator Key Trustee Server](#) on page 267.

Unlike OS and network-level encryption, HDFS transparent encryption is end-to-end. That is, it protects data at rest and in transit, which makes it more efficient than implementing a combination of OS-level and network-level encryption.

Data Redaction with Hadoop

Data redaction is the suppression of sensitive data, such as any personally identifiable information (PII). PII can be used on its own or with other information to identify or locate a single person, or to identify an individual in context. Enabling redaction allow you to transform PII to a pattern that does not contain any identifiable information. For example, you could replace all Social Security numbers (SSN) like 123-45-6789 with an unintelligible pattern like xxx-xx-xxxx, or replace only part of the SSN (xxx-xx-6789).

Although [encryption techniques](#) are available to protect Hadoop data, the underlying problem with using encryption is that an admin who has complete access to the cluster also access to unencrypted sensitive user data. Even users with appropriate ACLs on the data could have access to logs and queries where sensitive data might have leaked.

Data redaction provides compliance with industry regulations such as PCI and HIPAA, which require that access to PII be restricted to only those users whose jobs require such access. PII or other sensitive data must not be available through any other channels to users like cluster administrators or data analysts. However, if you already have permissions to access PII through queries, the query results will not be redacted. Redaction only applies to any incidental leak of data. Queries and query results must not show up in cleartext in logs, configuration files, UIs, or other unprotected areas.

Scope:

Data redaction in CDH targets sensitive SQL data and log files. Currently, you can enable or disable redaction for the whole cluster with a simple HDFS service-wide configuration change. Redaction is implemented with the assumption that sensitive information resides in the data itself, not the metadata. If you enable redaction for a file, only sensitive data inside the file is redacted. Metadata such as the name of the file or file owner is not redacted.

When data redaction is enabled, the following data is redacted:

- Logs in HDFS and any dependent cluster services. Log redaction is not available in Isilon-based clusters.
- Audit data sent to Cloudera Navigator
- SQL query strings displayed by Hue, Hive, and Impala.

For more information on enabling this feature, see [Sensitive Data Redaction](#) on page 437.

Password Redaction

Starting with Cloudera Manager and CDH 5.5, passwords will no longer be accessible in cleartext through the Cloudera Manager UI or in the configuration files stored on disk. For components such as HDFS, HBase, Hive, and so on, that use core Hadoop, the feature has been implemented by using Hadoop's `CredentialProvider` interface to encrypt and store passwords inside a secure `creds.jceks` keystore file. For components such as Hue and Impala, that do not use core Hadoop, instead of the password, we use a `password_script = /path/to/script/that/will/emit/password.sh` parameter that, when run, writes the password to `stdout`. Passwords contained within Cloudera Manager and Cloudera Navigator properties have been redacted internally in Cloudera Manager.

However, the database password contained in Cloudera Manager Server's `/etc/cloudera-scm-server/db.properties` file has not been redacted. The `db.properties` file is managed by customers and is populated manually when the Cloudera Manager Server database is being set up for the first time. Since this occurs before the Cloudera Manager Server has even started, encrypting the contents of this file is a completely different challenge as compared to that of redacting configuration files.

Password redaction (not including log and query redaction) is enabled by default for deployments with Cloudera Manager 5.5 (or higher) managing CDH 5.5 (or higher). There are no user-visible controls to enable or disable this feature. It is expected to work out of the box. The primary places where you will encounter the effects of password redaction are:

- In the Cloudera Manager Admin Console, on the **Processes** page for a given role instance, passwords in the linked configuration files have been replaced by `*****`.
- On the Cloudera Manager Server and Agent hosts, all configuration files in the `/var/run/cloudera-scm-agent/process` directory will have their passwords replaced by `*****`.

Exceptions:

- Solr does not have this feature enabled.
- The database password contained in Cloudera Manager Server's `/etc/cloudera-scm-server/db.properties` file has not been redacted.

Protecting Data In-Transit

For data-in-transit, implementing data protection and encryption is relatively easy. Wire encryption is built into the Hadoop stack, such as SSL, and typically does not require external systems. This data-in-transit encryption is built using session-level, one-time keys, by means of a session handshake with immediate and subsequent transmission. Thus, data-in-transit avoids much of the key management issues associated with data-at-rest due the temporal nature of the keys, but it does rely on proper authentication; a certificate compromise is an issue with authentication, but can compromise wire encryption. As the name implies, data-in-transit covers the secure transfer and intermediate storage of data. This applies to all process-to-process communication, within the same node or between nodes. There are three primary communication channels:

- **HDFS Transparent Encryption:** Data encrypted using [HDFS Transparent Encryption](#) on page 261 is protected end-to-end. Any data written to and from HDFS can only be encrypted or decrypted by the client. HDFS does not have access to the unencrypted data or the encryption keys. This supports both, at-rest encryption as well as in-transit encryption.
- **Data Transfer:** The first channel is data transfer, including the reading and writing of data blocks to HDFS. Hadoop uses a SASL-enabled wrapper around its native direct TCP/IP-based transport, called `DataTransportProtocol`, to secure the I/O streams within an DIGEST-MD5 envelope (For steps, see [Configuring Encrypted HDFS Data Transport](#) on page 341). This procedure also employs secured HadoopRPC (see Remote Procedure Calls) for the key exchange. The HttpFS REST interface, however, does not provide secure communication between the client and HDFS, only secured authentication using SPNEGO.

For the transfer of data between DataNodes during the shuffle phase of a MapReduce job (that is, moving intermediate results between the Map and Reduce portions of the job), Hadoop secures the communication channel with HTTP Secure (HTTPS) using Transport Layer Security (TLS). See [Encrypted Shuffle and Encrypted Web UIs](#) on page 250.

- **Remote Procedure Calls:** The second channel is system calls to remote procedures (RPC) to the various systems and frameworks within a Hadoop cluster. Like data transfer activities, Hadoop has its own native protocol for RPC, called HadoopRPC, which is used for Hadoop API client communication, intra-Hadoop services communication, as well as monitoring, heartbeats, and other non-data, non-user activity. HadoopRPC is SASL-enabled for secured transport and defaults to Kerberos and DIGEST-MD5 depending on the type of communication and security settings. For steps, see [Configuring Encrypted HDFS Data Transport](#) on page 341.
- **User Interfaces:** The third channel includes the various web-based user interfaces within a Hadoop cluster. For secured transport, the solution is straightforward; these interfaces employ HTTPS.

SSL/TLS Certificates Overview

Certificates can be signed in one three different ways:

Type	Usage Note
Public CA-signed certificates	Recommended. Using certificates signed by a trusted public CA simplifies deployment because the default Java client already trusts most public CAs. Obtain certificates from one of the trusted well-known (public) CAs, such as Symantec and Comodo, as detailed in Generate TLS Certificates on page 40
Internal CA-signed certificates	Obtain certificates from your organization's internal CA if your organization has its own. Using an internal CA can reduce costs (although cluster configuration may require establishing

Type	Usage Note
	the trust chain for certificates signed by an internal CA, depending on your IT infrastructure). See How to Configure TLS Encryption for Cloudera Manager on page 40 for information about establishing trust as part of configuring a Cloudera Manager cluster.
Self-signed certificates	Not recommended for production deployments. Using self-signed certificates requires configuring each client to trust the specific certificate, in addition to generating and distributing the certificates. However, to use self-signed certificates in a non-production (test, lab, proof-of-concept) deployment, see Using Self-Signed Certificates for TLS on page 226.



Note: Wildcard domain certificates and certificates using the SubjectAlternativeName extension are not supported at this time.

For more information on setting up SSL/TLS certificates, continue reading the topics at [TLS/SSL Certificates Overview](#) on page 205.

TLS Encryption Levels for Cloudera Manager

Transport Layer Security (TLS) provides encryption and authentication in communication between the Cloudera Manager Server and Agents. Encryption prevents snooping, and authentication helps prevent problems caused by malicious servers or agents.

Cloudera Manager supports three levels of TLS security.

- Level 1 (Good) - This level encrypts communication between the browser and Cloudera Manager, and between Agents and the Cloudera Manager Server. See [Configuring TLS Encryption Only for Cloudera Manager](#) on page 212 followed by [Level 1: Configuring TLS Encryption for Cloudera Manager Agents](#) on page 216 for instructions. Level 1 encryption prevents snooping of commands and controls ongoing communication between Agents and Cloudera Manager.
- Level 2 (Better) - This level encrypts communication between the Agents and the Server, and provides strong verification of the Cloudera Manager Server certificate by Agents. See [Level 2: Configuring TLS Verification of Cloudera Manager Server by the Agents](#) on page 217. Level 2 provides Agents with additional security by verifying trust for the certificate presented by the Cloudera Manager Server.
- Level 3 (Best) - This includes encrypted communication between the Agents and the Server, strong verification of the Cloudera Manager Server certificate by the Agents, *and* authentication of Agents to the Cloudera Manager Server using self-signed or CA-signed certs. See [Level 3: Configuring TLS Authentication of Agents to the Cloudera Manager Server](#) on page 219. Level 3 TLS prevents cluster Servers from being spoofed by untrusted Agents running on a host. Cloudera recommends that you configure Level 3 TLS encryption for untrusted network environments before enabling Kerberos authentication. This provides secure communication of keytabs between the Cloudera Manager Server and verified Agents across the cluster.

TLS/SSL Encryption for CDH Components

Cloudera recommends securing a cluster using Kerberos authentication before enabling encryption such as SSL on a cluster. If you enable SSL for a cluster that does not already have Kerberos authentication configured, a warning will be displayed.

Hadoop services differ in their use of SSL as follows:

- HDFS, MapReduce, and YARN daemons act as both SSL servers and clients.
- HBase daemons act as SSL servers only.
- Oozie daemons act as SSL servers only.
- Hue acts as an SSL client to all of the above.

Daemons that act as SSL servers load the keystores when starting up. When a client connects to an SSL server daemon, the server transmits the certificate loaded at startup time to the client, which then uses its truststore to validate the server's certificate.

For information on setting up SSL/TLS for CDH services, see [Configuring TLS/SSL Encryption for CDH Services](#) on page 228.

Data Protection within Hadoop Projects

The table below lists the various encryption capabilities that can be leveraged by CDH components and Cloudera Manager.

Project	Encryption for Data-in-Transit	Encryption for Data-at-Rest (HDFS Encryption + Navigator Encrypt + Navigator Key Trustee)
HDFS	SASL (RPC), SASL (DataTransferProtocol)	Yes
MapReduce	SASL (RPC), HTTPS (encrypted shuffle)	Yes
YARN	SASL (RPC)	Yes
Accumulo	Partial - Only for RPCs and Web UI (Not directly configurable in Cloudera Manager)	Yes
Flume	TLS (Avro RPC)	Yes
HBase	SASL - For web interfaces, inter-component replication, the HBase shell and the REST, Thrift 1 and Thrift 2 interfaces	Yes
HiveServer2	SASL (Thrift), SASL (JDBC), TLS (JDBC, ODBC)	Yes
Hue	TLS	Yes
Impala	TLS or SASL between impalad and clients, but not between daemons	
Oozie	TLS	Yes
Pig	N/A	Yes
Search	TLS	Yes
Sentry	SASL (RPC)	Yes
Spark	None	Yes
Sqoop	Partial - Depends on the RDBMS database driver in use	Yes
Sqoop2	Partial - You can encrypt the JDBC connection depending on the RDBMS database driver	Yes
ZooKeeper	SASL (RPC)	No
Cloudera Manager	TLS - Does not include monitoring	Yes
Cloudera Navigator	TLS - <i>Also see Cloudera Manager</i>	Yes
Backup and Disaster Recovery	TLS - <i>Also see Cloudera Manager</i>	Yes

Overview of Authorization Mechanisms for an Enterprise Data Hub

Authorization is concerned with who or what has access or control over a given resource or service. Within Hadoop, there are many resources and services ranging from computing frameworks to client applications and, of course, an unbounded amount and number of types of data, given the flexibility and scalability of the underlying HDFS storage

system. However, the authorization risks facing enterprises stems directly from Hadoop's storage and computing breadth and richness. Many audiences can and want to use, contribute, and view the data within Hadoop, and this trend will accelerate within organizations as Hadoop becomes critical to the data management infrastructure. Initially, an IT organization might employ a Hadoop cluster for ETL processing with only a handful of developers within the organization requesting access to the data and services. Yet as other teams, like line-of-business analysts, recognize the utility of Hadoop and its data, they may demand new computing needs, like interactive SQL or search, as well as request the addition of business-sensitive data within Hadoop. Since Hadoop merges together the capabilities of multiple varied, and previously separate IT systems as an enterprise data hub that stores and works on all data within an organization, it requires multiple authorization controls with varying granularities. In such cases, Hadoop management tools simplify setup and maintenance by:

- Tying all users to groups, which can be specified in existing LDAP or AD directories.
- Providing role-based access control for similar interaction methods, like batch and interactive SQL queries. For example, Apache Sentry permissions apply to Hive (HiveServer2) and Impala.

Authorization Mechanisms in Hadoop

Authorization for data access in Hadoop typically manifests in three forms:

- POSIX-style permissions on files and directories
- Access Control Lists (ACL) for management of services and resources
- Role-Based Access Control (RBAC) for certain services with advanced access controls to data.

Accordingly, CDH currently provides the following forms of access control:

- Traditional POSIX-style permissions for directories and files, where each directory and file is assigned a single owner and group. Each assignment has a basic set of permissions available; file permissions are simply read, write, and execute, and directories have an additional permission to determine access to child directories. (LINKTO)
- [Extended Access Control Lists](#) (ACLs) for HDFS that provide fine-grained control of permissions for HDFS files by allowing you to set different permissions for specific named users or named groups.
- Apache HBase uses ACLs to authorize various operations (`READ`, `WRITE`, `CREATE`, `ADMIN`) by column, column family, and column family qualifier. HBase ACLs are granted and revoked to both users and groups.
- Role-based access control with Apache Sentry. As of Cloudera Manager 5.1.x, Sentry permissions can be configured using either policy files or the database-backed Sentry service.
 - The Sentry service is the preferred way to set up Sentry permissions. See [The Sentry Service](#) on page 359 for more information.
 - For the policy file approach to configuring Sentry, see [Sentry Policy File Authorization](#) on page 392.



Important: Cloudera does not support Apache Ranger or Hive's native authorization frameworks for configuring access control in Hive. Use the Cloudera-supported Apache Sentry instead.

POSIX Permissions

The majority of services within the Hadoop ecosystem, from client applications like the CLI shell to tools written to use the Hadoop API, directly access data stored within HDFS. HDFS uses POSIX-style permissions for directories and files; each directory and file is assigned a single owner and group. Each assignment has a basic set of permissions available; file permissions are simply read, write, and execute, and directories have an additional permission to determine access to child directories.

Ownership and group membership for a given HDFS asset determines a user's privileges. If a given user fails either of these criteria, they are denied access. For services that may attempt to access more than one file, such as MapReduce, Cloudera Search, and others, data access is determined separately for each file access attempt. File permissions in HDFS are managed by the NameNode.

Access Control Lists

Hadoop also maintains general access controls for the services themselves in addition to the data within each service and in HDFS. Service access control lists (ACL) are typically defined within the global `hadoop-policy.xml` file and range from NameNode access to client-to-DataNode communication. In the context of MapReduce and YARN, user and group identifiers form the basis for determining permission for job submission or modification.

In addition, with MapReduce and YARN, jobs can be submitted using queues controlled by a scheduler, which is one of the components comprising the resource management capabilities within the cluster. Administrators define permissions to individual queues using ACLs. ACLs can also be defined on a job-by-job basis. Like HDFS permissions, local user accounts and groups must exist on each executing server, otherwise the queues will be unusable except by superuser accounts.

Apache HBase also uses ACLs for data-level authorization. HBase ACLs authorize various operations (READ, WRITE, CREATE, ADMIN) by column, column family, and column family qualifier. HBase ACLs are granted and revoked to both users and groups. Local user accounts are required for proper authorization, similar to HDFS permissions.

Apache ZooKeeper also maintains ACLs to the information stored within the DataNodes of a ZooKeeper data tree.

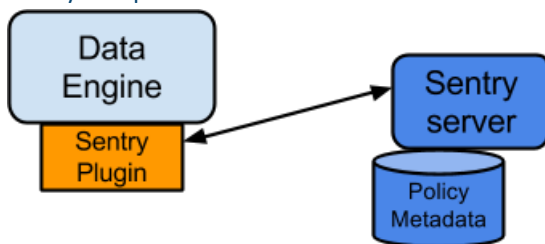
Role-Based Access Control with Apache Sentry (incubating)

For finer-grained access to data accessible using schema -- that is, data structures described by the Apache Hive Metastore and used by computing engines like Hive and Impala, as well as collections and indices within Cloudera Search -- CDH supports Apache Sentry (incubating), which offers a role-based privilege model for this data and its given schema. Apache Sentry is a granular, role-based authorization module for Hadoop. Sentry provides the ability to control and enforce precise levels of privileges on data for authenticated users and applications on a Hadoop cluster. Sentry currently works out of the box with Apache Hive, Hive Metastore/HCatalog, Apache Solr, Cloudera Impala and HDFS (limited to Hive table data).

Sentry is designed to be a pluggable authorization engine for Hadoop components. It allows you to define authorization rules to validate a user or application's access requests for Hadoop resources. Sentry is highly modular and can support authorization for a wide variety of data models in Hadoop.

Architecture Overview

Sentry Components



There are three components involved in the authorization process:

- **Sentry Server**

The Sentry RPC server manages the authorization metadata. It supports interfaces to securely retrieve and manipulate the metadata.

- **Data Engine**

This is a data processing application such as Hive or Impala that needs to authorize access to data or metadata resources. The data engine loads the Sentry plugin and all client requests for accessing resources are intercepted and routed to the Sentry plugin for validation.

- **Sentry Plugin**

The Sentry plugin runs in the data engine. It offers interfaces to manipulate authorization metadata stored in the Sentry server, and includes the authorization policy engine that evaluates access requests using the authorization metadata retrieved from the server.

Key Concepts

- Authentication - Verifying credentials to reliably identify a user
- Authorization - Limiting the user's access to a given resource
- User - Individual identified by underlying authentication system
- Group - A set of users, maintained by the authentication system
- Privilege - An instruction or rule that allows access to an object
- Role - A set of privileges; a template to combine multiple access rules
- Authorization models - Defines the objects to be subject to authorization rules and the granularity of actions allowed. For example, in the SQL model, the objects can be databases or tables, and the actions are `SELECT`, `INSERT`, `CREATE` and so on. For the Search model, the objects are indexes, collections and documents; the access modes are query, update and so on.

User Identity and Group Mapping

Sentry relies on underlying authentication systems such as Kerberos or LDAP to identify the user. It also uses the group mapping mechanism configured in Hadoop to ensure that Sentry sees the same group mapping as other components of the Hadoop ecosystem.

Consider users Alice and Bob who belong to an Active Directory (AD) group called `finance-department`. Bob also belongs to a group called `finance-managers`. In Sentry, you first create roles and then grant privileges to these roles. For example, you can create a role called Analyst and grant `SELECT` on tables Customer and Sales to this role.

The next step is to join these authentication entities (users and groups) to authorization entities (roles). This can be done by granting the Analyst role to the `finance-department` group. Now Bob and Alice who are members of the `finance-department` group get `SELECT` privilege to the Customer and Sales tables.

Role-based access control (RBAC) is a powerful mechanism to manage authorization for a large set of users and data objects in a typical enterprise. New data objects get added or removed, users join, move, or leave organisations all the time. RBAC makes managing this a lot easier. Hence, as an extension of the discussed previously, if Carol joins the Finance Department, all you need to do is add her to the `finance-department` group in AD. This will give Carol access to data from the Sales and Customer tables.

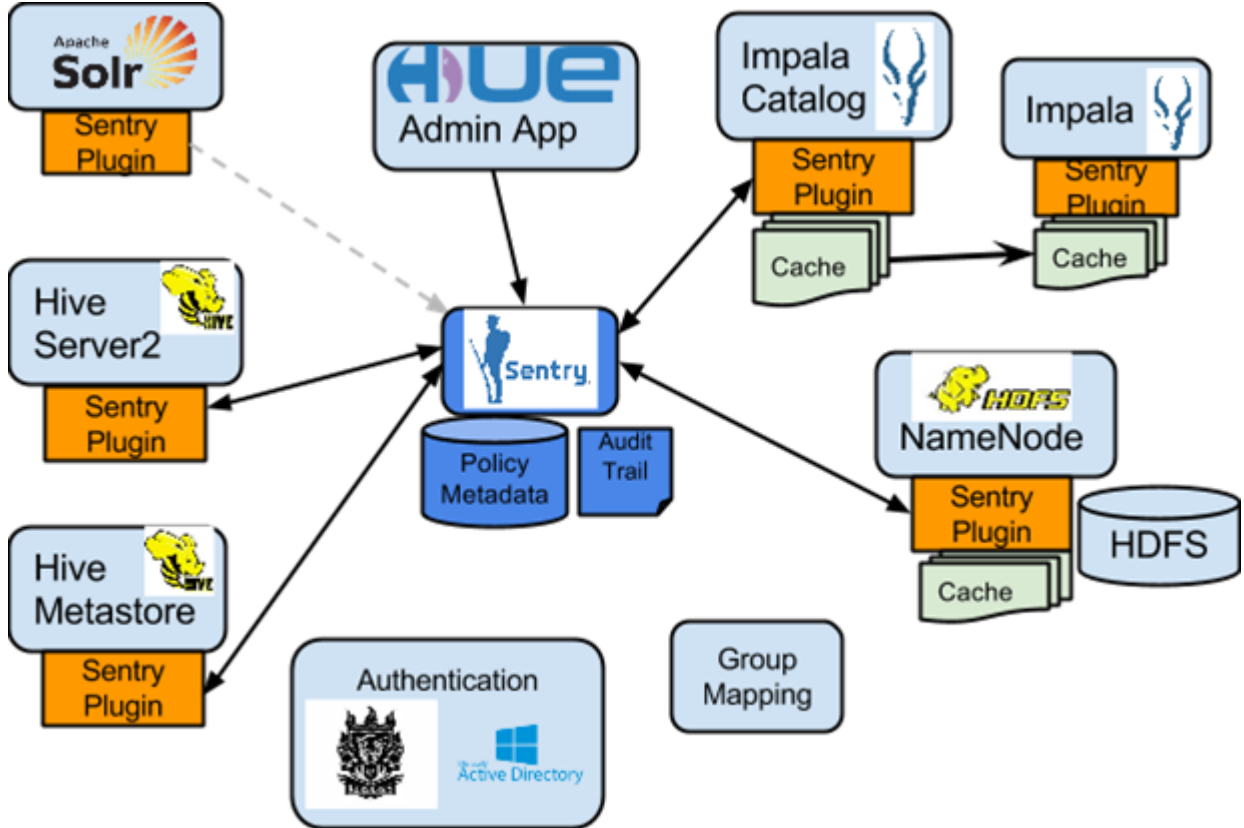
Unified Authorization

Another important aspect of Sentry is the unified authorization. The access control rules once defined, work across multiple data access tools. For example, being granted the Analyst role in the previous example will allow Bob, Alice, and others in the `finance-department` group to access table data from SQL engines such as Hive and Impala, as well as using MapReduce, Pig applications or metadata access using HCatalog.

As illustrated above, Apache Sentry works with multiple Hadoop components. At the heart you have the Sentry Server which stores authorization metadata and provides APIs for tools to retrieve and modify this metadata securely.

Note that the Sentry server only facilitates the metadata. The actual authorization decision is made by a policy engine which runs in data processing applications such as Hive or Impala. Each component loads the Sentry plugin which includes the service client for dealing with the Sentry service and the policy engine to validate the authorization request.

Sentry Integration with the Hadoop Ecosystem



As illustrated above, Apache Sentry works with multiple Hadoop components. At the heart you have the Sentry Server which stores authorization metadata and provides APIs for tools to retrieve and modify this metadata securely.

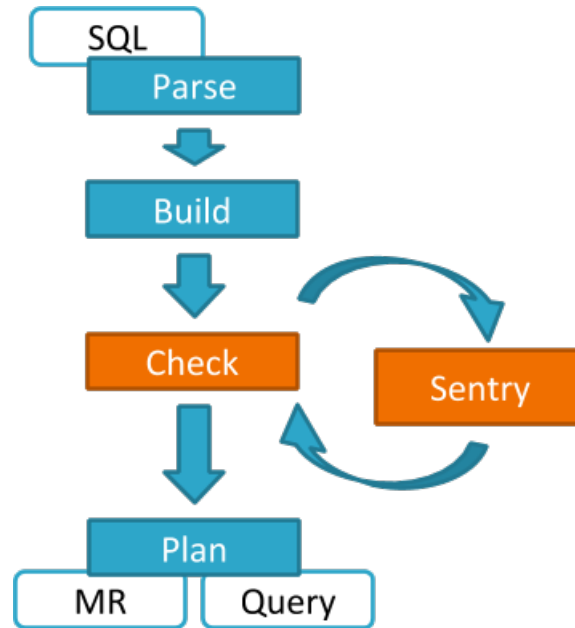
Note that the Sentry server only facilitates the metadata. The actual authorization decision is made by a policy engine which runs in data processing applications such as Hive or Impala. Each component loads the Sentry plugin which includes the service client for dealing with the Sentry service and the policy engine to validate the authorization request.

Hive and Sentry

Consider an example where Hive gets a request to access an object in a certain mode by a client. If Bob submits the following Hive query:

```
select * from production.sales
```

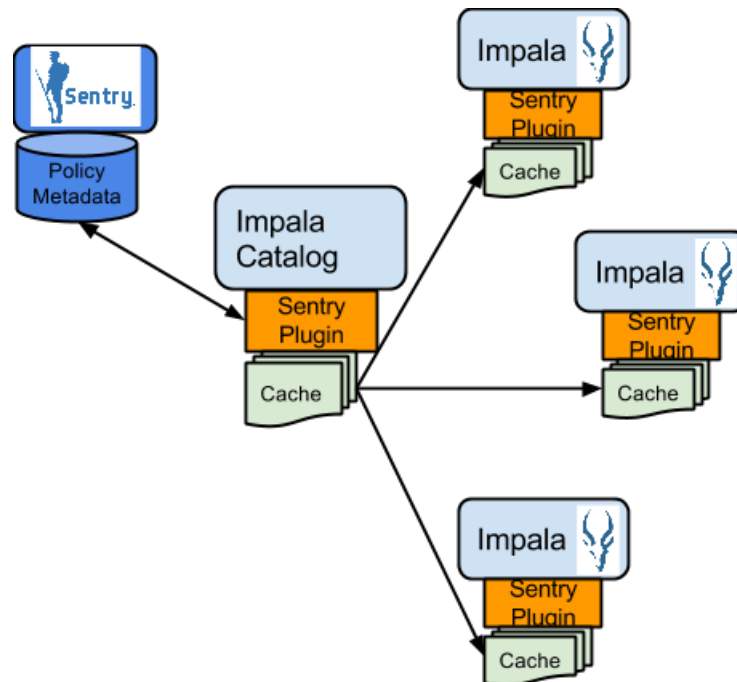
Hive will identify that user Bob is requesting `SELECT` access to the Sales table. At this point Hive will ask the Sentry plugin to validate Bob's access request. The plugin will retrieve Bob's privileges related to the Sales table and the policy engine will determine if the request is valid.



Hive works with both, the Sentry service and policy files. Cloudera recommends you use the Sentry service which makes it easier to manage user privileges. For more details and instructions, see [The Sentry Service](#) on page 359 or [Sentry Policy File Authorization](#) on page 392.

Impala and Sentry

Authorization processing in Impala is similar to that in Hive. The main difference is caching of privileges. Impala's Catalog server manages caching schema metadata and propagating it to all Impala server nodes. This Catalog server caches Sentry metadata as well. As a result, authorization validation in Impala happens locally and much faster.

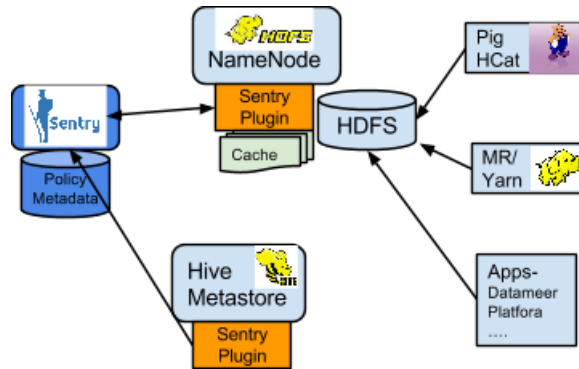


For detailed documentation, see [Enabling Sentry Authorization for Impala](#) on page 411.

Sentry-HDFS Synchronization

Sentry-HDFS authorization is focused on Hive warehouse data - that is, any data that is part of a table in Hive or Impala. The real objective of this integration is to expand the same authorization checks to Hive warehouse data being accessed

from any other components such as Pig, MapReduce or Spark. At this point, this feature does not replace HDFS ACLs. Tables that are not associated with Sentry will retain their old ACLs.



The mapping of Sentry privileges to HDFS ACL permissions is as follows:

- SELECT privilege -> Read access on the file.
- INSERT privilege -> Write access on the file.
- ALL privilege -> Read and Write access on the file.

The NameNode loads a Sentry plugin that caches Sentry privileges as well as Hive metadata. This helps HDFS to keep file permissions and Hive table privileges in sync. The Sentry plugin periodically polls the Sentry and Metastore to keep the metadata changes in sync.

For example, if Bob runs a Pig job that is reading from the Sales table data files, Pig will try to get the file handle from HDFS. At that point the Sentry plugin on the NameNode will figure out that the file is part of Hive data and overlay Sentry privileges on top of the file ACLs. As a result, HDFS will enforce the same privileges for this Pig client that Hive would apply for a SQL query.

For HDFS-Sentry synchronization to work, you *must* use the Sentry service, not policy file authorization. See [Synchronizing HDFS ACLs and Sentry Permissions](#) on page 385, for more details.

Search and Sentry

Sentry can apply a range of restrictions to various Search tasks, such as accessing data or creating collections. These restrictions are consistently applied, regardless of the way users attempt to complete actions. For example, restricting access to data in a collection restricts that access whether queries come from the command line, from a browser, or through the admin console.

With Search, Sentry stores its privilege policies in a policy file (for example, `sentry-provider.ini`) which is stored in an HDFS location such as `hdfs://ha-nn-uri/user/solr/sentry/sentry-provider.ini`.

Sentry with Search does not support multiple policy files for multiple databases. However, you must use a separate policy file for each Sentry-enabled service. For example, Hive and Search were using policy file authorization, using a combined Hive and Search policy file would result in an invalid configuration and failed authorization on both services.



Note: While Hive and Impala are compatible with the database-backed Sentry service, Search still uses Sentry's policy file authorization. Note that it is possible for a single cluster to use both, the Sentry service (for Hive and Impala as described above) and Sentry policy files (for Solr).

For detailed documentation, see [Enabling Sentry Authorization for Search using the Command Line](#) on page 422.

Authorization Administration

The Sentry server supports APIs to securely manipulate roles and privileges. Both Hive and Impala support SQL statements to manage privileges natively. Sentry assumes that HiveServer2 and Impala run as superusers, usually called `hive` and

impala. To initiate top-level permissions for Sentry, an admin must login as a superuser. You can use either Beeline or the Impala shell to execute the following sample statement:

```
GRANT ROLE Analyst TO GROUP finance-managers
```

Using Hue to Manage Sentry Permissions

Hue supports a Security app to manage Sentry authorization. This allows users to explore and change table permissions. Here is a [video blog](#) that demonstrates its functionality.

Integration with Authentication Mechanisms for Identity Management

Like many distributed systems, Hadoop projects and workloads often consist of a collection of processes working in concert. In some instances, the initial user process conducts authorization throughout the entirety of the workload or job's lifecycle. But for processes that spawn additional processes, authorization can pose challenges. In this case, the spawned processes are set to execute as if they were the authenticated user, that is, `setuid`, and thus only have the privileges of that user. The overarching system requires a mapping to the authenticated principal and the user account must exist on the local host system for the `setuid` to succeed.



Important:

- Cloudera strongly recommends *against* using Hadoop's `LdapGroupsMapping` provider. `LdapGroupsMapping` should only be used in cases where OS-level integration is not possible. Production clusters require an identity provider that works well with all applications, not just Hadoop. Hence, often the preferred mechanism is to use tools such as SSSD, VAS or Centrify to replicate LDAP groups.
- Cloudera does not support the use of Winbind in production environments. Winbind uses an inefficient approach to user/group mapping, which may lead to low performance or cluster failures as the size of the cluster, and the number of users and groups increases.

Irrespective of the mechanism used, user/group mappings must be applied consistently across all cluster hosts for ease with maintenance.

System and Service Authorization - Several Hadoop services are limited to inter-service interactions and are not intended for end-user access. These services do support authentication to protect against unauthorized or malicious users. However, any user or, more typically, another service that has login credentials and can authenticate to the service is authorized to perform all actions allowed by the target service. Examples include ZooKeeper, which is used by internal systems such as YARN, Cloudera Search, and HBase, and Flume, which is configured directly by Hadoop administrators and thus offers no user controls.

The authenticated Kerberos principals for these "system" services are checked each time they access other services such as HDFS, HBase, and MapReduce, and therefore must be authorized to use those resources. Thus, the fact that Flume does not have an explicit authorization model does not imply that Flume has unrestricted access to HDFS and other services; the Flume service principals still must be authorized for specific locations of the HDFS file system. Hadoop administrators can establish separate system users for a services such as Flume to segment and impose access rights to only the parts of the file system for a specific Flume application.

Authorization within Hadoop Projects

Project	Authorization Capabilities
HDFS	File Permissions, Sentry*
MapReduce	File Permissions, Sentry*
YARN	File Permissions, Sentry*
Accumulo	

Project	Authorization Capabilities
Flume	None
HBase	HBase ACLs
HiveServer2	File Permissions, Sentry
Hue	Hue authorization mechanisms (assigning permissions to Hue apps)
Impala	Sentry
Oozie	ACLs
Pig	File Permissions, Sentry*
Search	File Permissions, Sentry
Sentry	N/A
Spark	File Permissions, Sentry*
Sqoop	N/A
Sqoop2	None
ZooKeeper	ACLs
Cloudera Manager	Cloudera Manager roles
Cloudera Navigator	Cloudera Navigator roles
Backup and Disaster Recovery	N/A

* Sentry HDFS plug-in; when enabled, Sentry enforces its own access permissions over files that are part of tables defined in the Hive Metastore.

Overview of Data Management Mechanisms for an Enterprise Data Hub

For the data in the cluster, it is critical to understand where the data is coming from and how it's being used. The goal of auditing is to capture a complete and immutable record of all activity within a system. Auditing plays a central role in three key activities within the enterprise:

- First, auditing is part of a system’s security regime and can explain what happened, when, and to whom or what in case of a breach or other malicious intent. For example, if a rogue administrator deletes a user’s data set, auditing provides the details of this action, and the correct data may be retrieved from backup.
- The second activity is compliance, and auditing participates in satisfying the core requirements of regulations associated with sensitive or personally identifiable data (PII), such as the Health Insurance Portability and Accountability Act (HIPAA) or the Payment Card Industry (PCI) Data Security Standard. Auditing provides the touchpoints necessary to construct the trail of who, how, when, and how often data is produced, viewed, and manipulated.
- Lastly, auditing provides the historical data and context for data forensics. Audit information leads to the understanding of how various populations use different data sets and can help establish the access patterns of these data sets. This examination, such as trend analysis, is broader in scope than compliance and can assist content and system owners in their data optimization efforts.

The risks facing auditing are the reliable, timely, and tamper-proof capture of all activity, including administrative actions. Until recently, the native Hadoop ecosystem has relied primarily on using log files. Log files are unacceptable for most audit use cases in the enterprise as real-time monitoring is impossible, and log mechanics can be unreliable - a system crash before or during a write commit can compromise integrity and lead to data loss.

Cloudera Navigator is a fully integrated data management and security tool for the Hadoop platform. Data management and security capabilities are critical for enterprise customers that are in highly regulated industries and have stringent compliance requirements. This topic only provides an overview of some of the auditing and metadata management capabilities that Cloudera Navigator offers. For complete details, see [Cloudera Data Management](#).

Cloudera Navigator

The following sections describe some of the categories of functionalities Cloudera Navigator provides for auditing, metadata management and lineage.

Auditing

While Hadoop has historically lacked centralized cross-component audit capabilities, products such as Cloudera Navigator add secured, real-time audit components to key data and access frameworks. Cloudera Navigator allows administrators to configure, collect, and view audit events, to understand who accessed what data and how. Cloudera Navigator also allows administrators to generate reports that list the HDFS access permissions granted to groups. Cloudera Navigator tracks access permissions and actual accesses to all entities in HDFS, Hive, HBase, Impala, Sentry, and Solr, and the Cloudera Navigator Metadata Server itself to help answer questions such as - who has access to which entities, which entities were accessed by a user, when was an entity accessed and by whom, what entities were accessed using a service, which device was used to access, and so on. Cloudera Navigator auditing supports tracking access to:

- HDFS entities accessed by HDFS, Hive, HBase, Impala, and Solr services
- HBase and Impala
- Hive metadata
- Sentry
- Solr
- Cloudera Navigator Metadata Server

Data collected from these services also provides visibility into usage patterns for users, ability to see point-in-time permissions and how they have changed (leveraging Sentry), and review and verify HDFS permissions. Cloudera Navigator also provides out-of-the-box integration with leading enterprise metadata, lineage, and SIEM applications. For details on how Cloudera Navigator handles auditing, see [Cloudera Navigator Auditing Architecture](#).

The screenshot displays the Cloudera Navigator interface for viewing audit events. The top navigation bar includes 'Search', 'Audits', 'Analytics', 'Policies', 'Administration', and 'admin'. The main content area is titled 'Audit Events' and shows a table of recent audit events. The table has the following columns: Timestamp, Username, IP Address, Service Name, Operation, and Resource. The events listed include actions like 'savedSearch', 'authentication', 'getFileinfo', and 'listStatus' performed by users such as 'admin', 'accumulo', and 'oozie' on services like 'Navigator' and 'HDFS-1'. The interface also includes a search bar, a 'Save As Report' button, and a 'Filters' dropdown.

Metadata Management

For metadata and data discovery, Cloudera Navigator features complete metadata storage. First, it consolidates the technical metadata for all data inside Hadoop into a single, searchable interface and allows for automatic tagging of

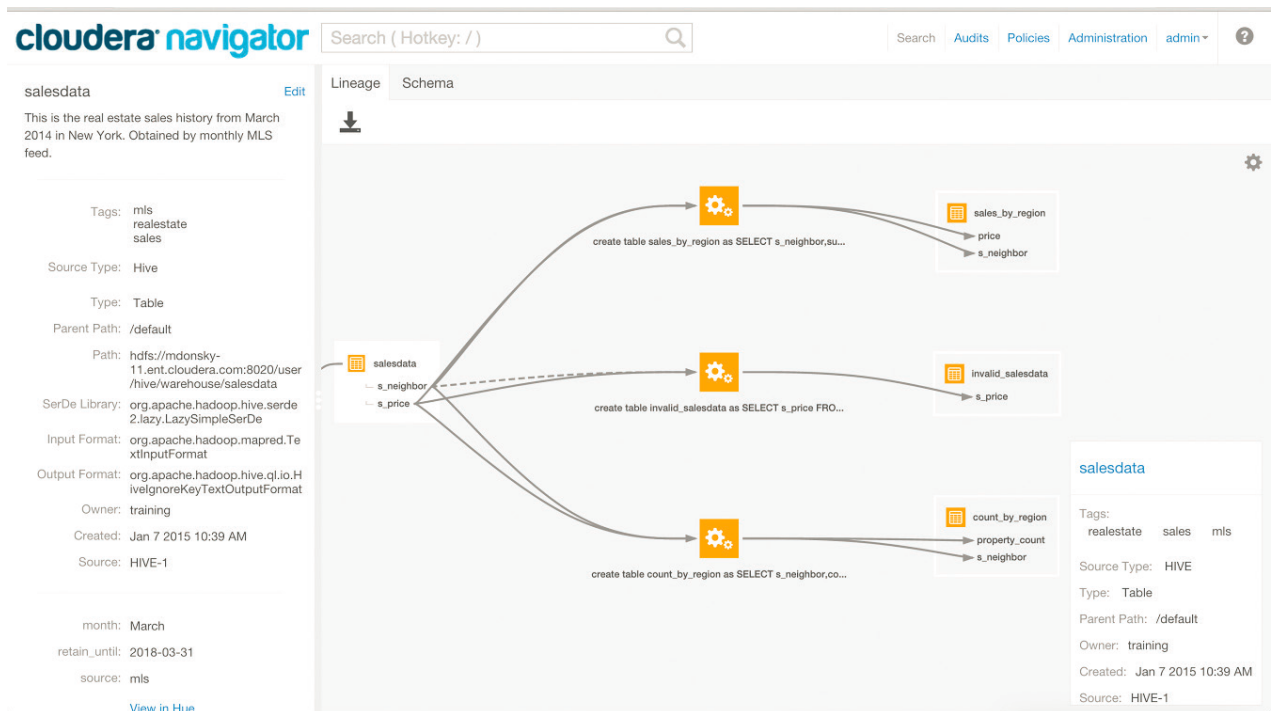
data based on the external sources entering the cluster. For example, if there is an external ETL process, data can be automatically tagged as such when it enters Hadoop. Second, it supports user-based tagging to augment files, tables, and individual columns with custom business context, tags, and key/value pairs. Combined, this allows data to be easily discovered, classified, and located to not only support governance and compliance, but also user discovery within Hadoop.

Cloudera Navigator also includes metadata policy management that can trigger actions (such as the autoclassification of metadata) for specific datasets based on arrival or scheduled intervals. This allows users to easily set, monitor, and enforce data management policies, while also integrating with common third-party tools.

For details on how Cloudera Navigator handles metadata, see [Cloudera Navigator Metadata Architecture](#).

Lineage

Cloudera Navigator provides an automatic collection and easy visualization of upstream and downstream data lineage to verify reliability. For each data source, it shows, down to the column-level within that data source, what the precise upstream data sources were, the transforms performed to produce it, and the impact that data has on downstream artifacts. Cloudera Navigator supports tracking the lineage of HDFS files, datasets, and directories, Hive tables and columns, MapReduce and YARN jobs, Hive queries, Impala queries, Pig scripts, Oozie workflows, Spark jobs, and Sqoop jobs. For details, see [Cloudera Navigator Lineage Diagrams](#).



Integration within an EDH

The monitoring and reporting of Hadoop systems, while critical elements to its enterprise usage, are only a part of an enterprise's complete audit infrastructure and data policy. Often these enterprise tools and policies require that all audit information route through a central interface to aid comprehensive reporting, and Hadoop-specific audit data can be integrated with these existing enterprise SIEM applications and other tools. For example, Cloudera Navigator exposes Hadoop audit data through several delivery methods:

- Using syslog, thus acting as a mediator between the raw event streams in Hadoop and the SIEM tools.
- Using a REST API for custom enterprise tools.
- You can also simply export the data to a file, such as a comma-delimited text file.

Auditing in Hadoop Projects

The table below depicts the auditing capabilities of Cloudera Manager and CDH components.

Project	Auditing Capabilities
HDFS	Events captured by Cloudera Navigator (including security events*)
MapReduce	Inferred through HDFS
YARN	Inferred through HDFS
Accumulo	Log Files - Partial inclusion of security events; does not include non-bulk writes
Flume	Log Files
HBase	Audit events captured by Cloudera Navigator (including security events*)
HiveServer2	Audit events captured by Cloudera Navigator
Hue	Inferred through underlying components
Impala	Audit events captured by Cloudera Navigator
Oozie	Log Files
Pig	Inferred through HDFS
Search	Log Files
Sentry	Audit events captured by Cloudera Navigator
Spark	Inferred through HDFS
Sqoop	Log Files
Sqoop2	Log Files (including security events*)
ZooKeeper	Log Files
Cloudera Manager	Audit events captured by Cloudera Navigator (partial capture of security events*)
Cloudera Navigator	Audit events captured by Cloudera Navigator itself
Backup and Disaster Recovery	None

* - Security events include a machine readable log of the following activities:

- User data read
- User data written
- Permission changes
- Configuration changes
- Login attempts
- Escalation of privileges
- Session Tracking
- Key Operations (Key Trustee)

How to Configure TLS Encryption for Cloudera Manager

When you configure authentication and authorization on a cluster, Cloudera Manager Server sends sensitive information over the network to cluster hosts, such as Kerberos keytabs and configuration files that contain passwords. To secure this transfer, you must configure TLS encryption between Cloudera Manager Server and all cluster hosts.

TLS encryption is also used to secure client connections to the Cloudera Manager Admin Interface, using HTTPS.

Cloudera Manager also supports TLS authentication. Without certificate authentication, a malicious user can add a host to Cloudera Manager by installing the Cloudera Manager Agent software and configuring it to communicate with Cloudera Manager Server. To prevent this, you must install certificates on each agent host and configure Cloudera Manager Server to trust those certificates.

This guide shows how to configure and enable TLS encryption and certificate authentication for Cloudera Manager. The provided examples use an internal certificate authority (CA) to sign all TLS certificates, so this guide also shows you how to establish trust with the CA. (For certificates signed by a trusted public CA, establishing trust is not necessary, because the Java Development Kit (JDK) already trusts them.)

Use this guide to enable TLS encryption and certificate authentication for Cloudera Manager:

Generate TLS Certificates



Important:

- You must use the Oracle JDK `keytool` utility. Do not use other JDK (such as OpenJDK) command line tools for this procedure. If you have multiple JDKs, set the `PATH` variable such that the Oracle JDK is first. For example:

```
$ export JAVA_HOME=/usr/java/jdk1.7.0_67-cloudera
$ export PATH=$JAVA_HOME/bin:$PATH
```

- Use the same password for the `-keypass` and `-storepass` values. Cloudera Manager does not support using different passwords for the key and keystore.

Before configuring Cloudera Manager Server and all Cloudera Manager Agents to use TLS encryption, generate the server and agent certificates:

Generate the Cloudera Manager Server Certificate

The following procedure assumes that an internal certificate authority (CA) is used, and shows how to establish trust for that internal CA. If you are using a trusted public CA (such as Symantec, GeoTrust, Comodo, and others), you do not need to explicitly establish trust for the issued certificates, unless you are using an older JDK and a newer public CA. Older JDKs might not trust newer public CAs by default.

1. On the Cloudera Manager Server host, create the `/opt/cloudera/security/pki` directory:

```
$ sudo mkdir -p /opt/cloudera/security/pki
```

If you choose to use a different directory, make sure you use the same directory on all cluster hosts to simplify management and maintenance.

2. Use the `keytool` utility to generate a Java keystore and certificate signing request (CSR). Replace the `OU`, `O`, `L`, `ST`, and `C` entries with the values for your environment. When prompted, use the same password for the `keystore`

password and key password. Cloudera Manager does not support using different passwords for the key and keystore.

```
$JAVA_HOME/bin/keytool -genkeypair -alias $(hostname -f) -keyalg RSA -keystore
/opt/cloudera/security/pki/$(hostname -f).jks -keysize 2048 -dname "CN=$(hostname
-f),OU=Engineering,O=Cloudera,L=Palo Alto,ST=California,C=US" -ext san=dns:$(hostname
-f)
```

```
$JAVA_HOME/bin/keytool -certreq -alias $(hostname -f) -keystore
/opt/cloudera/security/pki/$(hostname -f).jks -file /opt/cloudera/security/pki/$(hostname
-f).csr -ext san=dns:$(hostname -f)
```

3. Submit the CSR file (for example, `cm01.example.com-server.csr`) to your certificate authority to obtain a server certificate. If possible, obtain the certificate in PEM (Base64 ASCII) format. The certificate file is in PEM format if it looks like this:

```
-----BEGIN CERTIFICATE-----
MIIDAzCCAesCAQAwgY0xCzAJBgNVBAYTAlVTMRMwEQYDVQQLIEwpcDYWxpZm9ybmlh
MRIWEAYDVQQHEwlQYWxvIEFsdG8xETAPBgNVBAoTCENsb3VkZXJhMRQwEgYDVQQL
...
tudY0C32LjGjW0g5ALlin9Oylu2xRKGAVfapbzAZ2rcht1Czc7mtaT6BXgW8S+Db
0HhuObn1/8TL4Ho9G+KlJB3MWik2oEbOvQt0rBidMr9qaNX86m0i7pouXZelZ5c5
UnDPtrhW6A==
-----END CERTIFICATE-----
```

If your issued certificate is in binary (DER) format, convert it to PEM format.

4. After you receive the signed certificate, copy it to `/opt/cloudera/security/pki/$(hostname -f)-server.cert.pem`.
5. Copy the root and intermediate CA certificates to `/opt/cloudera/security/pki/rootca.cert.pem` and `/opt/cloudera/security/pki/intca.cert.pem` on the Cloudera Manager Server host. If you have a concatenated file containing the root CA and an intermediate CA certificate, split the file along the `END CERTIFICATE/BEGIN CERTIFICATE` boundary into individual files. If there are multiple intermediate CA certificates, use unique file names such as `intca-1.cert.pem`, `intca-2.cert.pem`, and so on.
6. On the Cloudera Manager Server host, copy the JDK `cacerts` file to `jssecacerts`:

```
$ sudo cp $JAVA_HOME/jre/lib/security/cacerts $JAVA_HOME/jre/lib/security/jssecacerts
```

If you do not have the `$JAVA_HOME` variable set, replace it with the path to the Oracle JDK (for example, `/usr/java/jdk1.7.0_67-cloudera/`).



Note: The Oracle JDK uses the `jssecacerts` file for its default truststore if it exists. Otherwise, it uses the `cacerts` file. Creating the `jssecacerts` file allows you to trust an internal CA without modifying the `cacerts` file that is included with the JDK.

7. Import the root CA certificate into the JDK truststore. If you do not have the `$JAVA_HOME` variable set, replace it with the path to the Oracle JDK.

```
$ sudo keytool -importcert -alias rootca -keystore $JAVA_HOME/jre/lib/security/jssecacerts
-file /opt/cloudera/security/pki/rootca.cert.pem -storepass changeit
```

The default password for the `cacerts` file is `changeit`. Cloudera recommends changing this password by running:

```
keytool -storepasswd -keystore $JAVA_HOME/jre/lib/security/cacerts
```

8. Copy the `jssecacerts` file from the Cloudera Manager Server host to all other cluster hosts. Make sure you copy the file to the correct location (`$JAVA_HOME/jre/lib/security/jssecacerts`), because the Oracle JDK expects it there.

How to Configure TLS Encryption for Cloudera Manager

9. On the Cloudera Manager Server host, append the intermediate CA certificate to the signed server certificate, and then import it into the keystore. Make sure that you use the append operator (>>) and not the overwrite operator (>):

```
$ sudo cat /opt/cloudera/security/pki/intca.cert.pem >>
/opt/cloudera/security/pki/${hostname -f}-server.cert.pem
$ sudo keytool -importcert -alias ${hostname -f}-server \
-file /opt/cloudera/security/pki/${hostname -f}-server.cert.pem \
-keystore /opt/cloudera/security/pki/${hostname -f}-server.jks
```

If you see a message like the following, enter `yes` to continue:

```
... is not trusted. Install reply anyway? [no]: yes
```

You *must* see the following response verifying that the certificate has been properly imported against its private key.

```
Certificate reply was installed in keystore
```

If you do not see this response, [contact Cloudera Support](#).

Generate the Cloudera Manager Agent Certificates

Complete the following procedure on each Cloudera Manager Agent host. The provided examples continue to use an internal certificate authority (CA) to sign the agent certificates.

1. On all Cloudera Manager Agent hosts, create the `/opt/cloudera/security/pki` directory:

```
$ sudo mkdir -p /opt/cloudera/security/pki
```

If you choose to use a different directory, make sure you use the same directory on all cluster hosts to simplify management and maintenance.

2. On all Cloudera Manager Agent hosts, create a Java Keystore and private key as follows:

```
$ keytool -genkeypair -alias ${hostname -f}-agent -keyalg RSA -keystore \
/opt/cloudera/security/pki/${hostname -f}-keystore.jks -keysize 2048 -dname \
"CN=${hostname -f},OU=Engineering,O=Cloudera,L=Palo Alto,ST=California,C=US" \
-storepass password -keypass password
```

Use the same password for the `-keypass` and `-storepass` values. Cloudera Manager does not support using different passwords for the key and keystore.

3. On all Cloudera Manager Agent hosts, generate the certificate signing request (CSR) and submit it to a CA. Use the `keytool` extended attributes to specify both `serverAuth` and `clientAuth` options:

```
$ keytool -certreq -alias ${hostname -f}-agent \
-keystore /opt/cloudera/security/pki/${hostname -f}-agent.jks \
-file /opt/cloudera/security/pki/${hostname -f}-agent.csr \
-ext EKU=serverAuth,clientAuth \
-storepass password -keypass password
```

For security purposes, many commercial CAs ignore requested extensions in a CSR. Make sure that you inform the CA that you require certificates with both server and client authentication options.

4. For each signed certificate you receive, copy it to `/opt/cloudera/security/pki/${hostname -f}-agent.cert.pem` on the correct host.
5. Inspect the certificates to verify that both server and client authentication options are present:

```
$ openssl x509 -in /opt/cloudera/security/pki/${hostname -f}-agent.cert.pem -noout -text
```

Look for output similar to the following:

```
X509v3 Extended Key Usage:
    TLS Web Server Authentication, TLS Web Client Authentication
X509v3 Subject Alternative Name:
DNS:your.hosts.name.com
If the certificate does not have the DNS field, re-submit the CSR to the CA, and request
that they generate a certificate that keeps the Subject Alternative Name field intact.
```

If the certificate does not have both TLS Web Server Authentication and TLS Web Client Authentication listed in the X509v3 Extended Key Usage section, re-submit the CSR to the CA, and request that they generate a certificate that can be used for both server and client authentication.

6. Copy the root and intermediate CA certificates to `/opt/cloudera/security/pki/rootca.cert.pem` and `/opt/cloudera/security/pki/intca.cert.pem` on each Cloudera Manager Agent host. If you have a concatenated file containing the root CA and an intermediate CA certificate, split the file along the `END CERTIFICATE/BEGIN CERTIFICATE` boundary into individual files. If there are multiple intermediate CA certificates, use unique file names such as `intca-1.cert.pem`, `intca-2.cert.pem`, and so on.
7. On each Cloudera Manager Agent host, append the intermediate CA certificate to the signed certificate, and then import it into the keystore. Make sure that you use the append operator (`>>`) and not the overwrite operator (`>`):

```
$ sudo cat /opt/cloudera/security/pki/intca.cert.pem >>
/opt/cloudera/security/pki/$(hostname -f)-agent.cert.pem
$ sudo keytool -importcert -alias $(hostname -f)-agent \
-file /opt/cloudera/security/pki/$(hostname -f)-agent.cert.pem \
-keystore /opt/cloudera/security/pki/$(hostname -f)-agent.jks
```

If you see a message like the following, enter `yes` to continue:

```
... is not trusted. Install reply anyway? [no]: yes
```

You *must* see the following response verifying that the certificate has been properly imported against its private key.

```
Certificate reply was installed in keystore
```

If you do not see this response, [contact Cloudera Support](#).

8. On each Cloudera Manager Agent host, create symbolic links (symlink) for the certificate and keystore files:

```
$ ln -s /opt/cloudera/security/pki/$(hostname -f)-agent.cert.pem
/opt/cloudera/security/pki/agent.cert.pem
$ ln -s /opt/cloudera/security/pki/$(hostname -f)-agent.jks
/opt/cloudera/security/pki/agent.jks
```

This allows you to use the same `/etc/cloudera-scm-agent/config.ini` file on all agent hosts rather than maintaining a file for each agent.

Configuring TLS Encryption for the Cloudera Manager Admin Console

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

Use the following procedure to enable TLS encryption for the Cloudera Manager Server admin interface. Make sure you have generated the server certificate as described in [Generate the Cloudera Manager Server Certificate](#) on page 40.

Step 1: Enable HTTPS for the Cloudera Manager Admin Console

1. Log in to the Cloudera Manager Admin Console.
2. Select **Administration > Settings**.

How to Configure TLS Encryption for Cloudera Manager

3. Select the **Security** category.
4. Configure the following TLS settings:

Property	Description
Cloudera Manager TLS/SSL Server JKS Keystore File Location	The complete path to the keystore file. In this example, the path is <code>/opt/cloudera/security/pki/cm01.example.com-server.jks</code> . Replace <code>cm01.example.com</code> with the Cloudera Manager Server hostname.
Cloudera Manager TLS/SSL Server JKS Keystore File Password	The password for the <code>/opt/cloudera/security/jks/cm01.example.com-server.jks</code> keystore.
Use TLS Encryption for Admin Console	Check this box to enable TLS encryption for Cloudera Manager.

5. Click **Save Changes** to save the settings.

Step 2: Specify SSL Truststore Properties for Cloudera Management Services

When enabling TLS for the Cloudera Manager Server admin interface, you must set the Java truststore location and password in the Cloudera Management Services configuration. Otherwise, roles such as Host Monitor and Service Monitor cannot connect to Cloudera Manager Server and will not start.

Configure the path and password for the `$JAVA_HOME/jre/lib/security/jssecacerts` truststore that you created earlier. Make sure that you copied this file to all cluster hosts, including the Cloudera Management Service hosts.

1. Open the Cloudera Manager Administration Console and go to the **Cloudera Management Service** service.
2. Click the **Configuration** tab.
3. Select **Scope > Cloudera Management Service (Service-Wide)**.
4. Select **Category > Security**.
5. Edit the following TLS/SSL properties according to your cluster configuration.

Property	Description
TLS/SSL Client Truststore File Location	The path to the client truststore file used in HTTPS communication. This truststore contains certificates of trusted servers, or of Certificate Authorities trusted to identify servers. For this example, set the value to <code>\$JAVA_HOME/jre/lib/security/jssecacerts</code> . Replace <code>\$JAVA_HOME</code> with the path to the Oracle JDK.
TLS/SSL Client Truststore File Password	The password for the truststore file.

6. Click **Save Changes** to commit the changes.

Step 3: Restart Cloudera Manager and Services

You must restart both Cloudera Manager Server and the Cloudera Management Service for TLS encryption to work. Otherwise, the Cloudera Management Services (such as Host Monitor and Service Monitor) cannot communicate with Cloudera Manager Server.

1. Restart the Cloudera Manager Server by running `service cloudera-scm-server restart` on the Cloudera Manager Server host.
2. After the restart completes, connect to the Cloudera Manager Admin Console using the HTTPS URL (for example: `https://cm01.example.com:7183`). If you used an internal CA-signed certificate, you must configure your browser to trust the certificate. Otherwise, you will see a warning in your browser any time you access the Cloudera Manager Administration Console. By default, certificates issued by public commercial CAs are trusted by most browsers, and no additional configuration is necessary if your certificate is signed by one of them.

- Restart the Cloudera Management Service (**Cloudera Management Service > Actions > Restart**).

Configuring TLS Encryption for Cloudera Manager Agents

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

Use the following procedure to encrypt the communication between Cloudera Manager Server and Cloudera Manager Agents:

Step 1: Enable TLS Encryption for Agents in Cloudera Manager

Configure the TLS properties for Cloudera Manager Agents.

- Log in to the Cloudera Manager Admin Console.
- Select **Administration > Settings**.
- Select the **Security** category.
- Select the **Use TLS Encryption for Agents** option.
- Click Save Changes.

Step 2: Enable TLS on Cloudera Manager Agent Hosts

To enable TLS between the Cloudera Manager agents and Cloudera Manager, you must specify values for the TLS properties in the `/etc/cloudera-scm-agent/config.ini` configuration file on all agent hosts.

- On each agent host, open the `/etc/cloudera-scm-agent/config.ini` configuration file and set the `use_tls` parameter in the `[Security]` section as follows:

```
use_tls=1
```

Alternatively, you can edit the `config.ini` file on one host, and then copy it to the other hosts because this file by default does not contain host-specific information. If you have modified properties such as `listening_hostname` or `listening_ip` address in `config.ini`, you must edit the file individually on each host.

Step 3: Restart Cloudera Manager Server and Agents

Restart the Cloudera Manager Server with the following command to activate the TLS configuration settings.

```
$ sudo service cloudera-scm-server restart
```

On each agent host, restart the Cloudera Manager agent service:

```
$ sudo service cloudera-scm-agent restart
```

Step 4: Verify that the Cloudera Manager Server and Agents are Communicating

In the Cloudera Manager Admin Console, go to **Hosts > All Hosts**. If you see successful heartbeats reported in the **Last Heartbeat** column after restarting the agents, TLS encryption is working properly.

Enabling Server Certificate Verification on Cloudera Manager Agents

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

If you have completed the previous sections, communication between Cloudera Manager server and the agents is encrypted, but the certificate authenticity is not verified. For full security, you must configure the agents to verify the Cloudera Manager server certificate. If you are using a server certificate signed by an internal certificate authority (CA), you must configure the agents to trust that CA:

How to Configure TLS Encryption for Cloudera Manager

1. On each agent host, open the `/etc/cloudera-scm-agent/config.ini` configuration file, and then uncomment and set the following property:

```
verify_cert_file=/opt/cloudera/security/pki/rootca.cert.pem
```

Alternatively, you can edit the `config.ini` file on one host, and then copy it to the other hosts because this file by default does not contain host-specific information. If you have modified properties such as `listening_hostname` or `listening_ip` address in `config.ini`, you must edit the file individually on each host.

2. Restart the Cloudera Manager agents. On each agent host, run the following command:

```
$ sudo service cloudera-scm-agent restart
```

3. Restart the Cloudera Management Service. On the **Home > Status** tab, click



to the right of the Cloudera Management Service and select **Restart**.

4. Verify that the Cloudera Manager server and agents are communicating. In the Cloudera Manager Admin Console, go to **Hosts > All Hosts**. If you see successful heartbeats reported in the **Last Heartbeat** column after restarting the agents and management service, TLS verification is working properly. If not, check the agent log (`/var/log/cloudera-scm-agent/cloudera-scm-agent.log`) for errors.

Configuring Agent Certificate Authentication



Important: Repeat this procedure on each agent host.

Without certificate authentication, a malicious user can add a host to Cloudera Manager by installing the Cloudera Manager agent software and configuring it to communicate with Cloudera Manager Server. To prevent this, you must configure Cloudera Manager to trust the agent certificates.

Step 1: Export the Private Key to a File

On each Cloudera Manager Agent host, use the `keytool` utility to export the private key and certificate to a PKCS12 file, which can then be split up into individual key and certificate files using the `openssl` command:

1. Export the private key and certificate:

```
$ keytool -importkeystore -srckeystore /opt/cloudera/security/pki/${hostname -f}-agent.jks \
-srcstorepass password -srckeypass password -destkeystore \
/opt/cloudera/security/pki/${hostname -f}-agent.p12 \
-deststoretype PKCS12 -srccalias ${hostname -f}-agent -deststorepass password -destkeypass \
password
```

2. Use the `openssl` command to export the private key into its own file:

```
$ openssl pkcs12 -in /opt/cloudera/security/pki/${hostname -f}-agent.p12 -passin \
pass:password -nocerts \
-out /opt/cloudera/security/pki/${hostname -f}-agent.key -passout pass:password
```

3. Create a symbolic link for the `.key` file:

```
$ ln -s /opt/cloudera/security/pki/${hostname -f}-agent.key \
/opt/cloudera/security/pki/agent.key
```

This allows you to use the same `/etc/cloudera-scm-agent/config.ini` file on all agent hosts rather than maintaining a file for each agent.

Step 2: Create a Password File

The Cloudera Manager agent obtains the password from a text file, not from a command line parameter or environment variable. The password file allows you to use file permissions to protect the password. For example, run the following commands on each Cloudera Manager Agent host, or run them on one host and copy the file to the other hosts:

```
$ echo "password" > /etc/cloudera-scm-agent/agentkey.pw
$ sudo chown root:root /etc/cloudera-scm-agent/agentkey.pw
$ sudo chmod 440 /etc/cloudera-scm-agent/agentkey.pw
```

Replace `password` with the password you created in [Step 1: Export the Private Key to a File](#) on page 46.

Step 3: Configure the Agent to Use Private Keys and Certificates

On a Cloudera Manager Agent, open the `/etc/cloudera-scm-agent/config.ini` configuration file and edit the following properties.

Property	Example Value	Description
<code>client_key_file</code>	<code>/etc/cloudera-scm-agent/agentkey</code>	Path to the private key file.
<code>client_keypw_file</code>	<code>/etc/cloudera-scm-agent/agentkey.pw</code>	Path to the private key password file.
<code>client_cert_file</code>	<code>/etc/cloudera-scm-agent/agentcert</code>	Path to the client certificate file.

Copy the file to all other cluster hosts. If you have modified properties such as `listening_hostname` or `listening_ip` address in `config.ini`, you must edit the file individually on each host.

Step 4: Enable Agent Certificate Authentication

1. Log in to the Cloudera Manager Admin Console.
2. Select **Administration > Settings**.
3. Click the **Security** category.
4. Configure the following TLS settings:

Setting	Description
Use TLS Authentication of Agents to Server	Select this option to enable TLS authentication of agents to the server.
Cloudera Manager TLS/SSL Certificate Trust Store File	Specify the full filesystem path to the <code>jssecacerts</code> file located on the Cloudera Manager Server host. For example, <code>/usr/java/jdk1.7.0_67-cloudera/jre/lib/security/jssecacerts</code> .
Cloudera Manager TLS/SSL Certificate Trust Store Password	Specify the password for the <code>jssecacerts</code> truststore.

5. Click **Save Changes** to save the settings.

Step 5: Restart Cloudera Manager Server and Agents

1. On the Cloudera Manager server host, restart the Cloudera Manager server:

```
$ sudo service cloudera-scm-server restart
```

How to Configure TLS Encryption for Cloudera Manager

2. On every agent host, restart the Cloudera Manager agent:

```
$ sudo service cloudera-scm-agent restart
```

Step 6: Verify that Cloudera Manager Server and Agents are Communicating

In the Cloudera Manager Admin Console, go to **Hosts > All Hosts**. If you see successful heartbeats reported in the **Last Heartbeat** column after restarting the agents and server, TLS certificate authentication is working properly. If not, check the agent log (`/var/log/cloudera-scm-agent/cloudera-scm-agent.log`) for errors.

Configuring Authentication

The purpose of authentication in Hadoop, as in other systems, is simply to prove that a user or service is who he or she claims to be.

Typically, authentication in enterprises is managed through a single distributed system, such as a Lightweight Directory Access Protocol (LDAP) directory. LDAP authentication consists of straightforward username/password services backed by a variety of storage systems, ranging from file to database.

A common enterprise-grade authentication system is Kerberos. Kerberos provides strong security benefits including capabilities that render intercepted authentication packets unusable by an attacker. It virtually eliminates the threat of impersonation by never sending a user's credentials in cleartext over the network.

Several components of the Hadoop ecosystem are converging to use Kerberos authentication with the option to manage and store credentials in LDAP or AD. For example, Microsoft's Active Directory (AD) is an LDAP directory that also provides Kerberos authentication for added security.

Before you use this guide to configure Kerberos on your cluster, ensure you have a working KDC (MIT KDC or Active Directory), set up. You can then use Cloudera Manager's Kerberos wizard to automate several aspects of Kerberos configuration on your cluster.



Important:

- You can use either Cloudera Manager or the following command-line instructions to complete this configuration.
- This information applies specifically to CDH 5.6.x. If you use an earlier version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

Configuring Authentication in Cloudera Manager

Before You Begin

When you configure authentication and authorization on a cluster, Cloudera Manager Server sends sensitive information over the network to cluster hosts, such as Kerberos keytabs and configuration files that contain passwords. To secure this transfer, you must configure TLS encryption between Cloudera Manager Server and all cluster hosts. For instruction on enabling TLS encryption for Cloudera Manager, see [How to Configure TLS Encryption for Cloudera Manager](#) on page 40.

Why Use Cloudera Manager to Implement Kerberos Authentication?

If you do not use Cloudera Manager to implement Hadoop security, you must manually create and deploy the Kerberos principals and keytabs on *every* host in your cluster. If you have a large number of hosts, this can be a time-consuming and error-prone process. After creating and deploying the keytabs, you must also manually configure properties in the `core-site.xml`, `hdfs-site.xml`, `mapred-site.xml`, and `taskcontroller.cfg` files on *every* host in the cluster to enable and configure Hadoop security in HDFS and MapReduce. You must also manually configure properties in the `oozie-site.xml` and `hue.ini` files on certain cluster hosts in order to enable and configure Hadoop security in Oozie and Hue.

Cloudera Manager enables you to automate all of those manual tasks. Cloudera Manager can automatically create and deploy a keytab file for the `hdfs` user and a keytab file for the `mapred` user on every host in your cluster, as well as keytab files for the `oozie` and `hue` users on select hosts. The `hdfs` keytab file contains entries for the `hdfs` principal and a `host` principal, and the `mapred` keytab file contains entries for the `mapred` principal and a `host` principal. The `host` principal will be the same in both keytab files. The `oozie` keytab file contains entries for the `oozie` principal and a `HTTP` principal. The `hue` keytab file contains an entry for the `hue` principal. Cloudera Manager can also

Configuring Authentication

automatically configure the appropriate properties in the `core-site.xml`, `hdfs-site.xml`, `mapred-site.xml`, and `taskcontroller.cfg` files on every host in the cluster, and the appropriate properties in `oozie-site.xml` and `hue.ini` for select hosts. Lastly, Cloudera Manager can automatically start up the NameNode, DataNode, Secondary NameNode, JobTracker, TaskTracker, Oozie Server, and Hue roles once all the appropriate configuration changes have been made.

Ways to Configure Kerberos Authentication Using Cloudera Manager

You can use *one* of the following ways to set up Kerberos authentication on your cluster using Cloudera Manager:

- Cloudera Manager 5.1 introduced a new wizard to automate the procedure to set up Kerberos on a cluster. Using the KDC information you enter, the wizard will create new principals and keytab files for your CDH services. The wizard can be used to deploy the `krb5.conf` file cluster-wide, and automate other manual tasks such as stopping all services, deploying client configuration and restarting all services on the cluster.


If you want to use the Kerberos wizard, follow the instructions at [Enabling Kerberos Authentication Using the Wizard](#) on page 59.

- If you do not want to use the Kerberos wizard, follow the instructions at [Enabling Kerberos Authentication Without the Wizard](#) on page 75.

Cloudera Manager User Accounts

Minimum Required Role: [User Administrator](#) (also provided by **Full Administrator**)

Access to Cloudera Manager features is controlled by user accounts. A user account identifies how a user is authenticated and determines what privileges are granted to the user.

When you are logged in to the Cloudera Manager Admin Console, the username you are logged in as is located at the far right of the top navigation bar—for example, if you are logged in as `admin` you will see  `admin`.

A user with the User Administrator or Full Administrator role manages user accounts through the **Administration > Users** page. View active user sessions on the **User Sessions** tab.

User Authentication

Cloudera Manager provides several mechanisms for authenticating users. You can configure Cloudera Manager to authenticate users against the Cloudera Manager database or against an external authentication service. The external authentication service can be an LDAP server (Active Directory or an OpenLDAP compatible directory), or you can specify another external service. Cloudera Manager also supports using the Security Assertion Markup Language (SAML) to enable single sign-on.

If you are using LDAP or another external service, you can configure Cloudera Manager so that it can use both methods of authentication (internal database and external service), and you can determine the order in which it performs these searches. If you select an external authentication mechanism, Full Administrator users can always authenticate against the Cloudera Manager database. This prevents locking everyone out if the authentication settings are misconfigured, such as with a bad LDAP URL.

With external authentication, you can restrict login access to members of specific groups, and can specify groups whose members are automatically given Full Administrator access to Cloudera Manager.

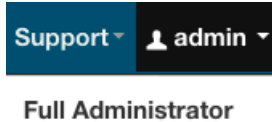
Users accounts in the Cloudera Manager database page show **Cloudera Manager** in the User Type column. User accounts in an LDAP directory or other external authentication mechanism show **External** in the User Type column.

User Roles

User accounts include the user's role, which determines the Cloudera Manager features visible to the user and the actions the user can perform. All tasks in the Cloudera Manager documentation indicate which role is required to perform the task. For more information about user roles, see [Cloudera Manager User Roles](#) on page 344.

Determining the Role of the Currently Logged in User

1. Click the logged-in username at the far right of the top navigation bar. The role displays under the username. For example:



Changing the Logged-In Internal User Password

1. Click the logged-in username at the far right of the top navigation bar and select **Change Password**.
2. Enter the current password and a new password twice, and then click **OK**.

Adding an Internal User Account

1. Select **Administration > Users**.
2. Click the **Add User** button.
3. Enter a username and password.
4. In the Role drop-down menu, select a role for the new user.
5. Click **Add**.

Assigning User Roles

1. Select **Administration > Users**.
2. Check the checkbox next to one or more usernames.
3. Select **Actions for Selected > Assign User Roles**.
4. In the drop-down menu, select the role.
5. Click the **Assign Role** button.

Changing an Internal User Account Password

1. Select **Administration > Users**.
2. Click the **Change Password** button next to a username with User Type **Cloudera Manager**.
3. Type the new password and repeat it to confirm.
4. Click the **Update** button to make the change.

Deleting Internal User Accounts

1. Select **Administration > Users**.
2. Check the checkbox next to one or more usernames with User Type **Cloudera Manager**.
3. Select **Actions for Selected > Delete**.
4. Click the **OK** button. (There is no confirmation of the action.)

Viewing User Sessions

1. Select **Administration > Users**.
2. Click the tab **User Sessions**.

Configuring External Authentication for Cloudera Manager

Minimum Required Role: [User Administrator](#) (also provided by **Full Administrator**)



Important: This feature is available only with a Cloudera Enterprise license; it is not available in Cloudera Express. For information on Cloudera Enterprise licenses, see [Managing Licenses](#).

Cloudera Manager supports user authentication against an internal database and against an external service. The following sections describe how to configure the supported external services.

Configuring Authentication Using Active Directory

1. Select **Administration > Settings**.
2. In the left-hand column, select the **External Authentication** category.
3. In the **Authentication Backend Order** field, select the order in which Cloudera Manager should attempt its authentication. You can choose to authenticate users using just one of the methods (using Cloudera Manager's own database is the default), or you can set it so that if the user cannot be authenticated by the first method, it will attempt using the second method.
4. For **External Authentication Type**, select **Active Directory**.
5. In the **LDAP URL** property, provide the URL of the Active Directory server.
6. In the **Active Directory NT Domain** property, provide the NT domain to authenticate against.

LDAP URL and **Active Directory NT Domain** are the only settings required to allow anyone in AD to log in to Cloudera Manager. For example, if you set LDAP URL to `ldap://adserver.example.com` and the Active Directory NT Domain to `ADREALM.EXAMPLE.COM`, AD users should now be able to log into Cloudera Manager using just their username, such as `sampleuser`. They no longer require the complete string: `sampleuser@ADREALM.EXAMPLE.COM`.

7. In the **LDAP User Groups** property, optionally provide a comma-separated list of case-sensitive LDAP group names. If this list is provided, only users who are members of one or more of the groups in the list will be allowed to log into Cloudera Manager. If this property is left empty, all authenticated LDAP users will be able to log into Cloudera Manager. For example, if there is a group called `CN=ClouderaManagerUsers,OU=Groups,DC=corp,DC=com`, add the group name `ClouderaManagerUsers` to the **LDAP User Groups** list to allow members of that group to log in to Cloudera Manager.
8. To automatically assign a [role](#) to users when they log in, provide a comma-separated list of LDAP group names in the following properties:
 - LDAP Full Administrator Groups
 - LDAP User Administrator Groups
 - LDAP Cluster Administrator Groups
 - LDAP BDR Administrator Groups
 - LDAP Configurator Groups
 - LDAP Navigator Administrator Groups
 - LDAP Operator Groups
 - LDAP Limited Operator Groups
 - LDAP Auditor Groups

If you specify groups in these properties, users must also be a member of at least one of the groups specified in the **LDAP User Groups** property or they will not be allowed to log in. If these properties are left empty, users will be assigned to the Read-Only role and any other role assignment must be performed manually by an Administrator.



Note: A user that is added to an LDAP group will not automatically be assigned the corresponding role in the internal Cloudera Manager database. Hence, the Users page in Cloudera Manager will display such users' roles as Read-Only, as this page only queries the Cloudera Manager database, and not LDAP.

9. [Restart](#) the Cloudera Manager Server.

Configuring Authentication Using an OpenLDAP-compatible Server

For an OpenLDAP-compatible directory, you have several options for searching for users and groups:

- You can specify a single base Distinguished Name (DN) and then provide a "Distinguished Name Pattern" to use to match a specific user in the LDAP directory.
- Search filter options let you search for a particular user based on somewhat broader search criteria – for example Cloudera Manager users could be members of different groups or organizational units (OUs), so a single pattern

won't find all those users. Search filter options also let you find all the groups to which a user belongs, to help determine if that user should have login or admin access.

1. Select **Administration > Settings**.
2. In the left-hand column, select the **External Authentication** category.
3. In the **Authentication Backend Order** field, select the order in which Cloudera Manager should attempt its authentication. You can choose to authenticate users using just one of the methods (using Cloudera Manager's own database is the default), or you can set it so that if the user cannot be authenticated by the first method, it will attempt using the second method.
4. For **External Authentication Type**, select **LDAP**.
5. In the **LDAP URL** property, provide the URL of the LDAP server and (optionally) the base Distinguished Name (DN) (the search base) as part of the URL — for example `ldap://ldap-server.corp.com/dc=corp,dc=com`.
6. If your server does not allow anonymous binding, provide the user DN and password to be used to bind to the directory. These are the **LDAP Bind User Distinguished Name** and **LDAP Bind Password** properties. By default, Cloudera Manager assumes anonymous binding.
7. Use one of the following methods to search for users and groups:

- You can search using User or Group search filters, using the **LDAP User Search Base**, **LDAP User Search Filter**, **LDAP Group Search Base** and **LDAP Group Search Filter** settings. These allow you to combine a base DN with a search filter to allow a greater range of search targets.

For example, if you want to authenticate users who may be in one of multiple OUs, the search filter mechanism will allow this. You can specify the User Search Base DN as `dc=corp,dc=com` and the user search filter as `uid={0}`. Then Cloudera Manager will search for the user anywhere in the tree starting from the Base DN. Suppose you have two OUs—`ou=Engineering` and `ou=Operations`—Cloudera Manager will find User "foo" if it exists in either of these OUs, that is, `uid=foo,ou=Engineering,dc=corp,dc=com` or `uid=foo,ou=Operations,dc=corp,dc=com`.

You can use a user search filter along with a DN pattern, so that the search filter provides a fallback if the DN pattern search fails.

The Groups filters let you search to determine if a DN or username is a member of a target group. In this case, the filter you provide can be something like `member={0}` where `{0}` will be replaced with the **DN** of the user you are authenticating. For a filter requiring the username, `{1}` may be used, as `memberUId={1}`. This will return a list of groups the user belongs to, which will be compared to the list in the group properties discussed in [step 8 of Configuring Authentication Using Active Directory](#) on page 52.

OR

- Alternatively, specify a single base Distinguished Name (DN) and then provide a "Distinguished Name Pattern" in the **LDAP Distinguished Name Pattern** property.

Use `{0}` in the pattern to indicate where the username should go. For example, to search for a distinguished name where the `uid` attribute is the username, you might provide a pattern similar to `uid={0},ou=People,dc=corp,dc=com`. Cloudera Manager substitutes the name provided at login into this pattern and performs a search for that specific user. So if a user provides the username "foo" at the Cloudera Manager login page, Cloudera Manager will search for the DN `uid=foo,ou=People,dc=corp,dc=com`.

If you provided a base DN along with the URL, the pattern only needs to specify the rest of the DN pattern. For example, if the URL you provide is `ldap://ldap-server.corp.com/dc=corp,dc=com`, and the pattern is `uid={0},ou=People`, then the search DN will be `uid=foo,ou=People,dc=corp,dc=com`.

8. [Restart](#) the Cloudera Manager Server.

Configuring Cloudera Manager to Use LDAPS

If the LDAP server certificate has been signed by a trusted Certificate Authority (that is, VeriSign, GeoTrust, and so on), steps 1 and 2 below may not be necessary.

Configuring Authentication

1. Copy the CA certificate file to the Cloudera Manager Server host.
2. Import the CA certificate(s) from the CA certificate file to the local truststore. The default truststore is located in the `$JAVA_HOME/jre/lib/security/cacerts` file. This contains the default CA information shipped with the JDK. Create an alternate default file called `jssecacerts` in the same location as the `cacerts` file. You can now safely append CA certificates for any private or public CAs not present in the default `cacerts` file, while keeping the original file intact.

For our example, we will follow this recommendation by copying the default `cacerts` file into the new `jssecacerts` file, and then importing the CA certificate to this alternate truststore.

```
$ cp $JAVA_HOME/jre/lib/security/cacerts \  
$JAVA_HOME/jre/lib/jssecacerts
```

```
$ /usr/java/latest/bin/keytool -import -alias nt_domain_name \  
-keystore /usr/java/latest/jre/lib/security/jssecacerts -file path_to_cert
```



Note:

- The default password for the `cacerts` store is **changeit**.
- The alias can be any name (not just the domain name).

3. Configure the **LDAP URL** property to use `ldaps://ldap_server` instead of `ldap://ldap_server`.
4. [Restart](#) the Cloudera Manager Server.

Configuring Authentication Using an External Program

You can configure Cloudera Manager to use an external authentication program of your own choosing. Typically, this may be a custom script that interacts with a custom authentication service. Cloudera Manager will call the external program with the username as the first command line argument. The password is passed over `stdin`. Cloudera Manager assumes the program will return the following exit codes identifying the user role for a successful authentication:

- 0 - Read-Only
- 1 - Full Administrator
- 2 - Limited Operator
- 3 - Operator
- 4 - Configurator
- 5 - Cluster Administrator
- 6 - BDR Administrator
- 7 - Navigator Administrator
- 8 - User Administrator
- 9 - Auditor
- 10 - Key Administrator

and a negative value is returned for a failure to authenticate.

To configure authentication using an external program:

1. Select **Administration > Settings**.
2. In the left-hand column, select the **External Authentication** category.
3. In the **Authentication Backend Order** field, select the order in which Cloudera Manager should attempt its authentication. You can choose to authenticate users using just one of the methods (using Cloudera Manager's own database is the default), or you can set it so that if the user cannot be authenticated by the first method, it will attempt using the second method.
4. For **External Authentication Type**, select **External Program**.
5. Provide a path to the external program in the **External Authentication Program Path** property.

Configuring Authentication Using SAML

Cloudera Manager supports the Security Assertion Markup Language (SAML), an XML-based open standard data format for exchanging authentication and authorization data between parties, in particular, between an identity provider (IDP) and a service provider (SP). The SAML specification defines three roles: the principal (typically a user), the IDP, and the SP. In the use case addressed by SAML, the principal (user agent) requests a service from the service provider. The service provider requests and obtains an identity assertion from the IDP. On the basis of this assertion, the SP can make an access control decision—in other words it can decide whether to perform some service for the connected principal.

The primary SAML use case is called web browser single sign-on (SSO). A user wielding a user agent (usually a web browser) requests a web resource protected by a SAML SP. The SP, wishing to know the identity of the requesting user, issues an authentication request to a SAML IDP through the user agent. In the context of this terminology, Cloudera Manager operates as a SP. This topic discusses the Cloudera Manager part of the configuration process; it assumes that you are familiar with SAML and SAML configuration in a general sense, and that you have a functioning IDP already deployed.



Note:

- Cloudera Manager supports both SP- and IDP-initiated SSO.
- The logout action in Cloudera Manager will send a single-logout request to the IDP.
- SAML authentication has been tested with specific configurations of CA Single Sign-On (formerly SiteMinder) and Shibboleth. While SAML is a standard, there is a great deal of variability in configuration between different IDP products, so it is possible that other IDP implementations, or other configurations of CA Single Sign-On and Shibboleth, may not interoperate with Cloudera Manager.
- To bypass SSO if SAML configuration is incorrect or not working, you can login using a Cloudera Manager local account using the URL: `http://cm_host:7180/cm/localLogin`

Setting up Cloudera Manager to use SAML requires the following steps.

Preparing Files

You will need to prepare the following files and information, and provide these to Cloudera Manager:

- A Java keystore containing a private key for Cloudera Manager to use to sign/encrypt SAML messages. For guidance on creating Java keystores, see [Creating Java Keystores and Truststores](#) on page 207.
- The SAML metadata XML file from your IDP. This file must contain the public certificates needed to verify the sign/encrypt key used by your IDP per the SAML Metadata Interoperability Profile. For example, if you are using the Shibboleth IdP, the metadata file is available at: `https://<IdPHOST>:8080/idp/shibboleth`.



Note: For guidance on how to obtain the metadata XML file from your IDP, either contact your IDP administrator or consult the documentation for the version of the IDP you are using.

- The entity ID that should be used to identify the Cloudera Manager instance
- How the user ID is passed in the SAML authentication response:
 - As an attribute. If so, what identifier is used.
 - As the NameID.
- The method by which the Cloudera Manager role will be established:
 - From an attribute in the authentication response:
 - What identifier will be used for the attribute
 - What values will be passed to indicate each role
 - From an external script that will be called for each use:

Configuring Authentication

- The script takes user ID as \$1
- The script sets an exit code to reflect successful authentication of the assigned role:
 - 0 - Full Administrator
 - 1 - Read-Only
 - 2 - Limited Operator
 - 3 - Operator
 - 4 - Configurator
 - 5 - Cluster Administrator
 - 6 - BDR Administrator
 - 7 - Navigator Administrator
 - 8 - User Administrator
 - 9 - Auditor
 - 10 - Key Administrator

and a negative value is returned for a failure to authenticate.

Configuring Cloudera Manager

1. Select **Administration > Settings**.
2. In the left-hand column, select the **External Authentication** category.
3. Set the **External Authentication Type** property to **SAML** (the Authentication Backend Order property is ignored for SAML).
4. Set the **Path to SAML IDP Metadata File** property to point to the IDP metadata file.
5. Set the **Path to SAML Keystore File** property to point to the Java keystore prepared earlier.
6. In the **SAML Keystore Password** property, set the keystore password.
7. In the **Alias of SAML Sign/Encrypt Private Key** property, set the alias used to identify the private key for Cloudera Manager to use.
8. In the **SAML Sign/Encrypt Private Key Password** property, set the private key password.
9. Set the **SAML Entity ID** property if:
 - There is more than one Cloudera Manager instance being used with the same IDP (each instance needs a different entity ID).
 - Entity IDs are assigned by organizational policy.
10. In the **Source of User ID in SAML Response** property, set whether the user ID will be obtained from an attribute or the NameID.

If an attribute will be used, set the attribute name in the **SAML attribute identifier for user ID** property. The default value is the normal OID used for user IDs and so may not need to be changed.
11. In the **SAML Role assignment mechanism** property, set whether the role assignment will be done from an attribute or an external script.
 - If an attribute will be used:
 - In the **SAML attribute identifier for user role** property, set the attribute name if necessary. The default value is the normal OID used for OrganizationalUnits and so may not need to be changed.
 - In the **SAML Attribute Values for Roles** property, set which attribute values will be used to indicate the user role.
 - If an external script will be used, set the path to that script in the **Path to SAML Role Assignment Script** property. Make sure that the script is executable (an executable binary is fine - it doesn't need to be a shell script).
12. Save the changes. Cloudera Manager will run a set of validations that ensure it can find the metadata XML and the keystore, and that the passwords are correct. If you see a validation error, correct the problem before proceeding.

3. [Restart](#) the Cloudera Manager Server.

Configuring the IDP

After the Cloudera Manager Server is restarted, it will attempt to redirect to the IDP login page instead of showing the normal CM page. This may or may not succeed, depending on how the IDP is configured. In either case, the IDP will need to be configured to recognize CM before authentication will actually succeed. The details of this process are specific to each IDP implementation - refer to your IDP documentation for details. If you are using the Shibboleth IDP, information on configuring the IDP to communicate with a Service Provider is available [here](#).

1. Download the Cloudera Manager's SAML metadata XML file from `http://hostname:7180/saml/metadata`.
2. Inspect the metadata file and ensure that any URLs contained in the file can be resolved by users' web browsers. The IDP will redirect web browsers to these URLs at various points in the process. If the browser cannot resolve them, authentication will fail. If the URLs are incorrect, you can manually fix the XML file or set the Entity Base URL in the CM configuration to the right value, and then re-download the file.
3. Provide this metadata file to your IDP using whatever mechanism your IDP provides.
4. Ensure that the IDP has access to whatever public certificates are necessary to validate the private key that was provided to Cloudera Manager earlier.
5. Ensure that the IDP is configured to provide the User ID and Role using the attribute names that Cloudera Manager was configured to expect, if relevant.
6. Ensure the changes to the IDP configuration have taken effect (a restart may be necessary).

Verifying Authentication and Authorization

1. Return to the Cloudera Manager Admin Console and refresh the login page.
2. Attempt to log in with credentials for a user that is entitled. The authentication should complete and you should see the **Home > Status** tab.
3. If authentication fails, you will see an IDP provided error message. Cloudera Manager is not involved in this part of the process, and you must ensure the IDP is working correctly to complete the authentication.
4. If authentication succeeds but the user is not authorized to use Cloudera Manager, they will be taken to an error page by Cloudera Manager that explains the situation. If an user who should be authorized sees this error, then you will need to verify their role configuration, and ensure that it is being properly communicated to Cloudera Manager, whether by attribute or external script. The Cloudera Manager log will provide details on failures to establish a user's role. If any errors occur during role mapping, Cloudera Manager will assume the user is unauthorized.

Kerberos Concepts - Principals, Keytabs and Delegation Tokens

This section describes how Hadoop uses Kerberos principals and keytabs for user authentication. It also briefly describes how Hadoop uses delegation tokens to authenticate jobs at execution time, to avoid overwhelming the KDC with authentication requests for each job.

Kerberos Principals

A user in Kerberos is called a *principal*, which is made up of three distinct components: the primary, instance, and realm. A *Kerberos principal* is used in a Kerberos-secured system to represent a unique identity. The first component of the principal is called the *primary*, or sometimes the user component. The primary component is an arbitrary string and may be the operating system username of the user or the name of a service. The primary component is followed by an optional section called the *instance*, which is used to create principals that are used by users in special roles or to define the host on which a service runs, for example. An instance, if it exists, is separated from the primary by a slash and then the content is used to disambiguate multiple principals for a single user or service. The final component of the principal is the realm. The *realm* is similar to a domain in DNS in that it logically defines a related group of objects, although rather than hostnames as in DNS, the Kerberos realm defines a group of principals. Each realm can have its own settings including the location of the KDC on the network and supported encryption algorithms. Large organizations commonly create distinct realms to delegate administration of a realm to a group within the enterprise. Realms, by convention, are written in uppercase characters.

Kerberos assigns tickets to Kerberos principals to enable them to access Kerberos-secured Hadoop services. For the Hadoop daemon principals, the principal names should be of the format

Configuring Authentication

`service/fully.qualified.domain.name@YOUR-REALM.COM`. Here, *service* in the `service/fully.qualified.domain.name@YOUR-REALM.COM` principal refers to the username of an existing Unix account that is used by Hadoop daemons, such as `hdfs` or `mapred`.

Human users who want to access the Hadoop cluster also need to have Kerberos principals of the format, `username@YOUR-REALM.COM`; in this case, *username* refers to the username of the user's Unix account, such as `joe` or `jane`. Single-component principal names (such as `joe@YOUR-REALM.COM`) are typical for client user accounts. Hadoop does not support more than two-component principal names.

Kerberos Keytabs

A **keytab** is a file containing pairs of Kerberos principals and an encrypted copy of that principal's key. A keytab file for a Hadoop daemon is unique to each host since the principal names include the hostname. This file is used to authenticate a principal on a host to Kerberos without human interaction or storing a password in a plain text file. Because having access to the keytab file for a principal allows one to act as that principal, access to the keytab files should be tightly secured. They should be readable by a minimal set of users, should be stored on local disk, and should not be included in host backups, unless access to those backups is as secure as access to the local host.

Delegation Tokens

Users in a Hadoop cluster authenticate themselves to the NameNode using their Kerberos credentials. However, once the user is authenticated, each job subsequently submitted must also be checked to ensure it comes from an authenticated user. Since there could be a time gap between a job being submitted and the job being executed, during which the user could have logged off, user credentials are passed to the NameNode using delegation tokens that can be used for authentication in the future.

Delegation tokens are a secret key shared with the NameNode, that can be used to impersonate a user to get a job executed. While these tokens can be renewed, new tokens can only be obtained by clients authenticating to the NameNode using Kerberos credentials. By default, delegation tokens are only valid for a day. However, since jobs can last longer than a day, each token specifies a JobTracker as a *renewer* which is allowed to renew the delegation token once a day, until the job completes, or for a maximum period of 7 days. When the job is complete, the JobTracker requests the NameNode to cancel the delegation token.

Token Format

The NameNode uses a random `masterKey` to generate delegation tokens. All active tokens are stored in memory with their expiry date (`maxDate`). Delegation tokens can either expire when the current time exceeds the expiry date, or, they can be canceled by the owner of the token. Expired or canceled tokens are then deleted from memory. The `sequenceNumber` serves as a unique ID for the tokens. The following section describes how the Delegation Token is used for authentication.

```
TokenID = {ownerID, renewerID, issueDate, maxDate, sequenceNumber}
TokenAuthenticator = HMAC-SHA1(masterKey, TokenID)
Delegation Token = {TokenID, TokenAuthenticator}
```

Authentication Process

To begin the authentication process, the client first sends the `TokenID` to the NameNode. The NameNode uses this `TokenID` and the `masterKey` to once again generate the corresponding `TokenAuthenticator`, and consequently, the Delegation Token. If the NameNode finds that the token already exists in memory, and that the current time is less than the expiry date (`maxDate`) of the token, then the token is considered valid. If valid, the client and the NameNode will then authenticate each other by using the `TokenAuthenticator` that they possess as the secret key, and MD5 as the protocol. Since the client and NameNode do not actually exchange `TokenAuthenticators` during the process, even if authentication fails, the tokens are not compromised.

Token Renewal

Delegation tokens must be renewed periodically by the designated renewer (`renewerID`). For example, if a JobTracker is the designated renewer, the JobTracker will first authenticate itself to the NameNode. It will then send the token to be authenticated to the NameNode. The NameNode verifies the following information before renewing the token:

- The JobTracker requesting renewal is the same as the one identified in the token by `renewerID`.
- The TokenAuthenticator generated by the NameNode using the `TokenID` and the `masterKey` matches the one previously stored by the NameNode.
- The current time must be less than the time specified by `maxDate`.

If the token renewal request is successful, the NameNode sets the new expiry date to `min(current time+renew period, maxDate)`. If the NameNode was restarted at any time, it will have lost all previous tokens from memory. In this case, the token will be saved to memory once again, this time with a new expiry date. Hence, designated renewers must renew all tokens with the NameNode after a restart, and before relaunching any failed tasks.

A designated renewer can also revive an expired or canceled token as long as the current time does not exceed `maxDate`. The NameNode cannot tell the difference between a token that was canceled, or has expired, and one that was erased from memory due to a restart, since only the `masterKey` persists in memory. The `masterKey` must be updated regularly.

Enabling Kerberos Authentication Using the Wizard

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)



Important: Ensure you have secured communication between the Cloudera Manager Server and Agents before you enable Kerberos on your cluster. Kerberos keytabs are sent from the Cloudera Manager Server to the Agents, and must be encrypted to prevent potential misuse of leaked keytabs. For instructions on securing this transfer with TLS encryption, see [How to Configure TLS Encryption for Cloudera Manager](#) on page 40.

This guide describes how to use Cloudera Manager and the Kerberos wizard (introduced in Cloudera Manager 5.1.0) to automate many of the manual tasks of implementing Kerberos security on your CDH cluster.

- **Prerequisites** - These instructions assume you know how to install and configure Kerberos, you already have a working Kerberos key distribution center (KDC) and realm setup, and that you've installed the following Kerberos client packages on all cluster hosts and hosts that will be used to access the cluster, depending on the OS in use.

OS	Packages Required
RHEL 7 Compatible, RHEL 6 Compatible, RHEL 5 Compatible	<ul style="list-style-type: none"> • <code>openldap-clients</code> on the Cloudera Manager Server host • <code>krb5-workstation, krb5-libs</code> on ALL hosts
SLES	<ul style="list-style-type: none"> • <code>openldap2-client</code> on the Cloudera Manager Server host • <code>krb5-client</code> on ALL hosts
Ubuntu or Debian	<ul style="list-style-type: none"> • <code>ldap-utils</code> on the Cloudera Manager Server host • <code>krb5-user</code> on ALL hosts
Windows	<ul style="list-style-type: none"> • <code>krb5-workstation, krb5-libs</code> on ALL hosts

Furthermore, Oozie and Hue require that the realm support renewable tickets. Cloudera Manager supports setting up kerberized clusters with MIT KDC and Active Directory.



Important: If you want to integrate Kerberos directly with Active Directory, ensure you have support from your AD administration team to do so. This includes any future support required to troubleshoot issues such as Kerberos TGT/TGS ticket renewal, access to KDC logs for debugging and so on.

For more information about using Active Directory, refer the section below on **Considerations when using an Active Directory KDC** and the [Microsoft AD documentation](#).

For more information about installing and configuring MIT KDC, see:

- [MIT Kerberos Home](#)
- [MIT Kerberos Documentation](#)
- **Support**
 - Cloudera supports the version of Kerberos that ships with each [supported operating system](#).

Considerations when using an Active Directory KDC

Performance:

As your cluster grows, so will the volume of Authentication Service (AS) and Ticket Granting Service (TGS) interaction between the services on each cluster server. Consider evaluating the volume of this interaction against the Active Directory domain controllers you have configured for the cluster before rolling this feature out to a production environment. If cluster performance suffers, over time it might become necessary to dedicate a set of AD domain controllers to larger deployments.

Network Proximity:

By default, Kerberos uses UDP for client/server communication. Often, AD services are in a different network than project application services such as Hadoop. If the domain controllers supporting a cluster for Kerberos are not in the same subnet, or they're separated by a firewall, consider using the `udp_preference_limit = 1` setting in the `[libdefaults]` section of the `krb5.conf` used by cluster services. Cloudera strongly recommends *against* using AD domain controller (KDC) servers that are separated from the cluster by a WAN connection, as latency in this service will significantly impact cluster performance.

Process:

Troubleshooting the cluster's operations, especially for Kerberos-enabled services, will need to include AD administration resources. Evaluate your organizational processes for engaging the AD administration team, and how to escalate in case a cluster outage occurs due to issues with Kerberos authentication against AD services. In some situations it might be necessary to [enable Kerberos event logging](#) to address desktop and KDC issues within windows environments.

Also note that if you decommission any Cloudera Manager roles or nodes, the related AD accounts will need to be deleted manually. This is required because Cloudera Manager will not delete existing entries in Active Directory.

Step 1: Install Cloudera Manager and CDH

If you have not already done so, Cloudera strongly recommends that you install and configure the Cloudera Manager Server and Cloudera Manager Agents and CDH to set up a fully-functional CDH cluster *before* you begin doing the following steps to implement Hadoop security features.

Overview of the User Accounts and Groups in CDH and Cloudera Manager to Support Security

User Accounts and Groups in CDH and Cloudera Manager Required to Support Security:

When you install the CDH packages and the Cloudera Manager Agents on your cluster hosts, Cloudera Manager takes some steps to provide system security such as creating the following Unix accounts and setting directory permissions as shown in the following table. These Unix accounts and directory permissions work with the Hadoop Kerberos security requirements.



Note: Cloudera Manager 5.3 introduces a new *single user mode*. In single user mode, the Cloudera Manager Agent and *all the processes run by services managed by Cloudera Manager* are started as a single configured user and group. See [Single User Mode Requirements](#) for more information.

This User	Runs These Roles
hdfs	NameNode, DataNodes, and Secondary Node
mapred	JobTracker and TaskTrackers (MR1) and Job History Server (YARN)
yarn	ResourceManager and NodeManagers (YARN)
oozie	Oozie Server
hue	Hue Server, Beeswax Server, Authorization Manager, and Job Designer

The `hdfs` user also acts as the HDFS superuser.

When you install the Cloudera Manager Server on the server host, a new Unix user account called `cloudera-scm` is created automatically to support security. The Cloudera Manager Server uses this account to create host principals and deploy the keytabs on your cluster.

Depending on whether you installed CDH and Cloudera Manager at the same time or not, use one of the following sections for information on configuring directory ownerships on cluster hosts:

If you installed CDH and Cloudera Manager at the Same Time

If you have a new installation and you installed CDH and Cloudera Manager at the same time, when you started the Cloudera Manager Agents on your cluster hosts, the Cloudera Manager Agent on each host automatically configured the directory owners shown in the following table to support security. Assuming the owners are configured as shown, the Hadoop daemons can then automatically set the permissions for each of the directories specified by the properties shown below to make sure they are properly restricted. It's critical that the owners are configured exactly as shown below, so do not change them:

Directory Specified in this Property	Owner
<code>dfs.name.dir</code>	<code>hdfs:hadoop</code>
<code>dfs.data.dir</code>	<code>hdfs:hadoop</code>
<code>mapred.local.dir</code>	<code>mapred:hadoop</code>
<code>mapred.system.dir</code> in HDFS	<code>mapred:hadoop</code>
<code>yarn.nodemanager.local-dirs</code>	<code>yarn:yarn</code>
<code>yarn.nodemanager.log-dirs</code>	<code>yarn:yarn</code>
<code>oozie.service.StoreService.jdbc.url</code> (if using Derby)	<code>oozie:oozie</code>
<code>[[database]] name</code>	<code>hue:hue</code>
<code>javax.jdo.option.ConnectionURL</code>	<code>hue:hue</code>

If you Installed and Used CDH Before Installing Cloudera Manager

If you have been using HDFS and running MapReduce jobs in an existing installation of CDH before you installed Cloudera Manager, you must manually configure the owners of the directories shown in the table above. Doing so enables the Hadoop daemons to automatically set the permissions for each of the directories. It's critical that you manually configure the owners exactly as shown above.

Step 2: If You are Using AES-256 Encryption, Install the JCE Policy File

If you are using CentOS or Red Hat Enterprise Linux 5.5 or higher, which use AES-256 encryption by default for tickets, you must install the [Java Cryptography Extension \(JCE\) Unlimited Strength Jurisdiction Policy File](#) on all cluster and Hadoop user hosts. There are 2 ways to do this:

- In the Cloudera Manager Admin Console, navigate to the **Hosts** page. Both, the **Add New Hosts to Cluster** wizard and the **Re-run Upgrade Wizard** will give you the option to have Cloudera Manager install the JCE Policy file for you.
- You can follow the JCE Policy File installation instructions in the `README.txt` file included in the `jce_policy-x.zip` file.

Alternatively, you can configure Kerberos to not use AES-256 by removing `aes256-cts:normal` from the `supported_encetypes` field of the `kdc.conf` or `krb5.conf` file. Note that after changing the `kdc.conf` file, you'll need to restart both the KDC and the `kadmin` server for those changes to take effect. You may also need to recreate or change the password of the relevant principals, including potentially the Ticket Granting Ticket principal (`krbtgt/REALM@REALM`). If AES-256 is still used after all of those steps, it's because the `aes256-cts:normal` setting existed when the Kerberos database was created. To fix this, create a new Kerberos database and then restart both the KDC and the `kadmin` server.

To verify the type of encryption used in your cluster:

1. **For MIT KDC:** On the local KDC host, type this command in the `kadmin.local` or `kadmin` shell to create a test principal:

```
kadmin: addprinc test
```

For Active Directory: Create a new AD account with the name, `test`.

2. On a cluster host, type this command to start a Kerberos session as test:

```
$ kinit test
```

3. On a cluster host, type this command to view the encryption type in use:

```
$ klist -e
```

If AES is being used, output like the following is displayed after you type the `klist` command (note that AES-256 is included in the output):

```
Ticket cache: FILE:/tmp/krb5cc_0
Default principal: test@Cloudera Manager
Valid starting Expires Service principal
05/19/11 13:25:04 05/20/11 13:25:04 krbtgt/Cloudera Manager@Cloudera Manager
Etype (skey, tkt): AES-256 CTS mode with 96-bit SHA-1 HMAC, AES-256 CTS mode with
96-bit SHA-1 HMAC
```

Step 3: Get or Create a Kerberos Principal for the Cloudera Manager Server

In order to create and deploy the host principals and keytabs on your cluster, the Cloudera Manager Server must have the correct Kerberos principal. Specifically, the Cloudera Manager Server must have a Kerberos principal that has privileges to create other accounts.

To get or create the Kerberos principal for the Cloudera Manager Server, you can do either of the following:

- Ask your Kerberos administrator to create a Kerberos administrator principal for the Cloudera Manager Server.
- Create the Kerberos principal for the Cloudera Manager Server yourself by using the following instructions in this step.

If for some reason, you cannot create a Cloudera Manager administrator principal on your KDC with the privileges to create other principals and keytabs for CDH services, then these will need to be created manually, and then retrieved by Cloudera Manager. See, [Using a Custom Kerberos Keytab Retrieval Script](#) on page 72.

Creating the Cloudera Manager Principal

The following instructions illustrate an example of creating the Cloudera Manager Server principal for MIT KDC and Active Directory KDC. (If you are using another version of Kerberos, refer to your Kerberos documentation for instructions.)

If you are using Active Directory:

1. Create an Organizational Unit (OU) in your AD setup where all the principals used by your CDH cluster will reside.
2. Add a new user account to Active Directory, for example, `<username>@YOUR-REALM.COM`. The password for this user should be set to never expire.
3. Use AD's Delegate Control wizard to allow this new user to **Create, Delete and Manage User Accounts**.

If you are using MIT KDC:

Typically, principals with the second component of `admin` in the principal name (for example, `username/admin@YOUR-LOCAL-REALM.com`) have administrator privileges. This is why `admin` is shown in the following example.



Note: If you are running `kadmin` and the Kerberos Key Distribution Center (KDC) on the same host, use `kadmin.local` in the following steps. If the Kerberos KDC is running on a remote host, you must use `kadmin` instead of `kadmin.local`.


In the `kadmin.local` or `kadmin` shell, type the following command to create the Cloudera Manager Server principal, replacing `YOUR-LOCAL-REALM.COM` with the name of your realm:

```
kadmin: addprinc -pw <Password> cloudera-scm/admin@YOUR-LOCAL-REALM.COM
```

Step 4: Enabling Kerberos Using the Wizard

Minimum Required Role: [Full Administrator](#)

To start the Kerberos wizard:

1. Go to the Cloudera Manager Admin Console and click  to the right of the cluster for which you want to enable Kerberos authentication.
2. Select **Enable Kerberos**.

Before you Begin Using the Wizard

The Welcome page lists the following action items that you should complete before you begin to secure the cluster using this wizard:

- Set up a working KDC. Cloudera Manager supports authentication with MIT KDC and Active Directory.
- Configure the KDC to allow renewable tickets with non-zero ticket lifetimes.

Active Directory allows renewable tickets with non-zero lifetimes by default. You can verify this by checking **Domain Security Settings > Account Policies > Kerberos Policy** in Active Directory.

For MIT KDC, make sure you have the following lines in the `kdc.conf`.

```
max_life = 1d
max_renewable_life = 7d
```

- If you are using Active Directory, make sure LDAP over TLS/SSL (LDAPS) is enabled for the Domain Controllers.
- Install the following packages on your cluster depending on the OS in use.

OS	Packages Required
RHEL 7 Compatible, RHEL 6 Compatible, RHEL 5 Compatible	<ul style="list-style-type: none"> • <code>openldap-clients</code> on the Cloudera Manager Server host • <code>krb5-workstation, krb5-libs</code> on ALL hosts
SLES	<ul style="list-style-type: none"> • <code>openldap2-client</code> on the Cloudera Manager Server host • <code>krb5-client</code> on ALL hosts
Ubuntu or Debian	<ul style="list-style-type: none"> • <code>ldap-utils</code> on the Cloudera Manager Server host • <code>krb5-user</code> on ALL hosts
Windows	<ul style="list-style-type: none"> • <code>krb5-workstation, krb5-libs</code> on ALL hosts

- Create an account for Cloudera Manager that has the permissions to create other accounts in the KDC. This should have been completed as part of [Step 3: Get or Create a Kerberos Principal for the Cloudera Manager Server](#) on page 62.



Important:

If you have enabled YARN Resource Manager HA in your non-secure cluster, you should clear the StateStore znode in ZooKeeper before enabling Kerberos. To do this:

1. Go to the Cloudera Manager Admin Console home page, click to the right of the YARN service and select **Stop**.
2. When you see a **Finished** status, the service has stopped.
3. Go to the YARN service and select **Actions > Format State Store**.
4. When the command completes, click **Close**.

Once you are able to check all the items on this list, click **Continue**.

KDC Information

On this page, select the KDC type you are using, MIT KDC or Active Directory, and complete the fields as applicable to enable Cloudera Manager to generate principals/accounts for the CDH services running on the cluster.



Note:

- If you are using AD and have multiple Domain Controllers behind a Load Balancer, enter the name of the Load Balancer in the **KDC Server Host** field and any *one* of the Domain Controllers in **Active Directory Domain Controller Override**. Hadoop daemons will use the Load Balancer for authentication, but Cloudera Manager will use the override for creating accounts.
- If you have multiple Domain Controllers (in case of AD) or MIT KDC servers, only enter the name of any *one* of them in the **KDC Server Host** field. Cloudera Manager will use that server only for creating accounts. If you choose to use Cloudera Manager to manage `krb5.conf`, you can specify the rest of the Domain Controllers using Safety Valve as explained below.
- Make sure the entries for the **Kerberos Encryption Types** field matches what your KDC supports.

Click **Continue** to proceed.

KRB5 Configuration

Manage `krb5.conf` through Cloudera Manager allows you to choose whether Cloudera Manager should deploy the `krb5.conf` on your cluster or not. If left unchecked, you must ensure that the `krb5.conf` is deployed on all hosts in the cluster, including the Cloudera Manager Server's host.

If you check **Manage krb5.conf through Cloudera Manager**, this page will let you configure the properties that will be emitted in it. In particular, the safety valves on this page can be used to configure cross-realm authentication. More information can be found at [Configuring a Cluster-dedicated MIT KDC with Cross-Realm Trust](#) on page 180.



Note: Cloudera Manager is unable to use a non-default realm. You must specify the default realm.

Click **Continue** to proceed.

Import KDC Account Manager Credentials

Enter the username and password for the user that can create principals for CDH cluster in the KDC. This is the user/principal you created in [Step 3: Get or Create a Kerberos Principal for the Cloudera Manager Server](#) on page 62. Cloudera Manager encrypts the username and password into a keytab and uses it as needed to create new principals.



Note: The username entered should have the realm portion in upper-case only as shown in the example in the UI.

Click **Continue** to proceed.

(Optional) Configuring Custom Kerberos Principals

Starting with Cloudera Manager 5.4, you can configure custom service principals for CDH services. Before you begin making configuration changes, see [Configuring a Cluster with Custom Kerberos Principals](#) on page 70 for some additional configuration changes required and limitations.

Configure HDFS DataNode Ports

On this page, specify the privileged ports needed by the DataNode's Transceiver Protocol and the HTTP Web UI in a secure cluster.

Use the checkbox to confirm you are ready to restart the cluster. Click **Continue**.

Enabling Kerberos

This page lets you track the progress made by the wizard as it first stops all services on your cluster, deploys the `krb5.conf`, generates keytabs for other CDH services, deploys client configuration and finally restarts all services. Click **Continue**.

Congratulations

The final page lists the cluster(s) for which Kerberos has been successfully enabled. Click **Finish** to return to the Cloudera Manager Admin Console home page.

Step 5: Create the HDFS Superuser

To be able to create home directories for users, you will need access to the HDFS superuser account. (CDH automatically created the HDFS superuser account on each cluster host during CDH installation.) When you enabled Kerberos for the HDFS service, you lost access to the default HDFS superuser account using `sudo -u hdfs` commands. Cloudera recommends you use a different user account as the superuser, not the default `hdfs` account.

Designating a Non-Default Superuser Group

To designate a different group of superusers instead of using the default `hdfs` account, follow these steps:

1. Go to the Cloudera Manager Admin Console and navigate to the HDFS service.
2. Click the **Configuration** tab.
3. Select **Scope > HDFS (Service-Wide)**.
4. Select **Category > Security**.
5. Locate the **Superuser Group** property and change the value to the appropriate group name for your environment. For example, `<superuser>`.

Configuring Authentication

6. Click **Save Changes** to commit the changes.
7. Restart the HDFS service.

To enable your access to the superuser account now that Kerberos is enabled, you must now create a Kerberos principal or an Active Directory user whose first component is *<superuser>*:

If you are using Active Directory

Add a new user account to Active Directory, *<superuser>@YOUR-REALM.COM*. The password for this account should be set to never expire.

If you are using MIT KDC

1. In the `kadmin.local` or `kadmin` shell, type the following command to create a Kerberos principal called *<superuser>*:

```
kadmin: addprinc <superuser>@YOUR-LOCAL-REALM.COM
```

This command prompts you to create a password for the *<superuser>* principal. You should use a strong password because having access to this principal provides superuser access to all of the files in HDFS.

2. To run commands as the HDFS superuser, you must obtain Kerberos credentials for the *<superuser>* principal. To do so, run the following command and provide the appropriate password when prompted.

```
$ kinit <superuser>@YOUR-LOCAL-REALM.COM
```

Step 6: Get or Create a Kerberos Principal for Each User Account

Now that Kerberos is configured and enabled on your cluster, you and every other Hadoop user must have a Kerberos principal or keytab to obtain Kerberos credentials to be allowed to access the cluster and use the Hadoop services. In the next step of this procedure, you will need to create your own Kerberos principals in order to verify that Kerberos security is working on your cluster. If you and the other Hadoop users already have a Kerberos principal or keytab, or if your Kerberos administrator can provide them, you can skip ahead to the next step.

The following instructions explain how to create a Kerberos principal for a user account.

If you are using Active Directory

Add a new AD user account for each new user that should have access to the cluster. You do not need to make any changes to existing user accounts.

If you are using MIT KDC

1. In the `kadmin.local` or `kadmin` shell, use the following command to create user principals by replacing `YOUR-LOCAL-REALM.COM` with the name of your realm, and replacing `USERNAME` with a username:

```
kadmin: addprinc USERNAME@YOUR-LOCAL-REALM.COM
```

2. Enter and re-enter a password when prompted.

Step 7: Prepare the Cluster for Each User

Before you and other users can access the cluster, there are a few tasks you must do to prepare the hosts for each user.

1. Make sure all hosts in the cluster have a Linux user account with the same name as the first component of that user's principal name. For example, the Linux account `joe` should exist on every box if the user's principal name is `joe@YOUR-REALM.COM`. You can use LDAP for this step if it is available in your organization.



Note: Each account must have a user ID that is greater than or equal to 1000. In the `/etc/hadoop/conf/taskcontroller.cfg` file, the default setting for the `banned.users` property is `mapred, hdfs, and bin` to prevent jobs from being submitted using those user accounts. The default setting for the `min.user.id` property is 1000 to prevent jobs from being submitted with a user ID less than 1000, which are conventionally Unix super users.

2. Create a subdirectory under `/user` on HDFS for each user account (for example, `/user/joe`). Change the owner and group of that directory to be the user.

```
$ hadoop fs -mkdir /user/joe
$ hadoop fs -chown joe /user/joe
```



Note: `sudo -u hdfs` is not included in the commands above. This is because it is not required if Kerberos is enabled on your cluster. You will, however, need to have Kerberos credentials for the HDFS super user in order to successfully run these commands. For information on gaining access to the HDFS super user account, see [Step 14: Create the HDFS Superuser Principal](#) on page 83

Step 8: Verify that Kerberos Security is Working

After you have Kerberos credentials, you can verify that Kerberos security is working on your cluster by trying to run MapReduce jobs. To confirm, try launching a sleep or a pi job from the provided Hadoop examples (`/usr/lib/hadoop/hadoop-examples.jar`).



Note:

This section assumes you have a fully-functional CDH cluster and you have been able to access HDFS and run MapReduce jobs before you followed these instructions to configure and enable Kerberos on your cluster. If you have not already done so, you should at a minimum use the Cloudera Manager Admin Console to generate a client configuration file to enable you to access the cluster. For instructions, see [Deploying Client Configuration Files](#).

To verify that Kerberos security is working:

1. Acquire Kerberos credentials for your user account.

```
$ kinit USERNAME@YOUR-LOCAL-REALM.COM
```

2. Enter a password when prompted.
3. Submit a sample pi calculation as a test MapReduce job. Use the following command if you use a package-based setup for Cloudera Manager:

```
$ hadoop jar /usr/lib/hadoop-0.20-mapreduce/hadoop-examples.jar pi 10 10000
Number of Maps = 10
Samples per Map = 10000
...
Job Finished in 38.572 seconds
Estimated value of Pi is 3.14120000000000000000
```

If you have a parcel-based setup, use the following command instead:

```
$ hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-0.20-mapreduce/hadoop-examples.jar pi
10 10000
Number of Maps = 10
Samples per Map = 10000
...
Job Finished in 30.958 seconds
Estimated value of Pi is 3.14120000000000000000
```

You have now verified that Kerberos security is working on your cluster.



Important:

Running a MapReduce job will fail if you do not have a valid Kerberos ticket in your credentials cache. You can examine the Kerberos tickets currently in your credentials cache by running the `klist` command. You can obtain a ticket by running the `kinit` command and either specifying a keytab file containing credentials, or entering the password for your principal. If you do not have a valid ticket, you will receive an error such as:

```
11/01/04 12:08:12 WARN ipc.Client:
Exception encountered while connecting to the server :
javax.security.sasl.SaslException:GSS initiate failed
[Caused by GSSException: No valid credentials provided (Mechanism level:
Failed to find any
Kerberos tgt)]
Bad connection to FS. command aborted. exception: Call to
nn-host/10.0.0.2:8020 failed on local exception:
java.io.IOException:javax.security.sasl.SaslException: GSS initiate
failed
[Caused by GSSException: No valid credentials provided
(Mechanism level: Failed to find any Kerberos tgt)]
```

Step 9: (Optional) Enable Authentication for HTTP Web Consoles for Hadoop Roles

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)

Authentication for access to the HDFS, MapReduce, and YARN roles' web consoles can be enabled using a configuration option for the appropriate service. To enable this authentication:

1. From the **Clusters** tab, select the service (HDFS, MapReduce, or YARN) for which you want to enable authentication.
2. Click the **Configuration** tab.
3. Select **Scope** > *service name* **Service-Wide**.
4. Select **Category** > **Security**.
5. Type `Enable Kerberos` in the Search box.
6. Select **Enable Kerberos Authentication for HTTP Web-Consoles**.
7. Click **Save Changes** to commit the changes.
8. When the command finishes, restart all roles of that service.

Enabling SPNEGO as an Authentication Backend for Hue

1. In Cloudera Manager, set the authentication backend to `SpnegoDjangoBackend`.
 - a. Go to the Cloudera Manager Admin Console. From the **Clusters** tab, select the Hue service.
 - b. Click the **Configuration** tab.
 - c. Select **Scope** > **Service-Wide**.
 - d. Select **Category** > **Security**.
 - e. Locate the **Authentication Backend** property and select `desktop.auth.backend.SpnegoDjangoBackend`.
 - f. Click **Save Changes**.
2. Restart the Hue service.
3. On the host running the Hue Kerberos Ticket Renewer, switch to the `KT_RENEWER` process directory. For example:

```
cd /var/run/cloudera-scm-agent/process/`ls -lrt /var/run/cloudera-scm-agent/process/
| awk '{print $9}' | grep KT_RENEWER | tail -1`/
```

4. Verify that the Hue keytab includes the HTTP principal.

```
klist -kte ./hue.keytab

Keytab name: FILE:./hue.keytab
KVNO Timestamp Principal
-----
1 03/09/15 20:20:35 hue/host-10-16-8-168.openstacklocal@EXAMPLE.CLOUDERA.COM
(aes128-cts-hmac-sha1-96)
1 03/09/15 20:20:36 HTTP/host-10-16-8-168.openstacklocal@EXAMPLE.CLOUDERA.COM
(aes128-cts-hmac-sha1-96)
```

5. Copy the hue.keytab file to /var/lib/hue and change ownership to the hue user and group.

```
$ cp ./hue.keytab /var/lib/hue/
$ chown hue:hue /var/lib/hue/hue.keytab
```

6. Go to the Cloudera Manager Admin Console. From the **Clusters** tab, select the Hue service.

7. Click the **Configuration** tab.

8. Select **Scope** > **Service-Wide**.

9. Select **Category** > **Advanced**.

10. Locate the **Hue Service Environment Advanced Configuration Snippet (Safety Valve)** property and add the following line:

```
KRB5_KTNAME=/var/lib/hue/hue.keytab
```

11. Click **Save Changes** to commit the changes.

12. Restart the Hue service.

Enabling Kerberos Authentication for Single User Mode or Non-Default Users

The steps described in this topic are only applicable in the following cases:

- You are running the Cloudera Manager in the [single user mode](#). In this case, configure all the services described in the table below.

OR

- You are running one or more CDH services with non-default users. This means if you have modified the default value for the **System User** property for any service in Cloudera Manager, you must only perform the command (as described below) corresponding to that service, to be able to successfully run jobs with the non-default user.

MapReduce	<p>Configure the <code>mapred.system.dir</code> directory to be owned by the <code>mapred</code> user.</p> <pre>sudo -u hdfs hadoop fs -chown mapred:hadoop \${mapred.system.dir}</pre> <p>By default, <code>mapred.system.dir</code> is <code>/tmp/mapred/system</code>.</p>
HBase	<p>Give the <code>hbase</code> user ownership of the HBase root directory:</p> <pre>sudo -u hdfs hadoop fs -chown -R hbase \${hbase.rootdir}</pre> <p>By default, <code>hbase.rootdir</code> is <code>/hbase</code>.</p>
Hive	<p>Give the <code>hive</code> user ownership of the <code>/user/hive</code> directory.</p> <pre>sudo -u hdfs hadoop fs -chown hive /user/hive</pre>

Configuring Authentication

YARN	<p>For every NodeManager host, for each path in <code>yarn.nodemanager.local-dirs</code>, run:</p> <pre>rm -rf \${yarn.nodemanager.local-dirs}/usercache/*</pre> <p>This removes the <code>/usercache</code> directory that contains intermediate data stored for previous jobs.</p>
------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Configuring a Cluster with Custom Kerberos Principals

By default, the Cloudera Manager [Kerberos wizard](#) configures CDH services to use the same Kerberos principals as the default process users. For example, the `hdfs` principal for the HDFS service, `hive` principal for the Hive service, and so on. The advantage to this is that when Kerberos is enabled, no HDFS directory permissions need to be changed for the new principals. However, starting with Cloudera Manager 5.4, you can configure custom service principals for CDH services.

Important Considerations

- Using different Kerberos principals for different services will make it easier to track the HDFS directories being accessed by each service.
- If you are using `ShellBasedUnixGroupsMapping` to obtain user-group mappings, ensure you have the UNIX accounts for the principals present on all hosts of the cluster.

Configuring Directory Permissions

Configure the following HDFS directories to give their corresponding custom service principals `read`, `write` and `execute` permissions.

Service	HDFS Directory
Accumulo	<ul style="list-style-type: none">• HDFS Directory• <code>/user/principal</code>
HBase	HBase Root Directory
Hive	<ul style="list-style-type: none">• Hive Warehouse Directory• <code>/user/principal</code>
Impala	<code>/user/principal</code>
MapReduce v1	<code>/tmp/mapred</code>
Oozie	Oozie ShareLib Root Directory
Solr	HDFS Data Directory
Spark on YARN	<ul style="list-style-type: none">• <code>/user/principal</code>• Spark History Location• Spark Jar Location
Sqoop2	<code>/user/principal</code>

Configuring CDH Services

The following services will require additional settings if you are using custom principals:

- **HDFS** - If you have enabled synchronization of HDFS and Sentry permissions, add the Hive and Impala principals to the **Sentry Authorization Provider Group** property.
 1. Go to the HDFS service.
 2. Click Configuration.
 3. Select **Scope > HDFS Service-Wide**.

4. Select **Category > Security**.
 5. Locate the **Sentry Authorization Provider Group** property and add the custom Hive and Impala principals.
 6. Click **Save Changes**.
- **YARN** - The principals used by YARN daemons should be part of `hadoop` group so that they are allowed to read JobHistory Server data.
 - **Impala** - If you are running the Hue service with a custom principal, configure Impala to allow the Hue principal to impersonate other users.
 1. Go to the Impala service.
 2. Click Configuration.
 3. Select **Scope > Impala (Service-Wide)**.
 4. Select **Category > Policy File-Based Sentry**.
 5. Locate the **Proxy User Configuration** property and add the custom Hue principal.
 6. Click **Save Changes**.
 - **Hive** - If the Sentry service is enabled, allow the Kerberos principals used by Hive, Impala, Hue, HDFS and the Service Monitor to bypass Sentry authorization in the Hive metastore.
 1. Go to the Hive service.
 2. Click Configuration.
 3. Select **Scope > Impala (Service-Wide)**.
 4. Select **Category > Policy File-Based Sentry**.
 5. Locate the **Bypass Sentry Authorization Users** property and add the custom Hive, Impala, Hue and HDFS principals to the list.
 6. Click **Save Changes**.
 - **Spark on YARN** - The principal used by the Spark service should be part of the `spark` group.
 - **Sentry** - Allow the Hive, Impala, Hue and HDFS principals to connect to the Sentry service.
 1. Go to the Sentry service.
 2. Click Configuration.
 3. Search for the **Allowed Connecting Users** property and add the custom Hive, Impala, Hue and HDFS principals to the list.
 4. Search for the **Admin Groups** property and include the groups to which the Hive, Impala, and Hue principals belong.
 5. Click **Save Changes**.
 - **Cloudera Management Service** - Configure the Reports Manager principal and the Navigator principal for HDFS as HDFS superusers.
 1. Go to the Cloudera Management Service.
 2. Click Configuration.
 3. Search for *kerberos*.
 4. Locate the **Reports Manager Kerberos Principal** property and set it to a principal with administrative and superuser privileges on all HDFS services.
 5. Locate the **Navigator Kerberos Principal for HDFS** property and set it to a principal with administrative and superuser privileges on all HDFS services.
 6. Click **Save Changes**.

Incompatibilities

The following features do not work with custom principals:

- Llama must always use the default Kerberos principal `llama`.
- If you are using MapReduce v1, the Activity Monitor and Cloudera Navigator should use the same principal as the Hue service.

Configuring Authentication

- If you are using the Java KeyStore KMS or KeyTrustee KMS with a custom principal, you will need to add the proxy user for the custom principal to the `kms-site.xml` safety valve.

For example, if you've replaced the default `oozie` principal with `oozieprinc`, add the `hadoop.kms.proxyuser.oozieprinc.groups` and `hadoop.kms.proxyuser.oozieprinc.hosts` properties to the `kms-site.xml` safety valve.

Viewing and Regenerating Kerberos Principals

Minimum Required Role: [Full Administrator](#)

As soon as you enable Hadoop secure authentication for HDFS and MapReduce service instances, Cloudera Manager starts creating the Kerberos principals for each of the role instances. The amount of time this process will take depends on the number of hosts and HDFS and MapReduce role instances on your cluster. The process can take from a few seconds for a small cluster to several minutes for a larger cluster. After the process is completed, you can use the Cloudera Manager Admin Console to view the list of Kerberos principals that Cloudera Manager has created for the cluster. Make sure there are principals for each of the hosts and HDFS and MapReduce role instances on your cluster. If there are no principals after 10 minutes, then there is most likely a problem with the principal creation. See the [Troubleshooting Authentication Issues](#) on page 194 section below for more information. If necessary, you can use Cloudera Manager to regenerate the principals.

If you make a global configuration change in your cluster, such as changing the encryption type, you must use the following instructions to regenerate the principals for your cluster.



Important:

- Regenerate principals using the following steps in the Cloudera Manager Admin Console and not directly using `kadmin` shell.
- Do not regenerate the principals for your cluster unless you have made a global configuration change. Before regenerating, be sure to read [Configuring a Cluster-dedicated MIT KDC with Cross-Realm Trust](#) on page 180 to avoid making your existing host keytabs invalid.
- If you are using Active Directory, delete the AD accounts with the `userPrincipalName` (or login names) that you want to manually regenerate before continuing with the steps below. This is required because Cloudera Manager will not delete existing entries in Active Directory.

To view and regenerate the Kerberos principals for your cluster:

1. Select **Administration > Security**.
2. The currently configured Kerberos principals are displayed under the **Kerberos Credentials** tab. If you are running HDFS, the `hdfs/hostname` and `host/hostname` principals are listed. If you are running MapReduce, the `mapred/hostname` and `host/hostname` principals are listed. The principals for other running services are also listed.
3. Only if necessary, select the principals you want to regenerate.
4. Click **Regenerate**.

The Security Inspector

The Security Inspector uses the Host Inspector to run a security-related set of commands on the hosts in your cluster. It reports on things such as how Java is configured for encryption and on the default realms configured on each host:

1. Select **Administration > Security**.
2. Click **Security Inspector**. Cloudera Manager begins several tasks to inspect the managed hosts.
3. After the inspection completes, click **Download Result Data** or **Show Inspector Results** to review the results.

Using a Custom Kerberos Keytab Retrieval Script

The Cloudera Manager [Kerberos setup procedure](#) requires you to create an administrator account for the Cloudera Manager user. Cloudera Manager then connects to your KDC and uses this admin account to generate principals and

keytabs for the remaining CDH services. If for some reason, you cannot create a Cloudera Manager administrator account on your KDC with the privileges to create other principals and keytabs for CDH services, then these will need to be created manually.

Cloudera Manager gives you the option to use a custom script to retrieve keytabs from the local filesystem. To use a custom Kerberos keytab retrieval script:

1. The KDC administrators should create the required principals and keytabs, and store them securely on the Cloudera Manager Server host.
2. Create the keytab retrieval script. Your script should take two arguments: a full principal name for which it should retrieve a keytab, and a destination to which it can write the keytab. The script must be executable by the Cloudera Manager admin user, `cloudera-scm`. Depending on the principal name input by Cloudera Manager, the script should locate the corresponding keytab on the Cloudera Manager Server host (stored in step 1), and copy it into a location accessible to the `cloudera-scm` user. Here is a simple example:

```
#!/bin/bash

# Cloudera Manager will input a destination path
DEST="$1"

# Cloudera Manager will input the principal name in the format: <service>/<fqdn>@REALM
PRINC="$2"

# Assuming the '<service>_<fqdn>@REALM.keytab' naming convention for keytab files
IN=$(echo $PRINC | sed -e 's/\//_/')
SRC="/keytabs/${IN}.keytab"

# Copy the keytab to the destination input by Cloudera Manager
cp -v $SRC $DEST
```

Note that the script will change according to the keytab naming convention followed by your organization.

3. Configure the location for the script in Cloudera Manager:
 1. Go to the Cloudera Manager Admin console.
 2. Select **Administration > Settings**.
 3. Select **Category > Kerberos**.
 4. Locate the **Custom Kerberos Keytab Retrieval Script** and set it to point to the script created in step 2.
 5. Click **Save Changes** to commit the changes.
4. Once the **Custom Kerberos Keytab Retrieval Script** property is set, whenever Cloudera Manager needs a keytab, it will ignore all other Kerberos configuration and run the keytab retrieval script to copy the required keytab to the desired destination.
5. Cloudera Manager can now distribute the keytab to the services that need access to it.



Note: The Cloudera Navigator web server accesses HDFS and Hue using the keytabs corresponding to those principals; however the custom script does not move these additional keytabs to the Navigator Metadata Server. To complete the setup for Navigator, move keytabs for HDFS and Hue principals to the Navigator home directory on the Navigator Metadata Server host manually (typically `/var/lib/cloudera-scm-navigator`).

Mapping Kerberos Principals to Short Names

Kerberos user principals typically have the format `username@REALM`, whereas Hadoop usernames are typically just `username`. To translate Kerberos principals to Hadoop usernames, Hadoop uses rules defined in the `hadoop.security.auth_to_local` property. The default setting strips the `@REALM` portion from the Kerberos principal, where `REALM` is the Kerberos realm defined by the `default_realm` setting in the NameNode `krb5.conf` file.

If you configure your cluster's Kerberos realm to trust other realms, such as a trust between your cluster's realm and a central Active Directory or MIT Kerberos realm, you must identify the trusted realms in Cloudera Manager so it can

automatically generate the appropriate rules. If you do not do so, user accounts in those realms cannot access the cluster.

To specify trusted realms using Cloudera Manager:

1. Go to the **HDFS Service > Configuration** tab.
2. Select **Scope > HDFS (Service-Wide)**.
3. Select **Category > Security**.
4. In the Search field, type `Kerberos Realms` to find the **Trusted Kerberos Realms** and **Additional Rules to Map Kerberos Principals to Short Names** settings.
5. Add realms that are trusted by the cluster's Kerberos realm. Realm names, including Active Directory realms, must be specified in uppercase letters (for example, `CORP.EXAMPLE.COM`). To add multiple realms, use the **+** button.
6. Click **Save Changes**.

The auto-generated mapping rules strip the Kerberos realm (for example, `@CORP.EXAMPLE.COM`) for each realm specified in the **Trusted Kerberos Realms** setting. To customize the mapping rules, specify additional rules in the **Additional Rules to Map Kerberos Principals to Short Names** setting, one rule per line. Only enter rules in this field; Cloudera Manager automatically surrounds the rules with the appropriate XML tags for the generated `core-site.xml` file. For more information on creating custom rules, including how to translate mixed-case Kerberos principals to lowercase Hadoop usernames, see [Mapping Rule Syntax](#) on page 117.

If you specify custom mapping rules for a Kerberos realm using the **Additional Rules to Map Kerberos Principals to Short Names** setting, ensure that the same realm is not specified in the **Trusted Kerberos Realms** setting. If it is, the auto-generated rule (which only strips the realm from the principal and does no additional transformations) takes precedent, and the custom rule is ignored.

For these changes to take effect, you must restart the cluster and redeploy the client configuration. On the Cloudera Manager **Home > Status** tab, click the cluster-wide button and select **Deploy Client Configuration**.

Moving Kerberos Principals to Another OU Within Active Directory

If you have a Kerberized cluster configured with an Active Directory KDC, you can use the following steps to move the Kerberos principals from one AD Organizational Unit (OU) to another.

1. Create the new OU on the Active Directory Server.
2. Use AD's Delegate Control wizard to set the permissions on the new OU such that the configured Cloudera Manager admin account has the ability to **Create, Delete and Manage User Accounts** within this OU.
3. [Stop the cluster](#).
4. [Stop the Cloudera Management Service](#).
5. In Active Directory, move all the Cloudera Manager and CDH components' user accounts to the new OU.
6. Go to Cloudera Manager and go to **Administration > Security**.
7. Go to the **Kerberos Credentials** tab and click **Configuration**.
8. Select **Scope > Settings**.
9. Select **Category > Kerberos**.
10. Locate the **Active Directory Suffix** property and edit the value to reflect the new OU name.
11. Click **Save Changes** to commit the changes.

Using Auth-to-Local Rules to Isolate Cluster Users

By default, the Hadoop auth-to-local rules map a principal of the form `<username>/<hostname>@<REALM>` to `<username>`. This means if there are multiple clusters in the same realm, then principals associated with hosts of one cluster would map to the same user in all other clusters.

For example, if you have two clusters, `cluster1-host-[1..4].example.com` and `cluster2-host-[1..4].example.com`, that are part of the same Kerberos realm, `EXAMPLE.COM`, then the `cluster2` principal, `hdfs/cluster2-host1.example.com@EXAMPLE.COM`, will map to the `hdfs` user even on `cluster1` hosts.

To prevent this, use auth-to-local rules as follows to ensure only principals containing hostnames of `cluster1` are mapped to legitimate users.

1. Go to the **HDFS Service > Configuration** tab.
2. Select **Scope > HDFS (Service-Wide)**.
3. Select **Category > Security**.
4. In the Search field, type `Additional Rules` to find the **Additional Rules to Map Kerberos Principals to Short Names** settings.
5. Additional mapping rules can be added to the **Additional Rules to Map Kerberos Principals to Short Names** property. These rules will be inserted before the rules generated from the list of trusted realms (configured above) and before the default rule.

```
RULE:[2:$1/$2@$0](hdfs/cluster1-host1.example.com@EXAMPLE.COM)s/(.*)@EXAMPLE.COM/hdfs/
RULE:[2:$1/$2@$0](hdfs/cluster1-host2.example.com@EXAMPLE.COM)s/(.*)@EXAMPLE.COM/hdfs/
RULE:[2:$1/$2@$0](hdfs/cluster1-host3.example.com@EXAMPLE.COM)s/(.*)@EXAMPLE.COM/hdfs/
RULE:[2:$1/$2@$0](hdfs/cluster1-host4.example.com@EXAMPLE.COM)s/(.*)@EXAMPLE.COM/hdfs/
RULE:[2:$1/$2@$0](hdfs.*@EXAMPLE.COM)s/(.*)@EXAMPLE.COM/nobody/
```

In the example, the principal `hdfs/<hostname>@REALM` is mapped to the `hdfs` user if `<hostname>` is one of the cluster hosts. Otherwise it gets mapped to `nobody`, thus ensuring that principals from other clusters do not have access to `cluster1`.

If the cluster hosts can be represented with a regular expression, that expression can be used to make the configuration easier and more conducive to scaling. For example:

```
RULE:[2:$1/$2@$0](hdfs/cluster1-host[1-4].example.com@EXAMPLE.COM)s/(.*)@EXAMPLE.COM/hdfs/
RULE:[2:$1/$2@$0](hdfs.*@EXAMPLE.COM)s/(.*)@EXAMPLE.COM/nobody/
```

6. Click **Save Changes**.
7. Restart the HDFS service and any dependent services.

Enabling Kerberos Authentication Without the Wizard

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator, Full Administrator**)

Note that certain steps in the following procedure to configure Kerberos security may not be completed without `Full Administrator` role privileges.



Important: Ensure you have secured communication between the Cloudera Manager Server and Agents before you enable Kerberos on your cluster. Kerberos keytabs are sent from the Cloudera Manager Server to the Agents, and must be encrypted to prevent potential misuse of leaked keytabs. For instructions on securing this transfer with TLS encryption, see [How to Configure TLS Encryption for Cloudera Manager](#) on page 40.

- **Prerequisites** - These instructions assume you know how to install and configure Kerberos, you already have a working Kerberos key distribution center (KDC) and realm setup, and that you've installed the following Kerberos client packages on all cluster hosts and hosts that will be used to access the cluster, depending on the OS in use.

OS	Packages Required
RHEL 7 Compatible, RHEL 6 Compatible, RHEL 5 Compatible	<ul style="list-style-type: none"> • <code>openldap-clients</code> on the Cloudera Manager Server host • <code>krb5-workstation, krb5-libs</code> on ALL hosts
SLES	<ul style="list-style-type: none"> • <code>openldap2-client</code> on the Cloudera Manager Server host • <code>krb5-client</code> on ALL hosts
Ubuntu or Debian	<ul style="list-style-type: none"> • <code>ldap-utils</code> on the Cloudera Manager Server host • <code>krb5-user</code> on ALL hosts
Windows	<ul style="list-style-type: none"> • <code>krb5-workstation, krb5-libs</code> on ALL hosts

Furthermore, Oozie and Hue require that the realm support renewable tickets. Cloudera Manager supports setting up kerberized clusters with MIT KDC and Active Directory.



Important: If you want to integrate Kerberos directly with Active Directory, ensure you have support from your AD administration team to do so. This includes any future support required to troubleshoot issues such as Kerberos TGT/TGS ticket renewal, access to KDC logs for debugging and so on.

For more information about using Active Directory, refer the section below on **Considerations when using an Active Directory KDC** and the [Microsoft AD documentation](#).

For more information about installing and configuring MIT KDC, see:

- [MIT Kerberos Home](#)
- [MIT Kerberos Documentation](#)
- **Support**
 - Cloudera supports the version of Kerberos that ships with each [supported operating system](#).

Here are the general steps to using Cloudera Manager to configure Hadoop security on your cluster, each of which is described in more detail in the following sections:

Step 1: Install Cloudera Manager and CDH

If you have not already done so, Cloudera strongly recommends that you install and configure the Cloudera Manager Server and Cloudera Manager Agents and CDH to set up a fully-functional CDH cluster *before* you begin performing the steps to implement Kerberos authentication.

User Accounts and Groups in CDH and Cloudera Manager Required to Support Security:

When you install the CDH packages and the Cloudera Manager Agents on your cluster hosts, Cloudera Manager takes some steps to provide system security such as creating the following Unix accounts and setting directory permissions as shown in the following table. These Unix accounts and directory permissions work with the Hadoop Kerberos security requirements.



Note: Cloudera Manager 5.3 introduces a new *single user mode*. In single user mode, the Cloudera Manager Agent and *all the processes run by services managed by Cloudera Manager* are started as a single configured user and group. See [Single User Mode Requirements](#) for more information.

This User	Runs These Roles
hdfs	NameNode, DataNodes, and Secondary Node
mapred	JobTracker and TaskTrackers (MR1) and Job History Server (YARN)
yarn	ResourceManager and NodeManagers (YARN)
oozie	Oozie Server
hue	Hue Server, Beeswax Server, Authorization Manager, and Job Designer

The `hdfs` user also acts as the HDFS superuser.

When you install the Cloudera Manager Server on the server host, a new Unix user account called `cloudera-scm` is created automatically to support security. The Cloudera Manager Server uses this account to create and deploy the host principals and keytabs on your cluster.

Depending on whether you installed CDH and Cloudera Manager at the same time or not, use one of the following sections for information on configuring directory ownerships on cluster hosts:

If you installed CDH and Cloudera Manager at the Same Time

If you have a new installation and you installed CDH and Cloudera Manager at the same time, when you started the Cloudera Manager Agents on your cluster hosts, the Cloudera Manager Agent on each host automatically configured the directory owners shown in the following table to support security. Assuming the owners are configured as shown, the Hadoop daemons can then automatically set the permissions for each of the directories specified by the properties shown below to make sure they are properly restricted. It's critical that the owners are configured exactly as shown below, so do not change them:

Directory Specified in this Property	Owner
<code>dfs.name.dir</code>	<code>hdfs:hadoop</code>
<code>dfs.data.dir</code>	<code>hdfs:hadoop</code>
<code>mapred.local.dir</code>	<code>mapred:hadoop</code>
<code>mapred.system.dir</code> in HDFS	<code>mapred:hadoop</code>
<code>yarn.nodemanager.local-dirs</code>	<code>yarn:yarn</code>
<code>yarn.nodemanager.log-dirs</code>	<code>yarn:yarn</code>
<code>oozie.service.StoreService.jdbc.url</code> (if using Derby)	<code>oozie:oozie</code>
<code>[[database]] name</code>	<code>hue:hue</code>
<code>javax.jdo.option.ConnectionURL</code>	<code>hue:hue</code>

If you Installed and Used CDH Before Installing Cloudera Manager

If you have been using HDFS and running MapReduce jobs in an existing installation of CDH before you installed Cloudera Manager, you must manually configure the owners of the directories shown in the table above. Doing so enables the Hadoop daemons to automatically set the permissions for each of the directories. It's critical that you manually configure the owners exactly as shown above.

Step 2: If You are Using AES-256 Encryption, Install the JCE Policy File

If you are using CentOS or RHEL 5.5 or higher, which use AES-256 encryption by default for tickets, you must install the [Java Cryptography Extension \(JCE\) Unlimited Strength Jurisdiction Policy File](#) on all cluster and Hadoop user hosts. There are 2 ways to do this:

- In the Cloudera Manager Admin Console, navigate to the **Hosts** page. Both, the **Add New Hosts to Cluster** wizard and the **Re-run Upgrade Wizard** will give you the option to have Cloudera Manager install the JCE Policy file for you.
- You can follow the JCE Policy File installation instructions in the `README.txt` file included in the `jce_policy-x.zip` file.

Alternatively, you can configure Kerberos to not use AES-256 by removing `aes256-cts:normal` from the `supported_enctypes` field of the `kdc.conf` or `krb5.conf` file. Note that after changing the `kdc.conf` file, you'll need to restart both the KDC and the `kadmin` server for those changes to take affect. You may also need to recreate or change the password of the relevant principals, including potentially the Ticket Granting Ticket principal (for example, `krbtgt/EXAMPLE.COM@EXAMPLE.COM`). If AES-256 is still used after all of those steps, it's because the `aes256-cts:normal` setting existed when the Kerberos database was created. To fix this, create a new Kerberos database and then restart both the KDC and the `kadmin` server.

To verify the type of encryption used in your cluster:

1. On the local KDC host, type this command in the `kadmin.local` or `kadmin` shell to create a test principal:

```
kadmin: addprinc test
```

Configuring Authentication

2. On a cluster host, type this command to start a Kerberos session as the test principal:

```
$ kinit test
```

3. After successfully running the previous command, type this command to view the encryption type in use:

```
$ klist -e
```

If AES is being used, output like the following is displayed after you type the `klist` command (note that AES-256 is included in the output):

```
Ticket cache: FILE:/tmp/krb5cc_0
Default principal: test@EXAMPLE.COM
Valid starting Expires Service principal
05/19/11 13:25:04 05/20/11 13:25:04 krbtgt/EXAMPLE.COM@EXAMPLE.COM
Etype (skey, tkt): AES-256 CTS mode with 96-bit SHA-1 HMAC, AES-256 CTS mode with
96-bit SHA-1 HMAC
```

Step 3: Get or Create a Kerberos Principal for the Cloudera Manager Server

In order to create and deploy the host principals and keytabs on your cluster, the Cloudera Manager Server must have the correct Kerberos principal. Specifically, the Cloudera Manager Server must have a Kerberos principal that has administrator privileges. Typically, principals with the second component of `admin` in the principal name (for example, `username/admin@EXAMPLE.COM`) have administrator privileges. This is why `admin` is shown in the following instructions and examples.

To get or create the Kerberos principal for the Cloudera Manager Server, you can do either of the following:

- Ask your Kerberos administrator to create a Kerberos administrator principal for the Cloudera Manager Server.
- Create the Kerberos principal for the Cloudera Manager Server yourself by using the following instructions in this step.

Creating the Cloudera Manager Principal

If you are using Active Directory

1. Create an Organizational Unit (OU) in your AD where all the principals used by your CDH cluster will reside.
2. Add a new AD user, for example, `<username>@EXAMPLE.COM`. The password for this user should be set to never expire.
3. Use AD's Delegate Control wizard to allow this new user to **Create, Delete and Manage User Accounts**.

If you are using MIT KDC

The instructions in this section illustrate an example of creating the Cloudera Manager Server principal for MIT Kerberos. (If you are using another version of Kerberos, refer to your Kerberos documentation for instructions.)



Note: If you are running `kadmin` and the Kerberos Key Distribution Center (KDC) on the same host, use `kadmin.local` in the following steps. If the Kerberos KDC is running on a remote host, you must use `kadmin` instead of `kadmin.local`.

In the `kadmin.local` or `kadmin` shell, type the following command to create the Cloudera Manager Server principal, replacing `EXAMPLE.COM` with the name of your realm:

```
kadmin: addprinc -pw <Password> cloudera-scm/admin@EXAMPLE.COM
```

Step 4: Import KDC Account Manager Credentials

1. In the Cloudera Manager Admin Console, select **Administration > Security**.
2. Go to the **Kerberos Credentials** tab and click **Import Kerberos Account Manager Credentials**.

3. In the **Import Kerberos Account Manager Credentials** dialog box, enter the username and password for the user that can create principals for CDH cluster in the KDC. This is the user/principal you created in [Step 3: Get or Create a Kerberos Principal for the Cloudera Manager Server](#) on page 78. Cloudera Manager encrypts the username and password into a keytab and uses it as needed to create new principals.



Note: The username entered should have the realm portion in upper-case only as shown in the example in the UI.

Click **Close** when complete.

Step 5: Configure the Kerberos Default Realm in the Cloudera Manager Admin Console

Minimum Required Role: [Full Administrator](#)



Important: Hadoop is unable to use a non-default realm. The Kerberos default realm is configured in the `libdefaults` property in the `/etc/krb5.conf` file on every host in the cluster:

```
[libdefaults]
  default_realm = EXAMPLE.COM
```

1. In the Cloudera Manager Admin Console, select **Administration > Settings**.
2. Click the **Security** category, and enter the Kerberos realm for the cluster in the **Kerberos Security Realm** field (for example, `EXAMPLE.COM` or `HADOOP.EXAMPLE.COM`) that you configured in the `krb5.conf` file.
3. Click **Save Changes**.

Step 6: Stop All Services

Minimum Required Role: [Operator](#) (also provided by **Configurator**, **Cluster Administrator**, **Full Administrator**)

Before you enable security in CDH, you must stop all Hadoop daemons in your cluster and then change some configuration properties. You must stop all daemons in the cluster because after one Hadoop daemon has been restarted with the configuration properties set to enable security. Daemons running without security enabled will be unable to communicate with that daemon. This requirement to stop all daemons makes it impossible to do a rolling upgrade to enable security on a Hadoop cluster.

Stop all running services, and the Cloudera Management service, as follows:

Stopping All Services

1. On the **Home > Status** tab, click



to the right of the cluster name and select **Stop**.

2. Click **Stop** in the confirmation screen. The **Command Details** window shows the progress of stopping services.

When **All services successfully stopped** appears, the task is complete and you can close the **Command Details** window.

Stopping the Cloudera Management Service

1. On the **Home > Status** tab, click



to the right of **Cloudera Management Service** and select **Stop**.

2. Click **Stop** to confirm. The **Command Details** window shows the progress of stopping the roles.
3. When **Command completed with n/n successful subcommands** appears, the task is complete. Click **Close**.

Step 7: Enable Hadoop Security

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)

To enable Hadoop security for the cluster, you enable it on an HDFS service. After you do so, the Cloudera Manager Server automatically enables Hadoop security on the MapReduce and YARN services associated with that HDFS service.

1. Go to the **HDFS Service > Configuration** tab.
2. In the Search field, type **Hadoop Secure** to show the Hadoop security properties (found under the **Service-Wide > Security** category).
3. Click the value for the **Hadoop Secure Authentication** property and select the `kerberos` option to enable Hadoop security on the selected HDFS service.
4. Click the value for the **Hadoop Secure Authorization** property and select the checkbox to enable service-level authorization on the selected HDFS service. You can specify comma-separated lists of users and groups authorized to use Hadoop services or perform admin operations using the following properties under the **Service-Wide > Security** section:
 - **Authorized Users:** Comma-separated list of users authorized to use Hadoop services.
 - **Authorized Groups:** Comma-separated list of groups authorized to use Hadoop services.
 - **Authorized Admin Users:** Comma-separated list of users authorized to perform admin operations on Hadoop.
 - **Authorized Admin Groups:** Comma-separated list of groups authorized to perform admin operations on Hadoop.



Important: For Cloudera Manager's Monitoring services to work, the `hue` user should always be added as an authorized user.

5. In the Search field, type **DataNode Transceiver** to find the **DataNode Transceiver Port** property.
6. Click the value for the **DataNode Transceiver Port** property and specify a privileged port number (below 1024). Cloudera recommends 1004.



Note: If there is more than one DataNode Role Group, you must specify a privileged port number for each DataNode Transceiver Port property.

7. In the Search field, type **DataNode HTTP** to find the **DataNode HTTP Web UI Port** property and specify a privileged port number (below 1024). Cloudera recommends 1006.



Note: These port numbers for the two DataNode properties must be below 1024 in order to provide part of the security mechanism to make it impossible for a user to run a MapReduce task that impersonates a DataNode. The port numbers for the NameNode and Secondary NameNode can be anything you want, but the default port numbers are good ones to use.

8. In the Search field type **Data Directory Permissions** to find the **DataNode Data Directory Permissions** property.
9. Reset the value for the **DataNode Data Directory Permissions** property to the default value of 700 if not already set to that.
10. Make sure you have changed the **DataNode Transceiver Port**, **DataNode Data Directory Permissions** and **DataNode HTTP Web UI Port** properties for every DataNode role group.
11. Click **Save Changes** to save the configuration settings.

To enable ZooKeeper security:

1. Go to the **ZooKeeper Service > Configuration** tab and click **View and Edit**.
2. Click the value for **Enable Kerberos Authentication** property.
3. Click **Save Changes** to save the configuration settings.

To enable HBase security:

1. Go to the **HBase Service > Configuration** tab and click **View and Edit**.

2. In the Search field, type **HBase Secure** to show the Hadoop security properties (found under the **Service-Wide > Security** category).
3. Click the value for the **HBase Secure Authorization** property and select the checkbox to enable authorization on the selected HBase service.
4. Click the value for the **HBase Secure Authentication** property and select `kerberos` to enable authorization on the selected HBase service.
5. Click **Save Changes** to save the configuration settings.

(CDH 4.3 or later) To enable Solr security:

1. Go to the **Solr Service > Configuration** tab and click **View and Edit**.
2. In the Search field, type **Solr Secure** to show the Solr security properties (found under the **Service-Wide > Security** category).
3. Click the value for the **Solr Secure Authentication** property and select `kerberos` to enable authorization on the selected Solr service.
4. Click **Save Changes** to save the configuration settings.



Note: If you use the Cloudera Manager Admin Console to generate a client configuration file after you enable Hadoop security on your cluster, the generated configuration file will not contain the Kerberos principal and keytab file that end users need to authenticate. Users must obtain Kerberos principal and keytab file from your Kerberos administrator and then run the `kinit` command themselves.

Step 8: Wait for the Generate Credentials Command to Finish

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator, Full Administrator**)

After you enable security for any of the services in Cloudera Manager, a command called Generate Credentials will be triggered automatically. You can watch the progress of the command on the top right corner of the screen that shows the running commands. Wait for this command to finish (indicated by a grey box containing "0" in it).

Step 9: Enable Hue to Work with Hadoop Security using Cloudera Manager

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

If you are using a Hue service, you must add a role instance of Kerberos Ticket Renewer to the Hue service to enable Hue to work properly with the secure Hadoop cluster using Cloudera Manager. The Kerberos Ticket Renewer role must be located on the same host as the Hue Server role. You add can the necessary Kerberos Ticket Renewer role instances using Cloudera Manager.

The Hue Kerberos Ticket Renewer service will only renew tickets for the Hue service, for the principal `hue/<hostname>@<YOUR-REALM.COM>`. The Hue principal is then used to impersonate other users for applications within Hue such as the Job Browser, File Browser and so on.

Other services, such as HDFS and MapReduce, do not use the Hue Kerberos Ticket Renewer. They obtain tickets at startup and use those tickets to obtain Delegation Tokens for various access privileges. Each service handles its own ticket renewal as needed.

To add a Kerberos Ticket Renewer role instance using Cloudera Manager:

1. Go to the Hue service.
2. Click the **Instances** tab.
3. Click the **Add Role Instances** button.
4. Assign the Kerberos Ticket Renewer role instance to the same host as the Hue server.
5. When the wizard is finished, the status will display **Finished** and the Kerberos Ticket Renewer role instance is configured. The Hue service will now work with the secure Hadoop cluster.
6. Repeat these steps for each Hue Server role.

Troubleshooting the Kerberos Ticket Renewer:

Configuring Authentication

If the Hue Kerberos Ticket Renewer does not start, check your KDC configuration and the ticket renewal property, `maxrenewlife`, for the `hue/<hostname>` and `krbtgt` principals to ensure they are renewable. If not, running the following commands on the KDC will enable renewable tickets for these principals.

```
kadmin.local: modprinc -maxrenewlife 90day krbtgt/YOUR_REALM.COM
kadmin.local: modprinc -maxrenewlife 90day +allow_renewable hue/<hostname>@YOUR-REALM.COM
```

Step 10: (Flume Only) Use Substitution Variables for the Kerberos Principal and Keytab

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)

As described in Flume [security configuration](#), if you are using Flume on a secure cluster you must configure the HDFS sink with the following configuration options in the `flume.conf` file:

- `hdfs.kerberosPrincipal` - fully-qualified principal.
- `hdfs.kerberosKeytab` - location on the local host of the keytab containing the user and host keys for the above principal

Here is an example of an HDFS sink configuration in the `flume.conf` file (the majority of the HDFS sink configuration options have been omitted):

```
agent.sinks.sink-1.type = HDFS
agent.sinks.sink-1.hdfs.kerberosPrincipal = flume/_HOST@YOUR-REALM.COM
agent.sinks.sink-1.hdfs.kerberosKeytab = /etc/flume-ng/conf/flume.keytab
agent.sinks.sink-1.hdfs.proxyUser = webloggs
```

Since Cloudera Manager generates the Flume keytab files for you, and the locations of the keytab files cannot be known beforehand, substitution variables are required for Flume. Cloudera Manager provides two Flume substitution variables called `$KERBEROS_PRINCIPAL` and `$KERBEROS_KEYTAB` to configure the principal name and the keytab file path respectively on each host.

Here is an example of using the substitution variables to configure the options shown in the previous example:

```
agent.sinks.sink-1.type = hdfs
agent.sinks.sink-1.hdfs.kerberosPrincipal = $KERBEROS_PRINCIPAL
agent.sinks.sink-1.hdfs.kerberosKeytab = $KERBEROS_KEYTAB
agent.sinks.sink-1.hdfs.proxyUser = webloggs
```

Use the following instructions to have Cloudera Manager add these variables to the `flume.conf` file on every host that Cloudera Manager manages.

To use the Flume substitution variables for the Kerberos principal and keytab:

1. Go to the **Flume service > Configuration** page in Cloudera Manager.
2. Click **Agent**.
3. In the **Configuration File** property, add the configuration options with the substitution variables. For example:

```
agent.sinks.sink-1.type = hdfs
agent.sinks.sink-1.hdfs.kerberosPrincipal = $KERBEROS_PRINCIPAL
agent.sinks.sink-1.hdfs.kerberosKeytab = $KERBEROS_KEYTAB
agent.sinks.sink-1.hdfs.proxyUser = webloggs
```

4. Click **Save**.

Step 11: (CDH 4.0 and 4.1 only) Configure Hue to Use a Local Hive Metastore

If using Hue and the **Bypass Hive Metastore Server** option is not selected (metastore bypass is disabled by default), then Hue will not be able to communicate with Hive with CDH 4.0 or CDH 4.1. This is not a problem with CDH 4.2 or higher.

If you are using CDH 4.0 or 4.1, you can workaround this issue following the instructions in the Known Issues section of the [Cloudera Manager 4 Release Notes](#).

Step 12: Start All Services

Minimum Required Role: [Operator](#) (also provided by **Configurator**, **Cluster Administrator**, **Full Administrator**)

Start all services on your cluster:

Starting All Services

1. On the **Home > Status** tab, click



to the right of the cluster name and select **Start**.

2. Click **Start** that appears in the next screen to confirm. The **Command Details** window shows the progress of starting services.

When **All services successfully started** appears, the task is complete and you can close the **Command Details** window.

Starting the Cloudera Management Service

1. On the **Home > Status** tab, click



to the right of **Cloudera Management Service** and select **Start**.

2. Click **Start** to confirm. The **Command Details** window shows the progress of starting the roles.
3. When **Command completed with n/n successful subcommands** appears, the task is complete. Click **Close**.

Step 13: Deploy Client Configurations

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)

1. On the **Home > Status** tab, click



to the right of the cluster name and select **Deploy Client Configuration**.

2. Click **Deploy Client Configuration**.

Step 14: Create the HDFS Superuser Principal

To be able to create home directories for users, you will need access to the HDFS superuser account. (CDH automatically created the HDFS superuser account on each cluster host during CDH installation.) When you enabled Kerberos for the HDFS service, you lost access to the default HDFS superuser account using `sudo -u hdfs` commands. Cloudera recommends you use a different user account as the superuser, not the default `hdfs` account.

Designating a Non-Default Superuser Group

To designate a different group of superusers instead of using the default `hdfs` account, follow these steps:

1. Go to the Cloudera Manager Admin Console and navigate to the HDFS service.
2. Click the **Configuration** tab.
3. Select **Scope > HDFS (Service-Wide)**.
4. Select **Category > Security**.
5. Locate the **Superuser Group** property and change the value to the appropriate group name for your environment. For example, `<superuser>`.
6. Click **Save Changes** to commit the changes.
7. Restart the HDFS service.

To enable your access to the superuser account now that Kerberos is enabled, you must now create a Kerberos principal or an Active Directory user whose first component is `<superuser>`:

Configuring Authentication

If you are using Active Directory

Add a new user account to Active Directory, `<superuser>@YOUR-REALM.COM`. The password for this account should be set to never expire.

If you are using MIT KDC

1. In the `kadmin.local` or `kadmin` shell, type the following command to create a Kerberos principal called `<superuser>`:

```
kadmin: addprinc <superuser>@YOUR-LOCAL-REALM.COM
```

This command prompts you to create a password for the `<superuser>` principal. You should use a strong password because having access to this principal provides superuser access to all of the files in HDFS.

2. To run commands as the HDFS superuser, you must obtain Kerberos credentials for the `<superuser>` principal. To do so, run the following command and provide the appropriate password when prompted.

```
$ kinit <superuser>@YOUR-LOCAL-REALM.COM
```

Step 15: Get or Create a Kerberos Principal for Each User Account

Now that Kerberos is configured and enabled on your cluster, you and every other Hadoop user must have a Kerberos principal or keytab to obtain Kerberos credentials to be allowed to access the cluster and use the Hadoop services. In the next step of this procedure, you will need to create your own Kerberos principals in order to verify that Kerberos security is working on your cluster. If you and the other Hadoop users already have a Kerberos principal or keytab, or if your Kerberos administrator can provide them, you can skip ahead to the next step.

The following instructions explain how to create a Kerberos principal for a user account.

If you are using Active Directory

Add a new AD user account, `<username>@EXAMPLE.COM` for each Cloudera Manager service that should use Kerberos authentication. The password for these service accounts should be set to never expire.

If you are using MIT KDC

1. In the `kadmin.local` or `kadmin` shell, use the following command to create a principal for your account by replacing `EXAMPLE.COM` with the name of your realm, and replacing `username` with a username:

```
kadmin: addprinc username@EXAMPLE.COM
```

2. When prompted, enter the password twice.

Step 16: Prepare the Cluster for Each User

Before you and other users can access the cluster, there are a few tasks you must do to prepare the hosts for each user.

1. Make sure all hosts in the cluster have a Unix user account with the same name as the first component of that user's principal name. For example, the Unix account `joe` should exist on every box if the user's principal name is `joe@YOUR-REALM.COM`. You can use LDAP for this step if it is available in your organization.



Note: Each account must have a user ID that is greater than or equal to 1000. In the `/etc/hadoop/conf/taskcontroller.cfg` file, the default setting for the `banned.users` property is `mapred,hdfs,bin` to prevent jobs from being submitted from those user accounts. The default setting for the `min.user.id` property is `1000` to prevent jobs from being submitted with a user ID less than 1000, which are conventionally Unix super users.

2. Create a subdirectory under `/user` on HDFS for each user account (for example, `/user/joe`). Change the owner and group of that directory to be the user.

```
$ hadoop fs -mkdir /user/joe
$ hadoop fs -chown joe /user/joe
```



Note: `sudo -u hdfs` is not included in the commands above. This is because it is not required if Kerberos is enabled on your cluster. You will, however, need to have Kerberos credentials for the HDFS super user in order to successfully run these commands. For information on gaining access to the HDFS super user account, see [Step 14: Create the HDFS Superuser Principal](#) on page 83

Step 17: Verify that Kerberos Security is Working

After you have Kerberos credentials, you can verify that Kerberos security is working on your cluster by trying to run MapReduce jobs. To confirm, try launching a sleep or a pi job from the provided Hadoop examples (`/usr/lib/hadoop/hadoop-examples.jar`).



Note:

This section assumes you have a fully-functional CDH cluster and you have been able to access HDFS and run MapReduce jobs before you followed these instructions to configure and enable Kerberos on your cluster. If you have not already done so, you should at a minimum use the Cloudera Manager Admin Console to generate a client configuration file to enable you to access the cluster. For instructions, see [Deploying Client Configuration Files](#).

To verify that Kerberos security is working:

1. Acquire Kerberos credentials for your user account.

```
$ kinit USERNAME@YOUR-LOCAL-REALM.COM
```

2. Enter a password when prompted.
3. Submit a sample pi calculation as a test MapReduce job. Use the following command if you use a package-based setup for Cloudera Manager:

```
$ hadoop jar /usr/lib/hadoop-0.20/hadoop-0.20.2*examples.jar pi 10 10000
Number of Maps = 10
Samples per Map = 10000
...
Job Finished in 38.572 seconds
Estimated value of Pi is 3.14120000000000000000
```

If you have a parcel-based setup, use the following command instead:

```
$ hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-0.20-mapreduce/hadoop-examples.jar pi
10 10000
Number of Maps = 10
Samples per Map = 10000
...
Job Finished in 30.958 seconds
Estimated value of Pi is 3.14120000000000000000
```

You have now verified that Kerberos security is working on your cluster.



Important:

Running a MapReduce job will fail if you do not have a valid Kerberos ticket in your credentials cache. You can examine the Kerberos tickets currently in your credentials cache by running the `klist` command. You can obtain a ticket by running the `kinit` command and either specifying a keytab file containing credentials, or entering the password for your principal. If you do not have a valid ticket, you will receive an error such as:

```
11/01/04 12:08:12 WARN ipc.Client:
Exception encountered while connecting to the server :
javax.security.sasl.SaslException:GSS initiate failed
[Caused by GSSException: No valid credentials provided (Mechanism level:
Failed to find any
Kerberos tgt)]
Bad connection to FS. command aborted. exception: Call to
nn-host/10.0.0.2:8020 failed on local exception:
java.io.IOException:javax.security.sasl.SaslException: GSS initiate
failed
[Caused by GSSException: No valid credentials provided
(Mechanism level: Failed to find any Kerberos tgt)]
```

Step 18: (Optional) Enable Authentication for HTTP Web Consoles for Hadoop Roles

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)

Authentication for access to the HDFS, MapReduce, and YARN roles' web consoles can be enabled using a configuration option for the appropriate service. To enable this authentication:

1. From the **Clusters** tab, select the service (HDFS, MapReduce, or YARN) for which you want to enable authentication.
2. Click the **Configuration** tab.
3. Select **Scope** > *serviceName* **Service-Wide**.
4. Select **Category** > **Security**.
5. Select **Enable Kerberos Authentication for HTTP Web-Consoles**.
6. Click **Save Changes** to commit the changes.
7. Once the command finishes, restart all roles of that service.

Configuring Authentication in the Cloudera Navigator Data Management Component

Cloudera Navigator data management component supports user authentication against Cloudera Manager user accounts and against an external LDAP or Active Directory service. External authentication enables you to assign Cloudera Navigator user roles to LDAP or Active Directory groups containing the appropriate users for each user role. Authentication with a Cloudera Manager user account requires either the Full Administrator or Navigator Administrator user role, and enables the user to use Cloudera Navigator features or to configure the external authentication service.

Configuring External Authentication for the Cloudera Navigator Data Management Component

Minimum Required Role: [Navigator Administrator](#) (also provided by **Full Administrator**)



Important: This feature is available only with a Cloudera Enterprise license; it is not available in Cloudera Express. For information on Cloudera Enterprise licenses, see [Managing Licenses](#).

Cloudera Navigator supports user authentication against Cloudera Manager user accounts and against an external service. The external service can be either LDAP or Active Directory. User authentication against Cloudera Manager user accounts requires users to have one of two Cloudera Manager user roles, either Full Administrator or Navigator Administrator. External authentication enables you to assign Cloudera Navigator user roles to LDAP or Active Directory groups to which the appropriate users belong.

For more information about Cloudera Manager user accounts, see [Cloudera Manager User Accounts](#) on page 50. For more information about Cloudera Navigator user roles, see [Cloudera Navigator Data Management Component User Roles](#) on page 347.

The following sections describe how to configure the supported external directory services.

Configuring Cloudera Navigator Authentication Using Active Directory



Important:

Cloudera Navigator has its own role-based access control and user management scheme. If you want to use LDAP/AD authentication, Cloudera Navigator roles must be explicitly assigned to AD users to allow them to log in to Navigator. To assign roles to AD users, log in to Cloudera Navigator for the first time using a Cloudera Manager admin user. Any *non-externally authenticated* Cloudera Manager user that has Full Administrator or Navigator Administrator privileges will have admin access to Cloudera Navigator. You can use this account to set up user groups and assign Cloudera Navigator roles to AD users.

Hence, Cloudera recommends that the **Authentication Backend Order** property be set initially to **Cloudera Manager then External**. Otherwise, the external authentication system will be checked first, and if the same user credentials also exist in the specified LDAP or Active Directory, the user will be authenticated there, and will not be authenticated as a Cloudera Manager administrator. Since no user roles will have been set up yet for the users in the external authentication system, the user's attempt to log in will fail. Once the groups and user roles for Cloudera Navigator are set up, the **Authentication Backend Order** can be changed to **External then Cloudera Manager** or **External Only**, if desired.

To configure Cloudera Navigator to use AD authentication:

1. Select **Clusters > Cloudera Management Service**.
2. Click the **Configuration** tab.
3. Select **Scope > Navigator Metadata Server**.
4. Select **Category > External Authentication**.
5. In the **Authentication Backend Order** field, select the order in which Cloudera Navigator should attempt its authentication. You can choose to authenticate users using just one of the methods (using Cloudera Manager user accounts is the default), or you can set it so that if the user cannot be authenticated by the first method, it will attempt using the second method.
6. In the **External Authentication Type** property, select **Active Directory**.
7. In the **LDAP URL** property, provide the URL of the LDAP/Active Directory server to authenticate against. The URL has the form: `[ldap|ldaps]://hostname:port`. If a port is not specified, the default LDAP port is used (389 for LDAP and 636 for LDAPS). For more details on the LDAP URL format, see [RFC 2255](#).
8. For the **Bind Distinguished Name** property, you require the `userPrincipalName` (UPN) of the user to bind as. For example, if the UPN is `sampleuser@EXAMPLE.COM`, the **Bind Distinguished Name** provided can be either `sampleuser` or the complete UPN, `sampleuser@EXAMPLE.COM`. This is used to connect to Active Directory for searching groups and to get other user information.
9. In the **LDAP Bind Password**, enter the password for the bind user entered above.
10. In the **Active Directory NT Domain** property, provide the NT domain to authenticate against.
11. Click **Save Changes** to commit the changes.
12. After changing the configuration settings, restart the **Navigator Metadata Service**: click the **Instances** tab on the **Cloudera Management Service** page, check **Navigator Metadata Service**, and click **Actions for Selected > Restart**.

Configuring Cloudera Navigator Authentication Using an OpenLDAP-compatible Server

For an OpenLDAP-compatible directory, you have several options for searching for users and groups:

- You can specify a single base Distinguished Name (DN) and then provide a "Distinguished Name Pattern" to use to match a specific user in the LDAP directory.

Configuring Authentication

- Search filter options let you search for a particular user based on somewhat broader search criteria – for example Cloudera Navigator users could be members of different groups or organizational units (OUs), so a single pattern won't find all those users. Search filter options also let you find all the groups to which a user belongs, to help determine if that user should be allowed to log in.

1. Select **Clusters > Cloudera Management Service**.
2. Click the **Configuration** tab.
3. Select **Scope > Navigator Metadata Server**.
4. Select **Category > External Authentication**.
5. In the **Authentication Backend Order** field, select the order in which Cloudera Navigator should attempt its authentication. You can choose to authenticate users using just one of the methods (using Cloudera Manager user accounts is the default), or you can set it so that if the user cannot be authenticated by the first method, it will attempt using the second method.
6. In the **External Authentication Type**, select **LDAP**.
7. In the **LDAP URL** property, provide the URL of the LDAP server and (optionally) the base Distinguished Name (DN) (the search base) as part of the URL — for example `ldap://ldap-server.corp.com/dc=corp,dc=com`.
8. In the **LDAP Bind User Distinguished Name** property, enter the user's sAMAccountName. This is used to connect to the LDAP server for searching groups and to get other user information.
9. In the **LDAP Bind Password** property, enter the password for the bind user entered above.
- 10 To use a single "Distinguished Name Pattern", provide a pattern in the **LDAP Distinguished Name Pattern** property.

Use `{0}` in the pattern to indicate where the username should go. For example, to search for a distinguished name where the `uid` attribute is the username, you might provide a pattern similar to `uid={0},ou=People,dc=corp,dc=com`. Cloudera Navigator substitutes the name provided at login into this pattern and performs a search for that specific user. So if a user provides the username "foo" at the Cloudera Navigator login page, Cloudera Navigator will search for the DN `uid=foo,ou=People,dc=corp,dc=com`.

If you provided a base DN along with the URL, the pattern only needs to specify the rest of the DN pattern. For example, if the URL you provide is `ldap://ldap-server.corp.com/dc=corp,dc=com`, and the pattern is `uid={0},ou=People`, then the search DN will be `uid=foo,ou=People,dc=corp,dc=com`.

- 11 You can also search using User or Group search filters, using the **LDAP User Search Base**, **LDAP User Search Filter**, **LDAP Group Search Base** and **LDAP Group Search Filter** settings. These allow you to combine a base DN with a search filter to allow a greater range of search targets.

For example, if you want to authenticate users who may be in one of multiple OUs, the search filter mechanism will allow this. You can specify the User Search Base DN as `dc=corp,dc=com` and the user search filter as `uid={0}`. Then Cloudera Navigator will search for the user anywhere in the tree starting from the Base DN. Suppose you have two OUs—`ou=Engineering` and `ou=Operations`—Cloudera Navigator will find User "foo" if it exists in either of these OUs, that is, `uid=foo,ou=Engineering,dc=corp,dc=com` or `uid=foo,ou=Operations,dc=corp,dc=com`.

You can use a user search filter along with a DN pattern, so that the search filter provides a fallback if the DN pattern search fails.

The Groups filters let you search to determine if a DN or username is a member of a target group. In this case, the filter you provide can be something like `member={0}` where `{0}` will be replaced with the **DN** of the user you are authenticating. For a filter requiring the username, `{1}` may be used, as `memberUId={1}`. This will return a list of groups to which the user belongs.

- 12 Click **Save Changes** to commit the changes.
- 13 After changing the configuration settings, restart the **Navigator Metadata Service**: click the **Instances** tab on the **Cloudera Management Service** page, check **Navigator Metadata Service**, and click **Actions for Selected > Restart**.

Configuring Cloudera Navigator to Use LDAPS

If the LDAP server certificate has been signed by a trusted Certificate Authority (that is, VeriSign, GeoTrust, and so on), steps 1 and 2 below may not be necessary.

1. Copy the CA certificate file to the Cloudera Navigator Server host.
2. Import the CA certificate(s) from the CA certificate file to the local truststore. The default truststore is located in the `$JAVA_HOME/jre/lib/security/cacerts` file. This contains the default CA information shipped with the JDK. Create an alternate default file called `jssecacerts` in the same location as the `cacerts` file. You can now safely append CA certificates for any private or public CAs not present in the default `cacerts` file, while keeping the original file intact.

For our example, we will follow this recommendation by copying the default `cacerts` file into the new `jssecacerts` file, and then importing the CA certificate to this alternate truststore.

```
$ cp $JAVA_HOME/jre/lib/security/cacerts \
  $JAVA_HOME/jre/lib/jssecacerts
```

```
$ /usr/java/latest/bin/keytool -import -alias nt_domain_name \
  -keystore /usr/java/latest/jre/lib/security/jssecacerts -file path_to_cert
```



Note:

- The default password for the `cacerts` store is **changeit**.
- The alias can be any name (not just the domain name).

3. Configure the **LDAP URL** property to use `ldaps://ldap_server` instead of `ldap://ldap_server`.

Configuring Cloudera Navigator Authentication Using SAML

Cloudera Navigator supports the Security Assertion Markup Language (SAML), an XML-based open standard data format for exchanging authentication and authorization data between parties, in particular, between an identity provider (IDP) and a service provider (SP). The SAML specification defines three roles: the principal (typically a user), the IDP, and the SP. In the use case addressed by SAML, the principal (user agent) requests a service from the service provider. The service provider requests and obtains an identity assertion from the IDP. On the basis of this assertion, the SP can make an access control decision—in other words it can decide whether to perform some service for the connected principal.

The primary SAML use case is called web browser single sign-on (SSO). A user wielding a user agent (usually a web browser) requests a web resource protected by a SAML SP. The SP, wishing to know the identity of the requesting user, issues an authentication request to a SAML IDP through the user agent. In the context of this terminology, Cloudera Navigator operates as a SP. This topic discusses the Cloudera Navigator part of the configuration process; it assumes that you are familiar with SAML and SAML configuration in a general sense, and that you have a functioning IDP already deployed.

Setting up Cloudera Navigator to use SAML requires the following steps.

Preparing Files

You will need to prepare the following files and information, and provide these to Cloudera Navigator:

- A Java keystore containing a private key for Cloudera Navigator to use to sign/encrypt SAML messages.
- The SAML metadata XML file from your IDP. This file must contain the public certificates needed to verify the sign/encrypt key used by your IDP per the SAML Metadata Interoperability Profile.
- The entity ID that should be used to identify the Navigator Metadata Server instance.
- How the user ID is passed in the SAML authentication response:
 - As an attribute. If so, what identifier is used.
 - As the NameID.
- The method by which the Cloudera Navigator role will be established:
 - From an attribute in the authentication response:
 - What identifier will be used for the attribute

Configuring Authentication

- What values will be passed to indicate each role
- From an external script that will be called for each use:
 - The script takes user ID as \$1
 - The script must assign an exit code to reflect successful authentication of the assigned role:
 - 0 - Full Administrator
 - 1 - User Administrator
 - 2 - Auditing Viewer
 - 4 - Lineage Viewer
 - 8 - Metadata Administrator
 - 16 - Policy Viewer
 - 32 - Policy Administrator
 - A negative value is returned for a failure to authenticate

To assign more than one role, add the numbers for the roles. For example, to assign the Policy Viewer and User Administrator roles, the exit code should be 17.

Configuring Cloudera Navigator

1. Select **Clusters > Cloudera Management Service**.
2. Click the **Configuration** tab.
3. Select **Scope > Navigator Metadata Server**.
4. Select **Category > External Authentication**.
5. Type `SAML` in the Search box.
6. Set the **External Authentication Type** property to **SAML** (the **Authentication Backend Order** property is ignored for SAML).
7. Set the **Path to SAML IDP Metadata File** property to point to the IDP metadata file.
8. Set the **Path to SAML Keystore File** property to point to the Java keystore file containing the Cloudera Navigator private key (prepared above).
9. In the **SAML Keystore Password** property, set the SAML keystore password.
10. In the **Alias of SAML Sign/Encrypt Private Key** property, set the alias used to identify the private key for Cloudera Navigator to use.
11. In the **SAML Sign/Encrypt Private Key Password** property, set the password for the sign/encrypt private key.
12. Set the **SAML Entity ID** property if:
 - There is more than one Cloudera Navigator instance being used with the same IDP (each instance needs a different entity ID).
 - Entity IDs are assigned by organizational policy.

The entity ID value should be unique to the current Navigator Metadata Server installation.

13. In the **Source of User ID in SAML Response** property, set whether the user ID will be obtained from an attribute or the NameID.

If an attribute will be used, set the attribute name in the **SAML Attribute Identifier for User ID** property. The default value is the normal OID used for user IDs and so may not need to be changed.

14. In the **SAML Role Assignment Mechanism** property, set whether the role assignment will be done from an attribute or an external script.
 - If an attribute will be used:
 - In the **SAML Attribute Identifier for User Role** property, set the attribute name if necessary. The default value is the normal OID used for OrganizationalUnits and so may not need to be changed.
 - In the **SAML Attribute Values for Roles** property, set which attribute values will be used to indicate the user role.

- If an external script will be used, set the path to that script in the **Path to SAML Role Assignment Script** property. Make sure that the script is executable (an executable binary is fine - it doesn't need to be a shell script).

15 Click **Save Changes** to commit the changes.

16 [Restart](#) the Navigator Metadata Server role.

Configuring the IDP

After the Cloudera Navigator is restarted, it will attempt to redirect to the IDP login page instead of showing the normal Cloudera Navigator login page. This may or may not succeed, depending on how the IDP is configured. In either case, the IDP will need to be configured to recognize Cloudera Navigator before authentication will actually succeed. The details of this process are specific to each IDP implementation - refer to your IDP documentation for details.

1. Download Cloudera Navigator's SAML metadata XML file from `http://hostname:7187/saml/metadata`.
2. Inspect the metadata file and ensure that any URLs contained in the file can be resolved by users' web browsers. The IDP will redirect web browsers to these URLs at various points in the process. If the browser cannot resolve them, authentication will fail. If the URLs are incorrect, you can manually fix the XML file or set the **SAML Entity Base URL** property in the Navigator Metadata Server configuration to the right value, and then re-download the file.
3. Provide this metadata file to your IDP using whatever mechanism your IDP provides.
4. Ensure that the IDP has access to whatever public certificates are necessary to validate the private key that was provided by Cloudera Navigator earlier.
5. Ensure that the IDP is configured to provide the User ID and Role using the attribute names that Cloudera Navigator was configured to expect, if relevant.
6. Ensure the changes to the IDP configuration have taken effect (a restart may be necessary).

Verifying Authentication and Authorization

1. Return to the Cloudera Navigator home page at: `http://hostname:7187/`.
2. Attempt to log in with credentials for a user that is entitled. The authentication should complete and you should see the Home page.
3. If authentication fails, you will see an IDP provided error message. Cloudera Navigator is not involved in this part of the process, and you must ensure the IDP is working correctly to complete the authentication.
4. If authentication succeeds but the user is not authorized to use Cloudera Navigator, they will be taken to an error page that explains the situation. If a user who should be authorized sees this error, then you will need to verify their role configuration, and ensure that it is being properly communicated to the Navigator Metadata Server, whether by attribute or external script. The Cloudera Navigator log will provide details on failures to establish a user's role. If any errors occur during role mapping, Cloudera Navigator will assume the user is unauthorized.

Bypassing SAML SSO

You can bypass SAML SSO by directly accessing the Cloudera Navigator login page at `http://hostname:7187/login.html`. If the Cloudera Management Service property **Authentication Backend Order**, is set to anything but **External Only**, a Cloudera Manager user will be able to log in to Cloudera Navigator.

Managing Users and Groups for the Cloudera Navigator Data Management Component

[Required Role:](#)

User Administrator

Full Administrator

These required roles refer to [Cloudera Navigator user roles](#). Users with the [Cloudera Manager user roles](#) Navigator Administrator or Full Administrator who log into the Cloudera Navigator Web UI with their Cloudera Manager credentials are logged into Cloudera Navigator with the Full Administrator Cloudera Navigator user role.

Configuring Authentication

Cloudera Navigator supports user authentication against Cloudera Manager user accounts and against an external LDAP or Active Directory service. External authentication enables you to assign Cloudera Navigator user roles to LDAP or Active Directory groups containing the appropriate users for each user role.

Assigning Cloudera Navigator User Roles to LDAP or Active Directory Groups

This section assumes that values for your LDAP or Active Directory directory service have been configured in Cloudera Manager as described in [Configuring External Authentication for Cloudera Navigator](#). This section also assumes that your LDAP or Active Directory service contains user groups that correspond to Cloudera Navigator user roles having the permissions you want each group of users to have. If not, you should assign your users to such groups now. The Cloudera Navigator user roles are as follows:

- Full Administrator
- User Administrator
- Auditing Viewer
- Lineage Viewer
- Metadata Administrator
- Policy Viewer
- Policy Administrator

Each of these roles and the permissions associated with them are described in [Cloudera Navigator User Roles](#).

To add or remove Cloudera Navigator user roles to LDAP or Active Directory user groups, you should know the names of the directory groups you want to configure, and then perform the following steps:

1. Do one of the following:

- Enter the URL of the Navigator UI in a browser: `http://Navigator_Metadata_Server_host:port/`, where *Navigator_Metadata_Server_host* is the name of the host on which you are running the Navigator Metadata Server role and *port* is the port configured for the role. The default port of the Navigator Metadata Server is 7187. To change the port, follow the instructions in [Configuring the Navigator Metadata Server Port](#).
- Do one of the following:
 - Select **Clusters > Cloudera Management Service > Cloudera Navigator**.
 - Navigate from the Navigator Metadata Server role:
 1. Do one of the following:
 - Select **Clusters > Cloudera Management Service > Cloudera Management Service**.
 - On the Status tab of the **Home > Status** tab, in **Cloudera Management Service** table, click the **Cloudera Management Service** link.
 2. Click the **Instances** tab.
 3. Click the **Navigator Metadata Server** role.
 4. Click the **Cloudera Navigator** link.

2. Log in to Cloudera Navigator with the credentials of a user having one or more of the following user roles:

- Cloudera Manager Full Administrator
- Cloudera Manager Navigator Administrator
- Cloudera Navigator Full Administrator
- Cloudera Navigator User Administrator

3. Click the **Administration** tab in the upper right.

4. Search for an LDAP or Active Directory group by entering its name (or the first portion of the name) in the search field and pressing **Enter** or **Return**.

- Select **All Groups** to search among all groups in the external directory.
- Select **Groups with Navigator Roles** to display only external directory groups that have already been assigned one or more Cloudera Navigator user roles.

5. From the LDAP or Active Directory groups displayed, select the group to which you want to assign a Cloudera Navigator user role or roles. If roles have already been assigned to the group, they are listed beneath the name of the group in the main panel.
6. Click **Manage Role Assignment** in the upper right.
7. Click the checkbox for each Cloudera Navigator user role you want assigned to that Active Directory or LDAP group. Uncheck any already-assigned roles that you want to remove from the group.
8. Click **Save**.

If a user's role assignments are changed, the changes take effect with the user's next new session, that is, the next time the user logs in to Cloudera Navigator.

Configuring Authentication in CDH Using the Command Line

The security features in CDH 5 enable Hadoop to prevent malicious user impersonation. The Hadoop daemons leverage Kerberos to perform user authentication on all remote procedure calls (RPCs). Group resolution is performed on the Hadoop master nodes, NameNode, JobTracker and ResourceManager to guarantee that group membership cannot be manipulated by users. Map tasks are run under the user account of the user who submitted the job, ensuring isolation there. In addition to these features, new authorization mechanisms have been introduced to HDFS and MapReduce to enable more control over user access to data.

The security features in CDH 5 meet the needs of most Hadoop customers because typically the cluster is accessible only to trusted personnel. In particular, Hadoop's current threat model assumes that users cannot:

1. Have `root` access to cluster machines.
2. Have `root` access to shared client machines.
3. Read or modify packets on the network of the cluster.



Note:

CDH 5 supports encryption of all user data sent over the network. For configuration instructions, see [Configuring Encrypted Shuffle, Encrypted Web UIs, and Encrypted HDFS Transport](#).

Note also that there is no built-in support for on-disk encryption.

Enabling Kerberos Authentication for Hadoop Using the Command Line



Important:

These instructions assume you know how to install and configure Kerberos, you already have a working Kerberos Key Distribution Center (KDC) and realm setup, and that you've installed the Kerberos user packages on all cluster machines and machines which will be used to access the cluster. Furthermore, Oozie and Hue require that the realm support renewable tickets. For more information about installing and configuring Kerberos, see:

- [MIT Kerberos Home](#)
- [MIT Kerberos Documentation](#)
- [Kerberos Explained](#)
- [Microsoft Kerberos Overview](#)
- [Microsoft Kerberos in Windows Server 2008](#)
- [Microsoft Kerberos in Windows Server 2003](#)

Kerberos security in CDH 5 has been tested with the following version of MIT Kerberos 5:

- krb5-1.6.1 on Red Hat Enterprise Linux 5 and CentOS 5

Kerberos security in CDH 5 is supported with the following versions of MIT Kerberos 5:

Configuring Authentication

- krb5-1.6.3 on SUSE Linux Enterprise Server (SLES) 11 Service Pack 1
- krb5-1.8.1 on Ubuntu
- krb5-1.8.2 on Red Hat Enterprise Linux 6 and CentOS 6
- krb5-1.9 on Red Hat Enterprise Linux 6.1



Note: The `krb5-server` package includes a `logrotate` policy file to rotate log files monthly. To take advantage of this, install the `logrotate` package. No additional configuration is necessary.

If you want to enable Kerberos SPNEGO-based authentication for the Hadoop web interfaces, see the [Hadoop Auth, Java HTTP SPNEGO Documentation](#).

Here are the general steps to configuring secure Hadoop, each of which is described in more detail in the following sections:

Step 1: Install CDH 5

Cloudera strongly recommends that you set up a fully-functional CDH 5 cluster before you begin configuring it to use Hadoop's security features. When a secure Hadoop cluster is not configured correctly, the resulting error messages are in a preliminary state, so it's best to start implementing security after you are sure your Hadoop cluster is working properly without security.

For information about installing and configuring Hadoop and CDH 5 components, and deploying them on a cluster, see [Cloudera Installation Guide](#).

Step 2: Verify User Accounts and Groups in CDH 5 Due to Security



Note: CDH 5 introduces a new version of MapReduce: MapReduce 2.0 (MRv2) built on the YARN framework. In this document, we refer to this new version as YARN. CDH 5 also provides an implementation of the previous version of MapReduce, referred to as MRv1 in this document.

- If you are using MRv1, see [Step 2a \(MRv1 only\): Verify User Accounts and Groups in MRv1](#) on page 94 for configuration information.
- If you are using YARN, see [Step 2b \(YARN only\): Verify User Accounts and Groups in YARN](#) on page 96 for configuration information.

Step 2a (MRv1 only): Verify User Accounts and Groups in MRv1



Note: If you are using YARN, skip this step and proceed to [Step 2b \(YARN only\): Verify User Accounts and Groups in YARN](#).

During CDH 5 package installation of MRv1, the following Unix user accounts are automatically created to support security:

This User	Runs These Hadoop Programs
<code>hdfs</code>	HDFS: NameNode, DataNodes, Secondary NameNode (or Standby NameNode if you are using HA)
<code>mapred</code>	MRv1: JobTracker and TaskTrackers

The `hdfs` user also acts as the HDFS superuser.

The `hadoop` user no longer exists in CDH 5. If you currently use the `hadoop` user to run applications as an HDFS super-user, you should instead use the new `hdfs` user, or create a separate Unix account for your application such as `myhadoopapp`.

MRv1: Directory Ownership in the Local File System

Because the HDFS and MapReduce services run as different users, you must be sure to configure the correct directory ownership of the following files on the local file system of each host:

File System	Directory	Owner	Permissions
Local	<code>dfs.namenode.name.dir</code> (<code>dfs.name.dir</code> is deprecated but will also work)	<code>hdfs:hdfs</code>	<code>drwx-----</code>
Local	<code>dfs.datanode.data.dir</code> (<code>dfs.data.dir</code> is deprecated but will also work)	<code>hdfs:hdfs</code>	<code>drwx-----</code>
Local	<code>mapred.local.dir</code>	<code>mapred:mapred</code>	<code>drwxr-xr-x</code>

See also [Deploying MapReduce v1 \(MRv1\) on a Cluster](#).

You must also configure the following permissions for the HDFS and MapReduce log directories (the default locations in `/var/log/hadoop-hdfs` and `/var/log/hadoop-0.20-mapreduce`), and the `$MAPRED_LOG_DIR/userlogs/` directory:

File System	Directory	Owner	Permissions
Local	<code>HDFS_LOG_DIR</code>	<code>hdfs:hdfs</code>	<code>drwxrwxr-x</code>
Local	<code>MAPRED_LOG_DIR</code>	<code>mapred:mapred</code>	<code>drwxrwxr-x</code>
Local	userlogs directory in <code>MAPRED_LOG_DIR</code>	<code>mapred:anygroup</code>	<i>permissions will be set automatically at daemon start time</i>

MRv1: Directory Ownership on HDFS

The following directories on HDFS must also be configured as follows:

File System	Directory	Owner	Permissions
HDFS	<code>mapreduce.jobtracker.system.dir</code> (<code>mapred.system.dir</code> is deprecated but will also work)	<code>mapred:hadoop</code>	<code>drwx-----</code>
HDFS	<code>/</code> (root directory)	<code>hdfs:hadoop</code>	<code>drwxr-xr-x</code>

MRv1: Changing the Directory Ownership on HDFS

- If Hadoop security is enabled, use `kinit hdfs` to obtain Kerberos credentials for the `hdfs` user by running the following commands before changing the directory ownership on HDFS:

```
$ sudo -u hdfs kinit -k -t hdfs.keytab hdfs/fully.qualified.domain.name@YOUR-REALM.COM
```

¹ In CDH 5, package installation and the Hadoop daemons will automatically configure the correct permissions for you if you configure the directory ownership correctly as shown in the table above.

² When starting up, MapReduce sets the permissions for the `mapreduce.jobtracker.system.dir` (or `mapred.system.dir`) directory in HDFS, assuming the user `mapred` owns that directory.

Configuring Authentication

If `kinit hdfs` does not work initially, run `kinit -R` after running `kinit` to obtain credentials. (For more information, see [Troubleshooting Authentication Issues](#) on page 194). To change the directory ownership on HDFS, run the following commands. Replace the example `/mapred/system` directory in the commands below with the HDFS directory specified by the `mapreduce.jobtracker.system.dir` (or `mapred.system.dir`) property in the `conf/mapred-site.xml` file:

```
$ sudo -u hdfs hadoop fs -chown mapred:hadoop /mapred/system
$ sudo -u hdfs hadoop fs -chown hdfs:hadoop /
$ sudo -u hdfs hadoop fs -chmod -R 700 /mapred/system
$ sudo -u hdfs hadoop fs -chmod 755 /
```

- In addition (whether or not Hadoop security is enabled) create the `/tmp` directory. For instructions on creating `/tmp` and setting its permissions, see [these instructions](#).

Step 2b (YARN only): Verify User Accounts and Groups in YARN



Note: If you are using MRv1, skip this step and proceed to [Step 3: If you are Using AES-256 Encryption, Install the JCE Policy File](#) on page 98.

During CDH 5 package installation of MapReduce 2.0 (YARN), the following Unix user accounts are automatically created to support security:

This User	Runs These Hadoop Programs
<code>hdfs</code>	HDFS: NameNode, DataNodes, Standby NameNode (if you are using HA)
<code>yarn</code>	YARN: ResourceManager, NodeManager
<code>mapred</code>	YARN: MapReduce Job History Server



Important: The HDFS and YARN daemons must run as different Unix users; for example, `hdfs` and `yarn`. The MapReduce Job History server must run as user `mapred`. Having all of these users share a common Unix group is recommended; for example, `hadoop`.

YARN: Directory Ownership in the Local File System

Because the HDFS and MapReduce services run as different users, you must be sure to configure the correct directory ownership of the following files on the local file system of each host:

File System	Directory	Owner	Permissions (see Footnote 1)
Local	<code>dfs.namenode.name.dir</code> (<code>dfs.name.dir</code> is deprecated but will also work)	<code>hdfs:hdfs</code>	<code>drwx-----</code>
Local	<code>dfs.datanode.data.dir</code> (<code>dfs.data.dir</code> is deprecated but will also work)	<code>hdfs:hdfs</code>	<code>drwx-----</code>
Local	<code>yarn.nodemanager.local-dirs</code>	<code>yarn:yarn</code>	<code>drwxr-xr-x</code>
Local	<code>yarn.nodemanager.log-dirs</code>	<code>yarn:yarn</code>	<code>drwxr-xr-x</code>
Local	<code>container-executor</code>	<code>root:yarn</code>	<code>--Sr-s---</code>

File System	Directory	Owner	Permissions (see Footnote 1)
Local	conf/container-executor.cfg	root:yarn	r-----



Important: Configuration changes to the Linux container executor could result in local NodeManager directories (such as `usercache`) being left with incorrect permissions. To avoid this, when making changes using either Cloudera Manager or the command line, first manually remove the existing NodeManager local directories from all configured local directories (`yarn.nodemanager.local-dirs`), and let the NodeManager recreate the directory structure.

You must also configure the following permissions for the HDFS, YARN and MapReduce log directories (the default locations in `/var/log/hadoop-hdfs`, `/var/log/hadoop-yarn` and `/var/log/hadoop-mapreduce`):

File System	Directory	Owner	Permissions
Local	HDFS_LOG_DIR	hdfs:hdfs	drwxrwxr-x
Local	\$YARN_LOG_DIR	yarn:yarn	drwxrwxr-x
Local	MAPRED_LOG_DIR	mapred:mapred	drwxrwxr-x

YARN: Directory Ownership on HDFS

The following directories on HDFS must also be configured as follows:

File System	Directory	Owner	Permissions
HDFS	/(root directory)	hdfs:hadoop	drwxr-xr-x
HDFS	yarn.nodemanager.remote-app-log-dir	yarn:hadoop	drwxrwxrwx
HDFS	mapreduce.jobhistory.intermediate-done-dir	mapred:hadoop	drwxrwxrwx
HDFS	mapreduce.jobhistory.done-dir	mapred:hadoop	drwxr-x---

YARN: Changing the Directory Ownership on HDFS

If Hadoop security is enabled, use `kinit hdfs` to obtain Kerberos credentials for the `hdfs` user by running the following commands:

```
$ sudo -u hdfs kinit -k -t hdfs.keytab hdfs/fully.qualified.domain.name@YOUR-REALM.COM
$ hadoop fs -chown hdfs:hadoop /
$ hadoop fs -chmod 755 /
```

If `kinit hdfs` does not work initially, run `kinit -R` after running `kinit` to obtain credentials. See [Troubleshooting Authentication Issues](#) on page 194. To change the directory ownership on HDFS, run the following commands:

```
$ sudo -u hdfs hadoop fs -chown hdfs:hadoop /
$ sudo -u hdfs hadoop fs -chmod 755 /
$ sudo -u hdfs hadoop fs -chown yarn:hadoop [yarn.nodemanager.remote-app-log-dir]
$ sudo -u hdfs hadoop fs -chmod 1777 [yarn.nodemanager.remote-app-log-dir]
$ sudo -u hdfs hadoop fs -chown mapred:hadoop [mapreduce.jobhistory.intermediate-done-dir]
$ sudo -u hdfs hadoop fs -chmod 1777 [mapreduce.jobhistory.intermediate-done-dir]
```

³ In CDH 5, package installation and the Hadoop daemons will automatically configure the correct permissions for you if you configure the directory ownership correctly as shown in the two tables above. See also [Deploying MapReduce v2 \(YARN\) on a Cluster](#).

Configuring Authentication

```
$ sudo -u hdfs hadoop fs -chown mapred:hadoop [mapreduce.jobhistory.done-dir]
$ sudo -u hdfs hadoop fs -chmod 750 [mapreduce.jobhistory.done-dir]
```

- In addition (whether or not Hadoop security is enabled) create the `/tmp` directory. For instructions on creating `/tmp` and setting its permissions, see [Step 7: If Necessary, Create the HDFS /tmp Directory](#).
- In addition (whether or not Hadoop security is enabled), change permissions on the `/user/history` Directory. See [Step 8: Create the history Directory and Set Permissions and Owner](#).

Step 3: If you are Using AES-256 Encryption, Install the JCE Policy File

If you are using CentOS/Red Hat Enterprise Linux 5.6 or higher, or Ubuntu, which use AES-256 encryption by default for tickets, you must install the [Java Cryptography Extension \(JCE\) Unlimited Strength Jurisdiction Policy File](#) on all cluster and Hadoop user machines. For JCE Policy File installation instructions, see the `README.txt` file included in the `jce_policy-x.zip` file.

Alternatively, you can configure Kerberos to not use AES-256 by removing `aes256-cts:normal` from the `supported_encetypes` field of the `kdc.conf` or `krb5.conf` file. After changing the `kdc.conf` file, you must restart both the KDC and the `kadmin` server for those changes to take affect. You may also need to re-create or change the password of the relevant principals, including potentially the Ticket Granting Ticket principal (`krbtgt/REALM@REALM`). If AES-256 is still used after completing steps, the `aes256-cts:normal` setting existed when the Kerberos database was created. To fix this, create a new Kerberos database and then restart both the KDC and the `kadmin` server.

To verify the type of encryption used in your cluster:

1. On the local KDC host, type this command to create a test principal:

```
$ kadmin -q "addprinc test"
```

2. On a cluster host, type this command to start a Kerberos session as test:

```
$ kinit test
```

3. On a cluster host, type this command to view the encryption type in use:

```
$ klist -e
```

If AES is being used, output like the following is displayed after you type the `klist` command; note that AES-256 is included in the output:

```
Ticket cache: FILE:/tmp/krb5cc_0
Default principal: test@SCM
Valid starting Expires Service principal
05/19/11 13:25:04 05/20/11 13:25:04 krbtgt/SCM@SCM
    Etype (skey, tkt): AES-256 CTS mode with 96-bit SHA-1 HMAC, AES-256 CTS mode with
96-bit SHA-1 HMAC
```

Step 4: Create and Deploy the Kerberos Principals and Keytab Files

A Kerberos principal is used in a Kerberos-secured system to represent a unique identity. Kerberos assigns tickets to Kerberos principals to enable them to access Kerberos-secured Hadoop services. For Hadoop, the principals should be of the format `username/fully.qualified.domain.name@YOUR-REALM.COM`. In this guide, the term `username` in the `username/fully.qualified.domain.name@YOUR-REALM.COM` principal refers to the username of an existing Unix account, such as `hdfs` or `mapred`.

A keytab is a file containing pairs of Kerberos principals and an encrypted copy of that principal's key. The keytab files are unique to each host since their keys include the hostname. This file is used to authenticate a principal on a host to Kerberos without human interaction or storing a password in a plain text file. Because having access to the keytab file for a principal allows one to act as that principal, access to the keytab files should be tightly secured. They should be readable by a minimal set of users, should be stored on local disk, and should not be included in machine backups, unless access to those backups is as secure as access to the local machine.

**Important:**

For both MRv1 and YARN deployments: *On every machine in your cluster, there must be a keytab file for the `hdfs` user and a keytab file for the `mapred` user. The `hdfs` keytab file must contain entries for the `hdfs` principal and a `HTTP` principal, and the `mapred` keytab file must contain entries for the `mapred` principal and a `HTTP` principal. On each respective machine, the `HTTP` principal will be the same in both keytab files.*

In addition, for YARN deployments only: *On every machine in your cluster, there must be a keytab file for the `yarn` user. The `yarn` keytab file must contain entries for the `yarn` principal and a `HTTP` principal. On each respective machine, the `HTTP` principal in the `yarn` keytab file will be the same as the `HTTP` principal in the `hdfs` and `mapred` keytab files.*

**Note:**

The following instructions illustrate an example of creating keytab files for MIT Kerberos. If you are using another version of Kerberos, refer to your Kerberos documentation for instructions. You may use either `kadmin` or `kadmin.local` to run these commands.

When to Use `kadmin.local` and `kadmin`

When creating the Kerberos principals and keytabs, you can use `kadmin.local` or `kadmin` depending on your access and account:

- If you have root access to the KDC machine, but you do not have a Kerberos admin account, use `kadmin.local`.
- If you do not have root access to the KDC machine, but you do have a Kerberos admin account, use `kadmin`.
- If you have both root access to the KDC machine and a Kerberos admin account, you can use either one.

To start `kadmin.local` (on the KDC machine) or `kadmin` from any machine, run this command:

```
$ sudo kadmin.local
```

OR:

```
$ kadmin
```

**Note:**

In this guide, `kadmin` is shown as the prompt for commands in the `kadmin` shell, but you can type the same commands at the `kadmin.local` prompt in the `kadmin.local` shell.

**Note:**

Running `kadmin.local` may prompt you for a password because it is being run via `sudo`. You should provide your Unix password. Running `kadmin` may prompt you for a password because you need Kerberos admin privileges. You should provide your Kerberos admin password.

To create the Kerberos principals



Important:

If you plan to use Oozie, Impala, or the Hue Kerberos ticket renewer in your cluster, you must configure your KDC to allow tickets to be renewed, and you must configure `krb5.conf` to request renewable tickets. Typically, you can do this by adding the `max_renewable_life` setting to your realm in `kdc.conf`, and by adding the `renew_lifetime` parameter to the `libdefaults` section of `krb5.conf`. For more information about renewable tickets, see the [Kerberos documentation](#).

Do the following steps for every host in your cluster. Run the commands in the `kadmin.local` or `kadmin` shell, replacing the `fully.qualified.domain.name` in the commands with the fully qualified domain name of each host. Replace `YOUR-REALM.COM` with the name of the Kerberos realm your Hadoop cluster is in.

1. In the `kadmin.local` or `kadmin` shell, create the `hdfs` principal. This principal is used for the NameNode, Secondary NameNode, and DataNodes.

```
kadmin: addprinc -randkey hdfs/fully.qualified.domain.name@YOUR-REALM.COM
```



Note:

If your Kerberos administrator or company has a policy about principal names that does not allow you to use the format shown above, you can work around that issue by configuring the `<kerberos principal>` to `<short name>` mapping that is built into Hadoop. For more information, see [Configuring the Mapping from Kerberos Principals to Short Names](#).

2. Create the `mapred` principal. If you are using MRv1, the `mapred` principal is used for the JobTracker and TaskTrackers. If you are using YARN, the `mapred` principal is used for the MapReduce Job History Server.

```
kadmin: addprinc -randkey mapred/fully.qualified.domain.name@YOUR-REALM.COM
```

3. **YARN only:** Create the `yarn` principal. This principal is used for the ResourceManager and NodeManager.

```
kadmin: addprinc -randkey yarn/fully.qualified.domain.name@YOUR-REALM.COM
```

4. Create the HTTP principal.

```
kadmin: addprinc -randkey HTTP/fully.qualified.domain.name@YOUR-REALM.COM
```



Important:

The HTTP principal must be in the format `HTTP/fully.qualified.domain.name@YOUR-REALM.COM`. The first component of the principal must be the literal string "HTTP". This format is standard for HTTP principals in SPNEGO and is hard-coded in Hadoop. It cannot be deviated from.

To create the Kerberos keytab files

**Important:**

The instructions in this section for creating keytab files require using the Kerberos `norandkey` option in the `xst` command. If your version of Kerberos does not support the `norandkey` option, or if you cannot use `kadmin.local`, then use [these alternate instructions](#) to create appropriate Kerberos keytab files. After using those alternate instructions to create the keytab files, continue with the next section [To deploy the Kerberos keytab files](#).

Do the following steps for every host in your cluster. Run the commands in the `kadmin.local` or `kadmin` shell, replacing the `fully.qualified.domain.name` in the commands with the fully qualified domain name of each host:

1. Create the `hdfs` keytab file that will contain the `hdfs` principal and `HTTP` principal. This keytab file is used for the NameNode, Secondary NameNode, and DataNodes.

```
kadmin: xst -norandkey -k hdfs.keytab hdfs/fully.qualified.domain.name
HTTP/fully.qualified.domain.name
```

2. Create the `mapred` keytab file that will contain the `mapred` principal and `HTTP` principal. If you are using MRv1, the `mapred` keytab file is used for the JobTracker and TaskTrackers. If you are using YARN, the `mapred` keytab file is used for the MapReduce Job History Server.

```
kadmin: xst -norandkey -k mapred.keytab mapred/fully.qualified.domain.name
HTTP/fully.qualified.domain.name
```

3. **YARN only:** Create the `yarn` keytab file that will contain the `yarn` principal and `HTTP` principal. This keytab file is used for the ResourceManager and NodeManager.

```
kadmin: xst -norandkey -k yarn.keytab yarn/fully.qualified.domain.name
HTTP/fully.qualified.domain.name
```

4. Use `klist` to display the keytab file entries; a correctly-created `hdfs` keytab file should look something like this:

```
$ klist -e -k -t hdfs.keytab
Keytab name: WRFILE:hdfs.keytab
slot KVNO Principal
-----
 1   7   HTTP/fully.qualified.domain.name@YOUR-REALM.COM (DES cbc mode with CRC-32)
 2   7   HTTP/fully.qualified.domain.name@YOUR-REALM.COM (Triple DES cbc mode with
HMAC/sha1)
 3   7   hdfs/fully.qualified.domain.name@YOUR-REALM.COM (DES cbc mode with CRC-32)
 4   7   hdfs/fully.qualified.domain.name@YOUR-REALM.COM (Triple DES cbc mode with
HMAC/sha1)
```

5. Continue with the next section [To deploy the Kerberos keytab files](#).

To deploy the Kerberos keytab files

On every node in the cluster, repeat the following steps to deploy the `hdfs.keytab` and `mapred.keytab` files. If you are using YARN, you will also deploy the `yarn.keytab` file.

1. On the host machine, copy or move the keytab files to a directory that Hadoop can access, such as `/etc/hadoop/conf`.

- a. If you are using MRv1:

```
$ sudo mv hdfs.keytab mapred.keytab /etc/hadoop/conf/
```

If you are using YARN:

```
$ sudo mv hdfs.keytab mapred.keytab yarn.keytab /etc/hadoop/conf/
```

- b.** Make sure that the `hdfs.keytab` file is only readable by the `hdfs` user, and that the `mapred.keytab` file is only readable by the `mapred` user.

```
$ sudo chown hdfs:hadoop /etc/hadoop/conf/hdfs.keytab
$ sudo chown mapred:hadoop /etc/hadoop/conf/mapred.keytab
$ sudo chmod 400 /etc/hadoop/conf/*.keytab
```



Note:

To enable you to use the same configuration files on every host, Cloudera recommends that you use the same name for the keytab files on every host.

- c. YARN only:** Make sure that the `yarn.keytab` file is only readable by the `yarn` user.

```
$ sudo chown yarn:hadoop /etc/hadoop/conf/yarn.keytab
$ sudo chmod 400 /etc/hadoop/conf/yarn.keytab
```



Important:

If the NameNode, Secondary NameNode, DataNode, JobTracker, TaskTrackers, HttpFS, or Oozie services are configured to use Kerberos HTTP SPNEGO authentication, and two or more of these services are running on the same host, then all of the running services must use the same HTTP principal and keytab file used for their HTTP endpoints.

Step 5: Shut Down the Cluster

To enable security in CDH, you must stop all Hadoop daemons in your cluster and then change some configuration properties. You must stop all daemons in the cluster because after one Hadoop daemon has been restarted with the configuration properties set to enable security, daemons running without security enabled will be unable to communicate with that daemon. This requirement to shut down all daemons makes it impossible to do a rolling upgrade to enable security on a Hadoop cluster.

To shut down the cluster, run the following command on every node in your cluster (as root):

```
$ for x in `cd /etc/init.d ; ls hadoop-*` ; do sudo service $x stop ; done
```

Step 6: Enable Hadoop Security

Cloudera recommends that all of the Hadoop configuration files throughout the cluster have the same contents.

To enable Hadoop security, add the following properties to the `core-site.xml` file *on every machine* in the cluster:

```
<property>
  <name>hadoop.security.authentication</name>
  <value>kerberos</value> <!-- A value of "simple" would disable security. -->
</property>

<property>
  <name>hadoop.security.authorization</name>
  <value>>true</value>
</property>
```

Enabling Service-Level Authorization for Hadoop Services

The `hadoop-policy.xml` file maintains access control lists (ACL) for Hadoop services. Each ACL consists of comma-separated lists of users and groups separated by a space. For example:

```
user_a,user_b group_a,group_b
```

If you only want to specify a set of users, add a comma-separated list of users followed by a blank space. Similarly, to specify only authorized groups, use a blank space at the beginning. A `*` can be used to give access to all users.

For example, to give users, `ann`, `bob`, and groups, `group_a`, `group_b` access to Hadoop's `DataNodeProtocol` service, modify the `security.datanode.protocol.acl` property in `hadoop-policy.xml`. Similarly, to give all users access to the `InterTrackerProtocol` service, modify `security.inter.tracker.protocol.acl` as follows:

```
<property>
  <name>security.datanode.protocol.acl</name>
  <value>ann,bob group_a,group_b</value>
  <description>ACL for DatanodeProtocol, which is used by datanodes to
  communicate with the namenode.</description>
</property>

<property>
  <name>security.inter.tracker.protocol.acl</name>
  <value>*</value>
  <description>ACL for InterTrackerProtocol, which is used by tasktrackers to
  communicate with the jobtracker.</description>
</property>
```

For more details, see [Service-Level Authorization in Hadoop](#).

Step 7: Configure Secure HDFS

When following the instructions in this section to configure the properties in the `hdfs-site.xml` file, keep the following important guidelines in mind:

- The properties for each daemon (NameNode, Secondary NameNode, and DataNode) must specify both the HDFS and HTTP principals, as well as the path to the HDFS keytab file.
- The Kerberos principals for the NameNode, Secondary NameNode, and DataNode are configured in the `hdfs-site.xml` file. The same `hdfs-site.xml` file with *all three* of these principals must be installed on every host machine in the cluster. That is, it is not sufficient to have the NameNode principal configured on the NameNode host machine only. This is because, for example, the DataNode must know the principal name of the NameNode in order to send heartbeats to it. Kerberos authentication is bi-directional.
- The special string `_HOST` in the properties is replaced at run-time by the fully-qualified domain name of the host machine where the daemon is running. This requires that reverse DNS is properly working on all the hosts configured this way. You may use `_HOST` only as the entirety of the second component of a principal name. For example, `hdfs/_HOST@YOUR-REALM.COM` is valid, but `hdfs._HOST@YOUR-REALM.COM` and `hdfs/_HOST.example.com@YOUR-REALM.COM` are not.
- When performing the `_HOST` substitution for the Kerberos principal names, the NameNode determines its own hostname based on the configured value of `fs.default.name`, whereas the DataNodes determine their hostnames based on the result of reverse DNS resolution on the DataNode hosts. Likewise, the JobTracker uses the configured value of `mapred.job.tracker` to determine its hostname whereas the TaskTrackers, like the DataNodes, use reverse DNS.
- The `dfs.datanode.address` and `dfs.datanode.http.address` port numbers for the DataNode *must* be below 1024, because this provides part of the security mechanism to make it impossible for a user to run a map task which impersonates a DataNode. The port numbers for the NameNode and Secondary NameNode can be anything you want, but the default port numbers are good ones to use.

Configuring Authentication

To configure secure HDFS

Add the following properties to the `hdfs-site.xml` file *on every machine* in the cluster. Replace these example values shown below with the correct settings for your site: *path to the HDFS keytab, YOUR-REALM.COM, fully qualified domain name of NN, and fully qualified domain name of 2NN*

```
<!-- General HDFS security config -->
<property>
  <name>dfs.block.access.token.enable</name>
  <value>true</value>
</property>

<!-- NameNode security config -->
<property>
  <name>dfs.namenode.keytab.file</name>
  <value>/etc/hadoop/conf/hdfs.keytab</value> <!-- path to the HDFS keytab -->
</property>
<property>
  <name>dfs.namenode.kerberos.principal</name>
  <value>hdfs/_HOST@YOUR-REALM.COM</value>
</property>
<property>
  <name>dfs.namenode.kerberos.internal.spnego.principal</name>
  <value>HTTP/_HOST@YOUR-REALM.COM</value>
</property>

<!-- Secondary NameNode security config -->
<property>
  <name>dfs.secondary.namenode.keytab.file</name>
  <value>/etc/hadoop/conf/hdfs.keytab</value> <!-- path to the HDFS keytab -->
</property>
<property>
  <name>dfs.secondary.namenode.kerberos.principal</name>
  <value>hdfs/_HOST@YOUR-REALM.COM</value>
</property>
<property>
  <name>dfs.secondary.namenode.kerberos.internal.spnego.principal</name>
  <value>HTTP/_HOST@YOUR-REALM.COM</value>
</property>

<!-- DataNode security config -->
<property>
  <name>dfs.datanode.data.dir.perm</name>
  <value>700</value>
</property>
<property>
  <name>dfs.datanode.address</name>
  <value>0.0.0.0:1004</value>
</property>
<property>
  <name>dfs.datanode.http.address</name>
  <value>0.0.0.0:1006</value>
</property>
<property>
  <name>dfs.datanode.keytab.file</name>
  <value>/etc/hadoop/conf/hdfs.keytab</value> <!-- path to the HDFS keytab -->
</property>
<property>
  <name>dfs.datanode.kerberos.principal</name>
  <value>hdfs/_HOST@YOUR-REALM.COM</value>
</property>

<!-- Web Authentication config -->
<property>
  <name>dfs.web.authentication.kerberos.principal</name>
  <value>HTTP/_HOST@YOUR_REALM</value>
</property>
```


To enable TLS/SSL for HDFS

Add the following property to `hdfs-site.xml` on *every machine* in your cluster.

```
<property>
<name>dfs.http.policy</name>
<value>HTTPS_ONLY</value>
</property>
```

Optional Step 8: Configuring Security for HDFS High Availability

CDH 5 supports the HDFS High Availability (HA) feature with Kerberos security enabled. There are two use cases that affect security for HA:

- If you are not using Quorum-based Storage (see [Software Configuration for Quorum-based Storage](#)), then no extra configuration for HA is necessary if automatic failover is not enabled. If automatic failover is enabled then access to ZooKeeper should be secured. See the [Software Configuration for Shared Storage Using NFS](#) documentation for details.
- If you are using Quorum-based Storage, then you must configure security for Quorum-based Storage by following the instructions in this section.

To configure security for Quorum-based Storage:

Add the following Quorum-based Storage configuration properties to the `hdfs-site.xml` file on all of the machines in the cluster:

```
<property>
  <name>dfs.journalnode.keytab.file</name>
  <value>/etc/hadoop/conf/hdfs.keytab</value> <!-- path to the HDFS keytab -->
</property>
<property>
  <name>dfs.journalnode.kerberos.principal</name>
  <value>hdfs/_HOST@YOUR-REALM.COM</value>
</property>
<property>
  <name>dfs.journalnode.kerberos.internal.spnego.principal</name>
  <value>HTTP/_HOST@YOUR-REALM.COM</value>
</property>
```

**Note:**

If you already have principals and keytabs created for the machines where the JournalNodes are running, then you should reuse those principals and keytabs in the configuration properties above. You will likely have these principals and keytabs already created if you are collocating a JournalNode on a machine with another HDFS daemon.

Optional Step 9: Configure secure WebHDFS

**Note:**

If you are not using WebHDFS, you can skip this step.

Security for WebHDFS is disabled by default. If you want use WebHDFS with a secure cluster, this is the time to enable and configure it.

To configure secure WebHDFS:

Configuring Authentication

1. If you have not already done so, enable WebHDFS by adding the following property to the `hdfs-site.xml` file *on every machine* in the cluster.

```
<property>
  <name>dfs.webhdfs.enabled</name>
  <value>>true</value>
</property>
```

2. Add the following properties to the `hdfs-site.xml` file *on every machine* in the cluster. Replace the example values shown below with the correct settings for your site.

```
<property>
  <name>dfs.web.authentication.kerberos.principal</name>
  <value>HTTP/_HOST@YOUR-REALM.COM</value>
</property>

<property>
  <name>dfs.web.authentication.kerberos.keytab</name>
  <value>/etc/hadoop/conf/HTTP.keytab</value> <!-- path to the HTTP keytab -->
</property>
```

Optional Step 10: Configuring a secure HDFS NFS Gateway

To deploy a Kerberized HDFS NFS gateway, add the following configuration properties to `hdfs-site.xml` on the NFS server.

```
<property>
  <name>dfs.nfs.keytab.file</name>
  <value>/etc/hadoop/conf/hdfs.keytab</value> <!-- path to the HDFS or NFS gateway keytab -->
</property>

<property>
  <name>dfs.nfs.kerberos.principal</name>
  <value>hdfs/_HOST@YOUR-REALM.COM</value>
</property>
```

Step 11: Set Variables for Secure DataNodes

In order to allow DataNodes to start on a secure Hadoop cluster, you must set the following variables on all DataNodes in `/etc/default/hadoop-hdfs-datanode`.

```
export HADOOP_SECURE_DN_USER=hdfs
export HADOOP_SECURE_DN_PID_DIR=/var/lib/hadoop-hdfs
export HADOOP_SECURE_DN_LOG_DIR=/var/log/hadoop-hdfs
export JSVC_HOME=/usr/lib/bigtop-utils/
```



Note:

Depending on the version of Linux you are using, you may not have the `/usr/lib/bigtop-utils` directory on your system. If that is the case, set the `JSVC_HOME` variable to the `/usr/libexec/bigtop-utils` directory by using this command:

```
export JSVC_HOME=/usr/libexec/bigtop-utils
```

Step 12: Start up the NameNode

You are now ready to start the NameNode. Use the `service` command to run the `/etc/init.d` script.

```
$ sudo service hadoop-hdfs-namenode start
```

You'll see some extra information in the logs such as:

```
10/10/25 17:01:46 INFO security.UserGroupInformation:
Login successful for user hdfs/fully.qualified.domain.name@YOUR-REALM.COM using keytab
file /etc/hadoop/conf/hdfs.keytab
```

and:

```
12/05/23 18:18:31 INFO http.HttpServer: Adding Kerberos (SPNEGO) filter to
getDelegationToken
12/05/23 18:18:31 INFO http.HttpServer: Adding Kerberos (SPNEGO) filter to
renewDelegationToken
12/05/23 18:18:31 INFO http.HttpServer: Adding Kerberos (SPNEGO) filter to
cancelDelegationToken
12/05/23 18:18:31 INFO http.HttpServer: Adding Kerberos (SPNEGO) filter to fsck
12/05/23 18:18:31 INFO http.HttpServer: Adding Kerberos (SPNEGO) filter to getimage
12/05/23 18:18:31 INFO http.HttpServer: Jetty bound to port 50070
12/05/23 18:18:31 INFO mortbay.log: jetty-6.1.26
12/05/23 18:18:31 INFO server.KerberosAuthenticationHandler: Login using keytab
/etc/hadoop/conf/hdfs.keytab, for principal
HTTP/fully.qualified.domain.name@YOUR-REALM.COM
12/05/23 18:18:31 INFO server.KerberosAuthenticationHandler: Initialized, principal
[HTTP/fully.qualified.domain.name@YOUR-REALM.COM] from keytab
[/etc/hadoop/conf/hdfs.keytab]
```

You can verify that the NameNode is working properly by opening a web browser to `http://machine:50070/` where *machine* is the name of the machine where the NameNode is running.

Cloudera also recommends testing that the NameNode is working properly by performing a metadata-only HDFS operation, which will now require correct Kerberos credentials. For example:

```
$ hadoop fs -ls
```

Information about the kinit Command



Important:

Running the `hadoop fs -ls` command will fail if you do not have a valid Kerberos ticket in your credentials cache. You can examine the Kerberos tickets currently in your credentials cache by running the `klist` command. You can obtain a ticket by running the `kinit` command and either specifying a keytab file containing credentials, or entering the password for your principal. If you do not have a valid ticket, you will receive an error such as:

```
11/01/04 12:08:12 WARN ipc.Client: Exception encountered while connecting
to the server : javax.security.sasl.SaslException:
GSS initiate failed [Caused by GSSException: No valid credentials
provided (Mechanism level: Failed to find any Kerberos tgt)]
Bad connection to FS. command aborted. exception: Call to
nn-host/10.0.0.2:8020 failed on local exception: java.io.IOException:
javax.security.sasl.SaslException: GSS initiate failed [Caused by
GSSException: No valid credentials provided (Mechanism level: Failed to
find any Kerberos tgt)]
```



Note:

The `kinit` command must either be on the path for user accounts running the Hadoop client, or else the `hadoop.kerberos.kinit.command` parameter in `core-site.xml` must be manually configured to the absolute path to the `kinit` command.

**Note:**

If you are running MIT Kerberos 1.8.1 or higher, a bug in versions of the Oracle JDK 6 Update 26 and higher causes Java to be unable to read the Kerberos credentials cache even after you have successfully obtained a Kerberos ticket using `kinit`. To work around this bug, run `kinit -R` after running `kinit` initially to obtain credentials. Doing so will cause the ticket to be renewed, and the credentials cache rewritten in a format which Java can read. For more information about this problem, see [Troubleshooting](#).

Step 12: Start up a DataNode

Begin by starting one DataNode only to make sure it can properly connect to the NameNode. Use the `service` command to run the `/etc/init.d` script.

```
$ sudo service hadoop-hdfs-datanode start
```

You'll see some extra information in the logs such as:

```
10/10/25 17:21:41 INFO security.UserGroupInformation:  
Login successful for user hdfs/fully.qualified.domain.name@YOUR-REALM.COM using keytab  
file /etc/hadoop/conf/hdfs.keytab
```

If you can get a single DataNode running and you can see it registering with the NameNode in the logs, then start up all the DataNodes. You should now be able to do all HDFS operations.

Step 14: Set the Sticky Bit on HDFS Directories

This step is optional but strongly recommended for security. In CDH 5, HDFS file permissions have support for the sticky bit. The sticky bit can be set on directories, preventing anyone except the superuser, directory owner, or file owner from deleting or moving the files within the directory. Setting the sticky bit for a file has no effect. This is useful for directories such as `/tmp` which previously had to be set to be world-writable. To set the sticky bit on the `/tmp` directory, run the following command:

```
$ sudo -u hdfs kinit -k -t hdfs.keytab hdfs/fully.qualified.domain.name@YOUR-REALM.COM  
$ sudo -u hdfs hadoop fs -chmod 1777 /tmp
```

After running this command, the permissions on `/tmp` will appear as shown below. (Note the "t" instead of the final "x".)

```
$ hadoop fs -ls /  
Found 2 items  
drwxrwxrwt - hdfs supergroup 0 2011-02-14 15:55 /tmp  
drwxr-xr-x - hdfs supergroup 0 2011-02-14 14:01 /user
```

Step 15: Start up the Secondary NameNode (if used)

At this point, you should be able to start the Secondary NameNode if you are using one:

```
$ sudo service hadoop-hdfs-secondarynamenode start
```

**Note:**

If you are using HDFS HA, do not use the Secondary NameNode. See [Configuring HDFS High Availability](#) for instructions on configuring and deploying the Standby NameNode.

You'll see some extra information in the logs such as:

```
10/10/26 12:03:18 INFO security.UserGroupInformation:
Login successful for user hdfs/fully.qualified.domain.name@YOUR-REALM using keytab file
/etc/hadoop/conf/hdfs.keytab
```

and:

```
12/05/23 18:33:06 INFO http.HttpServer: Adding Kerberos (SPNEGO) filter to getimage
12/05/23 18:33:06 INFO http.HttpServer: Jetty bound to port 50090
12/05/23 18:33:06 INFO mortbay.log: jetty-6.1.26
12/05/23 18:33:06 INFO server.KerberosAuthenticationHandler: Login using keytab
/etc/hadoop/conf/hdfs.keytab, for principal
HTTP/fully.qualified.domain.name@YOUR-REALM.COM
12/05/23 18:33:06 INFO server.KerberosAuthenticationHandler: Initialized, principal
[HTTP/fully.qualified.domain.name@YOUR-REALM.COM] from keytab
[/etc/hadoop/conf/hdfs.keytab]
```

You should make sure that the Secondary NameNode not only starts, but that it is successfully checkpointing.

If you're using the `service` command to start the Secondary NameNode from the `/etc/init.d` scripts, Cloudera recommends setting the property `fs.checkpoint.period` in the `hdfs-site.xml` file to a very low value (such as 5), and then monitoring the Secondary NameNode logs for a successful startup and checkpoint. Once you are satisfied that the Secondary NameNode is checkpointing properly, you should reset the `fs.checkpoint.period` to a reasonable value, or return it to the default, and then restart the Secondary NameNode.

You can make the Secondary NameNode perform a checkpoint by doing the following:

```
$ sudo -u hdfs hdfs secondarynamenode -checkpoint force
```

Note that this will not cause a running Secondary NameNode to checkpoint, but rather will start up a Secondary NameNode that will immediately perform a checkpoint and then shut down. This can be useful for debugging.



Note:

If you encounter errors during Secondary NameNode checkpointing, it may be helpful to enable Kerberos debugging output. For instructions, see [Enabling Debugging Output for the Sun Kerberos Classes](#).

Step 16: Configure Either MRv1 Security or YARN Security

At this point, you are ready to configure either MRv1 Security or YARN Security.

- If you are using MRv1, do the steps in [Configuring MRv1 Security](#) to configure, start, and test secure MRv1.
- If you are using YARN, do the steps in [Configuring YARN Security](#) to configure, start, and test secure YARN.

Configuring MRv1 Security

If you are using YARN, skip this section and see [Configuring YARN Security](#).

If you are using MRv1, do the following steps to configure, start, and test secure MRv1.

1. [Step 1: Configure Secure MRv1](#) on page 109
2. [Step 2: Start up the JobTracker](#) on page 111
3. [Step 3: Start up a TaskTracker](#) on page 111
4. [Step 4: Try Running a Map/Reduce Job](#) on page 111

Step 1: Configure Secure MRv1

Keep the following important information in mind when configuring secure MapReduce:

- The properties for Job Tracker and Task Tracker must specify the mapped principal, as well as the path to the mapped keytab file.

- The Kerberos principals for the Job Tracker and Task Tracker are configured in the `mapred-site.xml` file. The same `mapred-site.xml` file with *both* of these principals must be installed on every host machine in the cluster. That is, it is not sufficient to have the Job Tracker principal configured on the Job Tracker host machine only. This is because, for example, the TaskTracker must know the principal name of the JobTracker in order to securely register with the JobTracker. Kerberos authentication is bi-directional.
- Do not use `${user.name}` in the value of the `mapred.local.dir` or `hadoop.log.dir` properties in `mapred-site.xml`. Doing so can prevent tasks from launching on a secure cluster.
- Make sure that each user who will be running MRv1 jobs exists on all cluster nodes (that is, on every node that hosts any MRv1 daemon).
- Make sure the value specified for `mapred.local.dir` is identical in `mapred-site.xml` and `taskcontroller.cfg`. If the values are different, [this error message](#) is returned.
- Make sure the value specified in `taskcontroller.cfg` for `hadoop.log.dir` is the same as what the Hadoop daemons are using, which is `/var/log/hadoop-0.20-mapreduce` by default and can be configured in `mapred-site.xml`. If the values are different, [this error message](#) is returned.

To configure secure MapReduce:

1. Add the following properties to the `mapred-site.xml` file *on every machine* in the cluster:

```
<!-- JobTracker security configs -->
<property>
  <name>mapreduce.jobtracker.kerberos.principal</name>
  <value>mapred/_HOST@YOUR-REALM.COM</value>
</property>
<property>
  <name>mapreduce.jobtracker.keytab.file</name>
  <value>/etc/hadoop/conf/mapred.keytab</value> <!-- path to the MapReduce keytab -->
</property>

<!-- TaskTracker security configs -->
<property>
  <name>mapreduce.tasktracker.kerberos.principal</name>
  <value>mapred/_HOST@YOUR-REALM.COM</value>
</property>
<property>
  <name>mapreduce.tasktracker.keytab.file</name>
  <value>/etc/hadoop/conf/mapred.keytab</value> <!-- path to the MapReduce keytab -->
</property>

<!-- TaskController settings -->
<property>
  <name>mapred.task.tracker.task-controller</name>
  <value>org.apache.hadoop.mapred.LinuxTaskController</value>
</property>
<property>
  <name>mapreduce.tasktracker.group</name>
  <value>mapred</value>
</property>
```

2. Create a file called `taskcontroller.cfg` that contains the following information:

```
hadoop.log.dir=<Path to Hadoop log directory. Should be same value used to start the
TaskTracker. This is required to set proper permissions on the log files so that they
can be written to by the user's tasks and read by the TaskTracker for serving on the
web UI.>
mapreduce.tasktracker.group=mapred
banned.users=mapred,hdfs,bin
min.user.id=1000
```

**Note:**

In the `taskcontroller.cfg` file, the default setting for the `banned.users` property is `mapred, hdfs, and bin` to prevent jobs from being submitted using those user accounts. The default setting for the `min.user.id` property is `1000` to prevent jobs from being submitted with a user ID less than `1000`, which are conventionally Unix super users. Note that some operating systems such as CentOS 5 use a default value of `500` and above for user IDs, not `1000`. If this is the case on your system, change the default setting for the `min.user.id` property to `500`. If there are user accounts on your cluster that have a user ID less than the value specified for the `min.user.id` property, the TaskTracker returns an error code of `255`.

- The path to the `taskcontroller.cfg` file is determined relative to the location of the `task-controller` binary. Specifically, the path is `<path of task-controller binary>/../..../conf/taskcontroller.cfg`. If you installed the CDH 5 package, this path will always correspond to `/etc/hadoop/conf/taskcontroller.cfg`.

**Note:**

For more information about the `task-controller` program, see [Information about Other Hadoop Security Programs](#).

**Important:**

The same `mapred-site.xml` file and the same `hdfs-site.xml` file must both be installed on every host machine in the cluster so that the NameNode, Secondary NameNode, DataNode, Job Tracker and Task Tracker can all connect securely with each other.

Step 2: Start up the JobTracker

You are now ready to start the JobTracker.

If you're using the `/etc/init.d/hadoop-0.20-mapreduce-jobtracker` script, then you can use the `service` command to run it now:

```
$ sudo service hadoop-0.20-mapreduce-jobtracker start
```

You can verify that the JobTracker is working properly by opening a web browser to `http://machine:50030/` where `machine` is the name of the machine where the JobTracker is running.

Step 3: Start up a TaskTracker

You are now ready to start a TaskTracker.

If you're using the `/etc/init.d/hadoop-0.20-mapreduce-tasktracker` script, then you can use the `service` command to run it now:

```
$ sudo service hadoop-0.20-mapreduce-tasktracker start
```

Step 4: Try Running a Map/Reduce Job

You should now be able to run Map/Reduce jobs. To confirm, try launching a sleep or a pi job from the provided Hadoop examples (`/usr/lib/hadoop-0.20-mapreduce/hadoop-examples.jar`). Note that you will need Kerberos credentials to do so.

**Important:**

Remember that the user who launches the job must exist on every node.

Configuring Authentication

Configuring YARN Security

If you are using MRv1, skip this section and see [Configuring MRv1 Security](#).

If you are using YARN, do the following steps to configure, start, and test secure YARN.

1. [Configure Secure YARN](#).
2. [Start up the ResourceManager](#).
3. [Start up the NodeManager](#).
4. [Start up the MapReduce Job History Server](#).
5. [Try Running a Map/Reduce YARN Job](#).
6. [Step 6: \(Optional\) Configure YARN for Long-running Applications](#) on page 114

Step 1: Configure Secure YARN

Before you start:

- The Kerberos principals for the ResourceManager and NodeManager are configured in the `yarn-site.xml` file. The same `yarn-site.xml` file must be installed on every host machine in the cluster.
- Make sure that each user who will be running YARN jobs exists on all cluster nodes (that is, on every node that hosts any YARN daemon).

To configure secure YARN:

1. Add the following properties to the `yarn-site.xml` file *on every machine* in the cluster:

```
<!-- ResourceManager security configs -->
<property>
  <name>yarn.resourcemanager.keytab</name>
  <value>/etc/hadoop/conf/yarn.keytab</value> <!-- path to the YARN keytab -->
</property>
<property>
  <name>yarn.resourcemanager.principal</name>
  <value>yarn/_HOST@YOUR-REALM.COM</value>
</property>

<!-- NodeManager security configs -->
<property>
  <name>yarn.nodemanager.keytab</name>
  <value>/etc/hadoop/conf/yarn.keytab</value> <!-- path to the YARN keytab -->
</property>
<property>
  <name>yarn.nodemanager.principal</name>
  <value>yarn/_HOST@YOUR-REALM.COM</value>
</property>
<property>
  <name>yarn.nodemanager.container-executor.class</name>
  <value>org.apache.hadoop.yarn.server.nodemanager.LinuxContainerExecutor</value>
</property>
<property>
  <name>yarn.nodemanager.linux-container-executor.group</name>
  <value>yarn</value>
</property>

<!-- To enable TLS/SSL -->
<property>
  <name>yarn.http.policy</name>
  <value>HTTPS_ONLY</value>
</property>
```

2. Add the following properties to the `mapred-site.xml` file *on every machine* in the cluster:

```
<!-- MapReduce Job History Server security configs -->
<property>
  <name>mapreduce.jobhistory.address</name>
  <value>host:port</value> <!-- Host and port of the MapReduce Job History Server;
default port is 10020 -->
</property>
<property>
```



```

    <name>mapreduce.jobhistory.keytab</name>
    <value>/etc/hadoop/conf/mapred.keytab</value> <!-- path to the MAPRED keytab for the
    Job History Server -->
  </property>
</property>
  <name>mapreduce.jobhistory.principal</name>
  <value>mapred/_HOST@YOUR-REALM.COM</value>
</property>

<!-- To enable TLS/SSL -->
<property>
  <name>mapreduce.jobhistory.http.policy</name>
  <value>HTTPS_ONLY</value>
</property>

```

3. Create a file called `container-executor.cfg` for the Linux Container Executor program that contains the following information:

```

yarn.nodemanager.local-dirs=<comma-separated list of paths to local NodeManager
directories. Should be same values specified in yarn-site.xml. Required to validate
paths passed to container-executor in order.>
yarn.nodemanager.linux-container-executor.group=yarn
yarn.nodemanager.log-dirs=<comma-separated list of paths to local NodeManager log
directories. Should be same values specified in yarn-site.xml. Required to set proper
permissions on the log files so that they can be written to by the user's containers
and read by the NodeManager for log aggregation.
banned.users=hdfs,yarn,mapred,bin
min.user.id=1000

```

**Note:**

In the `container-executor.cfg` file, the default setting for the `banned.users` property is `hdfs, yarn, mapred, and bin` to prevent jobs from being submitted via those user accounts. The default setting for the `min.user.id` property is `1000` to prevent jobs from being submitted with a user ID less than 1000, which are conventionally Unix super users. Note that some operating systems such as CentOS 5 use a default value of 500 and above for user IDs, not 1000. If this is the case on your system, change the default setting for the `min.user.id` property to 500. If there are user accounts on your cluster that have a user ID less than the value specified for the `min.user.id` property, the NodeManager returns an error code of 255.

4. The path to the `container-executor.cfg` file is determined relative to the location of the container-executor binary. Specifically, the path is `<dirname of container-executor binary>/../etc/hadoop/container-executor.cfg`. If you installed the CDH 5 package, this path will always correspond to `/etc/hadoop/conf/container-executor.cfg`.

**Note:**

The `container-executor` program requires that the paths including and leading up to the directories specified in `yarn.nodemanager.local-dirs` and `yarn.nodemanager.log-dirs` to be set to 755 permissions as shown in [this table](#) on permissions on directories.

5. Verify that the ownership and permissions of the `container-executor` program corresponds to:

```

---Sr-s--- 1 root yarn 36264 May 20 15:30 container-executor

```

**Note:**

For more information about the Linux Container Executor program, see [Information about Other Hadoop Security Programs](#).

Configuring Authentication

Step 2: Start the ResourceManager

You are now ready to start the ResourceManager.



Note:

Make sure you always start ResourceManager before starting NodeManager.

If you're using the `/etc/init.d/hadoop-yarn-resourcemanager` script, then you can use the `service` command to run it now:

```
$ sudo service hadoop-yarn-resourcemanager start
```

You can verify that the ResourceManager is working properly by opening a web browser to `http://host:8088/` where `host` is the name of the machine where the ResourceManager is running.

Step 3: Start the NodeManager

You are now ready to start the NodeManager.

If you're using the `/etc/init.d/hadoop-yarn-nodemanager` script, then you can use the `service` command to run it now:

```
$ sudo service hadoop-yarn-nodemanager start
```

You can verify that the NodeManager is working properly by opening a web browser to `http://host:8042/` where `host` is the name of the machine where the NodeManager is running.

Step 4: Start the MapReduce Job History Server

You are now ready to start the MapReduce Job History Server.

If you're using the `/etc/init.d/hadoop-mapreduce-historyserver` script, then you can use the `service` command to run it now:

```
$ sudo service hadoop-mapreduce-historyserver start
```

You can verify that the MapReduce JobHistory Server is working properly by opening a web browser to `http://host:19888/` where `host` is the name of the machine where the MapReduce JobHistory Server is running.

Step 5: Try Running a Map/Reduce YARN Job

You should now be able to run Map/Reduce jobs. To confirm, try launching a sleep or a pi job from the provided Hadoop examples (`/usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar`). Note that you will need Kerberos credentials to do so.



Important:

Remember that the user who launches the job must exist on every node.

To try running a MapReduce job using YARN, set the `HADOOP_MAPRED_HOME` environment variable and then submit the job. For example:

```
$ export HADOOP_MAPRED_HOME=/usr/lib/hadoop-mapreduce
$ /usr/bin/hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar pi 10
10000
```

Step 6: (Optional) Configure YARN for Long-running Applications

Long-running applications such as Spark Streaming jobs will need additional configuration since the default settings only allow the `hdfs` user's delegation tokens a maximum lifetime of 7 days which is not always sufficient.

You can work around this by configuring the ResourceManager as a proxy user for the corresponding HDFS NameNode so that the ResourceManager can request new tokens when the existing ones are past their maximum lifetime. YARN will then be able to continue performing localization and log-aggregation on behalf of the `hdfs` user.

Set the following property in `yarn-site.xml` to true:

```
<property>
<name>yarn.resourcemanager.proxy-user-privileges.enabled</name>
<value>>true</value>
</property>
```

Configure the following properties in `core-site.xml` on the HDFS NameNode. You can use a more restrictive configuration by specifying hosts/groups instead of `*` as in the example below.

```
<property>
<name>hadoop.proxyuser.yarn.hosts</name>
<value>*</value>
</property>

<property>
<name>hadoop.proxyuser.yarn.groups</name>
<value>*</value>
</property>
```

FUSE Kerberos Configuration

This section describes how to use [FUSE](#) (Filesystem in Userspace) and CDH with Kerberos security on your Hadoop cluster. FUSE enables you to mount HDFS, which makes HDFS files accessible just as if they were UNIX files.

To use FUSE and CDH with Kerberos security, follow these guidelines:

- For each HDFS user, make sure that there is a UNIX user with the same name. If there isn't, some files in the FUSE mount point will appear to be owned by a non-existent user. Although this is harmless, it can cause confusion.
- When using Kerberos authentication, users must run `kinit` before accessing the FUSE mount point. Failure to do this will result in I/O errors when the user attempts to access the mount point. For security reasons, it is not possible to list the files in the mount point without first running `kinit`.
- When a user runs `kinit`, all processes that run as that user can use the Kerberos credentials. It is not necessary to run `kinit` in the same shell as the process accessing the FUSE mount point.

Using kadmin to Create Kerberos Keytab Files

If your version of Kerberos does not support the Kerberos `-norandkey` option in the `xst` command, or if you must use `kadmin` because you cannot use `kadmin.local`, then you can use the following procedure to create Kerberos keytab files. Using the `-norandkey` option when creating keytabs is optional and a convenience, but it is not required.



Important:

For both MRv1 and YARN deployments: *On every machine in your cluster, there must be a keytab file for the `hdfs` user and a keytab file for the `mapred` user. The `hdfs` keytab file must contain entries for the `hdfs` principal and an `HTTP` principal, and the `mapred` keytab file must contain entries for the `mapred` principal and an `HTTP` principal. On each respective machine, the `HTTP` principal will be the same in both keytab files.*

In addition, for YARN deployments only: *On every machine in your cluster, there must be a keytab file for the `yarn` user. The `yarn` keytab file must contain entries for the `yarn` principal and an `HTTP` principal. On each respective machine, the `HTTP` principal in the `yarn` keytab file will be the same as the `HTTP` principal in the `hdfs` and `mapred` keytab files.*

For instructions, see [To create the Kerberos keytab files](#) on page 116.

**Note:**

These instructions illustrate an example of creating keytab files for MIT Kerberos. If you are using another version of Kerberos, refer to your Kerberos documentation for instructions. You can use either `kadmin` or `kadmin.local` to run these commands.

To create the Kerberos keytab files

Do the following steps for every host in your cluster, replacing the `fully.qualified.domain.name` in the commands with the fully qualified domain name of each host:

1. Create the `hdfs` keytab file, which contains an entry for the `hdfs` principal. This keytab file is used for the NameNode, Secondary NameNode, and DataNodes.

```
$ kadmin
kadmin: xst -k hdfs-unmerged.keytab hdfs/fully.qualified.domain.name
```

2. Create the `mapred` keytab file, which contains an entry for the `mapred` principal. If you are using MRv1, the `mapred` keytab file is used for the JobTracker and TaskTrackers. If you are using YARN, the `mapred` keytab file is used for the MapReduce Job History Server.

```
kadmin: xst -k mapred-unmerged.keytab mapred/fully.qualified.domain.name
```

3. **YARN only:** Create the `yarn` keytab file, which contains an entry for the `yarn` principal. This keytab file is used for the ResourceManager and NodeManager.

```
kadmin: xst -k yarn-unmerged.keytab yarn/fully.qualified.domain.name
```

4. Create the `http` keytab file, which contains an entry for the HTTP principal.

```
kadmin: xst -k http.keytab HTTP/fully.qualified.domain.name
```

5. Use the `ktutil` command to merge the previously-created keytabs:

```
$ ktutil
ktutil: rkt hdfs-unmerged.keytab
ktutil: rkt http.keytab
ktutil: wkt hdfs.keytab
ktutil: clear
ktutil: rkt mapred-unmerged.keytab
ktutil: rkt http.keytab
ktutil: wkt mapred.keytab
ktutil: clear
ktutil: rkt yarn-unmerged.keytab
ktutil: rkt http.keytab
ktutil: wkt yarn.keytab
```

This procedure creates three new files: `hdfs.keytab`, `mapred.keytab` and `yarn.keytab`. These files contain entries for the `hdfs` and HTTP principals, the `mapred` and HTTP principals, and the `yarn` and HTTP principals respectively.

6. Use `klist` to display the keytab file entries. For example, a correctly-created `hdfs` keytab file should look something like this:

```
$ klist -e -k -t hdfs.keytab
Keytab name: WRFILE:hdfs.keytab
slot KVNO Principal
-----
1 7 HTTP/fully.qualified.domain.name@YOUR-REALM.COM (DES cbc mode with CRC-32)
2 7 HTTP/fully.qualified.domain.name@YOUR-REALM.COM (Triple DES cbc mode with
```

```
HMAC/sha1)
 3   7   hdfs/fully.qualified.domain.name@YOUR-REALM.COM (DES cbc mode with CRC-32)

 4   7   hdfs/fully.qualified.domain.name@YOUR-REALM.COM (Triple DES cbc mode with
HMAC/sha1)
```

- To verify that you have performed the merge procedure correctly, make sure you can obtain credentials as both the `hdfs` and `HTTP` principals using the single merged keytab:

```
$ kinit -k -t hdfs.keytab hdfs/fully.qualified.domain.name@YOUR-REALM.COM
$ kinit -k -t hdfs.keytab HTTP/fully.qualified.domain.name@YOUR-REALM.COM
```

If either of these commands fails with an error message such as "kinit: Key table entry not found while getting initial credentials", then something has gone wrong during the merge procedure. Go back to step 1 of this document and verify that you performed all the steps correctly.

- To continue the procedure of configuring Hadoop security in CDH 5, follow the instructions in the section [To deploy the Kerberos keytab files](#).

Configuring the Mapping from Kerberos Principals to Short Names

You configure the mapping from Kerberos principals to short names in the `hadoop.security.auth_to_local` property setting in the `core-site.xml` file. Kerberos has this support natively, and Hadoop's implementation reuses Kerberos's configuration language to specify the mapping.

A mapping consists of a set of rules that are evaluated in the order listed in the `hadoop.security.auth_to_local` property. The first rule that matches a principal name is used to map that principal name to a short name. Any later rules in the list that match the same principal name are ignored.

You specify the mapping rules on separate lines in the `hadoop.security.auth_to_local` property as follows:

```
<property>
  <name>hadoop.security.auth_to_local</name>
  <value>
    RULE:[<principal translation>](<acceptance filter>)<short name substitution>
    RULE:[<principal translation>](<acceptance filter>)<short name substitution>
    DEFAULT
  </value>
</property>
```

Mapping Rule Syntax

To specify a mapping rule, use the prefix string `RULE:` followed by three sections—principal translation, acceptance filter, and short name substitution—described in more detail below. The syntax of a mapping rule is:

```
RULE:[<principal translation>](<acceptance filter>)<short name substitution>
```

Principal Translation

The first section of a rule, `<principal translation>`, performs the matching of the principal name to the rule. If there is a match, the principal translation also does the initial translation of the principal name to a short name. In the `<principal translation>` section, you specify the number of components in the principal name and the pattern you want to use to translate those principal component(s) and realm into a short name. In Kerberos terminology, a principal name is a set of components separated by slash ("/") characters.

The principal translation is composed of two parts that are both specified within "[]" using the following syntax:

```
[<number of components in principal name>:<initial specification of short name>]
```

where:

<number of components in principal name> – This first part specifies the number of components in the principal name (not including the realm) and must be 1 or 2. A value of 1 specifies principal names that have a single component (for

Configuring Authentication

example, `hdfs`), and `2` specifies principal names that have two components (for example, `hdfs/fully.qualified.domain.name`). A principal name that has only one component will only match single-component rules, and a principal name that has two components will only match two-component rules.

<initial specification of short name> – This second part specifies a pattern for translating the principal component(s) and the realm into a short name. The variable `$0` translates the realm, `$1` translates the first component, and `$2` translates the second component.

Here are some examples of principal translation sections. These examples use `atm@YOUR-REALM.COM` and `atm/fully.qualified.domain.name@YOUR-REALM.COM` as principal name inputs:

This Principal Translation	Translates <code>atm@YOUR-REALM.COM</code> into this short name	Translates <code>atm/fully.qualified.domain.name@YOUR-REALM.COM</code> into this short name
<code>[1:\$1@\$0]</code>	<code>atm@YOUR-REALM.COM</code>	Rule does not match ¹
<code>[1:\$1]</code>	<code>atm</code>	Rule does not match ¹
<code>[1:\$1.foo]</code>	<code>atm.foo</code>	Rule does not match ¹
<code>[2:\$1/\$2@\$0]</code>	Rule does not match ²	<code>atm/fully.qualified.domain.name@YOUR-REALM.COM</code>
<code>[2:\$1/\$2]</code>	Rule does not match ²	<code>atm/fully.qualified.domain.name</code>
<code>[2:\$1@\$0]</code>	Rule does not match ²	<code>atm@YOUR-REALM.COM</code>
<code>[2:\$1]</code>	Rule does not match ²	<code>atm</code>

Footnotes:

¹Rule does not match because there are two components in principal name `atm/fully.qualified.domain.name@YOUR-REALM.COM`

²Rule does not match because there is one component in principal name `atm@YOUR-REALM.COM`

Acceptance Filter

The second section of a rule, (`<acceptance filter>`), matches the translated short name from the principal translation (that is, the output from the first section). The acceptance filter is specified in "()" characters and is a standard regular expression. A rule matches only if the specified regular expression matches the entire translated short name from the principal translation. That is, there's an implied `^` at the beginning of the pattern and an implied `$` at the end.

Short Name Substitution

The third and final section of a rule is the (`<short name substitution>`). If there is a match in the second section, the acceptance filter, the (`<short name substitution>`) section does a final translation of the short name from the first section. This translation is a `sed` replacement expression (`s/.../.../g`) that translates the short name from the first section into the final short name string. The short name substitution section is optional. In many cases, it is sufficient to use the first two sections only.

Converting Principal Names to Lowercase

In some organizations, naming conventions result in mixed-case usernames (for example, `John.Doe`) or even uppercase usernames (for example, `JDoe`) in Active Directory or LDAP. This can cause a conflict when the Linux username and HDFS home directory are lowercase.

To convert principal names to lowercase, append `/L` to the rule.

Example Rules

Suppose all of your service principals are either of the form `App.service-name/fully.qualified.domain.name@YOUR-REALM.COM` or

App.service-name@YOUR-REALM.COM, and you want to map these to the short name string service-name. To do this, your rule set would be:

```
<property>
  <name>hadoop.security.auth_to_local</name>
  <value>
    RULE:[1:$1](App\.*s/App\.(.*)/$1/g
    RULE:[2:$1](App\.*s/App\.(.*)/$1/g
    DEFAULT
  </value>
</property>
```

The first \$1 in each rule is a reference to the first component of the full principal name, and the second \$1 is a regular expression back-reference to text that is matched by (.*) .

In the following example, suppose your company's naming scheme for user accounts in Active Directory is FirstnameLastname (for example, JohnDoe), but user home directories in HDFS are /user/firstname.lastname. The following rule set converts user accounts in the CORP.EXAMPLE.COM domain to lowercase.

```
<property>
  <name>hadoop.security.auth_to_local</name>
  <value>
    RULE:[2:$1@$0](HTTP@QCORP.EXAMPLE.COM\E$)s/@\QCORP.EXAMPLE.COM\E$//
    RULE:[1:$1@$0](.*@QCORP.EXAMPLE.COM\E$)s/@\QCORP.EXAMPLE.COM\E$//L
    RULE:[2:$1@$0](.*@QCORP.EXAMPLE.COM\E$)s/@\QCORP.EXAMPLE.COM\E$//L
    DEFAULT
  </value>
</property>
```

In this example, the JohnDoe@CORP.EXAMPLE.COM principal becomes the johndoe HDFS user.

Default Rule

You can specify an optional default rule called DEFAULT (see example above). The default rule reduces a principal name down to its first component only. For example, the default rule reduces the principal names atm@YOUR-REALM.COM or atm/fully.qualified.domain.name@YOUR-REALM.COM down to atm, assuming that the default domain is YOUR-REALM.COM.

The default rule applies only if the principal is in the default realm.

If a principal name does not match any of the specified rules, the mapping for that principal name will fail.

Testing Mapping Rules

You can test mapping rules for a long principal name by running:

```
$ hadoop org.apache.hadoop.security.HadoopKerberosName name1 name2 name3
```

Enabling Debugging Output for the Sun Kerberos Classes

Initially getting a secure Hadoop cluster configured properly can be tricky, especially for those who are not yet familiar with Kerberos. To help with this, it can be useful to enable debugging output for the Sun Kerberos classes. To do so, set the HADOOP_OPTS environment variable to the following:

```
HADOOP_OPTS="-Dsun.security.krb5.debug=true"
```

Flume Authentication

Flume agents have the ability to store data on an HDFS filesystem configured with Hadoop security. The Kerberos system and protocols authenticate communications between clients and services. Hadoop clients include users and MapReduce jobs on behalf of users, and the services include HDFS and MapReduce. Flume acts as a Kerberos principal (user) and needs Kerberos credentials to interact with the Kerberos security-enabled service. Authenticating a user or

Configuring Authentication

a service can be done using a Kerberos keytab file. This file contains a key that is used to obtain a ticket-granting ticket (TGT). The TGT is used to mutually authenticate the client and the service using the Kerberos KDC.

The following sections describe how to use Flume 1.3.x and CDH 5 with Kerberos security on your Hadoop cluster:



Important:

To enable Flume to work with Kerberos security on your Hadoop cluster, make sure you perform the installation and configuration steps in [Configuring Hadoop Security in CDH 5](#).



Note:

These instructions have been tested with CDH 5 and MIT Kerberos 5 only. The following instructions describe an example of how to configure a Flume agent to be a client as the user `flume` to a secure HDFS service. This section does not describe how to secure the communications between Flume agents, which is not currently implemented.

Configuring Flume's Security Properties

Contents:

Writing as a single user for all HDFS sinks in a given Flume agent

The Hadoop services require a three-part principal that has the form of `username/fully.qualified.domain.name@YOUR-REALM.COM`. Cloudera recommends using `flume` as the first component and the fully qualified domain name of the host machine as the second. Assuming that Kerberos and security-enabled Hadoop have been properly configured on the Hadoop cluster itself, you must add the following parameters to the Flume agent's `flume.conf` configuration file, which is typically located at `/etc/flume-ng/conf/flume.conf`:

```
agentName.sinks.sinkName.hdfs.kerberosPrincipal =  
flume/fully.qualified.domain.name@YOUR-REALM.COM  
agentName.sinks.sinkName.hdfs.kerberosKeytab = /etc/flume-ng/conf/flume.keytab
```

where:

`agentName` is the name of the Flume agent being configured, which in this release defaults to the value "agent".
`sinkName` is the name of the HDFS sink that is being configured. The respective sink's `type` must be `HDFS`.

In the previous example, `flume` is the first component of the principal name, `fully.qualified.domain.name` is the second, and `YOUR-REALM.COM` is the name of the Kerberos realm your Hadoop cluster is in. The `/etc/flume-ng/conf/flume.keytab` file contains the keys necessary for `flume/fully.qualified.domain.name@YOUR-REALM.COM` to authenticate with other services.

Flume and Hadoop also provide a simple keyword, `_HOST`, that gets expanded to be the fully qualified domain name of the host machine where the service is running. This allows you to have one `flume.conf` file with the same `hdfs.kerberosPrincipal` value on all of your agent host machines.

```
agentName.sinks.sinkName.hdfs.kerberosPrincipal = flume/_HOST@YOUR-REALM.COM
```

Writing as different users across multiple HDFS sinks in a single Flume agent

In this release, support has been added for secure impersonation of Hadoop users (similar to "sudo" in UNIX). This is implemented in a way similar to how Oozie implements secure user impersonation.

The following steps to set up secure impersonation from Flume to HDFS assume your cluster is configured using Kerberos. (However, impersonation also works on non-Kerberos secured clusters, and Kerberos-specific aspects should be omitted in that case.)

1. Configure Hadoop to allow impersonation. Add the following configuration properties to your `core-site.xml`.

```
<property>
  <name>hadoop.proxyuser.flume.groups</name>
  <value>group1,group2</value>
  <description>Allow the flume user to impersonate any members of group1 and
group2</description>
</property>
<property>
  <name>hadoop.proxyuser.flume.hosts</name>
  <value>host1,host2</value>
  <description>Allow the flume user to connect only from host1 and host2 to impersonate
a user</description>
</property>
```

You can use the wildcard character `*` to enable impersonation of any user from any host. For more information, see [Secure Impersonation](#).

2. Set up a Kerberos keytab for the Kerberos principal and host Flume is connecting to HDFS from. This user must match the Hadoop configuration in the preceding step. For instructions, see [Configuring Hadoop Security in CDH 5](#).
3. Configure the HDFS sink with the following configuration options:
4. `hdfs.kerberosPrincipal` - fully-qualified principal. Note: `_HOST` will be replaced by the hostname of the local machine (only in-between the `/` and `@` characters)
5. `hdfs.kerberosKeytab` - location on the local machine of the keytab containing the user and host keys for the above principal
6. `hdfs.proxyUser` - the proxy user to impersonate

Example snippet (the majority of the HDFS sink configuration options have been omitted):

```
agent.sinks.sink-1.type = HDFS
agent.sinks.sink-1.hdfs.kerberosPrincipal = flume/_HOST@YOUR-REALM.COM
agent.sinks.sink-1.hdfs.kerberosKeytab = /etc/flume-ng/conf/flume.keytab
agent.sinks.sink-1.hdfs.proxyUser = webloggs

agent.sinks.sink-2.type = HDFS
agent.sinks.sink-2.hdfs.kerberosPrincipal = flume/_HOST@YOUR-REALM.COM
agent.sinks.sink-2.hdfs.kerberosKeytab = /etc/flume-ng/conf/flume.keytab
agent.sinks.sink-2.hdfs.proxyUser = applogs
```

In the above example, the flume Kerberos principal impersonates the user `weblogs` in `sink-1` and the user `applogs` in `sink-2`. This will only be allowed if the Kerberos KDC authenticates the specified principal (`flume` in this case), and the if NameNode authorizes impersonation of the specified proxy user by the specified principal.

Limitations

At this time, Flume does not support using multiple Kerberos principals or keytabs in the same agent. Therefore, if you want to create files as multiple users on HDFS, then impersonation must be configured, and exactly one principal must be configured in Hadoop to allow impersonation of all desired accounts. In addition, the same keytab path must be used across all HDFS sinks in the same agent. If you attempt to configure multiple principals or keytabs in the same agent, Flume will emit the following error message:

```
Cannot use multiple kerberos principals in the same agent. Must restart agent to use
new principal or keytab.
```

Configuring Kerberos for Flume Thrift Source and Sink Using Cloudera Manager

The Thrift source can be configured to start in secure mode by enabling Kerberos authentication. To communicate with a secure Thrift source, the Thrift sink should also be operating in secure mode.

1. Open the Cloudera Manager Admin Console and go to the **Flume** service.
2. Click the **Configuration** tab.
3. Select **Scope > Agent**.

4. Select **Category > Main**.
5. Edit the **Configuration File** property and add the Thrift source and sink properties listed in the tables below to the configuration file.

Table 1: Thrift Source Properties

Property	Description
kerberos	Set to <code>true</code> to enable Kerberos authentication. The <code>agent-principal</code> and <code>agent-keytab</code> properties are required for successful authentication. The Thrift source in secure mode, will accept connections only from Thrift sinks that have Kerberos-enabled and are successfully authenticated to the KDC.
agent-principal	The Kerberos principal used by the Thrift Source to authenticate to the KDC.
agent-keytab	The path to the keytab file used by the Thrift Source in combination with the <code>agent-principal</code> to authenticate to the KDC.

Table 2: Thrift Sink Properties

Property	Description
kerberos	Set to <code>true</code> to enable Kerberos authentication. In Kerberos mode, <code>client-principal</code> , <code>client-keytab</code> and <code>server-principal</code> are required for successful authentication and communication to a Kerberos enabled Thrift Source.
client-principal	The principal used by the Thrift Sink to authenticate to the Kerberos KDC.
client-keytab	The path to the keytab file used by the Thrift Sink in combination with the <code>client-principal</code> to authenticate to the KDC.
server-principal	The principal of the Thrift Source to which this Thrift Sink connects.

Make sure you are configuring these properties for **each** Thrift source and sink instance. For example, for agent `a1`, source `r1`, and sink `k1`, you would add the following properties:

```
# Kerberos properties for Thrift source s1
a1.sources.r1.kerberos=true
a1.sources.r1.agent-principal=<source_principal>
a1.sources.r1.agent-keytab=<path/to/source/keytab>

# Kerberos properties for Thrift sink k1
a1.sinks.k1.kerberos=true
a1.sinks.k1.client-principal=<sink_principal>
a1.sinks.k1.client-keytab=<path/to/sink/keytab>
a1.sinks.k1.server-principal=<path/to/source/keytab>
```

Configure these sets of properties for as many instances of the Thrift source and sink as needed to enable Kerberos.

6. Click **Save Changes** to commit the changes.
7. Restart the Flume service.

Configuring Kerberos for Flume Thrift Source and Sink Using the Command Line

The Thrift source can be configured to start in secure mode by enabling Kerberos authentication. To communicate with a secure Thrift source, the Thrift sink should also be operating in secure mode.

The following tables list the properties that must be configured in the `/etc/flume-ng/conf/flume.conf` file to enable Kerberos for Flume's Thrift source and sink instances.

Table 3: Thrift Source Properties

Property	Description
<code>kerberos</code>	Set to <code>true</code> to enable Kerberos authentication. The <code>agent-principal</code> and <code>agent-keytab</code> properties are required for successful authentication. The Thrift source in secure mode, will accept connections only from Thrift sinks that have Kerberos-enabled and are successfully authenticated to the KDC.
<code>agent-principal</code>	The Kerberos principal used by the Thrift Source to authenticate to the KDC.
<code>agent-keytab</code>	The path to the keytab file used by the Thrift Source in combination with the <code>agent-principal</code> to authenticate to the KDC.

Table 4: Thrift Sink Properties

Property	Description
<code>kerberos</code>	Set to <code>true</code> to enable Kerberos authentication. In Kerberos mode, <code>client-principal</code> , <code>client-keytab</code> and <code>server-principal</code> are required for successful authentication and communication to a Kerberos enabled Thrift Source.
<code>client-principal</code>	The principal used by the Thrift Sink to authenticate to the Kerberos KDC.
<code>client-keytab</code>	The path to the keytab file used by the Thrift Sink in combination with the <code>client-principal</code> to authenticate to the KDC.
<code>server-principal</code>	The principal of the Thrift Source to which this Thrift Sink connects.

Make sure you are configuring these properties for **each** Thrift source and sink instance. For example, for agent `a1`, source `r1`, and sink `k1`, you would add the following properties:

```
# Kerberos properties for Thrift source s1
a1.sources.r1.kerberos=true
a1.sources.r1.agent-principal=<source_principal>
a1.sources.r1.agent-keytab=<path/to/source/keytab>

# Kerberos properties for Thrift sink k1
a1.sinks.k1.kerberos=true
a1.sinks.k1.client-principal=<sink_principal>
a1.sinks.k1.client-keytab=<path/to/sink/keytab>
a1.sinks.k1.server-principal=<path/to/source/keytab>
```

Configure these sets of properties for as many instances of the Thrift source and sink as needed to enable Kerberos.

Flume Account Requirements

This section provides an overview of the account and credential requirements for Flume to write to a Kerberized HDFS. Note the distinctions between the Flume agent machine, DataNode machine, and NameNode machine, as well as the `flume` Unix user account versus the `flume` Hadoop/Kerberos user account.

- Each Flume agent machine that writes to HDFS (using a configured HDFS sink) needs a Kerberos principal of the form:

```
flume/fully.qualified.domain.name@YOUR-REALM.COM
```

Configuring Authentication

where `fully.qualified.domain.name` is the fully qualified domain name of the given Flume agent host machine, and `YOUR-REALM.COM` is the Kerberos realm.

- Each Flume agent machine that writes to HDFS does *not* need to have a `flume` Unix user account to write files owned by the `flume` Hadoop/Kerberos user. Only the keytab for the `flume` Hadoop/Kerberos user is required on the Flume agent machine.
- DataNode machines do *not* need Flume Kerberos keytabs and also do *not* need the `flume` Unix user account.
- TaskTracker (MRv1) or NodeManager (YARN) machines need a `flume` Unix user account *if and only if* MapReduce jobs are being run as the `flume` Hadoop/Kerberos user.
- The NameNode machine needs to be able to resolve the groups of the `flume` user. The groups of the `flume` user on the NameNode machine are mapped to the Hadoop groups used for authorizing access.
- The NameNode machine does *not* need a Flume Kerberos keytab.

Testing the Flume HDFS Sink Configuration

To test whether your Flume HDFS sink is properly configured to connect to your secure HDFS cluster, you must run data through Flume. An easy way to do this is to configure a Netcat source, a Memory channel, and an HDFS sink. Start Flume with that configuration, and use the `nc` command (available freely online and with many UNIX distributions) to send events to the Netcat source port. The resulting events should appear on HDFS in the configured location. If the events do not appear, check the Flume log at `/var/log/flume-ng/flume.log` for any error messages related to Kerberos.

Writing to a Secure HBase cluster

If you want to write to a secure HBase cluster, be aware of the following:

- Flume must be configured to use Kerberos security as documented above, and HBase must be configured to use Kerberos security as documented in [HBase Security Configuration](#).
- The `hbase-site.xml` file, which must be configured to use Kerberos security, must be in Flume's classpath or `HBASE_HOME/conf`.
- `org.apache.flume.sink.hbase.HBaseSink` supports secure HBase, but `org.apache.flume.sink.hbase.AsyncHBaseSink` does not.
- The Flume HBase Sink takes these two parameters:
- `kerberosPrincipal` – specifies the Kerberos principal to be used
- `kerberosKeytab` – specifies the path to the Kerberos keytab These are defined as:

```
agent.sinks.hbaseSink.kerberosPrincipal = flume/fully.qualified.domain.name@YOUR-REALM.COM
agent.sinks.hbaseSink.kerberosKeytab = /etc/flume-ng/conf/flume.keytab
```

- If HBase is running with the AccessController coprocessor, the `flume` user (or whichever user the agent is running as) must have permissions to write to the same table and the column family that the sink is configured to write to. You can grant permissions using the `grant` command from HBase shell as explained in [HBase Security Configuration](#).
- The Flume HBase Sink does not currently support impersonation; it will write to HBase as the user the agent is being run as.
- If you want to use HDFS Sink and HBase Sink to write to HDFS and HBase from the same agent respectively, both sinks have to use the same principal and keytab. If you want to use different credentials, the sinks have to be on different agents.
- Each Flume agent machine that writes to HBase (using a configured HBase sink) needs a Kerberos principal of the form:

```
flume/fully.qualified.domain.name@YOUR-REALM.COM
```

where `fully.qualified.domain.name` is the fully qualified domain name of the given Flume agent host machine, and `YOUR-REALM.COM` is the Kerberos realm.

HBase Authentication



Warning: Disabling security on a production HBase system is difficult and could cause data loss. Contact Cloudera Support if you need to disable security in HBase.

To configure HBase security, complete the following tasks:

- 1. Configure HBase Authentication:** You must establish a mechanism for HBase servers and clients to securely identify themselves with HDFS, ZooKeeper, and each other (called *authentication*). This ensures that, for example, a host claiming to be an HBase RegionServer or a particular HBase client are in fact who they claim to be.
- 2. Configure HBase Authorization:** You must establish rules for the resources that clients are allowed to access (called *authorization*). For more information, see [Configuring HBase Authorization](#) on page 432.

For more background information, see this [blog post](#).

The following sections describe how to use Apache HBase and CDH 5 with Kerberos security on your Hadoop cluster:

- [Configuring Kerberos Authentication for HBase](#) on page 125
- [Configuring Secure HBase Replication](#) on page 131
- [Configuring the HBase Client TGT Renewal Period](#) on page 131



Important:

To enable HBase to work with Kerberos security on your Hadoop cluster, make sure you perform the installation and configuration steps in [Configuring Hadoop Security in CDH 5](#) and [ZooKeeper Security Configuration](#).



Note:

These instructions have been tested with CDH and MIT Kerberos 5 only.



Important:

Although an HBase Thrift server can connect to a secured Hadoop cluster, access is not secured from clients to the HBase Thrift server.

Configuring Kerberos Authentication for HBase

Follow these steps to configure HBase authentication.

Configuring Kerberos Authentication for HBase Using Cloudera Manager

Cloudera Manager simplifies the task of configuring Kerberos authentication for HBase.

Configure HBase Servers to Authenticate with a Secure HDFS Cluster Using Cloudera Manager

Minimum Required Role: Security Administrator (also provided by **Full Administrator**)

1. If not done already, [enable Kerberos authentication](#) on your cluster using the Cloudera Manager wizard. This procedure will create principals for the Cloudera Manager server and associated CDH components such as HBase.
2. Go to the HBase service.
3. Click the **Configuration** tab.
4. Select **HBase Service-Wide**.
5. Select **Category > Security**.

Configuring Authentication

6. Ensure the Kerberos principal for the HBase service was generated. See [Viewing and Regenerating Kerberos Principals](#) on page 72.
7. Locate the **HBase Secure Authentication** property or search for it by typing its name in the Search box.
8. Select the **kerberos** option.
9. Click **Save Changes** to commit the changes.
10. Restart the role.
11. Restart the service.

Configure HBase Servers and Clients to Authenticate with a Secure ZooKeeper

To run a secure HBase, you must also use a secure ZooKeeper. To use a secure ZooKeeper, each HBase host machine (Master, RegionServer, and client) must have a principal that allows it to authenticate with your secure ZooKeeper ensemble.

Cloudera Manager automatically configures the following features when you enable Kerberos authentication, and no additional configuration is required.

- Authentication between HBase and ZooKeeper
- Authentication for the HBase REST and Thrift gateways
- `doAs` impersonation for the HBase Thrift gateway



Note:

If you use framed transport, you cannot use `doAs` impersonation, because SASL does not work with Thrift framed transport.

Configuring Kerberos Authentication for HBase Using the Command Line



Important:

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

Configure HBase Servers to Authenticate with a Secure HDFS Cluster Using the Command Line

To configure HBase servers to authenticate with a secure HDFS cluster, you must do the following:

Enabling HBase Authentication

Set the `hbase.security.authentication` property to `kerberos` in `hbase-site.xml` on every host acting as an HBase master, RegionServer, or client. In CDH 5, `hbase.rpc.engine` is automatically detected and does not need to be set.

```
<property>
  <name>hbase.security.authentication</name>
  <value>kerberos</value>
</property>
```

Configuring HBase Kerberos Principals

To run on a secure HDFS cluster, HBase must authenticate itself to the HDFS services. HBase acts as a Kerberos principal and needs Kerberos credentials to interact with the Kerberos-enabled HDFS daemons. You can authenticate a service by using a keytab file, which contains a key that allows the service to authenticate to the Kerberos Key Distribution Center (KDC).

1. Create a service principal for the HBase server using the following syntax. This principal is used to authenticate the HBase server with the HDFS services. Cloudera recommends using `hbase` as the username.

```
$ kadmin
kadmin: addprinc -randkey hbase/fully.qualified.domain.name@YOUR-REALM.COM
```

`fully.qualified.domain.name` is the host where the HBase server is running, and `YOUR-REALM` is the name of your Kerberos realm.

2. Create a keytab file for the HBase server.

```
$ kadmin
kadmin: xst -k hbase.keytab hbase/fully.qualified.domain.name
```

3. Copy the `hbase.keytab` file to the `/etc/hbase/conf` directory on the HBase server host. The owner of the `hbase.keytab` file should be the `hbase` user, and the file should have owner-only read permissions—that is, assign the file `0400` permissions and make it owned by `hbase:hbase`.

```
-r----- 1 hbase hbase 1343 2012-01-09 10:39 hbase.keytab
```

4. To test that the keytab file was created properly, try to obtain Kerberos credentials as the HBase principal using only the keytab file. Substitute your `fully.qualified.domain.name` and realm in the following command:

```
$ kinit -k -t /etc/hbase/conf/hbase.keytab
hbase/fully.qualified.domain.name@YOUR-REALM.COM
```

5. In the `/etc/hbase/conf/hbase-site.xml` configuration file on *all* cluster hosts running the HBase daemon, add the following lines:

```
<property>
  <name>hbase.regionserver.kerberos.principal</name>
  <value>hbase/_HOST@YOUR-REALM.COM</value>
</property>

<property>
  <name>hbase.regionserver.keytab.file</name>
  <value>/etc/hbase/conf/hbase.keytab</value>
</property>

<property>
  <name>hbase.master.kerberos.principal</name>
  <value>hbase/_HOST@YOUR-REALM.COM</value>
</property>

<property>
<name>hbase.master.keytab.file</name>
<value>/etc/hbase/conf/hbase.keytab</value>
</property>
```

Configure HBase Servers and Clients to Authenticate with a Secure ZooKeeper

To run a secure HBase, you must also use a secure ZooKeeper. To use a secure ZooKeeper, each HBase host machine (Master, RegionServer, and client) must have a principal that allows it to authenticate with your secure ZooKeeper ensemble.



Note: This section assumes the following:

- Your secure ZooKeeper is already configured according to the instructions in the [ZooKeeper Security Configuration](#) section and is *not* managed by HBase.
- You have successfully completed the previous steps and already have a principal and keytab file created and in place for every HBase server and client.

Configuring Authentication

Configure HBase JVMs (all Masters, RegionServers, and clients) to Use JAAS

1. On each host, set up a Java Authentication and Authorization Service (JAAS) by creating a `/etc/hbase/conf/zk-jaas.conf` file that contains the following:

```
Client {
  com.sun.security.auth.module.Krb5LoginModule required
  useKeyTab=true
  useTicketCache=false
  keyTab="/etc/hbase/conf/hbase.keytab"
  principal="hbase/fully.qualified.domain.name@<YOUR-REALM>" ;
};
```

2. Modify the `hbase-env.sh` file on HBase server and client hosts to include the following:

```
export HBASE_OPTS="$HBASE_OPTS
-Djava.security.auth.login.config=/etc/hbase/conf/zk-jaas.conf"
export HBASE_MANAGES_ZK=false
```

3. Restart the HBase cluster.

Configure the HBase Servers (Masters and RegionServers) to Use Authentication to Connect to ZooKeeper



Note: If you use Cloudera Manager, this configuration is done automatically and you do not need to perform these steps.

1. Update your `hbase-site.xml` on each HBase server host with the following properties:

```
<configuration>
  <property>
    <name>hbase.zookeeper.quorum</name>
    <value>$ZK_NODES</value>
  </property>
  <property>
    <name>hbase.cluster.distributed</name>
    <value>true</value>
  </property>
</configuration>
```

`$ZK_NODES` is the comma-separated list of hostnames of the ZooKeeper Quorum hosts that you configured according to the instructions in [ZooKeeper Security Configuration](#).

2. Add the following lines to the ZooKeeper configuration file `zoo.cfg`:

```
kerberos.removeHostFromPrincipal=true
kerberos.removeRealmFromPrincipal=true
```

3. Restart ZooKeeper.

Configure Authentication for the HBase REST and Thrift Gateways

By default, the REST gateway does not support impersonation, but accesses HBase as a statically configured user. The actual user who initiated the request is not tracked. With impersonation, the REST gateway user is a proxy user. The HBase server records the actual user who initiates each request and uses this information to apply authorization.



Note: If you use Cloudera Manager, this configuration is done automatically and you do not need to perform these steps.

1. Enable support for proxy users by adding the following properties to `hbase-site.xml`. Substitute the REST gateway proxy user for `$USER`, and the allowed group list for `$GROUPS`.

```
<property>
  <name>hbase.security.authorization</name>
  <value>true</value>
</property>
<property>
  <name>hadoop.proxyuser.$USER.groups</name>
  <value>$GROUPS</value>
</property>
<property>
  <name>hadoop.proxyuser.$USER.hosts</name>
  <value>$GROUPS</value>
</property>
```

2. Enable REST gateway impersonation by adding the following to the `hbase-site.xml` file for every REST gateway:

```
<property>
  <name>hbase.rest.authentication.type</name>
  <value>kerberos</value>
</property>
<property>
  <name>hbase.rest.authentication.kerberos.principal</name>
  <value>HTTP/fully.qualified.domain.name@<YOUR-REALM/>value>
</property>
<property>
  <name>hbase.rest.authentication.kerberos.keytab</name>
  <value>/etc/hbase/conf/hbase.keytab</value>
</property>
```

3. Add the following properties to `hbase-site.xml` for each Thrift gateway, replacing the Kerberos principal with a valid value:

```
<property>
  <name>hbase.thrift.keytab.file</name>
  <value>/etc/hbase/conf/hbase.keytab</value>
</property>
<property>
  <name>hbase.thrift.kerberos.principal</name>
  <value>hbase/fully.qualified.domain.name@<YOUR-REALM/>value>
</property>
<property>
  <name>hbase.thrift.security.qop</name>
  <value>auth</value>
</property>
```

The value for the property `hbase.thrift.security.qop` can be one of the following:

- `auth-conf`—Authentication, integrity, and confidentiality checking
- `auth-int`—Authentication and integrity checking
- `auth`—Authentication checking only

4. To use the Thrift API principal to interact with HBase, add the `hbase.thrift.kerberos.principal` to the `acl` table. For example, to provide administrative access to the Thrift API principal `thrift_server`, run an HBase Shell command like the following:

```
hbase> grant 'thrift_server', 'RWCA'
```

5. Optional: Configure HTTPS transport for Thrift by configuring the following parameters, substituting the placeholders with actual values:

```
<property>
```

```
<name>hbase.thrift.ssl.enabled</name>
<value>true</value>
</property>
<property>
  <name>hbase.thrift.ssl.keystore.store</name>
  <value>LOCATION_OF_KEYSTORE</value>
</property>
<property>
  <name>hbase.thrift.ssl.keystore.password</name>
  <value>KEYSTORE_PASSWORD</value>
</property>
<property>
  <name>hbase.thrift.ssl.keystore.keypassword</name>
  <value>LOCATION_OF_KEYSTORE_KEY_PASSWORD</value>
</property>
```

The Thrift gateway authenticates with HBase using the supplied credential. No authentication is performed by the Thrift gateway itself. All client access through the Thrift gateway uses the gateway's credential, and all clients have its privileges.

Configure `doAs` Impersonation for the HBase Thrift Gateway

**Note:**

If you use framed transport, you cannot use `doAs` impersonation, because SASL does not work with Thrift framed transport.

`doAs` Impersonation provides a flexible way to use the same client to impersonate multiple principals. `doAs` is supported only in Thrift 1, not Thrift 2.

Enable `doAs` support by adding the following properties to `hbase-site.xml` on each Thrift gateway:

```
<property>
  <name>hbase.regionserver.thrift.http</name>
  <value>true</value>
</property>
<property>
  <name>hbase.thrift.support.proxyuser</name>
  <value>true</value>
</property>
```

See the [demo client](#) for information on using `doAs` impersonation in your client applications.

Start HBase

If the configuration worked, you see something similar to the following in the HBase Master and RegionServer logs when you start the cluster:

```
INFO zookeeper.ZooKeeper: Initiating client connection,
connectString=ZK_QUORUM_SERVER:2181 sessionTimeout=180000 watcher=master:60000
INFO zookeeper.ClientCnxn: Opening socket connection to server /ZK_QUORUM_SERVER:2181
INFO zookeeper.RecoverableZooKeeper: The identifier of this process is
PID@ZK_QUORUM_SERVER
INFO zookeeper.Login: successfully logged in.
INFO client.ZooKeeperSaslClient: Client will use GSSAPI as SASL mechanism.
INFO zookeeper.Login: TGT refresh thread started.
INFO zookeeper.ClientCnxn: Socket connection established to ZK_QUORUM_SERVER:2181,
initiating session
INFO zookeeper.Login: TGT valid starting at:          Sun Apr 08 22:43:59 UTC 2012
INFO zookeeper.Login: TGT expires:                  Mon Apr 09 22:43:59 UTC 2012
INFO zookeeper.Login: TGT refresh sleeping until: Mon Apr 09 18:30:37 UTC 2012
INFO zookeeper.ClientCnxn: Session establishment complete on server ZK_QUORUM_SERVER:2181,
sessionid = 0x134106594320000, negotiated timeout = 180000
```

Configuring Secure HBase Replication

If you are using HBase Replication and you want to make it secure, read this section for instructions. Before proceeding, you should already have configured HBase Replication by following the instructions in the [HBase Replication](#) section of the *CDH 5 Installation Guide*.

To configure secure HBase replication, you must configure cross realm support for Kerberos, ZooKeeper, and Hadoop.

To configure secure HBase replication:

1. Create `krbtgt` principals for the two realms. For example, if you have two realms called `ONE.COM` and `TWO.COM`, you need to add the following principals: `krbtgt/ONE.COM@TWO.COM` and `krbtgt/TWO.COM@ONE.COM`. Add these two principals at both realms. Note that there must be at least one common encryption mode between these two realms.

```
kadmin: addprinc -e "<enc_type_list>" krbtgt/ONE.COM@TWO.COM
kadmin: addprinc -e "<enc_type_list>" krbtgt/TWO.COM@ONE.COM
```

2. Add rules for creating short names in Zookeeper. To do this, add a system level property in `java.env`, defined in the `conf` directory. Here is an example rule that illustrates how to add support for the realm called `ONE.COM`, and have two members in the principal (such as `service/instance@ONE.COM`):

```
-Dzookeeper.security.auth_to_local=RULE:[2:\$1@\$0](.*@\QONE.COM\E\$)s/@\QONE.COM\E\$//DEFAULT
```

The above code example adds support for the `ONE.COM` realm in a different realm. So, in the case of replication, you must add a rule for the primary cluster realm in the replica cluster realm. `DEFAULT` is for defining the default rule.

3. Add rules for creating short names in the Hadoop processes. To do this, add the `hadoop.security.auth_to_local` property in the `core-site.xml` file in the replica cluster. For example, to add support for the `ONE.COM` realm:

```
<property>
  <name>hadoop.security.auth_to_local</name>
  <value>
    RULE:[2:\$1@\$0](.*@\QONE.COM\E\$)s/@\QONE.COM\E\$//
    DEFAULT
  </value>
</property>
```

For more information about adding rules, see [Configuring the Mapping from Kerberos Principals to Short Names](#).

Configuring the HBase Client TGT Renewal Period

An HBase client user must also have a Kerberos principal which typically has a password that only the user knows. You should configure the `maxrenewlife` setting for the client's principal to a value that allows the user enough time to finish HBase client processes before the ticket granting ticket (TGT) expires. For example, if the HBase client processes require up to four days to complete, you should create the user's principal and configure the `maxrenewlife` setting by using this command:

```
kadmin: addprinc -maxrenewlife 4days
```

HCatalog Authentication

This section describes how to configure HCatalog in CDH 5 with Kerberos security in a Hadoop cluster:

- [Before You Start](#) on page 132
- [Step 1: Create the HTTP keytab file](#) on page 132
- [Step 2: Configure WebHCat to Use Security](#) on page 132

Configuring Authentication

- [Step 3: Create Proxy Users](#) on page 132
- [Step 4: Verify the Configuration](#) on page 133

For more information about HCatalog see [Installing and Using HCatalog](#).

Before You Start

Secure Web HCatalog requires a running [remote Hive metastore](#) service configured in secure mode. See [Hive MetaStoreServer Security Configuration](#) for instructions. Running secure WebHCat with an [embedded](#) repository is not supported.

Step 1: Create the HTTP keytab file

You need to create a keytab file for WebHCat. Follow these steps:

1. Create the file:

```
kadmin: addprinc -randkey HTTP/fully.qualified.domain.name@YOUR-REALM.COM
kadmin: xst -k HTTP.keytab HTTP/fully.qualified.domain.name
```

2. Move the file into the WebHCat configuration directory and restrict its access exclusively to the hcatalog user:

```
$ mv HTTP.keytab /etc/webhcat/conf/
$ chown hcatalog /etc/webhcat/conf/HTTP.keytab
$ chmod 400 /etc/webhcat/conf/HTTP.keytab
```

Step 2: Configure WebHCat to Use Security

Create or edit the WebHCat configuration file `webhcat-site.xml` in the configuration directory and set following properties:

Property	Value
templeton.kerberos.secret	Any random value
templeton.kerberos.keytab	<code>/etc/webhcat/conf/HTTP.keytab</code>
templeton.kerberos.principal	<code>HTTP/fully.qualified.domain.name@YOUR-REALM.COM</code>

Example configuration:

```
<property>
  <name>templeton.kerberos.secret</name>
  <value>SuPerS3c3tV@lue!</value>
</property>

<property>
  <name>templeton.kerberos.keytab</name>
  <value>/etc/webhcat/conf/HTTP.keytab</value>
</property>

<property>
  <name>templeton.kerberos.principal</name>
  <value>HTTP/fully.qualified.domain.name@YOUR-REALM.COM</value>
</property>
```

Step 3: Create Proxy Users

WebHCat needs access to your NameNode in order to work properly, and so you must configure Hadoop to allow impersonation from the hcatalog user. To do this, edit your `core-site.xml` configuration file and set the `hadoop.proxyuser.HTTP.hosts` and `hadoop.proxyuser.HTTP.groups` properties to specify the hosts from which HCatalog can do the impersonation and what users can be impersonated. You can use the value `*` for "any".

Example configuration:

```
<property>
  <name>hadoop.proxyuser.HTTP.hosts</name>
  <value>*</value>
</property>
<property>
  <name>hadoop.proxyuser.HTTP.groups</name>
  <value>*</value>
</property>
```

Step 4: Verify the Configuration

After restarting WebHcat you can verify that it is working by using `curl` (you may need to run `kinit` first):

```
$ curl --negotiate -i -u :
'http://fully.qualified.domain.name:50111/templeton/v1/ddl/database'
```

Hive Authentication

Here is a summary of the status of Hive authentication in CDH 5:

- HiveServer2 supports authentication of the Thrift client using Kerberos or user/password validation backed by LDAP. For configuration instructions, see [HiveServer2 Security Configuration](#).
- Earlier versions of HiveServer do not support Kerberos authentication for clients. However, the Hive MetaStoreServer does support Kerberos authentication for Thrift clients. For configuration instructions, see [Hive MetaStoreServer Security Configuration](#).

See also: [Using Hive to Run Queries on a Secure HBase Server](#) on page 140

For authorization, Hive uses Apache Sentry to enable role-based, fine-grained authorization for HiveServer2. See [Apache Sentry Overview](#).



Important: Cloudera does not support Apache Ranger or Hive's native authorization frameworks for configuring access control in Hive. Use the Cloudera-supported Apache Sentry instead.

HiveServer2 Security Configuration

HiveServer2 supports authentication of the Thrift client using either of these methods:

- Kerberos authentication
- LDAP authentication

If Kerberos authentication is used, authentication is supported between the Thrift client and HiveServer2, and between HiveServer2 and secure HDFS. If LDAP authentication is used, authentication is supported only between the Thrift client and HiveServer2.



Note: HiveServer2 cannot handle LDAP and Kerberos on the same interface. This is in contrast to when LDAP authentication is enabled on Impala in a Kerberized cluster. With Impala, you have the option to authenticate either with a Kerberos service ticket or with the LDAP username/password. However, this functionality is not available in HiveServer2. When LDAP is enabled on Hive, Kerberos authentication ceases to work.

Enabling Kerberos Authentication for HiveServer2

If you configure HiveServer2 to use Kerberos authentication, HiveServer2 acquires a Kerberos ticket during startup. HiveServer2 requires a principal and keytab file specified in the configuration. Client applications (for example, JDBC or Beeline) must have a valid Kerberos ticket before initiating a connection to HiveServer2.

Configuring Authentication

Configuring HiveServer2 for Kerberos-Secured Clusters

To enable Kerberos Authentication for HiveServer2, add the following properties in the `/etc/hive/conf/hive-site.xml` file:

```
<property>
  <name>hive.server2.authentication</name>
  <value>KERBEROS</value>
</property>
<property>
  <name>hive.server2.authentication.kerberos.principal</name>
  <value>hive/_HOST@YOUR-REALM.COM</value>
</property>
<property>
  <name>hive.server2.authentication.kerberos.keytab</name>
  <value>/etc/hive/conf/hive.keytab</value>
</property>
```

where:

- `hive.server2.authentication` in particular, is a client-facing property that controls the type of authentication HiveServer2 uses for connections to clients. In this case, HiveServer2 uses Kerberos to authenticate incoming clients.
- The `_HOST@YOUR-REALM.COM` value in the example above is the Kerberos principal for the host where HiveServer2 is running. The special string `_HOST` in the properties is replaced at run-time by the fully-qualified domain name of the host machine where the daemon is running. This requires that reverse DNS is properly working on all the hosts configured this way. Replace `YOUR-REALM.COM` with the name of the Kerberos realm your Hadoop cluster is in.
- The `/etc/hive/conf/hive.keytab` value in the example above is a keytab file for that principal.

If you configure HiveServer2 to use both Kerberos authentication and secure impersonation, JDBC clients and Beeline can specify an alternate session user. If these clients have proxy user privileges, HiveServer2 will impersonate the alternate user instead of the one connecting. The alternate user can be specified by the JDBC connection string `proxyUser=userName`

Configuring JDBC Clients for Kerberos Authentication with HiveServer2 (using the Apache Hive driver in Beeline)

JDBC-based clients must include `principal=<hive.server2.authentication.principal>` in the JDBC connection string. For example:

```
String url =
"jdbc:hive2://node1:10000/default;principal=hive/HiveServer2Host@YOUR-REALM.COM"
Connection con = DriverManager.getConnection(url);
```

where `hive` is the principal configured in `hive-site.xml` and `HiveServer2Host` is the host where HiveServer2 is running.

For JDBC clients using the Cloudera JDBC driver, refer the [Cloudera JDBC Driver for Hive](#) documentation. For ODBC clients, refer the [Cloudera ODBC Driver for Apache Hive](#) documentation.

Using Beeline to Connect to a Secure HiveServer2

Use the following command to start `beeline` and connect to a secure HiveServer2 process. In this example, the HiveServer2 process is running on `localhost` at port 10000:

```
$ /usr/lib/hive/bin/beeline
beeline> !connect
jdbc:hive2://localhost:10000/default;principal=hive/HiveServer2Host@YOUR-REALM.COM
0: jdbc:hive2://localhost:10000/default>
```

For more information about the Beeline CLI, see [Using the Beeline CLI](#).

For instructions on encrypting communication with the ODBC/JDBC drivers, see [Configuring Encrypted Communication Between HiveServer2 and Client Drivers](#) on page 235.

Using LDAP Username/Password Authentication with HiveServer2

As an alternative to Kerberos authentication, you can configure HiveServer2 to use user and password validation backed by LDAP. In this case, the client sends a user name and password during the connection initiation. HiveServer2 validates these credentials using an external LDAP service.

You can enable LDAP Authentication with HiveServer2 using Active Directory or OpenLDAP.



Important: When using LDAP username/password authentication with HiveServer2, make sure you have enabled encrypted communication between HiveServer2 and its client drivers to avoid sending cleartext passwords. For instructions, see [Configuring Encrypted Communication Between HiveServer2 and Client Drivers](#) on page 235. Also see [Configuring LDAPS Authentication with HiveServer2](#) on page 136.

Enabling LDAP Authentication with HiveServer2 using Active Directory

- **For managed clusters, use Cloudera Manager:**

1. In the Cloudera Manager Admin Console, click **Hive** in the list of components, and then select the **Configuration** tab.
2. Type "ldap" in the Search text box to locate the LDAP configuration fields.
3. Check **Enable LDAP Authentication**.
4. Enter the **LDAP URL** in the format `ldap[s]://<host>:<port>`
5. Enter the **Active Directory Domain** for your environment.
6. Click **Save Changes**.

- **For unmanaged clusters, use the command line:**

Add the following properties to the `hive-site.xml`:

```
<property>
  <name>hive.server2.authentication</name>
  <value>LDAP</value>
</property>
<property>
  <name>hive.server2.authentication.ldap.url</name>
  <value>LDAP_URL</value>
</property>
<property>
  <name>hive.server2.authentication.ldap.Domain</name>
  <value>AD_DOMAIN_ADDRESS</value>
</property>
```

Where:

The `LDAP_URL` value is the access URL for your LDAP server. For example, `ldap[s]://<host>:<port>`

Enabling LDAP Authentication with HiveServer2 using OpenLDAP

To enable the LDAP mode of authentication using OpenLDAP, include the following properties in the `hive-site.xml` file:

```
<property>
  <name>hive.server2.authentication</name>
  <value>LDAP</value>
</property>
<property>
  <name>hive.server2.authentication.ldap.url</name>
  <value>LDAP_URL</value>
```

Configuring Authentication

```
</property>
<property>
  <name>hive.server2.authentication.ldap.baseDN</name>
  <value>LDAP_BaseDN</value>
</property>
```

where:

- The `LDAP_URL` value is the access URL for your LDAP server.
- The `LDAP_BaseDN` value is the base LDAP DN for your LDAP server. For example, `ou=People,dc=example,dc=com`.

Configuring JDBC Clients for LDAP Authentication with HiveServer2

The JDBC client needs to use a connection URL as shown below. -

JDBC-based clients must include `user=LDAP_Userid;password=LDAP_Password` in the JDBC connection string. For example:

```
String url = "jdbc:hive2://node1:10000/default;user=LDAP_Userid;password=LDAP_Password"
Connection con = DriverManager.getConnection(url);
```

where the `LDAP_Userid` value is the user id and `LDAP_Password` is the password of the client user.

For ODBC Clients, refer the [Cloudera ODBC Driver for Apache Hive](#) documentation.

Enabling LDAP Authentication for HiveServer2 in Hue

LDAP authentication with HiveServer2 can be enabled by setting the following properties under the `[beeswax]` section in `hue.ini`.

<code>auth_username</code>	LDAP username of Hue user to be authenticated.
<code>auth_password</code>	LDAP password of Hue user to be authenticated.

These login details are only used by Hive to authenticate to LDAP. The Hive service trusts Hue to have already validated the user being impersonated, rather than simply passing on the credentials.

Configuring LDAPS Authentication with HiveServer2

HiveServer2 supports [LDAP username/password authentication](#) for clients. Clients send LDAP credentials to HiveServer2 which in turn verifies them against the configured LDAP provider, such as OpenLDAP or Microsoft Active Directory. Most implementations now support LDAPS (LDAP over TLS/SSL), an authentication protocol that uses TLS/SSL to encrypt communication between the LDAP service and its client (in this case, HiveServer2) to avoid sending LDAP credentials in cleartext.

To configure the LDAPS service with HiveServer2:

1. Import the LDAP server CA certificate or the server certificate into a truststore on the HiveServer2 host. If you import the CA certificate, HiveServer2 will trust any server with a certificate issued by the LDAP server's CA. If you only import the server certificate, HiveServer2 trusts only that server. See [Creating Java Keystores and Truststores](#) on page 207 for more details.
2. Make sure the truststore file is readable by the `hive` user.
3. Set the `hive.server2.authentication.ldap.url` configuration property in `hive-site.xml` to the LDAPS URL. For example, `ldaps://sample.myhost.com`.



Note: The URL scheme should be `ldaps` and *not* `ldap`.

- If this is a managed cluster, in Cloudera Manager, go to the Hive service and select **Configuration**. Under Category, select **Security**. In the right panel, search for **HiveServer2 TLS/SSL Certificate Trust Store File**, and add the path to the truststore file that you created in step 1.

If you are using an unmanaged cluster, set the environment variable `HADOOP_OPTS` as follows:

```
HADOOP_OPTS="-Djavax.net.ssl.trustStore=<trustStore-file-path>
-Djavax.net.ssl.trustStorePassword=<trustStore-password>"
```

- Restart HiveServer2.

Pluggable Authentication

Pluggable authentication allows you to provide a custom authentication provider for HiveServer2.

To enable pluggable authentication:

- Set the following properties in `/etc/hive/conf/hive-site.xml`:

```
<property>
  <name>hive.server2.authentication</name>
  <value>CUSTOM</value>
  <description>Client authentication types.
  NONE: no authentication check
  LDAP: LDAP/AD based authentication
  KERBEROS: Kerberos/GSSAPI authentication
  CUSTOM: Custom authentication provider
  (Use with property hive.server2.custom.authentication.class)
</description>
</property>

<property>
  <name>hive.server2.custom.authentication.class</name>
  <value>pluggable-auth-class-name</value>
  <description>
  Custom authentication class. Used when property
  'hive.server2.authentication' is set to 'CUSTOM'. Provided class
  must be a proper implementation of the interface
  org.apache.hive.service.auth.PasswdAuthenticationProvider. HiveServer2
  will call its Authenticate(user, passed) method to authenticate requests.
  The implementation may optionally extend the Hadoop's
  org.apache.hadoop.conf.Configured class to grab Hive's Configuration object.
  </description>
</property>
```

- Make the class available in the CLASSPATH of HiveServer2.

Trusted Delegation with HiveServer2

HiveServer2 determines the identity of the connecting user from the underlying authentication subsystem (Kerberos or LDAP). Any new session started for this connection runs on behalf of this connecting user. If the server is configured to proxy the user at the Hadoop level, then all MapReduce jobs and HDFS accesses will be performed with the identity of the connecting user. If Apache Sentry is configured, then this connecting userid can also be used to verify access rights to underlying tables, views and so on.

In CDH 4.5, a connecting user (for example, `hue`) with Hadoop-level superuser privileges, can request an alternate user for the given session. HiveServer2 will check if the connecting user has Hadoop-level privileges to proxy the requested userid (for example, `bob`). If it does, then the new session will be run on behalf of the alternate user, `bob`, requested by connecting user, `hue`.

To specify an alternate user for new connections, the JDBC client needs to add the `hive.server2.proxy.user=<alternate_user_id>` property to the JDBC connection URL. Note that the connecting

Configuring Authentication

user needs to have Hadoop-level proxy privileges over the alternate user. For example, if user `hue` requests access to run a session as user `bob`, the JDBC connection string should be as follows:

```
# Login as super user Hue
kinit hue -k -t hue.keytab hue@MY-REALM.COM

# Connect using following JDBC connection string
#
jdbc:hive2://myHost.myOrg.com:10000/default;principal=hive/_HOST@MY-REALM.COM;hive.server2.proxy.user=bob
```

HiveServer2 Impersonation



Note: This is not the recommended method to implement HiveServer2 impersonation. Cloudera recommends you use [Sentry](#) to implement this instead.

Impersonation support in HiveServer2 allows users to execute queries and access HDFS files as the connected user rather than the super user who started the HiveServer2 daemon. Impersonation allows admins to enforce an access policy at the file level using HDFS file and directory permissions.

To enable impersonation in HiveServer2:

1. Add the following property to the `/etc/hive/conf/hive-site.xml` file and set the value to `true`. (The default value is `false`.)

```
<property>
  <name>hive.server2.enable.impersonation</name>
  <description>Enable user impersonation for HiveServer2</description>
  <value>true</value>
</property>
```

2. In HDFS or MapReduce configurations, add the following property to the `core-site.xml` file:

```
<property>
  <name>hadoop.proxyuser.hive.hosts</name>
  <value>*</value>
</property>
<property>
  <name>hadoop.proxyuser.hive.groups</name>
  <value>*</value>
</property>
```

See also [File System Permissions](#).

Securing the Hive Metastore



Note: This is not the recommended method to protect the Hive Metastore. Cloudera recommends you use [Sentry](#) to implement this instead.

To prevent users from accessing the Hive metastore and the Hive metastore database using any method other than through HiveServer2, the following actions are recommended:

- Add a firewall rule on the metastore service host to allow access to the metastore port only from the HiveServer2 host. You can do this using [iptables](#).
- Grant access to the metastore database only from the metastore service host. This is specified for MySQL as:

```
GRANT ALL PRIVILEGES ON metastore.* TO 'hive'@'metastorehost';
```

where `metastorehost` is the host where the metastore service is running.

- Make sure users who are not admins cannot log on to the host on which HiveServer2 runs.

Disabling the Hive Security Configuration

Hive's security related metadata is stored in the configuration file `hive-site.xml`. The following sections describe how to disable security for the Hive service.

Disable Client/Server Authentication

To disable client/server authentication, set `hive.server2.authentication` to `NONE`. For example,

```
<property>
  <name>hive.server2.authentication</name>
  <value>NONE</value>
  <description>
    Client authentication types.
    NONE: no authentication check
    LDAP: LDAP/AD based authentication
    KERBEROS: Kerberos/GSSAPI authentication
    CUSTOM: Custom authentication provider
           (Use with property hive.server2.custom.authentication.class)
  </description>
</property>
```

Disable Hive Metastore security

To disable Hive Metastore security, perform the following steps:

- Set the `hive.metastore.sasl.enabled` property to `false` in all configurations, the metastore service side as well as for all clients of the metastore. For example, these might include HiveServer2, Impala, Pig and so on.
- Remove or comment the following parameters in `hive-site.xml` for the metastore service. Note that this is a server-only change.
 - `hive.metastore.kerberos.keytab.file`
 - `hive.metastore.kerberos.principal`

Disable Underlying Hadoop Security

If you also want to disable the underlying Hadoop security, remove or comment out the following parameters in `hive-site.xml`.

- `hive.server2.authentication.kerberos.keytab`
- `hive.server2.authentication.kerberos.principal`

Hive Metastore Server Security Configuration



Important:

This section describes how to configure security for the Hive metastore server. If you are using HiveServer2, see [HiveServer2 Security Configuration](#).

Here is a summary of Hive metastore server security in CDH 5:

- No additional configuration is required to run Hive on top of a security-enabled Hadoop cluster in standalone mode using a local or embedded metastore.
- HiveServer does not support Kerberos authentication for clients. While it is possible to run HiveServer with a secured Hadoop cluster, doing so creates a security hole since HiveServer does not authenticate the Thrift clients that connect to it. Instead, you can use HiveServer2 [HiveServer2 Security Configuration](#).

Configuring Authentication

- The Hive metastore server supports Kerberos authentication for Thrift clients. For example, you can configure a standalone Hive metastore server instance to force clients to authenticate with Kerberos by setting the following properties in the `hive-site.xml` configuration file used by the metastore server:

```
<property>
  <name>hive.metastore.sasl.enabled</name>
  <value>>true</value>
  <description>If true, the metastore thrift interface will be secured with SASL. Clients
  must authenticate with Kerberos.</description>
</property>

<property>
  <name>hive.metastore.kerberos.keytab.file</name>
  <value>/etc/hive/conf/hive.keytab</value>
  <description>The path to the Kerberos Keytab file containing the metastore thrift
  server's service principal.</description>
</property>

<property>
  <name>hive.metastore.kerberos.principal</name>
  <value>hive/_HOST@YOUR-REALM.COM</value>
  <description>The service principal for the metastore thrift server. The special string
  _HOST will be replaced automatically with the correct host name.</description>
</property>
```



Note:

The values shown above for the `hive.metastore.kerberos.keytab.file` and `hive.metastore.kerberos.principal` properties are examples which you will need to replace with the appropriate values for your cluster. Also note that the Hive keytab file should have its access permissions set to 600 and be owned by the same account that is used to run the Metastore server, which is the `hive` user by default.

- Requests to access the metadata are fulfilled by the Hive metastore impersonating the requesting user. This includes read access to the list of databases, tables, properties of each table such as their HDFS location, file type and so on. You can restrict access to the Hive metastore service by allowing it to impersonate only a subset of Kerberos users. This can be done by setting the `hadoop.proxyuser.hive.groups` property in `core-site.xml` on the Hive metastore host.

For example, if you want to give the `hive` user permission to impersonate members of groups `hive` and `user1`:

```
<property>
<name>hadoop.proxyuser.hive.groups</name>
<value>hive,user1</value>
</property>
```

In this example, the Hive metastore can impersonate users belonging to *only* the `hive` and `user1` groups. Connection requests from users not belonging to these groups will be rejected.

Using Hive to Run Queries on a Secure HBase Server

To use Hive to run queries on a secure HBase Server, you must set the following `HIVE_OPTS` environment variable:

```
env HIVE_OPTS="-hiveconf hbase.security.authentication=kerberos -hiveconf
hbase.master.kerberos.principal=hbase/_HOST@YOUR-REALM.COM -hiveconf
hbase.regionserver.kerberos.principal=hbase/_HOST@YOUR-REALM.COM -hiveconf
hbase.zookeeper.quorum=zookeeper1,zookeeper2,zookeeper3" hive
```

where:

- You replace `YOUR-REALM` with the name of your Kerberos realm
- You replace `zookeeper1`, `zookeeper2`, `zookeeper3` with the names of your ZooKeeper servers. The `hbase.zookeeper.quorum` property is configured in the `hbase-site.xml` file.

- The special string `_HOST` is replaced at run-time by the fully-qualified domain name of the host machine where the HBase Master or RegionServer is running. This requires that reverse DNS is properly working on all the hosts configured this way.

In the following, `_HOST` is the name of the host where the HBase Master is running:

```
-hiveconf hbase.master.kerberos.principal=hbase/_HOST@YOUR-REALM.COM
```

In the following, `_HOST` is the hostname of the HBase RegionServer that the application is connecting to:

```
-hiveconf hbase.regionserver.kerberos.principal=hbase/_HOST@YOUR-REALM.COM
```



Note:

You can also set the `HIVE_OPTS` environment variable in your shell profile.

HttpFS Authentication

This section describes how to configure HttpFS CDH 5 with Kerberos security on a Hadoop cluster:

- [Configuring the HttpFS Server to Support Kerberos Security](#) on page 141
- [Using curl to access a URL Protected by Kerberos HTTP SPNEGO](#) on page 142

For more information about HttpFS, see

<https://archive.cloudera.com/cdh5/cdh/5/hadoop/hadoop-hdfs-httpfs/index.html>.



Important:

To enable HttpFS to work with Kerberos security on your Hadoop cluster, make sure you perform the installation and configuration steps in [Configuring Hadoop Security in CDH 5](#).



Important:

If the NameNode, Secondary NameNode, DataNode, JobTracker, TaskTrackers, ResourceManager, NodeManagers, HttpFS, or Oozie services are configured to use Kerberos HTTP SPNEGO authentication, and two or more of these services are running on the same host, then all of the running services must use the same HTTP principal and keytab file used for their HTTP endpoints.

Configuring the HttpFS Server to Support Kerberos Security

1. Create an HttpFS service user principal that is used to authenticate with the Hadoop cluster. The syntax of the principal is: `httpfs/<fully.qualified.domain.name>@<YOUR-REALM>` where: `fully.qualified.domain.name` is the host where the HttpFS server is running `YOUR-REALM` is the name of your Kerberos realm

```
kadmin: addprinc -randkey httpfs/fully.qualified.domain.name@YOUR-REALM.COM
```

2. Create a HTTP service user principal that is used to authenticate user requests coming to the HttpFS HTTP web-services. The syntax of the principal is: `HTTP/<fully.qualified.domain.name>@<YOUR-REALM>` where: `'fully.qualified.domain.name'` is the host where the HttpFS server is running `YOUR-REALM` is the name of your Kerberos realm

```
kadmin: addprinc -randkey HTTP/fully.qualified.domain.name@YOUR-REALM.COM
```



Important:

The HTTP/ component of the HTTP service user principal must be upper case as shown in the syntax and example above.

3. Create keytab files with both principals.

```
$ kadmin
kadmin: xst -k httpfs.keytab httpfs/fully.qualified.domain.name
kadmin: xst -k http.keytab HTTP/fully.qualified.domain.name
```

4. Merge the two keytab files into a single keytab file:

```
$ ktutil
ktutil: rkt httpfs.keytab
ktutil: rkt http.keytab
ktutil: wkt httpfs-http.keytab
```

5. Test that credentials in the merged keytab file work. For example:

```
$ klist -e -k -t httpfs-http.keytab
```

6. Copy the `httpfs-http.keytab` file to the HttpFS configuration directory. The owner of the `httpfs-http.keytab` file should be the `httpfs` user and the file should have owner-only read permissions.

7. Edit the HttpFS server `httpfs-site.xml` configuration file in the HttpFS configuration directory by setting the following properties:

Property	Value
<code>httpfs.authentication.type</code>	<code>kerberos</code>
<code>httpfs.hadoop.authentication.type</code>	<code>kerberos</code>
<code>httpfs.authentication.kerberos.principal</code>	<code>HTTP/<HTTPFS-HOSTNAME>@<YOUR-REALM.COM></code>
<code>httpfs.authentication.kerberos.keytab</code>	<code>/etc/hadoop-httpfs/conf/httpfs-http.keytab</code>
<code>httpfs.hadoop.authentication.kerberos.principal</code>	<code>httpfs/<HTTPFS-HOSTNAME>@<YOUR-REALM.COM></code>
<code>httpfs.hadoop.authentication.kerberos.keytab</code>	<code>/etc/hadoop-httpfs/conf/httpfs-http.keytab</code>
<code>httpfs.authentication.kerberos.name.rules</code>	Use the value configured for 'hadoop.security.auth_to_local' in 'core-site.xml'



Important:

You must restart the HttpFS server to have the configuration changes take effect.

Using curl to access an URL Protected by Kerberos HTTP SPNEGO



Important:

Your version of `curl` must support GSS and be capable of running `curl -V`.

To configure curl to access an URL protected by Kerberos HTTP SPNEGO:

1. Run curl -V:

```
$ curl -V
curl 7.19.7 (universal-apple-darwin10.0) libcurl/7.19.7 OpenSSL/0.9.8l
zlib/1.2.3
Protocols: tftp ftp telnet dict ldap http file https ftps
Features: GSS-Negotiate IPv6 Largefile NTLM SSL libz
```

2. Login to the KDC using kinit.

```
$ kinit
Please enter the password for tucu@LOCALHOST:
```

3. Use curl to fetch the protected URL:

```
$ curl --cacert /path/to/truststore.pem --negotiate -u : -b ~/cookiejar.txt -c
~/cookiejar.txt https://localhost:14000/webhdfs/v1/?op=liststatus
```

where:

- The `--cacert` option is required if you are using TLS/SSL certificates that curl does not recognize by default.
- The `--negotiate` option enables SPNEGO in curl.
- The `-u :` option is required but the username is ignored (the principal that has been specified for `kinit` is used).
- The `-b` and `-c` options are used to store and send HTTP cookies.
- Cloudera does not recommend using the `-k` or `--insecure` option as it turns off curl's ability to verify the certificate.

Hue Authentication

The following sections describe how to configure Kerberos security, enable single sign-on with SAML and encrypt session communication between Hue and other CDH services.

- [Configuring Kerberos Authentication for Hue](#) on page 144
- [Integrating Hue with LDAP](#) on page 147
- [Configuring Hue for SAML](#) on page 151

**Important:**

To enable Hue to work with Kerberos security on your Hadoop cluster, make sure you perform the installation and configuration steps in [Configuring Hadoop Security in CDH 5](#).

Hue Security Enhancements

Enabling LDAP Authentication with HiveServer2 and Impala

LDAP authentication with HiveServer2 and Impala can be enabled by setting the following properties under their respective sections in `hue.ini`, `[beeswax]` and `[impala]`.

<code>auth_username</code>	LDAP username of Hue user to be authenticated.
<code>auth_password</code>	LDAP password of Hue user to be authenticated.

These login details are only used by Impala and Hive to authenticate to LDAP. The Impala and Hive services trust Hue to have already validated the user being impersonated, rather than simply passing on the credentials.

Configuring Authentication

Session Timeout

Session timeouts can be set by specifying the `t1` configuration property under the `[desktop]>[session]` section in `hue.ini`.

<code>t1</code>	The cookie containing the users' session ID will expire after this amount of time in seconds. Default: <code>60*60*24*14</code>
-----------------	----------------------------------------------------------------------------------------------------------------------------------------

Secure Cookies

Secure session cookies can be enable by specifying the `secure` configuration property under the `[desktop]>[session]` section in `hue.ini`. Additionally, you can set the `http-only` flag for cookies containing users' session IDs.

<code>secure</code>	The cookie containing the users' session ID will be secure. Should only be enabled with HTTPS. Default: <code>false</code>
<code>http-only</code>	The cookie containing the users' session ID will use the HTTP only flag. Default: <code>false</code>

Allowed HTTP Methods

You can specify the HTTP request methods that the server should respond to using the `http_allowed_methods` property under the `[desktop]` section in `hue.ini`.

<code>http_allowed_methods</code>	Default: <code>options, get, head, post, put, delete, connect</code>
-----------------------------------	----------------------------------------------------------------------

Restricting the Cipher List

Cipher list support with HTTPS can be restricted by specifying the `ssl_cipher_list` configuration property under the `[desktop]` section in `hue.ini`.

<code>ssl_cipher_list</code>	Default: <code>!aNULL:!eNULL:!LOW:!EXPORT:!SSLv2</code>
------------------------------	---------------------------------------------------------

URL Redirect Whitelist

Restrict the domains or pages to which Hue can redirect users. The `redirect_whitelist` property can be found under the `[desktop]` section in `hue.ini`.

<code>redirect_whitelist</code>	For example, to restrict users to your local domain and FQDN, the following value can be used: <code>^\./.*\$,^http://\./www.mydomain.com\./.*\$</code>
---------------------------------	----------------------------------------------------------------------------------------------------------------------------------------------------------------

Configuring Kerberos Authentication for Hue

You can configure Hue in CDH 5 to support Hadoop security on a cluster using Kerberos.

To configure the Hue server to support Hadoop security using Kerberos:

1. Create a Hue user principal in the same realm as the Hadoop cluster of the form:

```
kadmin: addprinc -randkey hue/hue.server.fully.qualified.domain.name@YOUR-REALM.COM
```


where: `hue` is the principal the Hue server is running as, `hue.server.fully.qualified.domain.name` is the fully-qualified domain name (FQDN) of your Hue server, `YOUR-REALM.COM` is the name of the Kerberos realm your Hadoop cluster is in

2. Create a keytab file for the Hue principal using the same procedure that you used to create the keytab for the `hdfs` or `mapred` principal for a specific host. You should name this file `hue.keytab` and put this keytab file in the directory `/etc/hue` on the machine running the Hue server. Like all keytab files, this file should have the most limited set of permissions possible. It should be owned by the user running the hue server (usually `hue`) and should have the permission `400`.
3. To test that the keytab file was created properly, try to obtain Kerberos credentials as the Hue principal using only the keytab file. Substitute your FQDN and realm in the following command:

```
$ kinit -k -t /etc/hue/hue.keytab
hue/hue.server.fully.qualified.domain.name@YOUR-REALM.COM
```

4. In the `/etc/hue/hue.ini` configuration file, add the following lines in the sections shown. Replace the `kinit_path` value, `/usr/kerberos/bin/kinit`, shown below with the correct path on the user's system.

```
[desktop]

[[kerberos]]
# Path to Hue's Kerberos keytab file
hue_keytab=/etc/hue/hue.keytab
# Kerberos principal name for Hue
hue_principal=hue/FQDN@REALM
# add kinit path for non root users
kinit_path=/usr/kerberos/bin/kinit

[beeswax]
# If Kerberos security is enabled, use fully-qualified domain name (FQDN)
## hive_server_host=<FQDN of Hive Server>
# Hive configuration directory, where hive-site.xml is located
## hive_conf_dir=/etc/hive/conf

[impala]
## server_host=localhost
# The following property is required when impalad and Hue
# are not running on the same host
## impala_principal=impala/impalad.hostname.domainname.com

[search]
# URL of the Solr Server
## solr_url=http://localhost:8983/solr/
# Requires FQDN in solr_url if enabled
## security_enabled=false

[hadoop]

[[hdfs_clusters]]

[[[default]]]
# Enter the host and port on which you are running the Hadoop NameNode
namenode_host=FQDN
hdfs_port=8020
http_port=50070
security_enabled=true

# Thrift plugin port for the name node
## thrift_port=10090

# Configuration for YARN (MR2)
# -----
[[yarn_clusters]]

[[[default]]]
# Enter the host on which you are running the ResourceManager
## resourcemanager_host=localhost
# Change this if your YARN cluster is Kerberos-secured
## security_enabled=false
```

Configuring Authentication

```
# Thrift plug-in port for the JobTracker
## thrift_port=9290

[liboozie]
# The URL where the Oozie service runs on. This is required in order for users to submit
jobs.
## oozie_url=http://localhost:11000/oozie
# Requires FQDN in oozie_url if enabled
## security_enabled=false
```



Important:

In the `/etc/hue/hue.ini` file, verify the following:

- Make sure the `jobtracker_host` property is set to the fully-qualified domain name of the host running the JobTracker. The JobTracker hostname must be fully-qualified in a secured environment.
- Make sure the `fs.defaultfs` property under each `[[hdfs_clusters]]` section contains the fully-qualified domain name of the file system access point, which is typically the NameNode.
- Make sure the `hive_conf_dir` property under the `[beeswax]` section points to a directory containing a valid `hive-site.xml` (either the original or a synced copy).
- Make sure the FQDN specified for HiveServer2 is the same as the FQDN specified for the `hue_principal` configuration property. Without this, HiveServer2 will not work with security enabled.

5. In the `/etc/hadoop/conf/core-site.xml` configuration file on all of your cluster nodes, add the following lines:

```
<!-- Hue security configuration -->
<property>
  <name>hue.kerberos.principal.shortname</name>
  <value>hue</value>
</property>
<property>
  <name>hadoop.proxyuser.hue.groups</name>
  <value>*</value> <!-- A group which all users of Hue belong to, or the wildcard value
  "*" -->
</property>
<property>
  <name>hadoop.proxyuser.hue.hosts</name>
  <value>hue.server.fully.qualified.domain.name</value>
</property>
```



Important:

Make sure you change the `/etc/hadoop/conf/core-site.xml` configuration file on *all* of your cluster nodes.

6. If Hue is configured to communicate to Hadoop using HttpFS, then you must add the following properties to `httpfs-site.xml`:

```
<property>
  <name>httpfs.proxyuser.hue.hosts</name>
  <value>fully.qualified.domain.name</value>
</property>
<property>
  <name>httpfs.proxyuser.hue.groups</name>
  <value>*</value>
</property>
```

7. Add the following properties to the Oozie server `oozie-site.xml` configuration file in the Oozie configuration directory:

```
<property>
  <name>oozie.service.ProxyUserService.proxyuser.hue.hosts</name>
  <value>*</value>
</property>
<property>
  <name>oozie.service.ProxyUserService.proxyuser.hue.groups</name>
  <value>*</value>
</property>
```

8. Restart the JobTracker to load the changes from the `core-site.xml` file.

```
$ sudo service hadoop-0.20-mapreduce-jobtracker restart
```

9. Restart Oozie to load the changes from the `oozie-site.xml` file.

```
$ sudo service oozie restart
```

- 10 Restart the NameNode, JobTracker, and all DataNodes to load the changes from the `core-site.xml` file.

```
$ sudo service hadoop-0.20-(namenode|jobtracker|datanode) restart
```

Integrating Hue with LDAP

When Hue is integrated with LDAP users can use their existing credentials to authenticate and inherit their existing groups transparently. There is no need to save or duplicate any employee password in Hue. There are several other ways to authenticate with Hue such as PAM, SPNEGO, OpenID, OAuth, SAML2 and so on. This topic details how you can configure Hue to authenticate against an LDAP directory server.

When authenticating using LDAP, Hue validates login credentials against an LDAP directory service if configured with the LDAP authentication backend:

```
[desktop]
[[auth]]
backend=desktop.auth.backend.LdapBackend
```

The LDAP authentication backend will automatically create users that don't exist in Hue by default. Hue needs to import users in order to properly perform the authentication. Passwords are never imported when importing users. If you want to disable automatic import set the `create_users_on_login` property under the `[desktop] > [[ldap]]` section of `hue.ini` to `false`.

```
[desktop]
[[ldap]]
create_users_on_login=false
```

The purpose of disabling the automatic import is to allow only a predefined list of manually imported users to login.

There are two ways to authenticate with a directory service through Hue:

- [Search Bind](#)
- [Direct Bind](#)

You can specify the authentication mechanism using the `search_bind_authentication` property under the `[desktop] > [[ldap]]` section of `hue.ini`.

<code>search_bind_authentication</code>	Uses search bind authentication by default. Set this property to <code>false</code> to use direct bind authentication. Default: <code>true</code>
-----------------------------------------	----------------------------------------------------------------------------------------------------------------------------------------------------------

Configuring Authentication

Search Bind

The search bind mechanism for authenticating will perform an [ldapsearch](#) against the directory service and bind using the found [distinguished name](#) (DN) and password provided. This is the default method of authentication used by Hue with LDAP.

The following configuration properties under the `[desktop] > [[ldap]] > [[[users]]]` section in `hue.ini` can be set to restrict the search process.

<code>user_filter</code>	General LDAP filter to restrict the search. Default: <code>"objectclass=*"</code>
<code>user_name_attr</code>	The attribute that will be considered the username to be searched against. Typical attributes to search for include: <code>uid</code> , <code>sAMAccountName</code> . Default: <code>sAMAccountName</code>

With the above configuration, the LDAP search filter will take on the form:

```
((&(objectClass=*)(sAMAccountName=<user entered username>))
```



Important: Setting `search_bind_authentication=true` in `hue.ini` tells Hue to perform an LDAP search using the bind credentials specified for the `bind_dn` and `bind_password` configuration properties. Hue will start searching the subtree starting from the base DN specified for the `base_dn` property. It will then search the base DN for an entry whose attribute, specified in `user_name_attr`, has the same value as the short name provided on login. The search filter, defined in `user_filter` will also be used to limit the search.

Direct Bind

The direct bind mechanism for authenticating will bind to the LDAP server using the username and password provided at login.

The following configuration properties can be used to determine how Hue binds to the LDAP server. These can be set under the `[desktop] > [[ldap]]` section of `hue.ini`.

<code>nt_domain</code>	The NT domain to connect to (only for use with Active Directory). This AD-specific property allows Hue to authenticate with AD without having to follow LDAP references to other partitions. This typically maps to the email address of the user or the user's ID in conjunction with the domain. If provided, Hue will use User Principal Names (UPNs) to bind to the LDAP service. Default: <code>mycompany.com</code>
<code>ldap_username_pattern</code>	Provides a template for the DN that will ultimately be sent to the directory service when authenticating. The <code><username></code> parameter will be replaced with the username provided at login. Default: <code>"uid=<username>,ou=People,dc=mycompany,dc=com"</code>



Important: Setting `search_bind_authentication=false` in `hue.ini` tells Hue to perform a direct bind to LDAP using the credentials provided (*not* `bind_dn` and `bind_password` specified in `hue.ini`). There are two ways direct bind works depending on whether the `nt_domain` property is specified in `hue.ini`:

- `nt_domain` is specified: This is used to connect to an Active Directory service. In this case, the User Principal Name (UPN) is used to perform a direct bind. Hue forms the UPN by concatenating the short name provided at login with the `nt_domain`. For example, `<short name>@<nt_domain>`. The `ldap_username_pattern` property is ignored.
- `nt_domain` is *not* specified: This is used to connect to all other directory services (can handle Active Directory, but `nt_domain` is the preferred way for AD). In this case, `ldap_username_pattern` is used and it should take on the form `cn=<username>,dc=example,dc=com` where `<username>` will be replaced with the username provided at login.

Importing LDAP Users and Groups

If an LDAP user needs to be part of a certain group and be given a particular set of permissions, you can import this user with the User Admin interface in Hue.

The screenshot shows the Hue User Admin interface. The top navigation bar includes 'HUE', 'Query Editors', 'Data Browsers', 'Workflows', 'Search', 'File Browser', 'Job Browser', and 'hdfs'. Below this, the 'User Admin' section is active, with sub-tabs for 'Users', 'Groups', and 'Permissions'. The main content area is titled 'Hue Users - Add/Sync LDAP user' and contains a form with the following fields:

- Username:** A text input field.
- Distinguished name:** A checkbox that is currently unchecked.
- Create home directory:** A checkbox that is currently checked.

At the bottom of the form, there are two buttons: 'Add/Sync user' (in blue) and 'Cancel' (in grey).

Groups can also be imported using the User Admin interface, and users can be added to this group. As in the image below, not only can groups be discovered using DN and [rDN](#) search, but users that are members of the group or members of its subordinate groups can be imported as well.

The screenshot shows the Hue User Admin interface with the 'Groups' tab selected. The page title is 'Hue Groups - Add/Sync LDAP group'. The form includes a 'Name' input field, a 'Distinguished name' checkbox, an 'Import new members' checkbox, an 'Import new members from all subgroups' checkbox, and a 'Create home directories' checkbox which is checked. At the bottom of the form are two buttons: 'Add/Sync group' and 'Cancel'.

You have the following options available when importing a user/group:

- **Distinguished name:** If checked, the username provided must be a full distinguished name (for example, `uid=hue,ou=People,dc=gethue,dc=com`). Otherwise, the **Username** provided should be a fragment of a [Relative Distinguished Name](#) (rDN) (for example, the username `hue` maps to the rDN `uid=hue`). Hue will perform an LDAP search using the same methods and configurations as described above. That is, Hue will take the provided username and create a search filter using the `user_filter` and `user_name_attr` configurations.
- **Create home directory:** If checked, when the user is imported, their home directory in HDFS will automatically be created if it doesn't already exist.



Important: When managing LDAP entries, the User Admin app will always perform an LDAP search and will always use `bind_dn`, `bind_password`, `base_dn`, as defined in `hue.ini`.

Synchronizing LDAP Users and Groups

Users and groups can be synchronized with the directory service using the User Admin interface or using a [command line utility](#). The image from the **Importing LDAP Users and Groups** section uses the words **Add/Sync** to indicate that when a user or group that already exists in Hue is being added, it will in fact be synchronized instead. In the case of importing users for a particular group, new users will be imported and existing users will be synchronized.



Note: Users that have been deleted from the directory service will not be deleted from Hue. Those users can be manually deactivated from Hue using the User Admin interface.

Attributes Synchronized

Currently, only the first name, last name, and email address are synchronized. Hue looks for the LDAP attributes `givenName`, `sn`, and `mail` when synchronizing. The `user_name_attr` configuration property is used to appropriately choose the username in Hue. For instance, if `user_name_attr` is set to `uid`, then the "uid" returned by the directory service will be used as the username of the user in Hue.

User Admin interface

The **Sync LDAP users/groups** button in the User Admin interface will automatically synchronize all users and groups.

Synchronize Using a Command-Line Interface

For example, to synchronize users and groups using a command-line interface:

```
<hue root>/build/env/bin/hue sync_ldap_users_and_groups
```

LDAPS/StartTLS support

Secure communication with LDAP is provided using the TLS/SSL and StartTLS protocols. They allow Hue to validate the directory service it is going to converse with. Hence, if a Certificate Authority certificate file is provided, Hue will communicate using LDAPS. You can specify the path to the CA certificate under :

```
[desktop]
[[ldap]]
ldap_cert=/etc/hue/ca.crt
```

The StartTLS protocol can be used as well:

```
[desktop]
[[ldap]]
use_start_tls=true
```

Configuring Hue for SAML

This section describes the configuration required to use Hue with SAML 2.0 (Security Assertion Markup Language) for single sign-on (SSO) authentication.

The [SAML 2.0 Web Browser SSO](#) profile has three components: a *Service Provider*, a *User Agent* and an *Identity Provider*. In this case, Hue is the Service Provider (SP), you can use an Identity Provider (IdP) of your choice, and your browser, representing you, the user, is the User Agent. When a user requests access to an application, Hue sends an authentication request from the browser to the Identity Provider which then authenticates the user and redirects the browser back to Hue.

Prerequisites

The instructions on this page assume that you have an Identity Provider set up and running. This [blog post](#) guides users through setting up SSO with Hue, using the SAML backend and [Shibboleth](#) as the Identity Provider.

Step 1: Install git, gcc, python-devel, swig, and openssl packages

For example, on RHEL systems, use the following commands:

```
yum install git gcc python-devel swig openssl
```



Important: Users with **CDH 5.5.x** and higher should disable cipher algorithms, MD5, RC4, and DH in the file, `java.security`. See [Troubleshooting](#) below.

Step 2: Install.djangosaml and pysaml2 libraries



Important: CDH 5.4.x and higher include the `djangosaml` and `pysaml2` libraries. Do not install.

The libraries, `djangosaml2` and `pysaml2`, support SAML in Hue. They depend on the `xmlsec1` package to be installed and executable by the user, Hue. First, install the `xmlsec1` package on your system.

RHEL, CentOS and SLES:

For RHEL, CentOS and SLES systems, the `xmlsec1` package is available for download from the EPEL repository. To install packages from the EPEL repository, first download the appropriate the RPM package to your machine, substituting

Configuring Authentication

the version in the package URL with the one required for your system. For example, use the following commands for CentOS 5 or RHEL 5:

```
rpm -Uvh http://download.fedoraproject.org/pub/epel/5/i386/epel-release-5-4.noarch.rpm
yum install xmlsec1
yum install xmlsec1-openssl
```

Oracle Linux:

For Oracle Linux systems, download the `xmlsec1` package from <http://www.aleksey.com/xmlsec/> and execute the following commands:

```
tar -xvzf xmlsec1-<version>.tar.gz
cd xmlsec1-<version>
./configure && make
sudo make install
```



Important: The `xmlsec1` package must be executable by the user, Hue.

You should now be able to install `djangosaml` and `pysaml2` on your machines.

```
build/env/bin/pip install -e git+https://github.com/abec/pysaml2@HEAD#egg=pysaml2
build/env/bin/pip install -e git+https://github.com/abec/djangosaml2@HEAD#egg=djangosaml2
```

Step 3: Update the Hue configuration file

To enable support for SAML, update the necessary parameters in Hue's configuration file, `hue.ini`. This table lists the available SAML parameters in `hue.ini` under the section, `[libsaml]`.

Parameter	Description
<code>xmlsec_binary</code>	Path to the <code>xmlsec_binary</code> , an executable to sign, verify, encrypt, and decrypt SAML requests and assertions. It must be executable by the user, Hue.
<code>create_users_on_login</code>	Boolean, that when True, creates users from OpenId, upon successful login.
<code>required_attributes</code>	Comma-separated list of attributes that Hue requests from the Identity Provider. For example, <code>uid, email</code> , and so on.
<code>optional_attributes</code>	Comma-separated list of optional attributes that Hue requests from the Identity Provider.
<code>metadata_file</code>	Path to the Identity Provider metadata that you copy to a local file. This file should be readable.
<code>key_file</code>	Path to the private key used to encrypt the metadata. File format <code>.PEM</code>
<code>cert_file</code>	Path to the X.509 certificate to be sent along with the encrypted metadata. File format <code>.PEM</code>
<code>user_attribute_mapping</code>	Map of Identity Provider attributes to Hue django user attributes. For example, <code>{'uid':'username', 'email':'email'}</code> .
<code>logout_requests_signed</code>	Boolean, that when True, signs Hue-initiated logout requests and provides a certificate.

For non-parcel based deployments, set the parameter, `redirect_whitelist`, in the `[[desktop]]` section of `hue.ini`, to the fully-qualified domain name of the SAML server so that Hue can redirect the SAML server for authentication. This is not required if you are using a CDH parcel-based deployment managed by Cloudera Manager.

<code>redirect_whitelist</code>	Value: <code>"^\\/.\$,^https:\\/\\/<fqdn_of_SAML_server>\\/. \$"</code>
---------------------------------	--------------------------------------------------------------------------------------

Step 3a: Update the SAML Metadata file

Update the metadata file pointed to by your Hue configuration file, `hue.ini`. Check your Identity Provider documentation for details on how to procure the XML metadata and paste it into the `<metadata_file_name>.xml` file at the location specified by the configuration parameter `metadata_file`.

For example, if your Identity Provider is Shibboleth, you would visit `https://<IdPHOST>:8443/idp/shibboleth`, copy the metadata content available there, and paste it into the Hue metadata file.



Note: You may have to edit the content copied from your Identity Provider's metadata file. For example, the IdP's port number (8443) might be missing from its URL.

Step 3b: Configure Paths to the Private key and Certificate

To enable communication between Hue and the Identity Provider, set `key_file` to the path of the private key and `cert_pem` to the path of the certificate file. The private key signs requests sent to the Identity Provider and the certificate file encrypts and decrypts messages from the Identity Provider.



Note: The key and certificate files specified by the `key_file` and `cert_file` parameters must be .PEM files.

Step 3c: Configure Hue to use SAML Backend

For SAML to allow user logins and create users, update the `backend` configuration property in `hue.ini` to use the SAML authentication backend. The `backend` property is in the `[[auth]]` sub-section under `[desktop]`.

```
backend=libsaml.backend.SAML2Backend
```

Here is an example configuration of the `[libsaml]` section from `hue.ini`.

```
xmlsec_binary=/usr/local/bin/xmlsec1
create_users_on_login=true
metadata_file=/etc/hue/saml/metadata.xml
key_file=/etc/hue/saml/key.pem
cert_file=/etc/hue/saml/cert.pem
logout_requests_signed=true
```

Step 4: Restart the Hue server

Use the following command to restart the Hue server.

```
sudo service hue restart
```

Troubleshooting

- **SSLError:** OpenSSL might fail in CDH 5.5.x and higher with this message:

```
SSLError: [Errno bad handshake] [('SSL routines', 'SSL3_CHECK_CERT_AND_ALGORITHM', 'dh key too small')]
```

To resolve, append the following code to the file, `/usr/java/<your jdk version>-cloudera/jre/lib/security/java.security`:

```
jdk.tls.disabledAlgorithms=MD5, RC4, DH
```

Configuring Authentication

- **Failed to decrypt:** The following error is an indication that you are using a slightly different SAML protocol from what Hue expects:

```
Error: ('failed to decrypt', -1)
```

To resolve: download and run this Python script, `fix-xmlsec.py` and set `xmlsec_binary=<path to script/fix-xmlsec.py>` in `hue.ini`.

This script repairs the known issue whereby `xmlsec1` is not compiled with the `<RetrievalMethod>` and cannot find the location of the encrypted key. SAML2 responses can place `<EncryptedKey>` outside of the `<EncryptedData>` tree and this script simply moves the `<EncryptedKey>` under the `<EncryptedData>`.

Impala Authentication

Authentication is the mechanism to ensure that only specified hosts and users can connect to Impala. It also verifies that when clients connect to Impala, they are connected to a legitimate server. This feature prevents spoofing such as **impersonation** (setting up a phony client system with the same account and group names as a legitimate user) and **man-in-the-middle attacks** (intercepting application requests before they reach Impala and eavesdropping on sensitive information in the requests or the results).

Impala supports authentication using either Kerberos or LDAP.



Note: Regardless of the authentication mechanism used, Impala always creates HDFS directories and data files owned by the same user (typically `impala`). To implement user-level access to different databases, tables, columns, partitions, and so on, use the Sentry authorization feature, as explained in [Enabling Sentry Authorization for Impala](#) on page 411.

Once you are finished setting up authentication, move on to authorization, which involves specifying what databases, tables, HDFS directories, and so on can be accessed by particular users when they connect through Impala. See [Enabling Sentry Authorization for Impala](#) on page 411 for details.

Enabling Kerberos Authentication for Impala

Impala supports Kerberos authentication. For more information on enabling Kerberos authentication, see the topic on Configuring Hadoop Security in the [CDH 5 Security Guide](#).

When using Impala in a managed environment, Cloudera Manager automatically completes Kerberos configuration. In an unmanaged environment, create a Kerberos principal for each host running `impalad` or `statedored`. Cloudera recommends using a consistent format, such as `impala/_HOST@Your-Realm`, but you can use any three-part Kerberos server principal.

In Impala 2.0 and later, `user()` returns the full Kerberos principal string, such as `user@example.com`, in a Kerberized environment.



Note: Regardless of the authentication mechanism used, Impala always creates HDFS directories and data files owned by the same user (typically `impala`). To implement user-level access to different databases, tables, columns, partitions, and so on, use the Sentry authorization feature, as explained in [Enabling Sentry Authorization for Impala](#) on page 411.

An alternative form of authentication you can use is LDAP, described in [Enabling LDAP Authentication for Impala](#) on page 157.

Requirements for Using Impala with Kerberos

On version 5 of Red Hat Enterprise Linux and comparable distributions, some additional setup is needed for the `impala-shell` interpreter to connect to a Kerberos-enabled Impala cluster:

```
sudo yum install python-devel openssl-devel python-pip
sudo pip-python install ssl
```



Important:

- If you plan to use Impala in your cluster, you must configure your KDC to allow tickets to be renewed, and you must configure `krb5.conf` to request renewable tickets. Typically, you can do this by adding the `max_renewable_life` setting to your realm in `kdc.conf`, and by adding the `renew_lifetime` parameter to the `libdefaults` section of `krb5.conf`.

For more information about renewable tickets, see the [Kerberos documentation](#).

- The Impala Web UI does not support Kerberos authentication.
- You cannot use the Impala resource management feature on a cluster that has Kerberos authentication enabled.

Start all `impalad` and `statedored` daemons with the `--principal` and `--keytab-file` flags set to the principal and full path name of the `keytab` file containing the credentials for the principal.

Impala supports the Cloudera ODBC driver and the Kerberos interface provided. To use Kerberos through the ODBC driver, the host type must be set depending on the level of the ODBC driver:

- `SecImpala` for the ODBC 1.0 driver.
- `SecBeeswax` for the ODBC 1.2 driver.
- Blank for the ODBC 2.0 driver or higher, when connecting to a secure cluster.
- `HS2NoSasl` for the ODBC 2.0 driver or higher, when connecting to a non-secure cluster.

To enable Kerberos in the Impala shell, start the `impala-shell` command using the `-k` flag.

To enable Impala to work with Kerberos security on your Hadoop cluster, make sure you perform the installation and configuration steps in [Authentication in the CDH 5 Security Guide](#).

Configuring Impala to Support Kerberos Security

Enabling Kerberos authentication for Impala involves steps that can be summarized as follows:

- Creating service principals for Impala and the HTTP service. Principal names take the form: `serviceName/fully.qualified.domain.name@KERBEROS.REALM`
- Creating, merging, and distributing key tab files for these principals.
- Editing `/etc/default/impala` (in cluster not managed by Cloudera Manager), or editing the **Security** settings in the Cloudera Manager interface, to accommodate Kerberos authentication.

Enabling Kerberos for Impala

1. Create an Impala service principal, specifying the name of the OS user that the Impala daemons run under, the fully qualified domain name of each node running `impalad`, and the realm name. For example:

```
$ kadmin
kadmin: addprinc -requires_preauth -randkey
impala/impala_host.example.com@TEST.EXAMPLE.COM
```

2. Create an HTTP service principal. For example:

```
kadmin: addprinc -randkey HTTP/impala_host.example.com@TEST.EXAMPLE.COM
```



Note: The HTTP component of the service principal must be uppercase as shown in the preceding example.

3. Create keytab files with both principals. For example:

```
kadmin: xst -k impala.keytab impala/impala_host.example.com
kadmin: xst -k http.keytab HTTP/impala_host.example.com
kadmin: quit
```

4. Use `ktutil` to read the contents of the two keytab files and then write those contents to a new file. For example:

```
$ ktutil
ktutil: rkt impala.keytab
ktutil: rkt http.keytab
ktutil: wkt impala-http.keytab
ktutil: quit
```

5. (Optional) Test that credentials in the merged keytab file are valid, and that the “renew until” date is in the future. For example:

```
$ klist -e -k -t impala-http.keytab
```

6. Copy the `impala-http.keytab` file to the Impala configuration directory. Change the permissions to be only read for the file owner and change the file owner to the `impala` user. By default, the Impala user and group are both named `impala`. For example:

```
$ cp impala-http.keytab /etc/impala/conf
$ cd /etc/impala/conf
$ chmod 400 impala-http.keytab
$ chown impala:impala impala-http.keytab
```

7. Add Kerberos options to the Impala defaults file, `/etc/default/impala`. Add the options for both the `impalad` and `statedored` daemons, using the `IMPALA_SERVER_ARGS` and `IMPALA_STATE_STORE_ARGS` variables. For example, you might add:

```
-kerberos_reinit_interval=60
-principal=impala_1/impala_host.example.com@TEST.EXAMPLE.COM
-keytab_file=/var/run/cloudera-scm-agent/process/3212-impala-IMPALAD/impala.keytab
```

For more information on changing the Impala defaults specified in `/etc/default/impala`, see [Modifying Impala Startup Options](#).



Note: Restart `impalad` and `statedored` for these configuration changes to take effect.

Enabling Kerberos for Impala with a Proxy Server

A common configuration for Impala with High Availability is to use a proxy server to submit requests to the actual `impalad` daemons on different hosts in the cluster. This configuration avoids connection problems in case of machine failure, because the proxy server can route new requests through one of the remaining hosts in the cluster. This configuration also helps with load balancing, because the additional overhead of being the “coordinator node” for each query is spread across multiple hosts.

Although you can set up a proxy server with or without Kerberos authentication, typically users set up a secure Kerberized configuration. For information about setting up a proxy server for Impala, including Kerberos-specific steps, see [Using Impala through a Proxy for High Availability](#).

Enabling Impala Delegation for Kerberos Users

See [Configuring Impala Delegation for Hue and BI Tools](#) on page 159 for details about the delegation feature that lets certain users submit queries using the credentials of other users.

Using TLS/SSL with Business Intelligence Tools

You can use Kerberos authentication, TLS/SSL encryption, or both to secure connections from JDBC and ODBC applications to Impala. See [Configuring Impala to Work with JDBC](#) and [Configuring Impala to Work with ODBC](#) for details.

Currently, the Hive JDBC driver does not support connections that use both Kerberos authentication and SSL encryption. To use both of these security features with Impala through a JDBC application, use the [Cloudera JDBC Connector](#) as the JDBC driver.

Enabling LDAP Authentication for Impala

Authentication is the process of allowing only specified named users to access the server (in this case, the Impala server). This feature is crucial for any production deployment, to prevent misuse, tampering, or excessive load on the server. Impala uses LDAP for authentication, verifying the credentials of each user who connects through `impala-shell`, Hue, a Business Intelligence tool, JDBC or ODBC application, and so on.



Note: Regardless of the authentication mechanism used, Impala always creates HDFS directories and data files owned by the same user (typically `impala`). To implement user-level access to different databases, tables, columns, partitions, and so on, use the Sentry authorization feature, as explained in [Enabling Sentry Authorization for Impala](#) on page 411.

An alternative form of authentication you can use is Kerberos, described in [Enabling Kerberos Authentication for Impala](#) on page 154.

Requirements for Using Impala with LDAP

Authentication against LDAP servers is available in Impala 1.2.2 and higher. Impala 1.4.0 adds support for secure LDAP authentication through SSL and TLS.

The Impala LDAP support lets you use Impala with systems such as Active Directory that use LDAP behind the scenes.

Kerberos Authentication for Connections Between Impala Components

Only client->Impala connections can be authenticated by LDAP.

You must use the Kerberos authentication mechanism for connections between internal Impala components, such as between the `impalad`, `statedored`, and `catalogd` daemons. See [Enabling Kerberos Authentication for Impala](#) on page 154 on how to set up Kerberos for Impala.

Server-Side LDAP Setup

These requirements apply on the server side when configuring and starting Impala:

To enable LDAP authentication, set the following startup options for `impalad`:

- `--enable_ldap_auth` enables LDAP-based authentication between the client and Impala.
- `--ldap_uri` sets the URI of the LDAP server to use. Typically, the URI is prefixed with `ldap://`. In Impala 1.4.0 and higher, you can specify secure SSL-based LDAP transport by using the prefix `ldaps://`. The URI can optionally specify the port, for example: `ldap://ldap_server.cloudera.com:389` or `ldaps://ldap_server.cloudera.com:636`. (389 and 636 are the default ports for non-SSL and SSL LDAP connections, respectively.)
- For `ldaps://` connections secured by SSL, `--ldap_ca_certificate="/path/to/certificate.pem"` specifies the location of the certificate in standard `.PEM` format. Store this certificate on the local filesystem, in a location that only the `impala` user and other trusted users can read.

Configuring Authentication

Support for Custom Bind Strings

When Impala connects to LDAP it issues a bind call to the LDAP server to authenticate as the connected user. Impala clients, including the Impala shell, provide the short name of the user to Impala. This is necessary so that Impala can use Sentry for role-based access, which uses short names.

However, LDAP servers often require more complex, structured usernames for authentication. Impala supports three ways of transforming the short name (for example, 'henry') to a more complicated string. If necessary, specify one of the following configuration options when starting the `impalad` daemon on each DataNode:

- `--ldap_domain`: Replaces the username with a string `username@ldap_domain`.
- `--ldap_baseDN`: Replaces the username with a “distinguished name” (DN) of the form: `uid=userid,ldap_baseDN`. (This is equivalent to a Hive option).
- `--ldap_bind_pattern`: This is the most general option, and replaces the username with the string `ldap_bind_pattern` where all instances of the string `#UID` are replaced with `userid`. For example, an `ldap_bind_pattern` of `"user=#UID,OU=foo,CN=bar"` with a username of `henry` will construct a bind name of `"user=henry,OU=foo,CN=bar"`.

For clusters not managed by Cloudera Manager, specify the option on the `impalad` command line. For clusters managed by Cloudera Manager 5.4.0 and higher, search for the configuration field names `ldap_domain`, `ldap_basedn`, or `ldap_bind_pattern`, fill in and save the appropriate field values, and restart the Impala service. Prior to Cloudera Manager 5.4.0, these values were filled in using the **Impala Daemon Command Line Argument Advanced Configuration Snippet (Safety Valve)** field.

These options are mutually exclusive; Impala does not start if more than one of these options is specified.

Secure LDAP Connections

To avoid sending credentials over the wire in cleartext, you must configure a secure connection between both the client and Impala, and between Impala and the LDAP server. The secure connection could use SSL or TLS.

Secure LDAP connections through SSL:

For SSL-enabled LDAP connections, specify a prefix of `ldaps://` instead of `ldap://`. Also, the default port for SSL-enabled LDAP connections is 636 instead of 389.

Secure LDAP connections through TLS:

[TLS](#), the successor to the SSL protocol, is supported by most modern LDAP servers. Unlike SSL connections, TLS connections can be made on the same server port as non-TLS connections. To secure all connections using TLS, specify the following flags as startup options to the `impalad` daemon:

- `--ldap_tls` tells Impala to start a TLS connection to the LDAP server, and to fail authentication if it cannot be done.
- `--ldap_ca_certificate="/path/to/certificate/pem"` specifies the location of the certificate in standard .PEM format. Store this certificate on the local filesystem, in a location that only the `impala` user and other trusted users can read.

LDAP Authentication for `impala-shell` Interpreter

To connect to Impala using LDAP authentication, you specify command-line options to the `impala-shell` command interpreter and enter the password when prompted:

- `-l` enables LDAP authentication.
- `-u` sets the user. Per Active Directory, the user is the short user name, not the full LDAP distinguished name. If your LDAP settings include a search base, use the `--ldap_bind_pattern` on the `impalad` daemon to translate the short user name from `impala-shell` automatically to the fully qualified name.
- `impala-shell` automatically prompts for the password.

For the full list of available `impala-shell` options, see [impala-shell Configuration Options](#).

LDAP authentication for JDBC applications: See [Configuring Impala to Work with JDBC](#) for the format to use with the JDBC connection string for servers using LDAP authentication.

Enabling LDAP for Impala in Hue

Enabling LDAP for Impala in Hue Using Cloudera Manager

1. Go to the Hue service.
2. Click the Configuration tab.
3. Select **Scope > Hue Server**.
4. Select **Category > Advanced**.
5. Add the following properties to the **Hue Server Advanced Configuration Snippet (Safety Valve) for hue_safety_valve_server.ini** property.

```
[impala]
auth_username=<LDAP username of Hue user to be authenticated>
auth_password=<LDAP password of Hue user to be authenticated>
```

6. Click **Save Changes**.

Enabling LDAP for Impala in Hue Using the Command Line

LDAP authentication for the Impala app in Hue can be enabled by setting the following properties under the `[impala]` section in `hue.ini`.

<code>auth_username</code>	LDAP username of Hue user to be authenticated.
<code>auth_password</code>	LDAP password of Hue user to be authenticated.

These login details are only used by Impala to authenticate to LDAP. The Impala service trusts Hue to have already validated the user being impersonated, rather than simply passing on the credentials.

Enabling Impala Delegation for LDAP Users

See [Configuring Impala Delegation for Hue and BI Tools](#) on page 159 for details about the delegation feature that lets certain users submit queries using the credentials of other users.

LDAP Restrictions for Impala

The LDAP support is preliminary. It currently has only been tested against Active Directory.

Using Multiple Authentication Methods with Impala

Impala 2.0 and later automatically handles both Kerberos and LDAP authentication. Each `impalad` daemon can accept both Kerberos and LDAP requests through the same port. No special actions need to be taken if some users authenticate through Kerberos and some through LDAP.

Prior to Impala 2.0, you had to configure each `impalad` to listen on a specific port depending on the kind of authentication, then configure your network load balancer to forward each kind of request to a DataNode that was set up with the appropriate authentication type. Once the initial request was made using either Kerberos or LDAP authentication, Impala automatically handled the process of coordinating the work across multiple nodes and transmitting intermediate results back to the coordinator node.

Configuring Impala Delegation for Hue and BI Tools

When users submit Impala queries through a separate application, such as Hue or a business intelligence tool, typically all requests are treated as coming from the same user. In Impala 1.2 and higher, Impala supports applications to pass along credentials for the users that connect to them, known as “delegation”, and to issue Impala queries with the privileges for those users. Currently, the delegation feature is available only for Impala queries submitted through application interfaces such as Hue and BI tools; for example, Impala cannot issue queries using the privileges of the HDFS user.

Configuring Authentication

The delegation feature is enabled by a startup option for `impalad`: `--authorized_proxy_user_config`. When you specify this option, users whose names you specify (such as `hue`) can delegate the execution of a query to another user. The query runs with the privileges of the delegated user, not the original user such as `hue`. The name of the delegated user is passed using the HiveServer2 configuration property `impala.doas.user`.

You can specify a list of users that the application user can delegate to, or `*` to allow a superuser to delegate to any other user. For example:

```
impalad --authorized_proxy_user_config 'hue=user1,user2;admin=*' ...
```



Note: Make sure to use single quotes or escape characters to ensure that any `*` characters do not undergo wildcard expansion when specified in command-line arguments.

See [Modifying Impala Startup Options](#) for details about adding or changing `impalad` startup options. See [this Cloudera blog post](#) for background information about the delegation capability in HiveServer2.

To set up authentication for the delegated users:

- On the server side, configure either user/password authentication through LDAP, or Kerberos authentication, for all the delegated users. See [Enabling LDAP Authentication for Impala](#) on page 157 or [Enabling Kerberos Authentication for Impala](#) on page 154 for details.
- On the client side, follow the instructions in the “Using User Name and Password” section in the [ODBC driver installation guide](#). Then search for “delegation” in that same installation guide to learn about the **Delegation UID** field and `DelegationUID` configuration keyword to enable the delegation feature for ODBC-based BI tools.

Enabling Delegation in Cloudera Manager

To enable delegation in Cloudera Manager:

1. Navigate to **Clusters > Impala > Configuration > Policy File-Based Sentry**.
2. In the **Proxy User Configuration** field, type the a semicolon-separated list of key=value pairs of authorized proxy users to the user(s) they can impersonate. The list of delegated users are delimited with a comma, e.g. **hue=user1, user2**.
3. Click **Save Changes** and then restart Impala service.

Llama Authentication



Note: Though Impala can be used together with YARN via simple configuration of Static Service Pools in Cloudera Manager, the use of the general-purpose component Llama for integrated resource management within YARN is no longer supported with CDH 5.5 / Impala 2.3 and higher.

This section describes how to configure Llama in CDH 5 with Kerberos security in a Hadoop cluster.



Note: Llama has been tested only in a Cloudera Manager deployment. For information on using Cloudera Manager to configure Llama and Impala, see [Installing Impala](#).

Configuring Llama to Support Kerberos Security

1. Create a Llama service user principal using the syntax: `llama/fully.qualified.domain.name@YOUR-REALM`. This principal is used to authenticate with the Hadoop cluster, where `fully.qualified.domain.name` is the host where Llama is running and `YOUR-REALM` is the name of your Kerberos realm:

```
$ kadmin
kadmin: addprinc -randkey
llama/fully.qualified.domain.name@YOUR-REALM
```

2. Create a keytab file with the Llama principal:

```
$ kadmin
kadmin: xst -k llama.keytab llama/fully.qualified.domain.name
```

3. Test that the credentials in the keytab file work. For example:

```
$ klist -e -k -t llama.keytab
```

4. Copy the `llama.keytab` file to the Llama configuration directory. The owner of the `llama.keytab` file should be the `llama` user and the file should have owner-only read permissions.
5. Edit the Llama `llama-site.xml` configuration file in the Llama configuration directory by setting the following properties:

Property	Value
<code>llama.am.server.thrift.security</code>	<code>true</code>
<code>llama.am.server.thrift.kerberos.keytab.file</code>	<code>llama/conf.keytab</code>
<code>llama.am.server.thrift.kerberos.server.principal.name</code>	<code>llama/fully.qualified.domain.name</code>
<code>llama.am.server.thrift.kerberos.notification.principal.name</code>	<code>impala</code>

6. Restart Llama to make the configuration changes take effect.

Oozie Authentication

This section describes how to configure Oozie CDH 5 with Kerberos security on a Hadoop cluster:

- [Configuring Kerberos Authentication for the Oozie Server](#) on page 162
- [Configuring Oozie HA with Kerberos](#) on page 163



Important:

To enable Oozie to work with Kerberos security on your Hadoop cluster, make sure you perform the installation and configuration steps in [Configuring Hadoop Security in CDH 5](#). Also note that when Kerberos security is enabled in Oozie, a web browser that supports Kerberos HTTP SPNEGO is required to access the Oozie web-console (for example, Firefox, Internet Explorer or Chrome).



Important:

If the NameNode, Secondary NameNode, DataNode, JobTracker, TaskTrackers, ResourceManager, NodeManagers, HttpFS, or Oozie services are configured to use Kerberos HTTP SPNEGO authentication, and two or more of these services are running on the same host, then all of the running services must use the same HTTP principal and keytab file used for their HTTP endpoints.

Configuring Kerberos Authentication for the Oozie Server

1. Create a Oozie service user principal using the syntax:

`oozie/<fully.qualified.domain.name>@<YOUR-REALM>`. This principal is used to authenticate with the Hadoop cluster. where: `fully.qualified.domain.name` is the host where the Oozie server is running `YOUR-REALM` is the name of your Kerberos realm.

```
kadmin: addprinc -randkey oozie/fully.qualified.domain.name@YOUR-REALM.COM
```

2. Create a HTTP service user principal using the syntax: `HTTP/<fully.qualified.domain.name>@<YOUR-REALM>`. This principal is used to authenticate user requests coming to the Oozie web-services. where: `fully.qualified.domain.name` is the host where the Oozie server is running `YOUR-REALM` is the name of your Kerberos realm.

```
kadmin: addprinc -randkey HTTP/fully.qualified.domain.name@YOUR-REALM.COM
```



Important:

The `HTTP/` component of the HTTP service user principal must be upper case as shown in the syntax and example above.

3. Create keytab files with both principals.

```
$ kadmin
kadmin: xst -k oozie.keytab oozie/fully.qualified.domain.name
kadmin: xst -k http.keytab HTTP/fully.qualified.domain.name
```

4. Merge the two keytab files into a single keytab file:

```
$ ktutil
ktutil: rkt oozie.keytab
ktutil: rkt http.keytab
ktutil: wkt oozie-http.keytab
```

5. Test that credentials in the merged keytab file work. For example:

```
$ klist -e -k -t oozie-http.keytab
```

6. Copy the `oozie-http.keytab` file to the Oozie configuration directory. The owner of the `oozie-http.keytab` file should be the `oozie` user and the file should have owner-only read permissions.
7. Edit the Oozie server `oozie-site.xml` configuration file in the Oozie configuration directory by setting the following properties:



Important: You must restart the Oozie server to have the configuration changes take effect.

Property	Value
<code>oozie.service.HadoopAccessorService.kerberos.enabled</code>	<code>true</code>
<code>local.realm</code>	<code><REALM></code>
<code>oozie.service.HadoopAccessorService.keytab.file</code>	<code>/etc/oozie/conf/oozie-http.keytab</code> for a package installation, or <code><EXPANDED_DIR>/conf/oozie-http.keytab</code> for a tarball installation

Property	Value
oozie.service.HadoopAccessorService.kerberos.principal	oozie/<fully.qualified.domain.name>@<YOUR-REALM.COM>
oozie.authentication.type	kerberos
oozie.authentication.kerberos.principal	HTTP/<fully.qualified.domain.name>@<YOUR-REALM.COM>
oozie.authentication.kerberos.name.rules	Use the value configured for <code>hadoop.security.auth_to_local</code> in <code>core-site.xml</code>

Configuring Oozie HA with Kerberos

In CDH 5, you can configure multiple active Oozie servers against the same database, providing high availability for Oozie. For instructions on setting up Oozie HA, see [Configuring Oozie for High Availability](#)

Let's assume you have three hosts running Oozie servers, `host1.example.com`, `host2.example.com`, `host3.example.com` and the Load Balancer running on `oozie.example.com`. The Load Balancer directs traffic to the Oozie servers: `host1`, `host2` and `host3`. For such a configuration, assuming your Kerberos realm is `EXAMPLE.COM`, create the following Kerberos principals:

- `oozie/host1.example.com@EXAMPLE.COM`
- `oozie/host2.example.com@EXAMPLE.COM`
- `oozie/host3.example.com@EXAMPLE.COM`
- `HTTP/host1.example.com@EXAMPLE.COM`
- `HTTP/host2.example.com@EXAMPLE.COM`
- `HTTP/host3.example.com@EXAMPLE.COM`
- `HTTP/oozie.example.com@EXAMPLE.COM`

On each of the hosts, `host1`, `host2` and `host3`, create a keytab file with its corresponding `oozie` and `HTTP` principals from the list above. All keytab files should also have the load balancer's `HTTP` principal. Hence, each keytab file should have 3 principals in all.

Edit the following property in the Oozie server configuration file, `oozie-site.xml`:

```
<property>
<name>oozie.authentication.kerberos.principal</name>
<value>*</value>
</property>
```

Search Authentication

This section describes how to configure Search in CDH 5 to enable authentication.

When authentication is enabled, only specified hosts and users can connect to Search. Authentication also verifies that clients connect to legitimate servers. This feature prevents spoofing such as impersonation and man-in-the-middle attacks. Search supports Kerberos and LDAP authentication.

Cloudera Search supports a variety of combinations of authentication protocols:

Table 5: Authentication Protocol Combinations

Solr Authentication	Use Case
No authentication	Insecure cluster
Kerberos only	The Hadoop cluster has Kerberos turned on and every user (or client) connecting to Solr has a Kerberos principal.

Solr Authentication	Use Case
Kerberos and LDAP	The Hadoop cluster has Kerberos turned on. External Solr users (or clients) don't have Kerberos principal but do have an identity in the LDAP server. Client authentication using LDAP requires that Kerberos is enabled for the cluster. Using LDAP alone is not supported.

Once you are finished setting up authentication, configure Sentry authorization. Authorization involves specifying which resources can be accessed by particular users when they connect through Search. See [Enabling Sentry Authorization for Search using the Command Line](#) on page 422 for details.

Enabling Kerberos Authentication for Search

Cloudera Search supports Kerberos authentication. All necessary packages are installed when you install Search. To enable Kerberos, create principals and keytabs and then modify default configurations.

The following instructions only apply to configuring Kerberos in an unmanaged environment. Kerberos configuration is automatically handled by Cloudera Manager if you are using Search in a Cloudera Manager environment.

To create principals and keytabs

Repeat this process on all Solr server hosts.

1. Create a Solr service user principal using the syntax: `solr/<fully.qualified.domain.name>@<YOUR-REALM>`. This principal is used to authenticate with the Hadoop cluster. where: `fully.qualified.domain.name` is the host where the Solr server is running `YOUR-REALM` is the name of your Kerberos realm.

```
$ kadmin
kadmin: addprinc -randkey solr/fully.qualified.domain.name@YOUR-REALM.COM
```

2. Create a HTTP service user principal using the syntax: `HTTP/<fully.qualified.domain.name>@<YOUR-REALM>`. This principal is used to authenticate user requests coming to the Solr web-services. where: `fully.qualified.domain.name` is the host where the Solr server is running `YOUR-REALM` is the name of your Kerberos realm.

```
kadmin: addprinc -randkey HTTP/fully.qualified.domain.name@YOUR-REALM.COM
```



Note:

The `HTTP/` component of the HTTP service user principal must be upper case as shown in the syntax and example above.

3. Create keytab files with both principals.

```
kadmin: xst -norandkey -k solr.keytab solr/fully.qualified.domain.name \
HTTP/fully.qualified.domain.name
```

4. Test that credentials in the merged keytab file work. For example:

```
$ klist -e -k -t solr.keytab
```

5. Copy the `solr.keytab` file to the Solr configuration directory. The owner of the `solr.keytab` file should be the `solr` user and the file should have owner-only read permissions.

To modify default configurations

Repeat this process on all Solr server hosts.

1. Ensure that the following properties appear in `/etc/default/solr` or `/opt/cloudera/parcels/CDH-*/etc/default/solr` and that they are uncommented. Modify these properties to match your environment. The relevant properties to be uncommented and modified are:

```
SOLR_AUTHENTICATION_TYPE=kerberos
SOLR_AUTHENTICATION_SIMPLE_ALLOW_ANON=true
SOLR_AUTHENTICATION_KERBEROS_KEYTAB=/etc/solr/conf/solr.keytab
```

```
SOLR_AUTHENTICATION_KERBEROS_PRINCIPAL=HTTP/localhost@LOCALHOST
SOLR_AUTHENTICATION_KERBEROS_NAME_RULES=DEFAULT
SOLR_AUTHENTICATION_JAAS_CONF=/etc/solr/conf/jaas.conf
```



Note: Modify the values for these properties to match your environment. For example, the `SOLR_AUTHENTICATION_KERBEROS_PRINCIPAL=HTTP/localhost@LOCALHOST` must include the principal instance and Kerberos realm for your environment. That is often different from `localhost@LOCALHOST`.

2. Set `hadoop.security.auth_to_local` to match the value specified by `SOLR_AUTHENTICATION_KERBEROS_NAME_RULES` in `/etc/default/solr` or `/opt/cloudera/parcels/CDH-*/etc/default/solr`.



Note: For information on how to configure the rules, see [Configuring the Mapping from Kerberos Principals to Short Names](#). For additional information on using Solr with HDFS, see [Configuring Solr for Use with HDFS](#).

3. If using applications that use the `solrj` library, set up the Java Authentication and Authorization Service (JAAS).

- a. Create a `jaas.conf` file in the Solr configuration directory containing the following settings. This file and its location must match the `SOLR_AUTHENTICATION_JAAS_CONF` value. Make sure that you substitute a value for `principal` that matches your particular environment.

```
Client {
  com.sun.security.auth.module.Krb5LoginModule required
  useKeyTab=true
  useTicketCache=false
  keyTab="/etc/solr/conf/solr.keytab"
  principal="solr/fully.qualified.domain.name@<YOUR-REALM>";
};
```

Enabling LDAP Authentication for Search

Before continuing, make sure that you have completed the steps in [Enabling Kerberos Authentication for Search](#) on page 164. Solr supports LDAP authentication for external Solr client including:

- Command-line tools
- curl
- Web browsers
- Solr Java clients

In some cases, Solr does not support LDAP authentication. Use Kerberos authentication instead in these cases. Solr does not support LDAP authentication with:

- Search indexing components including the MapReduce indexer, Lily HBase indexer, or Flume.
- Solr internal requests such as those for replication or querying.
- Hadoop delegation token management requests such as `GETDELEGATIONTOKEN` or `RENEWDELEGATIONTOKEN`.

Configuring LDAP Authentication for Solr using Cloudera Manager

You can configure LDAP-based authentication using Cloudera Manager at the Solr service level.

1. Go to the **Solr** service.
2. Click the **Configuration** tab.
3. Select **Scope > Solr**
4. Select **Category > Security**
5. Select **Enable LDAP**.
6. Enter the LDAP URI in the **LDAP URI** property.
7. Configure only one of following mutually exclusive parameters:

Configuring Authentication

- **LDAP BaseDN:** Replaces the username with a "distinguished name" (DN) of the form: `uid=userid,ldap_baseDN`. Typically used for OpenLDAP server installation.

-OR-

- **LDAP Domain:** Replaces the username with a string `username@ldap_domain`. Typically used for Active Directory server installation.

Configuring LDAP Authentication for Solr Using the Command Line

To enable LDAP authentication using the command line, configure the following environment variables in `/etc/default/solr`:

```
SOLR_AUTHENTICATION_HTTP_SCHEMES=Negotiate,Basic
SOLR_AUTHENTICATION_HTTP_DELEGATION_MGMT_SCHEMES=Negotiate
SOLR_AUTHENTICATION_HTTP_BASIC_HANDLER=ldap
SOLR_AUTHENTICATION_HTTP_NEGOTIATE_HANDLER=kerberos
SOLR_AUTHENTICATION_LDAP_PROVIDER_URL=ldap://www.example.com

# Specify value for only one of SOLR_AUTHENTICATION_LDAP_BASE_DN or
SOLR_AUTHENTICATION_LDAP_BIND_DOMAIN property.
SOLR_AUTHENTICATION_LDAP_BASE_DN=ou=Users,dc=example,dc=com
# SOLR_AUTHENTICATION_LDAP_BIND_DOMAIN=
# Required when using 'Start TLS' extension
# SOLR_AUTHENTICATION_LDAP_ENABLE_START_TLS=false
```

Securing LDAP Connections

You can secure communications using LDAP-based encryption.

To avoid sending credentials over the wire in clear-text, you must configure a secure connection between both the client and Solr, and between Solr and the LDAP server. The secure connection could use SSL or TLS.

Secure LDAP connections through SSL:

For SSL-enabled LDAP connections, specify a prefix of `ldaps://` instead of `ldap://`. Also, the default port for SSL-enabled LDAP connections is 636 instead of 389.

Secure LDAP connections through TLS:

[TLS](#), the successor to the SSL protocol, is supported by most modern LDAP servers. Unlike SSL connections, TLS connections can be made on the same server port as non-TLS connections. You can enable xxx using Cloudera Manager.

1. Go to the **Solr** service.
2. Click the **Configuration** tab.
3. Select **Scope > Solr**
4. Select **Category > Security**
5. Select **Enable LDAP TLS**.
6. Import the LDAP server security certificate in the Solr Trust Store file:
 - a. Enter the location for the Solr Trust Store File in **Solr TLS/SSL Certificate Trust Store File**.
 - b. Enter the password for the Solr Trust Store File in **Solr TLS/SSL Certificate Trust Store Password**.

LDAP Client Configuration

Some HTTP clients such as curl or the Apache Http Java client must be configured to use a particular scheme. For example:

- curl tool supports using Kerberos or username/password authentication. Kerberos is activated using the `--negotiate` flag and username/password based authentication is activated using the `--basic` and `-u` flags.
- Apache HttpClient library can be configured to use specific authentication scheme. For more information, see the [HTTP authentication](#) chapter of Apache's HttpClient Tutorial.

Typically, web browsers automatically choose a preferred authentication scheme. For more information, see the [HTTP authentication](#) topic in The Chromium Projects.

To use LDAP authentication with Solr Java clients, `HttpClientConfigurer` needs to be configured for Solr. This can either be done programmatically or using Java system properties.

For example, programmatic initialization might appear as:

SampleSolrClient.java

```
import org.apache.solr.client.solrj.impl.HttpClientUtil;
import org.apache.solr.client.solrj.impl.PreemptiveBasicAuthConfigurer;
import org.apache.solr.common.params.ModifiableSolrParams;

/**
 * This method initializes the Solr client to use LDAP authentication
 * This configuration is applicable to all Solr clients.
 * @param ldapUserName LDAP user name
 * @param ldapPassword LDAP user password
 */
public static void initialize(String ldapUserName, String ldapPassword) {
    HttpClientUtil.setConfigurer(new PreemptiveBasicAuthConfigurer());
    ModifiableSolrParams params = new ModifiableSolrParams();
    params.set(HttpClientUtil.PROP_BASIC_AUTH_USER, ldapUserName);
    params.set(HttpClientUtil.PROP_BASIC_AUTH_PASS, ldapPassword);
    // Configure the JVM default parameters.
    PreemptiveBasicAuthConfigurer.setDefaultSolrParams(params);
}
```

For configuration using system properties, configure the following system properties:

Table 6: System properties configuration for LDAP authentication

System property	Description
<code>solr.httpclient.configurer</code>	Fully qualified classname of <code>HttpClientConfigurer</code> implementation. For example, <code>org.apache.solr.client.solrj.impl.PreemptiveBasicAuthConfigurer</code> .
<code>solr.httpclient.config</code>	Http client configuration properties file path. For example, <code>ldap-credentials.properties</code> .

For example, the entry in `ldap-credentials.properties` might appear as:

ldap-credentials.properties

```
httpBasicAuthUser=user1
httpBasicAuthPassword=passwd
```

Using Kerberos with Search

The process of enabling Solr clients to authenticate with a secure Solr is specific to the client. This section demonstrates:

- Using Kerberos and curl
- Using solrctl
- Configuring SolrJ Library Usage
- This enables technologies including:
 - Command line solutions
 - Java applications
 - The MapReduceIndexerTool
- Configuring Flume Morphline Solr Sink Usage

Configuring Authentication

Secure Solr requires that the CDH components that it interacts with are also secure. Secure Solr interacts with HDFS, ZooKeeper and optionally HBase, MapReduce, and Flume.

Using Kerberos and curl

You can use Kerberos authentication with clients such as `curl`. To use `curl`, begin by acquiring valid Kerberos credentials and then run the desired command. For example, you might use commands similar to the following:

```
$ kinit -kt username.keytab username
$ curl --negotiate -u foo:bar http://solrserver:8983/solr/
```



Note: Depending on the tool used to connect, additional arguments may be required. For example, with `curl`, `--negotiate` and `-u` are required. The username and password specified with `-u` is not actually checked because Kerberos is used. As a result, any value such as `foo:bar` or even just `:` is acceptable. While any value can be provided for `-u`, note that the option is required. Omitting `-u` results in a 401 Unauthorized error, even though the `-u` value is not actually used.

Using solrctl

If you are using `solrctl` to manage your deployment in an environment that requires Kerberos authentication, you must have valid Kerberos credentials, which you can get using `kinit`. For more information on `solrctl`, see [solrctl Reference](#)

Configuring SolrJ Library Usage

If using applications that use the `solrj` library, begin by establishing a Java Authentication and Authorization Service (JAAS) configuration file.

Create a JAAS file:

- If you have already used `kinit` to get credentials, you can have the client use those credentials. In such a case, modify your `jaas-client.conf` file to appear as follows:

```
Client {
  com.sun.security.auth.module.Krb5LoginModule required
  useKeyTab=false
  useTicketCache=true
  principal="user/fully.qualified.domain.name@<YOUR-REALM>";
};
```

where `user/fully.qualified.domain.name@<YOUR-REALM>` is replaced with your credentials.

- You want the client application to authenticate using a keytab you specify:

```
Client {
  com.sun.security.auth.module.Krb5LoginModule required
  useKeyTab=true
  keyTab="/path/to/keytab/user.keytab"
  storeKey=true
  useTicketCache=false
  principal="user/fully.qualified.domain.name@<YOUR-REALM>";
};
```

where `/path/to/keytab/user.keytab` is the keytab file you want to use and

`user/fully.qualified.domain.name@<YOUR-REALM>` is the principal in that keytab you want to use.

Use the JAAS file to enable solutions:

- **Command line solutions**

Set the property when invoking the program. For example, if you were using a jar, you might use:

```
java -Djava.security.auth.login.config=/home/user/jaas-client.conf -jar app.jar
```

- **Java applications**

Set the Java system property `java.security.auth.login.config`. For example, if the JAAS configuration file is located on the filesystem as `/home/user/jaas-client.conf`. The Java system property `java.security.auth.login.config` must be set to point to this file. Setting a Java system property can be done programmatically, for example using a call such as:

```
System.setProperty("java.security.auth.login.config", "/home/user/jaas-client.conf");
```

- **The MapReduceIndexerTool**

The `MapReduceIndexerTool` uses SolrJ to pass the JAAS configuration file. Using the `MapReduceIndexerTool` in a secure environment requires the use of the `HADOOP_OPTS` variable to specify the JAAS configuration file. For example, you might issue a command such as the following:

```
HADOOP_OPTS="-Djava.security.auth.login.config=/home/user/jaas.conf" \  
hadoop jar MapReduceIndexerTool
```

- **Configuring the hbase-indexer CLI**

Certain `hbase-indexer` CLI commands such as `replication-status` attempt to read ZooKeeper hosts owned by HBase. To successfully use these commands in Search for CDH 5 in a secure environment, specify a JAAS configuration file with the HBase principal in the `HBASE_INDEXER_OPTS` environment variable. For example, you might issue a command such as the following:

```
HBASE_INDEXER_OPTS="-Djava.security.auth.login.config=/home/user/hbase-jaas.conf" \  
hbase-indexer replication-status
```

Configuring Flume Morphline Solr Sink Usage

Repeat this process on all Flume hosts:

1. If you have not created a keytab file, do so now at `/etc/flume-ng/conf/flume.keytab`. This file should contain the service principal `flume/<fully.qualified.domain.name>@<YOUR-REALM>`. See [Flume Authentication](#) on page 119 for more information.
2. Create a JAAS configuration file for flume at `/etc/flume-ng/conf/jaas-client.conf`. The file should appear as follows:

```
Client {  
  com.sun.security.auth.module.Krb5LoginModule required  
  useKeyTab=true  
  useTicketCache=false  
  keyTab="/etc/flume-ng/conf/flume.keytab"  
  principal="flume/<fully.qualified.domain.name>@<YOUR-REALM>";  
};
```

3. Add the flume JAAS configuration to the `JAVA_OPTS` in `/etc/flume-ng/conf/flume-env.sh`. For example, you might change:

```
JAVA_OPTS="-Xmx500m"
```

to:

```
JAVA_OPTS="-Xmx500m -Djava.security.auth.login.config=/etc/flume-ng/conf/jaas-client.conf"
```

Spark Authentication

Spark currently support two methods of authentication. Authentication can be configured using Kerberos *or* using a [shared secret](#). When using Spark on YARN, Cloudera recommends using Kerberos authentication since it is stronger security measure.

Configuring Kerberos Authentication for Spark Using the Command Line

Important:

- If you want to enable Spark event logging on a Kerberos-enabled cluster, you will need to enable Kerberos authentication for Spark as well, since Spark's event logs are written to HDFS.
- You can use Spark on a Kerberos-enabled cluster only in the YARN mode, *not* in the Standalone mode.

The following steps describe how to set up Kerberos authentication for Spark using the command line.

Create the Spark Principal and Keytab File

1. Create the `spark` principal and `spark.keytab` file:

```
kadmin: addprinc -randkey spark/fully.qualified.domain.name@YOUR-REALM.COM
kadmin: xst -k spark.keytab spark/fully.qualified.domain.name
```

2. Move the file into the Spark configuration directory and restrict its access exclusively to the `spark` user:

```
$ mv spark.keytab /etc/spark/conf/
$ chown spark /etc/spark/conf/spark.keytab
$ chmod 400 /etc/spark/conf/spark.keytab
```

For more details on creating Kerberos principals and keytabs, see [Step 4: Create and Deploy the Kerberos Principals and Keytab Files](#) on page 98.

Configure the Spark History Server to Use Kerberos

Open the Spark configuration file `/etc/spark/conf/spark-env.sh` file and add the following properties:

```
SPARK_HISTORY_OPTS=-Dspark.history.kerberos.enabled=true \
-Dspark.history.kerberos.principal=spark/fully.qualified.domain.name@YOUR-REALM.COM \
-Dspark.history.kerberos.keytab=/etc/spark/conf/spark.keytab
```

Running Spark Applications on a Secure Cluster

You can submit compiled Spark applications with the `spark-submit` script. Specify the following additional command-line options when running the `spark-submit` script on a secure cluster using the form: `--option value`.

Option	Description
<code>--keytab</code>	The full path to the file that contains the keytab for the principal. This keytab is copied to the node running the ApplicationMaster using the Secure Distributed Cache, for periodically renewing the login tickets and the delegation tokens. For information on setting up the principal and keytab, see Configuring a Cluster with Custom Kerberos Principals on page 70 and Spark Authentication on page 170.
<code>--principal</code>	Principal to be used to log in to the KDC, while running on secure HDFS.
<code>--proxy-user</code>	This property allows you to use the <code>spark-submit</code> script to impersonate client users when submitting jobs.

Configuring Spark Authentication With a Shared Secret Using Cloudera Manager

Minimum Required Role: Security Administrator (also provided by **Full Administrator**)

Authentication using a shared secret can be configured using the `spark.authenticate` configuration parameter. The authentication process is essentially a handshake between Spark and the other party to ensure they have the same shared secret and can be allowed to communicate. If the shared secret does not match, they will not be allowed to communicate.

If you are using Spark on YARN, [setting](#) `spark.authenticate` parameter to `true` will generate and distribute the shared secret to all applications communicating with Spark. For Cloudera Manager deployments, use the following instructions:

1. Go to the **Spark Service > Configuration** tab.
2. In the Search field, type `spark.authenticate` to find the **Spark Authentication** settings.
3. Check the checkbox for the **Spark Authentication** property.
4. Click **Save Changes**.

Configuring Spark on YARN for Long-Running Applications

For long-running applications, such as Spark Streaming jobs, to write to HDFS, you must configure [Kerberos authentication for Spark](#) for Spark, and pass the Spark principal and keytab to the `spark-submit` script using the `--principal` and `--keytab` parameters. The keytab is copied to the host running the ApplicationMaster, and the Kerberos login is renewed periodically by using the principal and keytab to generate the required delegation tokens for communication with HDFS.

To make sure the Spark keytab is delivered to the ApplicationMaster host securely, configure [TLS/SSL communication for YARN](#) and [HDFS encryption](#) on your cluster.

Sqoop 2 Authentication

This section describes how to configure Sqoop 2 with Kerberos security in a Hadoop cluster.



Note: Sqoop 2 is being deprecated. Cloudera recommends using Sqoop 1.

Create the Sqoop 2 Principal and Keytab File

You need to create a `sqoop2.keytab` file for Sqoop 2. Follow these steps:

1. Create the principal and keytab file:

```
kadmin: addprinc -randkey sqoop2/fully.qualified.domain.name@YOUR-REALM.COM
kadmin: xst -k sqoop2.keytab sqoop2/fully.qualified.domain.name
```

2. Move the file into the Sqoop 2 configuration directory and restrict its access exclusively to the `sqoop2` user:

```
$ mv sqoop2.keytab /etc/sqoop2/conf/
$ chown sqoop2 /etc/sqoop2/conf/sqoop2.keytab
$ chmod 400 /etc/sqoop2/conf/sqoop2.keytab
```

For more details on creating Kerberos principals and keytabs, see [Step 4: Create and Deploy the Kerberos Principals and Keytab Files](#) on page 98.

Configure Sqoop 2 to Use Kerberos

Edit the Sqoop 2 configuration file `sqoop.properties` file in the `/etc/sqoop2/conf` directory and add the following properties:

```
org.apache.sqoop.authentication.type=KERBEROS
org.apache.sqoop.authentication.handler=org.apache.sqoop.security.KerberosAuthenticationHandler
org.apache.sqoop.authentication.kerberos.principal=sqoop2/fully.qualified.domain.name@YOUR-REALM.COM
org.apache.sqoop.authentication.kerberos.keytab=/etc/sqoop2/conf/sqoop2.keytab
```

ZooKeeper Authentication

This section describes how to configure ZooKeeper in CDH 5 to enable Kerberos security:

- [Configuring ZooKeeper Server for Kerberos Authentication](#) on page 172
- [Configuring the ZooKeeper Client Shell to Support Kerberos Security](#) on page 173
- [Verifying the Configuration](#) on page 174



Important:

Prior to enabling ZooKeeper to work with Kerberos security on your cluster, make sure you first review the requirements in [Configuring Hadoop Security in CDH 5](#).

Configuring ZooKeeper Server for Kerberos Authentication

You can configure the ZooKeeper server for Kerberos authentication in Cloudera Manager or through the command line.

Using Cloudera Manager to Configure ZooKeeper Server for Kerberos Authentication

To set up the ZooKeeper server for Kerberos authentication in Cloudera Manager, complete the following steps:

1. In Cloudera Manager, open the ZooKeeper service.
2. Click the **Configuration** tab.
3. Enter **Kerberos** in the in the **Search** bar.
4. Find the **Enable Kerberos Authentication** property and select the check-box next to the ZooKeeper services that you want to configure for Kerberos authentication.

Using the Command Line to Configure ZooKeeper Server for Kerberos Authentication

Follow the steps below for each ZooKeeper server in the ensemble. To maintain consistency across ZooKeeper servers in the ensemble, use the same name for the keytab file you deploy to each server, for example, `zookeeper.keytab`. Each keytab file will contain its respective host's fully-qualified domain name (FQDN).

1. Create a service principal for the ZooKeeper server using the fully-qualified domain name (FQDN) of the host on which ZooKeeper server is running and the name of your Kerberos realm using the pattern `zookeeper/fqdn.example.com@YOUR-REALM`. This principal will be used to authenticate the ZooKeeper server with the Hadoop cluster. Create this service principal as follows:

```
kadmin: addprinc -randkey zookeeper/fqdn.example.com@YOUR-REALM
```

2. Create a keytab file for the ZooKeeper server:

```
$ kadmin
kadmin: xst -k zookeeper.keytab zookeeper/fqdn.example.com@YOUR-REALM
```



Note: For consistency across ZooKeeper Servers, use the same name for the keytab file you create for each subsequent ZooKeeper Server host system you configure using these steps, for example, `zookeeper.keytab`.

3. Copy the `zookeeper.keytab` file to the ZooKeeper configuration directory on the ZooKeeper server host, using the appropriate ZooKeeper configuration directory: `/etc/zookeeper/conf/`. The `zookeeper.keytab` file should be owned by the `zookeeper` user, with owner-only read permissions.
4. Add the following lines to the ZooKeeper configuration file `zoo.cfg`:

```
authProvider.1=org.apache.zookeeper.server.auth.SASLAuthenticationProvider
jaasLoginRenew=3600000
```

5. Set up the [Java Authentication and Authorization Service \(JAAS\)](#) by creating a `jaas.conf` file in the ZooKeeper configuration directory with the settings shown below, replacing `fqdn.example.com` with the ZooKeeper server's hostname.

```
Server {
  com.sun.security.auth.module.Krb5LoginModule required
  useKeyTab=true
  keyTab="/etc/zookeeper/conf/zookeeper.keytab"
  storeKey=true
  useTicketCache=false
  principal="zookeeper/fqdn.example.com
  @YOUR-REALM";
};
```

6. Add the following setting to the `java.env` file located in the ZooKeeper configuration directory, creating the file if necessary:

```
export JVMFLAGS="-Djava.security.auth.login.config=/etc/zookeeper/conf/jaas.conf"
```

7. Repeat these steps for each ZooKeeper server in the ensemble.
8. Restart the ZooKeeper server to have the configuration changes take effect. See [ZooKeeper Installation](#) for details.

Configuring the ZooKeeper Client Shell to Support Kerberos Security

1. If you want to use the ZooKeeper client shell `zookeeper-client` with Kerberos authentication, create a principal using the syntax: `zkcli@<YOUR-REALM>`. This principal is used to authenticate the ZooKeeper client shell to the ZooKeeper service. where: `YOUR-REALM` is the name of your Kerberos realm.

```
kadmin: addprinc -randkey zkcli@YOUR-REALM.COM
```

2. Create a keytab file for the ZooKeeper client shell.

```
$ kadmin
kadmin: xst -norandkey -k zkcli.keytab zkcli@YOUR-REALM.COM
```



Note:

Some versions of `kadmin` do not support the `-norandkey` option in the command above. If your version does not, you can omit it from the command. Note that doing so will result in a new password being generated every time you export a keytab, which will invalidate previously-exported keytabs.

Configuring Authentication

3. Set up JAAS in the configuration directory on the host where the ZooKeeper client shell is running. For a package installation, the configuration directory is `/etc/zookeeper/conf/`. For a tar ball installation, the configuration directory is `<EXPANDED_DIR>/conf`. Create a `jaas.conf` file containing the following settings:

```
Client {
  com.sun.security.auth.module.Krb5LoginModule required
  useKeyTab=true
  keyTab="/path/to/zkcli.keytab"
  storeKey=true
  useTicketCache=false
  principal="zkcli@YOUR-REALM";
};
```

4. Add the following setting to the `java.env` file located in the configuration directory. (Create the file if it does not already exist.)

```
export JVMFLAGS="-Djava.security.auth.login.config=/etc/zookeeper/conf/jaas.conf"
```

Verifying the Configuration

1. Make sure that you have restarted the ZooKeeper cluster with Kerberos enabled, as described above.
2. Start the client (where the hostname is the name of a ZooKeeper server):

```
zookeeper-client -server hostname:port
```

3. Create a protected znode from within the ZooKeeper CLI. Make sure that you substitute `YOUR-REALM` as appropriate.

```
create /znode1 znode1data sasl:zkcli@{{YOUR-REALM}}:cdwra
```

4. Verify the znode is created and the ACL is set correctly:

```
getAcl /znode1
```

The results from `getAcl` should show that the proper scheme and permissions were applied to the znode.


Hadoop Users in Cloudera Manager and CDH

A number of special users are created by default when installing and using CDH and Cloudera Manager. Given below is a list of users and groups as of the latest release. Also listed are the corresponding Kerberos principals and keytab files that should be created when you configure Kerberos security on your cluster.



Note: Cloudera Manager 5.3 introduces a new *single user mode*. In single user mode, the Cloudera Manager Agent and *all the processes run by services managed by Cloudera Manager* are started as a single configured user and group. See [Single User Mode Requirements](#) for more information.

Table 7: Users and Groups

Component (Version)	Unix User ID	Groups	Notes
Cloudera Manager (all versions)	cloudera-scm	cloudera-scm	<p>Cloudera Manager processes such as the Cloudera Manager Server and the monitoring roles run as this user.</p> <p>The Cloudera Manager keytab file must be named <code>cmf.keytab</code> since that name is hard-coded in Cloudera Manager.</p> <div style="border: 1px solid green; padding: 5px; margin-top: 10px;"> <p> Note: Applicable to clusters managed by Cloudera Manager only.</p> </div>
Apache Accumulo (Accumulo 1.4.3 and higher)	accumulo	accumulo	Accumulo processes run as this user.
Apache Avro			No special users.
Apache Flume (CDH 4, CDH 5)	flume	flume	The sink that writes to HDFS as this user must have write privileges.
Apache HBase (CDH 4, CDH 5)	hbase	hbase	The Master and the RegionServer processes run as this user.
HDFS (CDH 4, CDH 5)	hdfs	hdfs, hadoop	The NameNode and DataNodes run as this user, and the HDFS root directory as well as the directories used for edit logs should be owned by it.
Apache Hive (CDH 4, CDH 5)	hive	hive	<p>The HiveServer2 process and the Hive Metastore processes run as this user.</p> <p>A user must be defined for Hive access to its Metastore DB (for example, MySQL or Postgres) but it can be any identifier and does not correspond to a Unix uid. This is <code>javax.jdo.option.ConnectionUserName</code> in <code>hive-site.xml</code>.</p>
Apache HCatalog (CDH 4.2 and higher, CDH 5)	hive	hive	The WebHCat service (for REST access to Hive functionality) runs as the <code>hive</code> user.
HttpFS (CDH 4, CDH 5)	httpfs	httpfs	The HttpFS service runs as this user. See HttpFS Security Configuration for instructions on how to generate the merged <code>httpfs-http.keytab</code> file.
Hue (CDH 4, CDH 5)	hue	hue	Hue services run as this user.
Hue Load Balancer (Cloudera Manager 5.5 and higher)	apache	apache	The Hue Load balancer has a dependency on the <code>apache2</code> package that uses the <code>apache</code> user name. Cloudera Manager does not run processes using this user ID.
Impala	impala	impala, hive	Impala services run as this user.

Component (Version)	Unix User ID	Groups	Notes
Apache Kafka (Cloudera Distribution of Kafka 1.2.0)	kafka	kafka	Kafka services run as this user.
Java KeyStore KMS (CDH 5.2.1 and higher)	kms	kms	The Java KeyStore KMS service runs as this user.
Key Trustee KMS (CDH 5.3 and higher)	kms	kms	The Key Trustee KMS service runs as this user.
Key Trustee Server (CDH 5.4 and higher)	keytrustee	keytrustee	The Key Trustee Server service runs as this user.
Kudu	kudu	kudu	Kudu services run as this user.
Llama (CDH 5)	llama	llama	Llama runs as this user.
Apache Mahout			No special users.
MapReduce (CDH 4, CDH 5)	mapred	mapred, hadoop	Without Kerberos, the JobTracker and tasks run as this user. The LinuxTaskController binary is owned by this user for Kerberos.
Apache Oozie (CDH 4, CDH 5)	oozie	oozie	The Oozie service runs as this user.
Parquet			No special users.
Apache Pig			No special users.
Cloudera Search (CDH 4.3 and higher, CDH 5)	solr	solr	The Solr processes run as this user.
Apache Spark (CDH 5)	spark	spark	The Spark History Server process runs as this user.
Apache Sentry (incubating) (CDH 5.1 and higher)	sentry	sentry	The Sentry service runs as this user.
Apache Sqoop (CDH 4, CDH 5)	sqoop	sqoop	This user is only for the Sqoop1 Metastore, a configuration option that is not recommended.
Apache Sqoop2 (CDH 4.2 and higher, CDH 5)	sqoop2	sqoop, sqoop2	The Sqoop2 service runs as this user.
Apache Whirr			No special users.
YARN (CDH 4, CDH 5)	yarn	yarn, hadoop	Without Kerberos, all YARN services and applications run as this user. The LinuxContainerExecutor binary is owned by this user for Kerberos.
Apache ZooKeeper (CDH 4, CDH 5)	zookeeper	zookeeper	The ZooKeeper processes run as this user. It is not configurable.

Keytabs and Keytab File Permissions

**Note:**

The Kerberos principal names should be of the format, `username/fully.qualified.domain.name@YOUR-REALM.COM`, where the term `username` refers to the username of an existing UNIX account, such as `hdfs` or `mapred`. The table below lists the usernames to be used for the Kerberos principal names. For example, the Kerberos principal for Apache Flume would be `flume/fully.qualified.domain.name@YOUR-REALM.COM`.

For keytabs with multiple principals, Cloudera Manager merges them appropriately from individual keytabs. If you do not use Cloudera Manager, you must merge the keytabs manually.

Table 8: Clusters Managed by Cloudera Manager

Component (Unix User ID)	Service	Kerberos Principals	Filename (*.keytab)	Keytab File Owner	Keytab File Group	File Permission (octal)
Cloudera Manager (cloudera-scm)	NA	cloudera-scm	cmf	cloudera-scm	cloudera-scm	600
Cloudera Management Service (cloudera-scm)	cloudera-mgmt-REPORTSMANAGER	hdfs	headlamp	cloudera-scm	cloudera-scm	600
Cloudera Management Service (cloudera-scm)	cloudera-mgmt-SERVICEMONITOR, cloudera-mgmt-ACTIVITYMONITOR	hue	cmon	cloudera-scm	cloudera-scm	600
Cloudera Management Service (cloudera-scm)	cloudera-mgmt-HOSTMONITOR	N/A	N/A	N/A	N/A	N/A
Apache Accumulo (accumulo)	accumulo16-ACCUMULO16_MASTER	accumulo	accumulo16	cloudera-scm	cloudera-scm	600
	accumulo16-ACCUMULO16_TRACER					
	accumulo16-ACCUMULO16_MONITOR					
	accumulo16-ACCUMULO16_GC					
	accumulo16-ACCUMULO16_TSERVER					
Flume (flume)	flume-AGENT	flume	flume	cloudera-scm	cloudera-scm	600
HBase (hbase)	hbase-HBASETHRIFTSERVER	HTTP	HTTP	cloudera-scm	cloudera-scm	600
	hbase-REGIONSERVER	hbase	hbase			
	hbase-HBASERESTSERVER					
	hbase-MASTER					
HDFS (hdfs)	hdfs-NAMENODE	hdfs, HTTP	hdfs	cloudera-scm	cloudera-scm	600
	hdfs-DATANODE					
	hdfs-SECONDARYNAMENODE					

Component (Unix User ID)	Service	Kerberos Principals	Filename (*.keytab)	Keytab File Owner	Keytab File Group	File Permission (octal)
Hive (hive)	hive-HIVESERVER2	hive	hive	cloudera-scm	cloudera-scm	600
	hive-WEBHCAT	HTTP	HTTP			
	hive-HIVEMETASTORE	hive	hive			
HttpFS (httpfs)	hdfs-HTTPFS	httpfs	httpfs	cloudera-scm	cloudera-scm	600
Hue (hue)	hue-KT_RENEWER	hue	hue	cloudera-scm	cloudera-scm	600
Impala (impala)	impala-STATESTORE	impala	impala	cloudera-scm	cloudera-scm	600
	impala-CATALOGSERVER					
	impala-IMPALAD					
Java KeyStore KMS (kms)	kms-KMS	HTTP	kms	cloudera-scm	cloudera-scm	600
Apache Kafka (kafka)	kafka-KAFKA_BROKER	kafka	kafka	kafka	kafka	600
Key Trustee KMS (kms)	keytrustee-KMS_KEYTRUSTEE	HTTP	keytrustee	cloudera-scm	cloudera-scm	600
Llama (llama)	impala-LLAMA	llama, HTTP	llama	cloudera-scm	cloudera-scm	600
MapReduce (mapred)	mapreduce-JOBTRACKER	mapred, HTTP	mapred	cloudera-scm	cloudera-scm	600
	mapreduce-TASKTRACKER					
Oozie (oozie)	oozie-OOZIE_SERVER	oozie, HTTP	oozie	cloudera-scm	cloudera-scm	600
Search (solr)	solr-SOLR_SERVER	solr, HTTP	solr	cloudera-scm	cloudera-scm	600
Sentry (sentry)	sentry-SENTRY_SERVER	sentry	sentry	cloudera-scm	cloudera-scm	600
Spark (spark)	spark_on_yarn-SPARK_YARN_HISTORY_SERVER	spark	spark	cloudera-scm	cloudera-scm	600
YARN (yarn)	yarn-NODEMANAGER	yarn, HTTP	yarn	cloudera-scm	cloudera-scm	644
	yarn-RESOURCEMANAGER					600
	yarn-JOBHISTORY					600
ZooKeeper (zookeeper)	zookeeper-server	zookeeper	zookeeper	cloudera-scm	cloudera-scm	600

Table 9: CDH Clusters Not Managed by Cloudera Manager

Component (Unix User ID)	Service	Kerberos Principals	Filename (*.keytab)	Keytab File Owner	Keytab File Group	File Permission (octal)
Apache Accumulo (accumulo)	accumulo16-ACCUMULO16_MASTER	accumulo	accumulo16	accumulo	accumulo	600
	accumulo16-ACCUMULO16_TRACER					
	accumulo16-ACCUMULO16_MONITOR					
	accumulo16-ACCUMULO16_GC					
	accumulo16-ACCUMULO16_TSERVER					

Component (Unix User ID)	Service	Kerberos Principals	Filename (*.keytab)	Keytab File Owner	Keytab File Group	File Permission (octal)
Flume (flume)	flume-AGENT	flume	flume	flume	flume	600
HBase (hbase)	hbase-HBASETHRIFTSERVER	HTTP	HTTP	hbase	hbase	600
	hbase-REGIONSERVER	hbase	hbase			
	hbase-HBASERESTSERVER					
	hbase-MASTER					
HDFS (hdfs)	hdfs-NAMENODE	hdfs, HTTP	hdfs	hdfs	hdfs	600
	hdfs-DATANODE					
	hdfs- SECONDARYNAMENODE					
Hive (hive)	hive-HIVESERVER2	hive	hive	hive	hive	600
	hive-WEBHCAT	HTTP	HTTP			
	hive-HIVEMETASTORE	hive	hive			
HttpFS (httpfs)	hdfs-HTTPFS	httpfs	httpfs	httpfs	httpfs	600
Hue (hue)	hue-KT_RENEWER	hue	hue	hue	hue	600
Impala (impala)	impala-STATESTORE	impala	impala	impala	impala	600
	impala-CATALOGSERVER					
	impala-IMPALAD					
Llama (llama)	impala-LLAMA	llama, HTTP	llama	llama	llama	600
Java KeyStore KMS (kms)	kms-KMS	HTTP	kms	kms	kms	600
Apache Kafka (kafka)	kafka-KAFKA_BROKER	kafka	kafka	kafka	kafka	600
Key Trustee KMS (kms)	kms-KEYTRUSTEE	HTTP	kms	kms	kms	600
MapReduce (mapred)	mapreduce-JOBTRACKER	mapred, HTTP	mapred	mapred	hadoop	600
	mapreduce- TASKTRACKER					
Oozie (oozie)	oozie-OOZIE_SERVER	oozie, HTTP	oozie	oozie	oozie	600
Search (solr)	solr-SOLR_SERVER	solr, HTTP	solr	solr	solr	600
Sentry (sentry)	sentry-SENTRY_SERVER	sentry	sentry	sentry	sentry	600
Spark (spark)	spark_on_yarn-SPARK_YARN_HISTORY_SERVER	spark	spark	spark	spark	600
YARN (yarn)	yarn-NODEMANAGER	yarn, HTTP	yarn	yarn	hadoop	644
	yarn- RESOURCEMANAGER					600
	yarn-JOBHISTORY					600

Component (Unix User ID)	Service	Kerberos Principals	Filename (*.keytab)	Keytab File Owner	Keytab File Group	File Permission (octal)
ZooKeeper (zookeeper)	zookeeper-server	zookeeper	zookeeper	zookeeper	zookeeper	600

Configuring a Cluster-dedicated MIT KDC with Cross-Realm Trust

If you use Cloudera Manager to enable Hadoop security on your cluster, the Cloudera Manager Server will create several principals and then generate keytabs for those principals. Cloudera Manager will then deploy the keytab files on every host in the cluster. See [Hadoop Users in Cloudera Manager and CDH](#) on page 174 for a complete listing of the principals created by Cloudera Manager.



Note: The following instructions illustrate an example of creating and deploying the principals and keytab files for MIT Kerberos. (If you are using another version of Kerberos, refer to the Kerberos documentation for the version of the operating system you are using, for instructions.)

When to use `kadmin.local` and `kadmin`

When performing the Kerberos commands in this document, you can use `kadmin.local` or `kadmin` depending on your access and account:

- If you can log on to the KDC host directly, and have root access or a Kerberos admin account, use the `kadmin.local` command.
- When accessing the KDC from a remote host, use the `kadmin` command.

To start `kadmin.local` on the KDC host:

```
$ sudo kadmin.local
```

To run `kadmin` from any host:

```
$ kadmin
```



Note:

- In this guide, `kadmin` is shown as the prompt for commands in the `kadmin` shell, but you can type the same commands at the `kadmin.local` prompt in the `kadmin.local` shell.
- Running `kadmin.local` may prompt you for a password because it is being run using `sudo`. You should provide your Unix password. Running `kadmin` may prompt you for a password because you need Kerberos admin privileges. You should provide your Kerberos admin password.

Setting up a Cluster-Dedicated KDC and Default Realm for the Hadoop Cluster

Cloudera has tested the following configuration approaches to Kerberos security for clusters managed by Cloudera Manager. For administration teams that are just getting started with Kerberos security, we recommend starting with these approaches to the configuration of KDC services for a number of reasons.

The number of Service Principal Names (SPNs) that are created and managed by the Cloudera Manager server for a CDH cluster can be significant, so it is important to realize the potential impact on cluster uptime and overall operations if you choose to manage keytabs manually instead. The Cloudera Manager server manages the creation of service keytabs on the proper hosts based on the current configuration of the database. Manual keytab management can be error prone and introduce delays when deploying or moving services within the cluster, especially under time-sensitive conditions.

Cloudera Manager creates SPNs within a KDC that it can access with the `kadmin` command based on configuration of the `/etc/krb5.conf` file on the Cloudera Manager host. SPNs are created with the format `service-name/host.fqdn.name@EXAMPLE.COM` where `service-name` is the relevant CDH service name such as `hue` or `hbase` or `hdfs`.

If your site already has a working KDC, and any existing principals share the same name as any of the principals that Cloudera Manager creates, the Cloudera Manager Server generates a new randomized key for those principals, and consequently causes existing keytabs to become invalid.

This is why Cloudera recommends using a dedicated local MIT Kerberos KDC and realm for the Hadoop cluster. You can set up a one-way cross-realm trust from the cluster-dedicated KDC and realm to your existing central MIT Kerberos KDC, or to an existing Active Directory realm. Using this method, there is no need to create Hadoop service principals in the central MIT Kerberos KDC or in Active Directory, but principals (users) in the central MIT KDC or in Active Directory can be authenticated to Hadoop. The steps to implement this approach are as follows:

1. Install and configure a cluster-dedicated MIT Kerberos KDC that will be managed by Cloudera Manager for creating and storing the service principals for your Hadoop cluster.



Note: The `krb5-server` package includes a `logrotate` policy file to rotate log files monthly. To take advantage of this, install the `logrotate` package. No additional configuration is necessary.

2. See the example `kdc.conf` and `krb5.conf` files in [Sample Kerberos Configuration files: krb5.conf, kdc.conf, kadm5.acl](#) on page 194 for configuration considerations for the KDC and Kerberos clients.
3. Configure a default Kerberos realm for the cluster you want Cloudera Manager to manage and set up one-way cross-realm trust between the cluster-dedicated KDC and either your central KDC or Active Directory. Follow the appropriate instructions below for your deployment: [Using a Cluster-Dedicated KDC with a Central MIT KDC](#) on page 181 or [Using a Cluster-Dedicated MIT KDC with Active Directory](#) on page 183.

Cloudera strongly recommends the method above because:

- It requires minimal configuration in Active Directory.
- It is comparatively easy to script the creation of many principals and keytabs. A principal and keytab must be created for every daemon in the cluster, and in a large cluster this can be extremely onerous to do directly in Active Directory.
- There is no need to involve central Active Directory administrators in order to get service principals created.
- It allows for incremental configuration. The Hadoop administrator can completely configure and verify the functionality the cluster independently of integrating with Active Directory.

Using a Cluster-Dedicated KDC with a Central MIT KDC



Important: If you plan to use Oozie or the Hue Kerberos Ticket Renewer in your cluster, you must configure your KDC to allow tickets to be renewed, and you must configure `krb5.conf` to request renewable tickets. Typically, you can do this by adding the `max_renewable_life` setting to your realm in `kdc.conf`, and by adding the `renew_lifetime` parameter to the `libdefaults` section of `krb5.conf`. For more information about renewable tickets, see the [Kerberos documentation](#). This is demonstrated in the [Sample Kerberos Configuration files: krb5.conf, kdc.conf, kadm5.acl](#) on page 194.

1. In the `/var/kerberos/krb5kdc/kdc.conf` file on the local dedicated KDC server host, configure the default realm for the Hadoop cluster by substituting your Kerberos realm in the following `realms` property:

```
[realms]
HADOOP.EXAMPLE.COM = {
```

2. In the `/etc/krb5.conf` file on all cluster hosts and all Hadoop client user hosts, configure the default realm for the Hadoop cluster by substituting your Kerberos realm in the following `realms` property. Also specify the local dedicated KDC server hostname in the `/etc/krb5.conf` file (for example, `kdc01.example.com`).

```
[libdefaults]
  default_realm = HADOOP.EXAMPLE.COM
[realms]
  HADOOP.EXAMPLE.COM = {
    kdc = kdc01.hadoop.example.com:88
    admin_server = kdc01.hadoop.example.com:749
    default_domain = hadoop.example.com
  }
  EXAMPLE.COM = {
    kdc = kdc01.example.com:88
    admin_server = kdc01.example.com:749
    default_domain = example.com
  }
[domain_realm]
  .hadoop.example.com = HADOOP.EXAMPLE.COM
  hadoop.example.com = HADOOP.EXAMPLE.COM
  .example.com = EXAMPLE.COM
  example.com = EXAMPLE.COM
```

3. To set up the cross-realm trust in the cluster-dedicated KDC, type the following command in the `kadmin.local` or `kadmin` shell on the cluster-dedicated KDC host to create a `krbtgt` principal. Substitute your cluster-dedicated KDC realm for `HADOOP.EXAMPLE.COM`, and substitute your central KDC realm for `EXAMPLE.COM`. Enter a trust password when prompted. Note the password because you will need to enter the exact same password in the central KDC in the next step.

```
kadmin: addprinc krbtgt/HADOOP.EXAMPLE.COM@EXAMPLE.COM
```

4. Each of your Hadoop client users must also place this information in their local `core-site.xml` file. The easiest way to do so is by using the Cloudera Manager Admin Console to generate a [client configuration file](#).
5. To set up the cross-realm trust in the central KDC, type the same command in the `kadmin.local` or `kadmin` shell on the central KDC host to create the exact same `krbtgt` principal and password.

```
kadmin: addprinc krbtgt/HADOOP.EXAMPLE.COM@EXAMPLE.COM
```



Important: In order for a cross-realm trust to operate properly, both KDCs must have the same `krbtgt` principal and password, and both KDCs must be configured to use the same encryption type.

6. To properly translate principal names from the central KDC realm into the cluster-dedicated KDC realm for the Hadoop cluster, configure the **Trusted Kerberos Realms** property of the HDFS service.
 - a. Open the Cloudera Manager Admin Console.
 - b. Go to the HDFS service.
 - c. Click the **Configuration** tab.
 - d. Select **Scope > HDFS (Service Wide)**
 - e. Select **Category > Security**.
 - f. Type `Kerberos` in the **Search** box.
 - g. Edit the **Trusted Kerberos Realms** property to add the name of your central KDC realm. If you need to use more advanced mappings which do more than just allow principals from another domain, you may enter them in the **Additional Rules to Map Kerberos Principals to Short Names** property. For more information about name mapping rules, see [Configuring the Mapping from Kerberos Principals to Short Names](#) on page 117.
7. Each of your Hadoop client users must also place this information in their local `core-site.xml` file. The easiest way to do so is by using the Cloudera Manager Admin Console to generate a [client configuration file](#).

- Proceed to [Step 2: If You are Using AES-256 Encryption, Install the JCE Policy File](#) on page 77. Later in this procedure, you will restart the services to have the configuration changes in `core-site.xml` take effect.

Using a Cluster-Dedicated MIT KDC with Active Directory



Important: If you are using Cloudera Manager, ensure you have installed the `openldap-clients` package on the Cloudera Manager Server host before you begin configuring Kerberos authentication.

On the Active Directory Server

- On the Active Directory server host, type the following command to add the local realm trust to Active Directory:

```
netdom trust HADOOP.EXAMPLE.COM /Domain:EXAMPLE.COM /add /realm /passwordt:TrustPassword
```

- On the Active Directory server host, type the following command to set the proper encryption type:

Windows 2003 RC2

Windows 2003 server installations do not support AES encryption for Kerberos. Therefore RC4 should be used. Please see the [Microsoft reference documentation](#) for more information.

```
ktpass /MITRealmName HADOOP.EXAMPLE.COM /TrustEncryp RC4
```

Windows 2008

```
ksetup /SetEncTypeAttr HADOOP.EXAMPLE.COM <enc_type>
```

Where the `<enc_type>` parameter can be replaced with parameter strings for AES, DES, or RC4 encryption modes. For example, for AES encryption, replace `<enc_type>` with `AES256-CTS-HMAC-SHA1-96` or `AES128-CTS-HMAC-SHA1-96` and for RC4 encryption, replace with `RC4-HMAC-MD5`. See the [Microsoft reference documentation](#) for more information.



Important: Make sure that the encryption type you specify is supported on both your version of Windows Active Directory and your version of MIT Kerberos.

On the MIT KDC Server

- In the `/var/kerberos/krb5kdc/kdc.conf` file on the local dedicated KDC server host, configure the default realm for the Hadoop cluster by substituting your Kerberos realm in the following `realms` property:

```
[realms]
HADOOP.EXAMPLE.COM = {
```

- Each of your Hadoop client users must also place this information in their local `core-site.xml` file. The easiest way to do so is by using the Cloudera Manager Admin Console to generate a [client configuration file](#).
- On the local MIT KDC server host, type the following command in the `kadmin.local` or `kadmin` shell to add the cross-realm `krbtgt` principal:

```
kadmin: addprinc -e "<keysalt_list>" krbtgt/HADOOP.EXAMPLE.COM@EXAMPLE.COM
```

where the `<keysalt_list>` parameter specifies the types of keys and their salt to be used for encryption of the password for this cross-realm `krbtgt` principal. It can be set to AES, or RC4 keytypes with a salt value of `:normal`. Note that DES is deprecated and should no longer be used. You can specify multiple keysalt types using the parameter in the command above. Make sure that at least one of the encryption types corresponds to the encryption types found in the tickets granted by the KDC in the remote realm. For an example of the values to use, see the examples based on the Active Directory functional domain level, below.

Examples by Active Directory Domain or Forest "Functional level"

Active Directory will, based on the Domain or Forest functional level, use encryption types supported by that release of the Windows Server operating system. It is not possible to use AES encryption types with an AD 2003 functional level. If you notice that DES encryption types are being used when authenticating or requesting service tickets to Active Directory then it might be necessary to enable weak encryption types in the `/etc/krb5.conf`. See [Sample Kerberos Configuration files: krb5.conf, kdc.conf, kadm5.acl](#) on page 194 for an example.

- **Windows 2003**

```
kadmin: addprinc -e "rc4-hmac:normal" krbtgt/HADOOP.EXAMPLE.COM@EXAMPLE.COM
```

- **Windows 2008**

```
kadmin: addprinc -e "aes256-cts:normal aes128-cts:normal rc4-hmac:normal"
krbtgt/HADOOP.EXAMPLE.COM@EXAMPLE.COM
```



Note: The cross-realm `krbtgt` principal that you add in this step must have *at least one entry* that uses the same encryption type as the tickets that are issued by the remote KDC. If there are no matching encryption types, principals in the local realm can successfully access the Hadoop cluster, but principals in the remote realm are unable to.

On All Cluster Hosts

1. In the `/etc/krb5.conf` file on all cluster hosts and all Hadoop client user hosts, configure both Kerberos realms. Note that `default_realm` should be configured as the local MIT Kerberos realm for the cluster. Your `krb5.conf` may contain more configuration properties than those demonstrated below. This example is provided to clarify configuration parameters. See [Sample Kerberos Configuration files: krb5.conf, kdc.conf, kadm5.acl](#) on page 194 for more information.

```
[libdefaults]
default_realm = HADOOP.EXAMPLE.COM
[realms]
EXAMPLE.COM = {
    kdc = dc01.example.com:88
    admin_server = dc01.example.com:749
}
HADOOP.EXAMPLE.COM = {
    kdc = kdc01.hadoop.example.com:88
    admin_server = kdc01.hadoop.example.com:749
}
[domain_realm]
.hadoop.example.com = HADOOP.EXAMPLE.COM
hadoop.example.com = HADOOP.EXAMPLE.COM
.example.com = EXAMPLE.COM
example.com = EXAMPLE.COM
```

2. Use one of the following methods to properly translate principal names from the Active Directory realm into the cluster-dedicated KDC realm for the Hadoop cluster.
 - **Using Cloudera Manager:** Configure the **Trusted Kerberos realms** property of the HDFS service:
 1. Open the Cloudera Manager Admin Console.
 2. Go to the HDFS service.
 3. Click the **Configuration** tab.
 4. Select **Scope > HDFS (Service Wide)**
 5. Select **Category > Security**.
 6. Type `Kerberos` in the **Search** box.
 7. Edit the **Trusted Kerberos Realms** property to add the name of your central KDC realm. If you need to use more advanced mappings which do more than just allow principals from another domain, you may

enter them in the **Additional Rules to Map Kerberos Principals to Short Names** property. For more information about name mapping rules, see [Configuring the Mapping from Kerberos Principals to Short Names](#) on page 117.

- **Using the Command Line:** Configure the `hadoop.security.auth_to_local` setting in the `core-site.xml` file on all of the cluster hosts. The following example translates all principal names with the realm `EXAMPLE.COM` into the first component of the principal name only. It also preserves the standard translation for the default realm (the cluster realm).

```
<property>
  <name>hadoop.security.auth_to_local</name>
  <value>
    RULE:[1:$1@$0](^.*@EXAMPLE\.COM$)s/^(.*)@EXAMPLE\.COM$/$1/g
    RULE:[2:$1@$0](^.*@EXAMPLE\.COM$)s/^(.*)@EXAMPLE\.COM$/$1/g
    DEFAULT
  </value>
</property>
```

Integrating Hadoop Security with Active Directory

Considerations when using an Active Directory KDC

Performance:

As your cluster grows, so will the volume of Authentication Service (AS) and Ticket Granting Service (TGS) interaction between the services on each cluster server. Consider evaluating the volume of this interaction against the Active Directory domain controllers you have configured for the cluster before rolling this feature out to a production environment. If cluster performance suffers, over time it might become necessary to dedicate a set of AD domain controllers to larger deployments.

Network Proximity:

By default, Kerberos uses UDP for client/server communication. Often, AD services are in a different network than project application services such as Hadoop. If the domain controllers supporting a cluster for Kerberos are not in the same subnet, or they're separated by a firewall, consider using the `udp_preference_limit = 1` setting in the `[libdefaults]` section of the `krb5.conf` used by cluster services. Cloudera strongly recommends *against* using AD domain controller (KDC) servers that are separated from the cluster by a WAN connection, as latency in this service will significantly impact cluster performance.

Process:

Troubleshooting the cluster's operations, especially for Kerberos-enabled services, will need to include AD administration resources. Evaluate your organizational processes for engaging the AD administration team, and how to escalate in case a cluster outage occurs due to issues with Kerberos authentication against AD services. In some situations it might be necessary to [enable Kerberos event logging](#) to address desktop and KDC issues within windows environments.

Also note that if you decommission any Cloudera Manager roles or nodes, the related AD accounts will need to be deleted manually. This is required because Cloudera Manager will not delete existing entries in Active Directory.



Important: With CDH 5.1 and higher, clusters managed by Cloudera Manager 5.1 (and higher) do not require a local MIT KDC and are able to integrate directly with an Active Directory KDC. Cloudera recommends you use a direct-to-AD setup. For instructions, see [Enabling Kerberos Authentication Using the Wizard](#) on page 59.

If direct integration with AD is not currently possible, use the following instructions to configure a local MIT KDC to trust your AD server:

1. Run an MIT Kerberos KDC and realm local to the cluster and create all service principals in this realm.

2. Set up one-way cross-realm trust from this realm to the Active Directory realm. Using this method, there is no need to create service principals in Active Directory, but Active Directory principals (users) can be authenticated to Hadoop. See [Configuring a Local MIT Kerberos Realm to Trust Active Directory](#) on page 186.

Configuring a Local MIT Kerberos Realm to Trust Active Directory

On the Active Directory Server

1. Add the local realm trust to Active Directory with this command:

```
netdom trust YOUR-LOCAL-REALM.COMPANY.COM /Domain:AD-REALM.COMPANY.COM /add /realm /passwordt:<TrustPassword>
```

2. Set the proper encryption type with this command:

On Windows 2003 RC2:

```
ktpass /MITRealmName YOUR-LOCAL-REALM.COMPANY.COM /TrustEncryp <enc_type>
```

On Windows 2008:

```
ksetup /SetEncTypeAttr YOUR-LOCAL-REALM.COMPANY.COM <enc_type>
```

The <enc_type> parameter specifies AES, DES, or RC4 encryption. Refer to the documentation for your version of Windows Active Directory to find the <enc_type> parameter string to use.

3. Get and verify the list of encryption types set with this command:

On Windows 2008:

```
ksetup /GetEncTypeAttr YOUR-LOCAL-REALM.COMPANY.COM
```



Important: Make sure the encryption type you specify is supported on both your version of Windows Active Directory and your version of MIT Kerberos.

On the MIT KDC Server

Type the following command in the `kadmin.local` or `kadmin` shell to add the cross-realm `krbtgt` principal. Use the same password you used in the `netdom` command on the Active Directory Server.

```
kadmin: addprinc -e "<enc_type_list>"  
krbtgt/YOUR-LOCAL-REALM.COMPANY.COM@AD-REALM.COMPANY.COM
```

where the <enc_type_list> parameter specifies the types of encryption this cross-realm `krbtgt` principal will support: either AES, DES, or RC4 encryption. You can specify multiple encryption types using the parameter in the command above, what's important is that at least one of the encryption types corresponds to the encryption type found in the tickets granted by the KDC in the remote realm. For example:

```
kadmin: addprinc -e "rc4-hmac:normal des3-hmac-sha1:normal"  
krbtgt/YOUR-LOCAL-REALM.COMPANY.COM@AD-REALM.COMPANY.COM
```



Note: The cross-realm `krbtgt` principal that you add in this step must have *at least one entry* that uses the same encryption type as the tickets that are issued by the remote KDC. If no entries have the same encryption type, then the problem you will see is that authenticating as a principal in the local realm will allow you to successfully run Hadoop commands, but authenticating as a principal in the remote realm will not allow you to run Hadoop commands.

On All of the Cluster Hosts

1. Verify that both Kerberos realms are configured on all of the cluster hosts. Note that the default realm and the domain realm should remain set as the MIT Kerberos realm which is local to the cluster.

```
[realms]
AD-REALM.CORP.FOO.COM = {
  kdc = ad.corp.foo.com:88
  admin_server = ad.corp.foo.com:749
  default_domain = foo.com
}
CLUSTER-REALM.CORP.FOO.COM = {
  kdc = cluster01.corp.foo.com:88
  admin_server = cluster01.corp.foo.com:749
  default_domain = foo.com
}
```

2. To properly translate principal names from the Active Directory realm into local names within Hadoop, you must configure the `hadoop.security.auth_to_local` setting in the `core-site.xml` file on all of the cluster machines. The following example translates all principal names with the realm `AD-REALM.CORP.FOO.COM` into the first component of the principal name only. It also preserves the standard translation for the default realm (the cluster realm).

```
<property>
  <name>hadoop.security.auth_to_local</name>
  <value>
    RULE:[1:$1@$0](^.*@AD-REALM\.CORP\.FOO\.COM$)s/^.*(.*@AD-REALM\.CORP\.FOO\.COM$)/$1/g

    RULE:[2:$1@$0](^.*@AD-REALM\.CORP\.FOO\.COM$)s/^.*(.*@AD-REALM\.CORP\.FOO\.COM$)/$1/g

    DEFAULT
  </value>
</property>
```

For more information about name mapping rules, see [Configuring the Mapping from Kerberos Principals to Short Names](#) on page 117.

Integrating Hadoop Security with Alternate Authentication

One of the ramifications of enabling security on a Hadoop cluster is that every user who interacts with the cluster must have a Kerberos principal configured. For some of the services, specifically Oozie and Hadoop (for example, JobTracker and TaskTracker), it can be convenient to run a mixed form of authentication where Kerberos authentication is used for API or command line access while some other form of authentication (for example, SSO and LDAP) is used for accessing Web UIs. Using an alternate authentication deployment is considered an advanced topic because only a partial implementation is provided in this release: you will have to implement some of the code yourself.



Note: The following instructions assume you already have a Kerberos-enabled cluster.

Proceed as follows:

- [Configuring the AuthenticationFilter to use Kerberos](#) on page 188
- [Creating an AltKerberosAuthenticationHandler Subclass](#) on page 188
- [Enabling Your AltKerberosAuthenticationHandler Subclass](#) on page 188

See also the [Example Implementation for Oozie](#) on page 189.

Configuring the AuthenticationFilter to use Kerberos

First, you must do all of the steps in the Server Side Configuration section of the [Hadoop Auth, Java HTTP SPNEGO Documentation](#) to configure `AuthenticationFilter` to use Kerberos. You must configure `AuthenticationFilter` to use Kerberos before doing the steps below.

Creating an `AltKerberosAuthenticationHandler` Subclass

An `AuthenticationHandler` is installed on the server-side to handle authenticating clients and creating an `AuthenticationToken`.

1. Subclass the `org.apache.hadoop.security.authentication.server.AltKerberosAuthenticationHandler` class (in the `hadoop-auth` package).
2. When a client sends a request, the `authenticate` method will be called. For browsers, `AltKerberosAuthenticationHandler` will call the `alternateAuthenticate` method, which is what you need to implement to interact with the desired authentication mechanism. For non-browsers, `AltKerberosAuthenticationHandler` will follow the Kerberos SPNEGO sequence (this is provided for you).
3. The `alternateAuthenticate(HttpServletRequest request, HttpServletResponse response)` method in your subclass should follow these rules:
4. Return `null` if the authentication is still in progress; the `response` object can be used to interact with the client.
5. Throw an `AuthenticationException` if the authentication failed.
6. Return an `AuthenticationToken` if the authentication completed successfully.

Enabling Your `AltKerberosAuthenticationHandler` Subclass

You can enable the alternate authentication on Hadoop Web UIs, Oozie Web UIs, or both. You will need to include a JAR containing your subclass on the classpath of Hadoop or Oozie. All Kerberos-related configuration properties will still apply.

Enabling Your `AltKerberosAuthenticationHandler` Subclass on Hadoop Web UIs

1. Stop Hadoop by running the following command on every node in your cluster (as root):

```
$ for x in `cd /etc/init.d ; ls hadoop-*` ; do sudo service $x stop ; done
```

2. Set the following property in `core-site.xml`, where `org.my.subclass.of.AltKerberosAuthenticationHandler` is the classname of your subclass:

```
<property>  
  <name>hadoop.http.authentication.type</name>  
  <value>org.my.subclass.of.AltKerberosAuthenticationHandler</value>  
</property>
```

3. (Optional) You can also specify which user-agents you do not want to be considered as browsers by setting the following property as required (default value is shown). Note that all Java-based programs (such as Hadoop client) will use `java` as their user-agent.

```
<property>  
  <name>hadoop.http.authentication.alt-kerberos.non-browser.user-agents</name>  
  <value>java,curl,wget,perl</value>  
</property>
```

4. Copy the JAR containing your subclass into `/usr/lib/hadoop/lib/`.
5. Start Hadoop by running the following command:

```
$ for x in `cd /etc/init.d ; ls hadoop-*` ; do sudo service $x start ; done
```

Enabling Your AltKerberosAuthenticationHandler Subclass on Oozie Web UI

**Note:**

These instructions assume you have already performed the installation and configuration steps in [Oozie Security Configuration](#).

1. Stop the Oozie Server:

```
sudo /sbin/service oozie stop
```

2. Set the following property in oozie-site.xml, where

`org.my.subclass.of.AltKerberosAuthenticationHandler` is the classname of your subclass:

```
<property>
  <name>oozie.authentication.type</name>
  <value>org.my.subclass.of.AltKerberosAuthenticationHandler</value>
</property>
```

3. (Optional) You can also specify which user-agents you do not want to be considered as browsers by setting the following property as required (default value is shown). Note that all Java-based programs (such as Hadoop client) will use java as their user-agent.

```
<property>
  <name>oozie.authentication.alt-kerberos.non-browser.user-agents</name>
  <value>java,curl,wget,perl</value>
</property>
```

4. Copy the JAR containing your subclass into /var/lib/oozie.**5. Start the Oozie Server:**

```
sudo /sbin/service oozie start
```

Example Implementation for Oozie

**Warning:**

The example implementation is **NOT SECURE**. Its purpose is to be as simple as possible, as an example of how to write your own `AltKerberosAuthenticationHandler` subclass.

It should NOT be used in a production environment

An example implementation of `AltKerberosAuthenticationHandler` is included (though not built by default) with Oozie. Also included is a simple Login Server with two implementations. The first one will authenticate any user who is using a username and password that are identical, such as `foo:foo`. The second one can be configured against an LDAP server to use LDAP for authentication.

You can read comprehensive documentation on the example at [Creating Custom Authentication](#).



Important:

If you installed Oozie from the CDH packages and are deploying `oozie-login.war` alongside `oozie.war`, you will also need to run the following commands after you copy the `oozie-login.war` file to `/usr/lib/oozie/oozie-server` (if using YARN or `/usr/lib/oozie/oozie-server-0.20` if using MRv1) because it won't automatically be expanded:

```
jar xvf oozie-login.war
mkdir oozie-login
mv META-INF oozie-login/
mv WEB-INF oozie-login/
```

Authenticating Kerberos Principals in Java Code

This topic provides an example of how to authenticate a Kerberos principal in a Java application using the `org.apache.hadoop.security.UserGroupInformation` class.

The following code snippet authenticates the `cloudera` principal using the `cloudera.keytab` file:

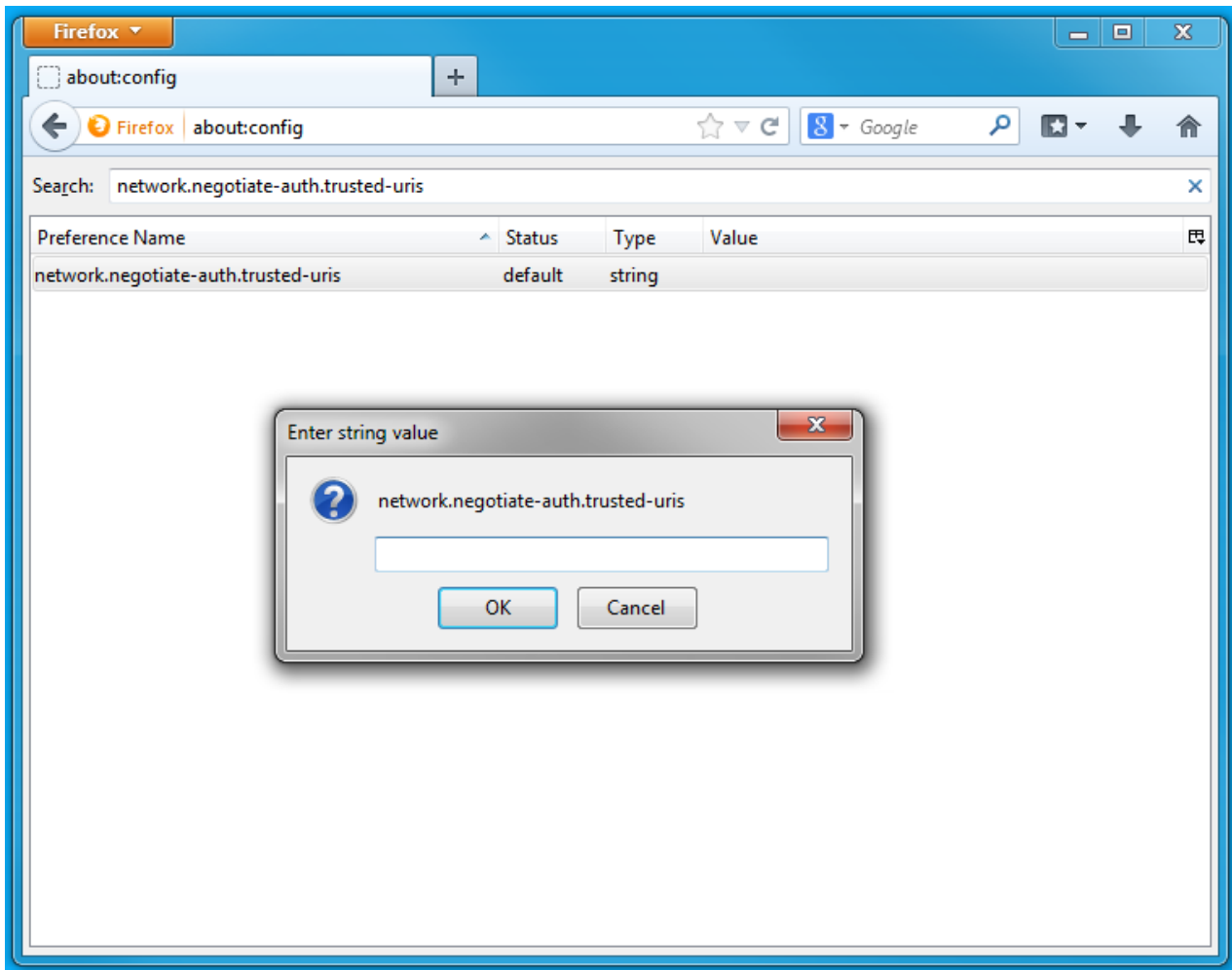
```
// Authenticating Kerberos principal
System.out.println("Principal Authentication: ");
final String user = "cloudera@CLOUDERA.COM";
final String keyPath = "cloudera.keytab";
UserGroupInformation.loginUserFromKeytab(user, keyPath);
```

Using a Web Browser to Access an URL Protected by Kerberos HTTP SPNEGO

To access a URL protected by Kerberos HTTP SPNEGO, use the following instructions for the browser you are using.

To configure Mozilla Firefox:

1. Open the low level Firefox configuration page by loading the `about:config` page.
2. In the **Search** text box, enter: `network.negotiate-auth.trusted-uris`
3. Double-click the `network.negotiate-auth.trusted-uris` preference and enter the hostname or the domain of the web server that is protected by Kerberos HTTP SPNEGO. Separate multiple domains and hostnames with a comma.
4. Click **OK**.

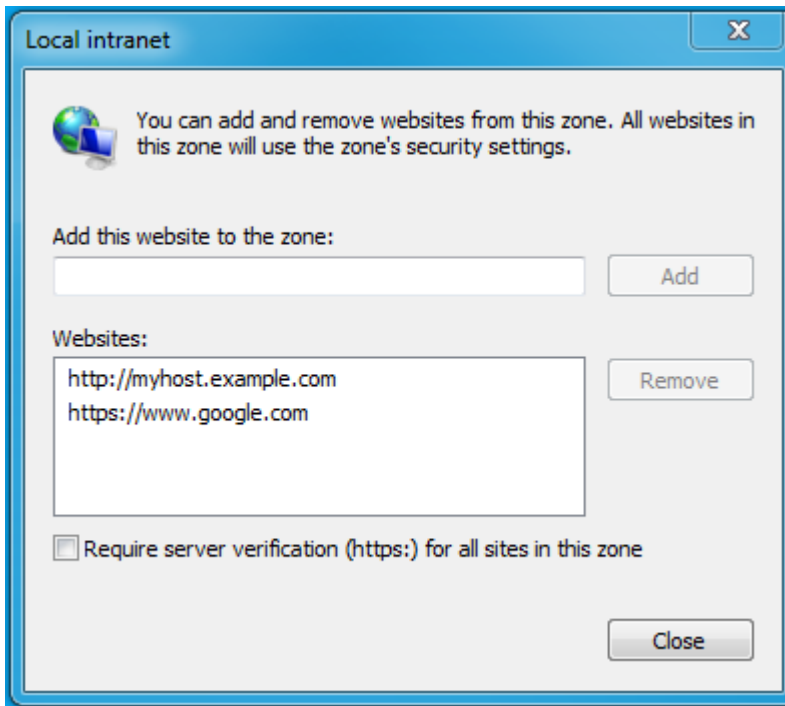


To configure Internet Explorer:

Follow the instructions given below to configure Internet Explorer to access URLs protected by

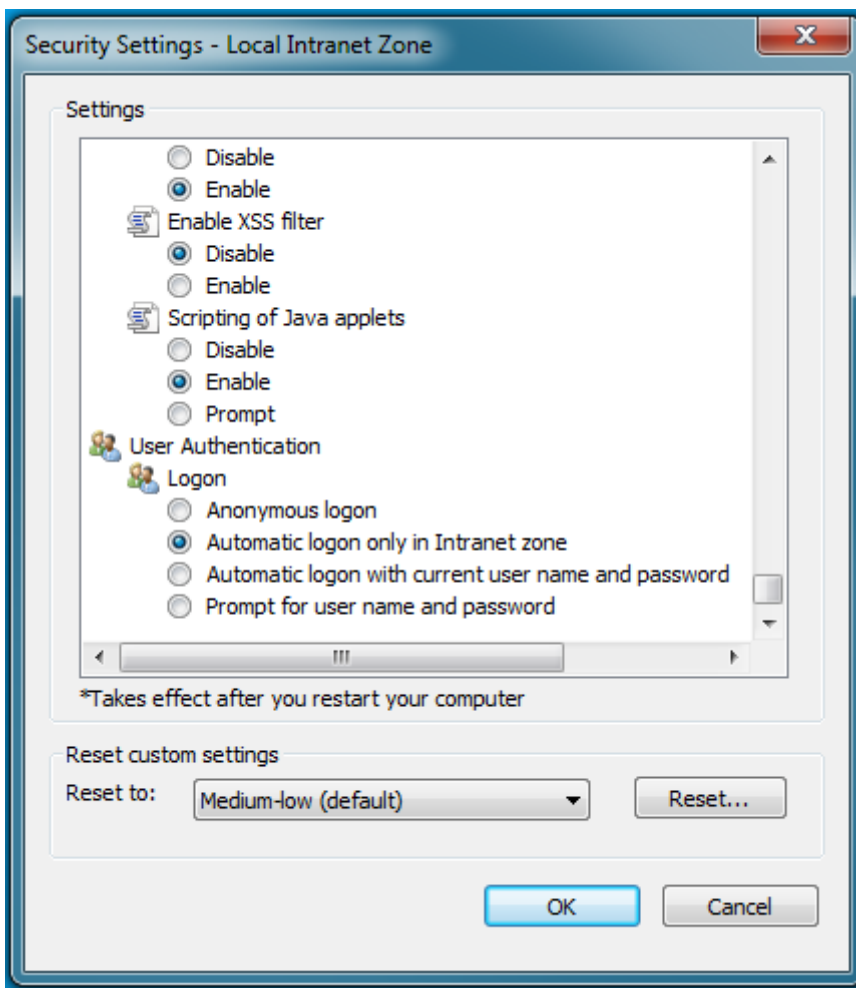
Configuring the Local Intranet Domain

1. Open Internet Explorer and click the Settings "gear" icon in the top-right corner. Select **Internet options**.
2. Select the **Security** tab.
3. Select the **Local Intranet** zone and click the **Sites** button.
4. Make sure that the first two options, **Include all local (intranet) sites not listed in other zones** and **Include all sites that bypass the proxy server** are checked.
5. Click **Advanced** and add the names of the domains that are protected by Kerberos HTTP SPNEGO, one at a time, to the list of websites. For example, `myhost.example.com`. Click **Close**.
6. Click **OK** to save your configuration changes.



Configuring Intranet Authentication

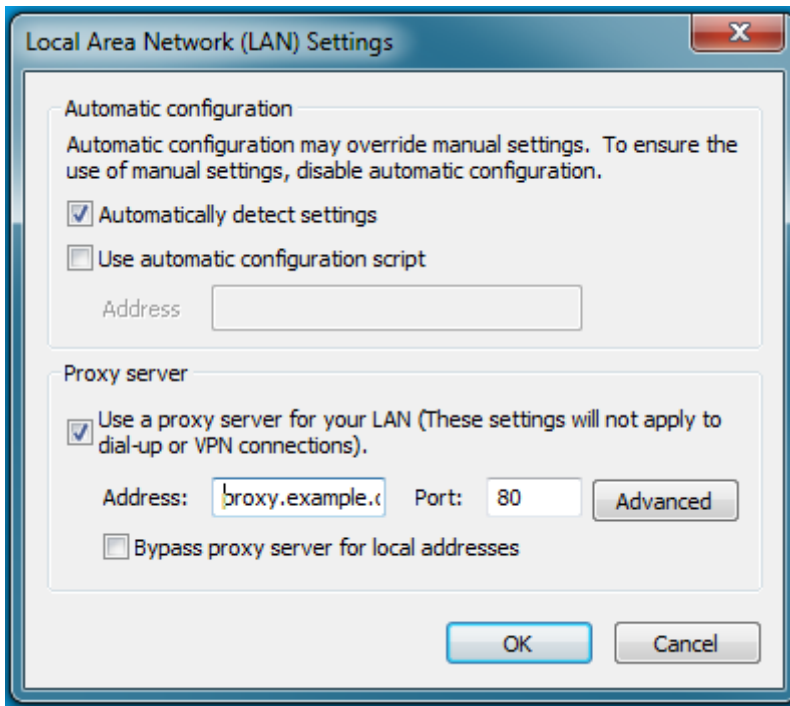
1. Click the Settings "gear" icon in the top-right corner. Select **Internet options**.
2. Select the **Security** tab.
3. Select the **Local Intranet** zone and click the **Custom level...** button to open the **Security Settings - Local Intranet Zone** dialog box.
4. Scroll down to the **User Authentication** options and select **Automatic logon only in Intranet zone**.
5. Click **OK** to save these changes.



Verifying Proxy Settings

You need to perform the following steps only if you have a proxy server already enabled.

1. Click the Settings "gear" icon in the top-right corner. Select **Internet options**.
2. Select the **Connections** tab and click **LAN Settings**.
3. Verify that the proxy server **Address** and **Port** number settings are correct.
4. Click **Advanced** to open the **Proxy Settings** dialog box.
5. Add the Kerberos-protected domains to the **Exceptions** field.
6. Click **OK** to save any changes.



To configure Google Chrome:

If you are using Windows, use the Control Panel to go to the **Internet Options** dialog box. Configuration changes required are the same as those described above for Internet Explorer.

On MacOS or Linux, add the `--auth-server-whitelist` parameter to the `google-chrome` command. For example, to run Chrome from a Linux prompt, run the `google-chrome` command as follows,

```
> google-chrome --auth-server-whitelist = "hostname/domain"
```

Troubleshooting Authentication Issues

Typically, if there are problems with security, Hadoop will display generic messages about the cause of the problem. This topic contains some sample Kerberos configuration files for your reference. It also has solutions to potential problems you might face when configuring a secure cluster:

Sample Kerberos Configuration files: krb5.conf, kdc.conf, kadm5.acl

kdc.conf:

```
[kdcdefaults]
kdc_ports = 88
kdc_tcp_ports = 88

[realms]
EXAMPLE.COM = {
  #master_key_type = aes256-cts
  max_renewable_life = 7d 0h 0m 0s
  acl_file = /var/kerberos/krb5kdc/kadm5.acl
  dict_file = /usr/share/dict/words
  admin_keytab = /var/kerberos/krb5kdc/kadm5.keytab
  # note that aes256 is ONLY supported in Active Directory in a domain / forrest operating
  # at a 2008 or greater functional level.
  # aes256 requires that you download and deploy the JCE Policy files for your JDK release
  # level to provide
```

```
# strong java encryption extension levels like AES256. Make sure to match based on the
# encryption configured within AD for
# cross realm auth, note that RC4 = arcfour when comparing windows and linux encyptes
supported_encyptes = aes256-cts:normal aes128-cts:normal arcfour-hmac:normal
default_principal_flags = +renewable, +forwardable
}
```

krb5.conf:

```
[logging]
default = FILE:/var/log/krb5libs.log
kdc = FILE:/var/log/krb5kdc.log
admin_server = FILE:/var/log/kadmind.log

[libdefaults]
default_realm = EXAMPLE.COM
dns_lookup_realm = false
dns_lookup_kdc = false
ticket_lifetime = 24h
renew_lifetime = 7d
forwardable = true
# udp_preference_limit = 1

# set udp_preference_limit = 1 when TCP only should be
# used. Consider using in complex network environments when
# troubleshooting or when dealing with inconsistent
# client behavior or GSS (63) messages.

# uncomment the following if AD cross realm auth is ONLY providing DES encrypted tickets
# allow-weak-crypto = true

[realms]
AD-REALM.EXAMPLE.COM = {
    kdc = AD1.ad-realm.example.com:88
    kdc = AD2.ad-realm.example.com:88
    admin_server = AD1.ad-realm.example.com:749
    admin_server = AD2.ad-realm.example.com:749
    default_domain = ad-realm.example.com
}
EXAMPLE.COM = {
    kdc = kdc1.example.com:88
    admin_server = kdc1.example.com:749
    default_domain = example.com
}

# The domain_realm is critical for mapping your host domain names to the kerberos realms
# that are servicing them. Make sure the lowercase left hand portion indicates any
# domains or subdomains
# that will be related to the kerberos REALM on the right hand side of the expression.
# REALMs will
# always be UPPERCASE. For example, if your actual DNS domain was test.com but your
# kerberos REALM is
# EXAMPLE.COM then you would have,

[domain_realm]
test.com = EXAMPLE.COM
#AD domains and realms are usually the same
ad-domain.example.com = AD-REALM.EXAMPLE.COM
ad-realm.example.com = AD-REALM.EXAMPLE.COM
```

kadm5.acl:

```
*/admin@HADOOP.COM *
cloudera-scm@HADOOP.COM * flume/*@HADOOP.COM
cloudera-scm@HADOOP.COM * hbase/*@HADOOP.COM
cloudera-scm@HADOOP.COM * hdfs/*@HADOOP.COM
cloudera-scm@HADOOP.COM * hive/*@HADOOP.COM
cloudera-scm@HADOOP.COM * httpfs/*@HADOOP.COM
cloudera-scm@HADOOP.COM * HTTP/*@HADOOP.COM
cloudera-scm@HADOOP.COM * hue/*@HADOOP.COM
```

```
cloudera-scm@HADOOP.COM * impala/*@HADOOP.COM
cloudera-scm@HADOOP.COM * mapred/*@HADOOP.COM
cloudera-scm@HADOOP.COM * oozie/*@HADOOP.COM
cloudera-scm@HADOOP.COM * solr/*@HADOOP.COM
cloudera-scm@HADOOP.COM * sqoop/*@HADOOP.COM
cloudera-scm@HADOOP.COM * yarn/*@HADOOP.COM
cloudera-scm@HADOOP.COM * zookeeper/*@HADOOP.COM
```

Potential Security Problems and Their Solutions

This Troubleshooting section contains sample Kerberos configuration files, `krb5.conf` and `kdc.conf` for your reference. It also has solutions to potential problems you might face when configuring a secure cluster:

Issues with Generate Credentials

Cloudera Manager uses a command called **Generate Credentials** to create the accounts needed by CDH for enabling authentication using Kerberos. The command is triggered automatically when you are using the Kerberos Wizard or making changes to your cluster that will require new Kerberos principals.

When configuring Kerberos, if CDH services do not start, and on the Cloudera Manager **Home > Status** tab you see a validation error, `Role is missing Kerberos keytab`, it means the **Generate Credentials** command failed. To see the output of the command, go to the **Home > Status** tab and click the **All Recent Commands** tab.

Here are some common error messages:

Problems	Possible Causes	Solutions
With Active Directory		
<code>ldap_sasl_interactive_bind_s: Can't contact LDAP server (-1)</code>	The Domain Controller specified is incorrect or LDAPS has not been enabled for it.	Verify the KDC configuration by going to the Cloudera Manager Admin Console and go to Administration > Settings > Kerberos . Also check that LDAPS is enabled for Active Directory.
<code>ldap_add: Insufficient access (50)</code>	The Active Directory account you are using for Cloudera Manager does not have permissions to create other accounts.	Use the Delegate Control wizard to grant permission to the Cloudera Manager account to create other accounts. You can also login to Active Directory as the Cloudera Manager user to check that it can create other accounts in your Organizational Unit.
With MIT KDC		
<code>kadmin: Cannot resolve network address for admin server in requested realm while initializing kadmin interface.</code>	The hostname for the KDC server is incorrect.	Check the <code>kdc</code> field for your default realm in <code>krb5.conf</code> and make sure the hostname is correct.

Running any Hadoop command fails after enabling security.

Description:

A user must have a valid Kerberos ticket in order to interact with a secure Hadoop cluster. Running any Hadoop command (such as `hadoop fs -ls`) will fail if you do not have a valid Kerberos ticket in your credentials cache. If you do not have a valid ticket, you will receive an error such as:

```
11/01/04 12:08:12 WARN ipc.Client: Exception encountered while connecting to the server
: javax.security.sasl.SaslException:
GSS initiate failed [Caused by GSSException: No valid credentials provided (Mechanism
level: Failed to find any Kerberos tgt)]
```

```
Bad connection to FS. command aborted. exception: Call to nn-host/10.0.0.2:8020 failed
on local exception: java.io.IOException:
javax.security.sasl.SaslException: GSS initiate failed [Caused by GSSException: No valid
credentials provided (Mechanism level: Failed to find any Kerberos tgt)]
```

Solution:

You can examine the Kerberos tickets currently in your credentials cache by running the `klist` command. You can obtain a ticket by running the `kinit` command and either specifying a keytab file containing credentials, or entering the password for your principal.

Java is unable to read the Kerberos credentials cache created by versions of MIT Kerberos 1.8.1 or higher.

Description:

If you are running MIT Kerberos 1.8.1 or higher, the following error will occur when you attempt to interact with the Hadoop cluster, even after successfully obtaining a Kerberos ticket using `kinit`:

```
11/01/04 12:08:12 WARN ipc.Client: Exception encountered while connecting to the server
: javax.security.sasl.SaslException:
GSS initiate failed [Caused by GSSException: No valid credentials provided (Mechanism
level: Failed to find any Kerberos tgt)]
Bad connection to FS. command aborted. exception: Call to nn-host/10.0.0.2:8020 failed
on local exception: java.io.IOException:
javax.security.sasl.SaslException: GSS initiate failed [Caused by GSSException: No valid
credentials provided (Mechanism level: Failed to find any Kerberos tgt)]
```

Because of a change [1] in the format in which MIT Kerberos writes its credentials cache, there is a bug [2] in the Oracle JDK 6 Update 26 and earlier that causes Java to be unable to read the Kerberos credentials cache created by versions of MIT Kerberos 1.8.1 or higher. Kerberos 1.8.1 is the default in Ubuntu Lucid and higher releases and Debian Squeeze and higher releases. (On RHEL and CentOS, an older version of MIT Kerberos which does not have this issue, is the default.)

Footnotes:

[1] MIT Kerberos change: <http://krbdev.mit.edu/rt/Ticket/Display.html?id=6206>

[2] Report of bug in Oracle JDK 6 Update 26 and lower:
http://bugs.sun.com/bugdatabase/view_bug.do?bug_id=6979329

Solution:

If you encounter this problem, you can work around it by running `kinit -R` after running `kinit` initially to obtain credentials. Doing so will cause the ticket to be renewed, and the credentials cache rewritten in a format which Java can read. To illustrate this:

```
$ klist
klist: No credentials cache found (ticket cache FILE:/tmp/krb5cc_1000)
$ hadoop fs -ls
11/01/04 13:15:51 WARN ipc.Client: Exception encountered while connecting to the server
: javax.security.sasl.SaslException:
GSS initiate failed [Caused by GSSException: No valid credentials provided (Mechanism
level: Failed to find any Kerberos tgt)]
Bad connection to FS. command aborted. exception: Call to nn-host/10.0.0.2:8020 failed
on local exception: java.io.IOException:
javax.security.sasl.SaslException: GSS initiate failed [Caused by GSSException: No valid
credentials provided (Mechanism level: Failed to find any Kerberos tgt)]
$ kinit
Password for atm@YOUR-REALM.COM:
$ klist
Ticket cache: FILE:/tmp/krb5cc_1000
Default principal: atm@YOUR-REALM.COM

Valid starting    Expires          Service principal
01/04/11 13:19:31 01/04/11 23:19:31  krbtgt/YOUR-REALM.COM@YOUR-REALM.COM
```

```
renew until 01/05/11 13:19:30
$ hadoop fs -ls
11/01/04 13:15:59 WARN ipc.Client: Exception encountered while connecting to the server
: javax.security.sasl.SaslException:
GSS initiate failed [Caused by GSSException: No valid credentials provided (Mechanism
level: Failed to find any Kerberos tgt)]
Bad connection to FS. command aborted. exception: Call to nn-host/10.0.0.2:8020 failed
on local exception: java.io.IOException:
javax.security.sasl.SaslException: GSS initiate failed [Caused by GSSException: No valid
credentials provided (Mechanism level: Failed to find any Kerberos tgt)]
$ kinit -R
$ hadoop fs -ls
Found 6 items
drwx----- - atm atm          0 2011-01-02 16:16 /user/atm/.staging
```



Note:

This workaround for Problem 2 requires the initial ticket to be renewable. Note that whether or not you can obtain renewable tickets is dependent upon a KDC-wide setting, as well as a per-principal setting for both the principal in question and the Ticket Granting Ticket (TGT) service principal for the realm. A non-renewable ticket will have the same values for its "valid starting" and "renew until" times. If the initial ticket is not renewable, the following error message is displayed when attempting to renew the ticket:

```
kinit: Ticket expired while renewing credentials
```

java.io.IOException: Incorrect permission

Description:

An error such as the following example is displayed if the user running one of the Hadoop daemons has a umask of 0002, instead of 0022:

```
java.io.IOException: Incorrect permission for
/var/folders/B3/B3d2vCm4F+mmWzVPB89W6E+++TI/-Tmp-/tmpYTil84/dfs/data/data1,
expected: rwxr-xr-x, while actual: rwxrwxr-x
    at org.apache.hadoop.util.DiskChecker.checkPermission(DiskChecker.java:107)
    at
org.apache.hadoop.util.DiskChecker.mkdirsWithExistsAndPermissionCheck(DiskChecker.java:144)
        at org.apache.hadoop.util.DiskChecker.checkDir(DiskChecker.java:160)
        at org.apache.hadoop.hdfs.server.datanode.DataNode.makeInstance(DataNode.java:1484)
    at
org.apache.hadoop.hdfs.server.datanode.DataNode.instantiateDataNode(DataNode.java:1432)
        at
org.apache.hadoop.hdfs.server.datanode.DataNode.instantiateDataNode(DataNode.java:1408)
            at org.apache.hadoop.hdfs.MinidFScluster.startDataNodes(MinidFScluster.java:418)
                at org.apache.hadoop.hdfs.MinidFScluster.<init>(MinidFScluster.java:279)
                at org.apache.hadoop.hdfs.MinidFScluster.<init>(MinidFScluster.java:203)
            at
org.apache.hadoop.test.MinidHadoopClusterManager.start(MinidHadoopClusterManager.java:152)
                at
org.apache.hadoop.test.MinidHadoopClusterManager.run(MinidHadoopClusterManager.java:129)
                at
org.apache.hadoop.test.MinidHadoopClusterManager.main(MinidHadoopClusterManager.java:308)
                    at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
                    at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:39)
                at
```

```

sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:25)
  at java.lang.reflect.Method.invoke(Method.java:597)
  at
org.apache.hadoop.util.ProgramDriver$ProgramDescription.invoke(ProgramDriver.java:68)
  at org.apache.hadoop.util.ProgramDriver.driver(ProgramDriver.java:139)
  at org.apache.hadoop.test.AllTestDriver.main(AllTestDriver.java:83)
  at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
  at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:39)
  at
sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:25)
  at java.lang.reflect.Method.invoke(Method.java:597)
  at org.apache.hadoop.util.RunJar.main(RunJar.java:186)

```

Solution:

Make sure that the `umask` for `hdfs` and `mapred` is `0022`.

A cluster fails to run jobs after security is enabled.

Description:

A cluster that was previously configured to *not* use security may fail to run jobs for certain users on certain TaskTrackers (MRv1) or NodeManagers (YARN) after security is enabled due to the following sequence of events:

1. A cluster is at some point in time configured *without* security enabled.
2. A user X runs some jobs on the cluster, which creates a local user directory on each TaskTracker or NodeManager.
3. Security is enabled on the cluster.
4. User X tries to run jobs on the cluster, and the local user directory on (potentially a subset of) the TaskTrackers or NodeManagers is owned by the wrong user or has overly-permissive permissions.

The bug is that after step 2, the local user directory on the TaskTracker or NodeManager should be cleaned up, but isn't.

If you're encountering this problem, you may see errors in the TaskTracker or NodeManager logs. The following example is for a TaskTracker on MRv1:

```

10/11/03 01:29:55 INFO mapred.JobClient: Task Id : attempt_201011021321_0004_m_000011_0,
  Status : FAILED
Error initializing attempt_201011021321_0004_m_000011_0:
java.io.IOException: org.apache.hadoop.util.Shell$ExitCodeException:
at org.apache.hadoop.mapred.LinuxTaskController.runCommand(LinuxTaskController.java:212)
  at
org.apache.hadoop.mapred.LinuxTaskController.initializeUser(LinuxTaskController.java:442)
  at
org.apache.hadoop.mapreduce.server.tasktracker.Localizer.initializeUserDirs(Localizer.java:272)
  at org.apache.hadoop.mapred.TaskTracker.localizeJob(TaskTracker.java:963)
  at org.apache.hadoop.mapred.TaskTracker.startNewTask(TaskTracker.java:2209)
  at org.apache.hadoop.mapred.TaskTracker$TaskLauncher.run(TaskTracker.java:2174)
Caused by: org.apache.hadoop.util.Shell$ExitCodeException:
at org.apache.hadoop.util.Shell.runCommand(Shell.java:250)
at org.apache.hadoop.util.Shell.run(Shell.java:177)
at org.apache.hadoop.util.Shell$ShellCommandExecutor.execute(Shell.java:370)
at org.apache.hadoop.mapred.LinuxTaskController.runCommand(LinuxTaskController.java:203)
... 5 more

```

Solution:

Delete the `mapred.local.dir` or `yarn.nodemanager.local-dirs` directories for that user across the cluster.

Configuring Authentication

The NameNode does not start and KrbException Messages (906) and (31) are displayed.

Description:

When you attempt to start the NameNode, a login failure occurs. This failure prevents the NameNode from starting and the following KrbException messages are displayed:

```
Caused by: KrbException: Integrity check on decrypted field failed (31) - PREAUTH_FAILED}}
```

and

```
Caused by: KrbException: Identifier doesn't match expected value (906)
```



Note:

These KrbException error messages are displayed only if you enable debugging output. See [Enabling Debugging Output for the Sun Kerberos Classes](#).

Solution:

Although there are several possible problems that can cause these two KrbException error messages to display, here are some actions you can take to solve the most likely problems:

- If you are using CentOS/Red Hat Enterprise Linux 5.6 or higher, or Ubuntu, which use AES-256 encryption by default for tickets, you must install the [Java Cryptography Extension \(JCE\) Unlimited Strength Jurisdiction Policy File](#) on all cluster and Hadoop user machines. For information about how to verify the type of encryption used in your cluster, see [Step 3: If you are Using AES-256 Encryption, Install the JCE Policy File](#) on page 98. Alternatively, you can change your `kdc.conf` or `krb5.conf` to not use AES-256 by removing `aes256-cts:normal` from the `supported_etypes` field of the `kdc.conf` or `krb5.conf` file. Note that after changing the `kdc.conf` file, you'll need to restart both the KDC and the kadmin server for those changes to take affect. You may also need to recreate or change the password of the relevant principals, including potentially the TGT principal (`krbtgt/REALM@REALM`).
- In the `[realms]` section of your `kdc.conf` file, in the realm corresponding to `HADOOP.LOCALDOMAIN`, add (or replace if it's already present) the following variable:

```
supported_etypes = des3-hmac-sha1:normal arcfour-hmac:normal des-hmac-sha1:normal  
des-cbc-md5:normal des-cbc-crc:normal des-cbc-crc:v4 des-cbc-crc:afs3
```

- Recreate the `hdfs` keytab file and `mapred` keytab file using the `-norandkey` option in the `xst` command (for details, see [Step 4: Create and Deploy the Kerberos Principals and Keytab Files](#) on page 98).

```
kadmin.local: xst -norandkey -k hdfs.keytab hdfs/fully.qualified.domain.name  
HTTP/fully.qualified.domain.name  
kadmin.local: xst -norandkey -k mapred.keytab mapred/fully.qualified.domain.name  
HTTP/fully.qualified.domain.name
```

The NameNode starts but clients cannot connect to it and error message contains enctype code 18.

Description:

The NameNode keytab file does not have an AES256 entry, but client tickets do contain an AES256 entry. The NameNode starts but clients cannot connect to it. The error message doesn't refer to "AES256", but does contain an enctype code "18".

Solution:

Make sure the "Java Cryptography Extension (JCE) Unlimited Strength Jurisdiction Policy File" is installed or remove `aes256-cts:normal` from the `supported_enctypes` field of the `kdc.conf` or `krb5.conf` file. For more information, see the [first suggested solution above for Problem 5](#).

For more information about the Kerberos encryption types, see <http://www.iana.org/assignments/kerberos-parameters/kerberos-parameters.xml>.

(MRv1 Only) Jobs won't run and TaskTracker is unable to create a local mapred directory.

Description:

The TaskTracker log contains the following error message:

```
11/08/17 14:44:06 INFO mapred.TaskController: main : user is atm
11/08/17 14:44:06 INFO mapred.TaskController: Failed to create directory
/var/log/hadoop/cache/mapred/mapred/local1/taskTracker/atm - No such file or directory
11/08/17 14:44:06 WARN mapred.TaskTracker: Exception while localization
java.io.IOException: Job initialization failed (20)
    at
    org.apache.hadoop.mapred.LinuxTaskController.initializeJob(LinuxTaskController.java:191)
        at org.apache.hadoop.mapred.TaskTracker$4.run(TaskTracker.java:1199)
        at java.security.AccessController.doPrivileged(Native Method)
        at javax.security.auth.Subject.doAs(Subject.java:396)
        at
    org.apache.hadoop.security.UserGroupInformation.doAs(UserGroupInformation.java:1127)
    at org.apache.hadoop.mapred.TaskTracker.initializeJob(TaskTracker.java:1174)
    at org.apache.hadoop.mapred.TaskTracker.localizeJob(TaskTracker.java:1089)
    at org.apache.hadoop.mapred.TaskTracker.startNewTask(TaskTracker.java:2257)
    at org.apache.hadoop.mapred.TaskTracker$TaskLauncher.run(TaskTracker.java:2221)
Caused by: org.apache.hadoop.util.Shell$ExitCodeException:
    at org.apache.hadoop.util.Shell.runCommand(Shell.java:255)
    at org.apache.hadoop.util.Shell.run(Shell.java:182)
    at org.apache.hadoop.util.Shell$ShellCommandExecutor.execute(Shell.java:375)
    at
    org.apache.hadoop.mapred.LinuxTaskController.initializeJob(LinuxTaskController.java:184)
    ... 8 more
```

Solution:

Make sure the value specified for `mapred.local.dir` is identical in `mapred-site.xml` and `taskcontroller.cfg`. If the values are different, the error message above is returned.

(MRv1 Only) Jobs will not run and TaskTracker is unable to create a Hadoop logs directory.

Description:

The TaskTracker log contains an error message similar to the following :

```
11/08/17 14:48:23 INFO mapred.TaskController: Failed to create directory
/home/atm/src/cloudera/hadoop/build/hadoop-0.23.2-cdh3u1-SNAPSHOT/logs1/userlogs/job_201108171441_0004
- No such file or directory
11/08/17 14:48:23 WARN mapred.TaskTracker: Exception while localization
java.io.IOException: Job initialization failed (255)
    at
    org.apache.hadoop.mapred.LinuxTaskController.initializeJob(LinuxTaskController.java:191)
        at org.apache.hadoop.mapred.TaskTracker$4.run(TaskTracker.java:1199)
        at java.security.AccessController.doPrivileged(Native Method)
        at javax.security.auth.Subject.doAs(Subject.java:396)
        at
    org.apache.hadoop.security.UserGroupInformation.doAs(UserGroupInformation.java:1127)
    at org.apache.hadoop.mapred.TaskTracker.initializeJob(TaskTracker.java:1174)
    at org.apache.hadoop.mapred.TaskTracker.localizeJob(TaskTracker.java:1089)
```

Configuring Authentication

```
    at org.apache.hadoop.mapred.TaskTracker.startNewTask(TaskTracker.java:2257)
    at org.apache.hadoop.mapred.TaskTracker$TaskLauncher.run(TaskTracker.java:2221)
Caused by: org.apache.hadoop.util.Shell$ExitCodeException:
    at org.apache.hadoop.util.Shell.runCommand(Shell.java:255)
    at org.apache.hadoop.util.Shell.run(Shell.java:182)
    at org.apache.hadoop.util.Shell$ShellCommandExecutor.execute(Shell.java:375)
    at
org.apache.hadoop.mapred.LinuxTaskController.initializeJob(LinuxTaskController.java:184)
    ... 8 more
```

Solution:

In MRv1, the default value specified for `hadoop.log.dir` in `mapred-site.xml` is `/var/log/hadoop-0.20-mapreduce`. The path must be owned and be writable by the `mapred` user. If you change the default value specified for `hadoop.log.dir`, make sure the value is identical in `mapred-site.xml` and `taskcontroller.cfg`. If the values are different, the error message above is returned.

After you enable cross-realm trust, you can run Hadoop commands in the local realm but not in the remote realm.

Description:

After you enable cross-realm trust, authenticating as a principal in the local realm will allow you to successfully run Hadoop commands, but authenticating as a principal in the remote realm will not allow you to run Hadoop commands. The most common cause of this problem is that the principals in the two realms either do not have the same encryption type, or the cross-realm principals in the two realms do not have the same password. This issue manifests itself because you are able to get Ticket Granting Tickets (TGTs) from both the local and remote realms, but you are unable to get a service ticket to allow the principals in the local and remote realms to communicate with each other.

Solution:

On the local MIT KDC server host, type the following command in the `kadmin.local` or `kadmin` shell to add the cross-realm `krbtgt` principal:

```
kadmin: addprinc -e "<enc_type_list>"
krbtgt/YOUR-LOCAL-REALM.COMPANY.COM@AD-REALM.COMPANY.COM
```

where the `<enc_type_list>` parameter specifies the types of encryption this cross-realm `krbtgt` principal will support: AES, DES, or RC4 encryption. You can specify multiple encryption types using the parameter in the command above, what's important is that at least one of the encryption types parameters corresponds to the encryption type found in the tickets granted by the KDC in the remote realm. For example:

```
kadmin: addprinc -e "aes256-cts:normal rc4-hmac:normal des3-hmac-sha1:normal"
krbtgt/YOUR-LOCAL-REALM.COMPANY.COM@AD-REALM.COMPANY.COM
```

(MRv1 Only) Jobs won't run and cannot access files in `mapred.local.dir`

Description:

The TaskTracker log contains the following error message:

```
WARN org.apache.hadoop.mapred.TaskTracker: Exception while localization
java.io.IOException: Job initialization failed (1)
```

Solution:

1. Add the `mapred` user to the `mapred` and `hadoop` groups on all hosts.
2. Restart all TaskTrackers.

Users are unable to obtain credentials when running Hadoop jobs or commands.

Description:

This error occurs because the ticket message is too large for the default UDP protocol. An error message similar to the following may be displayed:

```
13/01/15 17:44:48 DEBUG ipc.Client: Exception encountered while connecting to the server
: javax.security.sasl.SaslException:
GSS initiate failed [Caused by GSSException: No valid credentials provided (Mechanism
level: Fail to create credential.
(63) - No service creds)]
```

Solution:

Force Kerberos to use TCP instead of UDP by adding the following parameter to `libdefaults` in the `krb5.conf` file on the client(s) where the problem is occurring.

```
[libdefaults]
udp_preference_limit = 1
```

If you choose to manage `krb5.conf` through Cloudera Manager, this will automatically get added to `krb5.conf`.



Note:

When sending a message to the KDC, the library will try using TCP before UDP if the size of the ticket message is larger than the setting specified for the `udp_preference_limit` property. If the ticket message is smaller than `udp_preference_limit` setting, then UDP will be tried before TCP. Regardless of the size, both protocols will be tried if the first attempt fails.

Request is a replay [exceptions in the logs](#).

Description:

The root cause of this exception is that Kerberos uses a second-resolution timestamp to protect against replay attacks (where an attacker can record network traffic, and play back recorded requests later to gain elevated privileges). That is, incoming requests are cached by Kerberos for a little while, and if there are similar requests within a few seconds, Kerberos will be able to detect them as replay attack attempts. However, if there are multiple valid Kerberos requests coming in at the same time, these may also be misjudged as attacks for the following reasons:

- **Multiple services in the cluster are using the same Kerberos principal.** All secure clients that run on multiple machines should use unique Kerberos principals for each machine. For example, rather than connecting as a service principal `myservice@EXAMPLE.COM`, services should have per-host principals such as `myservice/host123.example.com@EXAMPLE.COM`.
- **Clocks not in sync:** All hosts should run NTP so that clocks are kept in sync between clients and servers.

While having different principals for each service, and clocks in sync helps mitigate the issue, there are, however, cases where even if all of the above are implemented, the problem still persists. In such a case, disabling the cache (*and the replay protection as a consequence*), will allow parallel requests to succeed. This compromise between usability and security can be applied by setting the `KRB5RCACHETYPE` environment variable to `none`.

Note that the `KRB5RCACHETYPE` is not automatically detected by Java applications. For Java-based components:

- Ensure that the cluster runs on JDK 8.
- To disable the replay cache, add `-Dsun.security.krb5.rcache=none` to the Java Opts/Arguments of the targeted JVM. For example, HiveServer2 or the Sentry service.

For more information, refer the [MIT KDC documentation](#).

Symptom: The following exception shows up in the logs for one or more of the Hadoop daemons:

```
2013-02-28 22:49:03,152 INFO ipc.Server (Server.java:doRead(571)) - IPC Server listener
on 8020: readAndProcess threw exception javax.security.sasl.SaslException: GSS initiate
failed [Caused by GSSException: Failure unspecified at GSS-API level (Mechanism 1
javax.security.sasl.SaslException: GSS initiate failed [Caused by GSSException: Failure
unspecified at GSS-API level (Mechanism level: Request is a replay (34))]
    at
com.sun.security.sasl.gsskerb.GssKrb5Server.evaluateResponse(GssKrb5Server.java:159)
    at org.apache.hadoop.ipc.Server$Connection.saslReadAndProcess(Server.java:1040)

    at org.apache.hadoop.ipc.Server$Connection.readAndProcess(Server.java:1213)
    at org.apache.hadoop.ipc.Server$Listener.doRead(Server.java:566)
    at org.apache.hadoop.ipc.Server$Listener$Reader.run(Server.java:363)
Caused by: GSSException: Failure unspecified at GSS-API level (Mechanism level: Request
is a replay (34))
    at sun.security.jgss.krb5.Krb5Context.acceptSecContext(Krb5Context.java:741)
    at sun.security.jgss.GSSContextImpl.acceptSecContext(GSSContextImpl.java:323)
    at sun.security.jgss.GSSContextImpl.acceptSecContext(GSSContextImpl.java:267)
    at
com.sun.security.sasl.gsskerb.GssKrb5Server.evaluateResponse(GssKrb5Server.java:137)
    ... 4 more
Caused by: KrbException: Request is a replay (34)
    at sun.security.krb5.KrbApReq.authenticate(KrbApReq.java:300)
    at sun.security.krb5.KrbApReq.<init>(KrbApReq.java:134)
    at sun.security.jgss.krb5.InitSecContextToken.<init>(InitSecContextToken.java:79)

    at sun.security.jgss.krb5.Krb5Context.acceptSecContext(Krb5Context.java:724)
    ... 7 more
```

In addition, this problem can manifest itself as performance issues for all clients in the cluster, including dropped connections, timeouts attempting to make RPC calls, and so on.

CDH services fail to start

Possible Causes: Check that the encryption types are matched between your KDC and `krb5.conf` on all hosts.

Solution: If you are using AES-256, follow the instructions at [Step 2: If You are Using AES-256 Encryption, Install the JCE Policy File](#) on page 62 to deploy the JCE policy file on all hosts.

Configuring Encryption

The goal of encryption is to ensure that only authorized users can view, use, or contribute to a data set. These security controls add another layer of protection against potential threats by end-users, administrators, and other malicious actors on the network. Data protection can be applied at a number of levels within Hadoop:

- **OS Filesystem-level** - Encryption can be applied at the Linux operating system filesystem level to cover all files in a volume. An example of this approach is [Cloudera Navigator Encrypt](#) on page 316 (formerly Gazzang zNcrypt) which is available for Cloudera customers licensed for Cloudera Navigator. Navigator Encrypt operates at the Linux volume level, so it can encrypt cluster data inside and outside HDFS, such as temp/spill files, configuration files and metadata databases (to be used only for data related to a CDH cluster). Navigator Encrypt must be used with [Cloudera Navigator Key Trustee Server](#) on page 297 (formerly Gazzang zTrustee).

CDH components, such as Impala, MapReduce, YARN, or HBase, also have the ability to encrypt data that lives temporarily on the local filesystem outside HDFS. To enable this feature, see [Configuring Encryption for Data Spills](#) on page 338.

- **Network-level** - Encryption can be applied to encrypt data just before it gets sent across a network and to decrypt it just after receipt. In Hadoop, this means coverage for data sent from client user interfaces as well as service-to-service communication like remote procedure calls (RPCs). This protection uses industry-standard protocols such as TLS/SSL.



Note: Cloudera Manager and CDH components support either TLS 1.0, TLS 1.1, or TLS 1.2, but not SSL 3.0. References to SSL only continue because of its widespread use in technical jargon.

- **HDFS-level** - Encryption applied by the HDFS client software. [HDFS Transparent Encryption](#) on page 261 operates at the HDFS folder level, allowing you to encrypt some folders and leave others unencrypted. HDFS transparent encryption cannot encrypt any data outside HDFS. To ensure reliable key storage (so that data is not lost), use Cloudera Navigator Key Trustee Server; the default Java keystore can be used for test purposes. For more information, see [Enabling HDFS Encryption Using Cloudera Navigator Key Trustee Server](#) on page 267.

Unlike OS and network-level encryption, HDFS transparent encryption is end-to-end. That is, it protects data at rest and in transit, which makes it more efficient than implementing a combination of OS-level and network-level encryption.

TLS/SSL Certificates Overview

This topic will guide you through the different certificate strategies that you can employ on your cluster to allow TLS/SSL clients to securely connect to servers using trusted certificates or certificates issued by trusted authorities. The set of certificates required depends upon the certificate provisioning strategy you implement. The following strategies are possible:

- **Public CA-signed certificates:** Using certificates signed by a trusted public certificate authority (CA) simplifies the integration procedure, because the default Java client already trusts the CA.
- **Internal CA-signed certificates:** Using certificates signed by an internal CA can also simplify integration if your infrastructure already trusts the internal CA. Otherwise, you must configure all hosts to trust your internal CA.
- **Self-signed certificates:** Using self-signed certificates complicates the deployment process because you must configure all clients of a particular service to trust the specific certificate used by that service.



Note: Wildcard domain certificates and certificates using the SubjectAlternativeName extension are not supported at this time.

When choosing an approach to certificate provisioning, keep in mind that TLS/SSL must be enabled for all core Hadoop services (HDFS, MapReduce, and YARN) as a group. For example, if you are running HDFS and YARN on your cluster, you cannot choose to enable TLS/SSL for HDFS, but not for YARN. You must enable it for both services, which implies that you must make certificates available to all daemon roles of both services. You will need to obtain a certificate for each host on which an HDFS or YARN daemon role is running.



Note: Cloudera Manager and CDH components support either TLS 1.0, TLS 1.1, or TLS 1.2, but not SSL 3.0. References to SSL only continue because of its widespread use in technical jargon.

Creating Certificates

The following sections will walk you through obtaining certificates from commercial Certificate Authorities and creating self-signed test certificates.

Using Keytool

`keytool` is a utility for creating and managing certificates and cryptographic keys, and is part of the standard JDK distribution. The `keytool` executable usually resides in `$JAVA_HOME/bin`.

`keytool` stores certificates and keys in a file known as a [keystore](#). While several different keystore types are supported, by default `keytool` uses the [Java KeyStore \(JKS\)](#) format.

Java-based services such as HDFS, MapReduce, and YARN use the JKS format by default. For this reason it is particularly convenient to use `keytool` for managing keys and certificates for these services. In the following topics, we assume you are using `keytool`.

For additional information on `keytool`, refer the [keytool documentation](#).

Using OpenSSL

Python-based services such as Hue expect certificates and keys to be stored in `PEM` format. When managing certificates and keys for such services, you may find it convenient to use the `openssl` tool.

Refer the [openssl](#) documentation for more information.

Obtaining a Production Certificate from a Commercial CA

Once you have decided on a certificate-provisioning strategy, and have determined which hosts require certificates, you will typically purchase the necessary certificates from a commercial Certificate Authority (CA). The procedure for applying for a certificate varies from one CA to another, but typically involves providing some form of proof that you are the legitimate owner of the domain name for which you are requesting a certificate, generating a key pair, and submitting a Certificate Signing Request (CSR) to the CA.

As noted above, you may find it convenient to use the Java `keytool` utility to generate your key pair and CSR, and to manage your certificates. The CA you choose will provide instructions for obtaining and installing a certificate; typically, there will be separate sets of instructions for different web and application servers. The instructions for Java-based servers (Tomcat, for example), will usually describe the following process comprising three `keytool` commands to obtain a certificate:

1. `keytool -genkeypair` to generate a public/private key pair and create the keystore.
2. `keytool -certreq` to create the CSR.
3. `keytool -importcert` to import the signed certificate into the keystore.

For example, to generate a public/private key pair for the domain name `node1.example.com`, you would use a command similar to the one shown below:

```
$ keytool -genkeypair -keystore node1.keystore -alias node1 \  
-dname "CN=node1.example.com,O=Hadoop" -keyalg RSA \  
-keysize 2048 -storepass changeme -keypass changeme
```

This command generates a pair of 2048-bit keys using the RSA key algorithm, one of several available. The keys are stored in a keystore file called `node1.keystore`, in a keystore entry identified by the alias `node1`. The keystore password (which protects the keystore as a whole) and the key password (which protects the private key stored in the `node1` entry) are set using the `-storepass` and `-keypass` options (respectively). `-keypass` must be set to the same password value as `-storepass` for Cloudera Manager to access the keystore.

To create a CSR, you would use a command similar to the following:

```
$ keytool -certreq -keystore node1.keystore -alias node1 \
-storepass changeme -keypass changeme -file node1.csr
```

This command generates the CSR, and stores it in a file called `node1.csr`. Once you've submitted your CSR to the CA, and received the CA's reply (containing the signed certificate), you will use the following `keytool -importcert` command to import the reply into your keystore:

```
$ keytool -importcert -keystore node1.keystore -alias node1 \
-storepass changeme -keypass changeme -trustcacerts -file node1.crt
```

Here we assume that the CA's reply is stored in the file `node1.crt`.



Important: This section describes a generic procedure using `keytool` to obtain a certificate from a commercial Certificate Authority. This procedure will differ from one CA to another and Cloudera recommends you consult your CA's documentation for more specifics.

Creating Self-Signed Test Certificates



Important: Cloudera strongly recommends against the use of self-signed certificates in production clusters.

It is also possible to create your own test certificates. These certificates are typically self-signed; that is, they are signed by your own private key, rather than that of an external CA. Such test certificates are useful during testing and bringup of a cluster.

To generate a self-signed certificate, use `keytool -genkeypair`. (In addition to creating a public/private key pair, this command wraps the public key into a self-signed certificate.) For example, the following command creates a self-signed test certificate for the host `node1.example.com`, and stores it in a keystore named `node1.keystore`:

```
$ keytool -genkeypair -keystore node1.keystore -keyalg RSA \
-alias node1 -dname "CN=node1.example.com,O=Hadoop" \
-storepass changeme -keypass changeme -validity <val_days>
```

By default, self-signed certificates are only valid for 90 days. To increase this period, replace `<val_days>` in the previous command's `-validity <val_days>` parameter to specify the number of days for which the certificate should be considered valid.

Creating Java Keystores and Truststores

Typically, a keystore is used in one of two distinct ways:

- The keystore contains private keys and certificates used by TLS/SSL servers to authenticate themselves to TLS/SSL clients. By convention, such files are referred to as keystores.
- When used as a *truststore*, the file contains certificates of trusted TLS/SSL servers, or of Certificate Authorities trusted to identify servers. There are no private keys in the truststore.



Note: The foregoing assumes that certificate-based authentication is being used in one direction only—that is, TLS/SSL servers are using certificates to authenticate themselves to clients. It is also possible for clients to authenticate themselves to servers using certificates. (This is known as mutual authentication.) Throughout this document, we assume that client certificates are not in use.

While all TLS/SSL clients must have access to a truststore, it is not always necessary to create and deploy truststores across a cluster. The standard JDK distribution includes a default truststore which is pre-provisioned with the root certificates of a number of well-known Certificate Authorities. If you do not provide a custom truststore, the Hadoop daemons load this default truststore. Therefore, if you are using certificates issued by a CA in the default truststore, you do not need to provide custom truststores. However, you must consider the following before you decide to use the default truststore:

- If you choose to use the default truststore, it is your responsibility to maintain it. You may need to remove the certificates of CAs you do not deem trustworthy, or add or update the certificates of CAs you trust. Use the `keytool` utility to perform these actions.

Security Considerations for Keystores and Truststores



Note: While the strategy for certificate deployment you select will ultimately depend upon the security policies you wish to implement, the following guidelines may prove useful.

Because keystores contain private keys, while truststores do not, the security requirements for keystores are more stringent. In particular:

- Hadoop TLS/SSL requires that truststores and the truststore password be stored, in plaintext, in a configuration file that is readable by all.
- Keystore and key passwords are stored, in plaintext, in a file that is readable only by members of the appropriate group.

These considerations should inform your choice of which keys and certificates to store in the keystores and truststores you will deploy across your cluster.

- Keystores should contain a minimal set of keys and certificates. A reasonable strategy would be to create a unique keystore for each host, which would contain only the keys and certificates needed by the Hadoop TLS/SSL services running on the host. In most cases, the keystore would contain a single key/certificate entry.

Modifying Keystores: CDH services and processes must be restarted in case changes are made to a keystore. However, this is relatively rare since keystores do not need to be updated when hosts are added or deleted from a cluster.

- Because truststores do not contain sensitive information, it is reasonable to create a single truststore for an entire cluster. On a production cluster, such a truststore would often contain a single CA certificate (or certificate chain), since you would typically choose to have all certificates issued by a single CA.



Important: *Do not use the same password for truststores and keystores/keys.*

Since truststore passwords are stored in the clear in files readable by all, doing so would compromise the security of the private keys in the keystore.

Creating Keystores

Once you have settled on a storage plan for your keys and certificates, you can use `keytool` to create or update the necessary keystores and truststores. To create a new keystore with a certificate see [Creating Certificates](#) on page 206.

In many cases, you will already have created the set of keystores that you need. If you have followed the approach of creating a separate keystore for each private key and certificate, and wish to maintain this arrangement when deploying the keystores, no additional steps are required to prepare the keystores for deployment. If you wish to reorganize your

keys and certificates into a different set of keystores, you can use `keytool -importkeystore` to transfer entries from one keystore to another.

Creating Truststores

The steps involved in preparing the truststores to be used in your deployment depend on whether you have decided to use the default Java truststore, or to create custom truststores:

- If you are using the default truststore, you may need to add CA certificates (or certificate chains) to the truststore, or delete them from the truststore.
- If you are creating custom truststores, you will need to build the truststores by importing trusted certificates into new truststores. The trusted certificates can be CA certificates (typically downloaded from the CA's website), or self-signed certificates that you have created. Note that Cloudera strongly recommends against using self-signed certificates in production.

As shown in the examples below, when creating a truststore you must select a password. All truststore passwords for a given service must be the same. In practice, this restriction rarely comes into play, since it is only relevant when you wish to create distinct custom truststores for each host.

The following sections provide examples of the steps required for several common scenarios:

Example 1: Adding a CA Certificate to the alternative Default Truststore

In this example, we assume that you have chosen to use the default Java truststore, but have obtained a certificate from a CA not included in the truststore. (This situation can also arise if the CA that issued your certificate has an entry in the default truststore, but the particular certificate product you purchased requires an alternate CA certificate chain.)

1. Locate the default truststore on your system. The default truststore is located in the `$JAVA_HOME/jre/lib/security/cacerts` file. This contains the default CA information shipped with the JDK. Create an alternate default file called `jssecacerts` in the same location as the `cacerts` file. You can now safely append CA certificates for any private or public CAs not present in the default `cacerts` file, while keeping the original file intact.

The alternate file will always be read unless the `javax.net.ssl.trustStore` flag is set in the arguments for the startup of the `java` process.

For our example, we will be following this recommendation by copying the default `cacerts` file into the new `jssecacerts` file.

```
$ cp $JAVA_HOME/jre/lib/security/cacerts \
  $JAVA_HOME/jre/lib/security/jssecacerts
```

If you use a copy of the `cacerts` file, remember the default keystore password is `changeit`.

2. Import the CA certificate into the default truststore. Assuming that the file `CA-root.cer` contains the CA's certificate, which you have previously downloaded from the CA's web site, the following command imports this certificate into the alternative default truststore.

```
$ keytool -importcert -file CA-root.cer -alias CAcert \
  -keystore /usr/java/default/jre/lib/security/jssecacerts \
  -storepass changeit
```

When you give this command, you will be prompted to confirm that you trust the certificate. Be sure to verify that the certificate is genuine before importing it.



Important: Test the trust relationship before you import any intermediary CA certificates. Trust should be derived from the root CA only. Import intermediary CA certificates only if necessary.

Note that any updates you make to the default truststore must be made on all hosts in the cluster.

Example 2: Creating a Custom Truststore Containing a Single CA Certificate Chain

In this example, we demonstrate how to use `keytool` to create a custom truststore. We assume all certificates were issued by a single CA, so a truststore containing the certificate chain for that CA will serve for all hosts in the cluster.

Our example certificate chain consists of a root certificate and a single intermediate certificate. We assume that you have downloaded these and saved them in the files `CA-root.cer` and `CA-intermediate.cer` (respectively). The steps below show the commands needed to build a custom truststore containing the root and intermediate certificates.

1. Import the root certificate and create the truststore:

```
$ keytool -importcert -keystore custom.truststore -alias CA-cert \  
-storepass trustchangeme -file CA-root.cer
```

You will be prompted to confirm that the root certificate is trustworthy. Be sure to verify that the certificate is genuine before you import it.

2. Import the intermediate certificate into the truststore created in Step 1:

```
$ keytool -importcert -keystore custom.truststore \  
-alias CA-intermediate -storepass trustchangeme \  
-file CA-intermediate.cer
```

Example 3: Creating a Custom Truststore Containing Self-Signed Test Certificates



Important: Cloudera strongly recommends against the use of self-signed certificates in production clusters.

This example is particularly relevant when setting up a test cluster. We assume that you have generated a set of self-signed test certificates for the hosts in the cluster, and wish to create a single truststore that can be deployed on all hosts. Because the certificates are self-signed, we cannot simply construct a truststore containing a single certificate chain, as in the previous example. When a client receives a self-signed certificate from a server during the TLS/SSL handshake, it must be able to find the server's certificate in the truststore, since no other signing certificate exists to establish trust. Therefore, the truststore must contain all the test certificates.

We assume that the test certificates reside in keystores named `node1.keystore` ... `node100.keystore`, which were created following the steps described in [Creating Self-Signed Test Certificates](#).

1. Export the test certificate for `node1.example.com`:

```
$ keytool -exportcert -keystore node1.keystore -alias node1 \  
-storepass changeme -file node1.cer
```

2. Import the test certificate into the custom truststore:

```
keytool -importcert -keystore custom.truststore -alias node1 \  
-storepass trustchangeme -file node1.cer -noprompt
```

Here we specify the `-noprompt` option to suppress the prompt asking you to confirm that the certificate is trustworthy. Since you created the certificate yourself, this confirmation is unnecessary.

3. Repeat Steps 1 and 2 for `node2.keystore` ... `node100.keystore`.

Private Key and Certificate Reuse Across Java Keystores and OpenSSL

This topic provides a quick tutorial on exporting/importing private keys for reuse from a Java keystore to OpenSSL and vice versa. Regardless of the procedure followed to create host private keys and certificates, sometimes it becomes necessary to reuse those private keys and certificates by other services on the same host. For example, if you used OpenSSL to create private keys and certificates for a service, you can reuse those keys for a Java-based service on the same host by converting them to the Java keystore format.

The documentation for [Configuring TLS Security for Cloudera Manager](#) describes both approaches to creating private keys, using Java keystore, and OpenSSL.

Why Reuse a Private Key?

Certificate authorities generally revoke previous generations of certificates issued to a host. Hence, a host cannot have 2 sets of CA-issued certificates and have both be valid. Once a certificate is issued to a host, it then becomes necessary to reuse the private key that requested the certificate, and the CA-issued certificate across different services, Java-based and otherwise.



Note: This following sections assume the default paths set up in [Configuring TLS Encryption Only for Cloudera Manager](#).

Conversion from Java Keystore to OpenSSL

First, use `keytool` to export the private key and certificate to a `PKCS12` file as a transitional file format that can then be split up into individual key and certificate files by the `openssl` command line. Replace `cmhost` and `hostname` in the commands below with the actual hostname of the server that is managing the certificate and keys.

```
$ keytool -importkeystore -srckeystore /opt/cloudera/security/jks/hostname-keystore.jks \
-srcstorepass password -srckeypass password -destkeystore /tmp/hostname-keystore.p12 \
-deststoretype PKCS12 -srcalias hostname -deststorepass password -destkeypass password
```

Now use `openssl` to split the `PKCS12` file created above into first, the certificate file, and then the private key file. While the CA-issued certificate can be used as is, the command has been provided here for completeness.

```
$ openssl pkcs12 -in /tmp/hostname-keystore.p12 -passin pass:password -nokeys \
-out /opt/cloudera/security/x509/hostname.pem

$ openssl pkcs12 -in /tmp/hostname-keystore.p12 -passin pass:password -nocerts \
-out /opt/cloudera/security/x509/hostname.key -passout pass:password
```

Note that the method above generates a key with a password. For services such as Impala and Hue that accept keys without passwords, you can use the following command:

```
$ openssl rsa -in /opt/cloudera/security/x509/hostname.key \
-passin pass:password -out /opt/cloudera/security/x509/hostname.pem
```

Conversion from OpenSSL to Java Keystore

First, convert the `openssl` private key and certificate files into a `PKCS12` file. The `PKCS12` file can then be imported into a Java keystore file. Replace `hostname` in the commands below with the FQDN for the host whose certificate is being imported.

```
$ openssl pkcs12 -export -in /opt/cloudera/security/x509/hostname.pem \
-inkey /opt/cloudera/security/x509/hostname.key -out /tmp/hostname.p12 \
-name hostname -passin pass:password -passout pass:password

$ keytool -importkeystore -srckeystore /tmp/hostname.p12 -srcstoretype PKCS12 \
-srcstorepass password -alias hostname -deststorepass password \
-destkeypass password -destkeystore /opt/cloudera/security/jks/hostname-keystore.jks
```

Configuring TLS Security for Cloudera Manager

Transport Layer Security (TLS) provides encryption and authentication in communication between the Cloudera Manager Server and Agents. Encryption prevents snooping, and authentication helps prevent problems caused by malicious

Configuring Encryption

servers or agents. If you are familiar with TLS encryption, and want to get started enabling it for Cloudera Manager, see [How to Configure TLS Encryption for Cloudera Manager](#) on page 40.

Cloudera Manager supports three levels of TLS security.

- Level 1 (Good) - This level encrypts communication between the browser and Cloudera Manager, and between Agents and the Cloudera Manager Server. See [Configuring TLS Encryption Only for Cloudera Manager](#) on page 212 followed by [Level 1: Configuring TLS Encryption for Cloudera Manager Agents](#) on page 216 for instructions. Level 1 encryption prevents snooping of commands and controls ongoing communication between Agents and Cloudera Manager.
- Level 2 (Better) - This level encrypts communication between the Agents and the Server, and provides strong verification of the Cloudera Manager Server certificate by Agents. See [Level 2: Configuring TLS Verification of Cloudera Manager Server by the Agents](#) on page 217. Level 2 provides Agents with additional security by verifying trust for the certificate presented by the Cloudera Manager Server.
- Level 3 (Best) - This includes encrypted communication between the Agents and the Server, strong verification of the Cloudera Manager Server certificate by the Agents, *and* authentication of Agents to the Cloudera Manager Server using self-signed or CA-signed certs. See [Level 3: Configuring TLS Authentication of Agents to the Cloudera Manager Server](#) on page 219. Level 3 TLS prevents cluster Servers from being spoofed by untrusted Agents running on a host. Cloudera recommends that you configure Level 3 TLS encryption for untrusted network environments before enabling Kerberos authentication. This provides secure communication of keytabs between the Cloudera Manager Server and verified Agents across the cluster.



Important:

- Cloudera strongly recommends that you set up a fully functional CDH cluster and Cloudera Manager before you configure the Cloudera Manager Server and Agents to use TLS.
- When TLS is enabled, Cloudera Manager continues to accept HTTP requests on port 7180 (default) but immediately redirects clients to port 7183 for HTTPS connectivity.
- You must finish configuring Level 1 and Level 2 TLS to configure Level 3 encryption. To enable TLS encryption for all connections between your Web browser running the Cloudera Manager Admin Console and the Cloudera Manager Server, see the first 2 steps of [Level 1: Configuring TLS Encryption for Cloudera Manager Agents](#) on page 216.
- When Level 3 TLS is configured, to add new hosts running Agents, you must manually deploy the Cloudera Manager agent and daemon packages for your platform, issue a new certificate for the host, configure `/etc/cloudera-scm-agent/config.ini` to use TLS/SSL, and then bring the host online.

Alternatively, you can disable TLS to add the host, configure the new host for TLS, and then re-enable with the proper configuration in place. Either approach is valid, based on your needs.

- For all hosts running Agents, Cloudera recommends that you first create the keystore in Java, and then export the key and certificate using `openssl` for the Agent or Hue.

For details on how HTTPS communication is handled Cloudera Manager Agents and Cloudera Management Services daemons, see [HTTPS Communication in Cloudera Manager](#) on page 224.

Configuring TLS Encryption Only for Cloudera Manager

Minimum Required Role: [Cluster Administrator](#) (also provided by [Full Administrator](#))



Important: The sequence of steps described in the following topics to configure Level 1 through 3 TLS will each build upon the steps of the previous level. The procedure and examples provided in these topics, are based on this concept.

Before enabling TLS security for Cloudera Manager, you must create a keystore, submit a certificate-signing request, and install the issued certificate for the Server. You do this using the Oracle JDK `keytool` command-line tool. If you are using a Private CA, append its certificate (and any required intermediary certificates) to the alternate default

truststore provided with the JDK for inherent trust. This process is described in detail in [Creating Truststores](#) on page 209.

The table below shows the paths for managing certificates in the following examples. These paths persist during any upgrades and should be removed manually if the host is removed from a CDH cluster. Note that the folders and filepaths listed here can reside anywhere on the system and must be created on every host, especially as later sections move on to creating certificates for each host.



Note: Set permissions on the paths such that `scm-user`, `hue`, Hadoop service users (or groups), and `root` users can read the private key, certificate, and keystore and truststore files.

Example Property Values	Description
<code>cmhost.sec.cloudera.com</code>	FQDN for Cloudera Manager Server host.
<code>/opt/cloudera/security</code>	Base location for security-related files.
<code>/opt/cloudera/security/x509</code>	Location for <code>openssl key/, cert/ and cacerts/</code> files to be used by the Cloudera Manager Agent and Hue.
<code>/opt/cloudera/security/jks</code>	Location for the Java-based <code>keystore/</code> and <code>truststore/</code> files for use by Cloudera Manager and Java-based cluster services.
<code>/opt/cloudera/security/CAcerts</code>	Location for CA certificates (root and intermediary/subordinate CAs). One PEM file per CA in the chain is required.



Important:

- You must use the Oracle JDK `keytool`. The following procedure requires use of the Cloudera-installed Oracle JDK (or JDK downloaded from Oracle). Do not use both `keytool` and OpenJDK, or varying versions of JDK command line tools like `keytool`. If necessary, set your `PATH` so that the Oracle JDK is first. For example:

```
$ export JAVA_HOME=/usr/java/jdk1.7.0_67-cloudera
$ export PATH=$JAVA_HOME/bin:$PATH
```

- Set `-keypass` to the same value as `-storepass`. Cloudera Manager assumes that the same password is used to access both the key and the keystore, and therefore, does not support separate values for `-keypass` and `-storepass`.

Step 1: Create the Cloudera Manager Server Keystore, Generate a Certificate Request, and Install the Certificate



Important: The following instructions use private CA-signed certificates. If you're testing using a self-signed certificate, see [Using Self-Signed Certificates for TLS](#) on page 226.

The following procedure assumes that a private Certificate Authority is used, and therefore trust must be established for that private CA. If a known public CA such as Verisign or GeoTrust is used, you may not need to explicitly establish trust for the issued certificates. Newer public CAs might not be present yet in the JDK default `cacerts` file. If you have problems with the import process (such as `keytool error: java.lang.Exception: Failed to establish chain from reply`), follow the steps for trusting private CAs below.

Configuring Encryption

1. Assuming the paths documented in the table above have been created, use `keytool` to generate a Java keystore and Certificate Signing Request (CSR) for the Cloudera Manager Server. Replace `cmhost` and `cmhost.sec.cloudera.com` in the commands below with your hostname and FQDN. For example:

```
$ keytool -genkeypair -alias cmhost -keyalg RSA -keystore \  
/opt/cloudera/security/jks/cmhost-keystore.jks -keysize 2048 -dname \  
"CN=cmhost.sec.cloudera.com,OU=Support,O=Cloudera,L=Denver,ST=Colorado,C=US" \  
-storepass password -keypass password
```

- `-alias` is a label used only in the keystore. In this example, the hostname is used for easy tracking and management of the key and certificate. Ensure that `-alias` is consistent across all your commands.
- `-keyalg` is the algorithm used to generate the key. Cloudera recommends you use RSA, which allows key lengths greater than 1024 bits for certificate requests. Other algorithms such as the DSA may not be supported by certain browsers, resulting in the `javax.net.ssl.SSLHandshakeException: no cipher suites in common` error.
- `-dname` allows you to provide the certificate subject as a single line. If not specified, you will be prompted for the values of the certificate subject information. In that case, use the host FQDN that agents and browsers will use to connect to in the subject **First and Last name (CN)** question prompt.



Note: The CN entry must match the hostname of the Cloudera Manager server, or you will get the `java.io.IOException: HTTPS hostname wrong` exception.

- `/opt/cloudera/security/jks/cmhost-keystore.jks` is an example path to the keystore where you store the keystore file and where the Cloudera Manager Server host can access it.
- `-keypass` must be set to the same password as `-storepass` to allow Cloudera Manager to access the keystore. This is because Cloudera Manager assumes that the same password is used to access the key and the keystore, and therefore, accepts only one password.

2. Generate a certificate signing request for the host (in this example, `cmhost`).

```
$ keytool -certreq -alias cmhost \  
-keystore /opt/cloudera/security/jks/cmhost-keystore.jks \  
-file /opt/cloudera/security/x509/cmhost.csr -storepass password \  
-keypass password
```

3. Submit the `.csr` file created by the `-certreq` command to your Certificate Authority to obtain a server certificate. When possible, work with certificates in the default Base64 (ASCII) format. You can easily modify Base64-encoded files from `.CER` or `.CRT` to `.PEM`. The file is in ASCII format if you see the opening and closing lines as follows:

```
-----BEGIN CERTIFICATE-----  
( the encoded certificate is represented by multiple lines of exactly 64 characters,  
except  
for the last line which can contain 64 characters or less.)  
-----END CERTIFICATE-----
```

If your issued certificate is in binary (DER) format, adjust the commands according to the `keytool` documentation.

4. Copy the root CA certificate and any intermediary or subordinate CA certificates to `/opt/cloudera/security/CAcerts/`.



Important: For a private CA, you must import the private CA and intermediary or subordinate CA certificates into an alternative default JDK truststore `jssecacerts`, *before* importing them to your Java keystore.

- a. Import the root CA certificate first, followed by any intermediary or subordinate CA certificates. Substitute `$JAVA_HOME` in the command below with the path for your Oracle JDK.

```
$ sudo cp $JAVA_HOME/jre/lib/security/cacerts $JAVA_HOME/jre/lib/security/jssecacerts
$ sudo keytool -importcert -alias RootCA -keystore $JAVA_HOME/jre/lib/security/jssecacerts \
-file /opt/cloudera/security/CAcerts/RootCA.cer -storepass changeit
$ sudo keytool -importcert -alias SubordinateCA -keystore \
$JAVA_HOME/jre/lib/security/jssecacerts \
-file /opt/cloudera/security/CAcerts/SubordinateCA.cer -storepass changeit
```

Repeat for as many subordinate or intermediary CA certificates as needed. The default `-storepass` for the `cacerts` file is `changeit`. After completing this step, copy the `jssecacerts` file created to the same path on all cluster hosts.

- b. Import the Private CA certificates into your Java keystore file. Import the root CA certificate first.

```
$ keytool -importcert -trustcacerts -alias RootCA -keystore \
/opt/cloudera/security/jks/<cmhost-keystore>.jks -file \
/opt/cloudera/security/CAcerts/RootCA.cer -storepass password
$ keytool -importcert -trustcacerts -alias SubordinateCA -keystore \
/opt/cloudera/security/jks/<cmhost-keystore>.jks -file \
/opt/cloudera/security/CAcerts/SubordinateCA.cer -storepass password
```

Repeat for as many subordinate/intermediary CA certificates as needed.

5. Copy the signed certificate file provided to a location where it can be used by the Cloudera Manager Agents (and Hue if necessary).

```
$ cp certificate-file.cer /opt/cloudera/security/x509/cmhost.pem
```

Install it with the following `keytool` command:

```
$ keytool -importcert -trustcacerts -alias cmhost \
-file /opt/cloudera/security/x509/cmhost.pem \
-keystore /opt/cloudera/security/jks/cmhost-keystore.jks -storepass password
```

You *must* see the following response verifying that the certificate has been properly imported against its private key.

```
Certificate reply was installed in keystore
```

Because the issued certificate has been imported by the Java keystore, the original certificate-signing request (.CSR) and certificate files are no longer needed by Java services on that host, and the certificate and private key are now accessed through the keystore.

However, you still must export the private key from the Java keystore to make the certificate usable by Hue and the Cloudera Manager Agent. For instructions on reusing certificates, see [Private Key and Certificate Reuse Across Java Keystores and OpenSSL](#) on page 210.

Step 2: Enable HTTPS for the Cloudera Manager Admin Console and Specify Server Keystore Properties

1. Log into the Cloudera Manager Admin Console.
2. Select **Administration > Settings**.
3. Click the **Security** category.
4. Configure the following TLS settings:

Property	Description
Path to TLS Keystore File	The complete path to the keystore file. In the example, this path would be: <code>/opt/cloudera/security/jks/cmhost-keystore.jks</code>
Keystore Password	The password for keystore: <code>password</code>
Use TLS Encryption for Admin Console	Check this box to enable TLS encryption for Cloudera Manager.

5. Click **Save Changes** to save the settings.

Step 3: Specify SSL Truststore Properties for Cloudera Management Services

When enabling TLS for the Cloudera Manager UI, you must set the Java truststore location and password in the Cloudera Management Services configuration. If this is not done, roles such as the Host Monitor and Service Monitor will be unable to connect to Cloudera Manager and will not start.

1. Open the Cloudera Manager Administration Console and go to the **Cloudera Management Service**.
2. Click the **Configuration** tab.
3. Select **Scope > Cloudera Management Service (Service-Wide)**.
4. Select **Category > Security**.
5. Edit the following TLS/SSL properties according to your cluster configuration.

Property	Description
TLS/SSL Client Truststore File Location	Path to the client truststore file used in HTTPS communication. The contents of this truststore can be modified without restarting the Cloudera Management Service roles. By default, changes to its contents are picked up within ten seconds.
TLS/SSL Client Truststore File Password	Password for the client truststore file.

6. Click **Save Changes** to commit the changes.
7. Restart the Cloudera Management Service. For more information, see [HTTPS Communication in Cloudera Manager](#) on page 224.

Step 4: Restart the Cloudera Manager Server

Restart the Cloudera Manager Server by running `service cloudera-scm-server restart` from the Cloudera Manager host command prompt.

You should now be able to connect to the Cloudera Manager Admin Console using an HTTPS browser connection. If a private CA certificate or self-signed certificate is used, you must establish trust in the browser for your certificate. This should be done for all browsers that will be used to access Cloudera Manager. By default, certificates issued by public commercial CAs should be trusted by the browsers accessing Cloudera Manager and other Java or OpenSSL-based services.

For more information on establishing trust for certificates, see [TLS/SSL Certificates Overview](#) on page 205 or the relevant [JDK documentation](#).

Level 1: Configuring TLS Encryption for Cloudera Manager Agents

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

Prerequisite:

You must have completed the steps described at [Configuring TLS Encryption Only for Cloudera Manager](#).

Step 1: Enable Agent Connections to Cloudera Manager to use TLS

In this step, you enable TLS properties for Cloudera Manager Agents and their connections to the Cloudera Manager Server. To configure agents to connect to Cloudera Manager over TLS, log into the Cloudera Manager Admin Console.



Note: If you are using a private certificate authority to sign certificate requests, see information on establishing trust for this CA in [Configuring TLS Encryption Only for Cloudera Manager](#) on page 212.

1. Log into the Cloudera Manager Admin Console.
2. Select **Administration > Settings**.
3. Click the **Security** category.
4. Configure the following TLS settings in the Cloudera Manager Server:

Property	Description
Use TLS Encryption for Agents	Enable TLS encryption for Agents connecting to the Server. The Agents will still connect to the defined agent listener port for Cloudera Manager (default: 7182). This property negotiates TLS connections to the service at this point.

5. Click Save Changes.

Step 2: Enable and Configure TLS on the Agent Hosts

To enable and configure TLS, you must specify values for the TLS properties in the `/etc/cloudera-scm-agent/config.ini` configuration file on all Agent hosts.

1. On the Agent host, open the `/etc/cloudera-scm-agent/config.ini` configuration file and edit the following property:

Property	Description
<code>use_tls</code>	Specify 1 to enable TLS on the Agent, or 0 (zero) to disable TLS.

2. Repeat this step on every Agent host. You can copy the Agent's `config.ini` file across all hosts since this file by default does not have host specific information within it. If you modify properties such as `listening_hostname` or `listening_ip` address in `config.ini`, you must configure the file individually for each host.

Step 3: Restart the Cloudera Manager Server

Restart the Cloudera Manager Server with the following command to activate the TLS configuration settings.

```
$ sudo service cloudera-scm-server restart
```

Step 4: Restart the Cloudera Manager Agents

On every Agent host, restart the Agent:

```
$ sudo service cloudera-scm-agent restart
```

Step 5: Verify that the Server and Agents are Communicating

In the Cloudera Manager Admin Console, open the **Hosts** page. If the Agents heartbeat successfully, TLS encryption is working properly.

Level 2: Configuring TLS Verification of Cloudera Manager Server by the Agents

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

Configuring Encryption

This level of TLS security requires that you provide a server certificate that is signed, either directly or through a chain, by a trusted root certificate authority (CA), to the Cloudera Manager Server. You must also provide the certificate of the CA that signed the Server certificate. For test environments, you can use a self-signed server certificate.



Note: If the Cloudera Manager Server certificate or the associated CA certificate is missing or expired, Agents will not communicate with the Cloudera Manager Server.

Step 1: Configure TLS encryption

If you have not done so, configure TLS encryption to use this level of security. For instructions, see [Configuring TLS Encryption Only for Cloudera Manager](#) on page 212 and [Level 1: Configuring TLS Encryption for Cloudera Manager Agents](#) on page 216.

Step 2: Copy the CA Certificate or Cloudera Manager Server's .pem file to the Agents

1. Agents can verify the Cloudera Manager Server using either the Server certificate or the associated root CA certificate. Pick any one of the following approaches to proceed:

- **Copy the Cloudera Manager Server .pem file to the Agent host**

1. For verification by the Agent, copy the Server .pem file (for example, `cmhost.pem`) to any directory on the Agent host. In the examples, this path is `/opt/cloudera/security/x509/cmhost.pem`.
2. On the Agent host, open the `/etc/cloudera-scm-agent/config.ini` configuration file and edit the following properties.

Property	Description
<code>verify_cert_file</code>	Point this property to the copied .pem file on the Agent host; in this example, <code>/opt/cloudera/security/x509/cmhost-cert.pem</code> .
<code>use_tls</code>	Set this property to 1.

OR

- **Copy the CA certificates to the Agent host**

1. If you have a CA-signed certificate, copy the root CA or intermediate CA certificates in PEM format to the Agent host. In the example, the CA certificates are copied to `/opt/cloudera/security/CAcerts/*`.
2. On the Agent host, open the `/etc/cloudera-scm-agent/config.ini` configuration file and edit the following properties.

Property	Description
<code>verify_cert_dir</code>	Point this property to the directory on the Agent host with the copied CA certificates; in the example, <code>/opt/cloudera/security/CAcerts/</code> .
<code>use_tls</code>	Set this property to 1.



Note: When configuring the `verify_cert_dir` property, ensure that the `openssl-perl` package is installed. The `openssl-perl` package comes with the `c_rehash` command required to generate the Subject Name hash values that will be linked to the certificates to make them usable. See the comments in the `config.ini` file for more information.

The following example is for RHEL-compatible systems. The package name for Debian-based systems is the same. After the package is installed, go to the CA certificate path and run the `c_rehash` command. This generates symbolic links to the certificate in that location, with `."` being the current path, as follows:

```
$ yum -y install openssl-perl
$ cd /opt/cloudera/security/CAcerts/
$ c_rehash .
  Doing .
  w2k8-1-root.pem => 4507f087.0
  w2k8-2-intermediary.pem => 082ba6df.0
$ ls -l
total 8.0K
lrwxrwxrwx 1 root root 23 Oct 6 22:44 082ba6df.0 ->
w2k8-2-intermediary.pem
lrwxrwxrwx 1 root root 15 Oct 6 22:44 4507f087.0 ->
w2k8-1-root.pem
-rw-r----- 1 root root 2.1K Oct 6 17:23 w2k8-1-root.pem
-rw-r----- 1 root root 2.8K Oct 6 17:23
w2k8-2-intermediary.pem
```

- Based on the approach you select in step 1, repeat the steps on every Agent host. You can copy the Agent's `config.ini` file across all hosts. However, if you modify properties such as `listening_hostname` or `listening_ip` address in `config.ini`, you must configure `config.ini` for each host individually.

Step 3: Restart the Cloudera Manager Agents

On every Agent host, restart the Agent:

```
$ sudo service cloudera-scm-agent restart
```

Step 4: Restart the Cloudera Management Services

To restart the Cloudera Management Service from the Cloudera Manager Admin Console:

- On the **Home** > **Status** tab, click



to the right of the service name and select **Restart**.

- Click **Start** on the next screen to confirm. When you see a **Finished** status, the service has restarted.

Step 5: Verify that the Server and Agents are communicating

In the Cloudera Manager Admin Console, open the **Hosts** page. If the Agents heartbeat successfully, the Server and Agents are communicating. If not, check the Agent log `/var/log/cloudera-scm-agent/cloudera-scm-agent.log`, which shows errors if the connection fails.

Level 3: Configuring TLS Authentication of Agents to the Cloudera Manager Server

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

This is the highest level of TLS security supported for Cloudera Manager Server-Agent communications, and requires you to create private keys and Certificate Signing Requests (CSR) for each cluster node. A Certificate Authority (CA)

Configuring Encryption

can then sign the CSR, thus providing a server certificate for each host. Agents will then need to authenticate themselves to Cloudera Manager using this server certificate.

This can be completed one of two ways, depending on the approach you choose to configuring TLS on your cluster.

- [Approach A](#) - Use OpenSSL to create private keys and request CA-signed certificates for every Agent on your cluster. Approach A is faster if you only need to enable TLS for Cloudera Manager Server-Agent communication.
- [Approach B](#) - Create a Java truststore file that contains the Agent and CA certificates, and authenticate Agents against this truststore file. If you plan to enable TLS communication for all CDH services cluster-wide, including Java-based components, consider using Approach B.

Steps for creating self-signed certificates are not included. Self-signed certificates are not recommended for production environments.



Note: Wildcard domain certificates and certificates using the SubjectAlternativeName extension are not supported at this time.

Step 1: Configure TLS encryption

If you have not already done so, you must configure TLS encryption to use this third level of security. For instructions, see [Configuring TLS Encryption Only for Cloudera Manager](#) on page 212 and [Configuring TLS Encryption for Cloudera Manager](#).

Step 2: Configure TLS Verification of Server Trust by Agents

If you have not already done so, you must configure TLS Verification of Server Trust by Agents. For instructions, see [Configuring TLS Authentication of Server to Agents](#).



Important:

Steps 3, 4, and 5 can be completed one of two ways, depending on the approach you choose to configuring TLS on your cluster.

- [Approach A](#) - Use OpenSSL to create private keys and request CA-signed certificates for every Agent on your cluster. Approach A is faster if you only need to enable TLS for Cloudera Manager Server-Agent communication.
- [Approach B](#) - Create a Java truststore file that contains the Agent and CA certificates, and authenticate Agents against this truststore file. If you plan to enable TLS communication for all CDH services cluster-wide, including Java-based components, consider using Approach B.

Irrespective of the path you select, it will still be possible to reuse OpenSSL private keys and certificates by exporting to a Java keystore and vice versa. For instructions, see [Private Key and Certificate Reuse Across Java Keystores and OpenSSL](#) on page 210.

After choosing an approach, follow steps 3-5 for all hosts in your cluster.

Approach A: Using OpenSSL to Create Private Keys and Request Agent Certificates

If the Cloudera Manager Server is running Management Services or CDH components (and therefore, has a Cloudera Manager Agent installed), you do not need to re-create a private key for the Server host. Follow the steps in [Private Key and Certificate Reuse Across Java Keystores and OpenSSL](#) on page 210 to reuse the host certificate. Follow steps 3-5 for all remaining cluster hosts.

Approach A Step 3: Generate the private key and certificate signing request for the Agent using OpenSSL.

Run the following command on the Agent, replacing `hostname` with your actual hostname. The `-subj` command line option allows you to provide the certificate subject as a single line. If you do not specify the certificate subject (`-subj`) as an argument, you will be prompted for the values of the certificate subject information. In that case, use the host

FQDN that Agents will use to connect from in the subject **First and Last name (CN)** question prompt. Country (C) requires a 2 letter country code. The "/" is replaced with "," in the actual CSR and private key file.

```
$ openssl req -subj
'/CN=hostname.sec.cloudera.com/OU=Support/O=Cloudera/L=Denver/ST=Colorado/C=US' \
-out /opt/cloudera/security/x509/hostname.csr -new -newkey rsa:2048 \
-keyout /opt/cloudera/security/x509/hostname.key -passout pass:password
```

password provides a password to protect the private key file. Keep the password in a safe place; you must provide a key password file to the Agent to complete configuration.

Approach A Step 4: Submit the certificate signing request to your CA and distribute the issued certificates.

The CSR file created (/opt/cloudera/security/x509/hostname.csr) is collected from cluster hosts for submission to the certificate authority (CA) for signing. In the example paths, you copy the issued CA-signed certificate file to /opt/cloudera/security/x509 on each cluster host. For easy management and tracking of files, name the files in the hostname.pem format, replacing hostname with the actual hostname.



Note: Certificate file extensions of .cer, .crt, and .pem are interchangeable. Rename the files so they have a .pem extension, and can therefore be used by the Agent and Hue (or any other Python-based component).

The CSR can be examined with the following command:

```
$ openssl req -text -noout -verify -in /opt/cloudera/security/x509/hostname.csr
```

The issued certificate file can be examined with the following command:

```
$ openssl x509 -in /opt/cloudera/security/x509/hostname.pem -text -noout
```

Approach A Step 5 (Optional): Import the OpenSSL private key and certificate into the per-host Java keystore.

Follow the steps in [Private Key and Certificate Reuse Across Java Keystores and OpenSSL](#) on page 210 for this step.



Important: If you are using Approach A, skip to [step 6](#) to continue.

Approach B: Creating a Java Keystore and Importing Signed Agent Certificates into it

If the Cloudera Manager Server is running Management Services or CDH components (and therefore, has a Cloudera Manager Agent installed), you do not need to re-create a private key for the Server host. Follow the steps in [Private Key and Certificate Reuse Across Java Keystores and OpenSSL](#) on page 210 to reuse the host certificate. Follow steps 3-5 for all remaining cluster hosts.

Approach B - Step 3: Create a Java Keystore and private key for a host

Create a Java Keystore and private key files for an Agent host as follows:

```
$ keytool -genkeypair -alias hostname -keyalg RSA -keystore \
/opt/cloudera/security/jks/hostname-keystore.jks -keysize 2048 -dname \
"CN=cmhost.sec.cloudera.com,OU=Support,O=Cloudera,L=Denver,ST=Colorado,C=US" \
-storepass password -keypass password
```

password provides a password to protect the private key file. Note the password in a safe place; you must provide a key password file to the Agent to complete configuration.

Configuring Encryption

Approach B - Step 4: Generate a certificate signing request and install the issued certificate into the Java Keystore

1. Generate a certificate signing request (CSR) and submit it to your CA for a signed certificate.

```
$ keytool -certreq -alias hostname \  
-keystore /opt/cloudera/security/jks/hostname-keystore.jks \  
-file /opt/cloudera/security/x509/hostname.csr \  
-storepass password -keypass password
```

2. If you are using a Private CA, first import the root CA certificate followed by the intermediary/subordinate CA certificates into the Java keystore created previously.

```
$ keytool -importcert -trustcacerts -alias RootCA -keystore \  
/opt/cloudera/security/jks/hostname-keystore.jks -file \  
/opt/cloudera/security/CAcerts/RootCA.cer -storepass password
```

Repeat the following for all subordinate/intermediary CA certificates presented.

```
$ keytool -importcert -trustcacerts -alias SubordinateCA -keystore \  
/opt/cloudera/security/jks/hostname-keystore.jks -file \  
/opt/cloudera/security/CAcerts/SubordinateCA.cer -storepass password
```

3. Copy the issued signed certificate file provided by your CA to the location from where it will be imported by the Cloudera Manager Agent and possibly Hue.

```
$ cp certificate-file.cer /opt/cloudera/security/x509/hostname.pem
```

4. Import the issued certificate file into the previously created Java keystore (.jks) with the following command:

```
$ keytool -import -trustcacerts -alias <hostname> \  
-keystore /opt/cloudera/security/jks/<hostname>-keystore.jks \  
-file /opt/cloudera/security/x509/<hostname>.pem -storepass password
```

Approach B - Step 5: Export the private key from the Java keystore and convert it with OpenSSL for reuse by Agent

Follow the steps in [Private Key and Certificate Reuse Across Java Keystores and OpenSSL](#) on page 210.

Step 6: Create a File that Contains the Password for the Key

The Agent reads the password from a text file, not from the command line. The password file allows you to use file permissions to protect the password. For our example the password file was created at, `/etc/cloudera-scm-agent/agentkey.pw`.

Step 7: Configure the Agent with its Private Key and Certificate

1. On the Agent host, open the `/etc/cloudera-scm-agent/config.ini` configuration file and edit the following properties:

Property	Description
<code>client_key_file</code>	Name of the client key file.
<code>client_keypw_file</code>	Name of the client key password file, <code>agentkey.pw</code> .
<code>client_cert_file</code>	Name of the client certificate file.

2. Repeat these steps on every Agent host.

Step 8: Verify that steps 3-7 Were Completed for every Agent Host in Your Cluster

Make sure each Agent's private key and certificate that you import into the Cloudera Manager Server's truststore is unique.

Step 9: Create a Truststore by Importing CA and Agent Certificates

Perform this step on the Cloudera Manager Server, where the new truststore is used to authenticate Agents.

Create a new truststore file (`/opt/cloudera/security/jks/truststore.jks`) and import the CA root and intermediary/subordinate certificates to this truststore. The new truststore functions like a keystore, containing only certificates and no private key.

1. Create a trusted keystore using the `keytool` command and import the root CA certificate to this truststore.

```
$ keytool -importcert -noprompt -keystore /opt/cloudera/security/jks/truststore.jks \
  -alias root_CA -file root.crt -storepass password
```

2. Import any remaining intermediary/subordinate CA certificates into the truststore.

```
$ keytool -importcert -noprompt -keystore /opt/cloudera/security/jks/truststore.jks
  -alias int_CA -file intermediate-CA.pem -storepass password
```

3. Save the `hostname.pem` certificate files from all cluster hosts in a single location. The Cloudera Manager Server can now import these host certificates (`hostname.pem`) into the new truststore.

```
$ keytool -keystore /opt/cloudera/security/jks/truststore.jks \
  -importcert -alias hostname -file hostname.pem -storepass password
```

Consider creating a `for` loop on a list of hostnames to speed up this process.

```
$ for HOST in `cat hostlist.txt`; do keytool -keystore
  /opt/cloudera/security/jks/truststore.jks \
  -importcert -alias $HOST -file $HOST.pem -storepass password
```

Step 10: Enable Agent Authentication and Configure the Cloudera Manager Server to Use the New Truststore

1. Log into the Cloudera Manager Admin Console.
2. Select **Administration > Settings**.
3. Click the **Security** category.
4. Configure the following TLS settings:

Setting	Description
Use TLS Authentication of Agents to Server	Select this option to enable TLS authentication of Agents to the Server.
Path to Truststore	Specify the full filesystem path to the truststore located on the Cloudera Manager Server host; in the example, <code>/opt/cloudera/security/jks/truststore.jks</code>
Truststore Password	Specify the password for the truststore.

5. Click **Save Changes** to save the settings.

Step 11: Restart the Cloudera Manager Server

```
$ sudo service cloudera-scm-server restart
```

Step 12: Restart the Cloudera Manager Agents

On every Agent host, restart the Agent:

```
$ sudo service cloudera-scm-agent restart
```

Step 13: Verify that the Server and Agents Are Communicating

In Cloudera Manager Admin Console, open the **Hosts** page. If the Agents heartbeat successfully, the Server and Agents are communicating. If they are not, you may see an error in the Server, such as a `null CA chain` error. This implies that either the truststore does not contain the Agent certificate, or the Agent is not presenting the certificate. Check all of your settings, and check the Server log to verify that TLS and Agent validation have been enabled correctly.

HTTPS Communication in Cloudera Manager

Both the Cloudera Manager Agent and the roles that make up the Cloudera Management Service use HTTPS to communicate with Cloudera Manager and CDH services. This topic aims to explain how the various aspects of HTTPS communication are handled by the Cloudera Manager Agents and the Cloudera Management Service roles.

Cloudera Manager Agents use HTTPS to communicate with HBase, HDFS, Impala, MapReduce, and YARN to collect monitoring data.

Cloudera Manager Agent

Configuring TLS communication between the Cloudera Manager Server and Agents is outlined in [Configuring TLS Security for Cloudera Manager](#) on page 211. You can configure the certificates available for server certificate verification using the `verify_cert_dir` parameter in the Agent `config.ini` file. See the comments in the `config.ini` file for a detailed explanation of this property. You can also use the existing value for the `verify_cert_file` parameter.

When the Cloudera Manager Agent communicates with CDH services using HTTPS:

- If `verify_cert_file` or `verify_cert_dir` are configured in the Agent `config.ini`, the Agent uses these settings to verify the server certificates. If these settings are not configured, no certificate verification occurs. If certificate verification is performed for the Cloudera Manager Server, it must also be performed for CDH daemons.
- An Agent never participates in mutual TLS authentication with any CDH service. Instead, each service has its own authentication scheme. Most services use Kerberos authentication, but Impala uses HTTP digest.

User Impact

This depends on how you use certificates.

- If you do not need certificate verification, do not configure `verify_cert_file` or `verify_cert_dir`. However, this leaves you vulnerable to man-in-the-middle attacks.
- If you are using a CA-signed certificate, configure the Agent accordingly. Adding new services or enabling TLS/SSL on a service requires no changes to the Agent configuration because the CA verifies the certificates used by any new servers brought online.
- If you are using self-signed certificates (recommended *only* on test clusters), the certificate for each new service that uses HTTPS must be available to the Agent. Modify the file pointed to by `verify_cert_file` (Agent restart required), or the directory pointed to by `verify_cert_dir`, to contain the new certificate.

Cloudera Management Services

Some Cloudera Management Service roles act as HTTPS clients when communicating with Cloudera Manager entities and CDH services.

You can verify server certificates in two ways:

- Configure a truststore through Cloudera Manager to perform certificate verification on the certificates of the servers with which it communicates. If this truststore is configured, it is used to verify server certificates.

OR

- If no truststore is configured through Cloudera Manager, the default Java truststore (`cacerts`) is used to verify certificates.

The following table shows Cloudera Management Service roles that act as HTTPS clients, and the corresponding Cloudera Manager entities that communicate with them as HTTPS servers. This table does not depict the entirety of the roles' communication, only communications over HTTPS.

Table 10: HTTPS Communication Between Cloudera Management Service Roles and Cloudera Manager Entities

Roles as HTTPS Clients	Communicating HTTPS Servers
Activity Monitor	<ul style="list-style-type: none"> • Cloudera Manager Server • JobTracker Web Server • Oozie server (may involve the load balancer in an HA configuration)
Host Monitor	<ul style="list-style-type: none"> • Cloudera Manager Server
Service Monitor	<ul style="list-style-type: none"> • Cloudera Manager Server • NameNode(s) Web Server(s) • Impala StateStore Web Server • YARN ResourceManager(s) Web Server(s) • YARN JobHistory Web Server • Oozie server (directly, not through the load balancer)
Event Server	<ul style="list-style-type: none"> • Cloudera Manager Server
Reports Manager	<ul style="list-style-type: none"> • Cloudera Manager Server • NameNode(s) Web Servers



Note: The Cloudera Navigator roles also act as HTTPS clients, but are outside the scope of this document.

The Cloudera Management Service roles communicate using HTTPS as follows:

- If the Cloudera Management Service **SSL Client Truststore File Location** parameter is configured, the roles use this truststore to verify server certificates. If this parameter is not set, the default Java truststore is used to verify certificates. Without using safety valves, you cannot verify certificates for some Cloudera Management Service roles but not for others. Nor can you verify certificates for only a subset of the HTTPS communication by a role.
- The Cloudera Management Service roles never participate in mutual TLS authentication with any CDH service or with the Cloudera Manager Server. Instead, each service has its own authentication scheme: Kerberos for most services, HTTP digest for Impala. For the Cloudera Manager Server, this authentication is session-based.

User Impact

This depends on how you use certificates:

- If you use a CA-signed certificate, configure the Cloudera Management Service **SSL Client Truststore File Location** parameter to point to a truststore that contains the CA certificate. Adding a new service or enabling TLS on an existing service requires no changes to the Cloudera Management Service configuration because the CA certificate verifies the certificates used by any new servers brought online. Alternatively, this CA-signed certificate can be added to the default Java truststore.
- If you are using self-signed certificates, the certificate for each new service that uses HTTPS must be available to the Agent. You must modify the truststore pointed to by the Cloudera Management Service **SSL Client Truststore File Location** parameter. Truststore changes are required on each host on which a Cloudera Management Service daemon is running. Changes to the truststore do not require a role restart, and should be picked up within 10 seconds by default.

If the Cloudera Management Service **SSL Client Truststore File Location** is not used, the certificate must be made available in the default Java truststore. The Cloudera Management Service role must be restarted for this change to take effect.

Troubleshooting TLS/SSL Connectivity

The `openssl` tool can be run from the host that is running the Cloudera Manager Agent or client service that should be inspected for connectivity issues. You should also test whether the certificate in use by the host is recognized by a trusted CA during the TLS/SSL negotiation.

Use the following command to inspect the connection.

```
$ openssl s_client -connect [host.fqdn.name]:[port]
```

For example:

```
$ openssl s_client -connect test1.sec.cloudera.com:7183
```

A return code 0 means `openssl` was able to establish trust of the server through its library of trusted public CAs. If the certificate was self-signed (recommended *only* on test clusters) or provided by a private CA it might be necessary to add the private CA or self-signed certificate to the truststore using the `openssl` command. Adding the path to the root CA, `-CAfile </path/to/root-ca.pem>`, should allow `openssl` to verify your self-signed or private CA-signed certificate as follows:

```
$ openssl s_client -connect test1.sec.cloudera.com:7183 -CAfile \
/opt/cloudera/security/CAcerts/RootCA.pem
```

Note that providing only the Root CA certificate is necessary to establish trust for this test. The result from the command is successful when you see the return code 0 as follows:

```
...
Verify return code: 0 (ok)
---
```

By default, the Cloudera Manager Server writes logs to the `/etc/cloudera-scm-server/cloudera-scm-server.log` file on startup. Successful start of the server process with the certificate will show logs similar to the following:

```
2014-10-06 21:33:47,515 INFO WebServerImpl:org.mortbay.log: jetty-6.1.26.cloudera.2
2014-10-06 21:33:47,572 INFO WebServerImpl:org.mortbay.log: Started
SslSelectChannelConnector@0.0.0.0:7183
2014-10-06 21:33:47,573 INFO WebServerImpl:org.mortbay.log: Started
SelectChannelConnector@0.0.0.0:7180
2014-10-06 21:33:47,573 INFO WebServerImpl:com.cloudera.server.cmf.WebServerImpl: Started
Jetty server.
```

Using Self-Signed Certificates for TLS

Self-signed certificates should not be used for production deployments. However, for testing and other non-production purposes, self-signed certificates let you quickly obtain the certificates needed for [Step 1: Create the Cloudera Manager Server Keystore, Generate a Certificate Request, and Install the Certificate](#) on page 213.

Replace paths, file names, aliases, and other examples in the commands below for your system.

1. Create the directory for the certificates:

```
$ mkdir -p /opt/cloudera/security/x509/ /opt/cloudera/security/jks/
$ cd /opt/cloudera/security/jks
```

Use `chmod/chown` to change ownership of the `/opt/cloudera/security/jks` directory to give Cloudera Manager access to the directory.

2. Generate the key pair and self-signed certificate, and store these in the keystore (`example.keystore`). Set `-keypass` to the same value as `-storepass` (Cloudera Manager does not support separate values for `keypass` and `storepass`.)



Note: For CN, use the FQDN of the Cloudera Manager Server host to avoid a `java.io.IOException: HTTPS hostname wrong exception`.

```
$ keytool -genkeypair -keystore example.keystore -keyalg RSA -alias cmhost \
-dname "CN=cmhost.sec.example.com,OU=Security,O=Example,L=Denver,ST=Colorado,C=US"
-storepass password -keypass password
```

3. Copy the default Java truststore (`cacerts`) to the alternate system truststore (`jssecacerts`). Self-signed certificates are appended to `jssecacerts` without modifying the default `cacerts` file.

```
$ sudo cp $JAVA_HOME/jre/lib/security/cacerts $JAVA_HOME/jre/lib/security/jssecacerts
```

4. Export the certificate from the keystore (`example.keystore`).

```
$ keytool -export -alias cmhost -keystore example.keystore -rfc -file selfsigned.cer
```

5. Copy the self-signed certificate (`selfsigned.cer`) to the `/opt/cloudera/security/x509/` directory.

```
$ cp selfsigned.cer /opt/cloudera/security/x509/cmhost.pem
```

6. Import the public key into the alternate system truststore (`jssecacerts`), so that any process that runs with Java on this machine will trust the key. The default password for the Java truststore is `changeit`. Do not use the password created for the keystore in Step 2.

```
$ keytool -import -alias cmhost -file /opt/cloudera/security/jks/selfsigned.cer \
-keystore $JAVA_HOME/jre/lib/security/jssecacerts -storepass changeit
```



Important: Repeat this process on each machine in the cluster.

7. Rename the keystore, such as from `example.keystore` to `cmhost-keystore.jks`:

```
$ mv /opt/cloudera/security/jks/example.keystore
/opt/cloudera/security/jks/cmhost-keystore.jks
```

You can also delete the certificate: it has been added to the keystore at `/opt/cloudera/security/x509/cmhost.pem`.

```
$ rm /opt/cloudera/security/selfsigned.cer
```

The self-signed certificate set up is complete. You can continue configuring TLS Level 1 as detailed in [Step 2: Enable HTTPS for the Cloudera Manager Admin Console and Specify Server Keystore Properties](#) on page 215.

Configuring TLS/SSL for the Cloudera Navigator Data Management Component



Important: The following instructions assume you have a Java keystore set up on the Navigator Metadata Server host.

To enable SSL communication between the Cloudera Navigator Metadata Server and its clients:

1. Open the Cloudera Manager Admin Console and go to the **Cloudera Management Service**.
2. Click the **Configuration** tab.
3. Select **Scope > Navigator Metadata Server**.

Configuring Encryption

4. Select **Category > Security**.
5. Edit the following properties according to your cluster configuration.

Property	Description
Enable TLS/SSL for Navigator Metadata Server	Encrypt communication between clients and Navigator Metadata Server using Transport Layer Security (TLS) (formerly known as Secure Socket Layer (SSL)).
TLS/SSL Keystore File Location	The path to the TLS/SSL keystore file containing the server certificate and private key used for TLS/SSL. Used when Navigator Metadata Server is acting as a TLS/SSL server. The keystore must be in JKS format.
TLS/SSL Keystore File Password	The password for the Navigator Metadata Server JKS keystore file.
TLS/SSL Keystore Key Password	The password that protects the private key contained in the JKS keystore used when Navigator Metadata Server is acting as a TLS/SSL server.

6. Click **Save Changes** to commit the changes.
7. Restart the Navigator Metadata server.



Note: Once you have enabled TLS/SSL, the Quick Links in Cloudera Manager pointing to the Cloudera Navigator UI will not work as they use HTTP, not HTTPS.

Configuring TLS/SSL for Cloudera Management Service Roles

To enable TLS/SSL communication between Cloudera Management Service roles, and CDH services and the Cloudera Manager Server:

1. Open the Cloudera Manager Administration Console and go to the **Cloudera Management Service**.
2. Click the **Configuration** tab.
3. Select **Scope > Cloudera Management Service (Service-Wide)**.
4. Select **Category > Security**.
5. Edit the following TLS/SSL properties according to your cluster configuration.

Property	Description
TLS/SSL Client Truststore File Location	Path to the client truststore file used in HTTPS communication. The contents of this truststore can be modified without restarting the Cloudera Management Service roles. By default, changes to its contents are picked up within ten seconds.
TLS/SSL Client Truststore File Password	Password for the client truststore file.

6. Click **Save Changes** to commit the changes.
7. Restart the Cloudera Management Service. For more information, see [HTTPS Communication in Cloudera Manager](#) on page 224.

Configuring TLS/SSL Encryption for CDH Services

In addition to configuring Cloudera Manager cluster to use TLS/SSL (as detailed, starting with [Configuring TLS Security for Cloudera Manager](#) on page 211), the various CDH services running on the cluster should also be configured to use TLS/SSL. The process of configuring TLS/SSL varies by component, so follow the steps below as needed for your system. Before trying to configure TLS/SSL, however, be sure your cluster meets the [prerequisites](#).



Note: TLS/SSL for Hadoop core services—HDFS, MapReduce, and YARN—must be enabled as a group. TLS/SSL for other components such as HBase, Hue, and Oozie can be enabled independently.

Prerequisites

- Cloudera recommends securing a cluster using Kerberos authentication before enabling encryption such as TLS/SSL on a cluster. If you enable TLS/SSL for a cluster that does not already have Kerberos authentication configured, a warning will be displayed.
- The following sections assume that you have created all the certificates required for TLS/SSL communication. If not, for information on how to do this, see [Creating Certificates](#).
- The certificates and keys to be deployed in your cluster should be organized into the appropriate set of keystores and truststores. For more information, see [Creating Java Keystores and Truststores](#) on page 207.



Note: Cloudera Manager and CDH components support either TLS 1.0, TLS 1.1, or TLS 1.2, but not SSL 3.0. References to SSL only continue because of its widespread use in technical jargon.

Hadoop Services as TLS/SSL Servers and Clients

Hadoop services differ in their use of TLS/SSL as follows:

- HDFS, MapReduce, and YARN daemons act as both TLS/SSL servers and clients.
- HBase daemons act as TLS/SSL servers only.
- Oozie daemons act as TLS/SSL servers only.
- Hue acts as an TLS/SSL client to all of the above.

Daemons that act as TLS/SSL servers load the keystores when starting up. When a client connects to an TLS/SSL server daemon, the server transmits the certificate loaded at startup time to the client, which then uses its truststore to validate the server's certificate.

Compatible Certificate Formats for Hadoop Components

Component	Compatible Certificate Format
HDFS	Java Keystore
MapReduce	Java Keystore
YARN	Java Keystore
Hue	PEM
Hive (for communication between Hive clients and HiveServer2)	Java Keystore
HBase	Java Keystore
Impala	PEM
Oozie	Java Keystore

Configuring TLS/SSL for HDFS, YARN and MapReduce

Required Role: [Configurator](#), [Cluster Administrator](#), or [Full Administrator](#)

TLS/SSL for the core Hadoop services—HDFS, MapReduce, and YARN—must be enabled as a group. Because most clusters run either MapReduce or YARN, not both, you will typically enable HDFS and YARN, or HDFS and MapReduce. Enabling TLS/SSL on HDFS is required before it can be enabled on either MapReduce or YARN.



Note: If you enable TLS/SSL for HDFS, you must also enable it for MapReduce or YARN.

The steps below include enabling Kerberos authentication for HTTP Web-Consoles. Enabling TLS/SSL for the core Hadoop services on a cluster without enabling authentication displays a warning.

Before You Begin

- Before enabling TLS/SSL, keystores containing certificates bound to the appropriate domain names will need to be accessible on all hosts on which at least one HDFS, MapReduce, or YARN daemon role is running.
- Since HDFS, MapReduce, and YARN daemons act as TLS/SSL clients as well as TLS/SSL servers, they must have access to truststores. In many cases, the most practical approach is to deploy truststores to all hosts in the cluster, as it may not be desirable to determine in advance the set of hosts on which clients will run.
- Keystores for HDFS, MapReduce and YARN must be owned by the `hadoop` group, and have permissions `0440` (that is, readable by owner and group). Truststores must have permissions `0444` (that is, readable by all)
- Cloudera Manager supports TLS/SSL configuration for HDFS, MapReduce and YARN at the service level. For each of these services, you must specify absolute paths to the keystore and truststore files. These settings apply to all hosts on which daemon roles of the service in question run. Therefore, the paths you choose must be valid on all hosts.

An implication of this is that the keystore file names for a given service must be the same on all hosts. If, for example, you have obtained separate certificates for HDFS daemons on hosts `node1.example.com` and `node2.example.com`, you might have chosen to store these certificates in files called `hdfs-node1.keystore` and `hdfs-node2.keystore` (respectively). When deploying these keystores, you must give them both the same name on the target host — for example, `hdfs.keystore`.

- Multiple daemons running on a host can share a certificate. For example, in case there is a DataNode and an Oozie server running on the same host, they can use the same certificate.

Configuring TLS/SSL for HDFS

1. Go to the **HDFS** service.
2. Click the **Configuration** tab.
3. Select **Scope > HDFS (Service-Wide)**.
4. Select **Category > Security**.
5. In the Search field, type **TLS/SSL** to show the TLS/SSL properties (found under the **Service-Wide > Security** category).
6. Edit the following properties according to your cluster configuration:

Property	Description
Hadoop TLS/SSL Server Keystore File Location	Path to the keystore file containing the server certificate and private key.
Hadoop TLS/SSL Server Keystore File Password	Password for the server keystore file.
Hadoop TLS/SSL Server Keystore Key Password	Password that protects the private key contained in the server keystore.

7. If you are not using the default truststore, configure TLS/SSL client truststore properties:




Important: The HDFS properties below define a cluster-wide default truststore that can be overridden by YARN and MapReduce (see the **Configuring TLS/SSL for YARN and MapReduce** section below).

Property	Description
Cluster-Wide Default TLS/SSL Client Truststore Location	Path to the client truststore file. This truststore contains certificates of trusted servers, or of Certificate Authorities trusted to identify servers.
Cluster-Wide Default TLS/SSL Client Truststore Password	Password for the client truststore file.

8. **(Optional)** Cloudera recommends you enable web UI authentication for the HDFS service. Web UI authentication uses SPNEGO. After enabling this, you cannot access the Hadoop web consoles without a valid Kerberos ticket and proper client-side configuration. For more information, see [Using a Web Browser to Access an URL Protected by Kerberos HTTP SPNEGO](#) on page 190.

To enable web UI authentication, enter **web consoles** in the Search field to bring up the **Enable Authentication for HTTP Web-Consoles** property (found under the **Service-Wide>Security** category). Check the property to enable web UI authentication.

Enable Authentication for HTTP Web-Consoles	Enables authentication for Hadoop HTTP web-consoles for all roles of this service. <div style="border: 1px solid black; padding: 5px; margin-top: 10px;">  Note: This is effective only if security is enabled for the HDFS service. </div>
----------------------------------------------------	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

9. Click **Save Changes**.
10. Follow the procedure described in the following **Configuring TLS/SSL for YARN and MapReduce** section, at the end of which you will be instructed to restart all the affected services (HDFS, MapReduce and YARN).

Configuring TLS/SSL for YARN or MapReduce

Perform the following steps to configure TLS/SSL for the YARN or MapReduce services:

- Go to the **YARN** or **MapReduce** service.
- Click the **Configuration** tab.
- Select **Scope > service name (Service-Wide)**.
- Select **Category > Security**.
- Locate the **<property name>** property or search for it by typing its name in the Search box.
- In the Search field, type **TLS/SSL** to show the TLS/SSL properties (found under the **Service-Wide > Security** category).
- Edit the following properties according to your cluster configuration:

Property	Description
Hadoop TLS/SSL Server Keystore File Location	Path to the keystore file containing the server certificate and private key.
Hadoop TLS/SSL Server Keystore File Password	Password for the server keystore file.
Hadoop TLS/SSL Server Keystore Key Password	Password that protects the private key contained in the server keystore.

8. Configure the following TLS/SSL client truststore properties for MRv1 or YARN only if you want to override the cluster-wide defaults set by the HDFS properties configured above.


Property	Description
TLS/SSL Client Truststore File Location	Path to the client truststore file. This truststore contains certificates of trusted servers, or of Certificate Authorities trusted to identify servers.

Configuring Encryption

Property	Description
TLS/SSL Client Truststore File Password	Password for the client truststore file.

9. Cloudera recommends you enable Web UI authentication for the service in question.

Enter **web consoles** in the Search field to bring up the **Enable Authentication for HTTP Web-Consoles** property (found under the **Service-Wide>Security** category). Check the property to enable web UI authentication.

Enable Authentication for HTTP Web-Consoles	Enables authentication for Hadoop HTTP web-consoles for all roles of this service. <div style="border: 1px solid green; padding: 5px;"> Note: This is effective only if security is enabled for the HDFS service.</div>
----------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

10 Click **Save Changes** to commit the changes.

11 Go to the **HDFS** service

12 Click the **Configuration** tab.

13 Type **Hadoop SSL Enabled** in the Search box.

14 Select the **Hadoop SSL Enabled** property to enable SSL communication for HDFS, MapReduce, and YARN.

Property	Description
Hadoop TLS/SSL Enabled	Enable TLS/SSL encryption for HDFS, MapReduce, and YARN web UIs, as well as encrypted shuffle for MapReduce and YARN.

15 Click **Save Changes** to commit the changes.

16 Restart all affected services (HDFS, MapReduce and YARN), as well as their dependent services.

Configuring TLS/SSL for HBase

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator, Full Administrator**)

Before You Begin

- Before enabling TLS/SSL, ensure that keystores containing certificates bound to the appropriate domain names will need to be accessible on all hosts on which at least one HBase daemon role is running.
- Keystores for HBase must be owned by the `hbase` group, and have permissions `0440` (that is, readable by owner and group).
- You must specify absolute paths to the keystore and truststore files. These settings apply to all hosts on which daemon roles of the HBase service run. Therefore, the paths you choose must be valid on all hosts.
- Cloudera Manager supports the TLS/SSL configuration for HBase at the service level. Ensure you specify absolute paths to the keystore and truststore files. These settings apply to all hosts on which daemon roles of the service in question run. Therefore, the paths you choose must be valid on all hosts.

An implication of this is that the keystore file names for a given service must be the same on all hosts. If, for example, you have obtained separate certificates for HBase daemons on hosts `node1.example.com` and `node2.example.com`, you might have chosen to store these certificates in files called `hbase-node1.keystore` and `hbase-node2.keystore` (respectively). When deploying these keystores, you must give them both the same name on the target host — for example, `hbase.keystore`.

Configuring TLS/SSL for HBase Web UIs

The steps for configuring and enabling TLS/SSL for HBase are similar to those for HDFS, YARN and MapReduce:

1. Go to the **HBase service**

- Click the **Configuration** tab.
- Select **Scope > HBASE (Service-Wide)**.
- Select **Category > Security**.
- In the Search field, type **TLS/SSL** to show the HBase TLS/SSL properties.
- Edit the following TLS/SSL properties according to your cluster configuration:

Table 11: HBase TLS/SSL Properties

Property	Description
HBase TLS/SSL Server JKS Keystore File Location	Path to the keystore file containing the server certificate and private key used for encrypted web UIs.
HBase TLS/SSL Server JKS Keystore File Password	Password for the server keystore file used for encrypted web UIs.
HBase TLS/SSL Server JKS Keystore Key Password	Password that protects the private key contained in the server keystore used for encrypted web UIs.

- Check the **Web UI TLS/SSL Encryption Enabled** property.

Web UI TLS/SSL Encryption Enabled	Enable TLS/SSL encryption for the HBase Master, RegionServer, Thrift Server, and REST Server web UIs.
------------------------------------------	-------------------------------------------------------------------------------------------------------

- Click **Save Changes** to commit the changes.
- Restart the HBase service.

Configuring TLS/SSL for HBase REST Server

- Go to the **HBase service**
- Click the **Configuration** tab.
- Select **Scope > HBase REST Server**.
- Select **Category > Security**.
- In the Search field, type **TLS/SSL REST** to show the HBase REST TLS/SSL properties.
- Edit the following TLS/SSL properties according to your cluster configuration:

Property	Description
Enable TLS/SSL for HBase REST Server	Encrypt communication between clients and HBase REST Server using Transport Layer Security (TLS).
HBase REST Server TLS/SSL Server JKS Keystore File Location	The path to the TLS/SSL keystore file containing the server certificate and private key used for TLS/SSL. Used when HBase REST Server is acting as a TLS/SSL server. The keystore must be in JKS format.file.
HBase REST Server TLS/SSL Server JKS Keystore File Password	The password for the HBase REST Server JKS keystore file.
HBase REST Server TLS/SSL Server JKS Keystore Key Password	The password that protects the private key contained in the JKS keystore used when HBase REST Server is acting as a TLS/SSL server.

- Click **Save Changes** to commit the changes.
- Restart the HBase service.

Configuring TLS/SSL for HBase Thrift Server

- Go to the **HBase service**

Configuring Encryption

2. Click the **Configuration** tab.
3. Select **Scope > HBase Thrift Server**.
4. Select **Category > Security**.
5. In the Search field, type **TLS/SSL Thrift** to show the HBase Thrift TLS/SSL properties.
6. Edit the following TLS/SSL properties according to your cluster configuration:

Property	Description
Enable TLS/SSL for HBase Thrift Server over HTTP	Encrypt communication between clients and HBase Thrift Server over HTTP using Transport Layer Security (TLS).
HBase Thrift Server over HTTP TLS/SSL Server JKS Keystore File Location	The path to the TLS/SSL keystore file containing the server certificate and private key used for TLS/SSL. Used when HBase Thrift Server over HTTP is acting as a TLS/SSL server. The keystore must be in JKS format.
HBase Thrift Server over HTTP TLS/SSL Server JKS Keystore File Password	The password for the HBase Thrift Server JKS keystore file.
HBase Thrift Server over HTTP TLS/SSL Server JKS Keystore Key Password	The password that protects the private key contained in the JKS keystore used when HBase Thrift Server over HTTP is acting as a TLS/SSL server.

7. Click **Save Changes** to commit the changes.
8. Restart the HBase service.

Configuring TLS/SSL for Flume Thrift Source and Sink

This topic describes how to enable TLS/SSL communication between Flume's Thrift source and sink.

The following tables list the properties that must be configured to enable TLS/SSL communication between Flume's Thrift source and sink instances.

Table 12: Thrift Source TLS/SSL Properties

Property	Description
ssl	Set to <code>true</code> to enable TLS/SSL encryption.
keystore	Path to a Java keystore file. Required for TLS/SSL.
keystore-password	Password for the Java keystore. Required for TLS/SSL.
keystore-type	The type of the Java keystore. This can be JKS or PKCS12.

Table 13: Thrift Sink TLS/SSL Properties

Property	Description
ssl	Set to <code>true</code> to enable TLS/SSL for this ThriftSink. When configuring TLS/SSL, you can optionally set the following <code>truststore</code> , <code>truststore-password</code> and <code>truststore-type</code> properties. If a custom truststore is not specified, Flume will use the default Java JSE truststore (typically <code>jssecacerts</code> or <code>cacerts</code> in the Oracle JRE) to verify the remote Thrift Source's TLS/SSL credentials.
truststore	(Optional) The path to a custom Java truststore file.
truststore-password	(Optional) The password for the specified truststore.

Property	Description
truststore-type	(Optional) The type of the Java truststore. This can be JKS or any other supported Java truststore type.

Make sure you are configuring TLS/SSL for *each* Thrift source and sink instance. For example, to the existing `flume.conf` file, for agent `a1`, source `r1`, and sink `k1`, you would add the following properties:

```
# TLS/SSL properties for Thrift source s1
a1.sources.r1.ssl=true
a1.sources.r1.keystore=<path/to/keystore>
a1.sources.r1.keystore-password=<keystore password>
a1.sources.r1.keystore-type=<keystore type>

# TLS/SSL properties for Thrift sink k1
a1.sinks.k1.ssl=true
a1.sinks.k1.truststore=<path/to/truststore>
a1.sinks.k1.truststore-password=<truststore password>
a1.sinks.k1.truststore-type=<truststore type>
```

Configure these sets of properties for more instances of the Thrift source and sink as required. You can use either Cloudera Manager or the command line to edit the `flume.conf` file.

Using Cloudera Manager

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)

1. Open the Cloudera Manager Admin Console and go to the **Flume** service.
2. Click the **Configuration** tab.
3. Select **Scope > Agent**.
4. Select **Category > Main**.
5. Edit the **Configuration File** property and add the Thrift source and sink properties for each Thrift source and sink instance as described above to the configuration file.
6. Click **Save Changes** to commit the changes.
7. Restart the Flume service.

Using the Command Line

Go to the `/etc/flume-ng/conf/flume.conf` file and add the Thrift source and sink properties for each Thrift source and sink instance as described above.

Configuring Encrypted Communication Between HiveServer2 and Client Drivers

Starting with CDH 5.5, encryption for HiveServer2 clients has been decoupled from the authentication mechanism. This means you can use either SASL QOP or TLS/SSL to encrypt traffic between HiveServer2 and its clients, irrespective of whether Kerberos is being used for authentication. Previously, the JDBC client drivers only supported SASL QOP encryption on Kerberos-authenticated connections.

SASL QOP encryption is better suited for encrypting RPC communication and may result in performance issues when dealing with large amounts of data. Move to using TLS/SSL encryption to avoid such issues.

Configuring Encrypted Client/Server Communication Using TLS/SSL

You can use either the Cloudera Manager or the command-line instructions described below to enable TLS/SSL encryption for JDBC/ODBC client connections to HiveServer2. For background information on setting up TLS/SSL truststores and keystores, see [TLS/SSL Certificates Overview](#) on page 205.



Note: Cloudera Manager and CDH components support either TLS 1.0, TLS 1.1, or TLS 1.2, but not SSL 3.0. References to SSL only continue because of its widespread use in technical jargon.

Configuring Encryption

Using Cloudera Manager

The steps for configuring and enabling TLS/SSL for Hive are as follows:

1. Open the Cloudera Manager Admin Console and go to the Hive service.
2. Click the **Configuration** tab.
3. Select **Scope > Hive (Service-Wide)**.
4. Select **Category > Security**.
5. In the Search field, type **TLS/SSL** to show the Hive properties.
6. Edit the following properties according to your cluster configuration.

Table 14: Hive TLS/SSL Properties

Property	Description
Enable TLS/SSL for HiveServer2	Enable support for encrypted client-server communication using Transport Layer Security (TLS) for HiveServer2 connections.
HiveServer2 TLS/SSL Server JKS Keystore File Location	Path to the TLS keystore.
HiveServer2 TLS/SSL Server JKS Keystore File Password	Password for the TLS keystore.

7. Click **Save Changes** to commit the changes.
8. Restart the Hive service.

Using the Command Line

- To enable TLS/SSL, add the following configuration parameters to `hive-site.xml` :

```
<property>
  <name>hive.server2.use.SSL</name>
  <value>>true</value>
  <description>enable/disable SSL </description>
</property>

<property>
  <name>hive.server2.keystore.path</name>
  <value>keystore-file-path</value>
  <description>path to keystore file</description>
</property>

<property>
  <name>hive.server2.keystore.password</name>
  <value>keystore-file-password</value>
  <description>keystore password</description>
</property>
```

- The keystore must contain the server's certificate.
- The JDBC client must add the following properties in the connection URL when connecting to a HiveServer2 using TLS/SSL:

```
;ssl=true[ ;sslTrustStore=<Trust-Store-Path>;trustStorePassword=<Trust-Store-password> ]
```

- Make sure one of the following is true:
 - *Either*: `sslTrustStore` points to the trust store file containing the server's certificate; for example:

```
jdbc:hive2://localhost:10000/default;ssl=true;\
sslTrustStore=/home/usrl/ssl/trust_store.jks;trustStorePassword=xyz
```

- *or*: the Trust Store arguments are set using the Java system properties `javax.net.ssl.trustStore` and `javax.net.ssl.trustStorePassword`; for example:

```
java -Djavax.net.ssl.trustStore=/home/usr1/ssl/trust_store.jks
-Djavax.net.ssl.trustStorePassword=xyz \
MyClass jdbc:hive2://localhost:10000/default;ssl=true
```

For more information on using self-signed certificates and the Trust Store, see the Oracle Java SE [keytool](#) page.

Configuring Encrypted Client/Server Communication Using SASL QOP

Traffic between the Hive JDBC or ODBC drivers and HiveServer2 can be encrypted using plain SASL QOP encryption which allows you to preserve data integrity (using checksums to validate message integrity) and confidentiality (by encrypting messages). This can be enabled by setting the `hive.server2.thrift.sasl.qop` property in `hive-site.xml`. For example,

```
<property>
<name>hive.server2.thrift.sasl.qop</name>
<value>auth-conf</value>
<description>Sasl QOP value; one of 'auth', 'auth-int' and 'auth-conf'</description>
</property>
```

Valid settings for the `value` field are:

- `auth`: Authentication only (default)
- `auth-int`: Authentication with integrity protection
- `auth-conf`: Authentication with confidentiality protection

The parameter value that you specify above in the HiveServer2 configuration, should match that specified in the Beeline client connection JDBC URL. For example:

```
!connect jdbc:hive2://ip-10-5-15-197.us-west-2.compute.internal:10000/default; \
principal=hive/_HOST@US-WEST-2.COMPUTE.INTERNAL;sasl.qop=auth-conf
```

Configuring TLS/SSL for Hue

This topic describes how to enable TLS/SSL communication for Hue:

Hue as a TLS/SSL Client

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)

Hue acts as an TLS/SSL client when communicating with Oozie, HBase and core Hadoop services. This means Hue may have to authenticate HDFS, MapReduce, and YARN daemons, as well as the HBase Thrift Server, and will need their certificates (or the relevant CA certificate) in its truststore.

Deploying the Hue Truststore:

You can create the Hue truststore by consolidating certificates of all TLS/SSL-enabled servers (or a single CA certificate chain) that Hue communicates with into one file. This will generally include certificates of all the HDFS, MapReduce and YARN daemons, and other TLS/SSL-enabled services such as Oozie..

The Hue truststore must be in PEM format whereas other services use JKS format by default. Hence, to populate the Hue truststore, you will need to extract the certificates from Hadoop's JKS keystores and convert them to PEM format. The following example assumes that `hadoop-server.keystore` contains the server certificate identified by alias `foo-1.example.com` and password `example123`.

```
$ keytool -exportcert -keystore hadoop-server.keystore -alias foo-1.example.com \
-storepass example123 -file foo-1.cert
$ openssl x509 -inform der -in foo-1.cert > foo-1.pem
```

Configuring Encryption

Once you've done this for each host in the cluster, you can concatenate the PEM files into one PEM file that can serve as the Hue truststore.

```
cat foo-1.pem foo-2.pem ... > huetrust.pem
```



Note: Ensure the final PEM truststore is deployed in a location that is accessible by the Hue service.

In Cloudera Manager, set TLS/SSL Truststore to the path of the consolidated PEM file, `huetrust.pem`:

1. Logon to Cloudera Manager Admin Console and go to the **Hue service**.
2. Click **Configuration**.
3. Select **Scope > Hue Server**.
4. Select **Category > Security**.
5. Find the property, `TLS/SSL Truststore`.
6. Enter the path to `huetrust.pem` on the host running the Hue web server.
7. Click **Save Changes**.
8. Restart the Hue service.

Hue as a TLS/SSL Server

Hue expects certificates and keys to be stored in PEM format. When managing certificates and keys for such services, using the `openssl` tool may be more convenient. To configure Hue to use HTTPS, you can generate a private key and certificate as described in [Creating Certificates](#) on page 206. Since Hue uses certificates in PEM format, you can reuse a host's existing Java keystore by converting it to the PEM format. For instructions, see [Conversion from Java Keystore to OpenSSL](#) on page 211.

Ensure secure session cookies for Hue have been enabled in `hue.ini` under `[desktop]>[[session]]`.

```
[desktop]
[[session]]
secure=true
```

Enabling TLS/SSL for the Hue Server at the Command Line

If you are not using Cloudera Manager, update the following properties in `hue.ini` under `[desktop]`.

```
[desktop]
ssl_certificate=/path/to/server.cert
ssl_private_key=/path/to/server.key
ssl_password=<private_key_password>
```

You can also store `ssl_password` more securely in a script and set this parameter instead:

```
ssl_password_script=<your_hue_passwords_script.sh>
```

For more, see [Storing Hue Passwords in a Script](#) on page 240.

Enabling TLS/SSL for the Hue Server in Cloudera Manager

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)

Perform the following steps in Cloudera Manager to enable TLS/SSL for the Hue web server.

1. Open the Cloudera Manager Admin Console and go to the **Hue service**.
2. Click **Configuration**.
3. Select **Scope > Hue Server**.
4. Select **Category > Security**.

5. Edit the following **TLS/SSL** properties according to your cluster configuration.

Property	Description
Enable HTTPS	Enable HTTPS for the Hue web server.
Local Path to TLS/SSL Certificate	Path to the TLS/SSL certificate on the host running the Hue web server.
Local Path to TLS/SSL Private Key	Path to the TLS/SSL private key on the host running the Hue web server.

If the private key has a password:

- a. Select **Scope > Hue-1 (Service-Wide)**.
- b. Select **Category > Advanced**.
- c. Locate the field, **Hue Service Advanced Configuration Snippet (Safety Valve) for hue_safety_valve.ini**.
- d. Add the TLS/SSL password parameter in the `[desktop]` section as follows:

```
[desktop]
ssl_password=<private_key_password>
```

You can also store `ssl_password` more securely in a script and set this parameter instead:

```
ssl_password_script=<your_hue_passwords_script.sh>
```

For more, see [Storing Hue Passwords in a Script](#) on page 240.

If more than one role group applies to this configuration, edit the value for the appropriate role group. See [Modifying Configuration Properties Using Cloudera Manager](#).

6. Click **Save Changes**.
7. Restart the Hue service.

For more details on configuring Hue with TLS/SSL, see this [blog post](#).

Enabling Hue TLS/SSL Communication with HiveServer2

In CDH 5.5.x and higher, HiveServer2 is enabled for TLS/SSL communication by default.

By providing a CA certificate, private key, and public certificate, Hue can communicate with HiveServer2 over TLS/SSL. You can now configure the following properties in the `[beeswax]` section under `[[ssl]]` in the Hue configuration file, `hue.ini`.

<code>enabled</code>	Choose to enable/disable TLS/SSL communication for this server. Default: <code>false</code>
<code>cacerts</code>	Path to Certificate Authority certificates. Default: <code>/etc/hue/cacerts.pem</code>
<code>validate</code>	Choose whether Hue should validate certificates received from the server. Default: <code>true</code>

Related Information

- [Configuring Encrypted Communication Between HiveServer2 and Client Drivers](#) on page 235

Enabling Hue TLS/SSL Communication with Impala

By providing a CA certificate, private key, and public certificate, Hue can communicate with Impala over TLS/SSL. You can configure the following properties in the `[impala]` section under `[[ssl]]` in the Hue configuration file, `hue.ini`.

Configuring Encryption

enabled	Choose to enable/disable TLS/SSL communication for this server. Default: false
cacerts	Path to Certificate Authority certificates. Default: /etc/hue/cacerts.pem
validate	Choose whether Hue should validate certificates received from the server. Default: true

Securing Database Connections using TLS/SSL

Connections vary depending on the database. Hue uses different clients to communicate with each database internally. Client specific options, such as secure connectivity, can be passed through the interface.

For example, for MySQL you can enable TLS/SSL communication by specifying the `options` configuration property under the `desktop>[[database]]` section in `hue.ini`. Here we identify the Certificate Authority (CA) certificate:

```
[desktop]
[[databases]]
...
options={"ssl":{"ca":"/tmp/ca-cert.pem"}}
```

You can also identify public and private keys, for example:

```
options={"ssl":{"ca":"/tmp/newcerts2/ca.pem", "key":"/tmp/newcerts2/client-key.pem",
"cert":"/tmp/newcerts2/client-cert.pem"}}
```

Storing Hue Passwords in a Script

In CDH 5.4, Hue added the ability to store passwords in a secure script and pull passwords from `stdout`. On startup, Hue runs one or more passwords scripts and grabs each password from `stdout`.

In `hue.ini`, add the suffix, `_script`, to any password property and set it equal to the script name. In Cludera Manager, set these properties in the configuration field, **Hue Service Advanced Configuration Snippet (Safety Valve) for hue_safety_valve.ini**. For example:

```
[desktop]
ldap_username=hueservice
ldap_password_script="/var/lib/hue/<your_hue_passwords_script.sh> ldap_password"
ssl_password_script="/var/lib/hue/<your_hue_passwords_script.sh> ssl_password"

[[ldap]]
bind_password_script="/var/lib/hue/<your_hue_passwords_script.sh> bind_password"

[[database]]
password_script="/var/lib/hue/<your_hue_passwords_script.sh> database"
```

Store the script in a directory that only the hue user can read, write, and execute. You can have one script per password or one script with parameters for all passwords. Here is an example of a script with parameters for multiple passwords:

```
#!/bin/bash

SERVICE=$1

if [[ ${SERVICE} == "ldap_password" ]]
then
    echo "password"
fi

if [[ ${SERVICE} == "ssl_password" ]]
then
    echo "password"
```



```

fi

if [[ ${SERVICE} == "bind_password" ]]
then
  echo "Password1"
fi

if [[ ${SERVICE} == "database_password" ]]
then
  echo "password"
fi

```



Note: The bind password parameter was added in CDH 5.4.6.

Configuring TLS/SSL for Impala

Impala supports TLS/SSL network encryption, between Impala and client programs, and between the Impala-related daemons running on different nodes in the cluster. This feature is important when you also use other features such as Kerberos authentication or Sentry authorization, where credentials are being transmitted back and forth.



Important:

- You can use either Cloudera Manager or the following command-line instructions to complete this configuration.
- This information applies specifically to CDH 5.6.x. If you use an earlier version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

Using Cloudera Manager

To configure Impala to listen for Beeswax and HiveServer2 requests on TLS/SSL-secured ports:

1. Open the Cloudera Manager Admin Console and go to the **Impala** service.
2. Click the **Configuration** tab.
3. Select **Scope > Impala (Service-Wide)**.
4. Select **Category > Security**.
5. Edit the following properties:

Table 15: Impala SSL Properties

Property	Description
Enable TLS/SSL for Impala Client Services	Encrypt communication between clients (like ODBC, JDBC, and the Impala shell) and the Impala daemon using Transport Layer Security (TLS) (formerly known as Secure Socket Layer (SSL)).
SSL/TLS Certificate for Clients	Local path to the X509 certificate that identifies the Impala daemon to clients during TLS/SSL connections. This file must be in PEM format.
SSL/TLS Private Key for Clients	Local path to the private key that matches the certificate specified in the Certificate for Clients. This file must be in PEM format.
SSL/TLS Private Key Password for Clients	A shell command for Impala to run on startup to retrieve the password for a password-protected private key file. The output of the command is truncated to a maximum of 1024 bytes, and any trailing whitespace (such as spaces or newline characters) is trimmed. If the command exits with an error, Impala does not start. If the password is incorrect, clients cannot connect to the server regardless of whether the public key is correct.

Property	Description
SSL/TLS CA Certificate	Must be specified for TLS/SSL encryption to be enabled for communication between internal Impala components.
SSL/TLS Certificate for <i>Impala component</i> Webserver	There are three of these configuration settings, one each for “Impala Daemon”, “Catalog Server”, and “Statestore”. Each of these Impala components has its own internal web server that powers the associated web UI with diagnostic information. The configuration setting represents the local path to the X509 certificate that identifies the web server to clients during TLS/SSL connections. This file must be in PEM format.

6. Click **Save Changes** to commit the changes.
7. Restart the Impala service.

For information on configuring TLS/SSL communication with the `impala-shell` interpreter, see [Configuring TLS/SSL Communication for the Impala Shell](#) on page 242.

Using the Command Line

To enable SSL for when client applications connect to Impala, add both of the following flags to the `impalad` startup options:

- `--ssl_server_certificate`: the full path to the server certificate, on the local filesystem.
- `--ssl_private_key`: the full path to the server private key, on the local filesystem.

In CDH 5.5 / Impala 2.3 and higher, Impala can also use SSL for its own internal communication between the `impalad`, `statedored`, and `catalogd` daemons. To enable this additional SSL encryption, set the `--ssl_server_certificate` and `--ssl_private_key` flags in the startup options for `impalad`, `catalogd`, and `statedored`, and also add the `--ssl_client_ca_certificate` flag for all three of those daemons.



Warning: Prior to CDH 5.5.2 / Impala 2.3.2, you could enable Kerberos authentication between Impala internal components, or SSL encryption between Impala internal components, but not both at the same time. This restriction has now been lifted. See [IMPALA-2598](#) to see the maintenance releases for different levels of CDH where the fix has been published.

If either of these flags are set, both must be set. In that case, Impala starts listening for Beeswax and HiveServer2 requests on SSL-secured ports only. (The port numbers stay the same; see [Ports Used by Impala](#) for details.)

Since Impala uses passphrase-less certificates in PEM format, you can reuse a host's existing Java keystore by converting it to the PEM format. For instructions, see [Conversion from Java Keystore to OpenSSL](#) on page 211.

Configuring TLS/SSL Communication for the Impala Shell

Typically, a client program has corresponding configuration properties in Cloudera Manager to verify that it is connecting to the right server. For example, with SSL enabled for Impala, you use the following options when starting the `impala-shell` interpreter:

- `--ssl`: enables TLS/SSL for `impala-shell`.
- `--ca_cert`: the local pathname pointing to the third-party CA certificate, or to a copy of the server certificate for self-signed server certificates.

If `--ca_cert` is not set, `impala-shell` enables TLS/SSL, but does not validate the server certificate. This is useful for connecting to a known-good Impala that is only running over TLS/SSL, when a copy of the certificate is not available (such as when debugging customer installations).

For `impala-shell` to successfully connect to an Impala cluster that has the minimum allowed TLS/SSL version set to 1.2 (`--ssl_minimum_version=tlsv1.2`), the Python version on the cluster that `impala-shell` runs on must be 2.7.9 or higher (or a vendor-provided Python version with the required support. Some vendors patched Python 2.7.5 versions on Red Hat Enterprise Linux 7 and derivatives).

Using TLS/SSL with Business Intelligence Tools

You can use Kerberos authentication, TLS/SSL encryption, or both to secure connections from JDBC and ODBC applications to Impala. See [Configuring Impala to Work with JDBC](#) and [Configuring Impala to Work with ODBC](#) for details.

Currently, the Hive JDBC driver does not support connections that use both Kerberos authentication and SSL encryption. To use both of these security features with Impala through a JDBC application, use the [Cloudera JDBC Connector](#) as the JDBC driver.

Configuring TLS/SSL for Oozie

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator, Full Administrator**)

Before You Begin

- Keystores for Oozie must be readable by the `oozie` user. This could be a copy of the Hadoop services' keystore with permissions `0440` and owned by the `oozie` group.
- Truststores must have permissions `0444` (that is, readable by all).
- Specify absolute paths to the keystore and truststore files. These settings apply to all hosts on which daemon roles of the Oozie service run. Therefore, the paths you choose must be valid on all hosts.
- In case there is a DataNode and an Oozie server running on the same host, they can use the same certificate.

For more information on obtaining signed certificates and creating keystores, see [TLS/SSL Certificates Overview](#) on page 205. You can also view the upstream documentation located [here](#).



Important:

- You can use either Cloudera Manager or the following command-line instructions to complete this configuration.
- This information applies specifically to CDH 5.6.x. If you use an earlier version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

Using Cloudera Manager

The steps for configuring and enabling Hadoop TLS/SSL for Oozie are as follows:

1. Open the Cloudera Manager Admin Console and go to the **Oozie service**.
2. Click the **Configuration** tab.
3. Select **Scope > All**.
4. Select **Category > All**.
5. In the Search field, type **TLS/SSL** to show the Oozie TLS/SSL properties.
6. Edit the following TLS/SSL properties according to your cluster configuration.

Table 16: Oozie TLS/SSL Properties

Property	Description
Enable TLS/SSL for Oozie	Check this field to enable TLS/SSL for Oozie.
Oozie TLS/SSL Server Keystore File Location	Location of the keystore file on the local file system.
Oozie TLS/SSL Server JKS Keystore File Password	Password for the keystore.

7. Click **Save Changes**.
8. Restart the Oozie service.

Configuring Encryption

Using the Command Line

To configure the Oozie server to use TLS/SSL:

1. Stop Oozie by running

```
sudo /sbin/service oozie stop
```

2. To enable TLS/SSL, set the MapReduce version that the Oozie server should work with using the `alternatives` command.



Note: The `alternatives` command is only available on RHEL systems. For SLES, Ubuntu and Debian systems, the command is `update-alternatives`.

For RHEL systems, to use YARN with TLS/SSL:

```
alternatives --set oozie-tomcat-conf /etc/oozie/tomcat-conf.https
```

For RHEL systems, to use MapReduce (MRv1) with TLS/SSL:

```
alternatives --set oozie-tomcat-conf /etc/oozie/tomcat-conf.https.mr1
```



Important:

The `OOZIE_HTTPS_KEYSTORE_PASS` variable must be the same as the password used when creating the keystore file. If you used a password other than `password`, you'll have to change the value of the `OOZIE_HTTPS_KEYSTORE_PASS` variable in this file.

3. Start Oozie by running

```
sudo /sbin/service oozie start
```

Connect to the Oozie Web UI using TLS/SSL (HTTPS)

Use `https://oozie.server.hostname:11443/oozie` though most browsers should automatically redirect you if you use `http://oozie.server.hostname:11000/oozie`.

Additional Considerations when Configuring TLS/SSL for Oozie HA

To allow clients to talk to Oozie servers (the target servers) through the load balancer using TLS/SSL, configure the load balancer to perform TLS/SSL pass-through. This allows clients to use the certificate provided by the target servers (so the load balancer will not need one). Consult your load balancer's documentation on how to configure this. Make sure to point the load balancer at the `https://HOST:HTTPS_PORT` addresses for your target servers. Clients can then connect to the load balancer at `https://LOAD_BALANCER_HOST:PORT`.

Configuring TLS/SSL for Solr

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)

Before You Begin

- The Solr service must be running.
- Keystores for Solr must be readable by the `solr` user. This could be a copy of the Hadoop services' keystore with permissions `0440` and owned by the `solr` group.
- Truststores must have permissions `0444` (that is, readable by all).
- Specify absolute paths to the keystore and truststore files. These settings apply to all hosts on which daemon roles of the Solr service run. Therefore, the paths you choose must be valid on all hosts.

- In case there is a DataNode and a Solr server running on the same host, they can use the same certificate.

For more information on obtaining signed certificates and creating keystores, see [TLS/SSL Certificates Overview](#) on page 205. You can also see the *Enabling SSL* section in the [Apache Solr 4.10 Reference Guide \(PDF\)](#).



Important:

- You can use either Cloudera Manager or the following command-line instructions to complete this configuration.
- This information applies specifically to CDH 5.6.x. If you use an earlier version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

Additional Considerations when Configuring TLS/SSL for Solr HA

To allow clients to talk to Solr servers (the target servers) through the load balancer using TLS/SSL, configure the load balancer to perform TLS/SSL pass-through. This allows clients to use the certificate provided by the target servers (so the load balancer will not need one). Consult your load balancer's documentation on how to configure this. Make sure to point the load balancer at the `https://HOST:HTTPS_PORT` addresses for your target servers. Clients can then connect to the load balancer at `https://LOAD_BALANCER_HOST:PORT`.

Configuring TLS/SSL for Solr Using Cloudera Manager

The steps for configuring and enabling Hadoop TLS/SSL for Search are as follows:

1. Open the Cloudera Manager Admin Console and go to the **Solr service**.
2. Click the **Configuration** tab.
3. Select **Scope > All**.
4. Select **Category > All**.
5. In the Search field, type **TLS/SSL** to show the Solr TLS/SSL properties.
6. Edit the following properties according to your cluster configuration.



Note: These values must be the same for all hosts running the Solr role.

Table 17: Solr TLS/SSL Properties

Property	Description
Enable TLS/SSL for Solr	Check this field to enable SSL for Solr.
Solr TLS/SSL Server Keystore File Location	The path to the TLS/SSL keystore file containing the server certificate and private key used for TLS/SSL. Used when Solr is acting as a TLS/SSL server. The keystore must be in JKS format.
Solr TLS/SSL Server JKS Keystore File Password	Password for the Solr JKS keystore.
Solr TLS/SSL Certificate Trust Store File	Required in case of self-signed or internal CA signed certificates. The location on disk of the trust store, in .jks format, used to confirm the authenticity of TLS/SSL servers that Solr might connect to. This is used when Solr is the client in a TLS/SSL connection. This trust store must contain the certificate(s) used to sign the service(s) being connected to. If this parameter is not provided, the default list of well-known certificate authorities is used instead.
Solr TLS/SSL Certificate Trust Store Password	The password for the Solr TLS/SSL Certificate Trust Store File. Note that this password is not required to access the trust store: this field can be left blank. This password provides optional integrity checking of the file. The contents of trust stores are certificates, and certificates are public information.

Configuring Encryption

7. Click **Save Changes** to commit the changes.
8. Restart the service.

Additional Considerations When Using a Load Balancer TLS/SSL for Solr HA

To configure a load balancer:

1. Go to the Solr service.
2. Click the **Configuration** tab.
3. Select **Scope > Solr**.
4. Enter the hostname and port number of the load balancer in the **Solr Load Balancer** property in the format *hostname:port number*.



Note:

When you set this property, Cloudera Manager regenerates the keytabs for Solr roles. The principal in these keytabs contains the load balancer hostname.

If there are services that depends on this Solr service, such as Hue, those services use the load balancer to communicate with Solr.

5. Click **Save Changes** to commit the changes.
6. Restart Solr and any dependent services or restart the entire cluster for this configuration to take effect.

Configuring TLS/SSL for Solr Using the Command Line

To configure the Search to use TLS/SSL:

1. Use `solrctl` to modify the `urlScheme` setting to specify `https`. For example:

```
solrctl --zk myZKEnsemble:2181/solr cluster --set-property urlScheme https
```

2. Stop Solr by running

```
sudo service solr-server stop
```

3. Edit `/etc/default/solr` to include the following environment variable settings:

```
SOLR_SSL_ENABLED=true
SOLR_KEYSTORE_PATH=<absolute_path_to_keystore_file>
SOLR_KEYSTORE_PASSWORD=<keystore_password>

#Following required only in case of self-signed or internal CA signed certificates
SOLR_TRUSTSTORE_PATH=<absolute_path_to_truststore_file>
SOLR_TRUSTSTORE_PASSWORD=<truststore_password>
```

4. Start Solr by running

```
sudo service solr-server start
```

Configuring TLS/SSL for the Key-Value Store Indexer Using Cloudera Manager

The steps for configuring and enabling Hadoop TLS/SSL for the Keystore Indexer are as follows:

1. Open the Cloudera Manager Admin Console and go to the **Key-Value Store Indexer**.
2. Click the **Configuration** tab.
3. Select **Scope > All**.
4. Select **Category > All**.
5. In the Search field, type **TLS/SSL** to show the Solr TLS/SSL properties.
6. Edit the following TLS/SSL properties according to your cluster configuration.



Note: These values must be the same for all hosts running the Key-Value Store Indexer role.

Table 18: Key-Value Store TLS/SSL Properties

Property	Description
HBase Indexer TLS/SSL Certificate Trust Store File	The location on disk of the trust store, in .jks format, used to confirm the authenticity of TLS/SSL servers that HBase Indexer might connect to. This is used when HBase Indexer is the client in a TLS/SSL connection. This trust store must contain the certificate(s) used to sign the service(s) being connected to. If this parameter is not provided, the default list of well-known certificate authorities is used instead.
HBase Indexer TLS/SSL Certificate Trust Store Password (Optional)	The password for the HBase Indexer TLS/SSL Certificate Trust Store File. Note that this password is not required to access the trust store: this field can be left blank. This password provides optional integrity checking of the file. The contents of trust stores are certificates, and certificates are public information.

7. Restart the service.

Configuring TLS/SSL for the Key-Value Store Indexer Using the Command Line

For every host running Key-Value Store Indexer server, specify Solr Trust Store details using the `HBASE_INDEXER_OPTS` environmental variable using following Java system properties:

- `-Djavax.net.ssl.trustStore=<absolute_path_to_truststore_file>`
- `-Djavax.net.ssl.trustStorePassword=<truststore_password>` (Optional)

Restart the Key-Value Store Indexer servers to apply these changes.

Configuring TLS/SSL for Flume Using Cloudera Manager

The steps for configuring and enabling Hadoop TLS/SSL for Flume are as follows:

1. Open the Cloudera Manager Admin Console and go to **Flume**.
2. Click the **Configuration** tab.
3. Select **Scope > All**.
4. Select **Category > All**.
5. In the Search field, type **TLS/SSL** to show the properties.
6. Edit the following SSL properties according to your cluster configuration.



Note: These values must be the same for all hosts running the Flume role.

Table 19: Key-Value Store SSL Properties

Property	Description
Flume TLS/SSL Certificate Trust Store File	The location on disk of the trust store, in .jks format, used to confirm the authenticity of TLS/SSL servers that Flume might connect to. This is used when Flume is the client in a TLS/SSL connection. This trust store must contain the certificate(s) used to sign the service(s) being connected to. If this parameter is not provided, the default list of well-known certificate authorities is used instead.
Flume TLS/SSL Certificate Trust Store Password (Optional)	The password for the Flume TLS/SSL Certificate Trust Store File. Note that this password is not required to access the trust store: this field can be left blank.

Configuring Encryption

Property	Description
	This password provides optional integrity checking of the file. The contents of trust stores are certificates, and certificates are public information.

7. Click **Save Changes** to commit the changes.
8. Restart the service.

Configuring TLS/SSL for Flume Using the Command Line

For every host running Flume agent, specify Solr Trust Store details using the `FLUME_AGENT_JAVA_OPTS` environmental variable using following Java system properties:

- `-Djavax.net.ssl.trustStore=<absolute_path_to_truststore_file>`
- `-Djavax.net.ssl.trustStorePassword=<truststore_password>` (Optional)

Restart the Flume agents to apply these changes.

Configuring TLS/SSL for HttpFS



Important:

- You can use either Cloudera Manager or the following command-line instructions to complete this configuration.
- This information applies specifically to CDH 5.6.x. If you use an earlier version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

Using Cloudera Manager

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)

The steps for configuring and enabling TLS/SSL for HttpFS using Cloudera Manager are as follows:

1. Go to the **HDFS service**
2. Click the **Configuration** tab.
3. Select **Scope** > **HttpFS**.
4. **Select** > **Security**.
5. Edit the following TLS/SSL properties according to your cluster configuration:

Table 20: HttpFS TLS/SSL Properties

Property	Description
Use TLS/SSL	Use TLS/SSL for HttpFS.
HttpFS Keystore File	Location of the keystore file used by the HttpFS role for TLS/SSL. Default: <code>/var/run/hadoop-httpfs/.keystore</code> . Note that the default location for the keystore file is on non-persistent disk.
HttpFS Keystore Password	Password of the keystore used by the HttpFS role for TLS/SSL. If the keystore password has a percent sign, it must be escaped. For example, for a password that is <code>pass%word</code> , use <code>pass%%word</code> .
HttpFS TLS/SSL Certificate Trust Store File	The location on disk of the truststore, in <code>.jks</code> format, used to confirm the authenticity of TLS/SSL servers that HttpFS might connect to. This is used when HttpFS is the client in a TLS/SSL connection.

Property	Description
HttpFS TLS/SSL Certificate Trust Store Password	The password for the HttpFS TLS/SSL Certificate Trust Store File. This password is not required to access the trust store; this field can be left blank. If the truststore password has a percent sign, it must be escaped. For example, for a password that is <code>pass%word</code> , use <code>pass%%word</code> .

6. Click **Save Changes**.
7. Restart the HDFS service.

Connect to the HttpFS Web UI using TLS/SSL (HTTPS)

Use `https://<httpfs_server_hostname>:14000/webhdfs/v1/`, though most browsers should automatically redirect you if you use `http://<httpfs_server_hostname>:14000/webhdfs/v1/`

Using the Command Line

Configure the HttpFS Server to use TLS/SSL (HTTPS)

1. Stop HttpFS by running

```
sudo /sbin/service hadoop-httpfs stop
```

2. To enable TLS/SSL, change which configuration the HttpFS server should work with using the `alternatives` command.



Note: The `alternatives` command is only available on RHEL systems. For SLES, Ubuntu and Debian systems, the command is `update-alternatives`.

For RHEL systems, to use TLS/SSL:

```
alternatives --set hadoop-httpfs-tomcat-conf /etc/hadoop-httpfs/tomcat-conf.https
```



Important:

The `HTTPFS_TLS/SSL_KEYSTORE_PASS` variable must be the same as the password used when creating the keystore file. If you used a password other than `password`, you'll have to change the value of the `HTTPFS_TLS/SSL_KEYSTORE_PASS` variable in `/etc/hadoop-httpfs/conf/httpfs-env.sh`.

3. Start HttpFS by running

```
sudo /sbin/service hadoop-httpfs start
```

Connect to the HttpFS Web UI using TLS/SSL (HTTPS)

Use `https://<httpfs_server_hostname>:14000/webhdfs/v1/`, though most browsers should automatically redirect you if you use `http://<httpfs_server_hostname>:14000/webhdfs/v1/`



Important:

If using a Self-Signed Certificate, your browser will warn you that it cannot verify the certificate or something similar. You will probably have to add your certificate as an exception.

Encrypted Shuffle and Encrypted Web UIs

Once you've enabled Kerberos, which provides for strong authentication, you can optionally enable network encryption if you so desire. CDH 5 supports the Encrypted Shuffle and Encrypted Web UIs feature that allows encryption of the MapReduce shuffle and web server ports using HTTPS with optional client authentication (also known as bi-directional HTTPS, or HTTPS with client certificates). It includes:

- Hadoop configuration setting for toggling the shuffle between HTTP and HTTPS.
- Hadoop configuration setting for toggling the Web UIs to use either HTTP or HTTPS.
- Hadoop configuration settings for specifying the keystore and truststore properties (location, type, passwords) that are used by the shuffle service, web server UIs and the reducers tasks that fetch shuffle data.
- A way to re-read truststores across the cluster (when a node is added or removed).

CDH 5 supports Encrypted Shuffle for both MRv1 and MRv2 (YARN), with common configuration properties used for both versions. The only configuration difference is in the parameters used to enable the features:

- For **MRv1**, setting the `hadoop.ssl.enabled` parameter in the `core-site.xml` file enables both the Encrypted Shuffle and the Encrypted Web UIs. In other words, the encryption toggling is coupled for the two features.
- For **MRv2**, setting the `hadoop.ssl.enabled` parameter enables the Encrypted Web UI feature; setting the `mapreduce.shuffle.ssl.enabled` parameter in the `mapred-site.xml` file enables the Encrypted Shuffle feature.

All other configuration properties apply to both the Encrypted Shuffle and Encrypted Web UI functionality.

When the Encrypted Web UI feature is enabled, all Web UIs for Hadoop components are served over HTTPS. If you configure the systems to require client certificates, browsers must be configured with the appropriate client certificates in order to access the Web UIs.



Important:

When the Web UIs are served over HTTPS, you must specify `https://` as the protocol; there is no redirection from `http://`. If you attempt to access an HTTPS resource over HTTP, your browser will probably show an empty screen with no warning.

Most components that run on top of MapReduce automatically use Encrypted Shuffle when it is configured.

Configuring Encrypted Shuffle and Encrypted Web UIs

To configure Encrypted Shuffle and Encrypted Web UIs, set the appropriate property/value pairs in the following:

- `core-site.xml` enables these features and defines the implementation
- `mapred-site.xml` enables Encrypted Shuffle for MRv2
- `ssl-server.xml` stores keystone and truststore settings for the server
- `ssl-client.xml` stores keystone and truststore settings for the client

core-site.xml Properties

To configure encrypted shuffle, set the following properties in the `core-site.xml` files of all nodes in the cluster:

Property	Default Value	Explanation
<code>hadoop.ssl.enabled</code>	<code>false</code>	For MRv1, setting this value to <code>true</code> enables both the Encrypted Shuffle and the Encrypted Web UI features. For MRv2, this property only enables the Encrypted WebUI; Encrypted Shuffle is enabled with a property in the <code>mapred-site.xml</code> file as described below.

Property	Default Value	Explanation
<code>hadoop.ssl.require.client.cert</code>	<code>false</code>	When this property is set to <code>true</code> , client certificates are required for all shuffle operations and all browsers used to access Web UIs. Cloudera recommends that this be set to <code>false</code> . See Client Certificates on page 255.
<code>hadoop.ssl.hostname.verifier</code>	<code>DEFAULT</code>	The hostname verifier to provide for <code>HttpsURLConnections</code> . Valid values are: <code>DEFAULT</code> , <code>STRICT</code> , <code>STRICT_I6</code> , <code>DEFAULT_AND_LOCALHOST</code> and <code>ALLOW_ALL</code> .
<code>hadoop.ssl.keystores.factory.class</code>	<code>org.apache.hadoop.security.ssl.FileBasedKeyStoresFactory</code>	The <code>KeyStoresFactory</code> implementation to use.
<code>hadoop.ssl.server.conf</code>	<code>ssl-server.xml</code>	Resource file from which ssl server keystore information is extracted. This file is looked up in the <code>classpath</code> ; typically it should be in the <code>/etc/hadoop/conf/</code> directory.
<code>hadoop.ssl.client.conf</code>	<code>ssl-client.xml</code>	Resource file from which ssl server keystore information is extracted. This file is looked up in the <code>classpath</code> ; typically it should be in the <code>/etc/hadoop/conf/</code> directory.

**Note:**

All these properties should be marked as `final` in the cluster configuration files.

Example

```
<configuration>
  ...
  <property>
    <name>hadoop.ssl.require.client.cert</name>
    <value>false</value>
    <final>true</final>
  </property>

  <property>
    <name>hadoop.ssl.hostname.verifier</name>
    <value>DEFAULT</value>
    <final>true</final>
  </property>

  <property>
    <name>hadoop.ssl.keystores.factory.class</name>
    <value>org.apache.hadoop.security.ssl.FileBasedKeyStoresFactory</value>
    <final>true</final>
  </property>
```

Configuring Encryption

```
<property>
  <name>hadoop.ssl.server.conf</name>
  <value>ssl-server.xml</value>
  <final>true</final>
</property>

<property>
  <name>hadoop.ssl.client.conf</name>
  <value>ssl-client.xml</value>
  <final>true</final>
</property>

<property>
  <name>hadoop.ssl.enabled</name>
  <value>true</value>
</property>
...
</configuration>
```

The cluster should be configured to use the Linux Task Controller in MRv1 and Linux container executor in MRv2 to run job tasks so that they are prevented from reading the server keystore information and gaining access to the shuffle server certificates. Refer to [Information about Other Hadoop Security Programs](#) for more information.

mapred-site.xml Property (MRv2 only)

To enable Encrypted Shuffle for MRv2, set the following property in the `mapred-site.xml` file on every node in the cluster:

Property	Default Value	Explanation
<code>mapreduce.shuffle.ssl.enabled</code>	<code>false</code>	If this property is set to <code>true</code> , encrypted shuffle is enabled. If this property is not specified, it defaults to the value of <code>hadoop.ssl.enabled</code> . This value can be <code>false</code> when <code>hadoop.ssl.enabled</code> is <code>true</code> but cannot be <code>true</code> when <code>hadoop.ssl.enabled</code> is <code>false</code> .

This property should be marked as `final` in the cluster configuration files.

Example:

```
<configuration>
...
  <property>
    <name>mapreduce.shuffle.ssl.enabled</name>
    <value>true</value>
    <final>true</final>
  </property>
...
</configuration>
```

Keystore and Truststore Settings

`FileBasedKeyStoresFactory` is the only `KeyStoresFactory` that is currently implemented. It uses properties in the `ssl-server.xml` and `ssl-client.xml` files to configure the keystores and truststores.

ssl-server.xml (Shuffle server and Web UI) Configuration

Use the following settings to configure the keystores and truststores in the `ssl-server.xml` file.

**Note:**

The `ssl-server.xml` should be owned by the `hdfs` or `mapred` Hadoop system user, belong to the `hadoop` group, and it should have 440 permissions. Regular users should not belong to the `hadoop` group.

Property	Default Value	Description
<code>ssl.server.keystore.type</code>	<code>jks</code>	Keystore file type
<code>ssl.server.keystore.location</code>	<code>NONE</code>	Keystore file location. The <code>mapred</code> user should own this file and have exclusive read access to it.
<code>ssl.server.keystore.password</code>	<code>NONE</code>	Keystore file password
<code>ssl.server.keystore.keypassword</code>	<code>NONE</code>	Key password
<code>ssl.server.truststore.type</code>	<code>jks</code>	Truststore file type
<code>ssl.server.truststore.location</code>	<code>NONE</code>	Truststore file location. The <code>mapred</code> user should own this file and have exclusive read access to it.
<code>ssl.server.truststore.password</code>	<code>NONE</code>	Truststore file password
<code>ssl.server.truststore.reload.interval</code>	<code>10000</code>	Truststore reload interval, in milliseconds

Example

```
<configuration>
<!-- Server Certificate Store -->
<property>
  <name>ssl.server.keystore.type</name>
  <value>jks</value>
</property>
<property>
  <name>ssl.server.keystore.location</name>
  <value>${user.home}/keystores/server-keystore.jks</value>
</property>
<property>
  <name>ssl.server.keystore.password</name>
  <value>serverfoo</value>
</property>
<property>
  <name>ssl.server.keystore.keypassword</name>
  <value>serverfoo</value>
</property>

<!-- Server Trust Store -->
<property>
  <name>ssl.server.truststore.type</name>
  <value>jks</value>
</property>
<property>
  <name>ssl.server.truststore.location</name>
  <value>${user.home}/keystores/truststore.jks</value>
</property>
<property>
  <name>ssl.server.truststore.password</name>
  <value>clientserverbar</value>
</property>
<property>
  <name>ssl.server.truststore.reload.interval</name>
  <value>10000</value>
</property>
</configuration>
```

ssl-client.xml (Reducer/Fetcher) Configuration

Use the following settings to configure the keystores and truststores in the `ssl-client.xml` file. This file should be owned by the `mapred` user for MRv1 and by the `yarn` user for MRv2; the file permissions should be 444 (read access for all users).

Property	Default Value	Description
<code>ssl.client.keystore.type</code>	<code>jks</code>	Keystore file type
<code>ssl.client.keystore.location</code>	<code>NONE</code>	Keystore file location. The <code>mapred</code> user should own this file and it should have default permissions.
<code>ssl.client.keystore.password</code>	<code>NONE</code>	Keystore file password
<code>ssl.client.keystore.keypassword</code>	<code>NONE</code>	Key password
<code>ssl.client.truststore.type</code>	<code>jks</code>	Truststore file type
<code>ssl.client.truststore.location</code>	<code>NONE</code>	Truststore file location. The <code>mapred</code> user should own this file and it should have default permissions.
<code>ssl.client.truststore.password</code>	<code>NONE</code>	Truststore file password
<code>ssl.client.truststore.reload.interval</code>	<code>10000</code>	Truststore reload interval, in milliseconds

Example

```
<configuration>
  <!-- Client certificate Store -->
  <property>
    <name>ssl.client.keystore.type</name>
    <value>jks</value>
  </property>
  <property>
    <name>ssl.client.keystore.location</name>
    <value>${user.home}/keystores/client-keystore.jks</value>
  </property>
  <property>
    <name>ssl.client.keystore.password</name>
    <value>clientfoo</value>
  </property>
  <property>
    <name>ssl.client.keystore.keypassword</name>
    <value>clientfoo</value>
  </property>

  <!-- Client Trust Store -->
  <property>
    <name>ssl.client.truststore.type</name>
    <value>jks</value>
  </property>
  <property>
    <name>ssl.client.truststore.location</name>
    <value>${user.home}/keystores/truststore.jks</value>
  </property>
  <property>
    <name>ssl.client.truststore.password</name>
    <value>clientserverbar</value>
  </property>
  <property>
    <name>ssl.client.truststore.reload.interval</name>
    <value>10000</value>
  </property>
</configuration>
```

```
</property>
</configuration>
```

Activating Encrypted Shuffle

When you have made the above configuration changes, activate Encrypted Shuffle by re-starting all TaskTrackers in MRv1 and all NodeManagers in YARN.



Important:

Encrypted shuffle has a significant performance impact. You should benchmark this before implementing it in production. In many cases, one or more additional cores are needed to maintain performance.

Client Certificates

Client Certificates are supported but they do not guarantee that the client is a reducer task for the job. The Client Certificate keystore file that contains the private key must be readable by all users who submit jobs to the cluster, which means that a rogue job could read those keystore files and use the client certificates in them to establish a secure connection with a Shuffle server. The JobToken mechanism that the Hadoop environment provides is a better protector of the data; each job uses its own JobToken to retrieve only the shuffle data that belongs to it. Unless the rogue job has a proper JobToken, it cannot retrieve Shuffle data from the Shuffle server.



Important:

If your certificates are signed by a certificate authority (CA), you must include the complete chain of CA certificates in the keystore that has the server's key.

Reloading Truststores

By default, each truststore reloads its configuration every 10 seconds. If a new truststore file is copied over the old one, it is re-read, and its certificates replace the old ones. This mechanism is useful for adding or removing nodes from the cluster, or for adding or removing trusted clients. In these cases, the client, TaskTracker or NodeManager certificate is added to (or removed from) all the truststore files in the system, and the new configuration is picked up without requiring that the TaskTracker in MRv1 and NodeManager in YARN daemons are restarted.



Note:

The keystores are not automatically reloaded. To change a keystore for a TaskTracker in MRv1 or a NodeManager in YARN, you must restart the TaskTracker or NodeManager daemon.

The reload interval is controlled by the `ssl.client.truststore.reload.interval` and `ssl.server.truststore.reload.interval` configuration properties in the `ssl-client.xml` and `ssl-server.xml` files described above.

Debugging



Important:

Enable debugging only for troubleshooting, and then only for jobs running on small amounts of data. Debugging is very verbose and slows jobs down significantly.

Configuring Encryption

To enable TLS/SSL debugging in the reducers, set `-Djavax.net.debug=all` in the `mapred.reduce.child.java.opts` property; for example:

```
<configuration>
  ...
  <property>
    <name>mapred.reduce.child.java.opts</name>
    <value>-Xmx200m -Djavax.net.debug=all</value>
  </property>
  ...
</configuration>
```

You can do this on a per-job basis, or by means of a cluster-wide setting in `mapred-site.xml`.

To set this property in TaskTrackers for MRv1, set it in `hadoop-env.sh`:

```
HADOOP_TASKTRACKER_OPTS="-Djavax.net.debug=all $HADOOP_TASKTRACKER_OPTS"
```

To set this property in NodeManagers for YARN, set it in `hadoop-env.sh`:

```
YARN_OPTS="-Djavax.net.debug=all $YARN_OPTS"
```

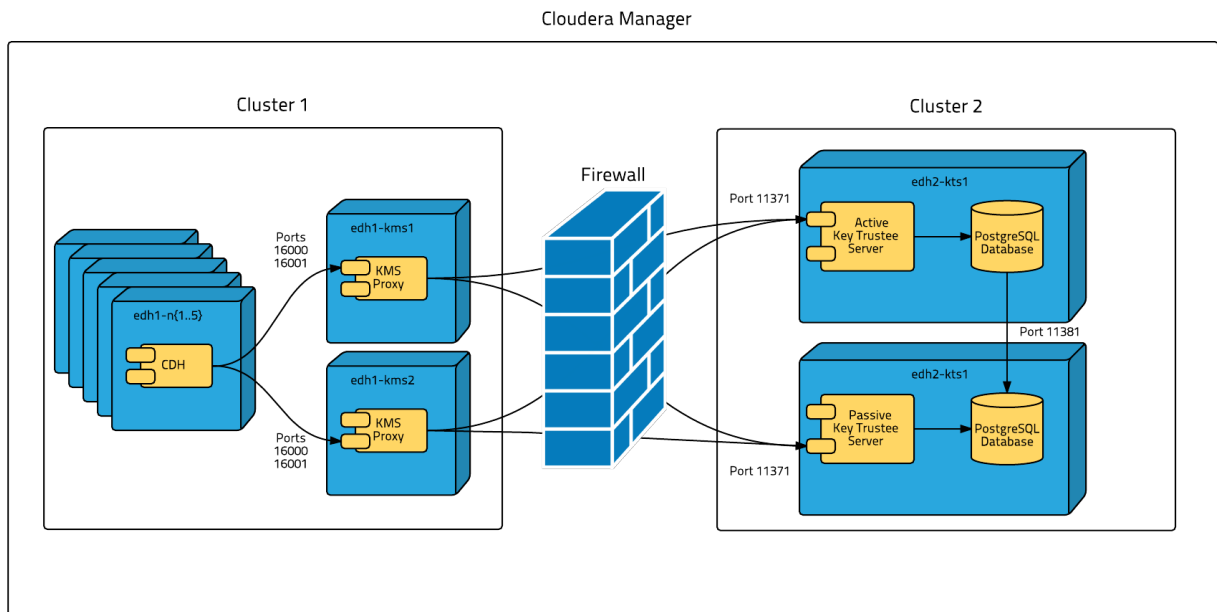
Deployment Planning for Data at Rest Encryption

Before deploying encryption for data at rest, familiarize yourself with the components, concepts, and architecture for encrypting data. See [Cloudera Navigator Data Encryption Overview](#) for more information.

For information on planning your encryption deployment, continue reading:

Data at Rest Encryption Reference Architecture

The following diagram illustrates the recommended architecture for deploying Cloudera Navigator encryption for data at rest:



To isolate Key Trustee Server from other enterprise data hub (EDH) services, deploy Key Trustee Server on dedicated hosts in a separate cluster in Cloudera Manager. Deploy Key Trustee KMS on dedicated hosts in the same cluster as the EDH services that require access to Key Trustee Server. This provides the following benefits:

- You can restart your EDH cluster without restarting Key Trustee Server, avoiding interruption to other clusters or clients that use the same Key Trustee Server instance.
- You can manage the Key Trustee Server upgrade cycle independently of other cluster components.
- You can limit access to the Key Trustee Server hosts to authorized key administrators only, reducing the attack surface of the system.
- Resource contention is reduced. Running Key Trustee Server and Key Trustee KMS services on dedicated hosts prevents other cluster services from reducing available resources (such as CPU and memory) and creating bottlenecks.

If you are using virtual machines for the Key Trustee Server or Key Trustee KMS hosts, see [Virtual Machine Considerations](#) on page 261.

Data at Rest Encryption Requirements

Encryption comprises several components, each with its own requirements.

Data at rest encryption protection can be applied at a number of levels within Hadoop:

- OS filesystem-level
- Network-level
- HDFS-level (protects both data at rest and in transit)

For more information on the components, concepts, and architecture for encrypting data at rest, see [Deployment Planning for Data at Rest Encryption](#) on page 256.

Product Compatibility Matrix

See [Product Compatibility Matrix for Cloudera Navigator Encryption](#) for the individual compatibility matrices for each Cloudera Navigator encryption component.

Entropy Requirements

Cryptographic operations require [entropy](#) to ensure randomness.

You can check the available entropy on a Linux system by running the following command:

```
$ cat /proc/sys/kernel/random/entropy_avail
```

The output displays the entropy currently available. Check the entropy several times to determine the state of the entropy pool on the system. If the entropy is consistently low (500 or less), you must increase it by installing `rng-tools` and starting the `rngd` service. Run the following commands on RHEL 6-compatible systems:

```
$ sudo yum install rng-tools
$ sudo echo 'EXTRAOPTIONS="-r /dev/urandom"' >> /etc/sysconfig/rngd
$ sudo service rngd start
$ sudo chkconfig rngd on
```

For RHEL 7, run the following commands:

```
$ sudo yum install rng-tools
$ cp /usr/lib/systemd/system/rngd.service /etc/systemd/system/
$ sed -i -e 's/ExecStart=/sbin\/rngd -f/ExecStart=/sbin\/rngd -f -r \/dev\/urandom/'
/etc/systemd/system/rngd.service
$ systemctl daemon-reload
$ systemctl start rngd
$ systemctl enable rngd
```

Make sure that the hosts running Key Trustee Server, Key Trustee KMS, and Navigator Encrypt have sufficient entropy to perform cryptographic operations.

Configuring Encryption

Key Trustee Server Requirements

Recommended Hardware and Supported Distributions

Cloudera recommends that the Key Trustee Server be installed on a dedicated server or virtual machine (VM) that is not used for any other purpose. The backing PostgreSQL database must be installed on the same host as the Key Trustee Server, and must not be shared with any other services. For high availability, the active and passive Key Trustee Servers must not share physical resources. See [Resource Planning for Data at Rest Encryption](#) on page 261 for more information.

The recommended minimum hardware specifications are as follows:

- Processor: 1 GHz 64-bit quad core
- Memory: 8 GB RAM
- Storage: 20 GB on moderate- to high-performance disk drives

For information on the supported Linux distributions, see [Table 2](#).

Cloudera Manager Requirements

Installing and managing Key Trustee Server using Cloudera Manager requires Cloudera Manager 5.4.0 and higher. Key Trustee Server does not require Cloudera Navigator Audit Server or Metadata Server.

umask Requirements

Key Trustee Server installation requires the default `umask` of `0022`.

Network Requirements

For new Key Trustee Server installations (5.4.0 and higher) and migrated upgrades (see [Migrate Apache Web Server to CherryPy](#) for more information), Key Trustee Server requires the following TCP ports to be opened for inbound traffic:

- 11371
Clients connect to this port over HTTPS.
- 11381 (PostgreSQL)
The passive Key Trustee Server connects to this port for database replication.

For upgrades that are not migrated to the CherryPy web server, the pre-upgrade port settings are preserved:

- 80
Clients connect to this port over HTTP to obtain the Key Trustee Server public key.
- 443 (HTTPS)
Clients connect to this port over HTTPS.
- 5432 (PostgreSQL)
The passive Key Trustee Server connects to this port for database replication.

TLS Certificate Requirements

To ensure secure network traffic, Cloudera recommends obtaining Transport Layer Security (TLS) certificates specific to the hostname of your Key Trustee Server. To obtain the certificate, generate a Certificate Signing Request (CSR) for the fully qualified domain name (FQDN) of the Key Trustee Server host. The CSR must be signed by a trusted Certificate Authority (CA). After the certificate has been verified and signed by the CA, the Key Trustee Server TLS configuration requires:

- The CA-signed certificate

- The private key used to generate the original CSR
- The intermediate certificate/chain file (provided by the CA)

Cloudera recommends not using self-signed certificates. If you use self-signed certificates, you must use the `--skip-ssl-check` parameter when registering Navigator Encrypt with the Key Trustee Server. This skips TLS hostname validation, which safeguards against certain network-level attacks. For more information regarding insecure mode, see [Table 24: Registration Options](#) on page 317.

Key Trustee KMS Requirements

Recommended Hardware and Supported Distributions

The recommended minimum hardware specifications are as follows:

- Processor: 2 GHz 64-bit quad core
- Memory: 16 GB RAM
- Storage: 40 GB on moderate- to high-performance disk drives

For information on the supported Linux distributions, see [Table 3](#).

The Key Trustee KMS workload is CPU-intensive. Cloudera recommends using machines with capabilities equivalent to your NameNode hosts, with Intel CPUs that support [AES-NI](#) for optimum performance.

Key HSM Requirements

The following are prerequisites for installing Navigator Key HSM:

- Oracle Java Runtime Environment (JRE) 7 or higher with Java Cryptography Extension (JCE) Unlimited Strength Jurisdiction Policy Files:
 - [JCE for Java SE 7](#)
 - [JCE for Java SE 8](#)
- A supported Linux distribution. See [Table 4](#).
- A supported HSM device:
 - SafeNet Luna
 - HSM firmware version: 6.2.1
 - HSM software version: 5.2.3-1
 - SafeNet KeySecure
 - HSM firmware version: 6.2.1
 - HSM software version: 8.0.1
 - Thales nSolo, nConnect
 - HSM firmware version: 11.4.0
 - Client software version: 2.28.9cam136
- Key Trustee Server 3.8 or higher



Important: You must install Key HSM on the same host as Key Trustee Server.

Root access is required to install Navigator Key HSM.

Configuring Encryption

Navigator Encrypt Requirements

Operating System Requirements

- Linux kernel 2.6.19 or higher (RHEL and CentOS can use 2.6.18-92 or higher)
- For supported Linux distributions, see [Table 5](#).



Note: Cloudera Enterprise, with the exception of Cloudera Navigator Encrypt, is supported on platforms with Security-Enhanced Linux (SELinux) enabled and in `enforcing` mode. Cloudera is not responsible for SELinux policy development, support, or enforcement. If you experience issues running Cloudera software with SELinux enabled, contact your OS provider for assistance.

If you are using SELinux in `enforcing` mode, Cloudera Support can request that you disable SELinux or change the mode to `permissive` to rule out SELinux as a factor when investigating reported issues.

Supported command-line interpreters:

- `sh` (Bourne)
- `bash` (Bash)
- `dash` (Debian)



Note: Navigator Encrypt does not support installation or use in `chroot` environments.

Network Requirements

For new Navigator Key Trustee Server (5.4.0 and higher) installations, Navigator Encrypt initiates TCP traffic over port 11371 (HTTPS) to the Key Trustee Server.

For upgrades and Key Trustee Server versions lower than 5.4.0, Navigator Encrypt initiates TCP traffic over ports 80 (HTTP) and 443 (HTTPS) to the Navigator Key Trustee Server.

Internet Access

You must have an active connection to the Internet to download many package dependencies, unless you have internal repositories or mirrors containing the dependent packages.

Maintenance Window

Data is not accessible during the encryption process. Plan for system downtime during installation and configuration.

Administrative Access

To enforce a high level of security, all Navigator Encrypt commands require administrative (`root`) access (including installation and configuration). If you do not have administrative privileges on your server, contact your system administrator before proceeding.

Package Dependencies

Navigator Encrypt requires these packages, which are resolved by your distribution package manager during installation:

- `dkms`
- `keyutils`
- `ecryptfs-utils`
- `libkeytrustee`
- `navencrypt-kernel-module`
- `openssl`

- `lsOf`
- `gcc`
- `cryptsetup`

These packages may have other dependencies that are also resolved by your package manager. Installation works with `gcc`, `gcc3`, and `gcc4`.

Resource Planning for Data at Rest Encryption

For production environments, you must configure high availability for Key Trustee Server and Key Trustee KMS.

For high availability, you must provision two dedicated Key Trustee Server hosts and at least two dedicated Key Trustee KMS hosts, for a minimum of four separate hosts. Do not run multiple Key Trustee Server or Key Trustee KMS services on the same physical host, and do not run these services on hosts with other cluster services. Doing so causes resource contention with other important cluster services and defeats the purpose of high availability. See [Data at Rest Encryption Reference Architecture](#) on page 256 for more information.

The Key Trustee KMS workload is CPU intensive. Cloudera recommends using machines with capabilities equivalent to your NameNode hosts, with Intel CPUs that support [AES-NI](#) for optimum performance.

Make sure that each host is secured and audited. Only authorized key administrators should have access to them. Red Hat provides security guides for RHEL:

- [RHEL 6 Security Guide](#)
- [RHEL 7 Security Guide](#)

For hardware sizing information, see [Data at Rest Encryption Requirements](#) on page 257 for recommendations for each Cloudera Navigator encryption component.

For Cloudera Manager deployments, deploy Key Trustee Server in its own dedicated cluster. Deploy Key Trustee KMS in each cluster that uses Key Trustee Server. See [Data at Rest Encryption Reference Architecture](#) on page 256 for more information.

Virtual Machine Considerations

If you are using virtual machines, make sure that the resources (such as virtual disks, CPU, and memory) for each Key Trustee Server and Key Trustee KMS host are allocated to separate physical hosts. Hosting multiple services on the same physical host defeats the purpose of high availability, because a single machine failure can take down multiple services.

To maintain the security of the cryptographic keys, make sure that all copies of the virtual disk (including any back-end storage arrays, backups, snapshots, and so on) are secured and audited with the same standards you apply to the live data.

HDFS Transparent Encryption

HDFS Encryption implements transparent, end-to-end encryption of data read from and written to HDFS, without requiring changes to application code. Because the encryption is end-to-end, data can be encrypted and decrypted only by the client. HDFS does not store or have access to unencrypted data or encryption keys. This supports both at-rest encryption (data on persistent media, such as a disk) and in-transit encryption (data traveling over a network).

Use Cases

Data encryption is required by a number of different government, financial, and regulatory entities. For example, the healthcare industry has HIPAA regulations, the card payment industry has PCI DSS regulations, and the United States government has FISMA regulations. Transparent encryption in HDFS makes it easier for organizations to comply with these regulations. Encryption can also be performed at the application-level, but by integrating it into HDFS, existing applications can operate on encrypted data without changes. This integrated architecture implements stronger encrypted file semantics and better coordination with other HDFS functions.

Architecture

Encryption Zones

An encryption zone is a directory in HDFS with all of its contents encrypted, that is, every file and subdirectory in it is encrypted. The files in this directory will be transparently encrypted upon write and transparently decrypted upon read. Each encryption zone is associated with a key, which is specified when the zone is created. Each file within an encryption zone also has its own encryption/decryption key, called the Data Encryption Key (DEK). These DEKs are never stored persistently unless they are encrypted with the encryption zone's key. This encrypted DEK is known as the EDEK. The EDEK is then stored persistently as part of the file's metadata on the NameNode.

A key can have multiple key versions, where each key version has its own distinct key material (that is, the portion of the key used during encryption and decryption). Key rotation is achieved by modifying the encryption zone's key, that is, bumping up its version. Per-file key rotation is then achieved by re-encrypting the file's DEK with the new encryption zone key to create new EDEKs. An encryption key can be fetched either by its key name, returning the latest version of the key, or by a specific key version.



Note: An encryption zone cannot be created on top of an existing directory. Each encryption zone begins as an empty directory and `distcp` can be used to add data to the zone.

Key Management Server

A new service, called the **Hadoop Key Management Server (KMS)**, needs to be added to your cluster to store, manage, and access encryption keys. The KMS service is a proxy that interfaces with a backing key store on behalf of HDFS daemons and clients. Both the backing key store and the KMS implement the Hadoop KeyProvider client API.

Encryption and decryption of EDEKs happens entirely on the KMS. More importantly, the client requesting creation or decryption of an EDEK never handles the EDEK's encryption key (that is, the encryption zone key). When a new file is created in an encryption zone, the NameNode asks the KMS to generate a new EDEK encrypted with the encryption zone's key. When reading a file from an encryption zone, the NameNode provides the client with the file's EDEK and the encryption zone key version that was used to encrypt the EDEK. The client then asks the KMS to decrypt the EDEK, which involves checking that the client has permission to access the encryption zone key version. Assuming that is successful, the client uses the DEK to decrypt the file's contents. All the steps for read and write take place automatically through interactions between the DFSClient, the NameNode, and the KMS.

Access to encrypted file data and metadata is controlled by normal HDFS filesystem permissions. Typically, the backing key store is configured to only allow end-user access to the encryption zone keys used to encrypt DEKs. This means that EDEKs can be safely stored and handled by HDFS, since the `hdfs` user will not have access to EDEK encryption keys. This means that if HDFS is compromised (for example, by gaining unauthorized access to a superuser account), a malicious user only gains access to the ciphertext and EDEKs. This does not pose a security threat since access to encryption zone keys is controlled by a separate set of permissions on the KMS and key store.

For more information about configuring the KMS, see [Configuring the Key Management Server \(KMS\)](#) on page 274.

Navigator Key Trustee

HDFS encryption can use a local Java KeyStore for key management. This is not sufficient for production environments where a more robust and secure key management solution is required. [Cloudera Navigator Key Trustee Server](#) on page 297 is a key store for managing encryption keys and other secure deposits.

In order to leverage the manageable, highly-available key management capabilities of the Navigator Key Trustee Server, Cloudera provides a custom KMS service, the Key Trustee KMS.

For more information, see [Enabling HDFS Encryption Using the Wizard](#) on page 266.

DistCp Considerations

A common use case for DistCp is to replicate data between clusters for backup and disaster recovery purposes. This is typically performed by the cluster administrator, who is an HDFS superuser. To retain this workflow when using HDFS encryption, a new virtual path prefix has been introduced, `/.reserved/raw/`, that gives superusers direct access to

the underlying block data in the filesystem. This allows superusers to `distcp` data without requiring access to encryption keys, and avoids the overhead of decrypting and re-encrypting data. It also means the source and destination data will be byte-for-byte identical, which would not have been true if the data was being re-encrypted with a new EDEK.



Warning:

When using `/.reserved/raw/` to `distcp` encrypted data, make sure you preserve extended attributes with the `-px` flag. This is because encrypted attributes such as the EDEK are exposed through extended attributes and *must* be preserved to be able to decrypt the file.

This means that if the `distcp` is initiated at or above the encryption zone root, it will automatically create a new encryption zone at the destination if it does not already exist. Hence, Cloudera recommends you first create identical encryption zones on the destination cluster to avoid any potential mishaps.

Copying between encrypted and unencrypted locations

By default, `distcp` compares checksums provided by the filesystem to verify that data was successfully copied to the destination. When copying between an unencrypted and encrypted location, the filesystem checksums will not match since the underlying block data is different.

In this case, you can specify the `-skipcrccheck` and `-update` flags to avoid verifying checksums.

Attack Vectors

Type of Exploit	Issue	Mitigation
Hardware Access Exploit		
These exploits assume the attacker has gained physical access to hard drives from cluster machines, that is, DataNodes and NameNodes.	<i>Access to swap files of processes containing DEKs.</i> This exploit does not expose cleartext, as it also requires access to encrypted block files.	It can be mitigated by disabling swap, using encrypted swap, or using <code>mlock</code> to prevent keys from being swapped out.
	<i>Access to encrypted block files.</i> This exploit does not expose cleartext, as it also requires access to the DEKs.	It can only be mitigated by restricting physical access to the cluster machines.
Root Access Exploits		
These exploits assume the attacker has gained root shell access to cluster machines running DataNodes and NameNodes. Many of these exploits cannot be addressed in HDFS, since a malicious root user has access to the in-memory state of processes holding encryption keys and cleartext. For these exploits, the only mitigation technique is carefully restricting and monitoring root shell access.	Access to encrypted block files. By itself, this does not expose cleartext, as it also requires access to encryption keys.	No mitigation required.
	Dump memory of client processes to obtain DEKs, delegation tokens, cleartext.	No mitigation.
	Recording network traffic to sniff encryption keys and encrypted data in transit. By itself, insufficient to read cleartext without the EDEK encryption key.	No mitigation required.
	Dump memory of DataNode process to obtain encrypted block data.	No mitigation required.

Type of Exploit	Issue	Mitigation
	By itself, insufficient to read cleartext without the DEK.	
	Dump memory of NameNode process to obtain encrypted data encryption keys.	No mitigation required.
	By itself, insufficient to read cleartext without the EDEK's encryption key and encrypted block files.	
HDFS Admin Exploits		
These exploits assume that the attacker has compromised HDFS, but does not have root or hdfs user shell access.	Access to encrypted block files.	No mitigation required.
	By itself, insufficient to read cleartext without the EDEK and EDEK encryption key.	
	Access to encryption zone and encrypted file metadata (including encrypted data encryption keys), using <code>-fetchImage</code> .	No mitigation required.
	By itself, insufficient to read cleartext without EDEK encryption keys.	
Root Access Exploits		
	A rogue user can collect keys to which they have access, and use them later to decrypt encrypted data.	This can be mitigated through periodic key rolling policies.

Optimizing for HDFS Data at Rest Encryption



Warning: To ensure that HDFS encryption functions as expected, the steps described in this section are *mandatory for production use*.

CDH implements the **Advanced Encryption Standard New Instructions** (AES-NI), which provide substantial performance improvements. To get these improvements, you need a recent version of `libcrypto.so` on HDFS and MapReduce client hosts -- that is, any host from which you originate HDFS or MapReduce requests. Many OS versions have an older version of the library that does not support AES-NI. The instructions that follow tell you what you need to do for each OS version that CDH supports.

RHEL/CentOS 6.5 or later

The installed version of `libcrypto.so` supports AES-NI, but you need to install the `openssl-devel` package on all clients:

```
$ sudo yum install openssl-devel
```

RHEL/CentOS 6.4 or earlier 6.x versions, or SLES 11

Download and extract a newer version of `libcrypto.so` from a CentOS 6.5 repository and install it on all clients in `/var/lib/hadoop/extra/native/`:

1. Download the latest version of the `openssl` package. For example:

```
$ wget
http://mirror.centos.org/centos/6/os/x86_64/Packages/openssl-1.0.1e-30.el6.x86_64.rpm
```

The `libcrypto.so` file in this package can be used on SLES 11 as well as RHEL/CentOS.

2. Decompress the files in the package, but **do not** install it:

```
$ rpm2cpio openssl-1.0.1e-30.el6.x86_64.rpm | cpio -idmv
```

3. If you are using parcels, create the `/var/lib/hadoop/extra/native/` directory:

```
$ sudo mkdir -p /var/lib/hadoop/extra/native
```

4. Copy the shared library into `/var/lib/hadoop/extra/native/`. Name the target file `libcrypto.so`, with no suffix at the end, exactly as in the command that follows.

```
$ sudo cp ./usr/lib64/libcrypto.so.1.0.1e /var/lib/hadoop/extra/native/libcrypto.so
```

RHEL/CentOS 5

In this case, you need to build `libcrypto.so` and copy it to all clients:

1. On one client, compile and install `openssl` from source:

```
$ wget http://www.openssl.org/source/openssl-1.0.1j.tar.gz
$ cd openssl-1.0.1j
$ ./config --shared --prefix=/opt/openssl-1.0.1j
$ sudo make install
```

2. If you are using parcels, create the `/var/lib/hadoop/extra/native/` directory:

```
$ sudo mkdir -p /var/lib/hadoop/extra/native
```

3. Copy the files into `/var/lib/hadoop/extra/native/`:

```
$ sudo cp /opt/openssl-1.0.1j/lib/libcrypto.so /var/lib/hadoop/extra/native
```

4. Copy the files to the remaining clients using a utility such as `rsync`

Debian Wheezy

The installed version of `libcrypto.so` supports AES-NI, but you need to install the `libssl-devel` package on all clients:

```
$ sudo apt-get install libssl-dev
```

Ubuntu Precise and Ubuntu Trusty

Install the `libssl-devel` package on all clients:

```
$ sudo apt-get install libssl-dev
```

Testing if encryption optimization works

To verify that a client host is ready to use the AES-NI instruction set optimization for HDFS encryption at rest, use the following command:

```
hadoop checknative
```

You should see a response such as the following:

```
14/12/12 13:48:39 INFO bzip2.Bzip2Factory: Successfully loaded & initialized native-bzip2
library system-native14/12/12 13:48:39 INFO zlib.ZlibFactory: Successfully loaded &
initialized native-zlib library
Native library checking:
hadoop: true /usr/lib/hadoop/lib/native/libhadoop.so.1.0.0
zlib: true /lib64/libz.so.1
snappy: true /usr/lib64/libsnappy.so.1
lz4: true revision:99
bzip2: true /lib64/libbz2.so.1
openssl: true /usr/lib64/libcrypto.so
```

If you see `true` in the `openssl` row, Hadoop has detected the right version of `libcrypto.so` and optimization will work. If you see `false` in this row, you do not have the right version.

Enabling HDFS Encryption Using the Wizard

To accommodate the security best practice of [separation of duties](#), enabling HDFS encryption using the wizard requires different Cloudera Manager user roles for different steps.

Launch the **Set up HDFS Data At Rest Encryption** wizard in one of the following ways:

- **Cluster** ▾ > **Set up HDFS Data At Rest Encryption**
Minimum Required Role: [Key Administrator](#) or [Cluster Administrator](#) (also provided by Full Administrator)
- **Administration** > **Security** > **Set up HDFS Data At Rest Encryption**
Minimum Required Role: [Key Administrator](#) or [Cluster Administrator](#) (also provided by Full Administrator)
- **HDFS service** > **Actions** > **Set up HDFS Data At Rest Encryption**
Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

On the first page of the wizard, select the root of trust for encryption keys:

- **Cloudera Navigator Key Trustee Server**
- **A file-based password-protected Java KeyStore**

Cloudera strongly recommends using Cloudera Navigator Key Trustee Server as the root of trust for production environments. The file-based Java KeyStore root of trust is insufficient to provide the security, scalability, and manageability required by most production systems. More specifically, the Java KeyStore KMS does not provide:

- Scalability, so you are limited to only one KMS, which can result in bottlenecks
- High Availability (HA)
- Recoverability, so if you lose the node where the Java KeyStore is stored, then you can lose access to all the encrypted data

Ultimately, the Java KeyStore does not satisfy the stringent security requirements of most organizations for handling master encryption keys.

Choosing a root of trust displays a list of steps required to enable HDFS encryption using that root of trust. Each step can be completed independently. The **Status** column indicates whether the step has been completed, and the **Notes** column provides additional context for the step. If your Cloudera Manager user account does not have sufficient privileges to complete a step, the **Notes** column indicates the required privileges.

Available steps contain links to wizards or documentation required to complete the step. If a step is unavailable due to insufficient privileges or a prerequisite step being incomplete, no links are present and the **Notes** column indicates the reason the step is unavailable.

Continue to the section for your selected root of trust for further instructions:

Enabling HDFS Encryption Using Cloudera Navigator Key Trustee Server

Enabling HDFS encryption using Key Trustee Server as the key store involves multiple components. For an overview of the components involved in encrypting data at rest, see [Cloudera Navigator Data Encryption Overview](#).

Before continuing, make sure the Cloudera Manager server host has access to the internal repository hosting the Key Trustee Server software. See [Setting Up an Internal Repository](#) for more information.

After selecting **Cloudera Navigator Key Trustee Server** as the root of trust, the following steps are displayed:

1. Enable Kerberos

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

For more information about enabling Kerberos, see [Enabling Kerberos Authentication Using the Wizard](#) on page 59.

2. Enable TLS/SSL

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

For more information about enabling TLS, see [Configuring TLS/SSL Encryption for CDH Services](#) on page 228.

3. Add a dedicated cluster for the Key Trustee Server

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

This step creates a new cluster in Cloudera Manager for the Key Trustee Server hosts to isolate them from other enterprise data hub (EDH) services for increased security and durability. For more information, see [Resource Planning for Data at Rest Encryption](#) on page 261 and [Data at Rest Encryption Reference Architecture](#) on page 256.

To complete this step:

1. Click **Add a dedicated cluster for the Key Trustee Server**.
2. Leave **Enable High Availability** checked to add two hosts to the cluster. Cloudera strongly recommends using high availability for Key Trustee Server. Failure to enable high availability can result in complete data loss in the case of catastrophic failure of a standalone Key Trustee Server. Click **Continue**.
3. Search for new hosts to add to the cluster, or select the **Currently Managed Hosts** tab to add existing hosts to the cluster. After selecting the hosts, click **Continue**.
4. Select the KEYTRUSTEE_SERVER parcel to install Key Trustee Server using parcels, or select **None** if you want to use packages. If you do not see a parcel available, click **More Options** and add the repository URL to the **Remote Parcel Repository URLs** list. After selecting a parcel or **None**, click **Continue**.

If you selected **None**, click **Continue** again, and skip to [4. Install Key Trustee Server binary using packages or parcels](#) on page 267.

5. After the KEYTRUSTEE_SERVER parcel is successfully downloaded, distributed, unpacked, and activated, click **Continue**.
6. Click **Continue** to complete this step and return to the main page of the wizard.

4. Install Key Trustee Server binary using packages or parcels

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)



Note: If you selected **None** on the parcel selection page in step [3. Add a dedicated cluster for the Key Trustee Server](#) on page 267, the step title is changed to **Install Parcel for Key Trustee Server**. If you are using packages, skip this step and see [Installing Key Trustee Server Using the Command Line](#) for package-based installation instructions. After installing Key Trustee Server using packages, continue to [5. Install Parcel for Key Trustee KMS](#) on page 268.

This step is completed automatically during [3. Add a dedicated cluster for the Key Trustee Server](#) on page 267 if you are using parcels. If the step is incomplete for any reason (such as the wizard being interrupted or a failure installing the parcel), complete it manually:

1. Click **Install Key Trustee Server binary using packages or parcels**.
2. Select the `KEYTRUSTEE_SERVER` parcel to install Key Trustee Server, or select **None** if you need to install Key Trustee Server manually using packages. If you do not see a parcel available, click **More Options** and add the repository URL to the **Remote Parcel Repository URLs** list. After selecting a parcel, click **Continue**.
3. After the `KEYTRUSTEE_SERVER` parcel is successfully downloaded, distributed, unpacked, and activated, click **Finish** to complete this step and return to the main page of the wizard.

5. Install Parcel for Key Trustee KMS

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

This step installs the Key Trustee KMS parcel. If you are using packages, skip this step and see [Installing Key Trustee KMS Using Packages](#) for instructions. After installing Key Trustee KMS using packages, continue to [6. Add a Key Trustee Server Service](#) on page 268.

To complete this step for parcel-based installations:

1. Click **Install Parcel for Key Trustee KMS**.
2. Select the `KEYTRUSTEE` parcel to install Key Trustee KMS. If you do not see a parcel available, click **More Options** and add the repository URL to the **Remote Parcel Repository URLs** list. After selecting a parcel, click **Continue**.
3. After the `KEYTRUSTEE` parcel is successfully downloaded, distributed, unpacked, and activated, click **Finish** to complete this step and return to the main page of the wizard.

6. Add a Key Trustee Server Service

Minimum Required Role: [Key Administrator](#) (also provided by **Full Administrator**)

This step adds the **Key Trustee Server** service to Cloudera Manager. To complete this step:

1. Click **Add a Key Trustee Server Service**.
2. Click **Continue**.
3. On the **Customize Role Assignments for Key Trustee Server** page, select the hosts for the **Active Key Trustee Server** and **Passive Key Trustee Server** roles. Make sure that the selected hosts are not used for other services (see [Resource Planning for Data at Rest Encryption](#) on page 261 for more information), and click **Continue**.
4. The **Entropy Considerations** page provides commands to install the `rng-tools` package to increase available entropy for cryptographic operations. For more information, see [Entropy Requirements](#) on page 257. After completing these commands, click **Continue**.
5. The **Synchronize Active and Passive Key Trustee Server Private Keys** page provides instructions for generating and copying the Active Key Trustee Server private key to the Passive Key Trustee Server. Cloudera recommends following security best practices and transferring the private key using offline media, such as a removable USB drive. For convenience (for example, in a development or testing environment where maximum security is not required), you can copy the private key over the network using the provided `rsync` command.

After you have synchronized the private keys, run the `ktadmin init` command on the Passive Key Trustee Server as described in the wizard. After the initialization is complete, check the box to indicate you have synchronized the keys and click **Continue** in the wizard.

6. The **Setup TLS for Key Trustee Server** page provides instructions on replacing the auto-generated self-signed certificate with a production certificate from a trusted Certificate Authority (CA). For more information, see [Managing Key Trustee Server Certificates](#) on page 307. Click **Continue** to view and modify the default certificate settings.
7. On the **Review Changes** page, you can view and modify the following settings:

- **Database Storage Directory** (`db_root`)
Default value: `/var/lib/keytrustee/db`

The directory on the local filesystem where the Key Trustee Server database is stored. Modify this value to store the database in a different directory.

- **Active Key Trustee Server TLS/SSL Server Private Key File (PEM Format)** (`ssl.privatekey.location`)

Default value: `/var/lib/keytrustee/.keytrustee/.ssl/ssl-cert-keytrustee-pk.pem`

The path to the Active Key Trustee Server TLS certificate private key. Accept the default setting to use the auto-generated private key. If you have a CA-signed certificate, change this path to the CA-signed certificate private key file. This file must be in [PEM](#) format.
- **Active Key Trustee Server TLS/SSL Server Certificate File (PEM Format)** (`ssl.cert.location`)

Default value: `/var/lib/keytrustee/.keytrustee/.ssl/ssl-cert-keytrustee.pem`

The path to the Active Key Trustee Server TLS certificate. Accept the default setting to use the auto-generated self-signed certificate. If you have a CA-signed certificate, change this to the path to the CA-signed certificate. This file must be in PEM format.
- **Active Key Trustee Server TLS/SSL Server CA Certificate (PEM Format)** (`ssl.cacert.location`)

Default value: (none)

The path to the file containing the CA certificate and any intermediate certificates (if any intermediate certificates exist, then they are *required* here) used to sign the Active Key Trustee Server certificate. If you have a CA-signed certificate, set this value to the path to the CA certificate or certificate chain file. This file must be in PEM format.
- **Active Key Trustee Server TLS/SSL Private Key Password** (`ssl.privatekey.password`)

Default value: (none)

The password for the Active Key Trustee Server private key file. Leave this blank if the file is not password-protected.
- **Passive Key Trustee Server TLS/SSL Server Private Key File (PEM Format)** (`ssl.privatekey.location`)

Default value: `/var/lib/keytrustee/.keytrustee/.ssl/ssl-cert-keytrustee-pk.pem`

The path to the Passive Key Trustee Server TLS certificate private key. Accept the default setting to use the auto-generated private key. If you have a CA-signed certificate, change this path to the CA-signed certificate private key file. This file must be in [PEM](#) format.
- **Passive Key Trustee Server TLS/SSL Server Certificate File (PEM Format)** (`ssl.cert.location`)

Default value: `/var/lib/keytrustee/.keytrustee/.ssl/ssl-cert-keytrustee.pem`

The path to the Passive Key Trustee Server TLS certificate. Accept the default setting to use the auto-generated self-signed certificate. If you have a CA-signed certificate, change this to the path to the CA-signed certificate. This file must be in PEM format.
- **Passive Key Trustee Server TLS/SSL Server CA Certificate (PEM Format)** (`ssl.cacert.location`)

Default value: (none)

The path to the file containing the CA certificate and any intermediate certificates (if any intermediate certificates exist, then they are *required* here) used to sign the Passive Key Trustee Server certificate. If you have a CA-signed certificate, set this value to the path to the CA certificate or certificate chain file. This file must be in PEM format.
- **Passive Key Trustee Server TLS/SSL Private Key Password** (`ssl.privatekey.password`)

Default value: (none)

The password for the Passive Key Trustee Server private key file. Leave this blank if the file is not password-protected.

After reviewing the settings and making any changes, click **Continue**.

8. After all commands complete successfully, click **Continue**. If the **Generate Key Trustee Server Keyring** appears stuck, make sure that the Key Trustee Server host has enough entropy. See [Entropy Requirements](#) on page 257 for more information.
9. Click **Finish** to complete this step and return to the main page of the wizard.

7. Add a Key Trustee KMS Service

Minimum Required Role: [Key Administrator](#) (also provided by **Full Administrator**)

This step adds a Key Trustee KMS service to the cluster. The Key Trustee KMS service is required to enable HDFS encryption to use Key Trustee Server for cryptographic key management. To complete this step:

1. Click **Add a Key Trustee KMS Service**.
2. Select an existing Key Trustee Server pair or specify an external Key Trustee Server pair. If you have an existing Key Trustee Server pair outside of Cloudera Manager control, select the **External Key Trustee Server** option and specify the fully-qualified domain names (FQDNs) of the Key Trustee Server pair. Click **Continue**.
3. Select cluster hosts for the Key Trustee KMS service. Cloudera recommends selecting at least two hosts for high availability. If you proceed with only one host, you can enable high availability later. See [Key Trustee KMS High Availability](#) for more information.

Make sure that the selected hosts are not used for other services (see [Resource Planning for Data at Rest Encryption](#) on page 261 for more information), and click **Continue**.

4. The **Entropy Considerations** page provides commands to install the `rng-tools` package to increase available entropy for cryptographic operations. For more information, see [Entropy Requirements](#) on page 257. After completing these commands, click **Continue**.
5. The **Setup Organization and Auth Secret** page generates the necessary commands to create an organization in Key Trustee Server. An organization is required to be able to register the Key Trustee KMS with Key Trustee Server. See [Managing Key Trustee Server Organizations](#) on page 304 for more information.

Enter an organization name and click **Generate Instruction**. Run the displayed commands to generate an organization and obtain the `auth_secret` value for the organization. Enter the secret in the `auth_secret` field and click **Continue**.

6. The **Setup Access Control List (ACL)** page allows you to generate ACLs for the Key Trustee KMS or to provide your own ACLs. To generate the recommended ACLs, enter the username and group responsible for managing cryptographic keys and click **Generate ACLs**. To specify your own ACLs, select the **Use Your Own kms-acls.xml File** option and enter the ACLs. For more information on the KMS Access Control List, see [Configuring KMS Access Control Lists](#) on page 281.

After generating or specifying the ACL, click **Continue**.

7. The **Setup TLS for Key Trustee KMS** page provides high-level instructions for configuring TLS communication between the Key Trustee KMS and the Key Trustee Server, as well as between the EDH cluster and the Key Trustee KMS. See [Configuring TLS/SSL for the KMS](#) on page 280 for more information.

Click **Continue**.

8. The **Review Changes** page lists all of the settings configured in this step. Click the ⓘ icon next to any setting for information about that setting. Review the settings and click **Continue**.
9. After the **First Run** commands have successfully completed, click **Continue**.
10. The **Synchronize Private Keys and HDFS Dependency** page provides instructions for copying the private key from one Key Management Server Proxy role to all other roles.



Warning: It is *very important* that you perform this step. Failure to do so leaves Key Trustee KMS in a state where keys are intermittently inaccessible, depending on which Key Trustee KMS host a client interacts with, because cryptographic key material encrypted by one Key Trustee KMS host cannot be decrypted by another. If you are already running multiple Key Trustee KMS hosts with different private keys, immediately [back up](#) all Key Trustee KMS hosts, and contact Cloudera Support for assistance correcting the issue.

To determine whether the Key Trustee KMS private keys are different, compare the MD5 hash of the private keys. On each Key Trustee KMS host, run the following command:

```
$ md5sum /var/lib/kms-keytrustee/keytrustee/.keytrustee/secring.gpg
```

If the outputs are different, contact Cloudera Support for assistance. Do not attempt to synchronize existing keys. If you overwrite the private key and do not have a backup, any keys encrypted by that private key are permanently inaccessible, and any data encrypted by those keys is permanently irretrievable. If you are configuring Key Trustee KMS high availability for the first time, continue synchronizing the private keys.

Cloudera recommends following security best practices and transferring the private key using offline media, such as a removable USB drive. For convenience (for example, in a development or testing environment where maximum security is not required), you can copy the private key over the network using the provided `rsync` command.

After you have synchronized the private keys, check the box to indicate you have done so and click **Continue**.

1. After the Key Trustee KMS service starts, click **Finish** to complete this step and return to the main page of the wizard.

8. Restart stale services and redeploy client configuration

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

This step restarts all services which were modified while enabling HDFS encryption. To complete this step:

1. Click **Restart stale services and redeploy client configuration**.
2. Click **Restart Stale Services**.
3. Make sure that **Re-deploy client configuration** is checked, and click **Restart Now**.
4. After all commands have completed, click **Finish**.

9. Validate Data Encryption

Minimum Required Role: [Key Administrator](#) or [Cluster Administrator](#) (also provided by Full Administrator)

This step launches a tutorial with instructions on creating an encryption zone and putting data into it to verify that HDFS encryption is enabled and working.

Enabling HDFS Encryption Using a Java KeyStore



Note: Cloudera strongly recommends using Cloudera Navigator Key Trustee Server as the root of trust for production environments. The file-based Java KeyStore root of trust is insufficient to provide the security, scalability, and manageability required by most production systems.

After selecting **A file-based password-protected Java KeyStore** as the root of trust, the following steps are displayed:

1. Enable Kerberos

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

For more information on enabling Kerberos, see [Enabling Kerberos Authentication Using the Wizard](#) on page 59.

Configuring Encryption

2. Enable TLS/SSL

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

For more information on enabling TLS, see [Configuring TLS Security for Cloudera Manager](#) on page 211.

3. Add a Java KeyStore KMS Service

Minimum Required Role: [Key Administrator](#) (also provided by **Full Administrator**)

This step adds the Java KeyStore KMS service to the cluster. The Java KeyStore KMS service uses a password-protected Java KeyStore for cryptographic key management. To complete this step:

1. Click **Add a Java KeyStore KMS Service**.
2. Select a cluster host for the Java KeyStore KMS service. Click **Continue**.
3. The **Setup TLS for Java KeyStore KMS** page provides high-level instructions for configuring TLS communication between the EDH cluster and the Java KeyStore KMS. See [Configuring TLS/SSL for the KMS](#) on page 280 for more information.

Click **Continue**.

4. The **Review Changes** page lists the Java KeyStore settings. Click the ⓘ icon next to any setting for information about that setting. Enter the location and password for the Java KeyStore and click **Continue**.
5. Click **Continue** to automatically configure the HDFS service to depend on the Java KeyStore KMS service.
6. Click **Finish** to complete this step and return to the main page of the wizard.

4. Restart stale services and redeploy client configuration

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

This step restarts all services which were modified while enabling HDFS encryption. To complete this step:

1. Click **Restart stale services and redeploy client configuration**.
2. Click **Restart Stale Services**.
3. Make sure that **Re-deploy client configuration** is checked, and click **Restart Now**.
4. After all commands have completed, click **Finish**.

5. Validate Data Encryption

Minimum Required Role: [Key Administrator](#) or [Cluster Administrator](#) (also provided by Full Administrator)

This step launches a tutorial with instructions on creating an encryption zone and putting data into it to verify that HDFS encryption is enabled and working.

Managing Encryption Keys and Zones

Interacting with the KMS and creating encryption zones requires the use of two new CLI commands: `hadoop key` and `hdfs crypto`. The following sections will help you get started with creating encryption keys and setting up encryption zones.

Before continuing, make sure that your KMS ACLs have been set up according to best practices. For more information, see [Configuring KMS Access Control Lists](#) on page 281.

Validating Hadoop Key Operations



Warning: If you are using or plan to use Cloudera Navigator Key HSM in conjunction with Cloudera Navigator Key Trustee Server, ensure that key names begin with alphanumeric characters and do not use special characters other than hyphen (-), period (.), or underscore (_). Using other special characters can prevent you from migrating your keys to an HSM. See [Integrating Key HSM with Key Trustee Server](#) on page 314 for more information.

Use `hadoop key create` to create a test key, and then use `hadoop key list` to retrieve the key list:

```
$ sudo -u <key_admin> hadoop key create keytrustee_test
$ hadoop key list
```

Creating Encryption Zones

Once a KMS has been set up and the NameNode and HDFS clients have been correctly configured, use the `hadoop key` and `hdfs crypto` command-line tools to create encryption keys and set up new encryption zones.

- Create an encryption key for your zone as the `keyadmin` for the user/group (regardless of the application that will be using the encryption zone):

```
$ sudo -u hdfs hadoop key create <key_name>
```

- Create a new empty directory and make it an encryption zone using the key created above.

```
$ sudo -u hdfs hadoop fs -mkdir /encryption_zone
$ sudo -u hdfs hdfs crypto -createZone -keyName <key_name> -path /encryption_zone
```

You can verify creation of the new encryption zone by running the `-listZones` command. You should see the encryption zone along with its key listed as follows:

```
$ sudo -u hdfs hdfs crypto -listZones
/encryption_zone <key_name>
```



Warning: Do not delete an encryption key as long as it is still in use for an encryption zone. This results in loss of access to data in that zone.



Important: The Java Keystore KMS default Truststore (for example, `org.apache.hadoop.crypto.key.JavaKeyStoreProvider`) does not support uppercase key names.

For more information and recommendations on creating encryption zones for each CDH component, see [Configuring CDH Services for HDFS Encryption](#) on page 290.

Adding Files to an Encryption Zone

Existing data can be encrypted by copying it copied into the new encryption zones using tools like `DistCp`. See [DistCp Considerations](#) on page 262 for information on using `DistCp` with encrypted data files.

You can add files to an encryption zone by copying them to the encryption zone using `distcp`. For example:

```
sudo -u hdfs hadoop distcp /user/dir /encryption_zone
```

Additional Information:

- For more information on KMS setup and high availability configuration, see [Configuring the Key Management Server \(KMS\)](#) on page 274.
- For instructions on securing the KMS using Kerberos, TLS/SSL communication and ACLs, see [Securing the Key Management Server \(KMS\)](#) on page 278.
- If you want to use the KMS to encrypt data used by other CDH services, see [Configuring CDH Services for HDFS Encryption](#) on page 290 for information on recommended encryption zones for each service.

Deleting Encryption Zones

To remove an encryption zone, delete the encrypted directory:

Configuring Encryption



Warning: This command deletes the entire directory and all of its contents. Ensure that the data is no longer needed before running this command.

```
$ sudo -u hdfs hadoop fs -rm -r -skipTrash /encryption_zone
```



Important: The Key Trustee KMS does not directly execute a key deletion (for example, it may perform a soft delete instead, or delay the actual deletion to prevent mistakes). In these cases, errors may occur when creating or deleting a key using the same name after it has already been deleted.

Backing Up Encryption Keys



Warning: It is *very important* that you regularly back up your encryption keys. Failure to do so can result in irretrievable loss of encrypted data.

If you are using the Java KeyStore KMS, make sure you regularly back up the Java KeyStore that stores the encryption keys. If you are using the Key Trustee KMS and Key Trustee Server, see [Backing Up and Restoring Key Trustee Server](#) on page 298 for instructions on backing up Key Trustee Server and Key Trustee KMS.

Configuring the Key Management Server (KMS)

Hadoop Key Management Server (KMS) is a cryptographic key management server based on the Hadoop **KeyProvider** API. It provides a KeyProvider implementation client that interacts with the KMS using the HTTP REST API. Both the KMS and its client support HTTP SPNEGO Kerberos authentication and TLS/SSL-secured communication. The KMS is a Java-based web application that uses a preconfigured Tomcat server bundled with the Hadoop distribution.

For instructions on securing the KMS, see [Securing the Key Management Server \(KMS\)](#) on page 278.

Cloudera provides two implementations of the Hadoop KMS:

- **Java KeyStore KMS** - The default Hadoop KMS included in CDH that uses a file-based Java KeyStore (JKS) for its backing keystore. For parcel-based installations, no additional action is required to install or upgrade the KMS. For package-based installations, you must install additional packages. For more information, see [Installing and Upgrading Java KeyStore KMS](#). Cloudera strongly recommends not using Java KeyStore KMS in production environments.
- **Key Trustee KMS** - A custom KMS that uses [Cloudera Navigator Key Trustee Server](#) for its backing keystore instead of the file-based Java KeyStore (JKS) used by the default Hadoop KMS. Cloudera strongly recommends using Key Trustee KMS in production environments to improve the security, durability, and scalability of your cryptographic key management. For more information about the architecture and components involved in encrypting data at rest for production environments, see [Cloudera Navigator Data Encryption Overview](#) and [Data at Rest Encryption Reference Architecture](#) on page 256. For instructions on installing and upgrading Key Trustee KMS, see:
 - [Installing Key Trustee KMS](#)
 - [Upgrading Key Trustee KMS](#)

Configuring KMS High Availability

For Key Trustee KMS high availability, see [Key Trustee KMS High Availability](#). Java KeyStore KMS does not support high availability.

Configuring the KMS Using Cloudera Manager

If you are using Cloudera Manager, you can view and edit the KMS configuration by navigating to the following pages, depending on the KMS implementation you are using:

- **Key Trustee KMS service > Configuration**
- **Java KeyStore KMS service > Configuration**

For more information on using Cloudera Manager to find and change configuration parameters, see [Modifying Configuration Properties Using Cloudera Manager](#).

For instructions about configuring the KMS and its clients using the command line for package-based installations, continue reading:

Configuring the KMS Cache Using Cloudera Manager

By default, the KMS caches keys to reduce the number of interactions with the key provider. You can disable the cache by setting the `hadoop.kms.cache.enable` property to `false`.

The cache is only used with the `getCurrentKey()`, `getKeyVersion()` and `getMetadata()` methods.

For the `getCurrentKey()` method, entries are cached for a maximum of 30000 milliseconds to prevent stale keys.

For the `getKeyVersion()` method, entries are cached with a default inactivity timeout of 600000 milliseconds (10 minutes).

You can configure the cache and its timeout values by adding the following properties to **KMS service > Configuration > Advanced > Key Management Server Proxy Advanced Configuration Snippet (Safety Valve) for kms-site.xml**:

```
<property>
  <name>hadoop.kms.cache.enable</name>
  <value>true</value>
</property>

<property>
  <name>hadoop.kms.cache.timeout.ms</name>
  <value>600000</value>
</property>

<property>
  <name>hadoop.kms.current.key.cache.timeout.ms</name>
  <value>30000</value>
</property>
```

See [Custom Configuration](#) for more information on adding custom properties using the **Advanced Configuration Snippet (Safety Valve)** feature.

Configuring the Audit Log Aggregation Interval Using the Command Line

Audit logs are generated for `GET_KEY_VERSION`, `GET_CURRENT_KEY`, `DECRYPT_EEK`, and `GENERATE_EEK` operations.

Entries are aggregated by user, key, and operation for a configurable interval, after which the number of aggregated operations by the user for a given key is written to the audit log.

The interval is configured in milliseconds by adding the `hadoop.kms.aggregation.delay.ms` property to **KMS service > Configuration > Advanced > Key Management Server Proxy Advanced Configuration Snippet (Safety Valve) for kms-site.xml**:

```
<property>
  <name>hadoop.kms.aggregation.delay.ms</name>
  <value>10000</value>
</property>
```

For more information about adding custom properties using the **Advanced Configuration Snippet (Safety Valve)** feature, see [Custom Configuration](#).

Configuring the Java KeyStore KMS Using the Command Line



Note: Because Key Trustee KMS is supported only in Cloudera Manager deployments, the following command line instructions apply only to Java KeyStore KMS. For instructions about configuring Key Trustee KMS, see [Configuring the KMS Using Cloudera Manager](#) on page 274.

For instructions about configuring the Java KeyStore KMS and its clients using the command line for package-based installations, continue reading:

Configuring the Java KeyStore KMS KeyProvider Using the Command Line

**Important:**

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

Configure the KMS backing KeyProvider properties in the `/etc/hadoop-kms/conf/kms-site.xml` configuration file:

```
<property>
  <name>hadoop.kms.key.provider.uri</name>
  <value>jceks://file@/${user.home}/kms.keystore</value>
</property>

<property>
  <name>hadoop.security.keystore.java-keystore-provider.password-file</name>
  <value>keystore_password_file</value>
</property>
```

If you do not specify the absolute path to the password file, you must include it in the Hadoop CLASSPATH.

Restart the KMS for configuration changes to take effect. See [Starting and Stopping the Java KeyStore KMS Using the Command Line](#) on page 277 for instructions.

Configuring the Java KeyStore KMS Cache Using the Command Line

**Important:**

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

By default, the KMS caches keys to reduce the number of interactions with the key provider. You can disable the cache by setting the `hadoop.kms.cache.enable` property to `false`.

The cache is only used with the `getCurrentKey()`, `getKeyVersion()` and `getMetadata()` methods.

For the `getCurrentKey()` method, entries are cached for a maximum of 30000 milliseconds to prevent stale keys.

For the `getKeyVersion()` method, entries are cached with a default inactivity timeout of 600000 milliseconds (10 minutes).

The cache and its timeout values are configured using the following properties in the `/etc/hadoop-kms/conf/kms-site.xml` configuration file:

```
<property>
  <name>hadoop.kms.cache.enable</name>
  <value>>true</value>
</property>

<property>
  <name>hadoop.kms.cache.timeout.ms</name>
  <value>600000</value>
</property>

<property>
  <name>hadoop.kms.current.key.cache.timeout.ms</name>
  <value>30000</value>
</property>
```

Configuring KMS Clients Using the Command Line

**Important:**

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

To configure KMS clients, set the `hadoop.security.key.provider.path` property in `core-site.xml` or `hdfs-site.xml`. Specify the value in the format `kms://<scheme>@<kms_hosts>:<port>/kms`. Replace `<scheme>` with `http` or `https`, depending on whether you have [configured TLS](#). Replace `<kms_hosts>` with a semicolon-separated list of the KMS hosts. Replace `<port>` with the port number on which the KMS is running (16000 by default).

For example, for a KMS running on `http://localhost:16000/kms`, the KeyProvider URI is `kms://http@localhost:16000/kms`. For high availability KMS (Key Trustee KMS only) running on `https://kms01.example.com:16000/kms` and `https://kms02.example.com:16000/kms`, the KeyProvider URI is `kms://https@kms01.example.com;kms02.example.com:16000/kms`.

See the following for an excerpt from `core-site.xml`:

```
<property>
  <name>hadoop.security.key.provider.path</name>
  <value>kms://https@kms01.example.com;kms02.example.com:16000/kms</value>
</property>
```

Starting and Stopping the Java KeyStore KMS Using the Command Line

**Important:**

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

To start or stop KMS use the `kms.sh` script. For example, to start the KMS:

```
$ sudo /usr/lib/hadoop-kms/sbin/kms.sh start
```

Running the script without parameters lists all possible parameters.

To use an `init` script to manage the KMS service, use your package manager to install the `hadoop-kms-server` package from the [CDH repository](#). For example, for RHEL 6:

```
$ sudo yum install hadoop-kms-server
```

After installation, use the `service hadoop-kms-server` command to manage the KMS service.

Configuring the Audit Log Aggregation Interval Using the Command Line

**Important:**

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

Audit logs are generated for `GET_KEY_VERSION`, `GET_CURRENT_KEY`, `DECRYPT_EEK`, and `GENERATE_EEK` operations.

Entries are aggregated by user, key, and operation for a configurable interval, after which the number of aggregated operations by the user for a given key is written to the audit log.

Configuring Encryption

The interval is configured in milliseconds using the `hadoop.kms.aggregation.delay.ms` property:

```
<property>
  <name>hadoop.kms.aggregation.delay.ms</name>
  <value>10000</value>
</property>
```

Configuring the Embedded Tomcat Server Using the Command Line



Important:

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

You can configure the embedded Tomcat server by using the `/etc/hadoop-kms/tomcat-conf/conf/server.xml.conf` file.

The following environment variables can be set in KMS `/etc/hadoop-kms/conf/kms-env.sh` script and can be used to alter the default ports and log directory:

- `KMS_HTTP_PORT`
- `KMS_ADMIN_PORT`
- `KMS_LOG`

[Restart the KMS](#) for the configuration changes to take effect.

Securing the Key Management Server (KMS)

Cloudera provides two implementations of the Hadoop KMS:

- **Java KeyStore KMS** - The default Hadoop KMS included in CDH that uses a file-based Java KeyStore (JKS) for its backing keystore. For parcel-based installations, no additional action is required to install or upgrade the KMS. For package-based installations, you must install additional packages. For more information, see [Installing and Upgrading Java KeyStore KMS](#). Cloudera strongly recommends not using Java Keystore KMS in production environments.
- **Key Trustee KMS** - A custom KMS that uses [Cloudera Navigator Key Trustee Server](#) for its backing keystore instead of the file-based Java KeyStore (JKS) used by the default Hadoop KMS. Cloudera strongly recommends using Key Trustee KMS in production environments to improve the security, durability, and scalability of your cryptographic key management. For more information about the architecture and components involved in encrypting data at rest for production environments, see [Cloudera Navigator Data Encryption Overview](#) and [Data at Rest Encryption Reference Architecture](#) on page 256. For instructions on installing and upgrading Key Trustee KMS, see:
 - [Installing Key Trustee KMS](#)
 - [Upgrading Key Trustee KMS](#)

This topic contains information on securing the KMS using Kerberos, TLS/SSL communication, and access control lists (ACLs) for operations on encryption keys. Cloudera Manager instructions can be performed for both Key Trustee KMS and Java KeyStore KMS deployments. Command-line instructions apply only to Java KeyStore KMS deployments. Key Trustee KMS is not supported outside of Cloudera Manager. See [Installing Key Trustee KMS](#) for more information.

Enabling Kerberos Authentication for the KMS

Enabling Kerberos Authentication for the KMS Using Cloudera Manager

Minimum Required Role: [Full Administrator](#)


To enable Kerberos for the KMS using Cloudera Manager:

1. Open the Cloudera Manager Admin Console and go to the KMS service.
2. Click **Configuration**.

3. Set the **Authentication Type** property to `kerberos`.
4. Click **Save Changes**.
5. Because Cloudera Manager does not automatically create the principal and keytab file for the KMS, you must run the **Generate Credentials** command manually. On the top navigation bar, go to **Administration > Security > Kerberos Credentials** and click **Generate Missing Credentials**.



Note: This does not create a new Kerberos principal if an existing HTTP principal exists for the KMS host.

6. Return to the Home page by clicking the Cloudera Manager logo.
7. Click  to invoke the cluster restart wizard.
8. Click **Restart Stale Services**.
9. Click **Restart Now**.
10. Click **Finish**.

Enabling Kerberos Authentication for the Java KeyStore KMS Using the Command Line



Important:

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

Configure `/etc/krb5.conf` with information for your KDC server. Create an HTTP principal and keytab file for the KMS.

Configure `/etc/hadoop-kms/conf/kms-site.xml` with the following properties:

```
<property>
  <name>hadoop.kms.authentication.type</name>
  <value>kerberos</value>
</property>

<property>
  <name>hadoop.kms.authentication.kerberos.keytab</name>
  <value>${user.home}/kms.keytab</value>
</property>

<property>
  <name>hadoop.kms.authentication.kerberos.principal</name>
  <value>HTTP/localhost</value>
</property>

<property>
  <name>hadoop.kms.authentication.kerberos.name.rules</name>
  <value>DEFAULT</value>
</property>
```

[Restart the KMS](#) service for the configuration changes to take effect.

Configuring the Java KeyStore KMS Proxyuser Using the Command Line



Important:

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

Each proxyuser must be configured in `/etc/hadoop-kms/conf/kms-site.xml` using the following properties:

```
<property>
  <name>hadoop.kms.proxyuser.#USER#.users</name>
  <value>*</value>
</property>

<property>
  <name>hadoop.kms.proxyuser.#USER#.groups</name>
  <value>*</value>
</property>

<property>
  <name>hadoop.kms.proxyuser.#USER#.hosts</name>
  <value>*</value>
</property>
```

where `#USER#` is the username of the proxyuser to be configured.

The `hadoop.kms.proxyuser.#USER#.users` property indicates the users that can be impersonated. The `hadoop.kms.proxyuser.#USER#.groups` property indicates the groups to which the users being impersonated must belong. At least one of these properties must be defined. If both are defined, the configured proxyuser can impersonate any user in the `users` list and any user belonging to a group listed in the `groups` list.

The `hadoop.kms.proxyuser.#USER#.hosts` property indicates the host from which the proxyuser can make impersonation requests. "*" means there are no restrictions for the `#USER#` regarding users, groups, or hosts.

Configuring TLS/SSL for the KMS

Configuring TLS/SSL for the KMS Using Cloudera Manager

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)


The steps for configuring and enabling Hadoop TLS/SSL for the KMS are as follows:

1. Go to the KMS service.
2. Click **Configuration**.
3. In the Search field, type **TLS/SSL** to show the KMS TLS/SSL properties (in the **Key Management Server Default Group > Security** category).
4. Edit the following TLS/SSL properties according to your cluster configuration.

Table 21: KMS TLS/SSL Properties

Property	Description
Enable TLS/SSL for Key Management Server	Encrypt communication between clients and Key Management Server using Transport Layer Security (TLS) (formerly known as Secure Socket Layer (TLS/SSL)).
Key Management Server TLS/SSL Server JKS Keystore File Location	The path to the TLS/SSL keystore file containing the server certificate and private key used for TLS/SSL. Used when Key Management Server is acting as a TLS/SSL server. The keystore must be in JKS format.
Key Management Server TLS/SSL Server JKS Keystore File Password	The password for the Key Management Server JKS keystore file.
Key Management Server Proxy TLS/SSL Certificate Trust Store File	The location on disk of the truststore, in .jks format, used to confirm the authenticity of TLS/SSL servers that Key Management Server Proxy might connect to. This is used when Key Management Server Proxy is the client in a TLS/SSL connection. This truststore must contain the certificates used to sign the services connected to. If this parameter is not provided, the default list of well-known certificate authorities is used instead.

Property	Description
Key Management Server Proxy TLS/SSL Certificate Trust Store Password	The password for the Key Management Server Proxy TLS/SSL Certificate Trust Store File. This password is not required to access the truststore; this field can be left blank. This password provides optional integrity checking of the file. The contents of truststores are certificates, and certificates are public information.

5. Click **Save Changes**.
6. Return to the Home page by clicking the Cloudera Manager logo.
7. Click  to invoke the cluster restart wizard.
8. Click **Restart Stale Services**.
9. Click **Restart Now**.
10. Click **Finish**.

For help troubleshooting TLS/SSL for KMS configuration issues, see [Troubleshooting TLS/SSL Connectivity](#) on page 226.

Configuring TLS/SSL for the Java KeyStore KMS Using the Command Line



Important:

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

To configure KMS to work over HTTPS, set the following properties in the `/etc/hadoop-kms/conf/kms_env.sh` script:

- `KMS_SSL_KEYSTORE_FILE`
- `KMS_SSL_KEYSTORE_PASS`
- `KMS_SSL_TRUSTSTORE_FILE`
- `KMS_SSL_TRUSTSTORE_PASS`

In the `/etc/hadoop-kms/tomcat-conf/conf/` directory, replace the `server.xml` file with the provided `ssl-server.xml` file.

Create a TLS/SSL certificate for the KMS. As the `kms` user, use the Java `keytool` command to create the TLS/SSL certificate:

```
$ keytool -genkey -alias tomcat -keyalg RSA
```

You are asked a series of questions in an interactive prompt. It creates the keystore file, which is named `.keystore` and located in the `kms` user home directory. The password you enter for the keystore must match the value of the `KMS_SSL_KEYSTORE_PASS` environment variable set in the `kms-env.sh` script in the configuration directory.

The answer to "What is your first and last name?" (CN) must be the hostname of the machine where the KMS will be running.



Note: Restart the KMS for the configuration changes to take effect.

Configuring KMS Access Control Lists

Hadoop KMS supports a range of ACLs that control access to keys and key operations on a granular basis. ACLs can be used, for instance, to only grant users access to certain keys. Restricting HDFS superusers from access to key material is an important design requirement. This prevents a malicious superuser from having access to all the key material and all the encrypted data, and thus being able to decrypt everything.

There are two categories of KMS ACLs:

1. KMS-wide: These ACLs specify the types of operations a user can perform. They are configured using the `hadoop.kms.acl.<OPERATION>` and `hadoop.kms.blacklist.<OPERATION>` parameters. The operations are as follows:

- CREATE
- DELETE
- ROLLOVER
- GET
- GET_KEYS
- GET_METADATA
- SET_KEY_MATERIAL
- GENERATE_EEK
- DECRYPT_EEK

2. Key-specific: These ACLs are set in a per-key basis. They are configured using the `default.key.acl.<OPERATION>`, `whitelist.key.acl.<OPERATION>`, and `key.acl.<key_name>.<OPERATION>` parameters. The operations and their programmatic equivalents are as follows:

- READ - `getKeyVersion`, `getKeyVersions`, `getMetadata`, `getKeysMetadata`, `getCurrentKey`
- MANAGEMENT - `createKey`, `deleteKey`, `rolloverNewVersion`
- GENERATE_EEK - `generateEncryptedKey`, `warmUpEncryptedKeys`
- DECRYPT_EEK - `decryptEncryptedKey`
- ALL - All of the above

The `default.key.acl.<OPERATION>` ACL applies to all keys for which an ACL has not been explicitly configured.

If no ACL is configured for a specific key, *and* no default ACL is configured for the requested operation, access is denied.



Note: The default ACL does not support the ALL operation qualifier.

The KMS supports both whitelist and blacklist ACLs. Blacklist entries override whitelist entries. A user or group accessing the KMS is first checked for inclusion in the ACL for the requested operation and then checked for exclusion in the blacklist for the operation before access is granted.

The group membership used by ACL entries relies on the configured group mapping mechanism for HDFS. By default, group membership is determined on the local Linux system running the KMS service. If you have configured HDFS to use LDAP for group mapping, the group membership for the ACL entries is determined using the configured LDAP settings. For more information about LDAP-based group membership, see [Configuring LDAP Group Mappings](#) on page 350.

The ACL syntax for both blacklist and whitelist entries is as follows:

- Users only:

```
user1,user2,userN
```



Note: There are no spaces following the commas separating the users in the list.

- Groups only:

```
nobody group1,group2,groupN
```



Note: There is a space between `nobody` and the comma-separated group list. The `nobody` user, if it exists, must not have privileges to log in to or interact with the system. If you are uncertain about its access privileges, specify a different nonexistent user in its place.

- Users and Groups:

```
user1,user2,userN group1,group2,groupN
```



Note: The comma-separated user list is separated from the comma-separated group list by a space.

Configuring KMS Access Control Lists Using Cloudera Manager

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)



Important: See related Known Issue and listed workaround: [KMS and Key Trustee ACLs do not work in Cloudera Manager 5.3](#).

The KMS installation wizard includes an option to generate the recommended ACLs. To view or edit the ACLs:

1. Go to the KMS service.
2. Click **Configuration**.
3. In the Search field, type `acl` to show the **Key Management Server Advanced Configuration Snippet (Safety Valve) for kms-acls.xml** (in the **Key Management Server Default Group** category).
4. Add or edit the ACL properties according to your cluster configuration. See [Recommended KMS Access Control List](#) on page 283 for example ACL entries.
5. Click **Save Changes**.
6. Return to the Home page by clicking the Cloudera Manager logo.
7. Click to invoke the cluster restart wizard.
8. Click **Restart Stale Services**.
9. Click **Restart Now**.
10. Click **Finish**.

Configuring Java KeyStore KMS Access Control Lists Using the Command Line



Important:

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

KMS ACLs are defined in the `/etc/hadoop-kms/conf/kms-acls.xml` configuration file. This file is hot-reloaded when it changes. See [Recommended KMS Access Control List](#) on page 283 for recommended ACL entries.

Recommended KMS Access Control List

Cloudera recommends the following ACL definition for secure production settings. Replace `keyadmin` and `keyadmingroup` with the user and group responsible for maintaining encryption keys.



Note: If you are entering the ACL using Cloudera Manager, omit the surrounding `<configuration>` and `</configuration>` tags; Cloudera Manager adds this automatically.

```
<configuration>
<!--
  KMS ACLs control which users can perform various actions on the KMS,
  and which users and groups have access to which keys.

  This file has the following sections:
  * ACLs for KMS operations
  ** Access to specific KMS operations
  ** Blacklists for those specific operations
  * ACLs for keys
  ** Default ACLs for keys
  ** Whitelist ACLs for keys
  ** Key-specific ACLs
-->
<!--
  KMS ACLs that govern access to specific key operations. If access is not
  granted for an operation here, then the operation is forbidden, even if
  a key ACL allows it.

  The ACL value should be either a username or a username and group name
  separated by whitespace.

  A value of "*" (for the username or groupname) indicates that
  all users are granted access to that operation. Any operation for which
  there is no ACL or an empty (zero-length) ACL is treated as having an
  ACL with a value of "*". To disallow all users, add an ACL with a
  value of " ", a single space.

  Note: This convention applies only to the KMS-level ACLs beginning with
  'hadoop.kms.acl'.
-->
<property>
  <name>hadoop.kms.acl.CREATE</name>
  <value>keyadmin keyadmingroup</value>
  <description>
    ACL for create-key operations.
    If the user is not in the GET ACL, the key material is not returned
    as part of the response.
  </description>
</property>
<property>
  <name>hadoop.kms.acl.DELETE</name>
  <value>keyadmin keyadmingroup</value>
  <description>
    ACL for delete-key operations.
  </description>
</property>
<property>
  <name>hadoop.kms.acl.ROLLOVER</name>
  <value>keyadmin keyadmingroup</value>
  <description>
    ACL for rollover-key operations.
    If the user does is not in the GET ACL, the key material is not returned
    as part of the response.
  </description>
</property>
<property>
  <name>hadoop.kms.acl.GET</name>
  <value></value>
  <description>
    ACL for get-key-version and get-current-key operations.
  </description>
</property>
</configuration>
```

```

</property>

<property>
  <name>hadoop.kms.acl.GET_KEYS</name>
  <value>keyadmin keyadmingroup</value>
  <description>
    ACL for get-keys operations.
  </description>
</property>

<property>
  <name>hadoop.kms.acl.SET_KEY_MATERIAL</name>
  <value></value>
  <description>
    Complementary ACL for CREATE and ROLLOVER operations to allow the client
    to provide the key material when creating or rolling a key.
  </description>
</property>

<property>
  <name>hadoop.kms.acl.GENERATE_EEK</name>
  <value>hdfs supergroup</value>
  <description>
    ACL for generateEncryptedKey CryptoExtension operations.
  </description>
</property>

<!--
  KMS blacklists to prevent access to operations. These settings override the
  permissions granted by the KMS ACLs above.

  The blacklist value should be either a username or a username and group name
  separated by whitespace.

  A blank value indicates that no user is blacklisted from the operation. A
  value of "*" (for either the username or groupname) indicates that all users
  are blacklisted from the operation. Any operation for which there is no
  blacklist will be treated as having a blacklist with an empty value.
-->

<!--
  In this template the hdfs user is blacklisted for everything except
  GET_METADATA, GET_KEYS, and GENERATE_EEK. The GET and SET_KEY_MATERIAL
  operations are blacklisted for all users since Hadoop users should not
  need to perform those operations, and access to the key material should
  be as restricted as possible.
-->

<property>
  <name>hadoop.kms.blacklist.CREATE</name>
  <value>hdfs supergroup</value>
</property>

<property>
  <name>hadoop.kms.blacklist.DELETE</name>
  <value>hdfs supergroup</value>
</property>

<property>
  <name>hadoop.kms.blacklist.ROLLOVER</name>
  <value>hdfs supergroup</value>
</property>

<property>
  <name>hadoop.kms.blacklist.GET</name>
  <value>*</value>
</property>

<property>
  <name>hadoop.kms.blacklist.GET_KEYS</name>
  <value></value>
</property>

```

```

<property>
  <name>hadoop.kms.blacklist.SET_KEY_MATERIAL</name>
  <value>*</value>
</property>

<property>
  <name>hadoop.kms.blacklist.DECRYPT_EEK</name>
  <value>hdfs supergroup</value>
</property>

<property>
  <name>keytrustee.kms.acl.UNDELETE</name>
  <value></value>
  <description>
    ACL that grants access to the UNDELETE operation on all keys.
    Only used by Key Trustee KMS.
  </description>
</property>

<property>
  <name>keytrustee.kms.acl.PURGE</name>
  <value></value>
  <description>
    ACL that grants access to the PURGE operation on all keys.
    Only used by Key Trustee KMS.
  </description>
</property>

<!--
  Default key ACLs that govern access to key operations for key-operation pairs
  that do not have a specific key ACL already. Specific key ACLs will override
  the default key ACLs

  The ACL value should be either a username or a username and group name
  separated by whitespace.

  An empty value for an ACL indicates that no user is granted access to that
  operation. A value of "*" (for the username or groupname) indicates that
  all users are granted access to that operation. Any operation for which
  there is no ACL will be treated as having an ACL with an empty value.
-->

<property>
  <name>default.key.acl.MANAGEMENT</name>
  <value></value>
  <description>
    Default ACL that grants access to the MANAGEMENT operation on all keys.
  </description>
</property>

<property>
  <name>default.key.acl.GENERATE_EEK</name>
  <value></value>
  <description>
    Default ACL that grants access to the GENERATE_EEK operation on all keys.
  </description>
</property>

<property>
  <name>default.key.acl.DECRYPT_EEK</name>
  <value></value>
  <description>
    Default ACL that grants access to the DECRYPT_EEK operation on all keys.
  </description>
</property>

<property>
  <name>default.key.acl.READ</name>
  <value></value>
  <description>
    Default ACL that grants access to the READ operation on all keys.
  </description>
</property>

```

```

<!--
Whitelist key ACLs that grant access to specific key operations. Any
permissions granted here will be added to whatever permissions are granted
by the specific key ACL or the default key ACL. Note that these whitelist
ACLs grant access to operations on specific keys. If the operations
themselves are not allowed because of the KMS ACLs/blacklists, then the
operation will not be permitted, regardless of the whitelist settings.

The ACL value should be either a username or a username and group name
separated by whitespace.

An empty value for an ACL indicates that no user is granted access to that
operation. A value of "*" (for the username or groupname) indicates that
all users are granted access to that operation. Any operation for which
there is no ACL will be treated as having an ACL with an empty value.
-->

<property>
  <name>whitelist.key.acl.MANAGEMENT</name>
  <value>keyadmin keyadmingroup</value>
  <description>
    Whitelist ACL for MANAGEMENT operations for all keys.
  </description>
</property>

<property>
  <name>whitelist.key.acl.READ</name>
  <value>hdfs supergroup</value>
  <description>
    Whitelist ACL for READ operations for all keys.
  </description>
</property>

<property>
  <name>whitelist.key.acl.GENERATE_EEK</name>
  <value>hdfs supergroup</value>
  <description>
    Whitelist ACL for GENERATE_EEK operations for all keys.
  </description>
</property>

<property>
  <name>whitelist.key.acl.DECRYPT_EEK</name>
  <value>keyadmin keyadmingroup</value>
  <description>
    Whitelist ACL for DECRYPT_EEK operations for all keys.
  </description>
</property>

<!--
Key ACLs that grant access to specific key operations. Any permissions
granted here are added to whatever permissions are granted by the whitelists.
The key ACL name should be key.acl.<keyname>.<OPERATION>.

The ACL value should be either a username or a username and group name
separated by whitespace.

An empty value for an ACL indicates that no user is granted access to that
operation. A value of "*" (for the username or groupname) indicates that
all users are granted access to that operation. Any key operation for which
there is no ACL will default to the default ACL for the operation.

Normally adding users or groups for a specific key and DECRYPT_EEK is
sufficient to allow access to data protected with HDFS data at rest
encryption.
-->

<!--
The following ACLs are required for proper functioning of services.
CM does not create keys or encryption zones, however our best practices
recommend encryption zones on certain directories. Below we assume that
the user has followed our recommended naming scheme and named the keys

```

```

according to our best practices: "hive-key" for the hive service,
"hbase-key" for the hbase service, etc. If the key names are different,
none of this will work out of the box, and you will need to edit these
ACLs to match your key names.
-->

<property>
  <name>key.acl.hive-key.DECRYPT_EEK</name>
  <value>hive hive</value>
  <description>
    Gives the hive user and the hive group access to the key named "hive-key".
    This allows the hive service to read and write files in /user/hive/.
    Also note that the impala user ought to be a member of the hive group
    in order to enjoy this same access.
  </description>
</property>

<property>
  <name>key.acl.hive-key.READ</name>
  <value>hive hive</value>
  <description>
    Required because hive compares key strengths when joining tables.
  </description>
</property>

<property>
  <name>key.acl.hbase-key.DECRYPT_EEK</name>
  <value>hbase hbase</value>
  <description>
    Gives the hbase user and hbase group access to the key named "hbase-key".
    This allows the hbase service to read and write files in /hbase.
  </description>
</property>

<property>
  <name>key.acl.solr-key.DECRYPT_EEK</name>
  <value>solr solr</value>
  <description>
    Gives the solr user and solr group access to the key named "solr-key".
    This allows the solr service to read and write files in /solr.
  </description>
</property>

<property>
  <name>key.acl.mapred-key.DECRYPT_EEK</name>
  <value>mapred,yarn hadoop</value>
  <description>
    Gives the mapred user and mapred group access to the key named "mapred-key".
    This allows mapreduce to read and write files in /user/history.
    This is required by YARN.
  </description>
</property>

<property>
  <name>key.acl.hue-key.DECRYPT_EEK</name>
  <value>oozie,hue oozie,hue</value>
  <description>
    Gives the appropriate users and groups access to the key named "hue-key".
    This allows hue and oozie to read and write files in /user/hue.
    Oozie is required here because it will attempt to access workflows in
    /user/hue/oozie/workspaces.
  </description>
</property>

<!-- This example is required if there are encryption zones on user's home
directories. -->
<!--
<property>
  <name>key.acl.username-key.DECRYPT_EEK</name>
  <value>username username,hive,hbase,solr,oozie,hue,yarn</value>
  <description>

```



```

    Designed to be placed on a key that protects the EZ /user/username,
    and assumes that the key name is also "username-key", this shows that
    a number of services may want to reach in to access data. Remove
    those are are not needed for your use-case.
  </description>
</property>
-->
</configuration>

```

Configuring Java KeyStore KMS Delegation Tokens Using the Command Line



Important:

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

Configure KMS delegation token secret manager using the following properties:

```

<property>
  <name>hadoop.kms.authentication.delegation-token.update-interval.sec</name>
  <value>86400</value>
  <description>
    How often the master key is rotated, in seconds. Default value 1 day.
  </description>
</property>

<property>
  <name>hadoop.kms.authentication.delegation-token.max-lifetime.sec</name>
  <value>604800</value>
  <description>
    Maximum lifetime of a delegation token, in seconds. Default value 7 days.
  </description>
</property>

<property>
  <name>hadoop.kms.authentication.delegation-token.renew-interval.sec</name>
  <value>86400</value>
  <description>
    Renewal interval of a delegation token, in seconds. Default value 1 day.
  </description>
</property>

<property>
  <name>hadoop.kms.authentication.delegation-token.removal-scan-interval.sec</name>
  <value>3600</value>
  <description>
    Scan interval to remove expired delegation tokens.
  </description>
</property>

```

Migrating Keys from a Java KeyStore to Cloudera Navigator Key Trustee Server

You can migrate keys from an existing Java KeyStore (JKS) to Key Trustee Server to improve security, durability, and scalability. If you are using the Java KeyStore KMS service, and want to use Key Trustee Server as the backing key store for [HDFS Transparent Encryption](#) on page 261, use the following procedure.

This procedure assumes that the Java KeyStore (JKS) is on the same host as the new Key Trustee KMS service.

1. Stop the Java KeyStore KMS service.
2. Add and configure the Key Trustee KMS service, and configure HDFS to use it for its **KMS Service** setting. For more information about how to install Key Trustee KMS, see [Installing Key Trustee KMS](#). Restart the HDFS service and redeploy client configuration for this to take effect: **Home** > **Cluster-wide** ▾ > **Deploy Client Configuration**

3. Add the following to the **Key Management Server Proxy Advanced Configuration Snippet (Safety Valve) for kms-site.xml (Key Trustee KMS Service > Configuration > Category > Advanced)**:

```
<property>
  <name>hadoop.kms.key.provider.uri</name>

  <value>keytrustee://file@/var/lib/kms-keytrustee/keytrustee/.keytrustee/,jceks://file@/path/to/kms.keystore</value>

  <description>URI of the backing KeyProvider for the KMS</description>
</property>

<property>
  <name>hadoop.security.keystore.java-keystore-provider.password-file</name>
  <value>/tmp/password.txt</value>
  <description>Java KeyStore password file</description>
</property>
```

If the Java KeyStore is *not* password protected, omit the `hadoop.security.keystore.java-keystore-provider.password-file` property.

4. Click **Save Changes** and restart the Key Trustee KMS service. If the Java KeyStore is *not* password protected, skip to step 7.
5. Create the file `/var/lib/keytrustee-kms/tomcat-deployment/webapps/kms/WEB-INF/classes/tmp/password.txt` and add the Java KeyStore password to it.
6. Change the ownership of `/var/lib/keytrustee-kms/tomcat-deployment/webapps/kms/WEB-INF/classes/tmp/password.txt` to `kms:kms`:

```
$ sudo chown kms:kms
/var/lib/keytrustee-kms/tomcat-deployment/webapps/kms/WEB-INF/classes/tmp/password.txt
```

7. From the host running the Key Trustee KMS service, if you have not configured Kerberos and TLS/SSL, run the following command:

```
$ curl -L -d "trusteeOp=migrate"
"http://kms01.example.com:16000/kms/v1/trustee/key/migrate?user.name=username&trusteeOp=migrate"
```

If you have configured Kerberos and TLS/SSL, use the following command instead:

```
$ curl --negotiate -u : -L -d "trusteeOp=migrate"
"https://kms01.example.com:16000/kms/v1/trustee/key/migrate?user.name=username&trusteeOp=migrate"
--cacert /path/to/kms/cert
```

8. Monitor `/var/log/kms-keytrustee/kms.log` and `/var/log/kms-keytrustee/kms-catalina.<date>.log` to verify that the migration is successful. You can also run `sudo -u <key_admin> hadoop key list` to verify that the keys are listed.
9. After you have verified that the migration is successful, remove the safety valve entry used in step 3 and restart the Key Trustee KMS service.

Configuring CDH Services for HDFS Encryption

The following topics contain recommendations for setting up [HDFS Transparent Encryption](#) on page 261 with various CDH services.

Hive

HDFS encryption has been designed so that files cannot be moved from one encryption zone to another or from encryption zones to unencrypted directories. Therefore, the landing zone for data when using the `LOAD DATA INPATH` command must always be inside the destination encryption zone.

To use HDFS encryption with Hive, ensure you are using *one* of the following configurations:

Single Encryption Zone

With this configuration, you can use HDFS encryption by having all Hive data inside the same encryption zone. In Cloudera Manager, configure the Hive Scratch Directory (`hive.exec.scratchdir`) to be inside the encryption zone.

Recommended HDFS Path: `/user/hive`

To use the auto-generated [KMS ACLs](#), make sure you name the encryption key `hive-key`.

For example, to configure a single encryption zone for the entire Hive warehouse, you can rename `/user/hive` to `/user/hive-old`, create an encryption zone at `/user/hive`, and then `distcp` all the data from `/user/hive-old` to `/user/hive`.

In Cloudera Manager, configure the Hive Scratch Directory (`hive.exec.scratchdir`) to be inside the encryption zone by setting it to `/user/hive/tmp`, ensuring that permissions are `1777` on `/user/hive/tmp`.

Multiple Encryption Zones

With this configuration, you can use encrypted databases or tables with different encryption keys. To read data from read-only encrypted tables, users must have access to a temporary directory that is encrypted at least as strongly as the table.

For example:

1. Configure two encrypted tables, `ezTbl1` and `ezTbl2`.
2. Create two new encryption zones, `/data/ezTbl1` and `/data/ezTbl2`.
3. Load data to the tables in Hive using `LOAD` statements.

For more information, see [Changed Behavior after HDFS Encryption is Enabled](#) on page 291.

Other Encrypted Directories

- **LOCALSCRATCHDIR:** The MapJoin optimization in Hive writes HDFS tables to a local directory and then uploads them to the distributed cache. To ensure these files are encrypted, either disable MapJoin by setting `hive.auto.convert.join` to `false`, or encrypt the *local* Hive Scratch directory (`hive.exec.local.scratchdir`) using [Cloudera Navigator Encrypt](#).
- **DOWNLOADED_RESOURCES_DIR:** JARs that are added to a user session and stored in HDFS are downloaded to `hive.downloaded.resources.dir` on the HiveServer2 local filesystem. To encrypt these JAR files, configure [Cloudera Navigator Encrypt](#) to encrypt the directory specified by `hive.downloaded.resources.dir`.
- **NodeManager Local Directory List:** Hive stores JARs and MapJoin files in the distributed cache. To use MapJoin or encrypt JARs and other resource files, the `yarn.nodemanager.local-dirs` YARN configuration property must be configured to a set of encrypted local directories on all nodes.

Changed Behavior after HDFS Encryption is Enabled

- Loading data from one encryption zone to another results in a copy of the data. `Distcp` is used to speed up the process if the size of the files being copied is higher than the value specified by `HIVE_EXEC_COPYFILE_MAXSIZE`. The minimum size limit for `HIVE_EXEC_COPYFILE_MAXSIZE` is 32 MB, which you can modify by changing the value for the `hive.exec.copyfile.maxsize` configuration property.
- When loading data to encrypted tables, Cloudera strongly recommends using a landing zone inside the same encryption zone as the table.
 - **Example 1: Loading unencrypted data to an encrypted table** - Use one of the following methods:
 - If you are loading new unencrypted data to an encrypted table, use the `LOAD DATA ...` statement. Because the source data is not inside the encryption zone, the `LOAD` statement results in a copy. For this reason, Cloudera recommends landing data that you need to encrypt inside the destination encryption zone. You can use `distcp` to speed up the copying process if your data is inside HDFS.
 - If the data to be loaded is already inside a Hive table, you can create a new table with a `LOCATION` inside an encryption zone as follows:

```
CREATE TABLE encrypted_table [STORED AS] LOCATION ... AS SELECT * FROM <unencrypted_table>
```

Configuring Encryption

The location specified in the `CREATE TABLE` statement must be inside an encryption zone. Creating a table pointing `LOCATION` to an unencrypted directory does not encrypt your source data. You must copy your data to an encryption zone, and then point `LOCATION` to that zone.

- **Example 2: Loading encrypted data to an encrypted table** - If the data is already encrypted, use the `CREATE TABLE` statement pointing `LOCATION` to the encrypted source directory containing the data. This is the fastest way to create encrypted tables.

```
CREATE TABLE encrypted_table [STORED AS] LOCATION ... AS SELECT * FROM
<encrypted_source_directory>
```

- Users reading data from encrypted tables that are read-only must have access to a temporary directory which is encrypted with at least as strong encryption as the table.
- Temporary data is now written to a directory named `.hive-staging` in each table or partition
- Previously, an `INSERT OVERWRITE` on a partitioned table inherited permissions for new data from the existing partition directory. With encryption enabled, permissions are inherited from the table.

KMS ACL Configuration for Hive

When Hive joins tables, it compares the encryption key strength for each table. For this operation to succeed, you must configure the [KMS ACLs](#) to allow the `hive` user and group `READ` access to the Hive key:

```
<property>
  <name>key.acl.hive-key.READ</name>
  <value>hive hive</value>
</property>
```

If you have restricted access to the `GET_METADATA` operation, you must grant permission for it to the `hive` user or group:

```
<property>
  <name>hadoop.kms.acl.GET_METADATA</name>
  <value>hive hive</value>
</property>
```

If you have disabled [HiveServer2 Impersonation](#) on page 138 (for example, to use [Apache Sentry](#)), you must configure the KMS ACLs to grant `DECRYPT_EEK` permissions to the `hive` user, as well as any user accessing data in the Hive warehouse.

Cloudera recommends creating a group containing all Hive users, and granting `DECRYPT_EEK` access to that group.

For example, suppose user `jdoue` (home directory `/user/jdoue`) is a Hive user and a member of the group `hive-users`. The encryption zone (EZ) key for `/user/jdoue` is named `jdoue-key`, and the EZ key for `/user/hive` is `hive-key`. The following ACL example demonstrates the required permissions:

```
<property>
  <name>key.acl.hive-key.DECRYPT_EEK</name>
  <value>hive hive-users</value>
</property>

<property>
  <name>key.acl.jdoue-key.DECRYPT_EEK</name>
  <value>jdoue,hive</value>
</property>
```

If you have enabled HiveServer2 impersonation, data is accessed by the user submitting the query or job, and the user account (`jdoue` in this example) may still need to access data in their home directory. In this scenario, the required permissions are as follows:

```
<property>
  <name>key.acl.hive-key.DECRYPT_EEK</name>
  <value>nobody hive-users</value>
```

```

</property>
<property>
  <name>key.acl.jdoe-key.DECRYPT_EEK</name>
  <value>jdoe</value>
</property>

```

Impala

Recommendations

- If HDFS encryption is enabled, configure Impala to encrypt data spilled to local disk.
- In releases lower than Impala 2.2.0 / CDH 5.4.0, Impala does not support the `LOAD DATA` statement when the source and destination are in different encryption zones. If you are running an affected release and need to use `LOAD DATA` with HDFS encryption enabled, copy the data to the table's encryption zone prior to running the statement.
- Use Cloudera Navigator to lock down the local directory where Impala UDFs are copied during execution. By default, Impala copies UDFs into `/tmp`, and you can configure this location through the `--local_library_dir` startup flag for the `impalad` daemon.
- Limit the rename operations for internal tables once encryption zones are set up. Impala cannot do an `ALTER TABLE RENAME` operation to move an internal table from one database to another, if the root directories for those databases are in different encryption zones. If the encryption zone covers a table directory but not the parent directory associated with the database, Impala cannot do an `ALTER TABLE RENAME` operation to rename an internal table, even within the same database.
- Avoid structuring partitioned tables where different partitions reside in different encryption zones, or where any partitions reside in an encryption zone that is different from the root directory for the table. Impala cannot do an `INSERT` operation into any partition that is not in the same encryption zone as the root directory of the overall table.
- If the data files for a table or partition are in a different encryption zone than the HDFS trashcan, use the `PURGE` keyword at the end of the `DROP TABLE` or `ALTER TABLE DROP PARTITION` statement to delete the HDFS data files immediately. Otherwise, the data files are left behind if they cannot be moved to the trashcan because of differing encryption zones. This syntax is available in Impala 2.3 / CDH 5.5 and higher.

Steps

Start every `impalad` process with the `--disk_spill_encryption=true` flag set. This encrypts all spilled data using AES-256-CFB. Set this flag by selecting the **Disk Spill Encryption** checkbox in the Impala configuration (**Impala service > Configuration > Category > Security**).



Important: Impala does not selectively encrypt data based on whether the source data is already encrypted in HDFS. This results in at most 15 percent performance degradation when data is spilled.

KMS ACL Configuration for Impala

Cloudera recommends making the `impala` user a member of the `hive` group, and following the ACL recommendations in [KMS ACL Configuration for Hive](#) on page 292.

HBase

Recommendations

Make `/hbase` an encryption zone. Do not create encryption zones as subdirectories under `/hbase`, because HBase may need to rename files across those subdirectories. When you create the encryption zone, name the key `hbase-key` for use with the [Recommended KMS Access Control List](#) on page 283.

Configuring Encryption

Steps

On a cluster without HBase currently installed, create the `/hbase` directory and make that an encryption zone.

On a cluster with HBase already installed, perform the following steps:

1. Stop the HBase service.
2. Move data from the `/hbase` directory to `/hbase-tmp`.
3. Create an empty `/hbase` directory and make it an encryption zone.
4. Distcp all data from `/hbase-tmp` to `/hbase`, preserving user-group permissions and extended attributes.
5. Start the HBase service and verify that it is working as expected.
6. Remove the `/hbase-tmp` directory.

KMS ACL Configuration for HBase

In the [KMS ACLs](#), grant the `hbase` user and group `DECRYPT_EEK` permission for the HBase key:

```
<property>
  <name>key.acl.hbase-key.DECRYPT_EEK</name>
  <value>hbase hbase</value>
  </description>
</property>
```

Search

Recommendations

Make `/solr` an encryption zone. When you create the encryption zone, name the key `solr-key` for use with the [Recommended KMS Access Control List](#) on page 283.

Steps

On a cluster without Solr currently installed, create the `/solr` directory and make that an encryption zone.

On a cluster with Solr already installed:

1. Create an empty `/solr-tmp` directory.
2. Make `/solr-tmp` an encryption zone.
3. DistCp all data from `/solr` into `/solr-tmp`.
4. Remove `/solr`, and rename `/solr-tmp` to `/solr`.

KMS ACL Configuration for Search

In the [KMS ACLs](#), grant the `solr` user and group `DECRYPT_EEK` permission for the Solr key:

```
<property>
  <name>key.acl.solr-key.DECRYPT_EEK</name>
  <value>solr solr</value>
  </description>
</property>
```

Sqoop

Recommendations

- **For Hive support:** Ensure that you are using Sqoop with the `--target-dir` parameter set to a directory that is inside the Hive encryption zone. For more details, see [Hive](#) on page 290.
- **For append/incremental support:** Make sure that the `sqoop.test.import.rootDir` property points to the same encryption zone as the `--target-dir` argument.
- **For HCatalog support:** No special configuration is required.

Hue

Recommendations

Make `/user/hue` an encryption zone because Oozie workflows and other Hue-specific data are stored there by default. When you create the encryption zone, name the key `hue-key` for use with the [Recommended KMS Access Control List](#) on page 283.

Steps

On a cluster without Hue currently installed, create the `/user/hue` directory and make it an encryption zone.

On a cluster with Hue already installed:

1. Create an empty `/user/hue-tmp` directory.
2. Make `/user/hue-tmp` an encryption zone.
3. DistCp all data from `/user/hue` into `/user/hue-tmp`.
4. Remove `/user/hue` and rename `/user/hue-tmp` to `/user/hue`.

KMS ACL Configuration for Hue

In the [KMS ACLs](#), grant the `hue` and `oozie` users and groups `DECRYPT_EEK` permission for the Hue key:

```
<property>
  <name>key.acl.hue-key.DECRYPT_EEK</name>
  <value>oozie,hue oozie,hue</value>
</property>
```

Spark

Recommendations

- By default, application event logs are stored at `/user/spark/applicationHistory`, which can be made into an encryption zone.
- Spark also optionally caches its JAR file at `/user/spark/share/lib` (by default), but encrypting this directory is not required.
- Spark does not encrypt shuffle data. To do so, configure the Spark local directory, `spark.local.dir` (in Standalone mode), to reside on an encrypted disk. For YARN mode, make the corresponding YARN configuration changes.

KMS ACL Configuration for Spark

In the [KMS ACLs](#), grant `DECRYPT_EEK` permission for the Spark key to the `spark` user and any groups that can submit Spark jobs:

```
<property>
  <name>key.acl.spark-key.DECRYPT_EEK</name>
  <value>spark spark-users</value>
</property>
```

MapReduce and YARN

MapReduce v1

Recommendations

MRv1 stores both history and logs on local disks by default. Even if you do configure history to be stored on HDFS, the files are not renamed. Hence, no special configuration is required.

Configuring Encryption

MapReduce v2 (YARN)

Recommendations

Make `/user/history` a single encryption zone, because history files are moved between the `intermediate` and `done` directories, and HDFS encryption does not allow moving encrypted files across encryption zones. When you create the encryption zone, name the key `mapred-key` to take advantage of auto-generated [KMS ACLs](#).

Steps

On a cluster with MRv2 (YARN) installed, create the `/user/history` directory and make that an encryption zone.

If `/user/history` already exists and is not empty:

1. Create an empty `/user/history-tmp` directory.
2. Make `/user/history-tmp` an encryption zone.
3. DistCp all data from `/user/history` into `/user/history-tmp`.
4. Remove `/user/history` and rename `/user/history-tmp` to `/user/history`.

KMS ACL Configuration for MapReduce

In the [KMS ACLs](#), grant `DECRYPT_EEK` permission for the MapReduce key to the `mapred` and `yarn` users and the `hadoop` group:

```
<property>
  <name>key.acl.mapred-key.DECRYPT_EEK</name>
  <value>mapred,yarn hadoop</value>
  </description>
</property>
```

Troubleshooting HDFS Encryption

This topic contains HDFS Encryption-specific troubleshooting information in the form of issues you might face when encrypting HDFS files/directories and their workarounds.

Retrieval of encryption keys fails

Description

You see the following error when trying to list encryption keys

```
user1@example-sles-4:~> hadoop key list
Cannot list keys for KeyProvider: KMSClientProvider[https://example-sles-2.example.com:16000/kms/v1/]: Retrieval of all keys failed.
```

Solution

Make sure your truststore has been updated with the relevant certificate(s), such as the Key Trustee server certificate.

DistCp between unencrypted and encrypted locations fails

Description

By default, DistCp compares checksums provided by the filesystem to verify that data was successfully copied to the destination. However, when copying between unencrypted and encrypted locations, the filesystem checksums will not match since the underlying block data is different.

Solution

Specify the `-skipcrccheck` and `-update distcp` flags to avoid verifying checksums.

Cannot move encrypted files to trash

Description

With HDFS encryption enabled, you cannot move encrypted files or directories to the trash directory.

Solution

To remove encrypted files/directories, use the following command with the `-skipTrash` flag specified to bypass trash.

```
rm -r -skipTrash /testdir
```

NameNode - KMS communication fails after long periods of inactivity

Description

Encrypted files and encryption zones cannot be created if a long period of time (by default, 20 hours) has passed since the last time the KMS and NameNode communicated.

Solution



Important: Upgrading your cluster to the latest CDH 5 release will fix this problem. For instructions, see [Upgrading from an Earlier CDH 5 Release to the Latest Release](#).

For lower CDH 5 releases, there are two possible workarounds to this issue :

- You can increase the KMS authentication token validity period to a very high number. Since the default value is 10 hours, this bug will only be encountered after 20 hours of no communication between the NameNode and the KMS. Add the following property to the `kms-site.xml` Safety Valve:

```
<property>
<name>hadoop.kms.authentication.token.validity</name>
<value>SOME VERY HIGH NUMBER</value>
</property>
```

- You can switch the KMS signature secret provider to the string secret provider by adding the following property to the `kms-site.xml` Safety Valve:

```
<property>
<name>hadoop.kms.authentication.signature.secret</name>
<value>SOME VERY SECRET STRING</value>
</property>
```

Cloudera Navigator Key Trustee Server

Cloudera Navigator Key Trustee Server is an enterprise-grade virtual safe-deposit box that stores and manages cryptographic keys and other security artifacts. With Navigator Key Trustee Server, encryption keys are separated from the encrypted data, ensuring that sensitive data is still protected in the event that unauthorized users gain access to the storage media.

Key Trustee Server protects these keys and other critical security objects from unauthorized access while enabling compliance with strict data security regulations. For added security, Key Trustee Server can integrate with a hardware security module (HSM). See [Cloudera Navigator Key HSM](#) on page 309 for more information.

In conjunction with the Key Trustee KMS, Navigator Key Trustee Server can serve as a backing key store for [HDFS Transparent Encryption](#) on page 261, providing enhanced security and scalability over the file-based Java KeyStore used by the default Hadoop Key Management Server.

[Cloudera Navigator Encrypt](#) on page 316 also uses Key Trustee Server for key storage and management.

For instructions on installing Key Trustee Server, see [Installing Cloudera Navigator Key Trustee Server](#).

For instructions on configuring Key Trustee Server, continue reading:

Backing Up and Restoring Key Trustee Server

Key Trustee Server high availability applies to read operations only. If either Key Trustee Server fails, the client automatically retries fetching keys from the functioning server. New write operations (for example, creating new encryption keys) are not allowed unless both Key Trustee Servers are operational.

If a Key Trustee Server fails catastrophically, you must restore it from backup to a new host with the same hostname and IP address as the failed host. Cloudera does not support PostgreSQL promotion to convert a passive Key Trustee Server to an active Key Trustee Server.

Cloudera strongly recommends regularly backing up Key Trustee Server databases and configuration files. Because these backups contain encryption keys and encrypted deposits, you must ensure that your backup repository is as secure as the Key Trustee Server.

You must also back up client configuration files and keys for Key Trustee Server clients, such as Key Trustee KMS and Navigator Encrypt clients.



Note: In an HA configuration, the backup need only be performed on one of the hosts for Key Trustee Server and the Key Trustee KMS. For Key Trustee Server, run the backup on the *active* server. For Key Trustee KMS, you can run the backup on any instance.

Backing Up Key Trustee Server

Use this procedure for both parcel-based and package-based installations.

If you have deployed [Cloudera Navigator Key Trustee Server High Availability](#), perform these steps on both the active and passive Key Trustee Servers. The following procedure references the default database port and location; if you modified these settings during installation, replace the database and port with your values.

1. Back up the Key Trustee Server database:

- For Key Trustee Server 3.8:

```
$ su - postgres
$ pg_dump -c -p 5432 keytrustee | zip --encrypt keytrustee-db.zip -
```

- For Key Trustee Server 5.4 and higher:

```
$ su - keytrustee
$ pg_dump -c -p 11381 keytrustee | zip --encrypt keytrustee-db.zip -
```

The `--encrypt` option prompts you to create a password used to encrypt the zip file. This password is required to decrypt the file.

For parcel-based installations, you must set environment variables after switching to the `keytrustee` user:

```
$ su - keytrustee
$ export PATH=$PATH:/opt/cloudera/parcels/KEYTRUSTEE_SERVER/PG_DB/opt/postgres/9.3/bin
$ export
LD_LIBRARY_PATH=/opt/cloudera/parcels/KEYTRUSTEE_SERVER/PG_DB/opt/postgres/9.3/lib
$ pg_dump -c -p 11381 keytrustee | zip --encrypt keytrustee-db.zip -
```

2. Back up the Key Trustee Server configuration directory (`/var/lib/keytrustee/.keytrustee`):

```
$ zip -r --encrypt keytrustee-conf.zip /var/lib/keytrustee/.keytrustee
```

The `--encrypt` option prompts you to create a password used to encrypt the zip file. This password is required to decrypt the file.

3. Move the backup files (`keytrustee-db.zip` and `keytrustee-conf.zip`) to a secure location.

Backing Up Key Trustee Server Clients

Cryptographic keys stored in Key Trustee Server are encrypted by clients before they are sent to Key Trustee Server. The primary clients for Key Trustee Server are Key Trustee KMS and Navigator Encrypt. Cloudera strongly recommends backing up regularly the configuration files and GPG keys for Key Trustee Server clients.



Warning: Failure to back up these files can result in irretrievable data loss. For example, encryption zone keys used for [HDFS Transparent Encryption](#) on page 261 are encrypted by the KMS before being stored in Key Trustee Server. A catastrophic failure of the KMS with no backup causes all HDFS data stored in encryption zones to become permanently irretrievable.

To prevent permanent data loss, regularly back up the following directories on each client that stores objects in Key Trustee Server:

Table 22: Key Trustee Server Client Configuration Directories

Key Trustee Server Client	Directories to Back Up
Key Trustee KMS	<code>/var/lib/kms-keytrustee</code>
Navigator Encrypt	<code>/etc/navencrypt</code>

Restoring Key Trustee Server

When restoring the Key Trustee Server database from backup, keep in mind that any keys or deposits created after the backup are not restored. If you are using Key Trustee Server high availability, you can restore the Active Key Trustee Server from the Passive Key Trustee Server. This restores all keys that were successfully written to the Passive Key Trustee Server before the failure.

The procedure to restore Key Trustee Server is different for parcel-based than for package-based installations. For more information about parcels, see [Parcels](#).

Restoring Key Trustee Server in Parcel-Based Installations



Note: These instructions apply to Key Trustee Servers deployed using parcels. For package-based deployments, skip to the [Restoring Key Trustee Server in Package-Based Installations](#) on page 300 section.

If you have deployed [Cloudera Navigator Key Trustee Server High Availability](#), perform these steps on both the active and passive Key Trustee Servers. The following procedures assume the default database port and location; if you modified these settings during installation, replace the database and port with your custom values.

If the Key Trustee Server host has failed completely, remove the host from the cluster and add a new host using Cloudera Manager:

1. Remove the failed host from the cluster. See [Deleting Hosts](#) for instructions.
2. Add a new host with the same hostname and IP address as the failed host to the cluster. See [Adding a Host to the Cluster](#) for instructions.



Important: Make sure that the replacement host uses the same operating system version as the failed host.

3. Install Key Trustee Server on the new host. See [Installing Cloudera Navigator Key Trustee Server](#) for instructions. Make sure to install the same Key Trustee Server version as the failed host.

After you have provisioned a new host and installed Key Trustee Server (or if you are restoring the database or configuration on the original host), restore the database and configuration directory. If you need to restore the Active Key Trustee Server from the Passive Key Trustee Server, skip to [Restoring Active Key Trustee Server from Passive Key Trustee Server](#) on page 301.

Configuring Encryption

If your backups were created manually using the `pg_dump` command, do the following:

1. Copy or move the backup files (*keytrustee-db.zip* and *keytrustee-conf.zip*) to the Key Trustee Server host.
2. Start the PostgreSQL server:

```
$ sudo ktadmin db --start --pg-rootdir /var/lib/keytrustee/db --background
```

3. Restore the Key Trustee Server database:

```
$ su - keytrustee
$ export PATH=$PATH:/opt/cloudera/parcels/KEYTRUSTEE_SERVER/PG_DB/opt/postgres/9.3/bin
$ export
LD_LIBRARY_PATH=/opt/cloudera/parcels/KEYTRUSTEE_SERVER/PG_DB/opt/postgres/9.3/lib
$ unzip -p /path/to/keytrustee-db.zip | psql -p 11381 -d keytrustee
```

If the zip file is encrypted, you are prompted for the password to decrypt the file.

4. Restore the Key Trustee Server configuration directory:

```
$ su - keytrustee
$ cd /var/lib/keytrustee
$ unzip /path/to/keytrustee-conf.zip
```

If the zip file is encrypted, you are prompted for the password to decrypt the file.

5. Stop the PostgreSQL server:

```
$ sudo ktadmin db --stop --pg-rootdir /var/lib/keytrustee/db
```

6. Start the Key Trustee Server service in Cloudera Manager (**Key Trustee Server service** > **Actions** > **Start**).
7. Remove the backup files (*keytrustee-db.zip* and *keytrustee-conf.zip*) from the Key Trustee Server host.

Restoring Key Trustee Server in Package-Based Installations

If you have deployed [Cloudera Navigator Key Trustee Server High Availability](#), perform these steps on both the active and passive Key Trustee Servers. The following procedures assume the default database port and location; if you modified these settings during installation, replace the database and port with your custom values.

If the Key Trustee Server host has failed completely, provision a new host with the same hostname and IP address as the failed host, and re-install Key Trustee Server. See [Installing Cloudera Navigator Key Trustee Server](#) for instructions.



Important: Make sure to install the same operating system and Key Trustee Server versions as the failed host.

After you have provisioned a new host and installed Key Trustee Server (or if you are restoring the database or configuration on the original host), restore the database and configuration directory. If you need to restore the Active Key Trustee Server from the Passive Key Trustee Server, skip to [Restoring Active Key Trustee Server from Passive Key Trustee Server](#) on page 301.

If your backups were created manually using the `pg_dump` command, do the following:

1. Copy or move the backup files (*keytrustee-db.zip* and *keytrustee-conf.zip*) to the Key Trustee Server host.
2. Change the file ownership on the backup files to `keytrustee:keytrustee`:

```
$ sudo chown keytrustee:keytrustee /path/to/keytrustee*.zip
```

3. Restore the Key Trustee Server database:

```
$ su - keytrustee
$ unzip -p /path/to/keytrustee-db.zip | psql -p 11381 -d keytrustee
```

If the zip file is encrypted, you are prompted for the password to decrypt the file.

4. Restore the Key Trustee Server configuration directory:

```
$ cd /var/lib/keytrustee
$ unzip /path/to/keytrustee-conf.zip
```

If the zip file is encrypted, you are prompted for the password to decrypt the file.

5. Start the Key Trustee Server service:

- RHEL 6-compatible: \$ sudo service keytrusteed start
- RHEL 7-compatible: \$ sudo systemctl start keytrusteed

6. Remove the backup files (*keytrustee-db.zip* and *keytrustee-conf.zip*) from the Key Trustee Server host.

Restoring Active Key Trustee Server from Passive Key Trustee Server

If the Active Key Trustee Server fails, and you do not have a backup, you can restore it from the Passive Key Trustee Server using the following procedure. You can also use this procedure if you need to restore keys that were successfully written to the Passive Key Trustee Server, but are not included in the most recent backup.

The following procedure assumes you have installed Key Trustee Server on the replacement host and (if you are using Cloudera Manager) added the Key Trustee Server service. For instructions on installing Key Trustee Server, see [Installing Cloudera Navigator Key Trustee Server](#)

1. Copy the Key Trustee Server database from the Passive Key Trustee Server host to the new Active Key Trustee Server host. Run the following command on the Passive Key Trustee Server host:

```
$ sudo rsync --exclude recovery.conf -a /var/lib/keytrustee/db
root@kts01.example.com:/var/lib/keytrustee/
```

Replace *kts01.example.com* with the hostname of the new Active Key Trustee Server.

2. Make sure that the *recovery.conf* file did not get copied to the Active Key Trustee Server (for example, if there was a typo in your *rsync* command). Run the following command on the Active Key Trustee Server host:

```
$ sudo ls -l /var/lib/keytrustee/db/recovery.conf
```

If the file exists on the Active Key Trustee Server host, delete it. Make sure you are on the Active Key Trustee Server host before deleting the file. Do not delete the *recovery.conf* file on the Passive Key Trustee Server host.

3. Copy the configuration directory from the Passive Key Trustee Server host to the new Active Key Trustee Server host. Run the following command on the Passive Key Trustee Server host:

```
$ sudo rsync --exclude .ssl --exclude '*.pid' -a /var/lib/keytrustee/.keytrustee
root@kts01.example.com:/var/lib/keytrustee/
```

Replace *kts01.example.com* with the hostname of the new Active Key Trustee Server.

4. Create the *logs* directory and make sure it is owned by the *keytrustee* user and group:

```
$ sudo mkdir /var/lib/keytrustee/logs
$ sudo chown keytrustee:keytrustee /var/lib/keytrustee/logs
```

5. (Cloudera Manager only) Generate the Key Trustee Server keyring: **Key Trustee Server service > Actions > Generate Key Trustee Server Keyring**

6. Set up the database on the Active Key Trustee Server host.

- **Using Cloudera Manager: Key Trustee Server service > Actions > Set Up Key Trustee Server Database**
- **Using the Command Line:**

```
$ sudo ktadmin --confdir /var/lib/keytrustee db --port 11381 --pg-rootdir
/var/lib/keytrustee/db --bootstrap --slave kts02.example.com
```

Configuring Encryption

Replace `kts02.example.com` with the hostname of the Passive Key Trustee Server.

7. Start the database.

- **Using Cloudera Manager: Key Trustee Server service > Instances > Active Database > Actions > Start this Active Database**
- **Using the Command Line:** Run the following command on the Active Key Trustee Server host:

```
$ sudo ktadmin --confdir /var/lib/keytrustee db --port 11381 --pg-rootdir /var/lib/keytrustee/db --bootstrap --slave kts02.example.com
```

Replace `kts02.example.com` with the hostname of the Passive Key Trustee Server.

8. Enable synchronous replication.

- **Using Cloudera Manager: Key Trustee Server service > Actions > Setup Enable Synchronous Replication in HA mode**
- **Using the Command Line:** Run the following command on the Active Key Trustee Server host:

```
$ sudo ktadmin --confdir /var/lib/keytrustee enable-synchronous-replication
```

9. Restart the Active Key Trustee Server.

- **Using Cloudera Manager: Key Trustee Server service > Actions > Restart**
- **Using the Command Line:** Run the following command on the Active Key Trustee Server host:

```
$ sudo service keytrusteed restart #RHEL 6-compatible  
$ sudo systemctl restart keytrusteed #RHEL 7-compatible
```

Initializing Standalone Key Trustee Server

If you are configuring high availability Key Trustee Servers, skip this step and proceed to [Cloudera Navigator Key Trustee Server High Availability](#).

Using Cloudera Manager



Note: These instructions apply to using Cloudera Manager only. For package-based deployments, skip to the [Using the Command Line](#) on page 302 section.



Important: If you are using SSH software other than OpenSSH, pre-create the SSH key before continuing:

```
$ sudo -u keytrustee ssh-keygen -t rsa -f /var/lib/keytrustee/.ssh/id_rsa
```

Add the Key Trustee Server service to your cluster, following the instructions in [Adding a Service](#). When customizing role assignments, assign only the Active Key Trustee Server and Active Database roles.



Important: You *must* assign the Key Trustee Server and Database roles to the same host.

Using the Command Line

To initialize a standalone Key Trustee Server, run the following commands on the Key Trustee Server:



Important: For Key Trustee Server 5.4.0 and higher, the `ktadmin init-master` command is deprecated. Use the `ktadmin init` command instead. If you are using SSH software other than OpenSSH, pre-create the SSH key before continuing:

```
$ sudo -u keytrustee ssh-keygen -t rsa /var/lib/keytrustee/.ssh/id_rsa
```

```
$ sudo ktadmin init --external-address keytrustee.example.com
$ sudo ktadmin db --bootstrap --port 11381 --pg-rootdir /var/lib/keytrustee/db
## For RHEL/CentOS 7, use 'sudo systemctl [stop|start] <service_name>' instead of 'sudo
service <service_name> [stop|start]' ##
$ sudo service keytrustee-db stop
$ sudo service keytrustee-db start
$ sudo service keytrusteed start
$ sudo chkconfig keytrustee-db on
$ sudo chkconfig keytrusteed on
```

Replace `keytrustee.example.com` with the fully-qualified domain name (FQDN) of the Key Trustee Server. Cloudera recommends using the default `/var/lib/keytrustee/db` directory for the PostgreSQL database.

To use a different port for the database, modify the `ktadmin init` and `ktadmin db` commands as follows:

```
$ sudo ktadmin init --external-address keytrustee.example.com --db-connect
postgresql://localhost:<port>/keytrustee?host=/tmp
$ sudo ktadmin db --bootstrap --port <port> --pg-rootdir /var/lib/keytrustee/db
```

If you specify a database directory other than `/var/lib/keytrustee/db`, create or edit the `/etc/sysconfig/keytrustee-db` file and add the following line:

```
ARGS="--pg-rootdir /path/to/db"
```

The `ktadmin init` command initializes the Key Trustee configuration directory (`/var/lib/keytrustee/.keytrustee` by default) and generates a self-signed certificate that Key Trustee Server uses for HTTPS communication.

The `ktadmin db --bootstrap` command initializes the database in the directory specified by the `--pg-rootdir` parameter.

The `sudo service keytrustee-db stop` and `sudo service keytrustee-db start` commands restart the Key Trustee Server database.

The `sudo service keytrusteed start` command starts Key Trustee Server.



Note: The `/etc/init.d/postgresql` script does not work when the PostgreSQL database is started by Key Trustee Server, and cannot be used to monitor the status of the database. Use `/etc/init.d/keytrustee-db` instead.

(Optional) Replace Self-Signed Certificate with CA-Signed Certificate



Important: Key Trustee Server certificates must be issued to the fully-qualified domain name (FQDN) of the Key Trustee Server host. If you are using CA-signed certificates, ensure that the generated certificates use the FQDN, and not the short name.

If you have a CA-signed certificate for Key Trustee Server, see [Managing Key Trustee Server Certificates](#) on page 307 for instructions on how to replace the self-signed certificate.

Configuring a Mail Transfer Agent for Key Trustee Server

The Key Trustee Server requires a mail transfer agent (MTA) to send email. Cloudera recommends Postfix, but you can use any MTA that meets your needs.

To configure Postfix for local delivery, run the following commands:

```
export
KEYTRUSTEE_SERVER_PK="/var/lib/keytrustee/.keytrustee/.ssl/ssl-cert-keytrustee-pk.pem"
export
KEYTRUSTEE_SERVER_CERT="/var/lib/keytrustee/.keytrustee/.ssl/ssl-cert-keytrustee.pem"
export
KEYTRUSTEE_SERVER_CA="/var/lib/keytrustee/.keytrustee/.ssl/ssl-cert-keytrustee-ca.pem"
export KEYTRUSTEE_SERVER_HOSTNAME="$(hostname -f)" # or adjust as required
postconf -e 'mailbox_command ='
postconf -e 'smtpd_sasl_local_domain ='
postconf -e 'smtpd_sasl_auth_enable = yes'
postconf -e 'smtpd_sasl_security_options = noanonymous'
postconf -e 'broken_sasl_auth_clients = yes'
postconf -e 'smtpd_recipient_restrictions =
permit_sasl_authenticated,permit_mynetworks,reject_unauth_destination'
postconf -e 'inet_interfaces = localhost'
postconf -e 'smtp_tls_security_level = may'
postconf -e 'smtpd_tls_security_level = may'
```

Start the Postfix service and ensure that it starts at boot:

```
$ service postfix restart
$ sudo chkconfig --level 235 postfix on
```

For information on installing Postfix or configuring a relay host, see the [Postfix documentation](#).

Verifying Cloudera Navigator Key Trustee Server Operations

Verify that the installation was successful by running the following command on all Key Trustee Servers.

```
curl -k https://keytrustee.example.com:11371/?a=fingerprint
```

Replace `keytrustee.example.com` with the fully qualified domain name (FQDN) of each Key Trustee Server you are validating. This command outputs a fingerprint similar to the following:

```
4096R/4EDC46882386C827E20DEEA2D850ACA33BEDB0D1
```

If high availability is enabled, the output should be identical for all Key Trustee Servers.

Managing Key Trustee Server Organizations

Organizations allow you to configure Key Trustee for use in a multi-tenant environment. Using the `keytrustee-orgtool` utility, you can create organizations and administrators for multiple organizations. Organization administrators can then approve or deny the registration of clients, depending on the registration method.

The `keytrustee-orgtool` Utility

`keytrustee-orgtool` is a command-line utility for administering organizations. The `keytrustee-orgtool` command must be run as the root user.

The following table explains the various `keytrustee-orgtool` commands and parameters. Run `keytrustee-orgtool --help` to view this information at the command line.

Table 23: Usage for `keytrustee-orgtool`

Operation	Usage	Description
Add	<code>keytrustee-orgtool add [-h] -n name -c contacts</code>	Adds a new organization and administrators for the organization.
List	<code>keytrustee-orgtool list</code>	Lists current organizations, including the authorization secret, all administrators, the organization creation date, and the organization expiration date.
Disable client	<code>keytrustee-orgtool disable-client [-h] --fingerprint fingerprint</code>	Disables a client that has already been activated by the organization administrator.
Enable client	<code>keytrustee-orgtool enable-client [-h] --fingerprint fingerprint</code>	Enables a client that has requested activation but has not yet been approved by the organization administrator.
Set authorization Code	<code>keytrustee-orgtool set-auth [-h] -n name -s secret</code>	Sets the authorization code to a new string, or to blank to allow automatic approvals without the code.

Create Organizations

Each new Key Trustee tenant needs its own organization. You can create new organizations using the `keytrustee-orgtool add` command. For example, to create a new organization for the Disaster Recovery group and add two administrators, Finn and Jake:

```
$ keytrustee-orgtool add -n disaster-recov -c finn@example.com,jake@example.com
```

When adding organizations:

- Do not use spaces or special characters in the organization name. Use hyphens or underscores instead.
- Do not use spaces between email addresses (when adding multiple administrators to an organization). Use a comma to separate email addresses, without any space (as shown in the example above).

Each contact email address added when creating the organization receives a [notification email](#), as detailed below.

Once an organization exists, use the `keytrustee-orgtool add` command to add new administrators to the organization. For example, to add an administrator to the `disaster-recov` organization:

```
keytrustee-orgtool add -n disaster-recov -c marceline@example.com
```



Note: You cannot remove contacts from an organization with the `keytrustee-orgtool` utility.

List Organizations

After creating an organization, verify its existence with the `keytrustee-orgtool list` command. This command lists details for all existing organizations. The following is the entry for the `disaster-recov` organization created in the example:

```
"disaster-recov": {
  "auth_secret": "/qFiICsyYqMLhdTznNY3Nw==",
  "contacts": [
```

Configuring Encryption

```
        "finn@example.com",
        "jake@example.com"
    ],
    "creation": "2013-12-02T09:55:21",
    "expiration": "9999-12-31T15:59:59",
    "key_info": null,
    "name": "disaster-recov",
    "state": 0,
    "uuid": "xY3Z8xCwMuKZMiTYJa0mZOdhMVdxhyCUOc6vSNc9I8X"
}
```

Change the Authorization Code

When an organization is created, an authorization code is automatically generated. When you run the `keytrustee-orgtool list` command, the code is displayed in the `auth_secret` field. To register with a Key Trustee Server, the client must have the authorization code along with the organization name. To set a new `auth_secret`, run the following command:

```
$ keytrustee-orgtool set-auth -n disaster-recov -s ThisISAs3cr3t!
```

Run the `keytrustee-orgtool list` command again, and confirm the updated `auth_secret` field:

```
"disaster-recov": {
  "auth_secret": "ThisISAs3cr3t!",
  "contacts": [
    "finn@example.com",
    "jake@example.com"
  ],
  "creation": "2013-12-02T09:55:21",
  "expiration": "9999-12-31T15:59:59",
  "key_info": null,
  "name": "disaster-recov",
  "state": 0,
  "uuid": "xY3Z8xCwMuKZMiTYJa0mZOdhMVdxhyCUOc6vSNc9I8X"
}
```

If you do not want to use an authorization code, set the `auth_secret` field to an empty string:

```
$ keytrustee-orgtool set-auth -n disaster-recov -s ""
```

Cloudera recommends requiring an authorization code.

Notification Email and GPG Keys

Whenever an administrator is added to an organization, the Key Trustee Server sends an automated email message (subject: "KeyTrustee Contact Registration") to the newly added administrator:

```
Hello, this is an automated message from your Cloudera keytrustee Server.

Welcome to Cloudera keytrustee! You have been listed as an administrator contact
for keytrustee services at your organization [test-org]. As an administrator,
you may be contacted to authorize the activation of new keytrustee clients.

We recommend that you register a GPG public key for secure administration of
your clients. To do so, visit the link below and follow the instructions.

https://keytrustee01.example.com:11371/?q=CnRV6u0nbm7zB07BQEpXCXsN0QJFBz684uC01cHMoWL

This link will expire in 12 hours, at Thu Sep  3 00:08:25 2015 UTC.
```

Cloudera recommends that an organization's administrators:

- Register the GPG public key by following the link contained in the notification email. Registering the GPG public key is optional, but if you choose to register your public key:

- Complete the process within 12 hours, before the link expires.
- Upload the entire key, including the BEGIN and END strings, as shown here:

```
-----BEGIN PGP PUBLIC KEY BLOCK-----
Version: GnuPG v1.2.1 (GNU/Linux)
Comment: For info see http://www.gnupg.org

mQGIBDkHP3URBACKWGsYh43pkXU9wj/X1G67K8/DSr185r7dNtHNfLL/ewil10k2
q8ssaWJn26QZPsDVqdUJModHfJ6kQTAt9NzQbgcVrxLYNfgeBsvkHF/POtnYcZRgL
tZ6syBBWs8JB4xt5V09iJSGAMPUQE8Jpdn2aRXPAPdoDw179LM8Rq6r+gwCg5ZZa
.
.
.
-----END PGP PUBLIC KEY BLOCK-----
```

- Import the Key Trustee Server's public GPG key to verify that the server is the sender.

The organization's administrators are notified by email when new clients are registered to the Key Trustee Server using the mail transfer agent (as discussed in [Configuring a Mail Transfer Agent for Key Trustee Server](#) on page 304). However, if the server does not have access to email, you can use a local system mail address, such as `username@hostname`, where `hostname` is the system hostname and `username` is a valid system user. If you use a local system email address, be sure to regularly monitor the email box.

Managing Key Trustee Server Certificates

Transport Layer Security (TLS) certificates are used to secure communication with Key Trustee Server. By default, Key Trustee Server generates self-signed certificates when it is first initialized. Cloudera strongly recommends using certificates signed by a trusted Certificate Authority (CA).

Generating a New Certificate

1. Generate a new certificate signing request (CSR):

```
$ openssl req -new -key keytrustee_private_key.pem -out new.csr
```

Replace *keytrustee_private_key.pem* with the filename of the private key. You can reuse the existing private key or generate a new private key in accordance with your company policies. For existing auto-generated self-signed certificates, the private key file is located at `/var/lib/keytrustee/.keytrustee/.ssl/ssl-cert-keytrustee-pk.pem`.

2. Generate a new certificate from the CSR:

- For a CA-signed certificate, submit the CSR to the CA, and they will provide a signed certificate.
- To generate a new self-signed certificate, run the following command:

```
$ openssl x509 -req -days 365 -in new.csr -signkey keytrustee_private_key.pem \
-out new_keytrustee_certificate.pem
```

Replacing Key Trustee Server Certificates

Use the following procedure if you need to replace an existing certificate for the Key Trustee Server. For example, you can use this procedure to replace the auto-generated self-signed certificate with a CA-signed certificate, or to replace an expired certificate.



Note: Key Trustee Server supports password-protected private keys, but not password-protected certificates.

1. After [Generating a New Certificate](#) on page 307, back up the original certificate and key files:

```
$ sudo cp -r /var/lib/keytrustee/.keytrustee/.ssl /var/lib/keytrustee/.keytrustee/.ssl.bak
```

2. (CA-Signed Certificates Only) Provide either the root or intermediate CA certificate:



Important: If you have separate root CA and intermediate CA certificate files, then you must concatenate them into a single file. If you need to convert from a JKS with combined certificate and private key file to a PEM file and separate private key file, refer to [Conversion from Java Keystore to OpenSSL](#) on page 211.

```
$ sudo mv /path/to/rootca.pem
/var/lib/keytrustee/.keytrustee/.ssl/ssl-cert-keytrustee-ca.pem
```

3. Make sure that the certificate and key files are owned by the `keytrustee` user and group, with file permissions set to 600:

```
$ sudo chown keytrustee:keytrustee /path/to/new/certificate.pem
/path/to/new/private_key.pem
$ sudo chmod 600 /path/to/new/certificate.pem /path/to/new/private_key.pem
```

4. Update the Key Trustee Server configuration with the location of the new certificate and key files:

- Using Cloudera Manager:

- Go to the Key Trustee Server service.
- Click the **Configuration** tab.
- Select **Category > Security**.
- Edit the following properties to specify the location of the new certificate and key files. If the private keys are not password protected, leave the password fields empty.
 - Active Key Trustee Server TLS/SSL Server Private Key File (PEM Format)**
 - Active Key Trustee Server TLS/SSL Server Certificate File (PEM Format)**
 - Active Key Trustee Server TLS/SSL Private Key Password**
 - Passive Key Trustee Server TLS/SSL Server Private Key File (PEM Format)**
 - Passive Key Trustee Server TLS/SSL Server Certificate File (PEM Format)**
 - Passive Key Trustee Server TLS/SSL Private Key Password**
- Click **Save Changes** to commit the changes.

- Using the command line:

- Edit `/var/lib/keytrustee/.keytrustee/keytrustee.conf` on the active and passive Key Trustee Server hosts and modify the `SSL_CERTIFICATE` and `SSL_PRIVATE_KEY` parameters as follows:

```
"SSL_CERTIFICATE": "/path/to/new/certificate.pem",
"SSL_PRIVATE_KEY": "/path/to/new/private_key.pem"
```

If the private key is password protected, add the following entry:

```
"SSL_PRIVATE_KEY_PASSWORD_SCRIPT": "/path/to/password_script [arguments]"
```

Replace `/path/to/password_script [arguments]` with the path to a script (and any necessary command arguments) that returns the password for the private key file. If you do not want to create a script, you can use a simple command, such as `echo -n password`. For example:

```
"SSL_PRIVATE_KEY_PASSWORD_SCRIPT": "/bin/echo -n password"
```

Keep in mind that this method can expose the private key password in plain text to anyone who can view the `/var/lib/keytrustee/.keytrustee/keytrustee.conf` file.

5. Restart Key Trustee Server:

- **Using Cloudera Manager:** Restart the Key Trustee Server service (**Key Trustee Server service > Actions > Restart**).
- **Using the Command Line:** Restart the Key Trustee Server daemon:
 - RHEL 6-compatible: `$ sudo service keytrusteed restart`
 - RHEL 7-compatible: `$ sudo systemctl restart keytrusteed`

6. If you are using the Key Trustee KMS service in Cloudera Manager for [HDFS Transparent Encryption](#) on page 261, update the Java KeyStore (JKS) used on the Key Trustee KMS host:

a. Download the new certificate to the Key Trustee KMS host:

```
$ echo -n | openssl s_client -connect keytrustee01.example.com:11371 \
| sed -ne '/-BEGIN CERTIFICATE-/,/-END CERTIFICATE-/p' > /tmp/keytrustee_certificate.pem
```

b. Delete the existing keystore entry for `keytrustee01.example.com`:

```
$ keytool -delete -alias key_trustee_alias_name -keystore /path/to/truststore -v
```

c. Add the new keystore entry for `keytrustee01.example.com`:

```
$ keytool -import -trustcacerts -alias keytrustee01.example.com \
-file /tmp/keytrustee_certificate.pem -keystore /path/to/truststore
```

d. Restart the Key Trustee KMS service in Cloudera Manager.

7. If you are using Key HSM, update the Key Trustee Server and Key HSM configuration:

a. Run the `keyhsm trust` command, using the path to the new certificate:

```
$ sudo keyhsm trust /path/to/new/key_trustee_server/cert
```

b. Run the `ktadmin keyhsm` command, using the `--client-certfile` and `--client-keyfile` options to specify the location of the new certificate file and private key:

```
$ sudo ktadmin keyhsm --server https://keyhsm01.example.com:9090 --client-certfile
/path/to/new/key_trustee_server/cert --client-keyfile
/path/to/new/key_trustee_server/private_key
```

Cloudera Navigator Key HSM

Cloudera Navigator Key HSM allows Cloudera Navigator Key Trustee Server to seamlessly integrate with a hardware security module (HSM). Key HSM enables Key Trustee Server to use an HSM as a root of trust for cryptographic keys, taking advantage of Key Trustee Server's policy-based key and security asset management capabilities while at the same time satisfying existing, internal security requirements regarding treatment of cryptographic materials.

For instructions on installing Key HSM, see [Installing Cloudera Navigator Key HSM](#).

For instructions on configuring Key HSM, continue reading:

Initializing Navigator Key HSM

Before initializing Navigator Key HSM, verify that the HSM is properly configured and accessible from the Key HSM host, and that the HSM client libraries are installed on the Key HSM host:

- **SafeNet Luna**
Install the SafeNet Luna client. No additional configuration is needed.
- **SafeNet KeySecure**

Configuring Encryption

Extract the KeySecure client tarball in the Key HSM library directory (`/usr/share/keytrustee-server-keyhsm/`).

- **Thales**

Install the Thales client service. Copy `nCipherKM.jar`, `jcetools.jar`, and `rsaprivenc.jar` from the installation media (usually located in `opt/nfast/java/classes` relative to the installation media mount point) to the Key HSM library directory (`/usr/share/keytrustee-server-keyhsm/`).

See your HSM product documentation for more information on installing and configuring your HSM and client libraries.



Note: When using an HSM with Key Trustee Server and Navigator Encrypt, encrypting many block devices may exceed the capacity of the HSM. A key is created in the HSM for each encrypted block device, so be sure that your HSM can support your encryption requirements.

To initialize Key HSM, use the `service keyhsm setup` command in conjunction with the name of the target HSM distribution:

```
$ sudo service keyhsm setup [keysecure|thales|luna]
```

For all HSM distributions, this first prompts for the IP address and port number that Key HSM listens on.



Important: If you have implemented Key Trustee Server high availability, initialize Key HSM on each Key Trustee Server.

Cloudera recommends using the loopback address (`127.0.0.1`) for the listener IP address and `9090` as the port number:

```
-- Configuring keyHsm General Setup --
Cloudera Recommends to use 127.0.0.1 as the listener port for Key HSM
Please enter Key HSM SSL listener IP address: [127.0.0.1]127.0.0.1
Will attempt to setup listener on 127.0.0.1
Please enter Key HSM SSL listener PORT number: 9090

validate Port:                               :[ Successful ]
```

If the setup utility successfully validates the listener IP address and port, you are prompted for additional information specific to your HSM. For HSM-specific instructions, continue to the [HSM-Specific Setup for Cloudera Navigator Key HSM](#) on page 311 section for your HSM.

The Key HSM keystore defaults to a strong, randomly-generated password. However, you can change the keystore password in the `application.properties` file:

```
keyhsm.keystore.password.set=yes
```

Then, run the `service keyhsm setup` command with the name of the HSM to which the keystore password applies. You will be prompted to enter the new keystore password, which must be a minimum of six characters in length:

```
$ sudo service keyhsm setup [keysecure|thales|luna]
```

After initial setup, the configuration is stored in the `/usr/share/keytrustee-server-keyhsm/application.properties` file, which contains human-readable configuration information for the Navigator Key HSM server.



Important: The truststore file created during Key HSM initialization must be stored at `/usr/share/keytrustee-server-keyhsm/`. There is no way to change the default location.

HSM-Specific Setup for Cloudera Navigator Key HSM

SafeNet KeySecure



Note: KeySecure was previously named DataSecure, but the Key HSM configuration process is the same for both.

After entering the Key HSM listener IP address and port, the HSM setup for SafeNet KeySecure prompts for login credentials, the IP address of the KeySecure HSM, and the port number:

```
-- Ingrian HSM Credential Configuration --
Please enter HSM login USERNAME: keyhsm
Please enter HSM login PASSWORD: *****

Please enter HSM IP Address or Hostname: 172.19.1.135
Please enter HSM Port number: 9020
```

If the connection is successful, the following message is displayed:

```
Valid address:                               :[ Successful ]
```

The KeySecure setup utility then prompts you whether to use SSL:

```
Use SSL? [Y/n] Y
```

If you choose to use SSL, Key HSM attempts to resolve the server certificate, and prompts you to trust the certificate:

```
[0]          Version: 3
            SerialNumber: 0
            IssuerDN: C=US,ST=TX,L=Austin,O=ACME,OU=Dev,
CN=172.19.1.135,E=webadmin@example.com
            Start Date: Thu Jan 29 09:55:57 EST 2015
            Final Date: Sat Jan 30 09:55:57 EST 2016
            SubjectDN: C=US,ST=TX,L=Austin,O=ACME,OU=Dev,
CN=172.19.1.135,E=webadmin@example.com
            Public Key: RSA Public Key
            modulus: abe4a8dcef92e145984309bd466b33b35562c7f875
1d1c406b1140e0584890272090424eb347647ba04b
34757cacc79652791427d0d8580a652c106bd26945
384b30b8107f8e15d2deba8a4e868bf17bb0207383
7cffef0ef16d5b5da5cfb4d3625c0affbda6320daf
7c6b6d8adfc563960fcd1207c059300feb6513408
79dd2d929a5b986517531be93f113c8db780c92ddf
30f5c8bf2b0bea60359b67be306c520358cc0c3fc3
65500d8abeeac99e53cc2b369b2031174e72e6fca1
f9a4639e09240ed6d4a73073885868e814839b09d5
6aa98a5a1e230b46cdb4818321f546ac15567c5968
33be47ef156a73e537fd09605482790714f4a276e5
f126f935
            public exponent: 10001

            Signature Algorithm: SHA256WithRSAEncryption
            Signature: 235168c68567b27a30b14ab443388039ff12357f
99ba439c6214e4529120d6ccb4a9b95ab25f81b4
7deb9354608df45525184e75e80eb0948eae3e15
c25c1d58c4f86cb9616dc5c68dfe35f718a0b6b5
56f520317eb5b96b30cd9d027a0e42f60de6dd24
5598dlfcea262b405266f484143a74274922884e
362192c4f6417643da2df6dd1a538d6d5921e78e
20a14e29ca1bb82b57c02000fa4907bd9f3c890a
bdae380c0b4dc68710deaeaf41576c0f767879a7
90f30a4b64a6afb3alace0f3ced17ae142ee6f18
5eff64e8b710606b28563dd99e8367a0d3cbab33
2e59c03cadce3a5f4e0aaa9d9165e96d062018f3
6a7e8e3075c40a95d61ebc8db43d77e7

            Extensions:
```

Configuring Encryption

```
critical(false) BasicConstraints: isCa(true)
critical(false) NetscapeCertType: 0xc0

Trust this server? [y/N] Y

Trusted server:                               :[ Successful ]
```

Thales HSM

By default, the Thales HSM client process listens on ports 9000 and 9001. The Cloudera Manager agent also listens on port 9000. To prevent a port conflict, you must change the Thales client ports. Cloudera recommends using ports 11500 and 11501.

To change the Thales client ports, run the following commands:

```
$ sudo /opt/nfast/bin/config-serverstartup --enable-tcp --enable-privileged-tcp
--port=11500 --privport=11501
$ sudo /opt/nfast/sbin/init.d-ncipher restart
```

To configure Key HSM to use the modified port, edit the `/usr/share/keytrustee-server-keyhsm/start.sh` file and add the `-DNFAST_SERVER_PORT=11500` Java system property. For example:

```
java -classpath "*/usr/safenet/lunaclient/jsp/lib/*:/opt/nfast/java/classes/*"
-Djava.library.path=/usr/safenet/lunaclient/jsp/lib/ -DNFAST_SERVER_PORT=11500
com.cloudera.app.run.Program $@
```

Before completing the Thales HSM setup, run the `nfkminfo` command to verify that the Thales HSM is properly configured:

```
$ sudo /opt/nfast/bin/nfkminfo
World generation 2
state           0x17270000 Initialised Usable Recovery !PINRecovery !ExistingClient
                RTC  NVRAM FTO !AlwaysUseStrongPrimes SEEDebug
```

If state reports `!Usable` instead of `Usable`, configure the Thales HSM before continuing. See the Thales product documentation for instructions.

After entering the Key HSM listener IP address and port, the HSM setup for Thales prompts for the OCS card password:

```
Please enter the OCS Card Password (input suppressed):

Configuration saved in 'application.properties' file
Configuration stored in: 'application.properties'. (Note: You can also use service keyHsm
settings to quickly view your current configuration)
```

Luna HSM



Important: If you have implemented Key Trustee Server high availability, ensure that the Luna client on each Key Trustee Server is configured with access to the same partition. See the Luna product documentation for instructions on configuring the Luna client.

Before completing the Luna HSM setup, run the `vtl verify` command (usually located at `/usr/safenet/lunaclient/bin/vtl`) to verify that the Luna HSM is properly configured.

After entering the Key HSM listener IP address and port, the HSM setup for Luna prompts for the slot number and password:

```
-- Configuring SafeNet Luna HSM --
Please enter SafeNetHSM Slot Number: 1
Please enter SafeNet HSM password (input suppressed):
Configuration stored in: 'application.properties'. (Note: You can also use service keyHsm
```



```
settings to quickly view your current configuration)
Configuration saved in 'application.properties' file
```

See the Luna product documentation for instructions on configuring your Luna HSM if you do not know what values to enter here.

Validating Key HSM Settings

After the setup completes, the Key HSM configuration is stored in `/usr/share/keytrustee-server-keyhsm/application.properties`.

You can view these settings using the `service keyhsm settings` command:

```
$ sudo service keyhsm settings

# keyHsm Server Configuration information:
keyhsm.management.address : 172.19.1.2
keyhsm.server.port : 9090
keyhsm.management.port : 9899
keyhsm.service.port : 19791
keyhsm.hardware : ncipher

# Module OCS Password
thales.ocs_password :
GIqhXDuZsJlOet137Lb+f+tqkYvKYDm/8StefpNqZWw1B+LfSY1B4eHd
endtYJio8qLjT+e7j2th5xf809t8FwFVguuyFW+6wdD
uNGvse1LY/itCwqF0ScM1B1Mnz4010xqC6y1PW71+0JjjkkqM5gJJb181sQFFaIGVM/pY=
```

These settings can be manually configured by modifying the `application.properties` file, with the exception of any passwords. These are encrypted by design, and can only be changed by re-running the setup utility.

Creating a Key Store with CA-Signed Certificate

Required Files

Before proceeding, ensure that you have the following three PEM files:

- Certificate Authority (CA) PEM file
- Signed PEM certificate
- Private key PEM file

The following example uses `ssl-cert-keyhsm-ca.pem`, `ssl-cert-keyhsm.pem`, and `ssl-cert-keyhsm-pk.pem`, respectively, to represent these files.

Create the Key Store

The following command accepts the `ssl-cert-keyhsm.pem` and `ssl-cert-keyhsm-pk.pem` files and converts them to a `.p12` file:

```
$ openssl pkcs12 -export -in ssl-cert-keyhsm.pem -inkey ssl-cert-keyhsm-pk.pem -out
mycert.p12 -name alias -CAfile ssl-cert-keyhsm-ca.pem -caname root -chain
```



Important: The certificate CN must match the fully-qualified domain name (FQDN) of the Key Trustee Server.

Managing the Navigator Key HSM Service

Use the `keyhsm` service for all basic server operations:

```
$ sudo service keyhsm
keyHsm service usage:
```

Configuring Encryption

```
setup <hsm name> - setup a new connection to an HSM
trust <path>     - add a trusted client certificate
validate        - validate that keyHSM is properly configured
settings        - display the current server configuration
start           - start the keyHSM proxy server
status          - show the current keyHSM server status
shutdown        - force keyHSM server to shut down
reload          - reload the server (without shutdown)
```

The `reload` command causes the application to restart internal services without ending the process itself. If you want to stop and start the process, use the `restart` command.

Logging and Audits

The Navigator Key HSM logs contain all log and audit information, and by default are stored in the `/var/log/keyhsm` directory.

Integrating Key HSM with Key Trustee Server

Using a hardware security module with Navigator Key Trustee Server requires Key HSM. This service functions as a driver to support interactions between Navigator Key Trustee Server and the hardware security module, and it must be installed on the same host system as Key Trustee Server. The steps below assume that both Key HSM and Key Trustee Server are set up and running. See [Installing Cloudera Navigator Key HSM](#) for details. Integrating Key HSM and Key Trustee Server involves the following steps:

1. [Check Existing Key Names](#) (for existing Key Trustee Server users only)
2. [Establish Trust from Key HSM to Key Trustee Server](#)
3. [Integrate Key HSM and Key Trustee Server](#)

Check Existing Key Names

During the process detailed below, you are prompted to migrate any existing keys from the Key Trustee Server to the HSM.



Warning: Migration fails if any existing keys do not adhere to the [constraints](#).

Successful migration depends on the existing keys conforming to the following constraints:

- Key names can begin with alpha-numeric characters only
- Key names can include only these special characters:
 - Hyphen -
 - Period .
 - Underscore _

If any existing key names in Key Trustee Server do not meet the [requirements listed above](#), the migration fails. To prepare for migration, check your key names and do the following if any of them are non-conforming:

- Decrypt any data using the non-conforming key
- Create a new key, named per the [requirements](#)
- Re-encrypt the data using the new key

After this, the migration from Key Trustee Server to the HSM will succeed during the process below.



Important: Keys are not available during migration, so you should perform these tasks during a maintenance window.

Establish Trust from Key HSM to Key Trustee Server

This step assumes that Key Trustee Server has a certificate for TLS (wire) encryption as detailed in [Managing Key Trustee Server Certificates](#) on page 307. Before you can run the commands in [the steps below](#), Key HSM service must explicitly trust the Key Trustee Server certificate (presented during TLS handshake). To establish this trust, run the following command:

```
$ sudo keyhsm trust /path/to/key_trustee_server/cert
```

The `/path/to/key_trustee_server/cert` in this command (and in the commands below) depends on whether the Key Trustee Server uses the **default certificate** (created by default during install), or uses a **custom certificate** (obtained from a commercial or internal CA). The two alternate paths are shown in the table below. The custom path is a common example but may differ from that shown.

Default	Custom
<code>/var/lib/keytrustee/.keytrustee/.ssl/ssl-cert-keytrustee.pem</code>	<code>/etc/pki/cloudera/certs/cert-file.crt</code>
<code>/var/lib/keytrustee/.keytrustee/.ssl/ssl-cert-keytrustee-pk.pem</code>	<code>/etc/pki/cloudera/private/private-key.key</code>



Note: The system requires TLS and Kerberos authentication throughout the system for security reasons. Connections attempted over SSL (1 through 3) and connections from untrusted clients are immediately terminated to prevent [POODLE](#) (Padding Oracle On Downgraded Legacy Encryption) exploits. See the [Cloudera Security Bulletin](#) for more information.

Integrate Key HSM and Key Trustee Server

The steps below assume that both Key HSM and the Key Trustee Server are on the same host system, as detailed in [Installing Cloudera Navigator Key HSM](#). These steps invoke commands on the Key HSM service and the Key Trustee Server, and they must be run on the host—they cannot be run remotely from another host.

1. Ensure the Key HSM service is running:

```
$ sudo service keyhsm start
```

2. Establish trust from Key Trustee Server to Key HSM specifying the path to the private key and certificate (Key Trustee Server is a client to Key HSM). This example shows how to use the `--client-certfile` and `--client-keyfile` options to specify the path to non-default certificate and key:

```
$ sudo ktadmin keyhsm --server https://keyhsm01.example.com:9090 \
--client-certfile /etc/pki/cloudera/certs/mycert.crt \
--client-keyfile /etc/pki/cloudera/certs/mykey.key --trust
```

For a password-protected Key Trustee Server private key, add the `--passphrase` argument to the command and enter the password when prompted:

```
$ sudo ktadmin keyhsm --passphrase \
--server https://keyhsm01.example.com:9090 \
--client-certfile /etc/pki/cloudera/certs/mycert.crt \
--client-keyfile /etc/pki/cloudera/certs/mykey.key --trust
```



Note: The preceding commands also create a certificate file for the Key HSM that is used by the Key Trustee Server. This certificate file is stored in `/var/lib/keytrustee/.keytrustee/.ssl/ssl-cert-keyhsm.pem`.

3. Restart Key Trustee Server:

- **Using Cloudera Manager:** Restart the Key Trustee Server service (**Key Trustee Server service > Actions > Restart**).
- **Using the Command Line:** Restart the Key Trustee Server daemon:
 - RHEL 6-compatible: `$ sudo service keytrusteed restart`
 - RHEL 7-compatible: `$ sudo systemctl restart keytrusteed`

Cloudera Navigator Encrypt

Cloudera Navigator Encrypt transparently encrypts and secures data at rest without requiring changes to your applications and ensures there is minimal performance lag in the encryption or decryption process. Advanced key management with [Cloudera Navigator Key Trustee Server](#) on page 297 and process-based access controls in Navigator Encrypt enable organizations to meet compliance regulations and ensure unauthorized parties or malicious actors never gain access to encrypted data.

For instructions on installing Navigator Encrypt, see [Installing Cloudera Navigator Encrypt](#).

For instructions on configuring Navigator Encrypt, continue reading:

Registering Cloudera Navigator Encrypt with Key Trustee Server

Prerequisites

Functioning Navigator Key Trustee Server

After [Installing Cloudera Navigator Encrypt](#) on a host, you must register the host with Navigator Key Trustee Server. If you have not yet installed Navigator Key Trustee Server, see [Installing Cloudera Navigator Key Trustee Server](#) for instructions.

Key Trustee Server Organization

To register with Key Trustee Server, you must have an existing organization. See [Managing Key Trustee Server Organizations](#) on page 304 for information on creating and viewing organizations on a Key Trustee Server.

Master Password

The Master Key is the primary Navigator Encrypt administrator access code and is configured by the Navigator Encrypt administrator during installation. The Master Key can take any one of three different forms:

- If you choose a passphrase (single), it must be between 15 and 32 characters long.
- If you choose passphrase (dual), both must be between 15 and 32 characters long.
- If you choose the RSA option, enter a path to the RSA key file, and if it has RSA passphrase, enter it for this private key.



Warning: It is *extremely* important that you keep your master password secret and safe. In the event that you lose your master password, **you will never be able to recover it**, leaving your encrypted data **irretrievably locked away**.

Registering with Key Trustee Server

After [Installing Cloudera Navigator Encrypt](#) on a host, you must register the host with Navigator Key Trustee Server to be able to encrypt and decrypt data. The following section lists the command options for registering your Navigator Encrypt client.



Note: Do not run Navigator Encrypt commands with the `screen` utility.

Example command:

```
$ sudo navencrypt register --server=https://keytrustee01.example.com:11371
--passive-server=https://keytrustee02.example.com:11371 --org=your_keytrustee_org
--auth=org_auth_token
```

Table 24: Registration Options

Command Option	Explanation
<code>--clientname=my_client_name</code>	User-defined unique name for this client to be used for administration and reports. You can verify your client name in the <code>/etc/navencrypt/keytrustee/clientname</code> file.
<code>--server=https://keytrustee01.example.com:11371</code>	Target Active Key Trustee Server for key storage. Replace <code>keytrustee01.example.com:11371</code> with the hostname and port of the Active Key Trustee Server. The default port is 11371.
<code>--passive-server=https://keytrustee02.example.com:11371</code>	Target Passive Key Trustee Server for key storage. Replace <code>keytrustee02.example.com:11371</code> with the hostname and port of the Passive Key Trustee Server. The default port is 11371.
<code>--org=your_keytrustee_org</code>	Key Trustee organization name configured by the Key Trustee Server administrator
<code>--auth=org_auth_token</code>	Organization authorization token, a pre-shared secret by the Navigator Key Trustee Server administrator
<code>--skip-ssl-check</code>	Skip SSL certificate verification. Use with self-signed certificates on the Navigator Key Trustee Server
<code>--trustee</code>	Add trustees for retrieval of the master key
<code>--votes</code>	Configure voting policy for trustees
<code>--recoverable</code>	Master Key will be uploaded without encrypting it with your local GPG Navigator Key Trustee
<code>--scheme "<scheme>"</code>	Key Trustee Server scheme that Navigator Encrypt uses for public key operations. Specify "http" or "https".
<code>--port</code>	Key Trustee Server port that Navigator Encrypt uses for public key operations.

Registering with Previous Versions of Key Trustee Server

By default, new installations of Navigator Key Trustee Server 5.4.0 use a single HTTPS port for key storage and public key operations. Previous versions and upgrades use separate ports for key storage and public key operations. For backward compatibility, Navigator Encrypt 3.7.0 introduces the `--scheme` and `--port` parameters for the `navencrypt register` command.

Configuring Encryption

For example, to register a version 3.7.0 Navigator Encrypt client with a version 3.8.0 Key Trustee Server using HTTPS over port 443 for key storage and HTTP over port 80 for public key operations, run the following command:

```
$ sudo navencrypt register --server=https://keytrustee.example.com:443
--org=key_trustee_org --auth=auth_token --scheme "http" --port 80
```

Navigator Encrypt versions lower than 3.7.0 do not support the `--scheme` and `--port` parameters. For these versions of Navigator Encrypt, you must ensure that the Key Trustee Server is configured to use port 443 (HTTPS) for key storage and port 80 (HTTP) for public key operations.

Navigator Encrypt versions lower than 3.8.0 do not support the `--passive-server` parameter.

Updating Key Trustee Server Ports

The `navencrypt register` command does not provide the ability to change the ports for existing registrations. If the Key Trustee Server ports are changed, you must update `/etc/navencrypt/keytrustee/ztrustee.conf` with the new port and scheme parameters (`HKP_PORT` and `HKP_SCHEME`, respectively).

For example, see the following `ztrustee.conf` excerpt from a registered client that has been upgraded to Navigator Encrypt 3.7.0:

```
{
  "LOCAL_FINGERPRINT": "2048R/182AAA838DC300AC334258D8E7F299BFB68A6F6F",
  "REMOTES": {
    "kts01.example.com": {
      "REMOTE_FINGERPRINT":
"4096R/AF6400E12DC149799CA8CE6BF1604C34D830DE20",
      "REMOTE_SERVER": "https://kts01.example.com",
      "DEFAULT": true,
      "SSL_INSECURE": false,
      "PROTOCOL": "json-encrypt"
    }
  }
}
```

In this example, the Key Trustee Server (`keytrustee.example.com`) is using the default configuration of port 443 (HTTPS) for key storage and port 80 (HTTP) for public key operations.

If the Key Trustee Server is then updated to use port 11371 (HTTPS) for both key storage and public key operations, you must update `ztrustee.conf` as follows (changes in **bold**):

```
{
  "LOCAL_FINGERPRINT": "2048R/182AAA838DC300AC334258D8E7F299BFB68A6F6F",
  "REMOTES": {
    "kts01.example.com": {
      "REMOTE_FINGERPRINT":
"4096R/AF6400E12DC149799CA8CE6BF1604C34D830DE20",
      "REMOTE_SERVER": "https://kts01.example.com:11371",
      ""HKP_PORT": 11371,
      ""HKP_SCHEME": "https",
      "DEFAULT": true,
      "SSL_INSECURE": false,
      "PROTOCOL": "json-encrypt"
    }
  }
}
```

Updating Navigator Encrypt for High Availability Key Trustee Server

If you registered a Navigator Encrypt client with a standalone Key Trustee Server, and then configured [high availability](#) for Key Trustee Server, you can edit `/etc/navencrypt/keytrustee/ztrustee.conf` to enable the client to take

advantage of the high availability features. The following example shows the contents of `ztrustee.conf` after adding the required `REMOTE_SERVERS` entry (changes in **bold**):

```
{
  "LOCAL_FINGERPRINT": "2048R/182AAA838DC300AC334258D8E7F299BFB68A6F6F",
  "REMOTES": {
    "kts01.example.com": {
      "REMOTE_FINGERPRINT":
"4096R/AF6400E12DC149799CA8CE6BF1604C34D830DE20",
      "REMOTE_SERVER": "https://kts01.example.com:11371",
      "HKP_PORT": 11371,
      "HKP_SCHEME": "https",
      "DEFAULT": true,
      "REMOTE_SERVERS": ["https://kts01.example.com:11371",
"https://kts02.example.com:11371"],
      "SSL_INSECURE": true,
      "PROTOCOL": "json-encrypt"
    }
  }
}
```

Configuration Files

The installer creates the `/etc/navencrypt` directory. All configuration settings are saved in this directory. **Do not** delete any file from `/etc/navencrypt`. These files provide the necessary information for the Navigator Encrypt application to function properly.



Warning: Perform backups of encrypted data, mount-points, and Navigator Encrypt configuration directories on a regular basis. To do this, ensure you have a backup of `/etc/navencrypt`. **Failure to backup this directory will make your backed up encrypted data unrecoverable in the event of data loss.**

Change Master Key by UUID

It is possible to re-use a previously used Master Key by its UUID. For example, if you currently have a Master key with a single passphrase, you can see the corresponding Navigator Key Trustee UUID in the `/etc/navencrypt/control` file:

```
$ cat /etc/navencrypt/control
{
  "app": {
    "name": "navencrypt",
    "version": "3.5"
  },
  "keys": {
    "master": {
      "type": "single-passphrase",
      "uuid": "qMAKRMdk4HVbhfzR79cp9w92YBmNHJ5nSLhfd8ZVo6L"
    },
    "targets": []
  }
}
```

You can copy the UUID (`qMAKRMdk4HVbhfzR79cp9w92YBmNHJ5nSLhfd8ZVo6L` in this example) and run `navencrypt key --change --new-master-key-uuid` to change a Master Key by using its UUID only:

```
$ sudo navencrypt key --change
--new-master-key-uuid=qMAKRMdk4HVbhfzR79cp9w92YBmNHJ5nSLhfd8ZVo6L
>> Type your OLD Master key
Type MASTER passphrase 1:
Type MASTER passphrase 2:
Verifying Master Key against Navigator Key Trustee (wait a moment)...
OK
Changing Master key (wait a moment)...
```

Configuring Encryption

```
* Setting up EXISTING MASTER key...
* Uploading CONTROL content...
* Re-encrypting local keys...
Master key successfully changed.
```



Note: The `navencrypt key` command fails if no volumes are encrypted or the kernel module is not loaded.

Preparing for Encryption Using Cloudera Navigator Encrypt

Before you can encrypt data, you must prepare a storage repository to hold the encrypted data and a mount point through which to access the encrypted data. The storage repository and mount point must exist before encrypting data using the `navencrypt-move` command.

Data stored and retrieved from the repository is encrypted and decrypted transparently.

Cloudera Navigator Encrypt *does not* support:

- Encrypting a directory that contains or is contained within a mount point for another service (including Navigator Encrypt and NFS). See [Encrypting Data](#) on page 326 for more information.
- Encrypting immutable files or directories containing immutable files.
- Installation or use in `chroot` environments, including creating `chroot` environments within an encrypted directory.
- Encrypting HDFS data files.

Navigator Encrypt Commands



Note: Do not run Navigator Encrypt commands with the `screen` utility.

The following table lists the commands used to encrypt data:

Table 25: Navigator Encrypt Commands

Command	Description
<code>navencrypt</code>	Manage, update, and verify your data.
<code>navencrypt-prepare</code>	Prepare your system for encryption by creating mount-points and specifying storage.
<code>navencrypt-prepare --undo</code>	Remove a mountpoint that is no longer in use.
<code>navencrypt-move</code>	Encrypt/decrypt your data to/from the encrypted filesystem.
<code>navencrypt-profile</code>	Generate process profile information in JSON format.
<code>navencrypt-module-setup</code>	Build or rebuild the Navigator Encrypt kernel module.

Preparing for Encryption



Note: When using an HSM with Key Trustee Server and Navigator Encrypt, encrypting many block devices may exceed the capacity of the HSM. A key is created in the HSM for each encrypted block device, so be sure that your HSM can support your encryption requirements.

To get an in-depth look at the details behind the `navencrypt-prepare` command, or to use a unique configuration, use the interactive prompt by executing `navencrypt-prepare` with no options. This launches an interactive console that guides you through the following operations:

- Creating internal encryption keys
- Registering internal keys in Navigator Key Trustee
- Registering mount point in `/etc/navencrypt/ztab`
- Mounting current mount point
- Establishing encryption method (`dm-crypt` for devices, `ecryptfs` for directories)

Using the console, you can choose how you want your data stored and accessed. Navigator Encrypt offers two different types of encryption:

- Block-level encryption with `dm-crypt`: Protect your data by encrypting the entire device. This option enables full disk encryption and is optimized for some system configurations. Block-level encryption can be used with logical devices such as a loop device.
- File-level encryption with `ecryptfs`: Protect your data by mounting an encrypted filesystem on top of an existing one. Enables transparent access to encrypted data without modifying your storage.



Note: As of August 2015, Filesystem-level encryption using eCryptfs is [deprecated](#). For new installations, use block-level encryption. For existing installations using eCryptfs, see [Migrating eCryptfs-Encrypted Data to dm-crypt](#) on page 328 for instructions on migrating data encrypted using eCryptfs to use dm-crypt.

See [Block-Level Encryption with dm-crypt](#) on page 322 and [Filesystem-Level Encryption with eCryptfs](#) on page 324 for more information.

To prepare for encryption, you must set a location to store the encrypted data and a mount point through which to access the data. The storage location and mount point must be created before encrypting data.



Note: If you are performing a file system check as part of your preparation work, then note that the crypto device must be mapped and active. Also, be aware that if you execute `fsck` in force mode (`-f`), there is a risk of data loss. If the force operation fails, it could cause file system corruption at the device header.

In the following example, we will use the directory `/navencrypt/encrypted-storage` for the encrypted storage and `/navencrypt/mount-point` for the mount point. If you have specific space/partition requirements, you can select a different directory, although Cloudera highly recommends that you place the encrypted directory on the same partition as the data you are planning to encrypt.

The syntax for the prepare command is as follows:

```
$ sudo navencrypt-prepare <data_storage_directory> <mount_point>
```

When specifying the storage path and the mount point path, *do not* use a trailing `/` in the path names. Both directories must exist prior to running the `navencrypt-prepare` command. They are not automatically created.

To create the encrypted partition, create the mount point and storage directories, and then use the `navencrypt-prepare` utility:

```
$ sudo mkdir -p /navencrypt/encrypted-storage /navencrypt/mount-point
$ sudo navencrypt-prepare /navencrypt/encrypted-storage /navencrypt/mount-point
```

For RHEL 7, run the following command after the `navencrypt-prepare` command completes:

```
$ sudo systemctl start navencrypt-mount
```

To demonstrate the difference between the two directories, this example uses different directories for the encrypted storage and the mount point. It is also possible to store and access the data using the same directory.

Configuring Encryption

To see the effects of these commands, run `df -h`. This command displays the partition information about your system. You should see an `ecryptfs` partition located at `/navencrypt/encrypted-storage`, and mounted at `/navencrypt/mount-point`.

After you have successfully prepared a client for encryption, you can encrypt and decrypt data using the commands described in [Encrypting and Decrypting Data Using Cloudera Navigator Encrypt](#) on page 325.

Block-Level Encryption with dm-crypt



Note: For best performance, Cloudera strongly recommends using block encryption with dm-crypt. See [Migrating eCryptfs-Encrypted Data to dm-crypt](#) on page 328 for instructions on migrating data encrypted using eCryptfs to use dm-crypt.

When choosing block-level encryption during the interactive console, you must specify two parameters:

1. The first parameter is the block device that you want to store the encrypted file system in. Because this device stores all of the encrypted data, it must be as large as or larger than the target data. The device must exist and be empty. Supported storage devices are:
 - Physical block devices (for example, a disk device)
 - Virtual block devices (for example, a block device created by [LVM](#))
 - Loop devices (see [Block-Level Encryption with a Loop Device](#) on page 322 for instructions on creating a loop device)
2. The second parameter is the mount point for the encrypted file system. This is the location where you can access the encrypted data stored in the first parameter. The mount point must already exist. It is not created by the `navencrypt-prepare` command.

The entire device in the first parameter is used for encrypted data.



Note: Do not manually unmount the encryption mount point (for example, using `umount`). If you do so, you must manually close the `dm-crypt` device using the following procedure:

1. Run `dmsetup table` to list the `dm-crypt` devices.
2. Run `cryptsetup luksClose <device_name>` to close the device for the unmounted mount point.

After choosing these two parameters and following the interactive console (discussed further in [Preparing for Encryption](#) on page 320), you are ready to encrypt your data. The following example shows successful output from a `navencrypt-prepare` command using dm-crypt for block-level encryption:

```
$ sudo /usr/sbin/navencrypt-prepare /dev/sdal /mnt/dm_encrypted
Type MASTER passphrase:
Encryption Type:  dmCrypt (LUKS)
Cipher:           aes
Key Size:         256
Random Interface: /dev/urandom
Filesystem:       ext4
Verifying MASTER key against Navigator Key Trustee (wait a moment) ... OK
Generation Encryption Keys with /dev/urandom ... OK
Preparing dmCrypt device (--use-urandom) ... OK
Creating ext4 filesystem ... OK
Registering Encryption Keys (wait a moment) ... OK
Mounting /dev/sdal ... OK
```

Block-Level Encryption with a Loop Device

A block-level encrypted device can be a physical device or a storage space treated as a device. See [Migrating eCryptfs-Encrypted Data to dm-crypt](#) on page 328 for instructions on migrating data encrypted using eCryptfs to use dm-crypt with a loop device.

To configure a loop device, use the `dd` command to create a storage space:



Warning: The space for the loop device is pre-allocated. After the loop device is created, the size cannot be changed. Make sure you are allocating enough storage for the current encrypted data as well as any future data.

If your disks are mounted individually with a filesystem on each disk, and your storage needs exceed the capacity of a single disk, you can create a loop device on each disk for which you want to allocate space for encrypted data. If you have consolidated storage (for example, using LVM), you can create a single loop device or multiple devices.

```
$ sudo dd if=/dev/zero of=/dmccrypt/storage bs=1G count=500
```

The `dd` command above creates a 500 GB file. Modify the `bs` and `count` values to generate the required file size.

After generating the file, run `losetup -f` to view unused loop devices. Use the available loop device with the `navencrypt-prepare -d` command, demonstrated below.

Specifically for loop devices, the `-d` parameter enables Navigator Encrypt to manage the loop device association. You no longer need to use the `losetup` command to associate the file with the loop device, and the loop device is automatically prepared at boot. For RHEL 7-compatible OS, you must run the following commands to ensure that a loop device is available at boot:

```
$ sudo bash -c 'echo "loop" > /etc/modules-load.d/loop.conf'
$ sudo bash -c 'echo "options loop max_loop=8" > /etc/modprobe.d/loop_options.conf'
```



Warning: Loop devices are not created by Navigator Encrypt. Instead, Navigator Encrypt assigns a datastore to a loop device when the `navencrypt-prepare --datastore` option is used. So, it is up to the system administrator to create persistent `/dev/loopX` devices, which are required to prepare a virtual block device. If the loop device being prepared is not persistent, then Navigator Encrypt will not mount the device upon a reboot.

The data storage directory name (`/dmccrypt/storage` in the previous example) must contain only alphanumeric characters, spaces, hyphens (`-`), or underscores (`_`). Other special characters are not supported.

The following example shows the output from a successful command:

```
$ losetup -f
/dev/loop0
$ sudo navencrypt-prepare -d /dmccrypt/storage /dev/loop0 /dmccrypt/mountpoint
Type MASTER passphrase:

Encryption Type:  dmCrypt (LUKS)
Cipher:           aes
Key Size:         256
Random Interface: OpenSSL
Filesystem:       ext4
Options:

Verifying MASTER key against KeyTrustee (wait a moment)  ... OK
Generation Encryption Keys with OpenSSL                  ... OK
Assigning '/dev/loop0' -> '/dmccrypt/storage'             ... OK
Preparing dmCrypt device                                 ... OK
Creating ext4 filesystem                                  ... OK
Registering Encryption Keys (wait a moment)               ... OK
Mounting /dev/loop0                                       ... OK
```

Configuring Encryption

For upgraded Navigator Encrypt clients that already use loop devices, you can enable Navigator Encrypt to manage the loop device file association (instead of configuring the system to run the `losetup` command at boot) by adding the `nav_datastore` option to the entry in `/etc/navencrypt/ztab`. For example:

```
# <target mount-dir>          <source device>      <type>      <options>
/dmccrypt/mountpoint         /dev/loop0          luks
key=keytrustee,nav_datastore='/dmccrypt/storage'
```



Important: Use caution when editing the `/etc/navencrypt/ztab` file. Entries are tab-separated (not space-separated). The `ztab` file must not contain empty lines.

After you have created the loop device, continue with the instructions in [Block-Level Encryption with dm-crypt](#) on page 322.

Filesystem-Level Encryption with eCryptfs



Note: As of August 2015, Filesystem-level encryption using eCryptfs is [deprecated](#). For best performance, Cloudera strongly recommends using [Block-Level Encryption with dm-crypt](#) on page 322 where possible. See [Migrating eCryptfs-Encrypted Data to dm-crypt](#) on page 328 for instructions on migrating data encrypted using eCryptfs to use dm-crypt.

RHEL 7 does not support eCryptfs. For new installations on RHEL 7, you must use [Block-Level Encryption with dm-crypt](#) on page 322. If you are planning on upgrading to RHEL 7 and are currently using eCryptfs, migrate to dm-crypt before upgrading.

When choosing file-level encryption during the interactive console, you must specify two parameters:

1. The first parameter is the storage directory you want to store the encrypted file system in. Because this directory will hold all of the encrypted data, it must be as large as or larger than the target data.
2. The second parameter is the mount point for the encrypted file system. This is the location where you can access the encrypted data stored in the location identified by the first parameter.

While the data is technically stored at the location identified by the first parameter, you can only access the data from the mount point identified by the second parameter. Consider this when choosing where to mount your data.

After choosing these two parameters and following the interactive console (discussed further in [Preparing for Encryption](#) on page 320), you are ready to encrypt your data.

Undo Operation

Navigator Encrypt 3.5 and higher supports a new command option, `navencrypt-prepare --undo`. This command reverses the operations from the regular `navencrypt-prepare` command by removing the device from Navigator Encrypt control and removing registered encryption keys.

The only parameter of the undo operation is the storage device used to store the encrypted file system (not the mount point). Here is an example showing `navencrypt-prepare` and `navencrypt-prepare --undo` operations:

```
$ sudo navencrypt-prepare /path/to/storage /path/to/mountpoint
Type MASTER passphrase:

Encryption Type:  eCryptfs
Cipher:           aes
Key Size:         256
Random Interface: OpenSSL
Filesystem:       ext4
Options:

Verifying MASTER key against Navigator Key Trustee (wait a moment)    ... OK
Generation Encryption Keys with OpenSSL                               ... OK
Registering Encryption Keys (wait a moment)                           ... OK
Mounting /path/to/mountpoint                                          ... OK
```

```
$ sudo navencrypt-prepare --undo /path/to/storage
Type MASTER passphrase:
Verifying MASTER key against Navigator Key Trustee (wait a moment)    ... OK
Unmounting /path/to/mountpoint                                     ... OK
```

Pass-through Mount Options for `navencrypt-prepare`

Navigator Encrypt 3.5 and higher provides the ability to specify options to pass to the `mount` command that is executed during `/etc/init.d/navencrypt-mount start` (`systemctl start navencrypt-mount` on RHEL 7). These options are specified with the `-o` option when preparing a mountpoint with the `navencrypt-prepare` command.

The following shows an example `navencrypt-prepare` command output when passing mount options with the `-o` option:

```
$ sudo navencrypt-prepare -o discard,resize /mnt/t2 /mnt/t2
Type MASTER passphrase:

Encryption Type:  eCryptfs
Cipher:           aes
Key Size:         256
Random Interface: OpenSSL
Filesystem:       ext4
Options:          discard,resize

Verifying MASTER key against Navigator Key Trustee(wait a moment)    ... OK
Generation Encryption Keys with OpenSSL                             ... OK
Registering Encryption Keys (wait a moment)                         ... OK
Mounting /mnt/t2                                                    ... OK
```

You can verify the results by viewing the `/etc/navencrypt/ztab` file:

```
$ cat /etc/navencrypt/ztab
/mnt/t2 /mnt/t2 eCryptfs key=keytrustee,cipher=aes,keysize=256,discard,resize
```

Options can be added or removed to existing mount points prepared with versions of Navigator Encrypt prior to 3.5 by editing the `/etc/navencrypt/ztab` file and adding the comma-separated options (no spaces) to the end of each line as seen in the previous example above.



Important: Use caution when editing the `/etc/navencrypt/ztab` file. Entries are tab-separated (not space-separated). The `ztab` file must not contain empty lines.

To see the mounted filesystems and options, run `mount`:

```
$ mount
/mnt/t2 on /mnt/t2 type eCryptfs
(rw,ecryptfs_sig=6de3db1e87077adb,ecryptfs_unlink_sigs,noauto,\
ecryptfs_cipher=aes,ecryptfs_key_bytes=32,discard,resize)
```

Pass-through mount options work for both `dm-crypt` and `eCryptfs`. For a list of available mount options, see the `man` pages for `cryptsetup` and `ecryptfs` respectively.

Encrypting and Decrypting Data Using Cloudera Navigator Encrypt



Warning: Before encrypting or decrypting any data, stop all processes (for example, MySQL, MongoDB, PostgreSQL, and so on) that have access to the target data. **Failure to do so could lead to data corruption.**

After the encrypted file system is created and initialized, it is ready to hold data. All encryption and decryption functionality is performed with a single command: `navencrypt-move`.

Configuring Encryption

Do not manually create directories or files under a Cloudera Navigator Encrypt mount point; use only the `navencrypt-move` command to encrypt and decrypt data. See [Preparing for Encryption Using Cloudera Navigator Encrypt](#) on page 320 for more information about mount points.

After encrypting a file or directory, all data written and read through the mount point is transparently encrypted and decrypted.

Before You Begin

Navigator Encrypt does not support encrypting data in certain environments, including the following:

- *Do not* attempt to encrypt a directory that contains or is contained within a mount point for another service (including Navigator Encrypt and NFS). For example:
 - If your encryption mount point is `/var/lib/navencrypt/mount`, do not attempt to encrypt `/var`, `/var/lib`, `/var/lib/navencrypt`, `/var/lib/navencrypt/mount`, or anything under `/var/lib/navencrypt/mount/`.
 - If you have mounted an NFS filesystem at `/mnt/home`, do not attempt to encrypt `/mnt`, `/mnt/home`, or anything under `/mnt/home`.
- Do not attempt to encrypt immutable files or directories containing immutable files.
- Do not use Navigator Encrypt within a `chroot` environment, or create a `chroot` environment within an encrypted directory.
- If your Key Trustee Server is managed by Cloudera Manager, *do not* encrypt the Cloudera Manager database with Navigator Encrypt; doing so prevents Cloudera Manager from starting.

Encrypting Data

Do not manually create directories or files under a Navigator Encrypt mount point; use only the `navencrypt-move` command to encrypt data.

Here is an example command to encrypt data, with an explanation for each option:


```
$ sudo navencrypt-move encrypt @<category> <directory_to_encrypt> <encrypted_mount_point>
```



Important: *Do not* run `navencrypt-move` commands simultaneously in multiple terminals. Doing so results in failure to encrypt or decrypt all of the specified data. No data is lost, as the source data is not removed, but you must re-run the failed operations sequentially.

Table 26: navencrypt-move Command Options

Command Option	Explanation
<code>navencrypt-move</code>	Main command interface for all actions that require moving data either to or from the encrypted file system. For more information see the <code>navencrypt-move</code> man page (<code>man navencrypt-move</code>).
<code>encrypt</code>	Identifies the cryptographic operation, in this case, encrypting data. The <code>decrypt</code> option is described later in Decrypting Data on page 327.

Command Option	Explanation
	 Note: By default, all Navigator Encrypt encryption commands require free space equal to twice the size of the encrypted data. If your environment does not have enough free space, add <code>--per-file</code> to the end of the command. This moves each file individually. Per-file encryption only requires free space equal to twice the size of the largest individual file, but is a slower operation.
@<category>	The access category that is applied to the data being encrypted. Encrypted data is protected by process-based access controls that restrict access to only the processes that you allow. You can use any naming convention you want (the @ symbol is required), but Cloudera recommends keeping it simple and memorable. For example, you can use a name referencing the data type being encrypted, such as <code>@mysql</code> for a MySQL deployment.
<directory to encrypt>	The data that you want to encrypt. This can be a single file or an entire directory. Navigator Encrypt starts after the system boots, so do not encrypt required system files and directories (such as the root partition, <code>/var</code> , and so on). Some examples of recommended data directories to encrypt are <code>/var/lib/mysql/data</code> , <code>/db/data</code> , and so on.
<encrypted mount-point>	Where you want to store the data. This is the path to the mount point specified during the navencrypt-prepare command.

When a file is encrypted, a symbolic link (symlink) is created which points to a mount point @<category> directory. The `navencrypt-move` command moves all specified data to the encrypted filesystem and replaces it with a symlink to the mount point for that encrypted filesystem.

Encrypting a directory is similar to encrypting a file. The following command encrypts a directory:

```
$ sudo /usr/sbin/navencrypt-move encrypt @mycategory /path/to/directory_to_encrypt/
/path/to/mount
```

In this command, a directory is specified instead of a filename, and a symlink is created for that particular directory. To see the effects of this command, run:

```
$ ls -l <directory_to_encrypt>
$ du -h <encrypted_storage_directory>
```

The output demonstrates the new filesystem layout. Everything that was in the target directory is now securely stored in the encrypted filesystem.

Decrypting Data

The decryption command requires only the path to the original data, which is now a symbolic link, as an argument. The following example demonstrates how to decrypt a file using the `navencrypt-move` command:

```
$ sudo /usr/sbin/navencrypt-move decrypt /path/to/encrypted/file
```



Important: Do not run `navencrypt-move` commands simultaneously in multiple terminals. Doing so results in failure to encrypt or decrypt all of the specified data. No data is lost, as the source data is not removed, but you must re-run the failed operations sequentially.

As with encryption, you can specify a directory instead of a file:

```
$ sudo /usr/sbin/navencrypt-move decrypt /path/to/encrypted/directory
```

Migrating eCryptfs-Encrypted Data to dm-crypt

As of August 2015, Filesystem-level encryption using eCryptfs is [deprecated](#). Use this procedure to migrate to dm-crypt.

RHEL 7 does not support eCryptfs. For new installations on RHEL 7, you must use [Block-Level Encryption with dm-crypt](#) on page 322. If you are planning on upgrading to RHEL 7 and are currently using eCryptfs, migrate to dm-crypt before upgrading.



Warning: Before encrypting or decrypting any data, stop all processes (for example, MySQL, MongoDB, PostgreSQL, and so on) that have access to the target data. **Failure to do so could lead to data corruption.**

1. Prepare an empty block device. This can be a physical block device (such as an unused disk) or a virtual block device (for example, a logical block device created by [LVM](#), or a loop device). For instructions on creating a loop device, see [Block-Level Encryption with a Loop Device](#) on page 322.
2. Stop any services which depend on the encrypted data to be moved.
3. Prepare a block-level encrypted mount point. See [Preparing for Encryption Using Cloudera Navigator Encrypt](#) on page 320 for details about the procedure.
4. [Add ACL rules](#) for the new encrypted mount point that match the ACL rules for the mount point you are migrating from. To view existing ACL rules, run `sudo navencrypt acl --print`.
5. Add an ACL rule for your preferred shell (for example, `/bin/bash`) to enable command-line utilities such as `mv` and `cp`:

```
$ sudo navencrypt acl --add --rule="ALLOW @category * /bin/bash"
```

6. Copy the encrypted data from the eCryptfs mount point to the dm-crypt mount point:

```
$ sudo cp -rp /ecryptfs/mountpoint/path/to/data /dmccrypt/mountpoint/path/to/data
```

7. Update any symbolic links referencing the encrypted data. The following example demonstrates updating a symbolic link for a PostgreSQL database that was originally encrypted using eCryptfs, but has been migrated to dm-crypt:

```
$ sudo ls -l /var/lib/db/data/base/16385
lrwxrwxrwx 1 root root 72 Jul 22 15:33 /var/lib/db/data/base/16385 ->
/ecryptfs/mountpoint/postgres/var/lib/db/data/base/16385
$ sudo ln -sif /dmccrypt/mountpoint/postgres/var/lib/db/data/base/16385
/var/lib/db/data/base/16385
$ sudo ls -l /var/lib/db/data/base/16385
lrwxrwxrwx 1 root root 72 Jul 22 15:33 /var/lib/db/data/base/16385 ->
/dmccrypt/mountpoint/postgres/var/lib/db/data/base/16385
```

8. Remove the ACL rule enabling command-line utilities:

```
$ sudo navencrypt acl --del --rule="ALLOW @category * /bin/bash"
```

9. Restart any services which depend on the encrypted data.
10. Verify that the data was successfully copied, then delete the original eCryptfs-encrypted data. Do not delete any data until you are certain that you no longer need the original data.

- a. Stop the `navencrypt-mount` service:

```
$ sudo service navencrypt-mount stop
```

- b. Remove the original mountpoint directory and the storage directory with the original encrypted data.
 c. Edit `/etc/navencrypt/ztab` and remove entries for the original encrypted directory where `eCryptfs` is listed as the `<type>`.



Important: Use caution when editing the `/etc/navencrypt/ztab` file. Entries are tab-separated (not space-separated). The `ztab` file must not contain empty lines.

- d. Start the `navencrypt-mount` service:

```
$ sudo service navencrypt-mount start
```

Cloudera Navigator Encrypt Access Control List

Managing the Access Control List

Cloudera Navigator Encrypt manages file system permissions with an access control list (ACL). This ACL is a security access control created by Cloudera that enables a predefined Linux process to access a file or directory managed by Navigator Encrypt.

The ACL uses rules to control process access to files. The rules specify whether a Linux process has access permissions to read from or write to a specific Navigator Encrypt path.

A rule is defined in the following order:

```
# TYPE @CATEGORY PATH PROCESS PARAMETERS
```

The following table defines the ACL rule components:

Table 27: ACL Rule Components

Component	Description
TYPE	Specifies whether to allow or deny a process. It can have either of the following values: <code>ALLOW</code> or <code>DENY</code> .
@CATEGORY	This is a user-defined shorthand, or container, for the encrypted dataset that the process will have access to. For example, if you are encrypting the directory <code>/var/lib/mysql</code> , you could use the category <code>@mysql</code> to indicate that this rule is granting access to a process on the MySQL data.
PATH	Specifies the rights permissions of a specific path. For example: <code>*,www/*.htaccess</code> . Omit the leading slash (<code>/</code>).
PROCESS	Specifies the process or command name for the rule.
PARAMETERS	Tells the process the parent-child process to be executed: <code>--shell</code> defines the script for Navigator Encrypt to allow for executable process. Supported shells are <code>/usr/bin/bash</code> , <code>/bin/bash</code> , <code>/usr/bin/dash</code> , and <code>/bin/bash</code> .

Component	Description
	<p><code>--children</code> defines for Navigator Encrypt which child processes to allow that are executed by a process/script.</p> <p>Example: <code>--shell=/bin/bash,</code> <code>--children=/bin/df,/bin/ls</code></p>

All rules are stored in an encrypted policy file together with a set of process signatures that are used by Navigator Encrypt to authenticate each Linux process. This file is encrypted with the Navigator Encrypt key you defined during installation.

Cloudera recommends using `permissive` mode to assist with the initial ACL rule creation for your environment. In `permissive` mode, Navigator Encrypt allows full access to the encrypted data by all processes, but logs them in `dmesg` as `action="denied"` messages. Consult these messages to identify required ACL rules. To set Navigator Encrypt to `permissive` mode, use the following command:

```
$ sudo /usr/sbin/navencrypt set --mode=permissive
```

To view the current mode, run `navencrypt status -d`. For more information on access modes, see [Access Modes](#).

`deny2allow`

After you generate the `action="denied"` messages, use the `navencrypt deny2allow` command to show which ACL rules are required, based on the `action="denied"` messages in `dmesg`. To show which ACL rules are required, perform the following steps:

1. Save the `dmesg` content to a file:

```
$ sudo dmesg > /tmp/dmesg.txt
```

2. Use the `dmesg.txt` file content as input to the `deny2allow` command to analyze the `action="denied"` messages and display a list of required ACL rules based on the `action="denied"` messages. Here is an example command and output:

```
$ sudo /usr/sbin/navencrypt deny2allow /tmp/dmesg.txt
ALLOW @mysql employees/* /usr/sbin/mysqld
ALLOW @mysql * /bin/bash
ALLOW @mysql * /bin/ls
```

If you need to clear the `dmesg` log and start fresh, run `dmesg -c`.

If a rule is displayed in the output from the command, it does not automatically mean the ACL rule must be added. You must determine which rules are actually needed. For example, the rule for `ls` would not typically be added as an ACL rule.

After the initial ACL rules are created, disable `permissive` mode with the following command:

```
$ sudo /usr/sbin/navencrypt set --mode=enforcing
```

Adding ACL Rules

Rules can be added one at a time using the command line or by specifying a policy file containing multiple rules. The following example shows how to add a single rule using the `navencrypt acl --add` command:

```
$ sudo /usr/sbin/navencrypt acl --add --rule="ALLOW @mysql * /usr/sbin/mysqld"
```

The following example shows how to add multiple rules using a policy file:

```
$ sudo /usr/sbin/navencrypt acl --add --file=/mnt/private/acl_rules
```

The contents of the policy file should contain one rule per line. For example:

```
ALLOW @mysql * /usr/sbin/mysqld
ALLOW @log * /usr/sbin/mysqld
ALLOW @apache * /usr/lib/apache2/mpm-prefork/apache2
```

Using a policy file is the fastest way to add multiple rules because it only requires the security key one time.

It is also possible to overwrite the entire current rules set with the option `--overwrite`. When this command is executed, all current rules are replaced by the ones specified in the file that contains the new set of rules. It is recommended to save a copy of your current set of rules by printing it with the option `--print`.

Here is an example command using the `--overwrite` option:

```
$ sudo /usr/sbin/navencrypt acl --overwrite --file=/mnt/private/acl_rules
```

Adding ACL Rules by Profile

If your environment requires more granular controls on the processes that can access the data, you can add extra controls by using profiles. Profiles set requirements on a process other than just having the correct fingerprint. They can include such things as process owner and group, required open files, and the current working directory. To see more about adding rules by profile, see [ACL Profile Rules](#) on page 333.

Deleting ACL Rules

Rules can be deleted in one of two ways:

1. Manually specifying the rule to delete using the command line.
2. Specifying the line number of the rule to delete.

The following example shows how to delete a rule by passing it as a parameter:

```
$ sudo /usr/sbin/navencrypt acl --del --rule="ALLOW @mysql * /usr/sbin/mysqld "
```

If you remove a MySQL `ALLOW` rule, the MySQL cache must be cleaned by executing the `FLUSH TABLES;` MySQL statement. Otherwise, it will still be possible to view data from encrypted table.

The following example shows how to delete a rule by specifying a line number:

```
$ sudo /usr/sbin/navencrypt acl --del --line 3
```

It is also possible to delete multiple ACL rules in a single command:

```
$ sudo /usr/sbin/navencrypt acl --del --line=1,3
```

See [Printing ACL Rules](#) on page 331 for information on determining line numbers.

Deleting ACL Rules by Profile

See [ACL Profile Rules](#) on page 333 for instructions on deleting rules by profile.

Printing ACL Rules

You can print the current Access Control List using the following command:

```
$ sudo /usr/sbin/navencrypt acl --print
```

Configuring Encryption

Save the ACL to a file with the `--file` option:

```
$ sudo /usr/sbin/navencrypt acl --print --file=policy-backup
```

To display additional information about the organization of the policy file, use the `--list` option:

```
$ sudo /usr/sbin/navencrypt acl --list
```

Universal ACL Rules

Universal ACLs will allow or deny a process access to all files or directories encrypted with Navigator Encrypt.

The rule `ALLOW @* * /process` allows the designated process to access anything from all encrypted categories.

The rule `ALLOW @data * *` allows all processes access to any path under the `@data` category.

The rule `ALLOW @* * *` allows all processes access to all encrypted categories. Cloudera does not recommend using this rule. Use it only in test environments.

Here is an example adding a universal ACL rule and then displaying it:

```
$ sudo /usr/sbin/navencrypt acl --add --rule="ALLOW @* * /usr/sbin/mysqld"
Type MASTER passphrase:
1 rule(s) were added
# navencrypt acl --listType MASTER passphrase:
# - Type      Category      Path      Profile  Process
1   ALLOW     @*              *         /usr/sbin/mysqld
```

Enabling Shell Scripts to Be Detected by ACL

All of the previous rules work for binary files. There may be times a script, such as a shell script, must be allowed to access the encrypted directory.

You can add the script as a rule by indicating the executable binary process of this script using the `--shell` option, for example:

```
ALLOW @scripts * /root/script.sh --shell=/bin/bash
```

The `--shell` option identifies which executable process is used to execute the script. Supported shells are `/usr/bin/bash`, `/bin/bash`, `/usr/bin/dash`, and `/bin/dash`.

If the script is altered, it will no longer be trusted by the ACL because the fingerprint has changed. If you edit the script you must invoke the update option to update the ACL with the new fingerprint.

In some cases, it may be necessary to grant permissions to sub-processes invoked by scripts. For example, it may be necessary to grant permissions to `/bin/bash` that also allow running the `/bin/df` command to allow the system administrator to check disk capacity through a script run using a `crontab` entry. By using the `--children` option, you can specify these permissions. For example:

```
ALLOW @scripts * /root/script.sh --shell=/bin/bash --children=/bin/df
```

The `--children` option tells Navigator Encrypt to allow the `/bin/df` binary process if it is executed by `/root/script.sh`.

To allow more than one sub-process, identify them with the `--children` option as comma-separated values. For example:

```
ALLOW @scripts * /root/script.sh --shell=/bin/bash --children=/bin/df,/bin/ls
```

To add shell-children sub-processes, execute the `navencrypt acl --add` command, for example:

```
$ sudo /usr/sbin/navencrypt acl --add --rule="ALLOW @mysql * /usr/bin/mysqld_safe \  
--shell=/bin/bash --children=/bin/df,/bin/ls"
```

ACL Profile Rules

If your environment requires more granular controls on the processes that can access the data, you can add extra controls by using profiles. Profiles set requirements on a process other than just having the correct fingerprint. They can include such things as process owner and group, required open files, and the current working directory.

A profile is generated by using the following command:

```
$ usr/sbin/navencrypt-profile --pid=<pid>
```

The output, by default, will be displayed on the screen. You can redirect the output to a file using the `>` or `>>` redirect operators. You can then edit the JSON output in the file to remove lines you do not want. By default, the profile includes the UID, the short name of the binary or script (identified as `comm`), and the full command line of the running process (including any parameters passed). You can generate information by using one of these flags:

- `-c, --with-cwd`
Output the current working directory
- `-e, --with-egid`
Output the egid
- `-g, --with-gid`
Output the gid
- `-u, --with-euid`
Output the euid

Example output from the `navencrypt-profile` command:

```
{  
  "uid": "0",  
  "comm": "NetworkManager",  
  "cmdline": "NetworkManager -pid-file=/var/run/NetworkManager/NetworkManager.pid",  
  "gid": "0",  
  "cwd": "/",  
  "fd0": "/dev/null",  
  "fd1": "/dev/null",  
  "fd2": "/dev/null"  
}
```

Some distributions do not support `euid` and `guid`. Make sure that your profile file is correct by executing the following command to verify the expected IDs:

```
$ ps -p <pid_of_process> -o euid,egid
```

If `cmdline` parameters are variable, such as appending a process start timestamp to a filename, then the process profile will not match on subsequent restarts of the process because the current profile will have an updated timestamp and access will be denied by the ACL. You can mark those parameters as variable inside the profile file. For example, if the `cmdline` of a process is something like this:

```
"cmdline": "NetworkManager -pid-file=/var/run/NetworkManager/NetworkManager.pid \  
-logfile=/var/log/NetworkManager/log-20130808152300.log"
```

Configuring Encryption

Where `log-20130505122300.log` is a variable `cmdline` parameter, before adding the process profile to the ACL, edit the process profile file and use `##` to specify that a particular parameter is variable:

```
"cmdline": "NetworkManager -pid-file=/var/run/NetworkManager/NetworkManager.pid
-logfile=##"
```

With the above configuration, the ACL will allow any value for the `-logfile` `cmdline` parameter.

To enable a profile in the ACL, use the additional parameter `--profile-file=<filename>` when adding the rule to the ACL:

```
$ sudo /usr/sbin/navencrypt acl --add --rule="ALLOW @mysql * /usr/sbin/mysqld" \
--profile-file=/path/to/profile/file
```

To display the profile portion of the rules, use the `--all` parameter with `navencrypt acl --list`:

```
$ sudo /usr/sbin/navencrypt acl --list --all
Type MASTER passphrase:
# - Type Category Path Profile Process
1 ALLOW @mysql * YES /usr/sbin/mysqld
PROFILE:
{"uid": "120", "comm": "mysqld", "cmdline": "mysqld" }
```

Maintaining Cloudera Navigator Encrypt

Manually Backing Up Navigator Encrypt

It is recommended that you back up Navigator Encrypt configuration directory after installation, and again after any configuration updates.

1. To manually back up the Navigator Encrypt configuration directory (`/etc/navencrypt`):

```
$ zip -r --encrypt nav-encrypt-conf.zip /etc/navencrypt
```

The `--encrypt` option prompts you to create a password used to encrypt the zip file. This password is also required to decrypt the file. Ensure that you protect the password by storing it in a secure location.

2. Move the backup file (`nav-encrypt-conf.zip`) to a secure location.



Warning: Failure to back up the configuration directory makes your backed-up encrypted data unrecoverable in the event of data loss.

Access Modes

Navigator Encrypt provides three different access modes:

- **Enforcing** is the default mode in which Navigator Encrypt validates access from all processes against the ACL. To protect your data, enforcing mode must be enabled.
- **Permissive** mode causes `action="denied"` messages to be logged in `dmesg`. It **does not** prevent access to the encrypted data. This mode is a dry-run feature to run and build ACL rules.
- **Admin** mode, as well as permissive mode, **does not** prevent access to the encrypted data. It allows any process to access the information because the ACL rules are not validated against the process. Admin mode does not cause `action="denied"` messages to be logged in `dmesg`.

To view the current access mode, run the following command:

```
$ sudo /usr/sbin/navencrypt status -d
```



Note: The `navencrypt status` command reports that the `navencrypt` module is not running if no volumes are encrypted or the kernel module is not loaded.

To change the access mode, use the following command:

```
$ sudo /usr/sbin/navencrypt set --mode={enforcing|permissive|admin}
```

You cannot change the Navigator Encrypt access mode unless the Navigator Encrypt module is running. To view the status of the Navigator Encrypt module, run `navencrypt status --module`.

To start the Navigator Encrypt module there must be at least one active mount-point. To verify the mount-points status, run the following command:

```
$ sudo /etc/init.d/navencrypt-mount status
```

For RHEL 7, use `systemctl` instead:

```
$ sudo systemctl status navencrypt-mount
```

Changing and Verifying the Master Key

You can perform two operations with the `navencrypt key` command: `change` and `verify`.

You can verify a key against the Navigator Encrypt module, the Navigator Key Trustee server, or both. For example:

```
$ sudo /usr/sbin/navencrypt key --verify
$ sudo /usr/sbin/navencrypt key --verify --only-module
$ sudo /usr/sbin/navencrypt key --verify --only-keytrustee
```



Note: The `navencrypt key` command fails if no volumes are encrypted or the kernel module is not loaded.

The master key can be changed in the event that another key-type authentication mechanism or a new master key is required. Valid master key types are single-passphrase, dual-passphrase, and RSA key files. To change the master key type, issue the following command and follow the interactive console:

```
$ sudo /usr/sbin/navencrypt key --change
```

You can use the `--trustees`, `--votes`, and `--recoverable` options for the new key as described in [Table 24: Registration Options](#) on page 317.

Updating ACL Fingerprints

All rules reference a process fingerprint (a SHA256 digest) that is used to authenticate the process into the file system. If the filesystem detects a fingerprint that is different from the one stored in the ACL, the Linux process is denied access and treated as an untrusted process.

Occasionally this process fingerprint must be updated, such as when software is upgraded. When the fingerprint must be updated, the Navigator Encrypt administrator re-authenticates the process on the ACL by executing the command `navencrypt acl --update`.

The following example demonstrates how to determine when a process fingerprint has been changed and must be updated:

```
$ sudo /usr/sbin/navencrypt acl --list
Type MASTER passphrase:
# - Type Category Path Profile Process
1 !! ALLOW @mysql * /usr/sbin/mysqld
```

Configuring Encryption

```
2 ALLOW @log * /usr/sbin/mysqld
3 !! ALLOW @apache * /usr/lib/apache2/mpm-prefork/
```

In the example above, the double exclamation (!!) characters indicate that a process fingerprint has changed and must be updated. Similarly, double E (EE) characters indicate a process read error. This error can be caused by a process that does not exist or that has permission issues.



Note:

For RHEL-compatible OSES, the `prelink` application may also be responsible for ACL fingerprint issues. Prelinking is intended to speed up a system by reducing the time a program needs to begin. Cloudera highly recommends disabling any automated `prelink` jobs, which are enabled by default in some systems. As an alternative, Cloudera recommends that you integrate a manual `prelink` run into your existing change control policies to ensure minimal downtime for applications accessing encrypted data.

To disable prelinking, modify the `/etc/sysconfig/prelink` file and change `PRELINKING=yes` to `PRELINKING=no`. Then, run the `/etc/cron.daily/prelink` script as root. Once finished, automatic prelinking is disabled.

For more information about how prelinking affects your system, see [prelink](#).

Managing Mount Points

The `/etc/init.d/navencrypt-mount` command mounts all mount points that were registered with the `navencrypt-prepare` command and are listed in the `/etc/navencrypt/ztab` file. The possible operations are:

- start
- stop
- status
- restart

For RHEL 7, use `systemctl [start|stop|status|restart] navencrypt-mount`.



Note: Do not manually unmount the encryption mount point (for example, using `umount`). If you do so, you must manually close the `dm-crypt` device using the following procedure:

1. Run `dmsetup table` to list the `dm-crypt` devices.
2. Run `cryptsetup luksClose <device_name>` to close the device for the unmounted mount point.

When executing the `stop` operation, the encrypted mount point is unmounted, and your data becomes inaccessible.

The following example shows how to execute `navencrypt-mount status` with some inactive mount points:

```
$ sudo /etc/init.d/navencrypt-mount status
```

The following example shows how to execute the `navencrypt-mount stop` command:

```
$ sudo /etc/init.d/navencrypt-mount stop
```

The following example shows how to execute the `navencrypt-mount start` command:

```
$ sudo /etc/init.d/navencrypt-mount start
```

Here is an example command used to manually mount a directory:

```
$ sudo /usr/sbin/mount.navencrypt /path/to_encrypted_data/ /path/to/mountpoint
```


This command can be executed only if the `navencrypt-prepare` command was previously executed.

Allowing Access to a Single Operation

Sometimes running a single operation on encrypted data is required, such as listing the encrypted files or copying a particular file. The `navencrypt exec` command provides an efficient method to accomplish this. In the following command example, permission is denied when the command is run without `navencrypt exec`:

```
$ ls /mnt/db_encrypted/mysql/
ls: cannot open directory /mnt/db_encrypted/mysql/: Permission denied
```

If the `ls` command is executed with `navencrypt exec`, access is allowed one time only:

```
$ sudo /usr/sbin/navencrypt exec "ls -l /mnt/db_encrypted/mysql/"
```

Navigator Encrypt Kernel Module Setup

If the kernel headers were not installed on your host, or if the wrong version of the kernel headers were installed, the Navigator Encrypt module was not built at installation time. To avoid reinstalling the system, install the correct headers and execute the `navencrypt-module-setup` command. This attempts to build the module and install it.

This method is also an efficient way to install any new Navigator Encrypt module feature or fix without otherwise modifying your current Navigator Encrypt environment.

Navigator Encrypt Configuration Directory Structure

The file and directory structure of `/etc/navencrypt` is as follows:

```
$ tree /etc/navencrypt/
/etc/navencrypt/
  control -> /etc/navencrypt/jSpi9SM65xUIIhraulNn8ZXmQhrrQ9e363EUz8HKiRs
  jSpi9SM65xUIIhraulNn8ZXmQhrrQ9e363EUz8HKiRs
  rules
  ztab
locust
  keytrustee
  clientname
  deposits
    dev.loop0
    media.31E5-79B9locustlocust[system ~]# . /etc/*release[system ~]# . /etc/*release

    mnt.a
    mnt.encrypted
    mnt.tomount
  pubring.gpg
  pubring.gpg~
  random_seed
  secring.gpg
  trustdb.gpg
  ztrustee.conf
```

The following files and folders are part of the created file structure:

- `control`
File that saves information about the mount points and corresponding Navigator Key Trustee keys.
- `rules`
File that contains the ACL rules. It is encrypted with the user-provided master key.
- `ztab`
File that contains information about all the mount points and their encryption type.



Important: Use caution when editing the `/etc/navencrypt/ztab` file. Entries are tab-separated (not space-separated). The `ztab` file must not contain empty lines.

- `keytrustee`

Directory where Navigator Key Trustee GPG keys are stored. These are generated during `navencrypt register` operations.

- `keytrustee/deposits`

Directory where the Navigator Encrypt keys are saved. These are encrypted with the user-provided master key.

Every mount point has an internal randomly-generated encryption passphrase.

Upgrading Navigator Encrypt Hosts

See [Best Practices for Upgrading Navigator Encrypt Hosts](#) for considerations when upgrading operating systems (OS) and kernels on hosts that have Navigator Encrypt installed.

Configuring Encryption for Data Spills

Certain CDH services have the ability to encrypt data that lives temporarily on the local filesystem outside HDFS. This usually includes data that may *spill* to disk when operations are too memory intensive and the service exceeds its allotted memory limit on a host. You can enable on-disk spill encryption for the following services:

MapReduce v2 (YARN)

MapReduce v2 allows you to encrypt intermediate files generated during encrypted shuffle and in case of data spills during the map and reduce stages. This can be enabled by setting the following properties in `mapred-site.xml`.

<code>mapreduce.job.encrypted-intermediate-data</code>	Enable or disable encryption for intermediate MapReduce spills. Default: <code>false</code>
<code>mapreduce.job.encrypted-intermediate-data-key-size-bits</code>	The key length used to encrypt data spilled to disk. Default: <code>128</code>
<code>mapreduce.job.encrypted-intermediate-data.buffer.kb</code>	The buffer size in kb for stream written to disk after encryption. Default: <code>128</code>



Note: Enabling encryption for intermediate data spills will restrict the number of attempts for a job to 1.

HBase

HBase does not write data outside HDFS, and does not require spill encryption.

Impala

Impala allows certain memory-intensive operations to be able to write temporary data to disk in case these operations come close to exceeding their memory limit on a host. For details, read [SQL Operations that Spill to Disk](#). To enable disk spill encryption in Impala:

1. Go to the Cloudera Manager Admin Console.

2. Click the **Configuration** tab.
3. Select **Scope > Impala Daemon**.
4. Select **Category > Security**.
5. Check the checkbox for the **Disk Spill Encryption** property.
6. Click **Save Changes** to commit the changes.

Hive

Hive jobs occasionally write data temporarily to local directories. If you enable HDFS encryption, then you must ensure that the following intermediate local directories are also protected:

- **LOCALSCRATCHDIR:** The MapJoin optimization in Hive writes HDFS tables to a local directory and then uploads them to the distributed cache. To ensure these files are encrypted, either disable MapJoin by setting `hive.auto.convert.join` to `false`, or encrypt the *local* Hive Scratch directory (`hive.exec.local.scratchdir`) using [Cloudera Navigator Encrypt](#).
- **DOWNLOADED_RESOURCES_DIR:** JARs that are added to a user session and stored in HDFS are downloaded to `hive.downloaded.resources.dir` on the HiveServer2 local filesystem. To encrypt these JAR files, configure [Cloudera Navigator Encrypt](#) to encrypt the directory specified by `hive.downloaded.resources.dir`.
- **NodeManager Local Directory List:** Hive stores JARs and MapJoin files in the distributed cache. To use MapJoin or encrypt JARs and other resource files, the `yarn.nodemanager.local-dirs` YARN configuration property must be configured to a set of encrypted local directories on all nodes.

For more information on Hive behavior with HDFS encryption enabled, see [Using HDFS Encryption with Hive](#).

Flume

Flume supports on-disk encryption for log files written by the Flume file channels. See [Using an On-disk Encrypted File Channel](#) on page 339.

Using an On-disk Encrypted File Channel

Flume supports on-disk encryption of data on the local disk. To implement this:

- Generate an encryption key to use for the Flume Encrypted File Channel
- Configure on-disk encryption by setting parameters in the `flume.conf` file



Important:

Flume on-disk encryption operates with a maximum strength of 128-bit AES encryption unless the JCE unlimited encryption cryptography policy files are installed. Please see this Oracle document for information about enabling strong cryptography:

<http://www.oracle.com/technetwork/java/javase/downloads/jce-6-download-429243.html>

Consult your security organization for guidance on the acceptable strength of your encryption keys. Cloudera has tested with AES-128, AES-192, and AES-256.

Generating Encryption Keys

Use the `keytool` program included with the Oracle JDK to create the AES encryption keys for use with Flume.

The command to generate a 128-bit key that uses the same password as the key store password is:

```
keytool -genseckey -alias key-1 -keyalg AES -keysize 128 -validity 9000 \
-keystore test.keystore -storetype jceks \
-storepass keyStorePassword
```

Configuring Encryption

The command to generate a 128-bit key that uses a different password from that used by the key store is:

```
keytool -genseckey -alias key-0 -keypass keyPassword -keyalg AES \  
-keysize 128 -validity 9000 -keystore test.keystore \  
-storetype jceks -storepass keyStorePassword
```

The key store and password files can be stored anywhere on the file system; both files should have `flume` as the owner and `0600` permissions.

Please note that `-keysize` controls the strength of the AES encryption key, in bits; 128, 192, and 256 are the allowed values.

Configuration

Flume on-disk encryption is enabled by setting parameters in the `/etc/flume-ng/conf/flume.conf` file.

Basic Configuration

The first example is a basic configuration with an alias called `key-0` that uses the same password as the key store:

```
agent.channels.ch-0.type = file  
agent.channels.ch-0.capacity = 10000  
agent.channels.ch-0.encryption.cipherProvider = AESCTRNOPADDING  
agent.channels.ch-0.encryption.activeKey = key-0  
agent.channels.ch-0.encryption.keyProvider = JCEKSFILE  
agent.channels.ch-0.encryption.keyProvider.keyStoreFile = /path/to/my.keystore  
agent.channels.ch-0.encryption.keyProvider.keyStorePasswordFile =  
/path/to/my.keystore.password  
agent.channels.ch-0.encryption.keyProvider.keys = key-0
```

In the next example, `key-0` uses its own password which may be different from the key store password:

```
agent.channels.ch-0.type = file  
agent.channels.ch-0.capacity = 10000  
agent.channels.ch-0.encryption.cipherProvider = AESCTRNOPADDING  
agent.channels.ch-0.encryption.activeKey = key-0  
agent.channels.ch-0.encryption.keyProvider = JCEKSFILE  
agent.channels.ch-0.encryption.keyProvider.keyStoreFile = /path/to/my.keystore  
agent.channels.ch-0.encryption.keyProvider.keyStorePasswordFile =  
/path/to/my.keystore.password  
agent.channels.ch-0.encryption.keyProvider.keys = key-0  
agent.channels.ch-0.encryption.keyProvider.keys.key-0.passwordFile =  
/path/to/key-0.password
```

Changing Encryption Keys Over Time

To modify the key, modify the configuration as shown below. This example shows how to change the configuration to use `key-1` instead of `key-0`:

```
agent.channels.ch-0.type = file  
agent.channels.ch-0.capacity = 10000  
agent.channels.ch-0.encryption.cipherProvider = AESCTRNOPADDING  
agent.channels.ch-0.encryption.activeKey = key-1  
agent.channels.ch-0.encryption.keyProvider = JCEKSFILE  
agent.channels.ch-0.encryption.keyProvider.keyStoreFile = /path/to/my.keystore  
agent.channels.ch-0.encryption.keyProvider.keyStorePasswordFile =  
/path/to/my.keystore.password  
agent.channels.ch-0.encryption.keyProvider.keys = key-0 key-1
```

The same scenario except that `key-0` and `key-1` have their own passwords is shown here:

```
agent.channels.ch-0.type = file  
agent.channels.ch-0.capacity = 10000  
agent.channels.ch-0.encryption.cipherProvider = AESCTRNOPADDING  
agent.channels.ch-0.encryption.activeKey = key-1  
agent.channels.ch-0.encryption.keyProvider = JCEKSFILE
```

```
agent.channels.ch-0.encryption.keyProvider.keyStoreFile = /path/to/my.keystore
agent.channels.ch-0.encryption.keyProvider.keyStorePasswordFile =
/path/to/my.keystore.password
agent.channels.ch-0.encryption.keyProvider.keys = key-0 key-1
agent.channels.ch-0.encryption.keyProvider.keys.key-0.passwordFile =
/path/to/key-0.password
agent.channels.ch-0.encryption.keyProvider.keys.key-1.passwordFile =
/path/to/key-1.password
```

Troubleshooting

If the unlimited strength JCE policy files are not installed, an error similar to the following is printed in the `flume.log`:

```
07 Sep 2012 23:22:42,232 ERROR [lifecycleSupervisor-1-0]
(org.apache.flume.channel.file.encryption.AESCTRNoPaddingProvider.getCipher:137) - Unable
to load key using transformation: AES/CTR/NoPadding; Warning: Maximum allowed key length
= 128 with the available JCE security policy files. Have you installed the JCE unlimited
strength jurisdiction policy files?
java.security.InvalidKeyException: Illegal key size
at javax.crypto.Cipher.a(DashoA13*..)
at javax.crypto.Cipher.a(DashoA13*..)
at javax.crypto.Cipher.a(DashoA13*..)
at javax.crypto.Cipher.init(DashoA13*..)
at javax.crypto.Cipher.init(DashoA13*..)
at
org.apache.flume.channel.file.encryption.AESCTRNoPaddingProvider.getCipher(AESCTRNoPaddingProvider.java:120)
at
org.apache.flume.channel.file.encryption.AESCTRNoPaddingProvider.access$200(AESCTRNoPaddingProvider.java:35)
at
org.apache.flume.channel.file.encryption.AESCTRNoPaddingProvider$AESCTRNoPaddingDecryptor.<init>(AESCTRNoPaddingProvider.java:94)
at
org.apache.flume.channel.file.encryption.AESCTRNoPaddingProvider$AESCTRNoPaddingDecryptor.<init>(AESCTRNoPaddingProvider.java:91)
at
org.apache.flume.channel.file.encryption.AESCTRNoPaddingProvider$DecryptorBuilder.build(AESCTRNoPaddingProvider.java:66)
at
org.apache.flume.channel.file.encryption.AESCTRNoPaddingProvider$DecryptorBuilder.build(AESCTRNoPaddingProvider.java:62)
at
org.apache.flume.channel.file.encryption.CipherProviderFactory.getDecrypter(CipherProviderFactory.java:47)
at org.apache.flume.channel.file.LogFileV3$SequentialReader.<init>(LogFileV3.java:257)
at
org.apache.flume.channel.file.LogFileFactory.getSequentialReader(LogFileFactory.java:110)
at org.apache.flume.channel.file.ReplayHandler.replayLog(ReplayHandler.java:258)
at org.apache.flume.channel.file.Log.replay(Log.java:339)
at org.apache.flume.channel.file.FileChannel.start(FileChannel.java:260)
at
org.apache.flume.lifecycle.LifecycleSupervisor$MonitorRunnable.run(LifecycleSupervisor.java:236)
at java.util.concurrent.Executors$RunnableAdapter.call(Executors.java:441)
at java.util.concurrent.FutureTask$Sync.innerRunAndReset(FutureTask.java:317)
at java.util.concurrent.FutureTask.runAndReset(FutureTask.java:150)
at
java.util.concurrent.ScheduledThreadPoolExecutor$ScheduledFutureTask.access$101(ScheduledThreadPoolExecutor.java:98)
at
java.util.concurrent.ScheduledThreadPoolExecutor$ScheduledFutureTask.runPeriodic(ScheduledThreadPoolExecutor.java:180)
at
java.util.concurrent.ScheduledThreadPoolExecutor$ScheduledFutureTask.run(ScheduledThreadPoolExecutor.java:204)
at java.util.concurrent.ThreadPoolExecutor$Worker.runTask(ThreadPoolExecutor.java:886)
at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:908)
at java.lang.Thread.run(Thread.java:662)
```

Configuring Encrypted HDFS Data Transport

This topic describes how to configure encrypted HDFS data transport using both Cloudera Manager, and the command line.

You must enable Kerberos before configuring encrypted HDFS data transport. See [Configuring Authentication](#) on page 49 for instructions.

Configuring Encryption

Using Cloudera Manager

Minimum Required Role: [Full Administrator](#)

To enable encryption of data transferred between DataNodes and clients, and among DataNodes, proceed as follows:

1. [Enable Hadoop security using Kerberos.](#)
2. Select the HDFS service.
3. Click the **Configuration** tab.
4. Select **Scope > HDFS (Service Wide)**
5. Select **Category > Security.**
6. Configure the following properties: (You can type the property name in the **Search** box to locate the property.)

Property	Description
Enable Data Transfer Encryption	Check this field to enable wire encryption.
Data Transfer Encryption Algorithm	Optionally configure the algorithm used to encrypt data.
Hadoop RPC Protection	Select privacy.

7. Click **Save Changes.**
8. Restart the HDFS service.

Using the Command Line



Important:

- You can use either Cloudera Manager or the following command-line instructions to complete this configuration.
- This information applies specifically to CDH 5.6.x. If you use an earlier version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

To enable encrypted data transport using the command line, proceed as follows:

1. Enable Kerberos authentication, following [these instructions](#).
2. Set the optional RPC encryption by setting `hadoop.rpc.protection` to "privacy" in the `core-site.xml` file in both client and server configurations.



Note:

If RPC encryption is not enabled, transmission of other HDFS data is also insecure.

3. Set `dfs.encrypt.data.transfer` to `true` in the `hdfs-site.xml` file on all server systems.
4. Restart all daemons.

Configuring Encrypted HBase Data Transport

This topic describes how to configure encrypted HBase data transport using both Cloudera Manager and the command line.

Using Cloudera Manager

Minimum Required Role: [Full Administrator](#)

To enable encryption of data transferred between HBase masters and RegionServers and between RegionServers and clients:

1. [Enable Hadoop security using Kerberos.](#)
2. [Configure Kerberos authentication for HBase.](#)
3. Select the HBase service.
4. Click the **Configuration** tab.
5. Select **Scope > HBase (Service Wide)**
6. Select **Category > Security**.
7. Search for the HBase Transport Security property. Select **privacy** to enable secure RPC transport. There are three settings: `authentication`, `integrity`, and `privacy`. Each level includes the prior ones. Selecting `privacy` also enables `authentication` and `integrity`.
8. Click **Save Changes**.
9. Restart the HBase service.

Using the Command Line



Important:

- You can use either Cloudera Manager or the following command-line instructions to complete this configuration.
- This information applies specifically to CDH 5.6.x. If you use an earlier version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

To enable encrypted data transport using the command line, proceed as follows:

1. Enable [Hadoop Security using Kerberos](#).
2. Enable [HBase security using Kerberos](#).
3. Enable RPC encryption by setting `hbase.rpc.protection` to "privacy" in the `hbase-site.xml` file.
4. Restart all daemons.

Configuring Authorization

Authorization is concerned with who or what has access or control over a given resource or service. Since Hadoop merges together the capabilities of multiple varied, and previously separate IT systems as an enterprise data hub that stores and works on all data within an organization, it requires multiple authorization controls with varying granularities. In such cases, Hadoop management tools simplify setup and maintenance by:

- Tying all users to groups, which can be specified in existing LDAP or AD directories.
- Providing role-based access control for similar interaction methods, like batch and interactive SQL queries. For example, Apache Sentry permissions apply to Hive (HiveServer2) and Impala.

CDH currently provides the following forms of access control:

- Traditional POSIX-style permissions for directories and files, where each directory and file is assigned a single owner and group. Each assignment has a basic set of permissions available; file permissions are simply read, write, and execute, and directories have an additional permission to determine access to child directories.
- [Extended Access Control Lists](#) (ACLs) for HDFS that provide fine-grained control of permissions for HDFS files by allowing you to set different permissions for specific named users or named groups.
- Apache HBase uses ACLs to authorize various operations (`READ`, `WRITE`, `CREATE`, `ADMIN`) by column, column family, and column family qualifier. HBase ACLs are granted and revoked to both users and groups.
- Role-based access control with [Apache Sentry](#). As of Cloudera Manager 5.1.x, Sentry permissions can be configured using either policy files or the database-backed Sentry service.
 - The Sentry service is the preferred way to set up Sentry permissions. See [The Sentry Service](#) on page 359 for more information.
 - For the policy file approach to configuring Sentry, see [Sentry Policy File Authorization](#) on page 392.



Important: Cloudera does not support Apache Ranger or Hive's native authorization frameworks for configuring access control in Hive. Use the Cloudera-supported Apache Sentry instead.

Cloudera Manager User Roles

Access to Cloudera Manager features is controlled by user accounts. For more information about user accounts, see [Cloudera Manager User Accounts](#). Among the properties of a user account is the *user role*, which determines the Cloudera Manager features visible to the user and the actions the user can perform. All the tasks in the Cloudera Manager documentation indicate which role is required to perform the task.



Note: The full set of roles are available with Cloudera Enterprise; Cloudera Express supports only the Read-Only and Full Administrator roles. When a Cloudera Enterprise Data Hub Edition trial license expires, only users with Read-Only and Full Administrator roles are allowed to log in. A Full Administrator must change the role of any other user to Read-Only or Full Administrator before that user can log in.

User Roles

A Cloudera Manager user account can be assigned one of the following roles with associated permissions:

- **Auditor**
 - View configuration and monitoring information in Cloudera Manager.
 - View audit events.
- **Read-Only**

- View configuration and monitoring information in Cloudera Manager.
- View service and monitoring information.
- View events and logs.
- View replication jobs and snapshot policies.
- View YARN applications and Impala queries.

The Read-Only role does not allow the user to:

- Add services or take any actions that affect the state of the cluster.
- Use the HDFS file browser.
- Use the HBase table browser.
- Use the Solr Collection Statistics browser.

- **Dashboard**

- Create, edit, or remove dashboards that belong to the user.
- Add an existing chart or create a new chart to add to a dashboard that belongs to the user.

Additionally, the Dashboard role can perform the same tasks as the Read-Only role.

- **Limited Operator**

- View configuration and monitoring information in Cloudera Manager.
- View service and monitoring information.
- Decommission hosts (except hosts running Cloudera Management Service roles).

The Limited Operator role does not allow the user to add services or take any other actions that affect the state of the cluster.

- **Operator**

- View configuration and monitoring information in Cloudera Manager.
- View service and monitoring information.
- Stop, start, and restart clusters, services (except the Cloudera Management Service), and roles.
- Decommission and recommission hosts (except hosts running Cloudera Management Service roles).
- Decommission and recommission roles (except Cloudera Management Service roles).
- Start, stop, and restart KMS.

The Operator role does not allow the user to add services, roles, or hosts, or take any other actions that affect the state of the cluster.

- **Configurator**

- View configuration and monitoring information in Cloudera Manager.
- Perform all Operator operations.
- Configure services (except the Cloudera Management Service).
- Enter and exit maintenance mode.
- Manage dashboards (including Cloudera Management Service dashboards).
- Start, stop, and restart KMS

- **Cluster Administrator** - View all data and perform all actions *except* the following:

- Administer Cloudera Navigator.
- View replication schedules and snapshot policies.
- View audit events.
- Manage user accounts and configuration of external authentication.
- Manage Full Administrator accounts.
- Configure HDFS encryption, administer Key Trustee Server, and manage encryption keys.
- Use the HDFS file browser, the HBase table browser, and the Solr Collection browser.

- **BDR Administrator**
 - View configuration and monitoring information in Cloudera Manager.
 - View service and monitoring information.
 - Perform replication and define snapshot operations.
 - Use the HDFS file browser, the HBase table browser, and the Solr Collection browser.
- **Navigator Administrator**
 - View configuration and monitoring information in Cloudera Manager.
 - View service and monitoring information.
 - Administer Cloudera Navigator.
 - View audit events.
 - Use the HDFS file browser, the HBase table browser, and the Solr Collection browser.
- **User Administrator**
 - View configuration and monitoring information in Cloudera Manager.
 - View service and monitoring information.
 - Manage user accounts and configuration of external authentication.
 - Use the HDFS file browser, the HBase table browser, and the Solr Collection browser.
- **Key Administrator**
 - View configuration and monitoring information in Cloudera Manager.
 - Configure HDFS encryption, administer Key Trustee Server, and manage encryption keys.
 - Start, stop, and restart KMS
 - Configure KMS ACLs
 - Use the HDFS file browser, the HBase table browser, and the Solr Collection browser.
- **Full Administrator** - Full Administrators have permissions to view all data and do all actions, including reconfiguring and restarting services, and administering other users.

Determining the Role of the Currently Logged in User

1. Click the logged-in username at the far right of the top navigation bar. The role displays under the username. For example:



Full Administrator

Removing the Full Administrator User Role

Minimum Required Role: [User Administrator](#) (also provided by **Full Administrator**)

In some organizations, security policies may prohibit the use of the Full Administrator role. The Full Administrator role is created during Cloudera Manager installation, but you can remove it as long as you have at least one remaining user account with User Administrator privileges.

To remove the Full Administrator user role, perform the following steps.

1. Add at least one user account with User Administrator privileges, or ensure that at least one such user account already exists.
2. Ensure that there is only a single user account with Full Administrator privileges.
3. While logged in as the single remaining Full Administrator user, select your own user account and either delete it or assign it a new user role.



Warning: After you delete the last Full Administrator account, you will be logged out immediately and will not be able to log in unless you have access to another user account. Also, it will no longer be possible to create or assign Full Administrators.

A consequence of removing the Full Administrator role is that some tasks may require collaboration between two or more users with different user roles. For example:

- If the machine that the Cloudera Navigator roles are running on needs to be replaced, the Cluster Administrator will want to move all the roles running on that machine to a different machine. The Cluster Administrator can move any non-Navigator roles by deleting and re-adding them, but would need a Navigator Administrator to perform the stop, delete, add, and start actions for the Cloudera Navigator roles.
- In order to take HDFS snapshots, snapshots must be enabled on the cluster by a Cluster Administrator, but the snapshots themselves must be taken by a BDR Administrator.

Cloudera Navigator Data Management Component User Roles

User roles determine the Cloudera Navigator features visible to the user and the actions the user can perform.

The menus displayed in the upper right indicate the user's access to Cloudera Navigator features, as determined by the roles associated with the [user's LDAP or Active Directory groups](#). For example, a user that belongs to a group with the Full Administrator role will see the **Search, Audits, Analytics, Policies, and Administration** tabs, while a user that belongs to a group with the Policy Administrator role will only see the **Search, Analytics (metadata), and Policies** tabs.

User Roles

A Cloudera Navigator data management component user account can be assigned one of the following user roles with associated permissions:

- **Auditing Viewer** - View audit events and audit analytics and create audit reports.
- **Full Administrator** - Full access, including role assignments to groups.
- **Lineage Viewer** - Search for entities, view metadata, and view lineage and metadata analytics.
- **Metadata Administrator** - Search for entities, view metadata, view lineage, view metadata analytics, and edit metadata.
- **Policy Administrator** - View, create, update, delete metadata and metadata policies, configure and perform command actions, and view metadata analytics.
- **Policy Viewer** - View metadata policies.
- **User Administrator** - Administer role assignments to groups.

Determining the Roles of the Currently Logged in User

To display the Cloudera Navigator user roles for the currently logged-in user:

1. In the upper right, select **username > My roles**. The **Roles** pop-up window displays all roles assigned to the LDAP or Active Directory groups to which the current user belongs.
2. Click **Close** to dismiss the window.

HDFS Extended ACLs

HDFS supports POSIX Access Control Lists (ACLs), as well as the traditional POSIX permissions model already supported. ACLs control access of HDFS files by providing a way to set different permissions for specific named users or named groups. They enhance the traditional permissions model by allowing users to define access control for arbitrary combinations of users and groups instead of a single owner/user or a single group.

Enabling HDFS Access Control Lists

By default, HDFS access control lists (ACLs) are disabled on a cluster. You can enable them using either Cloudera Manager or the command line.

Default ACLs are applied only to a directory (not to files), and have no direct effect on permission checks. Rather, they define the ACL that newly-created child files and directories receive automatically.



Important: Ensure that all users and groups resolve on the NameNode for ACLs to work as expected.

Enabling HDFS ACLs Using Cloudera Manager

1. Go to the Cloudera Manager Admin Console and navigate to the **HDFS** service.
2. Click the **Configuration** tab.
3. Select **Scope** > **Service_name (Service-Wide)**
4. Select **Category** > **Security**
5. Locate the **Enable Access Control Lists** property and select its checkbox to enable HDFS ACLs.
6. Click **Save Changes** to commit the changes.

Enabling HDFS ACLs Using the Command Line

To enable ACLs using the command line, set the `dfs.namenode.acls.enabled` property to `true` in the NameNode's `hdfs-site.xml`.

```
<property>
<name>dfs.namenode.acls.enabled</name>
<value>true</value>
</property>
```

Commands

To set and get file access control lists (ACLs), use the file system shell commands, `setfacl` and `getfacl`.

getfacl

```
hdfs dfs -getfacl [-R] <path>
```

```
<!-- COMMAND OPTIONS
```

```
<path>: Path to the file or directory for which ACLs should be listed.
```

```
-R: Use this option to recursively list ACLs for all files and directories.
```

```
-->
```

Examples:

```
<!-- To list all ACLs for the file located at /user/hdfs/file -->
```

```
hdfs dfs -getfacl /user/hdfs/file
```

```
<!-- To recursively list ACLs for /user/hdfs/file -->
```

```
hdfs dfs -getfacl -R /user/hdfs/file
```

setfacl

```
hdfs dfs -setfacl [-R] [-b|-k -m|-x <acl_spec> <path>]|[--set <acl_spec> <path>]
```

```
<!-- COMMAND OPTIONS
```

```
<path>: Path to the file or directory for which ACLs should be set.
```

```
-R: Use this option to recursively list ACLs for all files and directories.
```

```
-b: Revoke all permissions except the base ACLs for user, groups and others.
```

```
-k: Remove the default ACL.
```

```
-m: Add new permissions to the ACL with this option. Does not affect existing permissions.
```

```
-x: Remove only the ACL specified.
<acl_spec>: Comma-separated list of ACL permissions.
--set: Use this option to completely replace the existing ACL for the path specified.
       Previous ACL entries will no longer apply.
-->
```

Examples:

```
<!-- To give user ben read & write permission over /user/hdfs/file -->
hdfs dfs -setfacl -m user:ben:rw- /user/hdfs/file

<!-- To remove user alice's ACL entry for /user/hdfs/file -->
hdfs dfs -setfacl -x user:alice /user/hdfs/file

<!-- To give user hadoop read & write access, and group or others read-only access -->
hdfs dfs -setfacl --set user::rw-,user:hadoop:rw-,group::r--,other::r-- /user/hdfs/file
```

For more information on using HDFS ACLs, see the [HDFS Permissions Guide](#) on the Apache website.

HDFS Extended ACL Example

This example demonstrates how a user ("alice"), shares folder access with colleagues from another team ("hadoopdev"), so that the hadoopdev team can collaborate on the content of that folder; this is accomplished by updating the default extended ACL of that directory:

1. Make the files and sub-directories created within the content directory readable by team "hadoopdev":

```
$ hdfs dfs -setfacl -m group:hadoopdev:r-x /project
```

2. Set the default ACL setting for the parent directory:

```
$ hdfs dfs -setfacl -m default:group:hadoopdev:r-x /project
```

3. Create a sub-directory for the content you wish to share:

```
$ hdfs dfs -mkdir /project/dev
```

4. Inspect the new sub-directory ACLs to verify that HDFS has applied the new default values:

```
$ hdfs dfs -getfacl -R /project

file: /project
owner: alice
group: appdev
user::rwx
group::r-x
other::r-x
default:user::rwx
default:group::r-x
default:group:hadoopdev:r-x
default:mask::r-x
default:other::r-x

file: /project/dev
owner: alice
group: appdev
user::rwx
group::r-x
group:hadoopdev:r-x
mask:r-x
other:r-x
default:user::rwx
default:group::r-x
default:group:hadoopdev:r-x
default:mask::r-x
default:other::r-x
```



Note: At the time it is created, the default ACL is copied from the parent directory to the child directory. Subsequent changes to the parent directory default ACL do not change the ACLs of the existing child directories.

Configuring LDAP Group Mappings



Important:

- Cludera strongly recommends *against* using Hadoop's `LdapGroupsMapping` provider. `LdapGroupsMapping` should only be used in cases where OS-level integration is not possible. Production clusters require an identity provider that works well with all applications, not just Hadoop. Hence, often the preferred mechanism is to use tools such as SSSD, VAS or Centrify to replicate LDAP groups.
- Cludera does not support the use of Winbind in production environments. Winbind uses an inefficient approach to user/group mapping, which may lead to low performance or cluster failures as the size of the cluster, and the number of users and groups increases.

Irrespective of the mechanism used, user/group mappings must be applied consistently across all cluster hosts for ease with maintenance.

When configuring LDAP for group mappings in Hadoop, you must create the users and groups for your Hadoop services in LDAP. When using the default shell-based group mapping provider (`org.apache.hadoop.security.ShellBasedUnixGroupsMapping`), the requisite user and group relationships already exist because they are created during the installation procedure. When you switch to LDAP as the group mapping provider, you must re-create these relationships within LDAP.

Note that if you have modified the **System User** or **System Group** setting within Cludera Manager for any service, you must use those custom values to provision the users and groups in LDAP.

The table below lists users and their group members for CDH services:



Note: Cludera Manager 5.3 introduces a new *single user mode*. In single user mode, the Cludera Manager Agent and *all the processes run by services managed by Cludera Manager* are started as a single configured user and group. See [Single User Mode Requirements](#) for more information.

Table 28: Users and Groups

Component (Version)	Unix User ID	Groups	Notes
Cludera Manager (all versions)	cludera-scm	cludera-scm	<p>Cludera Manager processes such as the Cludera Manager Server and the monitoring roles run as this user.</p> <p>The Cludera Manager keytab file must be named <code>cmf.keytab</code> since that name is hard-coded in Cludera Manager.</p> <div style="border: 1px solid black; padding: 5px; margin-top: 10px;"> <p> Note: Applicable to clusters managed by Cludera Manager only.</p> </div>

Component (Version)	Unix User ID	Groups	Notes
Apache Accumulo (Accumulo 1.4.3 and higher)	accumulo	accumulo	Accumulo processes run as this user.
Apache Avro			No special users.
Apache Flume (CDH 4, CDH 5)	flume	flume	The sink that writes to HDFS as this user must have write privileges.
Apache HBase (CDH 4, CDH 5)	hbase	hbase	The Master and the RegionServer processes run as this user.
HDFS (CDH 4, CDH 5)	hdfs	hdfs, hadoop	The NameNode and DataNodes run as this user, and the HDFS root directory as well as the directories used for edit logs should be owned by it.
Apache Hive (CDH 4, CDH 5)	hive	hive	The HiveServer2 process and the Hive Metastore processes run as this user. A user must be defined for Hive access to its Metastore DB (for example, MySQL or Postgres) but it can be any identifier and does not correspond to a Unix uid. This is <code>javax.jdo.option.ConnectionUserName</code> in <code>hive-site.xml</code> .
Apache HCatalog (CDH 4.2 and higher, CDH 5)	hive	hive	The WebHCat service (for REST access to Hive functionality) runs as the <code>hive</code> user.
HttpFS (CDH 4, CDH 5)	httpfs	httpfs	The HttpFS service runs as this user. See HttpFS Security Configuration for instructions on how to generate the merged <code>httpfs-http.keytab</code> file.
Hue (CDH 4, CDH 5)	hue	hue	Hue services run as this user.
Hue Load Balancer (Cloudera Manager 5.5 and higher)	apache	apache	The Hue Load balancer has a dependency on the <code>apache2</code> package that uses the <code>apache</code> user name. Cloudera Manager does not run processes using this user ID.
Impala	impala	impala, hive	Impala services run as this user.
Apache Kafka (Cloudera Distribution of Kafka 1.2.0)	kafka	kafka	Kafka services run as this user.
Java KeyStore KMS (CDH 5.2.1 and higher)	kms	kms	The Java KeyStore KMS service runs as this user.
Key Trustee KMS (CDH 5.3 and higher)	kms	kms	The Key Trustee KMS service runs as this user.

Component (Version)	Unix User ID	Groups	Notes
Key Trustee Server (CDH 5.4 and higher)	keytrustee	keytrustee	The Key Trustee Server service runs as this user.
Kudu	kudu	kudu	Kudu services run as this user.
Llama (CDH 5)	llama	llama	Llama runs as this user.
Apache Mahout			No special users.
MapReduce (CDH 4, CDH 5)	mapred	mapred, hadoop	Without Kerberos, the JobTracker and tasks run as this user. The LinuxTaskController binary is owned by this user for Kerberos.
Apache Oozie (CDH 4, CDH 5)	oozie	oozie	The Oozie service runs as this user.
Parquet			No special users.
Apache Pig			No special users.
Cloudera Search (CDH 4.3 and higher, CDH 5)	solr	solr	The Solr processes run as this user.
Apache Spark (CDH 5)	spark	spark	The Spark History Server process runs as this user.
Apache Sentry (incubating) (CDH 5.1 and higher)	sentry	sentry	The Sentry service runs as this user.
Apache Sqoop (CDH 4, CDH 5)	sqoop	sqoop	This user is only for the Sqoop1 Metastore, a configuration option that is not recommended.
Apache Sqoop2 (CDH 4.2 and higher, CDH 5)	sqoop2	sqoop, sqoop2	The Sqoop2 service runs as this user.
Apache Whirr			No special users.
YARN (CDH 4, CDH 5)	yarn	yarn, hadoop	Without Kerberos, all YARN services and applications run as this user. The LinuxContainerExecutor binary is owned by this user for Kerberos.
Apache ZooKeeper (CDH 4, CDH 5)	zookeeper	zookeeper	The ZooKeeper processes run as this user. It is not configurable.



Important:

- You can use either Cloudera Manager or the following command-line instructions to complete this configuration.
- This information applies specifically to CDH 5.6.x. If you use an earlier version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

Using Cloudera Manager

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)

Make the following changes to the HDFS service's security configuration:

1. Open the Cloudera Manager Admin Console and go to the **HDFS** service.
2. Click the **Configuration** tab.
3. Select **Scope > HDFS (Service Wide)**
4. Select **Category > Security**.
5. Modify the following configuration properties using values from the table below:

Configuration Property	Value
Hadoop User Group Mapping Implementation	org.apache.hadoop.security.LdapGroupsMapping
Hadoop User Group Mapping LDAP URL	ldap://<server>
Hadoop User Group Mapping LDAP Bind User	Administrator@example.com
Hadoop User Group Mapping LDAP Bind User Password	***
Hadoop User Group Mapping Search Base	dc=example,dc=com

Although the above changes are sufficient to configure group mappings for Active Directory, some changes to the remaining default configurations might be required for OpenLDAP.

Using the Command Line

Add the following properties to the `core-site.xml` on the NameNode:

```
<property>
<name>hadoop.security.group.mapping</name>
<value>org.apache.hadoop.security.LdapGroupsMapping</value>
</property>

<property>
<name>hadoop.security.group.mapping.ldap.url</name>
<value>ldap://server</value>
</property>

<property>
<name>hadoop.security.group.mapping.ldap.bind.user</name>
<value>Administrator@example.com</value>
</property>

<property>
<name>hadoop.security.group.mapping.ldap.bind.password</name>
<value>****</value>
</property>

<property>
<name>hadoop.security.group.mapping.ldap.base</name>
<value>dc=example,dc=com</value>
</property>

<property>
<name>hadoop.security.group.mapping.ldap.search.filter.user</name>
<value>(&amp;!(objectClass=user)(sAMAccountName={0}))</value>
</property>

<property>
<name>hadoop.security.group.mapping.ldap.search.filter.group</name>
<value>(objectClass=group)</value>
</property>

<property>
<name>hadoop.security.group.mapping.ldap.search.attr.member</name>
<value>member</value>
</property>
```

Configuring Authorization

```
<property>
<name>hadoop.security.group.mapping.ldap.search.attr.group.name</name>
<value>cn</value>
</property>
```



Note: In addition:

- If you are using Sentry with Hive, you will also need to add these properties on the HiveServer2 node.
- If you are using Sentry with Impala, add these properties on all hosts

See [Users and Groups in Sentry](#) for more information.

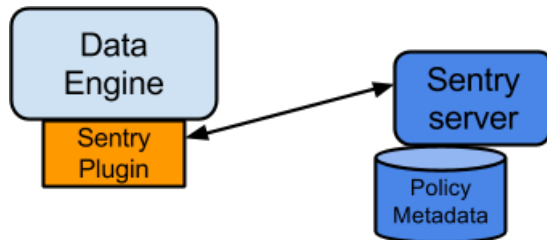
Authorization With Apache Sentry (Incubating)

Apache Sentry (incubating) is a granular, role-based authorization module for Hadoop. Sentry provides the ability to control and enforce precise levels of privileges on data for authenticated users and applications on a Hadoop cluster. Sentry currently works out of the box with Apache Hive, Hive Metastore/HCatalog, Apache Solr, Impala, and HDFS (limited to Hive table data).

Sentry is designed to be a pluggable authorization engine for Hadoop components. It allows you to define authorization rules to validate a user or application's access requests for Hadoop resources. Sentry is highly modular and can support authorization for a wide variety of data models in Hadoop.

Architecture Overview

Sentry Components



There are three components involved in the authorization process:

- **Sentry Server**

The Sentry RPC server manages the authorization metadata. It supports interfaces to securely retrieve and manipulate the metadata.

- **Data Engine**

This is a data processing application such as Hive or Impala that needs to authorize access to data or metadata resources. The data engine loads the Sentry plugin and all client requests for accessing resources are intercepted and routed to the Sentry plugin for validation.

- **Sentry Plugin**

The Sentry plugin runs in the data engine. It offers interfaces to manipulate authorization metadata stored in the Sentry server, and includes the authorization policy engine that evaluates access requests using the authorization metadata retrieved from the server.

Key Concepts

- **Authentication** - Verifying credentials to reliably identify a user

- Authorization - Limiting the user's access to a given resource
- User - Individual identified by underlying authentication system
- Group - A set of users, maintained by the authentication system
- Privilege - An instruction or rule that allows access to an object
- Role - A set of privileges; a template to combine multiple access rules
- Authorization models - Defines the objects to be subject to authorization rules and the granularity of actions allowed. For example, in the SQL model, the objects can be databases or tables, and the actions are `SELECT`, `INSERT`, `CREATE` and so on. For the Search model, the objects are indexes, collections and documents; the access modes are query, update and so on.

User Identity and Group Mapping

Sentry relies on underlying authentication systems, such as Kerberos or LDAP, to identify the user. It also uses the group mapping mechanism configured in Hadoop to ensure that Sentry sees the same group mapping as other components of the Hadoop ecosystem.

Consider a sample organization with users Alice and Bob who belong to an Active Directory (AD) group called `finance-department`. Bob also belongs to a group called `finance-managers`. In Sentry, you first create roles and then grant privileges to these roles. For example, you can create a role called Analyst and grant `SELECT` on tables Customer and Sales to this role.

The next step is to join these authentication entities (users and groups) to authorization entities (roles). This can be done by granting the Analyst role to the `finance-department` group. Now Bob and Alice who are members of the `finance-department` group get `SELECT` privilege to the Customer and Sales tables.

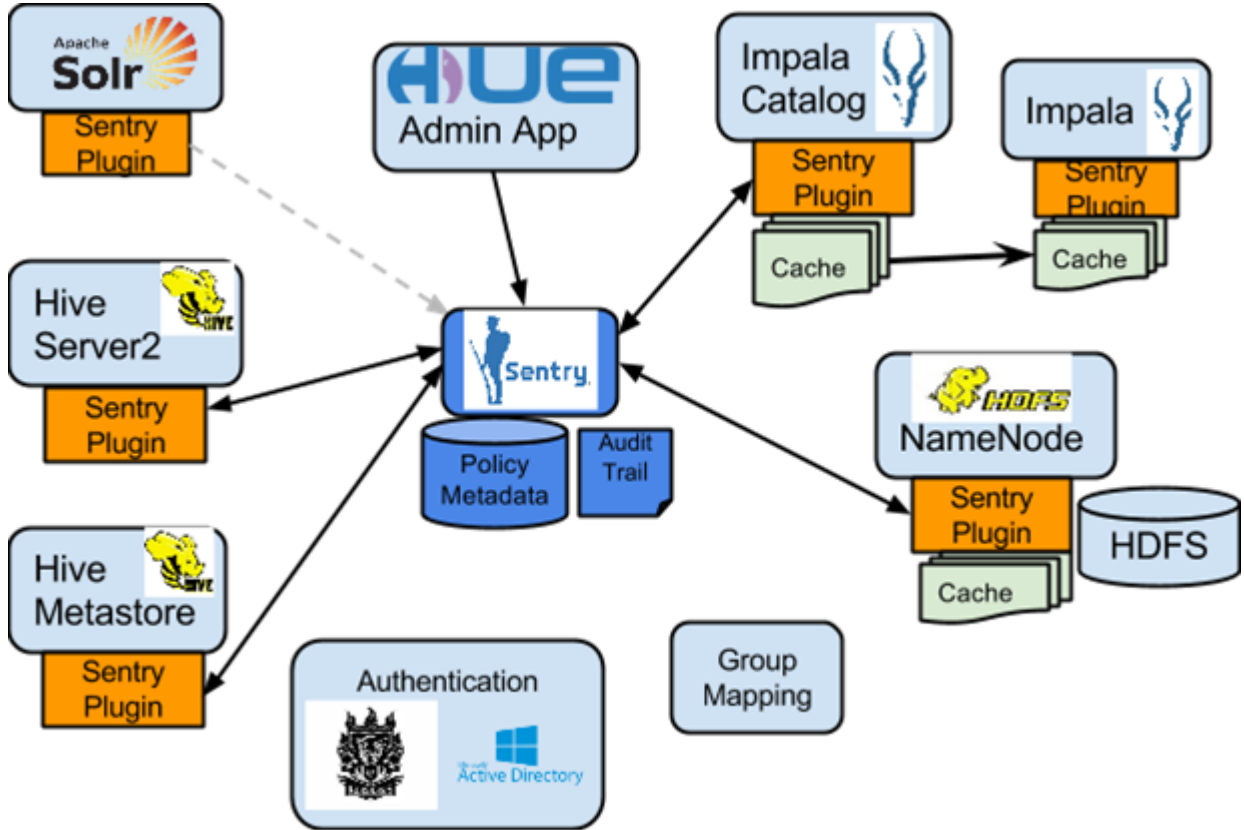
Role-Based Access Control

Role-based access control (RBAC) is a powerful mechanism to manage authorization for a large set of users and data objects in a typical enterprise. New data objects get added or removed, users join, move, or leave organisations all the time. RBAC makes managing this a lot easier. Hence, as an extension of the sample organization discussed previously, if a new employee Carol joins the Finance Department, all you need to do is add her to the `finance-department` group in AD. This will give Carol access to data from the Sales and Customer tables.

Unified Authorization

Another important aspect of Sentry is the unified authorization. The access control rules once defined, work across multiple data access tools. For example, being granted the Analyst role in the previous example will allow Bob, Alice, and others in the `finance-department` group to access table data from SQL engines such as Hive and Impala, as well as using MapReduce, Pig applications or metadata access using HCatalog.

Sentry Integration with the Hadoop Ecosystem



As illustrated above, Apache Sentry works with multiple Hadoop components. At the heart you have the Sentry Server which stores authorization metadata and provides APIs for tools to retrieve and modify this metadata securely.

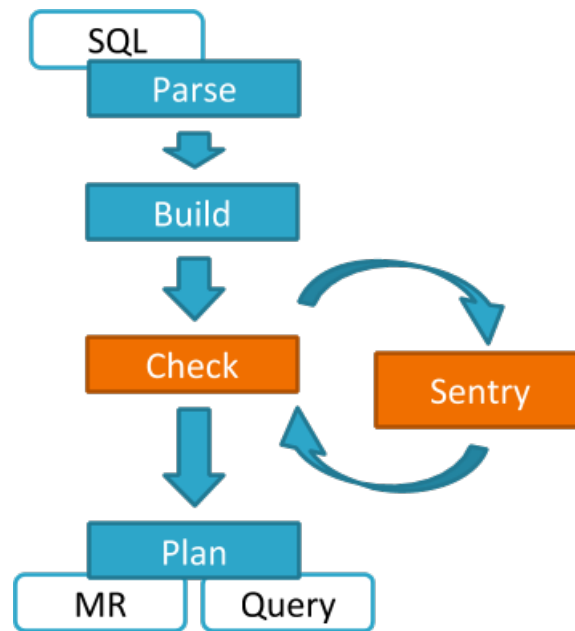
Note that the Sentry server only facilitates the metadata. The actual authorization decision is made by a policy engine which runs in data processing applications such as Hive or Impala. Each component loads the Sentry plugin which includes the service client for dealing with the Sentry service and the policy engine to validate the authorization request.

Hive and Sentry

Consider an example where Hive gets a request to access an object in a certain mode by a client. If Bob submits the following Hive query:

```
select * from production.sales
```

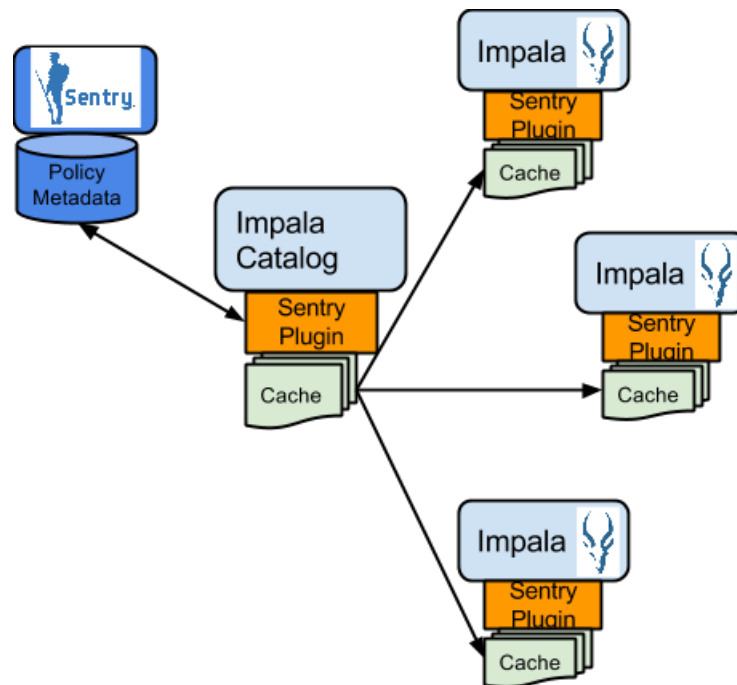
Hive will identify that user Bob is requesting `SELECT` access to the Sales table. At this point Hive will ask the Sentry plugin to validate Bob’s access request. The plugin will retrieve Bob’s privileges related to the Sales table and the policy engine will determine if the request is valid.



Hive works with both, the Sentry service and policy files. Cloudera recommends you use the Sentry service which makes it easier to manage user privileges. For more details and instructions, see [The Sentry Service](#) on page 359 or [Sentry Policy File Authorization](#) on page 392.

Impala and Sentry

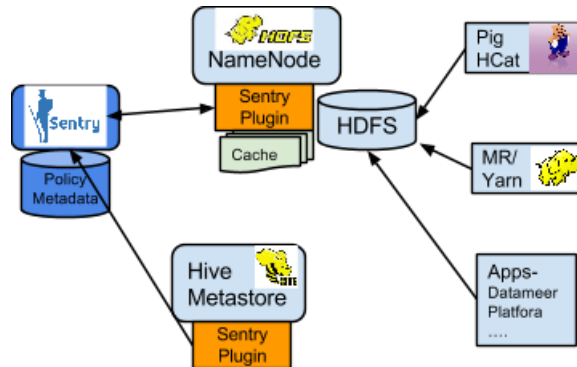
Authorization processing in Impala is similar to that in Hive. The main difference is caching of privileges. Impala's Catalog server manages caching schema metadata and propagating it to all Impala server nodes. This Catalog server caches Sentry metadata as well. As a result, authorization validation in Impala happens locally and much faster.



For detailed documentation, see [Enabling Sentry Authorization for Impala](#) on page 411.

Sentry-HDFS Synchronization

Sentry-HDFS authorization is focused on Hive warehouse data - that is, any data that is part of a table in Hive or Impala. The real objective of this integration is to expand the same authorization checks to Hive warehouse data being accessed from any other components such as Pig, MapReduce or Spark. At this point, this feature does not replace HDFS ACLs. Tables that are not associated with Sentry will retain their old ACLs.



The mapping of Sentry privileges to HDFS ACL permissions is as follows:

- SELECT privilege -> Read access on the file.
- INSERT privilege -> Write access on the file.
- ALL privilege -> Read and Write access on the file.

The NameNode loads a Sentry plugin that caches Sentry privileges as well as Hive metadata. This helps HDFS to keep file permissions and Hive table privileges in sync. The Sentry plugin periodically polls Sentry to keep the metadata changes in sync.

For example, if Bob runs a Pig job that is reading from the Sales table data files, Pig will try to get the file handle from HDFS. At that point the Sentry plugin on the NameNode will figure out that the file is part of Hive data and overlay Sentry privileges on top of the file ACLs. As a result, HDFS will enforce the same privileges for this Pig client that Hive would apply for a SQL query.

For HDFS-Sentry synchronization to work, you *must* use the Sentry service, not policy file authorization. See [Synchronizing HDFS ACLs and Sentry Permissions](#) on page 385, for more details.

Search and Sentry

Sentry can apply restrictions to various Search tasks including accessing data and creating collections. These restrictions are consistently applied, regardless of the way users attempt to complete actions. For example, restricting access to data in a collection restricts that access whether queries come from the command line, from a browser, or through the admin console.

With Search, Sentry restrictions can be stored in the database-backed Sentry service or in a policy file (for example, `sentry-provider.ini`) which is stored in an HDFS location such as

```
hdfs://ha-nn-uri/user/solr/sentry/sentry-provider.ini.
```

Sentry with Search does not support multiple policy files for multiple databases. If you choose to use policy files rather than database-backed Sentry service, you must use a separate policy file for each Sentry-enabled service. For example, if Hive and Search were using policy file authorization, using a combined Hive and Search policy file would result in an invalid configuration and failed authorization on both services.

Search works with both the Sentry service and policy files. Cloudera recommends you use the Sentry service, which makes it easier to manage user privileges. For more details and instructions, see [The Sentry Service](#) on page 359 or [Sentry Policy File Authorization](#) on page 392.

For detailed documentation, see [Enabling Sentry Authorization for Search using the Command Line](#) on page 422.

Authorization Administration

The Sentry Server supports APIs to securely manipulate roles and privileges. Both Hive and Impala support SQL statements to manage privileges natively. Sentry assumes that HiveServer2 and Impala run as superusers, usually called `hive` and `impala`. To initiate top-level permissions for Sentry, an admin must login as a superuser. You can use either Beeline or the Impala shell to execute the following sample statement:

```
GRANT ROLE Analyst TO GROUP finance_managers
```

Disabling Hive CLI

To execute Hive queries, you must use Beeline. Hive CLI is not supported with Sentry and therefore its access to the Hive Metastore must be disabled. This is especially necessary if the Hive metastore has sensitive metadata. To do this, set the **Hive Metastore Access Control and Proxy User Groups Override** property for the Hive service in Cloudera Manager. For example, to give the `hive` user permission to impersonate only members of the `hive` and `hue` groups, set the property to: `hive, hue`

If other user groups require access to the Hive Metastore, they can be added to the comma-separated list as needed. For example, setting this property to `hive, hue` blocks the Spark shell from accessing the metastore.

Using Hue to Manage Sentry Permissions

Hue supports a Security app to manage Sentry authorization. This allows users to explore and change table permissions. Here is a [video blog](#) that demonstrates its functionality.

The Sentry Service



Important: This is the documentation for the Sentry service introduced in CDH 5.1. If you want to use Sentry's previous policy file approach to secure your data, see [Sentry Policy File Authorization](#) on page 392.

The Sentry service is a RPC server that stores the authorization metadata in an underlying [relational database](#) and provides RPC interfaces to retrieve and manipulate privileges. It supports secure access to services using Kerberos. The service serves authorization metadata from the database backed storage; it does not handle actual privilege validation. The Hive and Impala services are clients of this service and will enforce Sentry privileges when configured to use Sentry.

The motivation behind introducing a new Sentry service is to make it easier to handle user privileges than the existing policy file approach. Providing a database instead, allows you to use the more traditional [GRANT/REVOKE](#) statements to modify privileges.

CDH 5.5 introduces column-level access control for tables in Hive and Impala. Previously, Sentry supported privilege granularity only down to a table. Hence, if you wanted to restrict access to a column of sensitive data, the workaround would be to first create view for a subset of columns, and then grant privileges on that view. To reduce the administrative overhead associated with such an approach, Sentry now allows you to assign the `SELECT` privilege on a subset of columns in a table.

For more information on installing, upgrading and configuring the Sentry service, see:

Prerequisites

- CDH 5.1.x (or higher) managed by Cloudera Manager 5.1.x (or higher). See the [Cloudera Manager Administration Guide](#) and [Cloudera Installation Guide](#) for instructions.
- HiveServer2 and the Hive Metastore running with strong authentication. For HiveServer2, strong authentication is either Kerberos or LDAP. For the Hive Metastore, only Kerberos is considered strong authentication (to override, see [Securing the Hive Metastore](#) on page 378).
- Impala 1.4.0 (or higher) running with strong authentication. With Impala, either Kerberos or LDAP can be configured to achieve strong authentication.

- Implement Kerberos authentication on your cluster. For instructions, see [Enabling Kerberos Authentication Using the Wizard](#) on page 59.

Terminologies

- An **object** is an entity protected by Sentry's authorization rules. The objects supported in the current release are `server`, `database`, `table`, and `URI`.
- A **role** is a collection of rules for accessing a given Hive object.
- A **privilege** is granted to a role to govern access to an object. With CDH 5.5, Sentry allows you to assign the `SELECT` privilege to columns (only for Hive and Impala). Supported privileges are:

Table 29: Valid privilege types and the objects they apply to

Privilege	Object
ALL	SERVER, TABLE, DB, URI, COLLECTION, CONFIG
INSERT	DB, TABLE
SELECT	DB, TABLE, COLUMN



Note: In Beeline, you can also grant `SELECT` and `INSERT` on `SERVER`.

- A user is an entity that is permitted by the authentication subsystem to access the Hive service. This entity can be a Kerberos principal, an LDAP `userid`, or an artifact of some other pluggable authentication system supported by HiveServer2.
- A group connects the authentication system with the authorization system. It is a collection of one or more users who have been granted one or more authorization roles. Sentry allows a set of roles to be configured for a group.
- A configured group provider determines a user's affiliation with a group. The current release supports HDFS-backed groups and locally configured groups.

Privilege Model

Sentry uses a role-based privilege model with the following characteristics.

- Allows any user to execute `show function`, `desc function`, and `show locks`.
- Allows the user to see only those tables and databases for which this user has privileges.
- Requires a user to have the necessary privileges on the URI to execute HiveQL operations that take in a location. Examples of such operations include `LOAD`, `IMPORT`, and `EXPORT`.
- Privileges granted on URIs are recursively applied to all subdirectories. That is, privileges only need to be granted on the parent directory.
- CDH 5.5 introduces column-level access control for tables in Hive and Impala. Previously, Sentry supported privilege granularity only down to a table. Hence, if you wanted to restrict access to a column of sensitive data, the workaround would be to first create view for a subset of columns, and then grant privileges on that view. To reduce the administrative overhead associated with such an approach, Sentry now allows you to assign the `SELECT` privilege on a subset of columns in a table.



Important:

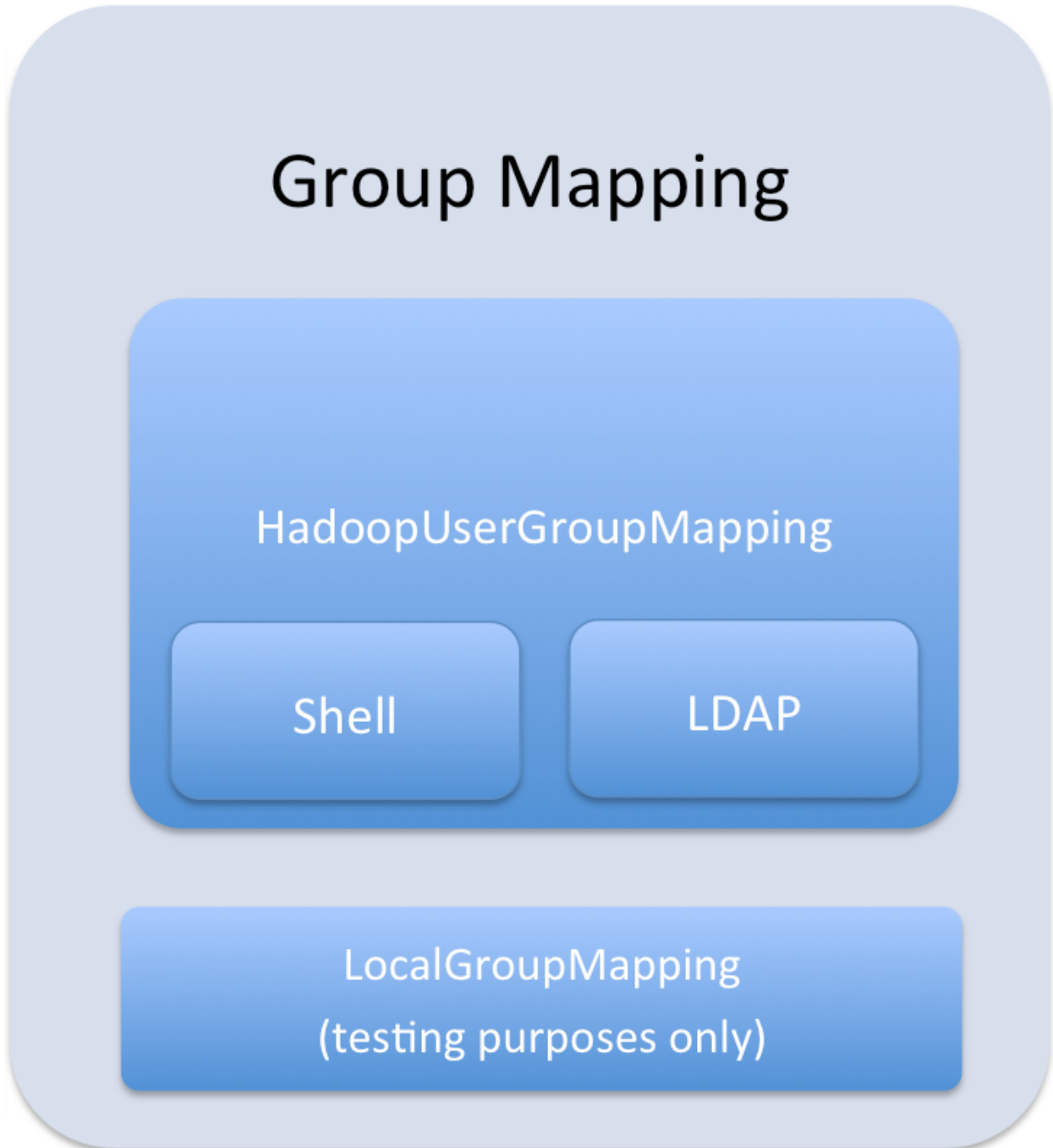
- When Sentry is enabled, you must use Beeline to execute Hive queries. Hive CLI is not supported with Sentry and must be disabled.
- When Sentry is enabled, a user with no privileges on a database will not be allowed to connect to HiveServer2. This is because the `use <database>` command is now executed as part of the connection to HiveServer2, which is why the connection fails. See [HIVE-4256](#).

For more information, see [Appendix: Authorization Privilege Model for Hive and Impala](#) on page 362.

User to Group Mapping

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator, Full Administrator**)

Group mappings in Sentry can be summarized as in the figure below.



The Sentry service only uses HadoopUserGroup mappings. You can refer [Configuring LDAP Group Mappings](#) on page 350 for details on configuring LDAP group mappings in Hadoop.

Important:

- Cloudera strongly recommends *against* using Hadoop's `LdapGroupsMapping` provider. `LdapGroupsMapping` should only be used in cases where OS-level integration is not possible. Production clusters require an identity provider that works well with all applications, not just Hadoop. Hence, often the preferred mechanism is to use tools such as SSSD, VAS or Centrify to replicate LDAP groups.
- Cloudera does not support the use of Winbind in production environments. Winbind uses an inefficient approach to user/group mapping, which may lead to low performance or cluster failures as the size of the cluster, and the number of users and groups increases.

Irrespective of the mechanism used, user/group mappings must be applied consistently across all cluster hosts for ease with maintenance.

Appendix: Authorization Privilege Model for Hive and Impala

Privileges can be granted on different objects in the Hive warehouse. Any privilege that can be granted is associated with a level in the object hierarchy. If a privilege is granted on a container object in the hierarchy, the base object automatically inherits it. For instance, if a user has `ALL` privileges on the database scope, then (s)he has `ALL` privileges on all of the base objects contained within that scope.

Important:

Note that because of this object hierarchy, it is possible for a user to read data from a database that the user does not have access to. For example, you have two roles:

- *role1* - full access to *database1* and *database2*
- *role2* - full access to *database1*, no access to *database2*

A user with *role1* can create a view in *database1* based on a table in *database2*. Because *role2* has access to *database1*, a user with *role2* can read the data in that view from *database2*.

Object Hierarchy in Hive

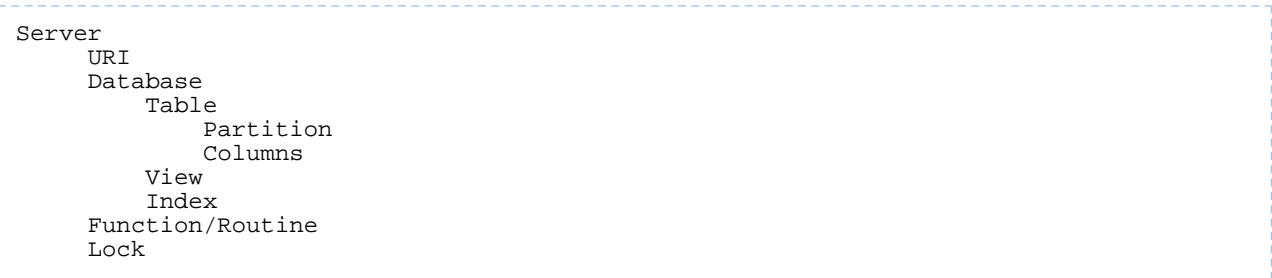


Table 30: Valid privilege types and objects they apply to

Privilege	Object
INSERT	DB, TABLE
SELECT	DB, TABLE, COLUMN
ALL	SERVER, TABLE, DB, URI

Note that when you grant `ALL` on a URI, those permissions extend into the subdirectories in that path. For example, if a role has `ALL` on the following URI:

- `hdfs://host:port/directory_A/directory_B`

That role will also have ALL on these directories:

- `hdfs://host:port/directory_A/directory_B/directory_C`
- `hdfs://host:port/directory_A/directory_B/directory_C/directory_D`
- `hdfs://host:port/directory_A/directory_B/directory_E`

URI permissions do not affect HDFS ACL's.

Table 31: Privilege hierarchy

Base Object	Granular privileges on object	Container object that contains the base object	Privileges on container object that implies privileges on the base object
DATABASE	ALL	SERVER	ALL
TABLE	INSERT	DATABASE	ALL
TABLE	SELECT	DATABASE	ALL
COLUMN	SELECT	DATABASE	ALL
VIEW	SELECT	DATABASE	ALL

Table 32: Privilege table for Hive & Impala operations

Operation	Scope	Privileges Required	URI
CREATE DATABASE	SERVER	ALL	
DROP DATABASE	DATABASE	ALL	
CREATE TABLE	DATABASE	ALL	
DROP TABLE	TABLE	ALL	
CREATE VIEW -This operation is allowed if you have column-level <code>SELECT</code> access to the columns being used.	DATABASE; SELECT on TABLE;	ALL	
ALTER VIEW -This operation is allowed if you have column-level <code>SELECT</code> access to the columns being used.	VIEW/TABLE	ALL	
DROP VIEW	VIEW/TABLE	ALL	
ALTER TABLE .. ADD COLUMNS	TABLE	ALL	
ALTER TABLE .. REPLACE COLUMNS	TABLE	ALL	
ALTER TABLE .. CHANGE column	TABLE	ALL	
ALTER TABLE .. RENAME	TABLE	ALL	
ALTER TABLE .. SET TBLPROPERTIES	TABLE	ALL	
ALTER TABLE .. SET FILEFORMAT	TABLE	ALL	
ALTER TABLE .. SET LOCATION	TABLE	ALL	URI

Operation	Scope	Privileges Required	URI
ALTER TABLE .. ADD PARTITION	TABLE	ALL	
ALTER TABLE .. ADD PARTITION location	TABLE	ALL	URI
ALTER TABLE .. DROP PARTITION	TABLE	ALL	
ALTER TABLE .. PARTITION SET FILEFORMAT	TABLE	ALL	
SHOW CREATE TABLE	TABLE	SELECT	
SHOW PARTITIONS	TABLE	SELECT/INSERT	
SHOW TABLES -Output includes all the tables for which the user has table-level privileges and all the tables for which the user has some column-level privileges.	TABLE	SELECT/INSERT	
SHOW GRANT ROLE -Output includes an additional field for any column-level privileges.	TABLE	SELECT/INSERT	
DESCRIBE TABLE -Output shows <i>all</i> columns if the user has table level-privileges or <code>SELECT</code> privilege on at least one table column	TABLE	SELECT/INSERT	
LOAD DATA	TABLE	INSERT	URI
SELECT -You can grant the <code>SELECT</code> privilege on a view to give users access to specific columns of a table they do not otherwise have access to. -See Column-level Authorization on page 380 for details on allowed column-level operations.	VIEW/TABLE; COLUMN	SELECT	
INSERT OVERWRITE TABLE	TABLE	INSERT	
CREATE TABLE .. AS SELECT -This operation is allowed if you have column-level <code>SELECT</code> access to the columns being used.	DATABASE; SELECT on TABLE	ALL	
USE <dbName>	Any		
CREATE FUNCTION	SERVER	ALL	
ALTER TABLE .. SET SERDEPROPERTIES	TABLE	ALL	

Operation	Scope	Privileges Required	URI
ALTER TABLE .. PARTITION SET SERDEPROPERTIES	TABLE	ALL on DATABASE	
Hive-Only Operations			
INSERT OVERWRITE DIRECTORY	TABLE	INSERT	URI
Analyze TABLE	TABLE	SELECT + INSERT	
IMPORT TABLE	DATABASE	ALL	URI
EXPORT TABLE	TABLE	SELECT	URI
ALTER TABLE TOUCH	TABLE	ALL on DATABASE	
ALTER TABLE TOUCH PARTITION	TABLE	ALL on DATABASE	
ALTER TABLE .. CLUSTERED BY SORTED BY	TABLE	ALL on DATABASE	
ALTER TABLE .. ENABLE/DISABLE	TABLE	ALL on DATABASE	
ALTER TABLE .. PARTITION ENABLE/DISABLE	TABLE	ALL on DATABASE	
ALTER TABLE .. PARTITION.. RENAME TO PARTITION	TABLE	ALL on DATABASE	
MSCK REPAIR TABLE	TABLE	ALL	
ALTER DATABASE	DATABASE	ALL	
DESCRIBE DATABASE	DATABASE	SELECT/INSERT	
SHOW COLUMNS -Output for this operation filters columns to which the user does not have explicit <code>SELECT</code> access	TABLE	SELECT/INSERT	
CREATE INDEX	TABLE	ALL	
DROP INDEX	TABLE	ALL	
SHOW INDEXES	TABLE	SELECT/INSERT	
GRANT PRIVILEGE	Allowed only for Sentry admin users		
REVOKE PRIVILEGE	Allowed only for Sentry admin users		
SHOW GRANT	Allowed only for Sentry admin users		
SHOW TBLPROPERTIES	TABLE	SELECT/INSERT	
DESCRIBE TABLE .. PARTITION	TABLE	SELECT/INSERT	
ADD ARCHIVE[S]	Not Allowed		
ADD FILE[S]	Not Allowed		
ADD JAR[S]	Not Allowed		
DELETE JAR[S]	Not Allowed		

Operation	Scope	Privileges Required	URI
DFS	Not Allowed		
LIST JAR[S]	Not Allowed		
SHOW CREATE VIEW	VIEW	SELECT	
Impala-Only Operations			
EXPLAIN SELECT	TABLE; COLUMN	SELECT	
EXPLAIN INSERT	TABLE; COLUMN	INSERT	
INVALIDATE METADATA	SERVER	ALL	
INVALIDATE METADATA <table name>	TABLE	SELECT/INSERT	
REFRESH <table name>	TABLE	SELECT/INSERT	
DROP FUNCTION	SERVER	ALL	
COMPUTE STATS	TABLE	ALL	
SHOW CREATE VIEW	VIEW / TABLE(S)	SELECT	

[Installing and Upgrading the Sentry Service](#)

This topic describes how to install and upgrade the Sentry service. If you are migrating from Sentry policy files to the database-backed Sentry service, see [Migrating from Sentry Policy Files to the Sentry Service](#) on page 369.

[Adding the Sentry Service](#)

Use one of the following sections to add/install the Sentry service:

[Adding the Sentry Service Using Cloudera Manager](#)

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

1. On the **Home > Status** tab, click



to the right of the cluster name and select **Add a Service**. A list of service types display. You can add one type of service at a time.

2. Select the **Sentry** service and click **Continue**.
3. Customize role assignments for Sentry. The wizard evaluates the hardware configurations of the available hosts and selects the best hosts for each role. If you are happy with the preselected hosts, click **Continue** and move to the next step. Otherwise, you can change the hosts that are assigned to the roles. The **View By Host** button allows you to view a list of hosts.

To change the host for a role, click the hostname under the role. A window appears with a list of hosts to choose from. Note that you can only select one host for the Sentry Server. You can search for a host in the Search field or you can filter the list by entering a range of hosts in the Search field. You can search for hosts in the following ways:

- Range of hostnames (without the domain portion)

Range Definition	Matching Hosts
10.1.1.[1-4]	10.1.1.1, 10.1.1.2, 10.1.1.3, 10.1.1.4
host[1-3].company.com	host1.company.com, host2.company.com, host3.company.com
host[07-10].company.com	host07.company.com, host08.company.com, host09.company.com, host10.company.com

- IP addresses
- Rack name

Click **Search** to filter the list and click a hostname to select the host. Click **OK** to close the window. The hostname that you selected appears under the role.

Click **Continue** to move to the next page in the wizard.

4. Configure database settings. You can use either an embedded or a custom database.

a. Choose the database type:

- Leave the default setting of **Use Embedded Database** to have Cloudera Manager create and configure required databases. Make a note of the auto-generated passwords.
- Select **Use Custom Databases** to specify external databases.
 1. Enter the database host, database type, database name, username, and password for the database that you created when you set up the database. See the [Creating Databases](#) documentation for Sentry Server database requirements.

b. Click **Test Connection** to confirm that Cloudera Manager can communicate with the database using the information you have supplied. If the test succeeds in all cases, click **Continue**; otherwise check and correct the information you have provided for the database and then try the test again. (For some servers, if you are using the embedded database, you will see a message saying the database will be created at a later step in the installation process.) The Review Changes page displays.

5. Click **Continue** then click **Finish**. You are returned to the [Home](#) page.

6. Verify the new service is started properly by checking the health status for the new service. If the Health Status is **Good**, then the service started properly.

7. To use the Sentry service, begin by enabling [Hive](#) and [Impala](#) for the service.

Installing Sentry Using the Command Line

Use the following the instructions, depending on your operating system, to install the latest version of Sentry.



Important: Configuration files

- If you install a newer version of a package that is already on the system, configuration files that you have modified will remain intact.
- If you uninstall a package, the package manager renames any configuration files you have modified from `<file>` to `<file>.rpmsave`. If you then re-install the package (probably to install a new version) the package manager creates a new `<file>` with applicable defaults. You are responsible for applying any changes captured in the original configuration file to the new configuration file. In the case of Ubuntu and Debian upgrades, you will be prompted if you have made changes to a file for which there is a new version; for details, see [Automatic handling of configuration files by dpkg](#).

OS	Command
RHEL	<code>\$ sudo yum install sentry</code>
SLES	<code>\$ sudo zypper install sentry</code>
Ubuntu or Debian	<code>\$ sudo apt-get update;</code> <code>\$ sudo apt-get install sentry</code>

Starting the Sentry Service

Perform the following steps to start the Sentry service on your cluster.

1. Set the `SENTRY_HOME` and `HADOOP_HOME` parameters.

Configuring Authorization

2. Create the Sentry database schema using the Sentry schematool. Sentry, by default, does not initialize the schema. The schematool is a built-in way for you to deploy the backend schema required by the Sentry service. For example, the following command uses the schematool to initialize the schema for a MySQL database.

```
bin/sentry --command schema-tool --conffile <sentry-site.xml> --dbType mysql --initSchema
```

Alternatively, you can set the `sentry.verify.schema.version` configuration property to `false`. However, this is not recommended.

3. Start the Sentry service.

```
bin/sentry --command service --conffile <sentry-site.xml>
```

Upgrading the Sentry Service

Use one of the following sections to upgrade the Sentry service:

Upgrading the Sentry Service Using Cloudera Manager

If you have a cluster managed by Cloudera Manager, go to [Upgrading CDH and Managed Services Using Cloudera Manager](#) and follow the instructions depending on the version of CDH you are upgrading to. If you are upgrading from CDH 5.1, you will notice an extra step in the procedure to upgrade the Sentry database schema.

For command-line instructions, continue reading.

Upgrading the Sentry Service Using the Command Line

1. Stop the Sentry service by identifying the PID of the Sentry Service and use the `kill` command to end the process:

```
ps -ef | grep sentry  
kill -9 <PID>
```

Replace `<PID>` with the PID of the Sentry Service.

2. Remove the previous version of Sentry.

OS	Command
RHEL	<code>\$ sudo yum remove sentry</code>
SLES	<code>\$ sudo zypper remove sentry</code>
Ubuntu or Debian	<code>\$ sudo apt-get remove sentry</code>

3. Install the new version of Sentry.

OS	Command
RHEL	<code>\$ sudo yum install sentry</code>
SLES	<code>\$ sudo zypper install sentry</code>
Ubuntu or Debian	<code>\$ sudo apt-get update;</code> <code>\$ sudo apt-get install sentry</code>

4. (From CDH 5.1 to CDH 5.x) Upgrade Sentry Database Schema

Use the Sentry `schematool` to upgrade the database schema as follows:

```
bin/sentry --command schema-tool --conffile <sentry-site.xml> --dbType <db-type>  
--upgradeSchema
```

Where `<db-type>` should be either `mysql`, `postgres` or `oracle`.

5. Start the Sentry Service

- a. Set the `SENTRY_HOME` and `HADOOP_HOME` parameters.
- b. Run the following command:

```
bin/sentry --command service --conffile <sentry-site.xml>
```

Migrating from Sentry Policy Files to the Sentry Service

Minimum Required Role: [Cluster Administrator](#) (also provided by **Full Administrator**)

The following steps describe how you can upgrade from Sentry's policy file-based approach to the new database-backed Sentry service.

1. If you haven't already done so, upgrade your cluster to the latest version of CDH and Cloudera Manager. Refer the [Cloudera Manager Administration Guide](#) for instructions.
2. Disable the existing Sentry policy file for any Hive or Impala services on the cluster. To do this:
 - a. Go to the Hive or Impala service.
 - b. Click the **Configuration** tab.
 - c. Select **Scope** > **Service Name (Service-Wide)**.
 - d. Select **Category** > **Policy File Based Sentry**.
 - e. Deselect **Enable Sentry Authorization using Policy Files**. Cloudera Manager will throw a validation error if you attempt to configure the Sentry service while this property is checked.
 - f. Repeat for any remaining Hive or Impala services.
3. Add the new Sentry service to your cluster. For instructions, see [Adding the Sentry Service](#) on page 366.
4. To begin using the Sentry service, see [Enabling the Sentry Service Using Cloudera Manager](#) on page 369 and [Configuring Impala as a Client for the Sentry Service](#) on page 377.
5. Use the command-line interface Beeline to issue grants to the Sentry service to match the contents of your old policy file(s). For more details on the Sentry service and examples on using Grant/Revoke statements to match your policy file, see [Hive SQL Syntax for Use with Sentry](#) on page 380.

Configuring the Sentry Service

This topic describes how to enable the Sentry service for Hive and Impala, and configuring the Hive metastore to communicate with the service.

Enabling the Sentry Service Using Cloudera Manager

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator, Full Administrator**)

Before Enabling the Sentry Service

- Ensure you satisfy all the [Prerequisites](#) on page 359 for the Sentry service.



Important: If you are going to enable [HDFS/Sentry synchronization](#), you do not need to perform the following step to explicitly set permissions for the Hive warehouse directory. With synchronization enabled, all Hive databases and tables will automatically be owned by `hive:hive`, and Sentry permissions on tables are translated to HDFS ACLs for the underlying table files.

The Hive warehouse directory (`/user/hive/warehouse` or any path you specify as `hive.metastore.warehouse.dir` in your `hive-site.xml`) must be owned by the Hive user and group.

- **Using the default Hive warehouse:** Permissions on the warehouse directory must be set as follows (see following Note for caveats):
 - **771** on the directory itself (for example, `/user/hive/warehouse`)
 - **771** on all subdirectories (for example, `/user/hive/warehouse/mysubdir`)

- All files and subdirectories should be owned by `hive:hive`

For example:

```
$ sudo -u hdfs hdfs dfs -chmod -R 771 /user/hive/warehouse
$ sudo -u hdfs hdfs dfs -chown -R hive:hive /user/hive/warehouse
```

If you have enabled Kerberos on your cluster, you must `kinit` as the `hdfs` user before you set permissions.

For example:

```
sudo -u hdfs kinit -kt <hdfs.keytab> hdfs
sudo -u hdfs hdfs dfs -chmod -R 771 /user/hive/warehouse
$ sudo -u hdfs hdfs dfs -chown -R hive:hive /user/hive/warehouse
```



Note:

- If you set `hive.warehouse.subdir.inherit.perms` to `true` in `hive-site.xml`, the permissions on the subdirectories will be set when you set permissions on the warehouse directory itself.
- If a user has access to any object in the warehouse, that user will be able to execute `use default`. This ensures that `use default` commands issued by legacy applications work when Sentry is enabled.
- The instructions described above for modifying permissions on the Hive warehouse directory override the recommendations in the Hive section of the CDH 5 Installation Guide.

- **Using a non-default Hive warehouse:** If you would like to use a different directory as the Hive warehouse, update the `hive.metastore.warehouse.dir` property, and make sure you set the required permissions on the new directory. For example, if the new warehouse directory is `/data`, set the permissions as follows:

```
$ hdfs dfs -chown hive:hive /data
$ hdfs dfs -chmod 771 /data
```

Note that when you update the default Hive warehouse, previously created tables will not be moved over automatically. Therefore, tables created before the update will remain at `/user/hive/warehouse/<old_table>`. However, after the update, any new tables created in the default location will be found at `/data/<new_table>`.

For Sentry/HDFS sync to work as expected, add the new warehouse URL to the list of [Sentry Synchronization Path Prefixes](#).

- Disable impersonation for HiveServer2 in the Cloudera Manager Admin Console:
 1. Go to the Hive service.
 2. Click the **Configuration** tab.
 3. Select **Scope > HiveServer2**.
 4. Select **Category > Main**.
 5. Uncheck the **HiveServer2 Enable Impersonation** checkbox.
 6. Click **Save Changes** to commit the changes.
- If you are using MapReduce, enable the Hive user to submit MapReduce jobs.
 1. Open the Cloudera Manager Admin Console and go to the MapReduce service.
 2. Click the **Configuration** tab.
 3. Select **Scope > TaskTracker**.
 4. Select **Category > Security**.
 5. Set the **Minimum User ID for Job Submission** property to zero (the default is 1000).

6. Click **Save Changes** to commit the changes.
 7. Repeat steps 1-6 for *every* TaskTracker role group for the MapReduce service that is associated with Hive, if more than one exists.
 8. Restart the MapReduce service.
- If you are using YARN, enable the Hive user to submit YARN jobs.
 1. Open the Cloudera Manager Admin Console and go to the YARN service.
 2. Click the **Configuration** tab.
 3. Select **Scope > NodeManager**.
 4. Select **Category > Security**.
 5. Ensure the **Allowed System Users** property includes the `hive` user. If not, add `hive`.
 6. Click **Save Changes** to commit the changes.
 7. Repeat steps 1-6 for *every* NodeManager role group for the YARN service that is associated with Hive, if more than one exists.
 8. Restart the YARN service.
 - Block the external applications from accessing the Hive metastore:
 1. In the Cloudera Manager Admin Console, select the **Hive** service.
 2. On the Hive service page, click the **Configuration** tab.
 3. In the search well on the right half of the Configuration page, search for `Hive Metastore Access Control` and `Proxy User Groups Override` to locate the `hadoop.proxyuser.hive.groups` parameter and click the plus sign.
 4. Enter `hive` into the text box and click the plus sign again.
 5. Enter `hue` into the text box.
 6. Click **Save Changes**.

Setting this parameter blocks access to the Hive metastore for non-service users. This effectively disables Hive CLI, Spark, and Sqoop applications from interacting with the Hive service. These application will still run, but after setting this parameter as described here, they will no longer be able to access the Hive metastore and all Hive queries will fail. Users running these tools must be part of the `hive` or `hue` groups to access the Hive service. To allow greater access, additional user groups must be added to the proxy list.



Important: Ensure you have unchecked the **Enable Sentry Authorization using Policy Files** configuration property for *both* Hive and Impala under the **Policy File Based Sentry** category before you proceed.

Enabling the Sentry Service for Hive

1. Go to the Hive service.
2. Click the **Configuration** tab.
3. Select **Scope > Hive (Service-Wide)**.
4. Select **Category > Main**.
5. Locate the **Sentry Service** property and select `Sentry`.
6. Click **Save Changes** to commit the changes.
7. Restart the Hive service.

Enabling Sentry on Hive service places several HiveServer2 properties on a restricted list properties that cannot be modified at runtime by clients. See [HiveServer2 Restricted Properties](#) on page 377.

Enabling the Sentry Service for Impala

1. Enable the Sentry service for Hive (as instructed above).
2. Go to the Impala service.

Configuring Authorization

3. Click the **Configuration** tab.
4. Select **Scope > Impala (Service-Wide)**.
5. Select **Category > Main**.
6. Locate the **Sentry Service** property and select `Sentry`.
7. Click **Save Changes** to commit the changes.
8. Restart Impala.

Enabling the Sentry Service for Hue

Hue uses a Security app to make it easier to interact with Sentry. When you set up Hue to manage Sentry permissions, make sure that users and groups are set up correctly. Every Hue user connecting to Sentry must have an equivalent OS-level user account on all hosts so that Sentry can authenticate Hue users. Each OS-level user should also be part of an OS-level group with the same name as the corresponding user's group in Hue.

For more information on using the Security app, see the related [blog post](#).

Enable the Sentry service as follows:

1. Enable the Sentry service for Hive and Impala (as instructed above).
2. Go to the Hue service.
3. Click the **Configuration** tab.
4. Select **Scope > Hue (Service-Wide)**.
5. Select **Category > Main**.
6. Locate the **Sentry Service** property and select `Sentry`.
7. Click **Save Changes** to commit the changes.
8. Restart Hue.

Add the Hive and Hue Groups to Sentry's Admin Groups

1. Go to the Sentry service.
2. Click the **Configuration** tab.
3. Select **Scope > Sentry (Service-Wide)**.
4. Select **Category > Main**.
5. Locate the **Admin Groups** property and add the `hive` and `hue` groups to the list. If an end user is in one of these admin groups, that user has administrative privileges on the Sentry Server.
6. Click **Save Changes** to commit the changes.

Enabling the Sentry Service Using the Command Line



Important:

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

Before Enabling the Sentry Service

-



Important: If you are going to enable [HDFS/Sentry synchronization](#), you do not need to perform the following step to explicitly set permissions for the Hive warehouse directory. With synchronization enabled, all Hive databases and tables will automatically be owned by `hive:hive`, and Sentry permissions on tables are translated to HDFS ACLs for the underlying table files.

The Hive warehouse directory (`/user/hive/warehouse` or any path you specify as `hive.metastore.warehouse.dir` in your `hive-site.xml`) must be owned by the Hive user and group.

- **Using the default Hive warehouse:** Permissions on the warehouse directory must be set as follows (see following Note for caveats):
 - **771** on the directory itself (for example, `/user/hive/warehouse`)
 - **771** on all subdirectories (for example, `/user/hive/warehouse/mysubdir`)
 - All files and subdirectories should be owned by `hive:hive`

For example:

```
$ sudo -u hdfs hdfs dfs -chmod -R 771 /user/hive/warehouse
$ sudo -u hdfs hdfs dfs -chown -R hive:hive /user/hive/warehouse
```

If you have enabled Kerberos on your cluster, you must `kinit` as the `hdfs` user before you set permissions.

For example:

```
sudo -u hdfs kinit -kt <hdfs.keytab> hdfs
sudo -u hdfs hdfs dfs -chmod -R 771 /user/hive/warehouse
$ sudo -u hdfs hdfs dfs -chown -R hive:hive /user/hive/warehouse
```



Note:

- If you set `hive.warehouse.subdir.inherit.perms` to `true` in `hive-site.xml`, the permissions on the subdirectories will be set when you set permissions on the warehouse directory itself.
- If a user has access to any object in the warehouse, that user will be able to execute `use default`. This ensures that `use default` commands issued by legacy applications work when Sentry is enabled.
- The instructions described above for modifying permissions on the Hive warehouse directory override the recommendations in the Hive section of the CDH 5 Installation Guide.

- **Using a non-default Hive warehouse:** If you would like to use a different directory as the Hive warehouse, update the `hive.metastore.warehouse.dir` property, and make sure you set the required permissions on the new directory. For example, if the new warehouse directory is `/data`, set the permissions as follows:

```
$ hdfs dfs -chown hive:hive /data
$ hdfs dfs -chmod 771 /data
```

Note that when you update the default Hive warehouse, previously created tables will not be moved over automatically. Therefore, tables created before the update will remain at `/user/hive/warehouse/<old_table>`. However, after the update, any new tables created in the default location will be found at `/data/<new_table>`.

For Sentry/HDFS sync to work as expected, add the new warehouse URL to the list of [Sentry Synchronization Path Prefixes](#).

- HiveServer2 impersonation must be turned off.
- If you are using MapReduce, you must enable the Hive user to submit MapReduce jobs. You can ensure that this is true by setting the minimum user ID for job submission to 0. Edit the `taskcontroller.cfg` file and set `min.user.id=0`.

If you are using YARN, you must enable the Hive user to submit YARN jobs, add the user `hive` to the `allowed.system.users` configuration property. Edit the `container-executor.cfg` file and add `hive` to the `allowed.system.users` property. For example,

```
allowed.system.users = nobody,impala,hive
```



Important: You must restart the cluster and HiveServer2 after changing these values.

- Block the Hive CLI user from accessing the Hive metastore by setting the following property in the cluster's `core-site.xml` file:

```
<property>
  <name>hadoop.proxyuser.hive.groups</name>
  <value>hive,hue</value>
  <description>Sets groups from which the hive user can impersonate other
  users.</description>
</property>
```

Setting this parameter blocks access to the Hive metastore for the user running the Hive CLI if they are not part of the `hive` or the `hue` groups. The Hive CLI can still run, but after setting this parameter as described here, the `hive` user can impersonate only members of the `hive` or the `hue` groups. If you are using Sqoop, the Sqoop user must also have access to the Hive metastore.

- Add the `hive`, `impala` and `hue` groups to Sentry's `sentry.service.admin.group` in the `sentry-site.xml` file. If an end user is in one of these admin groups, that user has administrative privileges on the Sentry Server.

```
<property>
  <name>sentry.service.admin.group</name>
  <value>hive,impala,hue</value>
</property>
```

Configuring the Sentry Server

Configure the following properties in `sentry-site.xml` on the Sentry Server host.

```
<property>
  <name>sentry.verify.schema.version</name>
  <value> </value>
  <description>
    value: true, false
    true Sentry store will verify the schema version in backed DB with expected version
    in jar.
    The service won't start if there's a mismatch
  </description>
</property>

<property>
  <name>sentry.service.server-max-threads</name>
  <value> </value>
  <description> Number of threads 500 Max worker threads to serve client
  requests</description>
</property>

<property>
  <name>sentry.service.server-min-threads</name>
  <value> </value>
  <description>Number of threads 10 Min worker threads to serve client
  requests</description>
</property>

<property>
  <name>sentry.service.allow.connect</name>
  <value> </value>
  <description>comma separated list of users - List of users that are allowed to
  connect to the service (eg Hive, Impala) </description>
</property>

<property>
  <name>sentry.store.jdbc.url</name>
```

```

    <value> </value>
    <description>JDBC connection URL for the backed DB</description>
  </property>

  <property>
    <name>sentry.store.jdbc.user</name>
    <value>sentry</value>
    <description>Userid for connecting to backend db </description>
  </property>

  <property>
    <name>sentry.store.jdbc.password</name>
    <value>Sentry</value>
    <description>Sentry password for backend JDBC user </description>
  </property>

  <property>
    <name>sentry.service.server.keytab</name>
    <value></value>
    <description>Keytab for service principal</description>
  </property>

  <property>
    <name>sentry.service.server.rpcport</name>
    <value>8038</value>
    <description> TCP port number for service</description>
  </property>

  <property>
    <name>sentry.service.server.rpcaddress</name>
    <value>0.0.0.0</value>
    <description> TCP interface for service to bind to</description>
  </property>

  <property>
    <name>sentry.store.jdbc.driver</name>
    <value>org.apache.derby.jdbc.EmbeddedDriver</value>
    <description>Backend JDBC driver - org.apache.derby.jdbc.EmbeddedDriver (only when
dbtype = derby) JDBC Driver class for the backed DB</description>
  </property>

  <property>
    <name>sentry.service.admin.group</name>
    <value> </value>
    <description>Comma separates list of groups. List of groups allowed to make policy
updates</description>
  </property>

  <property>
    <name>sentry.store.group.mapping</name>
    <value>org.apache.sentry.provider.common.HadoopGroupMappingService</value>
    <description>
Group mapping class for Sentry service. org.apache.sentry.provider.file.LocalGroupMapping
service can be used for local group mapping. </description>
  </property>

  <property>
    <name>sentry.store.group.mapping.resource</name>
    <value> </value>
    <description> Policy file for group mapping. Policy file path for local group mapping,
when sentry.store.group.mapping is set to LocalGroupMapping Service class.</description>

  </property>

  <property>
    <name>sentry.service.security.mode</name>
    <value>kerberos</value>
    <description>Options: kerberos, none. Authentication mode for Sentry service.
Currently supports Kerberos and trusted mode </description>
  </property>

  <property>
    <name>sentry.service.server.principal</name>

```

Configuring Authorization

```
<value> </value>
<description>Service Kerberos principal</description>
</property>
```

Configuring HiveServer2 for the Sentry Service

Configure the following properties in `sentry-site.xml` on the HiveServer2 host.

```
<property>
  <name>hive.sentry.server</name>
  <value>server1</value>
</property>
<property>
  <name>sentry.service.server.principal</name>
  <value>sentry/_HOST@EXAMPLE.COM</value>
</property>
<property>
  <name>sentry.service.security.mode</name>
  <value>kerberos</value>
</property>
<property>
  <name>sentry.hive.provider.backend</name>
  <value>org.apache.sentry.provider.db.SimpleDBProviderBackend</value>
</property>
<property>
  <name>sentry.service.client.server.rpc-address</name>
  <value>example.cloudera.com</value>
</property>
<property>
  <name>sentry.service.client.server.rpc-port</name>
  <value>8038</value>
</property>
<property>
  <name>hive.sentry.provider</name>
  <value>org.apache.sentry.provider.file.HadoopGroupResourceAuthorizationProvider</value>
</property>
<property>
  <name>hive.sentry.failure.hooks</name>
  <value>com.cloudera.navigator.audit.hive.HiveSentryOnFailureHook</value>
</property>
```

Add the following properties to `hive-site.xml` to allow the Hive service to communicate with the Sentry service.

```
<property>
  <name>hive.security.authorization.task.factory</name>
  <value>org.apache.sentry.binding.hive.SentryHiveAuthorizationTaskFactoryImpl</value>
</property>
<property>
  <name>hive.server2.session.hook</name>
  <value>org.apache.sentry.binding.hive.HiveAuthzBindingSessionHook</value>
</property>
<property>
  <name>hive.sentry.conf.url</name>
  <value>file:///{{PATH/TO/DIR}}/sentry-site.xml</value>
</property>
```

Enabling Sentry on Hive service places several HiveServer2 properties on a restricted list properties that cannot be modified at runtime by clients. See [HiveServer2 Restricted Properties](#) on page 377.

Configuring the Hive Metastore for the Sentry Service

Add the following properties to `hive-site.xml` to allow the Hive metastore to communicate with the Sentry service.

```
<property>
  <name>hive.metastore.filter.hook</name>
  <value>org.apache.sentry.binding.metastore.SentryMetaStoreFilterHook</value>
</property>

<property>
```



```

    <name>hive.metastore.pre.event.listeners</name>
    <value>org.apache.sentry.binding.metastore.MetastoreAuthzBinding</value>
    <description>list of comma separated listeners for metastore events.</description>
</property>
<property>
  <name>hive.metastore.event.listeners</name>
  <value>org.apache.sentry.binding.metastore.SentryMetastorePostEventListener</value>

  <description>list of comma separated listeners for metastore, post
  events.</description>
</property>

```

Configuring Impala as a Client for the Sentry Service

Set the following configuration properties in `../impala-conf/sentry-site.xml` on the Catalog Server.

```

<property>
  <name>sentry.service.client.server.rpc-port</name>
  <value>8038</value>
</property>
<property>
  <name>sentry.service.client.server.rpc-address</name>
  <value>hostname</value>
</property>
<property>
  <name>sentry.service.client.server.rpc-connection-timeout</name>
  <value>200000</value>
</property>
<property>
  <name>sentry.service.security.mode</name>
  <value>kerberos</value>
</property>

```

You must also add the following configuration properties to Impala's `/etc/default/impala` file. For more information, see [Configuring Impala Startup Options through the Command Line](#).

- On the catalogd and the impalad.

```
--sentry_config=<absolute path to sentry service configuration file>
```

- On the impalad.

```
--server_name=<server name>
```

If the `--authorization_policy_file` flag is set, Impala will use the policy file-based approach. Otherwise, the database-backed approach will be used to implement authorization.

HiveServer2 Restricted Properties

Enabling Sentry on Hive service places several HiveServer2 properties on a restricted list properties that cannot be modified at runtime by clients. This list is denoted by the `hive.conf.restricted.list` property and these properties are only configurable on the server side. The list includes:

```

hive.enable.spark.execution.engine
hive.semantic.analyzer.hook
hive.exec.pre.hooks
hive.exec.scratchdir
hive.exec.local.scratchdir
hive.metastore.uris,
javax.jdo.option.ConnectionURL
hadoop.bin.path
hive.session.id
hive.aux.jars.path
hive.stats.dbconnectionstring
hive.scratch.dir.permission
hive.security.command.whitelist
hive.security.authorization.task.factory

```

Configuring Authorization

```
hive.entity.capture.transform
hive.access.conf.url
hive.sentry.conf.url
hive.access.subject.name
hive.sentry.subject.name
hive.sentry.active.role.set
```

Configuring Pig and HCatalog for the Sentry Service

Once you have the Sentry service up and running, and Hive has been configured to use the Sentry service, there are some configuration changes you must make to your cluster to allow Pig, MapReduce (using HCatLoader, HCatStorer) and WebHCat queries to access Sentry-secured data stored in Hive.

Since the Hive warehouse directory is owned by `hive:hive`, with its permissions set to `771`, with these settings, other user requests such as commands coming through Pig jobs, WebHCat queries, and MapReduce jobs, may fail. To give these users access, perform the following configuration changes:

- Use HDFS ACLs to define permissions on a specific directory or file of HDFS. This directory/file is generally mapped to a database, table, partition, or a data file.
- Users running these jobs should have the required permissions in Sentry to add new metadata or read metadata from the Hive Metastore Server. For instructions on how to set up the required permissions, see [Hive SQL Syntax for Use with Sentry](#) on page 380. You can use HiveServer2's command line interface, Beeline to update the Sentry database with the user privileges.

Examples:

- A user who is using Pig HCatLoader will require read permissions on a specific table or partition. In such a case, you can `GRANT` read access to the user in Sentry and set the ACL to read and execute, on the file being accessed.
- A user who is using Pig HCatStorer will require ALL permissions on a specific table. In this case, you `GRANT ALL` access to the user in Sentry and set the ACL to write and execute, on the table being used.

Securing the Hive Metastore

It's important that the Hive metastore be secured. If you want to override the Kerberos prerequisite for the Hive metastore, set the `sentry.hive.testing.mode` property to `true` to allow Sentry to work with weaker authentication mechanisms. Add the following property to the HiveServer2 and Hive metastore's `sentry-site.xml`:

```
<property>
  <name>sentry.hive.testing.mode</name>
  <value>true</value>
</property>
```

Impala does not require this flag to be set.



Warning: Cloudera strongly recommends against enabling this property in production. Use Sentry's testing mode only in test environments.

You can also set the property in Cloudera Manager. Go to the Hive service and open the **Configuration** tab. Search for the **Hive Service Advanced Configuration Snippet (Safety Valve) for sentry-site.xml**. Click the plus sign (+) to add a new property with the following values:

- **Name:** `sentry.hive.testing.mode`
- **Value:** `true`

You can turn on Hive metastore security using the instructions in [Cloudera Security](#). To secure the Hive metastore; see [Hive Metastore Server Security Configuration](#) on page 139.

Using User-Defined Functions with HiveServer2

The `ADD JAR` command does *not* work with HiveServer2 and the Beeline client when Beeline runs on a different host. As an alternative to `ADD JAR`, Hive's *auxiliary paths* functionality should be used. There are some differences in the

procedures for creating permanent functions and temporary functions when Sentry is enabled. For detailed instructions, see:

- [User-Defined Functions \(UDFs\) with HiveServer2 Using Cloudera Manager](#)
- OR
- [User-Defined Functions \(UDFs\) with HiveServer2 Using the Command Line](#)

Sentry Debugging and Failure Scenarios

This topic describes how Sentry deals with conflicting policies, how to debug Sentry authorization request failures and how different CDH components respond when the Sentry service fails. Before debugging, ensure you have read through the CDH Release Notes and the list of [Known Issues for Sentry](#).

Resolving Policy Conflicts

Sentry treats all policies independently. Hence, for any operation, if Sentry can find a policy that allows it, that operation will be allowed. Consider an example with a table, `test_db.test_tbl`, whose HDFS directory is located at `hdfs://user/hive/warehouse/test_db.db/test_tbl`, and grant the following conflicting privileges to a user with the role, `test_role`. That is, you are granting `ALL` privilege to the role `test_role` on the URI, but only the `SELECT` privilege on the table itself.

```
GRANT ALL ON URI 'hdfs:///user/hive/warehouse/test_db.db/test_tbl' to role test_role;
USE test_db;
GRANT SELECT ON TABLE test_tbl to role test_role;
```

With these privileges, all users with the role `test_role` will be able to carry out the `EXPORT TABLE` operation even though they should only have `SELECT` privileges on `test_db.test_tbl`:

```
EXPORT TABLE <another-table-user-can-read> TO
'hdfs:///user/hive/warehouse/test_db.db/test_tbl'
```

Debugging Failed Sentry Authorization Requests

Sentry logs all facts that lead up to authorization decisions at the debug level. If you do not understand why Sentry is denying access, the best way to debug is to temporarily turn on debug logging:

- In Cloudera Manager, add `log4j.logger.org.apache.sentry=DEBUG` to the logging settings for your service through the corresponding **Logging Safety Valve** field for Impala or HiveServer2.
- On systems not managed by Cloudera Manager, add `log4j.logger.org.apache.sentry=DEBUG` to the `log4j.properties` file on each host in the cluster, in the appropriate configuration directory for each service.

Specifically, look for exceptions and messages such as:

```
FilePermission server..., RequestPermission server..., result [true|false]
```

which indicate each evaluation Sentry makes. The `FilePermission` is from the policy file, while `RequestPermission` is the privilege required for the query. A `RequestPermission` will iterate over all appropriate `FilePermission` settings until a match is found. If no matching privilege is found, Sentry returns `false` indicating “Access Denied”.

Sentry Service Failure Scenarios

If the Sentry service fails and you attempt to access the Hive warehouse, Hive, Impala and HDFS will behave as follows:

- **Hive:** Queries to the Hive warehouse will fail with an authentication error.
- **Impala:** The Impala Catalog server caches Sentry privileges. If Sentry goes down, Impala queries will continue to work and will be authorized against this cached copy of the metadata. However, authorization DDLs such as `CREATE ROLE` or `GRANT ROLE` will fail.
- **HDFS/Sentry Synchronized Permissions:** Affected HDFS files will continue to use a cached copy of the synchronized ACLs for a configurable period of time, after which they will fall back to the Hive System User and the Hive System Group (for example, `hive:hive`). The timeout value can be modified by adding the

Configuring Authorization

`sentry.authorization-provider.cache-stale-threshold.ms` parameter to the `hdfs-site.xml` Safety Valve in Cloudera Manager. The default timeout value is 60 seconds, but you can increase this value from several minutes to a few hours, as needed to accommodate large clusters.

Hive SQL Syntax for Use with Sentry

Sentry permissions can be configured through `GRANT` and `REVOKE` statements issued either interactively or programmatically through the HiveServer2 SQL command line interface, Beeline (documentation available [here](#)). The syntax described below is very similar to the `GRANT` and `REVOKE` commands that are available in well-established relational database systems.

In HUE, the Sentry Admin that creates roles and grants privileges must belong to a group that has ALL privileges on the server. For example, you can create a role for the group that contains the hive or impala user, and grant ALL ON SERVER .. WITH GRANT OPTION to that role:

```
CREATE ROLE <admin role>;
GRANT ALL ON SERVER <server1> TO ROLE <admin role> WITH GRANT OPTION;
GRANT ROLE <admin role> TO GROUP <hive>;
```



Important:

- When Sentry is enabled, you must use Beeline to execute Hive queries. Hive CLI is not supported with Sentry and must be disabled. See [Disabling Hive CLI](#) for information on how to disable the Hive CLI.
- There are some differences in syntax between Hive and the corresponding Impala SQL statements. For Impala syntax, see [SQL Statements](#).
- No privilege is required to drop a function. Any user can drop a function.

Sentry supports column-level authorization with the `SELECT` privilege. Information about column-level authorization is in the [Column-level Authorization](#) on page 380 section of this page.

See the sections below for details about the supported statements and privileges:

Column-level Authorization

Sentry allows you to assign the `SELECT` privilege on a subset of columns in a table.

The following command grants a role the `SELECT` privilege on a column:

```
GRANT SELECT (<column name>) ON TABLE <table name> TO ROLE <role name>;
```

The following command can be used to revoke the `SELECT` privilege on a column:

```
REVOKE SELECT (<column name>) ON TABLE <table name> FROM ROLE <role name>;
```

Any new columns added to a table will be inaccessible by default, until explicitly granted access.

Actions allowed for users with `SELECT` privilege on a column:

Users whose roles have been granted the `SELECT` privilege on columns only, can perform operations which explicitly refer to those columns. Some examples are:

```
SELECT <column name> FROM TABLE <table name>;
```

In this case, Sentry will first check to see if the user has the required privileges to access the table. It will then further check to see whether the user has the `SELECT` privilege to access the column(s).

```
SELECT COUNT <column name> FROM TABLE <table name>;
```

Users are also allowed to use the `COUNT` function to return the number of values in the column.

```
SELECT <column name> FROM TABLE <table name> WHERE <column name> <operator> GROUP BY <column name>;
```

The above command will work as long as you refer only to columns to which you already have access.

- To list the column(s) to which the current user has `SELECT` access:

```
SHOW COLUMNS (FROM|IN) <table name> [(FROM|IN) <database name>];
```

Exceptions:

- If a user has `SELECT` access to all columns in a table, the following command will work. Note that this is an exception, not the norm. In all other cases, `SELECT` on all columns does *not* allow you to perform table-level operations.

```
SELECT * FROM TABLE <table name>;
```

Limitations:

- Column-level privileges can only be applied to tables, not partitions or views.
- **HDFS-Sentry Sync:** With HDFS-Sentry sync enabled, even if a user has been granted access to all columns of a table, they will not have access to the corresponding HDFS data files. This is because Sentry does not consider `SELECT` on all columns equivalent to explicitly being granted `SELECT` on the table.
- Column-level access control for access from Spark SQL is not supported by the HDFS-Sentry plug-in.

CREATE ROLE Statement

The `CREATE ROLE` statement creates a role to which privileges can be granted. Privileges can be granted to roles, which can then be assigned to users. A user that has been assigned a role will only be able to exercise the privileges of that role.

Only users that have administrative privileges can create or drop roles. By default, the `hive`, `impala` and `hue` users have admin privileges in Sentry.

```
CREATE ROLE <role name>;
```

Note that role names are case-insensitive.

DROP ROLE Statement

The `DROP ROLE` statement can be used to remove a role from the database. Once dropped, the role will be revoked for all users to whom it was previously assigned. Queries that are already executing will not be affected. However, since Hive checks user privileges before executing each query, active user sessions in which the role has already been enabled will be affected.

```
DROP ROLE <role name>;
```

GRANT ROLE Statement

The `GRANT ROLE` statement can be used to grant roles to groups. Only Sentry admin users can grant roles to a group.

```
GRANT ROLE <role name> [, <role name>]
  TO GROUP <group name> [,GROUP <group name>]
```

Configuring Authorization

Sentry only allows you to grant roles to groups that have alphanumeric characters and underscores (`_`) in the group name. If the group name contains a non-alphanumeric character that is not an underscore, you can put the group name in backticks (```) to execute the command. For example, Sentry will return an error for the following command:

```
GRANT ROLE test TO GROUP test-group;
```

To grant a role to this group, put the group name in backticks:

```
GRANT ROLE test TO GROUP `test-group`;
```

The following command, which contains an underscore, is also acceptable:

```
GRANT ROLE test TO GROUP test_group;
```

Operating system group names must be in lowercase letters. Although group names are case-insensitive to Sentry, Sentry modifies capital letters within group names to be lowercase. For example, Sentry will change `TestGroup` to `testgroup`. It is not possible to disable this normalization. Therefore, group information within the environment must be in lowercase letters.

REVOKE ROLE Statement

The `REVOKE ROLE` statement can be used to revoke roles from groups. Only Sentry admin users can revoke the role from a group.

```
REVOKE ROLE <role name> [, <role name>]  
FROM GROUP <group name> [,GROUP <group name>]
```

GRANT <Privilege> Statement

Use the `GRANT <Privilege>` statement to grant privileges on an object to a role. The statement uses the following syntax:

```
GRANT  
  <privilege> [, <privilege> ]  
ON <object type> <object name>  
TO ROLE <role name> [,ROLE <role name>]
```

You can grant the `SELECT` privilege on specific columns of a table. For example:

```
GRANT SELECT <column name> ON TABLE <table name> TO ROLE <role name>;
```

GRANT <Privilege> ON URIs (HDFS and S3A)

If the `GRANT` for Sentry URI does not specify the complete scheme, or the URI mentioned in Hive DDL statements does not have a scheme, Sentry automatically completes the URI by applying the default scheme based on the HDFS configuration provided in the `fs.defaultFS` property. Using the same HDFS configuration, Sentry can also auto-complete URIs in case the URI is missing a scheme and an authority component.

When a user attempts to access a URI, Sentry will check to see if the user has the required privileges. During the authorization check, if the URI is incomplete, Sentry will complete the URI using the default HDFS scheme. Note that Sentry does not check URI schemes for completion when they are being used to grant privileges. This is because users can `GRANT` privileges on URIs that do not have a complete scheme or do not already exist on the filesystem.

For example, in CDH 5.8 and later, the following `CREATE EXTERNAL TABLE` statement works even though the statement does not include the URI scheme.

```
GRANT ALL ON URI 'hdfs://namenode:XXX/path/to/table' TO ROLE <role name>;  
CREATE EXTERNAL TABLE foo LOCATION 'namenode:XXX/path/to/table' TO ROLE <role name>;
```

Similarly, the following `CREATE EXTERNAL TABLE` statement works even though it is missing scheme and authority components.

```
GRANT ALL ON URI 'hdfs://namenode:XXX/path/to/table' TO ROLE <role name>;
CREATE EXTERNAL TABLE foo LOCATION
  '/path/to/table'
```

Since Sentry supports both HDFS and Amazon S3, in CDH 5.8 and later, Cloudera recommends that you specify the fully qualified URI in `GRANT` statements to avoid confusion. If the underlying storage is a mix of S3 and HDFS, the risk of granting the wrong privileges increases. The following are examples of fully qualified URIs:

- **HDFS:** `hdfs://host:port/path/to/hdfs/table`
- **S3:** `s3a://host:port/path/to/s3/table`

REVOKE <Privilege> Statement

You can use the `REVOKE <Privilege>` statement to revoke previously-granted privileges that a role has on an object.

```
REVOKE
  <privilege> [, <privilege> ]
  ON <object type> <object name>
  FROM ROLE <role name> [,ROLE <role name>]
```

For example, you can revoke previously-granted `SELECT` privileges on specific columns of a table with the following statement:

```
REVOKE SELECT <column name> ON TABLE <table name> FROM ROLE <role name>;
```

GRANT <Privilege> ... WITH GRANT OPTION

You can delegate granting and revoking privileges to other roles. For example, a role that is granted a privilege `WITH GRANT OPTION` can `GRANT/REVOKE` the same privilege to/from other roles. Hence, if a role has the `ALL` privilege on a database and the `WITH GRANT OPTION` set, users granted that role can execute `GRANT/REVOKE` statements only for that database or child tables of the database.

```
GRANT
  <privilege>
  ON <object type> <object name>
  TO ROLE <role name>
  WITH GRANT OPTION
```

Only a role with `GRANT` option on a specific privilege or its parent privilege can revoke that privilege from other roles. Once the following statement is executed, all privileges with and without grant option are revoked.

```
REVOKE
  <privilege>
  ON <object type> <object name>
  FROM ROLE <role name>
```

Hive does not currently support revoking only the `WITH GRANT OPTION` from a privilege previously granted to a role. To remove the `WITH GRANT OPTION`, revoke the privilege and grant it again without the `WITH GRANT OPTION` flag.

SET ROLE Statement

Sentry enforces restrictions on queries based on the roles and privileges that the user has. A user can have multiple roles and a role can have multiple privileges.

The `SET ROLE` command enforces restrictions at the role level, not at the user level. When you use the `SET ROLE` command to make a role active, the role becomes current for the session. If a role is not current for the session, it is inactive and the user does not have the privileges assigned to that role. A user can only use the `SET ROLE` command for roles that have been granted to the user.

To list the roles that are current for the user, use the `SHOW CURRENT ROLES` command. By default, all roles that are assigned to the user are current.

Configuring Authorization

You can use the following `SET ROLE` commands:

SET ROLE NONE

Makes all roles for the user inactive. When no role is current, the user does not have any privileges and cannot execute a query.

SET ROLE ALL

Makes all roles that have been granted to the user active. All privileges assigned to those roles are applied. When the user executes a query, the query is filtered based on those privileges.

SET ROLE *role name*

Makes a single role active. The privileges assigned to that role are applied. When the user executes a query, the query is filtered based on the privileges assigned to that role.

SHOW Statement

- To list the database(s) for which the current user has database, table, or column-level access:

```
SHOW DATABASES;
```

- To list the table(s) for which the current user has table or column-level access:

```
SHOW TABLES;
```

- To list the column(s) to which the current user has `SELECT` access:

```
SHOW COLUMNS (FROM|IN) <table name> [(FROM|IN) <database name>];
```

- To list all the roles in the system (only for sentry admin users):

```
SHOW ROLES;
```

- To list all the roles in effect for the current user session:

```
SHOW CURRENT ROLES;
```

- To list all the roles assigned to the given *group name* (only allowed for Sentry admin users and others users that are part of the group specified by *group name*):

```
SHOW ROLE GRANT GROUP group name;
```

- The `SHOW` statement can also be used to list the privileges that have been granted to a role or all the grants given to a role for a particular object.

To list all the grants for the given *<role name>* (only allowed for Sentry admin users and other users that have been granted the role specified by *<role name>*). The following command will also list any column-level privileges:

```
SHOW GRANT ROLE <role name>;
```

- To list all the grants for a role on the given *<object name>* (only allowed for Sentry admin users and other users that have been granted the role specified by *<role name>*). The following command will also list any column-level privileges:

```
SHOW GRANT ROLE <role name> on <object type> <object name>;
```


Example: Using Grant/Revoke Statements to Match an Existing Policy File



Note: In the following example(s), `server1` refers to an alias Sentry uses for the associated Hive service. It does not refer to any physical server. This alias can be modified using the `hive.sentry.server` property in `hive-site.xml`. If you are using Cloudera Manager, modify the Hive property, **Server Name for Sentry Authorization**, in the **Service-Wide > Advanced** category.

Here is a sample policy file:

```
[groups]
# Assigns each Hadoop group to its set of roles
manager = analyst_role, junior_analyst_role
analyst = analyst_role
jranalyst = junior_analyst_role
customers_admin = customers_admin_role
admin = admin_role

[roles] # The URIs below define a landing skid which
# the user can use to import or export data from the system.
# Since the server runs as the user "hive" files in that directory
# must either have the group hive and read/write set or
# be world read/write.
analyst_role = server=server1->db=analyst1, \
  server=server1->db=jranalyst1->table=*->action=select
  server=server1->uri=hdfs://ha-nn-uri/landing/analyst1
junior_analyst_role = server=server1->db=jranalyst1, \
  server=server1->uri=hdfs://ha-nn-uri/landing/jranalyst1

# Implies everything on server1.
admin_role = server=server1
```

The following sections show how you can use the new GRANT statements to assign privileges to roles (and assign roles to groups) to match the sample policy file above.

Grant privileges to analyst_role:

```
CREATE ROLE analyst_role;
GRANT ALL ON DATABASE analyst1 TO ROLE analyst_role;
GRANT SELECT ON DATABASE jranalyst1 TO ROLE analyst_role;
GRANT ALL ON URI 'hdfs://ha-nn-uri/landing/analyst1' \
  TO ROLE analyst_role;
```

Grant privileges to junior_analyst_role:

```
CREATE ROLE junior_analyst_role;
GRANT ALL ON DATABASE jranalyst1 TO ROLE junior_analyst_role;
GRANT ALL ON URI 'hdfs://ha-nn-uri/landing/jranalyst1' \
  TO ROLE junior_analyst_role;
```

Grant privileges to admin_role:

```
CREATE ROLE admin_role;
GRANT ALL ON SERVER server1 TO ROLE admin_role;
```

Grant roles to groups:

```
GRANT ROLE admin_role TO GROUP admin;
GRANT ROLE analyst_role TO GROUP analyst;
GRANT ROLE jranalyst_role TO GROUP jranalyst;
```

Synchronizing HDFS ACLs and Sentry Permissions

The HDFS-Sentry plugin allows you to configure synchronization of Sentry privileges with HDFS ACLs for specific HDFS directories.

Introduction

The integration of Sentry and HDFS permissions automatically keeps HDFS ACLs in sync with the privileges configured with Sentry. This feature offers the easiest way to share data between Hive, Impala and other components such as MapReduce, Spark, and Pig, while setting permissions for that data with just one set of rules through Sentry. It maintains the ability of Hive and Impala to set permissions on views, in addition to tables, while access to data outside of Hive and Impala (for example, reading files off HDFS) requires table permissions. HDFS permissions for some or all of the files that are part of tables defined in the Hive Metastore will now be controlled by Sentry.

This consists of three components:

- An HDFS NameNode plugin
- A Sentry-Hive Metastore plugin
- A Sentry Service plugin

With synchronization enabled, Sentry will translate permissions on databases and tables to the appropriate corresponding HDFS ACL on the underlying files in HDFS. For example, if a user group is assigned to a Sentry role that has SELECT permissions on a particular table, that user group will also have read access to the HDFS files that are part of that table. When you list those files in HDFS, this permission will be listed as an HDFS ACL. Or if a user group is assigned to a Sentry role that has SELECT permissions on a database, that user group will also have read access to the HDFS files that are part of that database. When you list those files in HDFS, those permissions will also be listed as an HDFS ACL.

Note that when Sentry was [enabled](#), the `hive` user/group was given ownership of all files/directories in the Hive warehouse (`/user/hive/warehouse`). Hence, the resulting synchronized Sentry permissions will reflect this fact. If you skipped that step, Sentry permissions will be based on the existing Hive warehouse ACLs. Sentry will not automatically grant ownership to the `hive` user.

The mapping of Sentry privileges to HDFS ACLs is as follows:

- SELECT privilege -> Read access on the file.
- INSERT privilege -> Write access on the file.
- ALL privilege -> Read and Write access on the file.

Note that you must explicitly specify the path prefix to the Hive warehouse (default: `user/hive/warehouse`) and any other directories that must be managed by Sentry. This procedure is described in the [Enabling the HDFS-Sentry Plugin Using Cloudera Manager](#) on page 388 section below.

**Important:**

- With synchronization enabled, your ability to set HDFS permissions for those files is disabled. Permissions for those particular files can be set only through Sentry, and when examined through HDFS these permissions appear as HDFS ACLs. A configurable set of users, such as `hive` and `impala`, will have full access to the files automatically. This ensures that a key requirement of using Sentry with Hive and Impala — giving these processes full access to regulate permissions on underlying data files — is met automatically.
- Tables and databases that are not associated with Sentry, that is, have no user with Sentry privileges to access them, will retain their old ACLs.
- Synchronized privileges are not persisted to HDFS. This means that when this feature is disabled, HDFS privileges will return to their original values.
- Setting HDFS ACLs on Sentry-managed paths will not affect the original HDFS ACLs. That is, if you set an ACL for a Hive object that also falls under the Sentry-managed path prefixes, no action will be taken. If the path does not point to a Hive object managed by Sentry, HDFS ACLs will be set as expected.

Removing HDFS ACLs from paths will work the same way. If you attempt to remove an ACL associated with a Hive object managed by Sentry, no action will be taken. In all other cases, the ACL will be removed as is expected behavior.

- With HDFS-Sentry sync enabled, if the NameNode plugin is unable to communicate with the Sentry Service, affected HDFS files will continue to use a cached copy of the synchronized ACLs for a configurable period of time, after which they will fall back to the Hive System User and the Hive System Group (for example, `hive:hive`). The timeout value can be modified by adding the `sentry.authorization-provider.cache-stale-threshold.ms` parameter to the `hdfs-site.xml` Safety Valve in Cloudera Manager. The default timeout value is 60 seconds, but you can increase this value from several minutes to a few hours, as needed to accommodate large clusters.
- Column-level access control for access from Spark SQL is not supported by the HDFS-Sentry plug-in.

Prompting HDFS ACL Changes

URIs do not have an impact on the HDFS-Sentry plugin. Therefore, you cannot manage all of your HDFS ACLs with the HDFS-Sentry plugin and you must continue to use standard HDFS ACLs for data outside of Hive.

HDFS ACL changes are triggered on:

- Hive DATABASE object LOCATION (HDFS) when a role is granted to the object
- Hive TABLE object LOCATION (HDFS) when a role is granted to the object

HDFS ACL changes are not triggered by:

- Hive URI LOCATION (HDFS) when a role is granted to a URI
- Hive SERVER object when a role is granted to the object. HDFS ACLs are not updated if a role is assigned to the SERVER. The privileges are inherited by child objects in standard Sentry interactions, but the plugin does not trickle the privileges down.
- Permissions granted on views. Views are not synchronized as objects in the HDFS file system.

Prerequisites

- CDH 5.3.0 or higher
- (Strongly Recommended) Implement Kerberos authentication on your cluster.

The following conditions must also be true when enabling Sentry-HDFS synchronization. Failure to comply with any of these will result in validation errors.

Configuring Authorization

- You must use the Sentry service, not policy file-based authorization.
- Enabling [HDFS Extended Access Control Lists \(ACLs\)](#) is required.
- There must be at least one Sentry service dependent on HDFS.
- The Sentry service must have at least one Sentry Server role.
- The Sentry service must have at least one dependent Hive service.
- The Hive service must have at least one Hive metastore role.

Enabling the HDFS-Sentry Plugin Using Cloudera Manager

1. Go to the HDFS service.
2. Click the **Configuration** tab.
3. Select **Scope > HDFS (Service-Wide)**.
4. Type `Check HDFS Permissions` in the Search box.
5. Select **Check HDFS Permissions**.
6. Select **Enable Sentry Synchronization**.
7. Locate the **Sentry Synchronization Path Prefixes** property or search for it by typing its name in the Search box.
8. Edit the **Sentry Synchronization Path Prefixes** property to list HDFS path prefixes where Sentry permissions should be enforced. Multiple HDFS path prefixes can be specified. By default, this property points to `/user/hive/warehouse` and must always be non-empty. If you are using a non-default location for the Hive warehouse, make sure you add it to the list of path prefixes. HDFS privilege synchronization will not occur for tables and databases located outside the HDFS regions listed here.



Important: Sentry will only manage paths that store Hive objects. If a path is listed under the **Sentry Synchronization Path Prefixes**, but there is no Hive object there, Sentry will not manage permissions for that path.

9. Click **Save Changes**.
10. Restart the cluster. Note that it may take an additional two minutes after cluster restart for privilege synchronization to take effect.

Enabling the HDFS-Sentry Plugin Using the Command Line



Important:

- You can use either Cloudera Manager or the following command-line instructions to complete this configuration.
- This information applies specifically to CDH 5.6.x. If you use an earlier version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

To enable the Sentry plugins on an unmanaged cluster, you must explicitly allow the `hdfs` user to interact with Sentry, and install the plugin packages as described in the following sections.

Allowing the `hdfs` user to connect with Sentry

For an unmanaged cluster, add `hdfs` to the `sentry.service.allow.connect` property in `sentry-site.xml`.

```
<property>
  <name>sentry.service.allow.connect</name>
  <value>impala,hive,hue,hdfs</value>
</property>
```

Installing the HDFS-Sentry Plugin**Note: Install Cloudera Repository**

Before using the instructions on this page to install the package, install the Cloudera `yum`, `zypper`/`YaST` or `apt` repository, and install or upgrade CDH 5 and make sure it is functioning correctly. For instructions, see [Installing the Latest CDH 5 Release](#).

Use the following the instructions, depending on your operating system, to install the `sentry-hdfs-plugin` package. The package must be installed (at a minimum) on the following hosts:

- The host running the NameNode and Secondary NameNode
- The host running the Hive Metastore
- The host running the Sentry Service

OS	Command
RHEL-compatible	\$ sudo yum install sentry-hdfs-plugin
SLES	\$ sudo zypper install sentry-hdfs-plugin
Ubuntu or Debian	\$ sudo apt-get install sentry-hdfs-plugin

Configuring the HDFS NameNode Plugin

Add the following properties to the `hdfs-site.xml` file on the NameNode host.

```
<property>
<name>dfs.namenode.acls.enabled</name>
<value>>true</value>
</property>

<property>
<name>dfs.namenode.authorization.provider.class</name>
<value>org.apache.sentry.hdfs.SentryAuthorizationProvider</value>
</property>

<property>
<name>dfs.permissions</name>
<value>>true</value>
</property>

<!-- Comma-separated list of HDFS path prefixes where Sentry permissions should be
enforced. -->
<!-- Privilege synchronization will occur only for tables located in HDFS regions
specified here. -->
<property>
<name>sentry.authorization-provider.hdfs-path-prefixes</name>
<value>/user/hive/warehouse</value>
</property>

<property>
<name>sentry.hdfs.service.security.mode</name>
<value>kerberos</value>
</property>

<property>
<name>sentry.hdfs.service.server.principal</name>
<value> SENTRY_SERVER_PRINCIPAL (for eg : sentry/_HOST@VPC.CLOUDERA.COM )</value>
</property>

<property>
<name>sentry.hdfs.service.client.server.rpc-port</name>
<value>SENTRY_SERVER_PORT</value>
</property>

<property>
<name>sentry.hdfs.service.client.server.rpc-address</name>
```

Configuring Authorization

```
<value>SENTRY_SERVER_HOST</value>
</property>
```

Configuring the Hive Metastore Plugin

Add the following properties to `hive-site.xml` on the Hive Metastore Server host.

```
<property>
<name>sentry.metastore.plugins</name>
<value>org.apache.sentry.hdfs.MetastorePlugin</value>
</property>

<property>
<name>sentry.hdfs.service.client.server.rpc-port</name>
<value> SENTRY_SERVER_PORT </value>
</property>

<property>
<name>sentry.hdfs.service.client.server.rpc-address</name>
<value> SENTRY_SERVER_HOSTNAME </value>
</property>

<property>
<name>sentry.hdfs.service.client.server.rpc-connection-timeout</name>
<value>200000</value>
</property>

<property>
<name>sentry.hdfs.service.security.mode</name>
<value>kerberos</value>
</property>

<property>
<name>sentry.hdfs.service.server.principal</name>
<value> SENTRY_SERVER_PRINCIPAL (for eg : sentry/_HOST@VPC.CLOUDERA.COM )</value>
</property>
```

Configuring the Sentry Service Plugin

Add the following properties to the `sentry-site.xml` file on the NameNode host.

```
<property>
<name>sentry.service.processor.factories</name>
<value>org.apache.sentry.provider.db.service.thrift.SentryPolicyStoreProcessorFactory,
org.apache.sentry.hdfs.SentryHDFSServiceProcessorFactory</value>
</property>

<property>
<name>sentry.policy.store.plugins</name>
<value>org.apache.sentry.hdfs.SentryPlugin</value>
</property>
```



Important: Once all the configuration changes are complete, restart your cluster. Note that it may take an additional two minutes after cluster restart for privilege synchronization to take effect.

Testing the Sentry Synchronization Plugins

The following tasks will help you ensure that Sentry-HDFS synchronization has been enabled and configured correctly:

For a folder that has been enabled for the plugin, such as the Hive warehouse, try accessing the files in that folder outside Hive and Impala. For this, you should know what tables and databases those HDFS files belong to and the Sentry permissions on those tables. Attempt to view or modify the Sentry permissions settings over those tables using one of the following tools:

- **(Recommended)** Hue's Security application
- HiveServer2 CLI
- Impala CLI

- Access the tables and databases directly in HDFS. For example:
 - List files inside the folder and verify that the file permissions shown in HDFS (including ACLs) match what was configured in Sentry.
 - Run a MapReduce, Pig or Spark job that accesses those files. Pick any tool besides HiveServer2 and Impala

Using the Sentry Web Server

The Sentry webserver can be used to view reported metrics which can prove useful for debugging. To enable the Sentry webserver for reporting metrics and secure it using Kerberos authentication, perform the following steps:

1. Go to the Sentry service in Cloudera Manager.
2. Click the **Configuration** tab.
3. Select **Scope > Sentry (Service-Wide)**.
4. Select **Category > Advanced**.
5. Locate the **Sentry Service Advanced Configuration Snippet (Safety Valve) for sentry-site.xml** property and add the following properties:
 - a. To enable the Sentry webserver:

```
<!-- Enable the Sentry web server -->
<property>
<name>sentry.service.web.enable</name>
<value>>true</value>
</property>
```

- b. Metrics for the Sentry service can now be reported using either JMX or console. To obtain the metrics in JSON format, you can use the Sentry Web Server which by default, listens on port 51000. To enable reporting of metrics:

```
<!-- Port on which the Sentry web server listens -->
<property>
<name>sentry.service.web.port</name>
<value>51000</value>
</property>

<!-- Tool being used to report metrics; jmx or console -->
<property>
<name>sentry.service.reporter</name>
<value>jmx</value>
</property>
```

- c. Kerberos authentication must be enabled for the Sentry web server to restrict who can access the debug webpage for the Sentry service. To enable:

```
<!-- Set Kerberos authentication properties -->
<property>
<name>sentry.service.web.authentication.type</name>
<value>KERBEROS</value>
</property>

<property>
<name>sentry.service.web.authentication.kerberos.principal</name>
<value>HTTP/<fully.qualified.domain.name>@YOUR-REALM</value>
</property>

<property>
<name>sentry.service.web.authentication.kerberos.keytab</name>
<value>/path/to/keytab/file</value>
</property>

<!-- Define comma-separated list of users allowed to connect to the web server -->
<property>
<name>sentry.service.web.authentication.allow.connect.users</name>
<value>user_a,user_b</value>
</property>
```

6. Click **Save Changes** to commit the changes.

Sentry Policy File Authorization



Important: This is the documentation for configuring Sentry using the policy file approach. Cloudera recommends you use the database-backed Sentry service introduced in CDH 5.1 to secure your data. See [The Sentry Service](#) on page 359 for more information.

Sentry enables role-based, fine-grained authorization for HiveServer2, Impala, and Cloudera Search.

The following topics provide instructions on how to install, upgrade, and configure policy file authorization.

Prerequisites

Sentry depends on an underlying authentication framework to reliably identify the requesting user. It requires:

- CDH 4.3.0 or higher.
- HiveServer2 and the Hive Metastore running with strong authentication. For HiveServer2, strong authentication is either Kerberos or LDAP. For the Hive Metastore, only Kerberos is considered strong authentication (to override, see [Securing the Hive Metastore](#) on page 411).
- Impala 1.2.1 (or higher) running with strong authentication. With Impala, either Kerberos or LDAP can be configured to achieve strong authentication. Auditing of authentication failures is supported only with CDH 4.4.0 and Impala 1.2.1 or higher.
- Implement Kerberos authentication on your cluster. This is to prevent a user bypassing the authorization and gaining direct access to the underlying data.

Terminologies

- An **object** is an entity protected by Sentry's authorization rules. The objects supported in the current release are `server`, `database`, `table`, and `URI`.
- A **role** is a collection of rules for accessing a given Hive object.
- A **privilege** is granted to a role to govern access to an object. With CDH 5.5, Sentry allows you to assign the `SELECT` privilege to columns (only for Hive and Impala). Supported privileges are:

Table 33: Valid privilege types and the objects they apply to

Privilege	Object
ALL	SERVER, TABLE, DB, URI, COLLECTION, CONFIG
INSERT	DB, TABLE
SELECT	DB, TABLE, COLUMN



Note: In Beeline, you can also grant `SELECT` and `INSERT` on `SERVER`.

- A user is an entity that is permitted by the authentication subsystem to access the Hive service. This entity can be a Kerberos principal, an LDAP `userid`, or an artifact of some other pluggable authentication system supported by HiveServer2.
- A group connects the authentication system with the authorization system. It is a collection of one or more users who have been granted one or more authorization roles. Sentry allows a set of roles to be configured for a group.
- A configured group provider determines a user's affiliation with a group. The current release supports HDFS-backed groups and locally configured groups.

Privilege Model

Sentry uses a role-based privilege model with the following characteristics.

- Allows any user to execute `show function`, `desc function`, and `show locks`.
- Allows the user to see only those tables and databases for which this user has privileges.
- Requires a user to have the necessary privileges on the URI to execute HiveQL operations that take in a location. Examples of such operations include `LOAD`, `IMPORT`, and `EXPORT`.
- Privileges granted on URIs are recursively applied to all subdirectories. That is, privileges only need to be granted on the parent directory.
- CDH 5.5 introduces column-level access control for tables in Hive and Impala. Previously, Sentry supported privilege granularity only down to a table. Hence, if you wanted to restrict access to a column of sensitive data, the workaround would be to first create view for a subset of columns, and then grant privileges on that view. To reduce the administrative overhead associated with such an approach, Sentry now allows you to assign the `SELECT` privilege on a subset of columns in a table.

**Important:**

- When Sentry is enabled, you must use Beeline to execute Hive queries. Hive CLI is not supported with Sentry and must be disabled.
- When Sentry is enabled, a user with no privileges on a database will not be allowed to connect to HiveServer2. This is because the `use <database>` command is now executed as part of the connection to HiveServer2, which is why the connection fails. See [HIVE-4256](#).

For more information, see [Authorization Privilege Model for Hive and Impala](#) on page 399.

Granting Privileges



Note: In the following example(s), `server1` refers to an alias Sentry uses for the associated Hive service. It does not refer to any physical server. This alias can be modified using the `hive.sentry.server` property in `hive-site.xml`. If you are using Cloudera Manager, modify the Hive property, **Server Name for Sentry Authorization**, found under the **Service-Wide > Advanced** category.

For example, a rule for the `Select` privilege on table `customers` from database `sales` would be formulated as follows:

```
server=server1->db=sales->table=customer->action=Select
```

To assign the `Select` privilege to the `sales_read` role on the `Id` column from the `customers` table, the rule would be as follows:

```
sales_read = server=server1->db=sales->table=customers->column=Id->action=select
```

Each object must be specified as a hierarchy of the containing objects, from server to table, followed by the privilege granted for that object. A role can contain multiple such rules, separated by commas. For example, a role might contain the `Select` privilege for the `customer` and `items` tables in the `sales` database, and the `Insert` privilege for the `sales_insights` table in the `reports` database. You would specify this as follows:

```
sales_reporting =
server=server1->db=sales->table=customer->action=Select,
server=server1->db=sales->table=items->action=Select,
server=server1->db=reports->table=sales_insights->action=Insert
```

User to Group Mapping

You can configure Sentry to use either Hadoop groups or groups defined in the policy file. By default, Sentry looks up groups locally, but it can be configured to look up Hadoop groups using LDAP (for Active Directory). User/group information for Sentry, Hive and Impala must be made available for lookup on the following hosts:

- Sentry - Groups are looked up on the host the Sentry Server runs on.

Configuring Authorization

- Hive - Groups are looked up on the hosts running HiveServer2 and the Hive Metastore.
- Impala - Groups are looked up on the Catalog Server and on all of the Impala daemon hosts.

Group mappings in Sentry can be summarized as in the figure below:



Important: You can use either Hadoop groups or local groups, but not both at the same time. Local groups are traditionally used for a quick proof-of-concept, while Hadoop groups are more commonly used in production. Refer [Configuring LDAP Group Mappings](#) on page 350 for details on configuring LDAP group mappings in Hadoop.

Policy File

The sections that follow contain notes on creating and maintaining the policy file, and using URIs to load external data and JARs.



Warning: An invalid policy file will be ignored while logging an exception. This will lead to a situation where users will lose access to all Sentry-protected data, since default Sentry behaviour is *deny* unless a user has been explicitly granted access. (Note that if only the per-DB policy file is invalid, it will invalidate only the policies in that file.)

Storing the Policy File

Considerations for storing the policy file(s) in HDFS include:

1. Replication count - Because the file is read for each query in Hive and read once every five minutes by all Impala daemons, you should increase this value; since it is a small file, setting the replication count equal to the number of client nodes in the cluster is reasonable.
2. Updating the file - Updates to the file are reflected immediately, so you should write them to a temporary copy of the file first, and then replace the existing file with the temporary one after all the updates are complete. This avoids race conditions caused by reads on an incomplete file.

Defining Roles

Keep in mind that role definitions are not cumulative; the definition that is further down in the file replaces the older one. For example, the following results in `role1` having `privilege2`, not `privilege1` and `privilege2`.

```
role1 = privilege1
role1 = privilege2
```

Role names are scoped to a specific file. For example, if you give `role1` the `ALL` privilege on `db1` in the global policy file and give `role1` `ALL` on `db2` in the per-db `db2` policy file, the user will be given both privileges.

URIs

Any command which references a URI such as `CREATE TABLE EXTERNAL`, `LOAD`, `IMPORT`, `EXPORT`, and more, in addition to `CREATE TEMPORARY FUNCTION` requires the `URI` privilege. This is an important security control because without this users could simply create an external table over an existing table they do not have access to and bypass Sentry.

URIs must start with either `hdfs://` or `file://`. If a URI starts with anything else, it will cause an exception and the policy file will be invalid.

When defining URIs for HDFS, you must also specify the NameNode. For example:

```
data_read = server=server1->uri=file:///path/to/dir,\
server=server1->uri=hdfs://namenode:port/path/to/dir
```



Important: Because the NameNode host and port must be specified, Cloudera strongly recommends you use High Availability (HA). This ensures that the URI will remain constant even if the NameNode changes.

Loading Data

Data can be loaded using a landing skid, either in HDFS or using a local/NFS directory where HiveServer2/Impala run. The following privileges can be used to grant a role access to a loading skid:

- **Load data from a local/NFS directory:**

```
server=server1->uri=file:///path/to/nfs/local/to/nfs
```

- **Load data from HDFS (MapReduce, Pig, and so on):**

```
server=server1->uri=hdfs://ha-nn-uri/data/landing-skid
```

In addition to the privilege in Sentry, the `hive` or `impala` user will require the appropriate file permissions to access the data being loaded. Groups can be used for this purpose. For example, create a group `hive-users`, and add the `hive` and `impala` users along with the users who will be loading data, to this group.

Configuring Authorization

The example `usermod` and `groupadd` commands below are only applicable to locally defined groups on the NameNode, JobTracker, and ResourceManager. If you use another system for group management, equivalent changes should be made in your group management system.

```
$ groupadd hive-users
$ usermod -G someuser,hive-users someuser
$ usermod -G hive,hive-users hive
```

External Tables

External tables require the `ALL@database` privilege in addition to the `URI` privilege. When data is being inserted through the `EXTERNAL TABLE` statement, or is referenced from an HDFS location outside the normal Hive database directories, the user needs appropriate permissions on the URIs corresponding to those HDFS locations. This means that the URI location must either be owned by the `hive:hive` user OR the `hive/impala` users must be members of the group that owns the directory.

You can configure access to the directory using a URI as follows:

```
[roles]
someuser_home_dir_role = server=server1->uri=hdfs://ha-nn-uri/user/someuser
```

You should now be able to create an external table:

```
CREATE EXTERNAL TABLE ...
LOCATION 'hdfs://ha-nn-uri/user/someuser/mytable';
```

Sample Sentry Configuration Files

This section provides a sample configuration.



Note: In the following example(s), `server1` refers to an alias Sentry uses for the associated Hive service. It does not refer to any physical server. This alias can be modified using the `hive.sentry.server` property in `hive-site.xml`. If you are using Cloudera Manager, modify the Hive property, **Server Name for Sentry Authorization**, found under the **Service-Wide > Advanced** category.

Policy Files

The following is an example of a policy file with a per-DB policy file. In this example, the first policy file, `sentry-provider.ini` would exist in HDFS; `hdfs://ha-nn-uri/etc/sentry/sentry-provider.ini` might be an appropriate location. The per-DB policy file is for the customer's database. It is located at `hdfs://ha-nn-uri/etc/sentry/customers.ini`.

sentry-provider.ini

```
[databases]
# Defines the location of the per DB policy file for the customers DB/schema
customers = hdfs://ha-nn-uri/etc/sentry/customers.ini

[groups]
# Assigns each Hadoop group to its set of roles
manager = analyst_role, junior_analyst_role
analyst = analyst_role
jranalyst = junior_analyst_role
customers_admin = customers_admin_role
admin = admin_role

[roles]
# The uris below define a define a landing skid which
# the user can use to import or export data from the system.
# Since the server runs as the user "hive" files in that directory
# must either have the group hive and read/write set or
# be world read/write.
```

```

analyst_role = server=server1->db=analyst1, \
  server=server1->db=jranalyst1->table=*->action=select
server=server1->uri=hdfs://ha-nn-uri/landing/analyst1
junior_analyst_role = server=server1->db=jranalyst1, \
  server=server1->uri=hdfs://ha-nn-uri/landing/jranalyst1

# Implies everything on server1 -> customers. Privileges for
# customers can be defined in the global policy file even though
# customers has its only policy file. Note that the Privileges from
# both the global policy file and the per-DB policy file
# are merged. There is no overriding.
customers_admin_role = server=server1->db=customers

# Implies everything on server1.
admin_role = server=server1

```

customers.ini

```

[groups]
manager = customers_insert_role, customers_select_role
analyst = customers_select_role

[roles]
customers_insert_role = server=server1->db=customers->table=*->action=insert
customers_select_role = server=server1->db=customers->table=*->action=select

```



Important: Sentry does not support using the `view` keyword in policy files. If you want to define a role against a view, use the keyword `table` instead. For example, to define the role `analyst_role` against the view `col_test_view`:

```

[roles]
analyst_role =
server=server1->db=default->table=col_test_view->action=select

```

Sentry Configuration File

The following is an example of a `sentry-site.xml` file.



Important: If you are using Cloudera Manager 4.6 (or earlier), make sure you **do not** store `sentry-site.xml` in `/etc/hive/conf`; that directory is regenerated whenever the Hive client configurations are redeployed. Instead, use a directory such as `/etc/sentry` to store the `sentry` file.

If you are using Cloudera Manager 4.7 (or higher), Cloudera Manager will create and deploy `sentry-site.xml` for you. See [The Sentry Service](#) on page 359 for more details on configuring Sentry with Cloudera Manager.

sentry-site.xml

```

<configuration>
  <property>
    <name>hive.sentry.provider</name>
    <value>org.apache.sentry.provider.file.HadoopGroupResourceAuthorizationProvider</value>
  </property>
  <property>
    <name>hive.sentry.provider.resource</name>
    <value>/path/to/authz-provider.ini</value>
  <!--
    If the hdfs-site.xml points to HDFS, the path will be in HDFS;

```

```
        alternatively you could specify a full path, e.g.:
        hdfs://namenode:port/path/to/authz-provider.ini
        file:///path/to/authz-provider.ini
    -->
</property>

<property>
  <name>sentry.hive.server</name>
  <value>server1</value>
</property>
</configuration>
```

Accessing Sentry-Secured Data Outside Hive/Impala

When Sentry is enabled, the `hive` user owns all data within the Hive warehouse. However, unlike traditional database systems the enterprise data hub allows for multiple engines to execute over the same dataset.



Note: Cloudera strongly recommends you use Hive/Impala SQL queries to access data secured by Sentry, as opposed to accessing the data files directly.

However, there are scenarios where fully vetted and reviewed jobs will also need to access the data stored in the Hive warehouse. A typical scenario would be a secured MapReduce transformation job that is executed automatically as an application user. In such cases it's important to know that the user executing this job will also have full access to the data in the Hive warehouse.

Scenario One: Authorizing Jobs

Problem

A reviewed, vetted, and automated job requires access to the Hive warehouse and cannot use Hive/Impala to access the data.

Solution

Create a group which contains `hive`, `impala`, and the user executing the automated job. For example, if the `etl` user is executing the automated job, you can create a group called `hive-users` which contains the `hive`, `impala`, and `etl` users.

The example `usermod` and `groupadd` commands below are only applicable to locally defined groups on the NameNode, JobTracker, and ResourceManager. If you use another system for group management, equivalent changes should be made in your group management system.

```
$ groupadd hive-users
$ usermod -G hive,impala,hive-users hive
$ usermod -G hive,impala,hive-users impala
$ usermod -G etl,hive-users etl
```

Once you have added users to the `hive-users` group, change directory permissions in the HDFS:

```
$ hadoop fs -chgrp -R hive:hive-users /user/hive/warehouse
$ hadoop fs -chmod -R 770 /user/hive/warehouse
```

Scenario Two: Authorizing Group Access to Databases

Problem

One group of users, `grp1` should have full access to the database, `db1`, outside of Sentry. The database, `db1` should not be accessible to any other groups, outside of Sentry. Sentry should be used for all other authorization needs.

Solution

Place the `hive` and `impala` users in `grp1`.

```
$ usermod -G hive,impala,grp1 hive
$ usermod -G hive,impala,grp1 impala
```

Then change group ownerships of all directories and files in `db1` to `grp1`, and modify directory permissions in the HDFS. This example is only applicable to local groups on a single host.

```
$ hadoop fs -chgrp -R hive:grp1 /user/hive/warehouse/db1.db
$ hadoop fs -chmod -R 770 /user/hive/warehouse/db1.db
```

Debugging Failed Sentry Authorization Requests

Sentry logs all facts that lead up to authorization decisions at the debug level. If you do not understand why Sentry is denying access, the best way to debug is to temporarily turn on debug logging:

- In Cloudera Manager, add `log4j.logger.org.apache.sentry=DEBUG` to the logging settings for your service through the corresponding **Logging Safety Valve** field for the Impala, Hive Server 2, or Solr Server services.
- On systems not managed by Cloudera Manager, add `log4j.logger.org.apache.sentry=DEBUG` to the `log4j.properties` file on each host in the cluster, in the appropriate configuration directory for each service.

Specifically, look for exceptions and messages such as:

```
FilePermission server..., RequestPermission server..., result [true|false]
```

which indicate each evaluation Sentry makes. The `FilePermission` is from the policy file, while `RequestPermission` is the privilege required for the query. A `RequestPermission` will iterate over all appropriate `FilePermission` settings until a match is found. If no matching privilege is found, Sentry returns `false` indicating “Access Denied”.

Authorization Privilege Model for Hive and Impala

Privileges can be granted on different objects in the Hive warehouse. Any privilege that can be granted is associated with a level in the object hierarchy. If a privilege is granted on a container object in the hierarchy, the base object automatically inherits it. For instance, if a user has `ALL` privileges on the database scope, then (s)he has `ALL` privileges on all of the base objects contained within that scope.

Object Hierarchy in Hive

```
Server
  URI
  Database
    Table
      Partition
      Columns
    View
    Index
  Function/Routine
  Lock
```

Table 34: Valid privilege types and objects they apply to

Privilege	Object
INSERT	DB, TABLE
SELECT	DB, TABLE, COLUMN
ALL	SERVER, TABLE, DB, URI

Note that when you grant `ALL` on a URI, those permissions extend into the subdirectories in that path. For example, if a role has `ALL` on the following URI:

- `hdfs://host:port/directory_A/directory_B`

That role will also have ALL on these directories:

- hdfs://host:port/directory_A/directory_B/directory_C
- hdfs://host:port/directory_A/directory_B/directory_C/directory_D
- hdfs://host:port/directory_A/directory_B/directory_E

URI permissions do not affect HDFS ACL's.

Table 35: Privilege hierarchy

Base Object	Granular privileges on object	Container object that contains the base object	Privileges on container object that implies privileges on the base object
DATABASE	ALL	SERVER	ALL
TABLE	INSERT	DATABASE	ALL
TABLE	SELECT	DATABASE	ALL
COLUMN	SELECT	DATABASE	ALL
VIEW	SELECT	DATABASE	ALL

Table 36: Privilege table for Hive & Impala operations

Operation	Scope	Privileges Required	URI
CREATE DATABASE	SERVER	ALL	
DROP DATABASE	DATABASE	ALL	
CREATE TABLE	DATABASE	ALL	
DROP TABLE	TABLE	ALL	
CREATE VIEW -This operation is allowed if you have column-level <code>SELECT</code> access to the columns being used.	DATABASE; SELECT on TABLE;	ALL	
ALTER VIEW -This operation is allowed if you have column-level <code>SELECT</code> access to the columns being used.	VIEW/TABLE	ALL	
DROP VIEW	VIEW/TABLE	ALL	
ALTER TABLE .. ADD COLUMNS	TABLE	ALL	
ALTER TABLE .. REPLACE COLUMNS	TABLE	ALL	
ALTER TABLE .. CHANGE column	TABLE	ALL	
ALTER TABLE .. RENAME	TABLE	ALL	
ALTER TABLE .. SET TBLPROPERTIES	TABLE	ALL	
ALTER TABLE .. SET FILEFORMAT	TABLE	ALL	
ALTER TABLE .. SET LOCATION	TABLE	ALL	URI

Operation	Scope	Privileges Required	URI
ALTER TABLE .. ADD PARTITION	TABLE	ALL	
ALTER TABLE .. ADD PARTITION location	TABLE	ALL	URI
ALTER TABLE .. DROP PARTITION	TABLE	ALL	
ALTER TABLE .. PARTITION SET FILEFORMAT	TABLE	ALL	
SHOW CREATE TABLE	TABLE	SELECT	
SHOW PARTITIONS	TABLE	SELECT/INSERT	
SHOW TABLES -Output includes all the tables for which the user has table-level privileges and all the tables for which the user has some column-level privileges.	TABLE	SELECT/INSERT	
SHOW GRANT ROLE -Output includes an additional field for any column-level privileges.	TABLE	SELECT/INSERT	
DESCRIBE TABLE -Output shows <i>all</i> columns if the user has table level-privileges or <code>SELECT</code> privilege on at least one table column	TABLE	SELECT/INSERT	
LOAD DATA	TABLE	INSERT	URI
SELECT -You can grant the <code>SELECT</code> privilege on a view to give users access to specific columns of a table they do not otherwise have access to. -See Column-level Authorization on page 380 for details on allowed column-level operations.	VIEW/TABLE; COLUMN	SELECT	
INSERT OVERWRITE TABLE	TABLE	INSERT	
CREATE TABLE .. AS SELECT -This operation is allowed if you have column-level <code>SELECT</code> access to the columns being used.	DATABASE; SELECT on TABLE	ALL	
USE <dbName>	Any		
CREATE FUNCTION	SERVER	ALL	
ALTER TABLE .. SET SERDEPROPERTIES	TABLE	ALL	

Operation	Scope	Privileges Required	URI
ALTER TABLE .. PARTITION SET SERDEPROPERTIES	TABLE	ALL on DATABASE	
Hive-Only Operations			
INSERT OVERWRITE DIRECTORY	TABLE	INSERT	URI
Analyze TABLE	TABLE	SELECT + INSERT	
IMPORT TABLE	DATABASE	ALL	URI
EXPORT TABLE	TABLE	SELECT	URI
ALTER TABLE TOUCH	TABLE	ALL on DATABASE	
ALTER TABLE TOUCH PARTITION	TABLE	ALL on DATABASE	
ALTER TABLE .. CLUSTERED BY SORTED BY	TABLE	ALL on DATABASE	
ALTER TABLE .. ENABLE/DISABLE	TABLE	ALL on DATABASE	
ALTER TABLE .. PARTITION ENABLE/DISABLE	TABLE	ALL on DATABASE	
ALTER TABLE .. PARTITION.. RENAME TO PARTITION	TABLE	ALL on DATABASE	
MSCK REPAIR TABLE	TABLE	ALL	
ALTER DATABASE	DATABASE	ALL	
DESCRIBE DATABASE	DATABASE	SELECT/INSERT	
SHOW COLUMNS -Output for this operation filters columns to which the user does not have explicit <code>SELECT</code> access	TABLE	SELECT/INSERT	
CREATE INDEX	TABLE	ALL	
DROP INDEX	TABLE	ALL	
SHOW INDEXES	TABLE	SELECT/INSERT	
GRANT PRIVILEGE	Allowed only for Sentry admin users		
REVOKE PRIVILEGE	Allowed only for Sentry admin users		
SHOW GRANT	Allowed only for Sentry admin users		
SHOW TBLPROPERTIES	TABLE	SELECT/INSERT	
DESCRIBE TABLE .. PARTITION	TABLE	SELECT/INSERT	
ADD ARCHIVE[S]	Not Allowed		
ADD FILE[S]	Not Allowed		
ADD JAR[S]	Not Allowed		
DELETE JAR[S]	Not Allowed		

Operation	Scope	Privileges Required	URI
DFS	Not Allowed		
LIST JAR[S]	Not Allowed		
SHOW CREATE VIEW	VIEW	SELECT	
Impala-Only Operations			
EXPLAIN SELECT	TABLE; COLUMN	SELECT	
EXPLAIN INSERT	TABLE; COLUMN	INSERT	
INVALIDATE METADATA	SERVER	ALL	
INVALIDATE METADATA <table name>	TABLE	SELECT/INSERT	
REFRESH <table name>	TABLE	SELECT/INSERT	
DROP FUNCTION	SERVER	ALL	
COMPUTE STATS	TABLE	ALL	
SHOW CREATE VIEW	VIEW / TABLE(S)	SELECT	

Installing and Upgrading Sentry for Policy File Authorization

Sentry stores the configuration as well as privilege policies in files. The `sentry-site.xml` file contains configuration options such as `group association provider`, `privilege policy file location`, and so on. The policy file contains the privileges and groups. It has a `.ini` file format and can be stored on a local file system or HDFS.

Sentry is plugged into Hive as session hooks which you [configure](#) in `hive-site.xml`. The `vsentry` package must be installed; it contains the required JAR files. You must also configure properties in the [Sentry Configuration File](#) on page 397.



Important:

If you have not already done so, install the Cloudera `yum`, `zypper/YaST` or `apt` repository before using the following commands. For instructions, see [Installing the Latest CDH 5 Release](#).

Installing Sentry

Use the following the instructions, depending on your operating system, to install the latest version of Sentry.



Important: Configuration files

- If you install a newer version of a package that is already on the system, configuration files that you have modified will remain intact.
- If you uninstall a package, the package manager renames any configuration files you have modified from `<file>` to `<file>.rpmsave`. If you then re-install the package (probably to install a new version) the package manager creates a new `<file>` with applicable defaults. You are responsible for applying any changes captured in the original configuration file to the new configuration file. In the case of Ubuntu and Debian upgrades, you will be prompted if you have made changes to a file for which there is a new version; for details, see [Automatic handling of configuration files by dpkg](#).

OS	Command
RHEL	<code>\$ sudo yum install sentry</code>

OS	Command
SLES	\$ sudo zypper install sentry
Ubuntu or Debian	\$ sudo apt-get update; \$ sudo apt-get install sentry

Upgrading Sentry

If you are upgrading from CDH 5.x to the latest CDH release, see [Installing Sentry](#) on page 403 to install the latest version.

Configuring Sentry Policy File Authorization Using Cloudera Manager

This topic describes how to configure Sentry policy files and enable policy file authorization for CDH services using Cloudera Manager.

Configuring User to Group Mappings

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)

Hadoop Groups

1. Go to the Hive service.
2. Click the **Configuration** tab.
3. Select **Scope > Hive (Service-Wide)**.
4. Select **Category > Policy File Based Sentry**.
5. Locate the **Sentry User to Group Mapping Class** property or search for it by typing its name in the Search box.
6. Set the **Sentry User to Group Mapping Class** property to `org.apache.sentry.provider.file.HadoopGroupResourceAuthorizationProvider`.
7. Click **Save Changes**.
8. Restart the Hive service.

Local Groups



Note: You can use either Hadoop groups or local groups, but not both at the same time. Use local groups if you want to do a quick proof-of-concept. For production, use Hadoop groups.

1. Define local groups in the `[users]` section of the [Policy File](#) on page 394. For example:

```
[users]
user1 = group1, group2, group3
user2 = group2, group3
```

2. Modify Sentry configuration as follows:
 - a. Go to the Hive service.
 - b. Click the **Configuration** tab.
 - c. Select **Scope > Hive (Service-Wide)**.
 - d. Select **Category > Policy File Based Sentry**.
 - e. Locate the **Sentry User to Group Mapping Class** property or search for it by typing its name in the Search box.
 - f. Set the **Sentry User to Group Mapping Class** property to `org.apache.sentry.provider.file.LocalGroupResourceAuthorizationProvider`.
 - g. Click **Save Changes**.
 - h. Restart the Hive service.

Enabling URIs for Per-DB Policy Files

The `ADD JAR` command does *not* work with HiveServer2 and the Beeline client when Beeline runs on a different host. As an alternative to `ADD JAR`, Hive's *auxiliary paths* functionality should be used as described in the following steps.



Important: Enabling URIs in per-DB policy files introduces a security risk by allowing the owner of the db-level policy file to grant himself/herself load privileges to anything the `hive` user has read permissions for in HDFS (including data in other databases controlled by different db-level policy files).

Add the following string to the Java configuration options for HiveServer2 during startup.

```
-Dsentry.allow.uri.db.policyfile=true
```

Using User-Defined Functions with HiveServer2

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator, Full Administrator**)

The `ADD JAR` command does *not* work with HiveServer2 and the Beeline client when Beeline runs on a different host. As an alternative to `ADD JAR`, Hive's *auxiliary paths* functionality should be used. There are some differences in the procedures for creating permanent functions and temporary functions. For detailed instructions, see [User-Defined Functions \(UDFs\) with HiveServer2 Using Cloudera Manager](#).

Enabling Policy File Authorization for Hive

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator, Full Administrator**)

1. Ensure the [Prerequisites](#) on page 392 have been satisfied.
2. The Hive warehouse directory (`/user/hive/warehouse` or any path you specify as `hive.metastore.warehouse.dir` in your `hive-site.xml`) must be owned by the Hive user and group.
 - Permissions on the warehouse directory must be set as follows (see following Note for caveats):
 - **771** on the directory itself (for example, `/user/hive/warehouse`)
 - **771** on all subdirectories (for example, `/user/hive/warehouse/mysubdir`)
 - All files and subdirectories should be owned by `hive:hive`

For example:

```
$ sudo -u hdfs hdfs dfs -chmod -R 771 /user/hive/warehouse
$ sudo -u hdfs hdfs dfs -chown -R hive:hive /user/hive/warehouse
```

If you have enabled Kerberos on your cluster, you must kinit as the `hdfs` user before you set permissions.

For example:

```
sudo -u hdfs kinit -kt <hdfs.keytab> hdfs
sudo -u hdfs hdfs dfs -chmod -R 771 /user/hive/warehouse
$ sudo -u hdfs hdfs dfs -chown -R hive:hive /user/hive/warehouse
```



Note:

- If you set `hive.warehouse.subdir.inherit.perms` to `true` in `hive-site.xml`, the permissions on the subdirectories will be set when you set permissions on the warehouse directory itself.
- If a user has access to any object in the warehouse, that user will be able to execute `use default`. This ensures that `use default` commands issued by legacy applications work when Sentry is enabled. Note that you can protect objects in the default database (or any other database) by means of a policy file.



Important: These instructions override the recommendations in the Hive section of the CDH 5 Installation Guide.

3. Disable impersonation for HiveServer2:
 - a. Go to the Hive service.
 - b. Click the **Configuration** tab.
 - c. Select **Scope > HiveServer2**.
 - d. Select **Category > All**.
 - e. Locate the **HiveServer2 Enable Impersonation** property or search for it by typing its name in the Search box.
 - f. Under the **HiveServer2** role group, deselect the **HiveServer2 Enable Impersonation** property.
 - g. Click **Save Changes** to commit the changes.
4. Create the Sentry policy file, `sentry-provider.ini`, as an HDFS file.
5. Enable the Hive user to submit MapReduce jobs.
 - a. Go to the MapReduce service.
 - b. Click the **Configuration** tab.
 - c. Select **Scope > TaskTracker**.
 - d. Select **Category > Security**.
 - e. Locate the **Minimum User ID for Job Submission** property or search for it by typing its name in the Search box.
 - f. Set the **Minimum User ID for Job Submission** property to 0 (the default is 1000).
 - g. Click **Save Changes** to commit the changes.
 - h. Repeat steps 5.a-5.d for *every* TaskTracker role group for the MapReduce service that is associated with Hive, if more than one exists.
 - i. Restart the MapReduce service.
6. Enable the Hive user to submit YARN jobs.
 - a. Go to the YARN service.
 - b. Click the **Configuration** tab.
 - c. Select **Scope > NodeManager**.
 - d. Select **Category > Security**.
 - e. Ensure the **Allowed System Users** property includes the `hive` user. If not, add `hive`.
 - f. Click **Save Changes** to commit the changes.
 - g. Repeat steps 6.a-6.d for *every* NodeManager role group for the YARN service that is associated with Hive, if more than one exists.
 - h. Restart the YARN service.
7. Go to the Hive service.
8. Click the **Configuration** tab.
9. Select **Scope > Hive (Service-Wide)**.
10. Select **Category > Policy File Based Sentry**.
11. Select **Enable Sentry Authorization Using Policy Files**.
12. Click **Save Changes** to commit the changes.
13. Add the Hive user group to Sentry's admin groups.
 - a. Go to the Sentry service.
 - b. Click the **Configuration** tab.
 - c. Select **Scope > Sentry (Service-Wide)**.
 - d. Select **Category > Main**.

- e. Locate the **Admin Groups** property and add the `hive` group to the list. If an end user is in one of these admin groups, that user has administrative privileges on the Sentry Server.
- f. Click **Save Changes** to commit the changes.

14 Restart the cluster and HiveServer2 after changing these values, whether you use Cloudera Manager or not.

Configuring Group Access to the Hive Metastore

Minimum Required Role: [Configurator](#) (also provided by **Cluster Administrator**, **Full Administrator**)

You can configure the Hive Metastore to reject connections from users not listed in the Hive group proxy list (in HDFS). If you do not configure this override, the Hive Metastore will use the value in the core-site HDFS configuration. To configure the Hive group proxy list:

1. Go to the Hive service.
2. Click the **Configuration** tab.
3. Select **Scope > Hive (Service-Wide)**.
4. Select **Category > Proxy**.
5. In the **Hive Metastore Access Control and Proxy User Groups Override** property, specify a list of groups whose users are allowed to access the Hive Metastore. If you do not specify "*" (wildcard), you will be warned if the groups do not include `hive` and `impala` (if the Impala service is configured) in the list of groups.
6. Click **Save Changes**.
7. Restart the Hive service.

Enabling Policy File Authorization for Impala

For a cluster managed by Cloudera Manager, perform the following steps to enable policy file authorization for Impala.

1. Enable Sentry's policy file based authorization for Hive. For details, see [Enabling Policy File Authorization for Hive](#) on page 405.
2. Go to the Cloudera Manager Admin Console and go to the Impala service.
3. Click the **Configuration** tab.
4. Select **Scope > Impala (Service-Wide)**.
5. Select **Category > Policy File-Based Sentry**.
6. Select **Enable Sentry Authorization Using Policy Files**.
7. Click **Save Changes** to commit the changes.
8. Restart the Impala service.
9. Add the Impala user group to Sentry's admin groups.
 - a. Go to the Sentry service.
 - b. Click the **Configuration** tab.
 - c. Select **Scope > Sentry (Service-Wide)**.
 - d. Select **Category > Main**.
 - e. Locate the **Admin Groups** property and add the `impala` group to the list. If an end user is in one of these admin groups, that user has administrative privileges on the Sentry Server.
 - f. Click **Save Changes** to commit the changes.

For more details, see [Starting the impalad Daemon with Sentry Authorization Enabled](#) on page 412.

Enabling Sentry Authorization for Solr

Minimum Required Role: [Full Administrator](#)

1. Ensure the following requirements are satisfied:
 - Cloudera Search 1.1.1 or higher or CDH 5 or higher.
 - A secure Hadoop cluster.
2. Create the policy file `sentry-provider.ini` as an HDFS file. When you create the policy file `sentry-provider.ini` follow the instructions in the Policy File section in [Configuring Sentry for Search \(CDH\)](#)

Configuring Authorization

4) or [Search Authentication](#) on page 163. The file must be owned by owned by the `solr` user in the `solr` group, with `perms=600`. By default Cloudera Manager assumes the policy file is in the HDFS location `/user/solr/sentry`. To configure the location:

- a. Go to the Solr service.
- b. Click the **Configuration** tab.
- c. Select **Scope > SOLR (Service-Wide)**.
- d. Select **Category > Policy File Based Sentry**.
- e. Locate the **Sentry Global Policy File** property.
- f. Modify the path in the **Sentry Global Policy File** property.
- g. Select **Enable Sentry Authorization**.
- h. Click **Save Changes**.

3. Restart the Solr service.

For more details, see [Enabling Sentry Authorization for Search using the Command Line](#) on page 422.

Configuring Sentry to Enable BDR Replication

Cloudera recommends the following steps when configuring Sentry and [data replication](#) is enabled.

- Group membership should be managed outside of Sentry (as typically OS groups, LDAP groups, and so on are managed) and replication for them also should be handled outside of Cloudera Manager.
- In Cloudera Manager, set up [HDFS replication](#) for the Sentry files of the databases that are being replicated (separately using Hive replication).
- On the source cluster:
 - Use a separate Sentry policy file for every database
 - Avoid placing any group or role info (except for server admin info) in the global Sentry policy file (to avoid manual replication/merging with the global file on the target cluster)
 - To avoid manual fix up of URI privileges, ensure that the URIs for the data are the same on both the source and target cluster
- On the target cluster:
 - In the global Sentry policy file, manually add the *DB name - DB file* mapping entries for the databases being replicated
 - Manually copy the server admin info from the global Sentry policy file on the source to the policy on the target cluster
 - For the databases being replicated, avoid adding more privileges (adding tables specific to target cluster may sometimes require adding extra privileges to allow access to those tables). If any target cluster specific privileges absolutely need to be added for a database, add them to the global Sentry policy file on the target cluster since the per database files would be overwritten periodically with source versions during scheduled replication.

Configuring Sentry Policy File Authorization Using the Command Line

This topic describes how to configure Sentry policy files and enable policy file authorization for unmanaged CDH services using the command line.

Configuring User to Group Mappings

Hadoop Groups

Set the `hive.sentry.provider` property in `sentry-site.xml`.

```
<property>
<name>hive.sentry.provider</name>
<value>org.apache.sentry.provider.file.HadoopGroupResourceAuthorizationProvider</value>
</property>
```


Local Groups

1. Define local groups in the `[users]` section of the [Policy File](#) on page 394. For example:

```
[users]
user1 = group1, group2, group3
user2 = group2, group3
```

2. Modify Sentry configuration as follows:

In `sentry-site.xml`, set `hive.sentry.provider` as follows:

```
<property>
<name>hive.sentry.provider</name>
<value>org.apache.sentry.provider.file.LocalGroupResourceAuthorizationProvider</value>
</property>
```

Enabling URIs for Per-DB Policy Files

The `ADD JAR` command does *not* work with HiveServer2 and the Beeline client when Beeline runs on a different host. As an alternative to `ADD JAR`, Hive's *auxiliary paths* functionality should be used as described in the following steps.



Important: Enabling URIs in per-DB policy files introduces a security risk by allowing the owner of the db-level policy file to grant himself/herself load privileges to anything the `hive` user has read permissions for in HDFS (including data in other databases controlled by different db-level policy files).

Add the following string to the Java configuration options for HiveServer2 during startup.

```
-Dsentry.allow.uri.db.policyfile=true
```

Using User-Defined Functions with HiveServer2

The `ADD JAR` command does *not* work with HiveServer2 and the Beeline client when Beeline runs on a different host. As an alternative to `ADD JAR`, Hive's *auxiliary paths* functionality should be used as described in the following steps. There are some differences in the procedures for creating permanent functions and temporary functions. For detailed instructions, see [User-Defined Functions \(UDFs\) with HiveServer2 Using the Command Line](#).

Enabling Policy File Authorization for Hive

Prerequisites

In addition to the [Prerequisites](#) on page 392 above, make sure that the following are true:

- The Hive warehouse directory (`/user/hive/warehouse` or any path you specify as `hive.metastore.warehouse.dir` in your `hive-site.xml`) must be owned by the Hive user and group.
 - Permissions on the warehouse directory must be set as follows (see following Note for caveats):
 - **771** on the directory itself (for example, `/user/hive/warehouse`)
 - **771** on all subdirectories (for example, `/user/hive/warehouse/mysubdir`)
 - All files and subdirectories should be owned by `hive:hive`

For example:

```
$ sudo -u hdfs hdfs dfs -chmod -R 771 /user/hive/warehouse
$ sudo -u hdfs hdfs dfs -chown -R hive:hive /user/hive/warehouse
```

If you have enabled Kerberos on your cluster, you must kinit as the `hdfs` user before you set permissions.

For example:

```
sudo -u hdfs kinit -kt <hdfs.keytab> hdfs
sudo -u hdfs hdfs dfs -chmod -R 771 /user/hive/warehouse
$ sudo -u hdfs hdfs dfs -chown -R hive:hive /user/hive/warehouse
```

**Note:**

- If you set `hive.warehouse.subdir.inherit.perms` to `true` in `hive-site.xml`, the permissions on the subdirectories will be set when you set permissions on the warehouse directory itself.
- If a user has access to any object in the warehouse, that user will be able to execute `use default`. This ensures that `use default` commands issued by legacy applications work when Sentry is enabled. Note that you can protect objects in the default database (or any other database) by means of a policy file.



Important: These instructions override the recommendations in the Hive section of the CDH 5 Installation Guide.

- HiveServer2 impersonation must be turned off.
- The Hive user must be able to submit MapReduce jobs. You can ensure that this is true by setting the minimum user ID for job submission to 0. Edit the `taskcontroller.cfg` file and set `min.user.id=0`.

To enable the Hive user to submit YARN jobs, add the user `hive` to the `allowed.system.users` configuration property. Edit the `container-executor.cfg` file and add `hive` to the `allowed.system.users` property. For example,

```
allowed.system.users = nobody,impala,hive
```

**Important:**

- You must restart the cluster and HiveServer2 after changing this value, whether you use Cloudera Manager or not.
- These instructions override the instructions under [Configuring MRv1 Security](#) on page 109
- These instructions override the instructions under [Configuring YARN Security](#) on page 112

- Add the Hive, Impala and Hue groups to Sentry's admin groups. If an end user is in one of these admin groups, that user has administrative privileges on the Sentry Server.

```
<property>
  <name>sentry.service.admin.group</name>
  <value>hive,impala,hue</value>
</property>
```

Configuration Changes Required

To enable Sentry, add the following properties to `hive-site.xml`:

```
<property>
<name>hive.server2.session.hook</name>
<value>org.apache.sentry.binding.hive.HiveAuthzBindingSessionHook</value>
</property>

<property>
<name>hive.sentry.conf.url</name>
<value></value>
<description>sentry-site.xml file location</description>
</property>

<property>
<name>hive.metastore.client.impl</name>
<value>org.apache.sentry.binding.metastore.SentryHiveMetaStoreClient</value>
```

```
<description>Sets custom Hive Metastore client which Sentry uses to filter out
metadata.</description>
</property>
```

Securing the Hive Metastore

It's important that the Hive metastore be secured. If you want to override the Kerberos prerequisite for the Hive metastore, set the `sentry.hive.testing.mode` property to `true` to allow Sentry to work with weaker authentication mechanisms. Add the following property to the HiveServer2 and Hive metastore's `sentry-site.xml`:

```
<property>
  <name>sentry.hive.testing.mode</name>
  <value>true</value>
</property>
```

Impala does not require this flag to be set.



Warning: Cloudera strongly recommends against enabling this property in production. Use Sentry's testing mode only in test environments.

You can turn on Hive metastore security using the instructions in [Cloudera Security](#). To secure the Hive metastore; see [Hive Metastore Server Security Configuration](#) on page 139.

Enabling Policy File Authorization for Impala

First, enable Sentry's policy file based authorization for Hive. For details, see [Enabling Policy File Authorization for Hive](#) on page 409.

See [Enabling Sentry Authorization for Impala](#) on page 411 for details on configuring Impala to work with Sentry policy files.

Enabling Sentry in Cloudera Search

See [Enabling Sentry in Cloudera Search for CDH 5](#) on page 425 for details on securing Cloudera Search with Sentry.

Enabling Sentry Authorization for Impala

Authorization determines which users are allowed to access which resources, and what operations they are allowed to perform. In Impala 1.1 and higher, you use the Sentry open source project for authorization. Sentry adds a fine-grained authorization framework for Hadoop. By default (when authorization is not enabled), Impala does all read and write operations with the privileges of the `impala` user, which is suitable for a development/test environment but not for a secure production environment. When authorization is enabled, Impala uses the OS user ID of the user who runs `impala-shell` or other client program, and associates various privileges with each user.



Note: Sentry is typically used in conjunction with Kerberos authentication, which defines which hosts are allowed to connect to each server. Using the combination of Sentry and Kerberos prevents malicious users from being able to connect by creating a named account on an untrusted machine. See [Enabling Kerberos Authentication for Impala](#) on page 154 for details about Kerberos authentication.

The Sentry Privilege Model

Privileges can be granted on different objects in the schema. Any privilege that can be granted is associated with a level in the object hierarchy. If a privilege is granted on a container object in the hierarchy, the child object automatically inherits it. This is the same privilege model as Hive and other database systems such as MySQL.

The object hierarchy for Impala covers Server, URI, Database, Table, and Column. (The Table privileges apply to views as well; anywhere you specify a table name, you can specify a view name instead.) Column-level authorization is available in CDH 5.5 / Impala 2.3 and higher, as described in [Column-level Authorization](#) on page 380. Previously, you constructed views to query specific columns and assigned privileges based on the views rather than the base tables.

A restricted set of privileges determines what you can do with each object:

SELECT privilege

Lets you read data from a table or view, for example with the `SELECT` statement, the `INSERT . . . SELECT` syntax, or `CREATE TABLE . . . LIKE`. Also required to issue the `DESCRIBE` statement or the `EXPLAIN` statement for a query against a particular table. Only objects for which a user has this privilege are shown in the output for `SHOW DATABASES` and `SHOW TABLES` statements. The `REFRESH` statement and `INVALIDATE METADATA` statements only access metadata for tables for which the user has this privilege.

INSERT privilege

Lets you write data to a table. Applies to the `INSERT` and `LOAD DATA` statements.

ALL privilege

Lets you create or modify the object. Required to run DDL statements such as `CREATE TABLE`, `ALTER TABLE`, or `DROP TABLE` for a table, `CREATE DATABASE` or `DROP DATABASE` for a database, or `CREATE VIEW`, `ALTER VIEW`, or `DROP VIEW` for a view. Also required for the URI of the “location” parameter for the `CREATE EXTERNAL TABLE` and `LOAD DATA` statements.

Privileges can be specified for a table or view before that object actually exists. If you do not have sufficient privilege to perform an operation, the error message does not disclose if the object exists or not.

Originally, privileges were encoded in a policy file, stored in HDFS. This mode of operation is still an option, but the emphasis of privilege management is moving towards being SQL-based. Although currently Impala does not have `GRANT` or `REVOKE` statements, Impala can make use of privileges assigned through `GRANT` and `REVOKE` statements done through Hive. The mode of operation with `GRANT` and `REVOKE` statements instead of the policy file requires that a special Sentry service be enabled; this service stores, retrieves, and manipulates privilege information stored inside the metastore database.

Starting the impalad Daemon with Sentry Authorization Enabled

To run the `impalad` daemon with authorization enabled, you add one or more options to the `IMPALA_SERVER_ARGS` declaration in the `/etc/default/impala` configuration file:

- The `-server_name` option turns on Sentry authorization for Impala. The authorization rules refer to a symbolic server name, and you specify the name to use as the argument to the `-server_name` option.
- If you specify just `-server_name`, Impala uses the Sentry service for authorization, relying on the results of `GRANT` and `REVOKE` statements issued through Hive. (This mode of operation is available in Impala 1.4.0 and higher.) Prior to Impala 1.4.0, or if you want to continue storing privilege rules in the policy file, also specify the `-authorization_policy_file` option as in the following item.
- Specifying the `-authorization_policy_file` option in addition to `-server_name` makes Impala read privilege information from a policy file, rather than from the metastore database. The argument to the `-authorization_policy_file` option specifies the HDFS path to the policy file that defines the privileges on different schema objects.

For example, you might adapt your `/etc/default/impala` configuration to contain lines like the following. To use the Sentry service rather than the policy file:

```
IMPALA_SERVER_ARGS=" \  
-server_name=server1 \  
. . .
```

Or to use the policy file, as in releases prior to Impala 1.4:

```
IMPALA_SERVER_ARGS=" \  
-authorization_policy_file=/user/hive/warehouse/auth-policy.ini \  
-server_name=server1 \  
. . .
```

The preceding examples set up a symbolic name of `server1` to refer to the current instance of Impala. This symbolic name is used in the following ways:

- In an environment managed by Cloudera Manager, the server name is specified through **Impala (Service-Wide) > Category > Advanced > Sentry Service** and **Hive > Service-Wide > Advanced > Sentry Service**. The values must

be the same for both, so that Impala and Hive can share the privilege rules. Restart the Impala and Hive services after setting or changing this value.

- In an environment not managed by Cloudera Manager, you specify this value for the `sentry.hive.server` property in the `sentry-site.xml` configuration file for Hive, as well as in the `-server_name` option for `impalad`.
If the `impalad` daemon is not already running, start it as described in [Starting Impala](#). If it is already running, restart it with the command `sudo /etc/init.d/impala-server restart`. Run the appropriate commands on all the nodes where `impalad` normally runs.
- If you use the mode of operation using the policy file, the rules in the `[roles]` section of the policy file refer to this same `server1` name. For example, the following rule sets up a role `report_generator` that lets users with that role query any table in a database named `reporting_db` on a node where the `impalad` daemon was started up with the `-server_name=server1` option:

```
[roles]
report_generator = server=server1->db=reporting_db->table=*->action=SELECT
```

When `impalad` is started with one or both of the `-server_name=server1` and `-authorization_policy_file` options, Impala authorization is enabled. If Impala detects any errors or inconsistencies in the authorization settings or the policy file, the daemon refuses to start.

Using Impala with the Sentry Service (CDH 5.1 or higher only)

When you use the Sentry service rather than the policy file, you set up privileges through `GRANT` and `REVOKE` statement in either Impala or Hive, then both components use those same privileges automatically. (Impala added the `GRANT` and `REVOKE` statements in Impala 2.0.0 / CDH 5.2.0.)

Hive already had `GRANT` and `REVOKE` statements prior to CDH 5.1, but those statements were not production-ready. CDH 5.1 is the first release where those statements use the Sentry framework and are considered GA level. If you used the Hive `GRANT` and `REVOKE` statements prior to CDH 5.1, you must set up these privileges with the CDH 5.1 versions of `GRANT` and `REVOKE` to take advantage of Sentry authorization.

For information about using the updated Hive `GRANT` and `REVOKE` statements, see [Sentry service](#) topic in the *CDH 5 Security Guide*.

Using Impala with the Sentry Policy File

The policy file is a file that you put in a designated location in HDFS, and is read during the startup of the `impalad` daemon when you specify both the `-server_name` and `-authorization_policy_file` startup options. It controls which objects (databases, tables, and HDFS directory paths) can be accessed by the user who connects to `impalad`, and what operations that user can perform on the objects.



Note:

In CDH 5 and higher, Cloudera recommends managing privileges through SQL statements, as described in [Using Impala with the Sentry Service \(CDH 5.1 or higher only\)](#) on page 413. If you are still using policy files, plan to migrate to the new approach some time in the future.

The location of the policy file is listed in the `auth-site.xml` configuration file. To minimize overhead, the security information from this file is cached by each `impalad` daemon and refreshed automatically, with a default interval of 5 minutes. After making a substantial change to security policies, restart all Impala daemons to pick up the changes immediately.

Policy File Location and Format

The policy file uses the familiar `.ini` format, divided into the major sections `[groups]` and `[roles]`. There is also an optional `[databases]` section, which allows you to specify a specific policy file for a particular database, as explained in [Using Multiple Policy Files for Different Databases](#) on page 417. Another optional section, `[users]`, allows you to

Configuring Authorization

override the OS-level mapping of users to groups; that is an advanced technique primarily for testing and debugging, and is beyond the scope of this document.

In the `[groups]` section, you define various categories of users and select which roles are associated with each category. The group and user names correspond to Linux groups and users on the server where the `impalad` daemon runs.

The group and user names in the `[groups]` section correspond to Linux groups and users on the server where the `impalad` daemon runs. When you access Impala through the `impalad` interpreter, for purposes of authorization, the user is the logged-in Linux user and the groups are the Linux groups that user is a member of. When you access Impala through the ODBC or JDBC interfaces, the user and password specified through the connection string are used as login credentials for the Linux server, and authorization is based on that user name and the associated Linux group membership.

In the `[roles]` section, you a set of roles. For each role, you specify precisely the set of privileges is available. That is, which objects users with that role can access, and what operations they can perform on those objects. This is the lowest-level category of security information; the other sections in the policy file map the privileges to higher-level divisions of groups and users. In the `[groups]` section, you specify which roles are associated with which groups. The group and user names correspond to Linux groups and users on the server where the `impalad` daemon runs. The privileges are specified using patterns like:

```
server=server_name->db=database_name->table=table_name->action=SELECT
server=server_name->db=database_name->table=table_name->action=CREATE
server=server_name->db=database_name->table=table_name->action=ALL
```

For the `server_name` value, substitute the same symbolic name you specify with the `impalad -server_name` option. You can use `*` wildcard characters at each level of the privilege specification to allow access to all such objects. For example:

```
server=impala-host.example.com->db=default->table=t1->action=SELECT
server=impala-host.example.com->db=*->table=*->action=CREATE
server=impala-host.example.com->db=*->table=audit_log->action=SELECT
server=impala-host.example.com->db=default->table=t1->action=*
```

When authorization is enabled, Impala uses the policy file as a *whitelist*, representing every privilege available to any user on any object. That is, only operations specified for the appropriate combination of object, role, group, and user are allowed; all other operations are not allowed. If a group or role is defined multiple times in the policy file, the last definition takes precedence.

To understand the notion of whitelisting, set up a minimal policy file that does not provide any privileges for any object. When you connect to an Impala node where this policy file is in effect, you get no results for `SHOW DATABASES`, and an error when you issue any `SHOW TABLES`, `USE database_name`, `DESCRIBE table_name`, `SELECT`, and or other statements that expect to access databases or tables, even if the corresponding databases and tables exist.

The contents of the policy file are cached, to avoid a performance penalty for each query. The policy file is re-checked by each `impalad` node every 5 minutes. When you make a non-time-sensitive change such as adding new privileges or new users, you can let the change take effect automatically a few minutes later. If you remove or reduce privileges, and want the change to take effect immediately, restart the `impalad` daemon on all nodes, again specifying the `-server_name` and `-authorization_policy_file` options so that the rules from the updated policy file are applied.

Examples of Policy File Rules for Security Scenarios

The following examples show rules that might go in the policy file to deal with various authorization-related scenarios. For illustration purposes, this section shows several very small policy files with only a few rules each. In your environment, typically you would define many roles to cover all the scenarios involving your own databases, tables, and applications, and a smaller number of groups, whose members are given the privileges from one or more roles.

A User with No Privileges

If a user has no privileges at all, that user cannot access any schema objects in the system. The error messages do not disclose the names or existence of objects that the user is not authorized to read.

This is the experience you want a user to have if they somehow log into a system where they are not an authorized Impala user. In a real deployment with a filled-in policy file, a user might have no privileges because they are not a member of any of the relevant groups mentioned in the policy file.

Examples of Privileges for Administrative Users

When an administrative user has broad access to tables or databases, the associated rules in the `[roles]` section typically use wildcards and/or inheritance. For example, in the following sample policy file, `db=*` refers to all databases and `db=*->table=*` refers to all tables in all databases.

Omitting the rightmost portion of a rule means that the privileges apply to all the objects that could be specified there. For example, in the following sample policy file, the `all_databases` role has all privileges for all tables in all databases, while the `one_database` role has all privileges for all tables in one specific database. The `all_databases` role does not grant privileges on URIs, so a group with that role could not issue a `CREATE TABLE` statement with a `LOCATION` clause. The `entire_server` role has all privileges on both databases and URIs within the server.

```
[groups]
supergroup = all_databases

[roles]
read_all_tables = server=server1->db=*->table=*->action=SELECT
all_tables = server=server1->db=*->table=*
all_databases = server=server1->db=*
one_database = server=server1->db=test_db
entire_server = server=server1
```

A User with Privileges for Specific Databases and Tables

If a user has privileges for specific tables in specific databases, the user can access those things but nothing else. They can see the tables and their parent databases in the output of `SHOW TABLES` and `SHOW DATABASES`, use the appropriate databases, and perform the relevant actions (`SELECT` and/or `INSERT`) based on the table privileges. To actually create a table requires the `ALL` privilege at the database level, so you might define separate roles for the user that sets up a schema and other users or applications that perform day-to-day operations on the tables.

The following sample policy file shows some of the syntax that is appropriate as the policy file grows, such as the `#` comment syntax, `\` continuation syntax, and comma separation for roles assigned to groups or privileges assigned to roles.

```
[groups]
cloudera = training_sysadmin, instructor
visitor = student

[roles]
training_sysadmin = server=server1->db=training, \
server=server1->db=instructor_private, \
server=server1->db=lesson_development
instructor = server=server1->db=training->table=*->action=*, \
server=server1->db=instructor_private->table=*->action=*, \
server=server1->db=lesson_development->table=lesson*
# This particular course is all about queries, so the students can SELECT but not INSERT
or CREATE/DROP.
student = server=server1->db=training->table=lesson_*->action=SELECT
```

Privileges for Working with External Data Files

When data is being inserted through the `LOAD DATA` statement, or is referenced from an HDFS location outside the normal Impala database directories, the user also needs appropriate permissions on the URIs corresponding to those HDFS locations.

In this sample policy file:

- The `external_table` role lets us insert into and query the Impala table, `external_table.sample`.
- The `staging_dir` role lets us specify the HDFS path `/user/cloudera/external_data` with the `LOAD DATA` statement. Remember, when Impala queries or loads data files, it operates on all the files in that directory, not just a single file, so any Impala `LOCATION` parameters refer to a directory rather than an individual file.

Configuring Authorization

- We included the IP address and port of the Hadoop name node in the HDFS URI of the `staging_dir` rule. We found those details in `/etc/hadoop/conf/core-site.xml`, under the `fs.default.name` element. That is what we use in any roles that specify URIs (that is, the locations of directories in HDFS).
- We start this example after the table `external_table.sample` is already created. In the policy file for the example, we have already taken away the `external_table_admin` role from the `cloudera` group, and replaced it with the lesser-privileged `external_table` role.
- We assign privileges to a subdirectory underneath `/user/cloudera` in HDFS, because such privileges also apply to any subdirectories underneath. If we had assigned privileges to the parent directory `/user/cloudera`, it would be too likely to mess up other files by specifying a wrong location by mistake.
- The `cloudera` under the `[groups]` section refers to the `cloudera` group. (In the demo VM used for this example, there is a `cloudera` user that is a member of a `cloudera` group.)

Policy file:

```
[groups]
cloudera = external_table, staging_dir

[roles]
external_table_admin = server=server1->db=external_table
external_table = server=server1->db=external_table->table=sample->action=*
staging_dir =
server=server1->uri=hdfs://127.0.0.1:8020/user/cloudera/external_data->action=*
```

impala-shell session:

```
[localhost:21000] > use external_table;
Query: use external_table
[localhost:21000] > show tables;
Query: show tables
Query finished, fetching results ...
+-----+
| name |
+-----+
| sample |
+-----+
Returned 1 row(s) in 0.02s

[localhost:21000] > select * from sample;
Query: select * from sample
Query finished, fetching results ...
+-----+
| x |
+-----+
| 1 |
| 5 |
| 150 |
+-----+
Returned 3 row(s) in 1.04s

[localhost:21000] > load data inpath '/user/cloudera/external_data' into table sample;
Query: load data inpath '/user/cloudera/external_data' into table sample
Query finished, fetching results ...
+-----+
| summary |
+-----+
| Loaded 1 file(s). Total files in destination location: 2 |
+-----+
Returned 1 row(s) in 0.26s
[localhost:21000] > select * from sample;
Query: select * from sample
Query finished, fetching results ...
+-----+
| x |
+-----+
| 2 |
| 4 |
| 6 |
| 8 |
+-----+
```



```

64738 |
49152 |
1     |
5     |
150  |
+-----+
Returned 9 row(s) in 0.22s

[localhost:21000] > load data inpath '/user/cloudera/unauthorized_data' into table
sample;
Query: load data inpath '/user/cloudera/unauthorized_data' into table sample
ERROR: AuthorizationException: User 'cloudera' does not have privileges to access:
hdfs://127.0.0.1:8020/user/cloudera/unauthorized_data

```

Separating Administrator Responsibility from Read and Write Privileges

Remember that to create a database requires full privilege on that database, while day-to-day operations on tables within that database can be performed with lower levels of privilege on specific table. Thus, you might set up separate roles for each database or application: an administrative one that could create or drop the database, and a user-level one that can access only the relevant tables.

For example, this policy file divides responsibilities between users in 3 different groups:

- Members of the `supergroup` group have the `training_sysadmin` role and so can set up a database named `training`.
- Members of the `cloudera` group have the `instructor` role and so can create, insert into, and query any tables in the `training` database, but cannot create or drop the database itself.
- Members of the `visitor` group have the `student` role and so can query those tables in the `training` database.

```

[groups]
supergroup = training_sysadmin
cloudera = instructor
visitor = student

[roles]
training_sysadmin = server=server1->db=training
instructor = server=server1->db=training->table=*->action=*
student = server=server1->db=training->table=*->action=SELECT

```

Using Multiple Policy Files for Different Databases

For an Impala cluster with many databases being accessed by many users and applications, it might be cumbersome to update the security policy file for each privilege change or each new database, table, or view. You can allow security to be managed separately for individual databases, by setting up a separate policy file for each database:

- Add the optional `[databases]` section to the main policy file.
- Add entries in the `[databases]` section for each database that has its own policy file.
- For each listed database, specify the HDFS path of the appropriate policy file.

For example:

```

[databases]
# Defines the location of the per-DB policy files for the 'customers' and 'sales'
databases.
customers = hdfs://ha-nn-uri/etc/access/customers.ini
sales = hdfs://ha-nn-uri/etc/access/sales.ini

```

To enable URIs in per-DB policy files, add the following string in the Cloudera Manager field **Impala Service Environment Advanced Configuration Snippet (Safety Valve)**:

```
JAVA_TOOL_OPTIONS="-Dsentry.allow.uri.db.policyfile=true"
```



Important: Enabling URIs in per-DB policy files introduces a security risk by allowing the owner of the db-level policy file to grant himself/herself load privileges to anything the `impala` user has read permissions for in HDFS (including data in other databases controlled by different db-level policy files).

Setting Up Schema Objects for a Secure Impala Deployment

Remember that in your role definitions, you specify privileges at the level of individual databases and tables, or all databases or all tables within a database. To simplify the structure of these rules, plan ahead of time how to name your schema objects so that data with different authorization requirements is divided into separate databases.

If you are adding security on top of an existing Impala deployment, remember that you can rename tables or even move them between databases using the `ALTER TABLE` statement. In Impala, creating new databases is a relatively inexpensive operation, basically just creating a new directory in HDFS.

You can also plan the security scheme and set up the policy file before the actual schema objects named in the policy file exist. Because the authorization capability is based on whitelisting, a user can only create a new database or table if the required privilege is already in the policy file: either by listing the exact name of the object being created, or a * wildcard to match all the applicable objects within the appropriate container.

Privilege Model and Object Hierarchy

Privileges can be granted on different objects in the schema. Any privilege that can be granted is associated with a level in the object hierarchy. If a privilege is granted on a container object in the hierarchy, the child object automatically inherits it. This is the same privilege model as Hive and other database systems such as MySQL.

The kinds of objects in the schema hierarchy are:

```
Server
URI
Database
Table
```

The server name is specified by the `-server_name` option when `impalad` starts. Specify the same name for all `impalad` nodes in the cluster.

URIs represent the HDFS paths you specify as part of statements such as `CREATE EXTERNAL TABLE` and `LOAD DATA`. Typically, you specify what look like UNIX paths, but these locations can also be prefixed with `hdfs://` to make clear that they are really URIs. To set privileges for a URI, specify the name of a directory, and the privilege applies to all the files in that directory and any directories underneath it.

In CDH 5.5 / Impala 2.3 and higher, you can specify privileges for individual columns, as described in [Column-level Authorization](#) on page 380. Formerly, to specify read privileges at this level, you created a view that queried specific columns and/or partitions from a base table, and gave `SELECT` privilege on the view but not the underlying table.

URIs must start with either `hdfs://` or `file://`. If a URI starts with anything else, it will cause an exception and the policy file will be invalid. When defining URIs for HDFS, you must also specify the NameNode. For example:

```
data_read = server=server1->uri=file:///path/to/dir, \
server=server1->uri=hdfs://namenode:port/path/to/dir
```



Warning:

Because the NameNode host and port must be specified, Cloudera strongly recommends you use High Availability (HA). This ensures that the URI will remain constant even if the NameNode changes.

```
data_read = server=server1->uri=file:///path/to/dir, \
server=server1->uri=hdfs://ha-nn-uri/path/to/dir
```

Table 37: Valid privilege types and objects they apply to

Privilege	Object
INSERT	DB, TABLE
SELECT	DB, TABLE, COLUMN
ALL	SERVER, TABLE, DB, URI



Note:

Although this document refers to the ALL privilege, currently if you use the policy file mode, you do not use the actual keyword ALL in the policy file. When you code role entries in the policy file:

- To specify the ALL privilege for a server, use a role like `server=server_name`.
- To specify the ALL privilege for a database, use a role like `server=server_name->db=database_name`.
- To specify the ALL privilege for a table, use a role like `server=server_name->db=database_name->table=table_name->action=*`.

Operation	Scope	Privileges Required	URI
CREATE DATABASE	SERVER	ALL	
DROP DATABASE	DATABASE	ALL	
CREATE TABLE	DATABASE	ALL	
DROP TABLE	TABLE	ALL	
CREATE VIEW -This operation is allowed if you have column-level SELECT access to the columns being used.	DATABASE; SELECT on TABLE;	ALL	
ALTER VIEW -This operation is allowed if you have column-level SELECT access to the columns being used.	VIEW/TABLE	ALL	
DROP VIEW	VIEW/TABLE	ALL	
ALTER TABLE .. ADD COLUMNS	TABLE	ALL	
ALTER TABLE .. REPLACE COLUMNS	TABLE	ALL	
ALTER TABLE .. CHANGE column	TABLE	ALL	
ALTER TABLE .. RENAME	TABLE	ALL	
ALTER TABLE .. SET TBLPROPERTIES	TABLE	ALL	
ALTER TABLE .. SET FILEFORMAT	TABLE	ALL	
ALTER TABLE .. SET LOCATION	TABLE	ALL	URI
ALTER TABLE .. ADD PARTITION	TABLE	ALL	

Operation	Scope	Privileges Required	URI
ALTER TABLE .. ADD PARTITION location	TABLE	ALL	URI
ALTER TABLE .. DROP PARTITION	TABLE	ALL	
ALTER TABLE .. PARTITION SET FILEFORMAT	TABLE	ALL	
SHOW CREATE TABLE	TABLE	SELECT/INSERT	
SHOW PARTITIONS	TABLE	SELECT/INSERT	
SHOW TABLES -Output includes all the tables for which the user has table-level privileges and all the tables for which the user has some column-level privileges.	TABLE	SELECT/INSERT	
SHOW GRANT ROLE -Output includes an additional field for any column-level privileges.	TABLE	SELECT/INSERT	
DESCRIBE TABLE -Output shows <i>all</i> columns if the user has table level-privileges or <code>SELECT</code> privilege on at least one table column	TABLE	SELECT/INSERT	
LOAD DATA	TABLE	INSERT	URI
SELECT -You can grant the <code>SELECT</code> privilege on a view to give users access to specific columns of a table they do not otherwise have access to. -See Column-level Authorization on page 380 for details on allowed column-level operations.	VIEW/TABLE; COLUMN	SELECT	
INSERT OVERWRITE TABLE	TABLE	INSERT	
CREATE TABLE .. AS SELECT -This operation is allowed if you have column-level <code>SELECT</code> access to the columns being used.	DATABASE; SELECT on TABLE	ALL	
USE <dbName>	Any		
CREATE FUNCTION	SERVER	ALL	
ALTER TABLE .. SET SERDEPROPERTIES	TABLE	ALL	

Operation	Scope	Privileges Required	URI
ALTER TABLE .. PARTITION SET SERDEPROPERTIES	TABLE	ALL	
EXPLAIN SELECT	TABLE; COLUMN	SELECT	
EXPLAIN INSERT	TABLE; COLUMN	INSERT	
INVALIDATE METADATA	SERVER	ALL	
INVALIDATE METADATA <table name>	TABLE	SELECT/INSERT	
REFRESH <table name> or REFRESH <table name> PARTITION (<partition_spec>)	TABLE	SELECT/INSERT	
DROP FUNCTION	SERVER	ALL	
COMPUTE STATS	TABLE	ALL	

Debugging Failed Sentry Authorization Requests

Sentry logs all facts that lead up to authorization decisions at the debug level. If you do not understand why Sentry is denying access, the best way to debug is to temporarily turn on debug logging:

- In Cloudera Manager, add `log4j.logger.org.apache.sentry=DEBUG` to the logging settings for your service through the corresponding **Logging Safety Valve** field for the Impala, Hive Server 2, or Solr Server services.
- On systems not managed by Cloudera Manager, add `log4j.logger.org.apache.sentry=DEBUG` to the `log4j.properties` file on each host in the cluster, in the appropriate configuration directory for each service.

Specifically, look for exceptions and messages such as:

```
FilePermission server..., RequestPermission server..., result [true|false]
```

which indicate each evaluation Sentry makes. The `FilePermission` is from the policy file, while `RequestPermission` is the privilege required for the query. A `RequestPermission` will iterate over all appropriate `FilePermission` settings until a match is found. If no matching privilege is found, Sentry returns `false` indicating “Access Denied”.

Managing Sentry for Impala through Cloudera Manager

To enable the Sentry service for Impala and Hive, set **Hive/Impala > Service-Wide > Sentry Service** parameter to the Sentry service. Then restart Impala and Hive. Simply adding Sentry service as a dependency and restarting enables Impala and Hive to use the Sentry service.

To set the server name to use when granting server level privileges, set the **Hive > Service-Wide > Advanced > Server Name for Sentry Authorization** parameter. When using Sentry with the Hive Metastore, you can specify the list of users that are allowed to bypass Sentry Authorization in Hive Metastore using **Hive > Service-Wide > Security > Bypass Sentry Authorization Users**. These are usually service users that already ensure all activity has been authorized.



Note: The **Hive/Impala > Service-Wide > Policy File Based Sentry** tab contains parameters only relevant to configuring Sentry using policy files. In particular, make sure that **Enable Sentry Authorization using Policy Files** parameter is unchecked when using the Sentry service. Cloudera Manager throws a validation error if you attempt to configure the Sentry service and policy file at the same time.

The DEFAULT Database in a Secure Deployment

Because of the extra emphasis on granular access controls in a secure deployment, you should move any important or sensitive information out of the `DEFAULT` database into a named database whose privileges are specified in the

policy file. Sometimes you might need to give privileges on the `DEFAULT` database for administrative reasons; for example, as a place you can reliably specify with a `USE` statement when preparing to drop a database.

Enabling Sentry Authorization for Search using the Command Line

Sentry enables role-based, fine-grained authorization for Cloudera Search. Sentry can apply a range of restrictions to various tasks, such as accessing data, managing configurations through config objects, or creating collections. Restrictions are consistently applied, regardless of the way users attempt to complete actions. For example, restricting access to data in a collection restricts that access whether queries come from the command line, from a browser, Hue, or through the admin console.

- You can use either Cloudera Manager or the following command-line instructions to complete this configuration.
- This information applies specifically to CDH 5.6.x. If you use an earlier version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

For information on enabling Sentry authorization using Cloudera Manager, see [Configuring Sentry Policy File Authorization Using Cloudera Manager](#) on page 404.

Follow the instructions below to configure Sentry under CDH 4.5 or higher or CDH 5. Sentry is included in the Search installation.



Note: Sentry for Search depends on Kerberos authentication. For additional information on using Kerberos with Search, see [Search Authentication](#) on page 163.

This document describes configuring Sentry for Cloudera Search. For information about alternate ways to configure Sentry or for information about installing Sentry for other services, see:

- [Enabling Sentry Authorization for Solr](#) for instructions for using Cloudera Manager to configure Search Authorization with Sentry.
- [Overview of Impala Security](#) for instructions on using Impala with Sentry.
- [Sentry Installation](#) for additional information on Sentry installation.

Using Roles and Privileges with Sentry

Sentry uses a role-based privilege model. A role is a set of rules for accessing a given Solr collection or Solr config. Access to each collection is governed by privileges: `Query`, `Update`, or `All (*)`. Access to each config provides access to all actions: `All (*)`.

- A rule for the `Query` privilege on collection `logs` would be formulated as follows:

```
collection=logs->action=Query
```

- A rule for the `All (*)` privilege on config `myConfig` would be formulated as follows:

```
config=myConfig->action=*
```

No action implies `*` and `*` is the only valid action. Because `config` objects only support `*`, the following `config` privilege is invalid:

```
config=myConfig->action=Update
```

Note that `config` objects cannot be combined with `collection` objects in a single privilege. For example, the following combinations are illegal:

```
config=myConfig->collection=myCollection->action=*
```

```
collection=myCollection->config=myConfig
```

You may specify these privileges separately. For example:

```
myRole = collection=myCollection->action=QUERY, config=myConfig->action=*
```

A role can contain multiple such rules, separated by commas. For example the `engineer_role` might contain the `Query` privilege for `hive_logs` and `hbase_logs` collections, and the `Update` privilege for the `current_bugs` collection. You would specify this as follows:

```
engineer_role = collection=hive_logs->action=Query, collection=hbase_logs->action=Query,
collection=current_bugs->action=Update
```

Using Users and Groups with Sentry

- A user is an entity that is permitted by the Kerberos authentication system to access the Search service.
- A group connects the authentication system with the authorization system. It is a set of one or more users who have been granted one or more authorization roles. Sentry allows a set of roles to be configured for a group.
- A configured group provider determines a user's affiliation with a group. The current release supports HDFS-backed groups and locally configured groups. For example,

```
dev_ops = dev_role, ops_role
```

Here the group `dev_ops` is granted the roles `dev_role` and `ops_role`. The members of this group can complete searches that are allowed by these roles.

User to Group Mapping

You can configure Sentry to use either Hadoop groups or groups defined in the policy file.



Important: You can use either Hadoop groups or local groups, but not both at the same time. Use local groups if you want to do a quick proof-of-concept. For production, use Hadoop groups.

To configure Hadoop groups:

Set the `sentry.provider` property in `sentry-site.xml` to `org.apache.sentry.provider.file.HadoopGroupResourceAuthorizationProvider`.

By default, this uses local shell groups. See the [Group Mapping](#) section of the HDFS Permissions Guide for more information.

In this case, Sentry uses the Hadoop configuration described in [Configuring LDAP Group Mappings](#) on page 350. Cloudera Manager automatically uses this configuration. In a deployment not managed by Cloudera Manager, manually set these configuration parameters in the `hadoop-conf` file that is passed to Solr.

OR

To configure local groups:

1. Define local groups in a `[users]` section of the Sentry Policy file. For example:

```
[users]
user1 = group1, group2, group3
user2 = group2, group3
```

Configuring Authorization

2. In `sentry-site.xml`, set `search.sentry.provider` as follows:

```
<property>
  <name>sentry.provider</name>
  <value>org.apache.sentry.provider.file.LocalGroupResourceAuthorizationProvider</value>
</property>
```

Using Policy Files with Sentry

The sections that follow contain notes on creating and maintaining the policy file.



Warning: An invalid configuration disables all authorization while logging an exception.

Storing the Policy File

Considerations for storing the policy file(s) include:

1. Replication count - Because the file is read for each query, you should increase this; 10 is a reasonable value.
2. Updating the file - Updates to the file are only reflected when the Solr process is restarted.

Defining Roles

Keep in mind that role definitions are not cumulative; the newer definition replaces the older one. For example, the following results in `role1` having `privilege2`, not `privilege1` and `privilege2`.

```
role1 = privilege1
role1 = privilege2
```

Sample Sentry Configuration

This section provides a sample configuration.



Note: Sentry with CDH Search does not support multiple policy files. Other implementations of Sentry such as Sentry for Hive do support different policy files for different databases, but Sentry for CDH Search has no such support for multiple policies.

Policy File

The following is an example of a CDH Search policy file. The `sentry-provider.ini` would exist in an HDFS location such as `hdfs://ha-nn-uri/user/solr/sentry/sentry-provider.ini`. This location must be readable by Solr.



Note: Use separate policy files for each Sentry-enabled service. Using one file for multiple services results in each service failing on the other services' entries. For example, with a combined Hive and Search file, Search would fail on Hive entries and Hive would fail on Search entries.

`sentry-provider.ini`

```
[groups]
# Assigns each Hadoop group to its set of roles
engineer = engineer_role
ops = ops_role
dev_ops = engineer_role, ops_role
hbase_admin = hbase_admin_role

[roles]
# The following grants all access to source_code.
# "collection = source_code" can also be used as syntactic
# sugar for "collection = source_code->action=*"
engineer_role = collection = source_code->action=*
```



```
# The following imply more restricted access.
ops_role = collection = hive_logs->action=Query
dev_ops_role = collection = hbase_logs->action=Query

#give hbase_admin_role the ability to create/delete/modify the hbase_logs collection
#as well as to update the config for the hbase_logs collection, called hbase_logs_config.
hbase_admin_role = collection=admin->action=*, collection=hbase_logs->action=*,
config=hbase_logs_config->action=*
```

Sentry Configuration File

Sentry stores the configuration as well as privilege policies in files. The `sentry-site.xml` file contains configuration options such as privilege policy file location. The [Policy File](#) on page 424 contains the privileges and groups. It has a `.ini` file format and should be stored on HDFS.

The following is an example of a `sentry-site.xml` file.

sentry-site.xml

```
<configuration>
  <property>
    <name>hive.sentry.provider</name>

<value>org.apache.sentry.provider.file.HadoopGroupResourceAuthorizationProvider</value>
  </property>

  <property>
    <name>sentry.solr.provider.resource</name>
    <value>/path/to/authz-provider.ini</value>
    <!--
      If the HDFS configuration files (core-site.xml, hdfs-site.xml)
      pointed to by SOLR_HDFS_CONFIG in /etc/default/solr
      point to HDFS, the path will be in HDFS;
      alternatively you could specify a full path,
      e.g. :hdfs://namenode:port/path/to/authz-provider.ini
    -->
  </property>
```

Enabling Sentry in Cloudera Search for CDH 5

You can enable Sentry using Cloudera Manager or by manually modifying files. For more information on enabling Sentry using Cloudera Manager, see [Configuring Sentry Policy File Authorization Using Cloudera Manager](#) on page 404 and [Enabling Sentry Authorization for Solr](#) on page 407.

Sentry is enabled with addition of two properties to `/etc/default/solr` or `/opt/cloudera/parcels/CDH-*/etc/default/solr`.

- If you are using configs, you must configure the proper `config=myConfig` permissions as described in [Using Roles and Privileges with Sentry](#) on page 422.
- In a Cloudera Manager deployment, these properties are added automatically when you click **Enable Sentry Authorization** in the Solr configuration page in Cloudera Manager.
- In a deployment not managed by Cloudera Manager, you must make these changes yourself. The variable `SOLR_AUTHORIZATION_SENTRY_SITE` specifies the path to `sentry-site.xml`. The variable `SOLR_AUTHORIZATION_SUPERUSER` specifies the first part of `SOLR_KERBEROS_PRINCIPAL`. This is `solr` for the majority of users, as `solr` is the default. Settings are of the form:

```
SOLR_AUTHORIZATION_SENTRY_SITE=/location/to/sentry-site.xml
SOLR_AUTHORIZATION_SUPERUSER=solr
```

To enable sentry collection-level authorization checking on a new collection, the `instancedir` for the collection must use a modified version of `solrconfig.xml` with Sentry integration. Each collection has a separate `solrconfig.xml` file, meaning you can define different behavior for each collection. The command `solrctl instancedir --generate` generates two versions of `solrconfig.xml`: the standard `solrconfig.xml` without sentry integration, and the

Configuring Authorization

sentry-integrated version called `solrconfig.xml.secure`. To use the sentry-integrated version, replace `solrconfig.xml` with `solrconfig.xml.secure` before creating the `instancedir`.

You can enable Sentry on an existing collection. The process varies depending on whether you are using a config or `instancedir`.

Enabling Sentry on Collections using configs

If you have a collection that is using a non-secured config, you can enable Sentry security on that collection by modifying the collection to use a secure config. The config in use must not be immutable, otherwise it cannot be changed. To update an existing non-immutable config:

1. Delete the existing config using the `solrctl config --delete` command. For example:

```
solrctl config --delete myManaged
```

2. Create a new non-immutable config using the `solrctl config --create` command. Use a sentry-enabled template such as `managedTemplateSecure`. The new config must have the same name as the config being replaced. For example:

```
solrctl config --create myManaged managedTemplateSecure -p immutable=false
```

3. Reload the collection using `solrctl collection --reload` command.

```
solrctl collection --reload myCollection
```

For a list of all available config templates, see [Included Immutable Config Templates](#).

Enabling Sentry on Collections using instancedirs

If you have a collection that is using a non-secured `instancedir` configuration, you can enable Sentry security on that collection by modifying the settings that are stored in `instancedir`. For example, you might have an existing collection named `foo` and a standard `solrconfig.xml`. By default, collections are stored in `instancedirs` that use the collection's name, which is `foo` in this case.

If your collection uses an unmodified `solrconfig.xml` file, you can enable Sentry by replacing the existing the `solrconfig.xml` file. If your collection uses a `solrconfig.xml` that contains modifications you want to preserve, you can attempt to use a `diff` tool to find an integrate changes in to the secure template.

To enable Sentry on an existing collection without preserving customizations



Warning: Executing the following commands replaces your existing `solrconfig.xml` file. Any customizations to this file will be lost.

```
# generate a fresh instancedir
solrctl instancedir --generate foosecure
# download the existing instancedir from ZK into subdirectory foo
solrctl instancedir --get foo foo
# replace the existing solrconfig.xml with the sentry-enabled one
cp foosecure/conf/solrconfig.xml.secure foo/conf/solrconfig.xml
# update the instancedir in ZK
solrctl instancedir --update foo foo
# reload the collection
solrctl collection --reload foo
```

To enable Sentry on an existing collection and preserve customizations

Generate a new `instancedir`, compare the differences between the default `solrconfig.xml` and `solrconfig.xml.secure` files, and then add the elements that are unique to `solrconfig.xml.secure` to the file that your environment is using.

1. Generate a fresh instancedir:

```
solrctl instancedir --generate foo
```

2. Compare the solrconfig.xml and solrconfig.xml.secure:

```
diff foo/conf/solrconfig.xml foo/conf/solrconfig.xml.secure
```

3. Add the elements that are unique to solrconfig.xml.secure to your existing solrconfig.xml file. You might complete this process by manually editing your existing solrconfig.xml file or by using a merge tool.

Note: If you have modified or specified additional request handlers, consider that Sentry:

- Supports protecting additional query request handlers by adding a search component, which should be shown in the diff.
- Supports protecting additional update request handlers with Sentry by adding an updateRequestProcessorChain, which should be shown in the diff.
- Does not support protecting modified or specified additional "special" request handlers like analysis handlers or admin handlers.

4. Reload the collection:

```
solrctl collection --reload foo
```

Providing Document-Level Security Using Sentry

For role-based access control of a collection, an administrator modifies a Sentry role so it has query, update, or administrative access, as described above.

Collection-level authorization is useful when the access control requirements for the documents in the collection are the same, but users may want to restrict access to a subset of documents in a collection. This finer-grained restriction could be achieved by defining separate collections for each subset, but this is difficult to manage, requires duplicate documents for each collection, and requires that these documents be kept synchronized.

Document-level access control solves this issue by associating authorization tokens with each document in the collection. This enables granting Sentry roles access to sets of documents in a collection.

Document-Level Security Model

Document-level security depends on a chain of relationships between users, groups, roles, and documents.

- Users are assigned to groups.
- Groups are assigned to roles.
- Roles are stored as "authorization tokens" in a specified field in the documents.

Document-level security supports restricting which documents can be viewed by which users. Access is provided by adding roles as "authorization tokens" to a specified document field. Conversely, access is implicitly denied by omitting roles from the specified field. In other words, in a document-level security enabled environment, a user might submit a query that matches a document; if the user is not part of a group that has a role has been granted access to the document, the result is not returned.

For example, Alice might belong to the administrators group. The administrators group may belong to the doc-mgmt role. A document could be ingested and the doc-mgmt role could be added at ingest time. In such a case, if Alice submitted a query that matched the document, Search would return the document, since Alice is then allowed to see any document with the "doc-mgmt" authorization token.

Similarly, Bob might belong to the guests group. The guests group may belong to the public-browser role. If Bob tried the same query as Alice, but the document did not have the public-browser role, Search would not return the result because Bob does not belong to a group that is associated with a role that has access.

Note that collection-level authorization rules still apply, if enabled. Even if Alice is able to view a document given document-level authorization rules, if she is not allowed to query the collection, the query will fail.

Roles are typically added to documents when those documents are ingested, either via the standard Solr APIs or, if using morphlines, the `setValues` morphline command.

Enabling Document-Level Security

Cloudera Search supports document-level security in Search for CDH 5.1 and later. Document-level security requires collection-level security. Configuring collection-level security is described earlier in this topic.

Document-level security is disabled by default, so the first step in using document-level security is to enable the feature by modifying the `solrconfig.xml.secure` file. Remember to replace the `solrconfig.xml` with this file, as described in [Enabling Sentry in Cloudera Search for CDH 5](#) on page 425.

To enable document-level security, change `solrconfig.xml.secure`. The default file contents are as follows:

```
<searchComponent name="queryDocAuthorization">
  <!-- Set to true to enabled document-level authorization -->

  <bool name="enabled">false</bool>

  <!-- Field where the auth tokens are stored in the document -->
  <str name="sentryAuthField">sentry_auth</str>

  <!-- Auth token defined to allow any role to access the document.
  Uncomment to enable. -->

  <!--<str name="allRolesToken">*</str>-->
</searchComponent>
```

- The `enabled` Boolean determines whether document-level authorization is enabled. To enable document level security, change this setting to `true`.
- The `sentryAuthField` string specifies the name of the field that is used for storing authorization information. You can use the default setting of `sentry_auth` or you can specify some other string to be used for assigning values during ingest.



Note: This field must exist as an explicit or dynamic field in the schema for the collection you are creating with document-level security. `sentry_auth` exists in the default `schema.xml`, which is automatically generated and can be found in the same directory as `solrconfig.xml`.

for the collection you are creating with document-level security. `Schema.xml` is in the generated configuration in the same directory as the `solrconfig.xml`

- The `allRolesToken` string represents a special token defined to allow any role access to the document. By default, this feature is disabled. To enable this feature, uncomment the specification and specify the token. This token should be different from the name of any sentry role to avoid collision. By default it is `"*"`. This feature is useful when first configuring document level security or it can be useful in granting all roles access to a document when the set of roles may change. See [Best Practices](#) for additional information.

Best Practices

Using `allRolesToken`

You may want to grant every user that belongs to a role access to certain documents. One way to accomplish this is to specify all known roles in the document, but this requires updating or re-indexing the document if you add a new role. Alternatively, an `allUser` role, specified in the Sentry `.ini` file, could contain all valid groups, but this role would need to be updated every time a new group was added to the system. Instead, specifying `allRolesToken` allows any user that belongs to a valid role to access the document. This access requires no updating as the system evolves.

In addition, `allRolesToken` may be useful for transitioning a deployment to use document-level security. Instead of having to define all the roles upfront, all the documents can be specified with `allRolesToken` and later modified as the roles are defined.

Consequences of Document-Level Authorization Only Affecting Queries

Document-level security does not prevent users from modifying documents or performing other update operations on the collection. Update operations are only governed by collection-level authorization.

Document-level security can be used to prevent documents being returned in query results. If users are not granted access to a document, those documents are not returned even if that user submits a query that matches those documents. This does not have affect attempted updates.

Consequently, it is possible for a user to not have access to a set of documents based on document-level security, but to still be able to modify the documents via their collection-level authorization update rights. This means that a user can delete all documents in the collection. Similarly, a user might modify all documents, adding their authorization token to each one. After such a modification, the user could access any document via querying. Therefore, if you are restricting access using document-level security, consider granting collection-level update rights only to those users you trust and assume they will be able to access every document in the collection.

Limitations on Query Size

By default queries support up to 1024 Boolean clauses. As a result, queries containing more that 1024 clauses may cause errors. Because authorization information is added by Sentry as part of a query, using document-level security can increase the number of clauses. In the case where users belong to many roles, even simple queries can become quite large. If a query is too large, an error of the following form occurs:

```
org.apache.lucene.search.BooleanQuery$TooManyClauses: maxClauseCount is set to 1024
```

To change the supported number of clauses, edit the `maxBooleanClauses` setting in `solrconfig.xml`. For example, to allow 2048 clauses, you would edit the setting so it appears as follows:

```
<maxBooleanClauses>2048</maxBooleanClauses>
```

For `maxBooleanClauses` to be applied as expected, make any change to this value to all collections and then restart the service. You must make this change to all collections because this option modifies a global Lucene property, affecting all Solr cores. If different `solrconfig.xml` files have different values for this property, the effective value is determined per host, based on the first Solr core to be initialized.

Enabling Secure Impersonation

Secure Impersonation is a feature that allows a user to make requests as another user in a secure way. For example, to allow the following impersonations:

- User `hue` can make requests as any user from any host.
- User `foo` can make requests as any member of group `bar`, from `host1` or `host2`.

Configure the following properties in `/etc/default/solr` or `/opt/cloudera/parcels/CDH-*/etc/default/solr`:

```
SOLR_SECURITY_ALLOWED_PROXYUSERS=hue,foo
SOLR_SECURITY_PROXYUSER_hue_HOSTS=*
SOLR_SECURITY_PROXYUSER_hue_GROUPS=*
SOLR_SECURITY_PROXYUSER_foo_HOSTS=host1,host2
SOLR_SECURITY_PROXYUSER_foo_GROUPS=bar
```

`SOLR_SECURITY_ALLOWED_PROXYUSERS` lists all of the users allowed to impersonate. For a user `x` in `SOLR_SECURITY_ALLOWED_PROXYUSERS`, `SOLR_SECURITY_PROXYUSER_x_HOSTS` list the hosts `x` is allowed to connect from to impersonate, and `SOLR_SECURITY_PROXYUSERS_x_GROUPS` lists the groups that the users is allowed to impersonate members of. Both `GROUPS` and `HOSTS` support the wildcard `*` and both `GROUPS` and `HOSTS` must be defined for a specific user.



Note: Cloudera Manager has its own management of secure impersonation for Hue. To add additional users for Secure Impersonation, use the environment variable `safety` value for Solr to set the environment variables as above. Be sure to include `hue` in `SOLR_SECURITY_ALLOWED_PROXYUSERS` if you want to use secure impersonation for hue.

Debugging Failed Sentry Authorization Requests

Sentry logs all facts that lead up to authorization decisions at the debug level. If you do not understand why Sentry is denying access, the best way to debug is to temporarily turn on debug logging:

- In Cloudera Manager, add `log4j.logger.org.apache.sentry=DEBUG` to the logging settings for your service through the corresponding **Logging Safety Valve** field for the Impala, Hive Server 2, or Solr Server services.
- On systems not managed by Cloudera Manager, add `log4j.logger.org.apache.sentry=DEBUG` to the `log4j.properties` file on each host in the cluster, in the appropriate configuration directory for each service.

Specifically, look for exceptions and messages such as:

```
FilePermission server..., RequestPermission server..., result [true|false]
```

which indicate each evaluation Sentry makes. The `FilePermission` is from the policy file, while `RequestPermission` is the privilege required for the query. A `RequestPermission` will iterate over all appropriate `FilePermission` settings until a match is found. If no matching privilege is found, Sentry returns `false` indicating “Access Denied” .

Appendix: Authorization Privilege Model for Search

The tables below refer to the request handlers defined in the generated `solrconfig.xml.secure`. If you are not using this configuration file, the below may not apply.

`admin` is a special collection in sentry used to represent administrative actions. A non-administrative request may only require privileges on the collection or config on which the request is being performed. This is called either `collection1` or `config1` in this appendix. An administrative request may require privileges on both the `admin` collection and `collection1`. This is denoted as `admin, collection1` in the tables below.



Note: If no privileges are granted, no access is possible. For example, accessing the Solr Admin UI requires the `QUERY` privilege. If no users are granted the `QUERY` privilege, no access to the Solr Admin UI is possible.

Table 38: Privilege table for non-administrative request handlers

Request Handler	Required Collection Privilege	Collections that Require Privilege
select	QUERY	collection1
query	QUERY	collection1
get	QUERY	collection1
browse	QUERY	collection1
tvrh	QUERY	collection1
clustering	QUERY	collection1
terms	QUERY	collection1
elevate	QUERY	collection1
analysis/field	QUERY	collection1
analysis/document	QUERY	collection1
update	UPDATE	collection1

Request Handler	Required Collection Privilege	Collections that Require Privilege
update/json	UPDATE	collection1
update/csv	UPDATE	collection1

Table 39: Privilege table for collections admin actions

Collection Action	Required Collection Privilege	Collections that Require Privilege
create	UPDATE	admin, collection1
delete	UPDATE	admin, collection1
reload	UPDATE	admin, collection1
createAlias	UPDATE	admin, collection1 <div style="border: 1px solid black; padding: 5px; margin-top: 10px;"> <p>Note: collection1 here refers to the name of the alias, not the underlying collection(s). For example, <code>http://YOUR-HOST:8983/solr/admin/collections?action=CREATEALIAS&name=collection1&collections-underlying=collection1</code></p> </div>
deleteAlias	UPDATE	admin, collection1 <div style="border: 1px solid black; padding: 5px; margin-top: 10px;"> <p>Note: collection1 here refers to the name of the alias, not the underlying collection(s). For example, <code>http://YOUR-HOST:8983/solr/admin/collections?action=DELETEALIAS&name=collection1</code></p> </div>
syncShard	UPDATE	admin, collection1
splitShard	UPDATE	admin, collection1
deleteShard	UPDATE	admin, collection1

Table 40: Privilege table for core admin actions

Collection Action	Required Collection Privilege	Collections that Require Privilege
create	UPDATE	admin, collection1
rename	UPDATE	admin, collection1
load	UPDATE	admin, collection1
unload	UPDATE	admin, collection1
status	UPDATE	admin, collection1
persist	UPDATE	admin

Collection Action	Required Collection Privilege	Collections that Require Privilege
reload	UPDATE	admin, collection1
swap	UPDATE	admin, collection1
mergeIndexes	UPDATE	admin, collection1
split	UPDATE	admin, collection1
prepRecover	UPDATE	admin, collection1
requestRecover	UPDATE	admin, collection1
requestSyncShard	UPDATE	admin, collection1
requestApplyUpdates	UPDATE	admin, collection1

Table 41: Privilege table for Info and AdminHandlers

Request Handler	Required Collection Privilege	Collections that Require Privilege
LukeRequestHandler	QUERY	admin
SystemInfoHandler	QUERY	admin
SolrInfoMBeanHandler	QUERY	admin
PluginInfoHandler	QUERY	admin
ThreadDumpHandler	QUERY	admin
PropertiesRequestHandler	QUERY	admin
LogginHandler	QUERY, UPDATE (or *)	admin
ShowFileRequestHandler	QUERY	admin

Table 42: Privilege table for Config Admin actions

Config Action	Required Collection Privilege	Collections that Require Privilege	Required Config Privilege	Configs that Require Privilege
CREATE	UPDATE	admin	*	config1
DELETE	UPDATE	admin	*	config1

Configuring HBase Authorization

After you have configured HBase authentication as described in the previous section, you must establish authorization rules for the resources that a client is allowed to access. HBase currently allows you to establish authorization rules at the table, column and cell-level. Cell-level authorization is fully supported since CDH 5.2.

Understanding HBase Access Levels

HBase access levels are granted independently of each other and allow for different types of operations at a given scope.

- **Read (R)** - can read data at the given scope
- **Write (W)** - can write data at the given scope
- **Execute (X)** - can execute coprocessor endpoints at the given scope
- **Create (C)** - can create tables or drop tables (even those they did not create) at the given scope
- **Admin (A)** - can perform cluster operations such as balancing the cluster or assigning regions at the given scope

The possible scopes are:

- **Superuser** - superusers can perform any operation available in HBase, to any resource. The user who runs HBase on your cluster is a superuser, as are any principals assigned to the configuration property `hbase.superuser` in `hbase-site.xml` on the HMaster.
- **Global** - permissions granted at `global` scope allow the admin to operate on all tables of the cluster.
- **Namespace** - permissions granted at `namespace` scope apply to all tables within a given namespace.
- **Table** - permissions granted at `table` scope apply to data or metadata within a given table.
- **ColumnFamily** - permissions granted at `ColumnFamily` scope apply to cells within that ColumnFamily.
- **Cell** - permissions granted at `Cell` scope apply to that exact cell coordinate. This allows for policy evolution along with data. To change an ACL on a specific cell, write an updated cell with new ACL to the precise coordinates of the original. If you have a multi-versioned schema and want to update ACLs on all visible versions, you'll need to write new cells for all visible versions. The application has complete control over policy evolution. The exception is `append` and `increment` processing. `Appends` and `increments` can carry an ACL in the operation. If one is included in the operation, then it will be applied to the result of the `append` or `increment`. Otherwise, the ACL of the existing cell being appended to or incremented is preserved.

The combination of access levels and scopes creates a matrix of possible access levels that can be granted to a user. In a production environment, it is useful to think of access levels in terms of what is needed to do a specific job. The following list describes appropriate access levels for some common types of HBase users. It is important not to grant more access than is required for a given user to perform their required tasks.

- **Superusers** - In a production system, only the HBase user should have superuser access. In a development environment, an administrator may need superuser access in order to quickly control and manage the cluster. However, this type of administrator should usually be a `Global Admin` rather than a superuser.
- **Global Admins** - A `global admin` can perform tasks and access every table in HBase. In a typical production environment, an admin should not have `Read` or `Write` permissions to data within tables.
 - A global admin with `Admin` permissions can perform cluster-wide operations on the cluster, such as balancing, assigning or unassigning regions, or calling an explicit major compaction. This is an operations role.
 - A global admin with `Create` permissions can create or drop any table within HBase. This is more of a DBA-type role.

In a production environment, it is likely that different users will have only one of `Admin` and `Create` permissions.



Warning:

In the current implementation, a `Global Admin` with `Admin` permission can grant himself `Read` and `Write` permissions on a table and gain access to that table's data. For this reason, only grant `Global Admin` permissions to trusted user who actually need them.

Also be aware that a `Global Admin` with `Create` permission can perform a `Put` operation on the ACL table, simulating a `grant` or `revoke` and circumventing the authorization check for `Global Admin` permissions. This issue (but not the first one) is fixed in CDH 5.3 and higher, as well as CDH 5.2.1. It is not fixed in CDH 4.x or CDH 5.1.x.

Due to these issues, be cautious with granting `Global Admin` privileges.

- **Table Admins** - A table admin can perform administrative operations only on that table. A table admin with `Create` permissions can create snapshots from that table or restore that table from a snapshot. A table admin with `Admin` permissions can perform operations such as splits or major compactions on that table.
- **Users** - Users can read or write data, or both. Users can also execute coprocessor endpoints, if given `Executable` permissions.



Important:

If you are using Kerberos principal names when setting ACLs for users, note that Hadoop uses only the first part (short) of the Kerberos principal when converting it to the user name. Hence, for the principal `ann/fully.qualified.domain.name@YOUR-REALM.COM`, HBase ACLs should only be set for user `ann`.

Table 43: Real-World Example of Access Levels

This table shows some typical job descriptions at a hypothetical company and the permissions they might require in order to get their jobs done using HBase.

Job Title	Scope	Permissions	Description
Senior Administrator	Global	Admin, Create	Manages the cluster and gives access to Junior Administrators.
Junior Administrator	Global	Create	Creates tables and gives access to Table Administrators.
Table Administrator	Table	Admin	Maintains a table from an operations point of view.
Data Analyst	Table	Read	Creates reports from HBase data.
Web Application	Table	Read, Write	Puts data into HBase and uses HBase data to perform operations.

Further Reading

- [Access Control Matrix](#)
- [Security - Apache HBase Reference Guide](#)

Enable HBase Authorization

HBase authorization is built on top of the Coprocessors framework, specifically `AccessController Coprocessor`.



Note: Once the Access Controller coprocessor is enabled, any user who uses the HBase shell will be subject to access control. Access control will also be in effect for native (Java API) client access to HBase.

Enable HBase Authorization Using Cloudera Manager

1. Go to **Clusters** and select the HBase cluster.
2. Select **Configuration**.
3. Search for **HBase Secure Authorization** and select it.
4. Search for **HBase Service Advanced Configuration Snippet (Safety Valve) for hbase-site.xml** and enter the following into it to enable `hbase.security.exec.permission.checks`. Without this option, all users will continue to have access to execute endpoint coprocessors. This option is not enabled when you enable HBase Secure Authorization for backward compatibility.

```
<property>
  <name>hbase.security.exec.permission.checks</name>
```

```
<value>true</value>
</property>
```

5. Optionally, search for and configure **HBase Coprocessor Master Classes** and **HBase Coprocessor Region Classes**.

Enable HBase Authorization Using the Command Line



Important:

- If you use Cloudera Manager, do not use these command-line instructions.
- This information applies specifically to CDH 5.6.x. If you use a lower version of CDH, see the documentation for that version located at [Cloudera Documentation](#).

To enable HBase authorization, add the following properties to the `hbase-site.xml` file *on every HBase server host (Master or RegionServer)*:

```
<property>
  <name>hbase.security.authorization</name>
  <value>true</value>
</property>
<property>
  <name>hbase.security.exec.permission.checks</name>
  <value>true</value>
</property>
<property>
  <name>hbase.coprocessor.master.classes</name>
  <value>org.apache.hadoop.hbase.security.access.AccessController</value>
</property>
<property>
  <name>hbase.coprocessor.region.classes</name>
  <value>org.apache.hadoop.hbase.security.token.TokenProvider,org.apache.hadoop.hbase.security.access.AccessController</value>
</property>
```

Configure Access Control Lists for Authorization

Now that HBase has the security coprocessor enabled, you can set ACLs using the HBase shell. Start the HBase shell as usual.



Important:

The host running the shell must be configured with a keytab file as described in [Configuring Kerberos Authentication for HBase](#).

The commands that control ACLs take the following form. Group names are prefixed with the @ symbol.

```
hbase> grant <user> <permissions> [ @<namespace> [ <table>[ <column family>[ <column
qualifier> ] ] ] ] # grants permissions

hbase> revoke <user> [ @<namespace> [ <table> [ <column family> [ <column qualifier> ]
] ] # revokes permissions

hbase> user_permission <table>
# displays existing permissions
```

In the above commands, fields encased in <> are variables, and fields in [] are optional. The `permissions` variable must consist of zero or more character from the set "RWCA".

- **R** denotes read permissions, which is required to perform `Get`, `Scan`, or `Exists` calls in a given scope.
- **w** denotes write permissions, which is required to perform `Put`, `Delete`, `LockRow`, `UnlockRow`, `IncrementColumnValue`, `CheckAndDelete`, `CheckAndPut`, `Flush`, or `Compact` in a given scope.

Configuring Authorization

- **x** denotes execute permissions, which is required to execute coprocessor endpoints.
- **C** denotes create permissions, which is required to perform `Create`, `Alter`, or `Drop` in a given scope.
- **A** denotes admin permissions, which is required to perform `Enable`, `Disable`, `Snapshot`, `Restore`, `Clone`, `Split`, `MajorCompact`, `Grant`, `Revoke`, and `Shutdown` in a given scope.

Access Control List Example Commands

```
grant 'user1', 'RWC'  
grant 'user2', 'RW', 'tableA'  
grant 'user3', 'C', '@my_namespace'
```

Be sure to review the information in [Understanding HBase Access Levels](#) on page 432 to understand the implications of the different access levels.

Sensitive Data Redaction

Data redaction is the suppression of sensitive data, such as any personally identifiable information (PII). PII can be used on its own or with other information to identify or locate a single person, or to identify an individual in context. Enabling redaction allow you to transform PII to a pattern that does not contain any identifiable information. For example, you could replace all Social Security numbers (SSN) like 123-45-6789 with an unintelligible pattern like XXX-XX-XXXX, or replace only part of the SSN (XXX-XX-6789).

Although [encryption techniques](#) are available to protect Hadoop data, the underlying problem with using encryption is that an admin who has complete access to the cluster also access to unencrypted sensitive user data. Even users with appropriate ACLs on the data could have access to logs and queries where sensitive data might have leaked.

Data redaction provides compliance with industry regulations such as PCI and HIPAA, which require that access to PII be restricted to only those users whose jobs require such access. PII or other sensitive data must not be available through any other channels to users like cluster administrators or data analysts. This is because redaction only applies to any incidental leaks of data. For example, if a user already has the required permissions to access PII through queries, then query results will not be redacted.

Password Redaction

Starting with Cloudera Manager and CDH 5.5, passwords will no longer be accessible in cleartext through the Cloudera Manager UI or in the configuration files stored on disk. For components such as HDFS, HBase, Hive, and so on, that use core Hadoop, the feature has been implemented by using Hadoop's `CredentialProvider` interface to encrypt and store passwords inside a secure `creds.jceks` keystore file. For components such as Hue and Impala, that do not use core Hadoop, instead of the password, we use a `password_script = /path/to/script/that/will/emit/password.sh` parameter that, when run, writes the password to `stdout`. Passwords contained within Cloudera Manager and Cloudera Navigator properties have been redacted internally in Cloudera Manager.

However, the database password contained in Cloudera Manager Server's `/etc/cloudera-scm-server/db.properties` file has not been redacted. The `db.properties` file is managed by customers and is populated manually when the Cloudera Manager Server database is being set up for the first time. Since this occurs before the Cloudera Manager Server has even started, encrypting the contents of this file is a completely different challenge as compared to that of redacting configuration files.

Password redaction (not including log and query redaction) is enabled by default for deployments with Cloudera Manager 5.5 (or higher) managing CDH 5.5 (or higher). There are no user-visible controls to enable or disable this feature. It is expected to work out of the box. The primary places where you will encounter the effects of password redaction are:

- In the Cloudera Manager Admin Console, on the **Processes** page for a given role instance, passwords in the linked configuration files have been replaced by `*****`.
- On the Cloudera Manager Server and Agent hosts, all configuration files in the `/var/run/cloudera-scm-agent/process` directory will have their passwords replaced by `*****`.

Exceptions:

- Solr does not have this feature enabled.
- The database password contained in Cloudera Manager Server's `/etc/cloudera-scm-server/db.properties` file has not been redacted.

Scope - Log and Query Redaction

Data redaction in CDH targets sensitive SQL data and log files. Currently, you can enable or disable redaction for the whole cluster with a simple HDFS service-wide configuration change. Redaction is implemented with the assumption that sensitive information resides in the data itself, not the metadata. If you enable redaction for a file, only sensitive data inside the file is redacted. Metadata such as the name of the file or file owner is not redacted.

Sensitive Data Redaction

When data redaction is enabled, the following data is redacted:

- Logs in HDFS and any dependent cluster services. Log redaction is not available in Isilon-based clusters.
- Audit data sent to Cloudera Navigator
- SQL query strings displayed by Hue, Hive, and Impala.

Redaction Rules

Redaction is based on pattern matching. Use regular expressions to define redaction rules that search for patterns of sensitive information such as Social Security numbers, credit card numbers, and dates.

Use Cloudera Manager to create redaction rules that have the following components:

- **Search** - A regular expression matched against the data. If the expression matches any part of the data, the match is replaced by the contents of the replace string. For example, to redact credit card numbers, your regular expression is `\d{4}[^\\w]\d{4}[^\\w]\d{4}[^\\w]\d{4}`.
- **Replace** - The string used to replace the redacted data. For example, to replace any matched credit card digits with Xs, the Replace string value would be `XXXX-XXXX-XXXX-XXXX`.
- **Trigger** - An optional field that specifies a simple string to be searched for in the data. The redactor searches for matches to the search regular expression only if the string is found. If no trigger is specified, redaction occurs when the Search regular expression is matched. Using the Trigger field improves performance: simple string matching is faster than regular expression matching.

The following redaction rules are preconfigured (*not* enabled) in Cloudera Manager. The order in which the rules are specified is relevant. For example, in the list of rules below, credit card numbers are redacted first, followed by SSNs, email addresses, and finally, hostnames.

Redaction Rule	Search Expression	Replace Expression
Credit Card numbers (with separator)	<code>\d{4}[^\\w]\d{4}[^\\w]\d{4}[^\\w]\d{4}</code>	<code>XXXX-XXXX-XXXX-XXXX</code>
Social Security numbers (with separator)	<code>\d{3}[^\\w]\d{2}[^\\w]\d{4}</code>	<code>XXX-XX-XXXX</code>
Email addresses	<code>\b([A-Za-z0-9] [A-Za-z0-9][A-Za-z0-9\-\._])\n*\n[A-Za-z0-9])@((\n[A-Za-z0-9] [A-Za-z0-9]\n[A-Za-z0-9\-\._])*\n[A-Za-z0-9])\n ([A-Za-z0-9][A-Za-z0-9\-\._])*\n[A-Za-z0-9])\n\b</code>	<code>email@redacted.host</code>
Hostnames	<code>\b((\n[A-Za-z] [A-Za-z][A-Za-z0-9\-_])\n*\n[A-Za-z0-9])\n ([A-Za-z0-9][A-Za-z0-9\-_])*\n[A-Za-z0-9])\n\b</code>	<code>HOSTNAME.REDACTED</code>

Cloudera Manager API Redaction

Cloudera Manager API does not have redaction enabled by default. You can configure redaction of the sensitive items by specifying a JVM parameter for Cloudera Manager. When you set this parameter, API calls to Cloudera Manager for configuration data do not include the sensitive information. For more information, see [Redacting Sensitive Information from the Exported Configuration](#).

Enabling Log and Query Redaction Using Cloudera Manager

Cloudera recommends using the new layout in Cloudera Manager, instead of the classic layout, to enable redaction. The new layout allows you to add preconfigured redaction rules and test your rules inline. To enable log and query redaction in Cloudera Manager:

1. Go to the **HDFS** service.
2. Click the **Configuration** tab.
3. In the Search box, type *redaction* to bring up the following redaction properties.

Property	Description
Enable Log and Query Redaction	Check this checkbox to enable log and query redaction for the cluster.
Log and Query Redaction Policy	List of rules for redacting sensitive information from log files and query strings. Choose a preconfigured rule or add a custom rule. Test your rules by entering sample text into the Test Redaction Rules text box and click Test Redaction . If no rules match, the text you entered is returned unchanged.

4. Optionally, enter a reason for the configuration changes.
5. Click **Save Changes** to commit the changes.
6. Restart the cluster.

Configuring the Cloudera Navigator Data Management Component to Redact PII



Warning: Cloudera strongly recommends you use the cluster-wide log and query redaction feature described in the [previous section](#). The Cloudera Navigator configuration described as follows will *not* redact Navigator data for all CDH components.

You can specify credit card number patterns and other PII to be masked in audit events, in the properties of entities displayed in lineage diagrams, and in information retrieved from the Audit Server database and the Metadata Server persistent storage. Redacting data other than credit card numbers is not supported out-of-the-box with this Cloudera Navigator property. You may use a different regular expression to redact Social Security numbers or other PII. Masking is not applied to audit events and lineage entities that existed before the mask was enabled.

Minimum Required Role: [Navigator Administrator](#) (also provided by **Full Administrator**)

1. Do one of the following:
 - Select **Clusters > Cloudera Management Service > Cloudera Management Service**.
 - On the Status tab of the **Home > Status** tab, in **Cloudera Management Service** table, click the **Cloudera Management Service** link.
2. Click the **Configuration** tab.
3. Expand the **Navigator Audit Server Default Group** category.
4. Click the **Advanced** category.
5. Configure the **PII Masking Regular Expression** property with a regular expression that matches the credit card number formats to be masked. The default expression is:

```
(4[0-9]{12}(?:[0-9]{3})?)|(5[1-5][0-9]{14})|(3[47][0-9]{13})
|(3(?:0[0-5]|[68][0-9])[0-9]{11})|(6(?:011|5[0-9]{2})[0-9]{12})|((?:2131|1800|35\\d{3})\\d{11})
```

which is constructed from the following subexpressions:

- Visa - (4[0-9]{12}(?:[0-9]{3})?)
- MasterCard - (5[1-5][0-9]{14})
- American Express - (3[47][0-9]{13})
- Diners Club - (3(?:0[0-5]|[68][0-9])[0-9]{11})
- Discover - (6(?:011|5[0-9]{2})[0-9]{12})
- JCB - ((?:2131|1800|35\\d{3})\\d{11})

Sensitive Data Redaction

If the property is left blank, PII information is not masked.

6. Click **Save Changes** to commit the changes.

Overview of Impala Security

Impala includes a fine-grained authorization framework for Hadoop, based on the Sentry open source project. Sentry authorization was added in Impala 1.1.0. Together with the Kerberos authentication framework, Sentry takes Hadoop security to a new level needed for the requirements of highly regulated industries such as healthcare, financial services, and government. Impala also includes an auditing capability; Impala generates the audit data, the Cloudera Navigator product consolidates the audit data from all nodes in the cluster, and Cloudera Manager lets you filter, visualize, and produce reports. The auditing feature was added in Impala 1.1.1.

The Impala security features have several objectives. At the most basic level, security prevents accidents or mistakes that could disrupt application processing, delete or corrupt data, or reveal data to unauthorized users. More advanced security features and practices can harden the system against malicious users trying to gain unauthorized access or perform other disallowed operations. The auditing feature provides a way to confirm that no unauthorized access occurred, and detect whether any such attempts were made. This is a critical set of features for production deployments in large organizations that handle important or sensitive data. It sets the stage for multi-tenancy, where multiple applications run concurrently and are prevented from interfering with each other.

The material in this section presumes that you are already familiar with administering secure Linux systems. That is, you should know the general security practices for Linux and Hadoop, and their associated commands and configuration files. For example, you should know how to create Linux users and groups, manage Linux group membership, set Linux and HDFS file permissions and ownership, and designate the default permissions and ownership for new files. You should be familiar with the configuration of the nodes in your Hadoop cluster, and know how to apply configuration changes or run a set of commands across all the nodes.

The security features are divided into these broad categories:

authorization

Which users are allowed to access which resources, and what operations are they allowed to perform? Impala relies on the open source Sentry project for authorization. By default (when authorization is not enabled), Impala does all read and write operations with the privileges of the `impala` user, which is suitable for a development/test environment but not for a secure production environment. When authorization is enabled, Impala uses the OS user ID of the user who runs `impala-shell` or other client program, and associates various privileges with each user. See [Enabling Sentry Authorization for Impala](#) on page 411 for details about setting up and managing authorization.

authentication

How does Impala verify the identity of the user to confirm that they really are allowed to exercise the privileges assigned to that user? Impala relies on the Kerberos subsystem for authentication. See [Enabling Kerberos Authentication for Impala](#) on page 154 for details about setting up and managing authentication.

auditing

What operations were attempted, and did they succeed or not? This feature provides a way to look back and diagnose whether attempts were made to perform unauthorized operations. You use this information to track down suspicious activity, and to see where changes are needed in authorization policies. The audit data produced by this feature is collected by the Cloudera Manager product and then presented in a user-friendly form by the Cloudera Manager product. See [Auditing Impala Operations](#) for details about setting up and managing auditing.

These other topics in the *Security Guide* cover how Impala integrates with security frameworks such as Kerberos, LDAP, and Sentry:

- [Impala Authentication](#) on page 154
- [Enabling Sentry Authorization for Impala](#) on page 411

Security Guidelines for Impala

The following are the major steps to harden a cluster running Impala against accidents and mistakes, or malicious attackers trying to access sensitive data:

- Secure the `root` account. The `root` user can tamper with the `impalad` daemon, read and write the data files in HDFS, log into other user accounts, and access other system services that are beyond the control of Impala.
- Restrict membership in the `sudoers` list (in the `/etc/sudoers` file). The users who can run the `sudo` command can do many of the same things as the `root` user.
- Ensure the Hadoop ownership and permissions for Impala data files are restricted.
- Ensure the Hadoop ownership and permissions for Impala log files are restricted.
- Ensure that the Impala web UI (available by default on port 25000 on each Impala node) is password-protected. See [Impala Web User Interface for Debugging](#) for details.
- Create a policy file that specifies which Impala privileges are available to users in particular Hadoop groups (which by default map to Linux OS groups). Create the associated Linux groups using the `groupadd` command if necessary.
- The Impala authorization feature makes use of the HDFS file ownership and permissions mechanism; for background information, see the [CDH HDFS Permissions Guide](#). Set up users and assign them to groups at the OS level, corresponding to the different categories of users with different access levels for various databases, tables, and HDFS locations (URIs). Create the associated Linux users using the `useradd` command if necessary, and add them to the appropriate groups with the `usermod` command.
- Design your databases, tables, and views with database and table structure to allow policy rules to specify simple, consistent rules. For example, if all tables related to an application are inside a single database, you can assign privileges for that database and use the `*` wildcard for the table name. If you are creating views with different privileges than the underlying base tables, you might put the views in a separate database so that you can use the `*` wildcard for the database containing the base tables, while specifying the precise names of the individual views. (For specifying table or database names, you either specify the exact name or `*` to mean all the databases on a server, or all the tables and views in a database.)
- Enable authorization by running the `impalad` daemons with the `-server_name` and `-authorization_policy_file` options on all nodes. (The authorization feature does not apply to the `statedored` daemon, which has no access to schema objects or data files.)
- Set up authentication using Kerberos, to make sure users really are who they say they are.

Securing Impala Data and Log Files

One aspect of security is to protect files from unauthorized access at the filesystem level. For example, if you store sensitive data in HDFS, you specify permissions on the associated files and directories in HDFS to restrict read and write permissions to the appropriate users and groups.

If you issue queries containing sensitive values in the `WHERE` clause, such as financial account numbers, those values are stored in Impala log files in the Linux filesystem and you must secure those files also. For the locations of Impala log files, see [Using Impala Logging](#).

All Impala read and write operations are performed under the filesystem privileges of the `impala` user. The `impala` user must be able to read all directories and data files that you query, and write into all the directories and data files for `INSERT` and `LOAD DATA` statements. At a minimum, make sure the `impala` user is in the `hive` group so that it can access files and directories shared between Impala and Hive. See [User Account Requirements](#) for more details.

Setting file permissions is necessary for Impala to function correctly, but is not an effective security practice by itself:

- The way to ensure that only authorized users can submit requests for databases and tables they are allowed to access is to set up Sentry authorization, as explained in [Enabling Sentry Authorization for Impala](#) on page 411. With authorization enabled, the checking of the user ID and group is done by Impala, and unauthorized access is blocked by Impala itself. The actual low-level read and write requests are still done by the `impala` user, so you must have appropriate file and directory permissions for that user ID.

- You must also set up Kerberos authentication, as described in [Enabling Kerberos Authentication for Impala](#) on page 154, so that users can only connect from trusted hosts. With Kerberos enabled, if someone connects a new host to the network and creates user IDs that match your privileged IDs, they will be blocked from connecting to Impala at all from that host.

Installation Considerations for Impala Security

Impala 1.1 comes set up with all the software and settings needed to enable security when you run the `impalad` daemon with the new security-related options (`-server_name` and `-authorization_policy_file`). You do not need to change any environment variables or install any additional JAR files. In a cluster managed by Cloudera Manager, you do not need to change any settings in Cloudera Manager.

Securing the Hive Metastore Database

It is important to secure the Hive metastore, so that users cannot access the names or other information about databases and tables through the Hive client or by querying the metastore database. Do this by turning on Hive metastore security, using the instructions in the [CDH 5 Security Guide](#) for securing different Hive components:

- Secure the Hive Metastore.
- In addition, allow access to the metastore only from the HiveServer2 server, and then disable local access to the HiveServer2 server.

Securing the Impala Web User Interface

The instructions in this section presume you are familiar with the [.htpasswd mechanism](#) commonly used to password-protect pages on web servers.

Password-protect the Impala web UI that listens on port 25000 by default. Set up a `.htpasswd` file in the `$IMPALA_HOME` directory, or start both the `impalad` and `statedored` daemons with the `--webserver_password_file` option to specify a different location (including the filename).

This file should only be readable by the Impala process and machine administrators, because it contains (hashed) versions of passwords. The username / password pairs are not derived from Unix usernames, Kerberos users, or any other system. The `domain` field in the password file must match the domain supplied to Impala by the new command-line option `--webserver_authentication_domain`. The default is `mydomain.com`.

Impala also supports using HTTPS for secure web traffic. To do so, set `--webserver_certificate_file` to refer to a valid `.pem` TLS/SSL certificate file. Impala will automatically start using HTTPS once the TLS/SSL certificate has been read and validated. A `.pem` file is basically a private key, followed by a signed TLS/SSL certificate; make sure to concatenate both parts when constructing the `.pem` file.

If Impala cannot find or parse the `.pem` file, it prints an error message and quits.



Note:

If the private key is encrypted using a passphrase, Impala will ask for that passphrase on startup, which is not useful for a large cluster. In that case, remove the passphrase and make the `.pem` file readable only by Impala and administrators.

When you turn on TLS/SSL for the Impala web UI, the associated URLs change from `http://` prefixes to `https://`. Adjust any bookmarks or application code that refers to those URLs.

Configuring Secure Access for Impala Web Servers

Cloudera Manager supports two methods of authentication for secure access to the Impala Catalog Server, Daemon, and StateStoreweb servers: password-based authentication and TLS/SSL certificate authentication.

Authentication for the three types of daemons can be configured independently.

Configuring Password Authentication

1. Navigate to **Clusters > Impala Service > Configuration**.
2. Search for "password" using the Search box in the **Configuration** tab. This should display the password-related properties (Username and Password properties) for the Impala Daemon, StateStore, and Catalog Server. If there are multiple role groups configured for Impala Daemon instances, the search should display all of them.
3. Enter a username and password into these fields.
4. Click **Save Changes**, and restart the Impala service.

Now when you access the Web UI for the Impala Daemon, StateStore, or Catalog Server, you are asked to log in before access is granted.

Configuring TLS/SSL Certificate Authentication

1. Create or obtain an TLS/SSL certificate.
2. Place the certificate, in `.pem` format, on the hosts where the Impala Catalog Server and StateStore are running, and on each host where an Impala Daemon is running. It can be placed in any location (path) you choose. If all the Impala Daemons are members of the same role group, then the `.pem` file must have the same path on every host.
3. Navigate to **Clusters > Impala Service > Configuration**.
4. Search for "certificate" using the Search box in the **Configuration** tab. This should display the certificate file location properties for the Impala Catalog Server, Impala Daemon, and StateStore. If there are multiple role groups configured for Impala Daemon instances, the search should display all of them.
5. In the property fields, enter the full path name to the certificate file.
6. Click **Save Changes**, and restart the Impala service.



Important: If Cloudera Manager cannot find the `.pem` file on the host for a specific role instance, that role will fail to start.

When you access the Web UI for the Impala Catalog Server, Impala Daemon, and StateStore, `https` will be used.

Miscellaneous Topics

This section comprises miscellaneous security guide topics that you may find useful once you have secured your cluster with authentication, encryption and authorization techniques.

Jsvc, Task Controller and Container Executor Programs

This section contains information about the following Hadoop security programs:

MRv1 and YARN: The jsvc Program

The `jsvc` program is part of the `bigtop-jsvc` package and installed in either `/usr/lib/bigtop-utils/jsvc` or `/usr/libexec/bigtop-utils/jsvc` depending on the particular Linux flavor.

`jsvc` ([more info](#)) is used to start the `DataNode` listening on low port numbers. Its entry point is the `SecureDataNodeStarter` class, which implements the `Daemon` interface that `jsvc` expects. `jsvc` is run as `root`, and calls the `SecureDataNodeStarter.init(...)` method while running as `root`. Once the `SecureDataNodeStarter` class has finished initializing, `jsvc` sets the effective UID to be the `hdfs` user, and then calls `SecureDataNodeStarter.start(...)`. `SecureDataNodeStarter` then calls the regular `DataNode` entry point, passing in a reference to the privileged resources it previously obtained.

MRv1 Only: The Linux TaskController Program

A `setuid` binary called `task-controller` is part of the `hadoop-0.20-mapreduce` package and is installed in either `/usr/lib/hadoop-0.20-mapreduce/sbin/Linux-amd64-64/task-controller` or `/usr/lib/hadoop-0.20-mapreduce/sbin/Linux-i386-32/task-controller`.

This `task-controller` program, which is used on MRv1 only, allows the `TaskTracker` to run tasks under the Unix account of the user who submitted the job in the first place. It is a `setuid` binary that must have a very specific set of permissions and ownership to function correctly. In particular, it must:

1. Be owned by `root`
2. Be owned by a group that contains only the user running the MapReduce daemons
3. Be `setuid`
4. Be group readable and executable

This corresponds to the ownership `root:mapred` and the permissions `4754`.

Here is the output of `ls` on a correctly-configured `Task-controller`:

```
-rwsr-xr-- 1 root mapred 30888 Mar 18 13:03 task-controller
```

The `TaskTracker` will check for this configuration on start up, and fail to start if the `Task-controller` is not configured correctly.

YARN Only: The Linux Container Executor Program

A `setuid` binary called `container-executor` is part of the `hadoop-yarn` package and is installed in `/usr/lib/hadoop-yarn/bin/container-executor`.

This `container-executor` program, which is used on YARN only and supported on GNU/Linux only, runs the containers as the user who submitted the application. It requires all user accounts to be created on the cluster hosts where the containers are launched. It uses a `setuid` executable that is included in the Hadoop distribution. The `NodeManager` uses this executable to launch and kill containers. The `setuid` executable switches to the user who has submitted the application and launches or kills the containers. For maximum security, this executor sets up restricted permissions and user/group ownership of local files and directories used by the containers such as the shared objects, jars,

intermediate files, and log files. As a result, only the application owner and NodeManager can access any of the local files/directories including those localized as part of the distributed cache.


Parcel Deployments

In a parcel deployment the `container-executor` file is located inside the parcel at `/opt/cloudera/parcels/CDH/lib/hadoop-yarn/bin/container-executor`. For the `/usr/lib` mount point, `setuid` should not be a problem. However, the parcel could easily be located on a different mount point. If you are using a parcel, make sure the mount point for the parcel directory is without the `nosuid` option.

The `container-executor` program must have a very specific set of permissions and ownership to function correctly. In particular, it must:

1. Be owned by root
2. Be owned by a group that contains only the user running the YARN daemons
3. Be `setuid`
4. Be group readable and executable. This corresponds to the ownership `root:yarn` and the permissions `6050`.

```
---Sr-s--- 1 root yarn 91886 2012-04-01 19:54 container-executor
```

 **Important:** Configuration changes to the Linux container executor could result in local NodeManager directories (such as `usercache`) being left with incorrect permissions. To avoid this, when making changes using either Cloudera Manager or the command line, first manually remove the existing NodeManager local directories from all configured local directories (`yarn.nodemanager.local-dirs`), and let the NodeManager recreate the directory structure.

Task-controller and Container-executor Error Codes

When you set up a secure cluster for the first time and debug problems with it, the `task-controller` or `container-executor` may encounter errors. These programs communicate these errors to the TaskTracker or NodeManager daemon via numeric error codes which will appear in the TaskTracker or NodeManager logs respectively (`/var/log/hadoop-mapreduce` or `/var/log/hadoop-yarn`). The following sections list the possible numeric error codes with descriptions of what they mean:

- [MRv1 ONLY: Task-controller Error Codes](#) on page 446
- [YARN ONLY: Container-executor Error Codes](#) on page 448

MRv1 ONLY: Task-controller Error Codes

The following table applies to the task-controller in MRv1.

Numeric Code	Name	Description
1	INVALID_ARGUMENT_NUMBER	<ul style="list-style-type: none"> • Incorrect number of arguments provided for the given task-controller command • Failure to initialize the job localizer
2	INVALID_USER_NAME	The user passed to the task-controller does not exist.
3	INVALID_COMMAND_PROVIDED	The task-controller does not recognize the command it was asked to execute.
4	SUPER_USER_NOT_ALLOWED_TO_RUN_TASKS	The user passed to the task-controller was the super user.

Numeric Code	Name	Description
5	INVALID_TT_ROOT	The passed TaskTracker root does not match the configured TaskTracker root (<code>mapred.local.dir</code>), or does not exist.
6	SETUID_OPER_FAILED	Either could not read the local groups database, or could not set UID or GID
7	UNABLE_TO_EXECUTE_TASK_SCRIPT	The task-controller could not execute the task launcher script.
8	UNABLE_TO_KILL_TASK	The task-controller could not kill the task it was passed.
9	INVALID_TASK_PID	The PID passed to the task-controller was negative or 0.
10	ERROR_RESOLVING_FILE_PATH	The task-controller could not resolve the path of the task launcher script file.
11	RELATIVE_PATH_COMPONENTS_IN_FILE_PATH	The path to the task launcher script file contains relative components (for example, <code>..</code>).
12	UNABLE_TO_STAT_FILE	The task-controller did not have permission to stat a file it needed to check the ownership of.
13	FILE_NOT_OWNED_BY_TASKTRACKER	A file which the task-controller must change the ownership of has the wrong the ownership.
14	PREPARE_ATTEMPT_DIRECTORIES_FAILED	The <code>mapred.local.dir</code> is not configured, could not be read by the task-controller, or could not have its ownership secured.
15	INITIALIZE_JOB_FAILED	The task-controller could not get, stat, or secure the job directory or job working working directory.
16	PREPARE_TASK_LOGS_FAILED	The task-controller could not find or could not change the ownership of the task log directory to the passed user.
17	INVALID_TT_LOG_DIR	The <code>hadoop.log.dir</code> is not configured.
18	OUT_OF_MEMORY	The task-controller could not determine the job directory path or the task launcher script path.
19	INITIALIZE_DISTCACHEFILE_FAILED	Could not get a unique value for, stat, or the local distributed cache directory.
20	INITIALIZE_USER_FAILED	Could not get, stat, or secure the per-user task tracker directory.
21	UNABLE_TO_BUILD_PATH	The task-controller could not concatenate two paths, most likely because it ran out of memory.
22	INVALID_TASKCONTROLLER_PERMISSIONS	The task-controller binary does not have the correct permissions set. See Information about Other Hadoop Security Programs .

Numeric Code	Name	Description
23	PREPARE_JOB_LOGS_FAILED	The task-controller could not find or could not change the ownership of the job log directory to the passed user.
24	INVALID_CONFIG_FILE	The taskcontroller.cfg file is missing, malformed, or has incorrect permissions.
255	Unknown Error	<p>There are several causes for this error. Some common causes are:</p> <ul style="list-style-type: none"> • There are user accounts on your cluster that have a user ID less than the value specified for the <code>min.user.id</code> property in the <code>taskcontroller.cfg</code> file. The default value is 1000 which is appropriate on Ubuntu systems, but may not be valid for your operating system. For information about setting <code>min.user.id</code> in the <code>taskcontroller.cfg</code> file, see this step. • Jobs won't run and the TaskTracker is unable to create a Hadoop logs directory. For more information, see (MRv1 Only) Jobs will not run and TaskTracker is unable to create a Hadoop logs directory on page 201. • This error is often caused by previous errors; look earlier in the log file for possible causes.

YARN ONLY: Container-executor Error Codes

The codes in the table apply to the container-executor in YARN, but are used by the LinuxContainerExecutor only.

Numeric Code	Name	Description
1	INVALID_ARGUMENT_NUMBER	<ul style="list-style-type: none"> • Incorrect number of arguments provided for the given task-controller command • Failure to initialize the container localizer
2	INVALID_USER_NAME	The user passed to the task-controller does not exist.
3	INVALID_COMMAND_PROVIDED	The container-executor does not recognize the command it was asked to run.
5	INVALID_NM_ROOT	The passed NodeManager root does not match the configured NodeManager root (<code>yarn.nodemanager.local-dirs</code>), or does not exist.
6	SETUID_OPER_FAILED	Either could not read the local groups database, or could not set UID or GID
7	UNABLE_TO_EXECUTE_CONTAINER_SCRIPT	The container-executor could not run the container launcher script.

Numeric Code	Name	Description
8	UNABLE_TO_SIGNAL_CONTAINER	The container-executor could not signal the container it was passed.
9	INVALID_CONTAINER_PID	The PID passed to the container-executor was negative or 0.
18	OUT_OF_MEMORY	The container-executor couldn't allocate enough memory while reading the container-executor.cfg file, or while getting the paths for the container launcher script or credentials files.
20	INITIALIZE_USER_FAILED	Couldn't get, stat, or secure the per-user NodeManager directory.
21	UNABLE_TO_BUILD_PATH	The container-executor couldn't concatenate two paths, most likely because it ran out of memory.
22	INVALID_CONTAINER_EXEC_PERMISSIONS	The container-executor binary does not have the correct permissions set. See Information about Other Hadoop Security Programs .
24	INVALID_CONFIG_FILE	The container-executor.cfg file is missing, malformed, or has incorrect permissions.
25	SETSID_OPER_FAILED	Could not set the session ID of the forked container.
26	WRITE_PIDFILE_FAILED	Failed to write the value of the PID of the launched container to the PID file of the container.
255	Unknown Error	<p>This error has several possible causes. Some common causes are:</p> <ul style="list-style-type: none"> • User accounts on your cluster have a user ID less than the value specified for the <code>min.user.id</code> property in the <code>container-executor.cfg</code> file. The default value is 1000 which is appropriate on Ubuntu systems, but may not be valid for your operating system. For information about setting <code>min.user.id</code> in the <code>container-executor.cfg</code> file, see this step. • This error is often caused by previous errors; look earlier in the log file for possible causes.

Sqoop, Pig, and Whirr Security Support Status

Here is a summary of the status of security in the other CDH 5 components:

- Sqoop 1 and Pig support security with no configuration required.
- Whirr does not support security in CDH 5.

Setting Up a Gateway Node to Restrict Cluster Access

Use the instructions that follow to set up and use a Hadoop cluster that is entirely firewalled off from outside access; the only exception will be one node which will act as a gateway. Client machines can access the cluster through the gateway using the REST API.

HttpFS will be used to allow REST access to HDFS, and Oozie will allow REST access for submitting and monitoring jobs.

Installing and Configuring the Firewall and Gateway

Follow these steps:

1. Choose a cluster node to be the gateway machine.
2. Install and configure the Oozie server by following the standard directions starting here: [Installing Oozie](#).
3. [Install HttpFS](#).
4. Start the Oozie server:

```
$ sudo service oozie start
```

5. Start the HttpFS server:

```
$ sudo service hadoop-httpfs start
```

6. Configure firewalls.

Block all access from outside the cluster.

- The gateway node should have ports 11000 (oozie) and 14000 (hadoop-httpfs) open.
- Optionally, to maintain access to the Web UIs for the cluster's JobTrackers, NameNode, and so on, open their HTTP ports: see [Ports Used by Components of CDH 5](#).

7. Optionally configure authentication in simple mode (default) or using Kerberos. See [HttpFS Authentication](#) on page 141 to configure Kerberos for HttpFS and [Oozie Authentication](#) on page 161 to configure Kerberos for Oozie.
8. Optionally encrypt communication using HTTPS for Oozie by following [these directions](#).

Accessing HDFS

With the Hadoop client:

All of the standard `hadoop fs` commands will work; just make sure to specify `-fs webhdfs://HOSTNAME:14000`. For example (where `GATEWAYHOST` is the hostname of the gateway machine):

```
$ hadoop fs -fs webhdfs://GATEWAYHOST:14000 -cat /user/me/myfile.txt
Hello World!
```

Without the Hadoop client:

You can run all of the standard `hadoop fs` commands by using the WebHDFS REST API and any program that can do GET, PUT, POST, and DELETE requests; for example:

```
$ curl "http://GATEWAYHOST:14000/webhdfs/v1/user/me/myfile.txt?op=OPEN&user.name=me"
Hello World!
```



Important: The `user.name` parameter is valid only if security is disabled. In a secure cluster, you must initiate a valid Kerberos session.

In general, the command will look like this:

```
$ curl "http://GATEWAYHOST/webhdfs/v1/PATH?[user.name=USER&]op=..."
```

You can find a full explanation of the commands in the [WebHDFS REST API documentation](#).

Submitting and Monitoring Jobs

The Oozie REST API supports the direct submission of jobs for MapReduce, Pig, and Hive; Oozie automatically creates a workflow with a single action. For any other action types, or to execute anything more complicated than a single job, you must create an actual workflow. Required files (JAR files, input data, and so on.) must already exist on HDFS; if they do not, you can use HttpFS to upload the files.

With the Oozie client:

All of the standard Oozie commands will work. You can find a full explanation of the commands in the documentation for the [command-line utilities](#).

Without the Oozie client:

You can run all of the standard Oozie commands by using the REST API and any program that can do GET, PUT, and POST requests. You can find a full explanation of the commands in the [Oozie Web Services API documentation](#).

Logging a Security Support Case

Before you log a support case, ensure you have either part or all of the following information to help Support investigate your case:

Kerberos Issues

- For Kerberos issues, your `krb5.conf` and `kdc.conf` files are valuable for support to be able to understand your configuration.
- If you are having trouble with client access to the cluster, provide the output for `klist -ef` after kiniting as the user account on the client host in question. Additionally, confirm that your ticket is renewable by running `kinit -R` after successfully kiniting.
- Specify if you are authenticating (kiniting) with a user outside of the Hadoop cluster's realm (such as Active Directory, or another MIT Kerberos realm).
- If using AES-256 encryption, ensure you have the [Unlimited Strength JCE Policy Files](#) deployed on all cluster and client nodes.

TLS/SSL Issues

- Specify whether you are using a private/commercial CA for your certificates, or if they are self-signed. Note that Cloudera strongly recommends against using self-signed certificates in production clusters.
- Clarify what services you are attempting to setup TLS/SSL for in your description.
- When troubleshooting TLS/SSL trust issues, provide the output of the following `openssl` command:

```
openssl s_client -connect host.fqdn.name:port
```

LDAP Issues

- Specify the LDAP service in use (Active Directory, OpenLDAP, one of Oracle Directory Server offerings, OpenDJ, etc)
- Provide a screenshot of the LDAP configuration screen you are working with if you are troubleshooting setup issues.
- Be prepared to troubleshoot using the `ldapsearch` command (requires the `openldap-clients` package) on the host where LDAP authentication or authorization issues are being seen.

Using Antivirus Software on CDH Hosts

If you use antivirus software on your servers, consider configuring it to skip scans on certain types of Hadoop-specific resources. It can take a long time to scan large files or directories with a large number of files. In addition, if your antivirus software locks files or directories as it scans them, those resources will be unavailable to your Hadoop processes during the scan, and can cause latency or unavailability of resources in your cluster. Consider skipping scans on the following types of resources:

- Scratch directories used by services such as Impala
- Log directories used by various Hadoop services
- Data directories which can grow to petabytes in size

The specific directory names and locations depend on the services your cluster uses and your configuration. In general, avoid scanning very large directories and filesystems. Instead, limit write access to these locations using security mechanisms such as access controls at the level of the operating system, HDFS, or at the service level.

Appendix: Apache License, Version 2.0

SPDX short identifier: Apache-2.0

Apache License
Version 2.0, January 2004
<http://www.apache.org/licenses/>

TERMS AND CONDITIONS FOR USE, REPRODUCTION, AND DISTRIBUTION

1. Definitions.

"License" shall mean the terms and conditions for use, reproduction, and distribution as defined by Sections 1 through 9 of this document.

"Licensor" shall mean the copyright owner or entity authorized by the copyright owner that is granting the License.

"Legal Entity" shall mean the union of the acting entity and all other entities that control, are controlled by, or are under common control with that entity. For the purposes of this definition, "control" means (i) the power, direct or indirect, to cause the direction or management of such entity, whether by contract or otherwise, or (ii) ownership of fifty percent (50%) or more of the outstanding shares, or (iii) beneficial ownership of such entity.

"You" (or "Your") shall mean an individual or Legal Entity exercising permissions granted by this License.

"Source" form shall mean the preferred form for making modifications, including but not limited to software source code, documentation source, and configuration files.

"Object" form shall mean any form resulting from mechanical transformation or translation of a Source form, including but not limited to compiled object code, generated documentation, and conversions to other media types.

"Work" shall mean the work of authorship, whether in Source or Object form, made available under the License, as indicated by a copyright notice that is included in or attached to the work (an example is provided in the Appendix below).

"Derivative Works" shall mean any work, whether in Source or Object form, that is based on (or derived from) the Work and for which the editorial revisions, annotations, elaborations, or other modifications represent, as a whole, an original work of authorship. For the purposes of this License, Derivative Works shall not include works that remain separable from, or merely link (or bind by name) to the interfaces of, the Work and Derivative Works thereof.

"Contribution" shall mean any work of authorship, including the original version of the Work and any modifications or additions to that Work or Derivative Works thereof, that is intentionally submitted to Licensor for inclusion in the Work by the copyright owner or by an individual or Legal Entity authorized to submit on behalf of the copyright owner. For the purposes of this definition, "submitted" means any form of electronic, verbal, or written communication sent to the Licensor or its representatives, including but not limited to communication on electronic mailing lists, source code control systems, and issue tracking systems that are managed by, or on behalf of, the Licensor for the purpose of discussing and improving the Work, but excluding communication that is conspicuously marked or otherwise designated in writing by the copyright owner as "Not a Contribution."

"Contributor" shall mean Licensor and any individual or Legal Entity on behalf of whom a Contribution has been received by Licensor and subsequently incorporated within the Work.

2. Grant of Copyright License.

Subject to the terms and conditions of this License, each Contributor hereby grants to You a perpetual, worldwide, non-exclusive, no-charge, royalty-free, irrevocable copyright license to reproduce, prepare Derivative Works of, publicly display, publicly perform, sublicense, and distribute the Work and such Derivative Works in Source or Object form.

3. Grant of Patent License.

Subject to the terms and conditions of this License, each Contributor hereby grants to You a perpetual, worldwide, non-exclusive, no-charge, royalty-free, irrevocable (except as stated in this section) patent license to make, have made, use, offer to sell, sell, import, and otherwise transfer the Work, where such license applies only to those patent claims

licensable by such Contributor that are necessarily infringed by their Contribution(s) alone or by combination of their Contribution(s) with the Work to which such Contribution(s) was submitted. If You institute patent litigation against any entity (including a cross-claim or counterclaim in a lawsuit) alleging that the Work or a Contribution incorporated within the Work constitutes direct or contributory patent infringement, then any patent licenses granted to You under this License for that Work shall terminate as of the date such litigation is filed.

4. Redistribution.

You may reproduce and distribute copies of the Work or Derivative Works thereof in any medium, with or without modifications, and in Source or Object form, provided that You meet the following conditions:

1. You must give any other recipients of the Work or Derivative Works a copy of this License; and
2. You must cause any modified files to carry prominent notices stating that You changed the files; and
3. You must retain, in the Source form of any Derivative Works that You distribute, all copyright, patent, trademark, and attribution notices from the Source form of the Work, excluding those notices that do not pertain to any part of the Derivative Works; and
4. If the Work includes a "NOTICE" text file as part of its distribution, then any Derivative Works that You distribute must include a readable copy of the attribution notices contained within such NOTICE file, excluding those notices that do not pertain to any part of the Derivative Works, in at least one of the following places: within a NOTICE text file distributed as part of the Derivative Works; within the Source form or documentation, if provided along with the Derivative Works; or, within a display generated by the Derivative Works, if and wherever such third-party notices normally appear. The contents of the NOTICE file are for informational purposes only and do not modify the License. You may add Your own attribution notices within Derivative Works that You distribute, alongside or as an addendum to the NOTICE text from the Work, provided that such additional attribution notices cannot be construed as modifying the License.

You may add Your own copyright statement to Your modifications and may provide additional or different license terms and conditions for use, reproduction, or distribution of Your modifications, or for any such Derivative Works as a whole, provided Your use, reproduction, and distribution of the Work otherwise complies with the conditions stated in this License.

5. Submission of Contributions.

Unless You explicitly state otherwise, any Contribution intentionally submitted for inclusion in the Work by You to the Licensor shall be under the terms and conditions of this License, without any additional terms or conditions. Notwithstanding the above, nothing herein shall supersede or modify the terms of any separate license agreement you may have executed with Licensor regarding such Contributions.

6. Trademarks.

This License does not grant permission to use the trade names, trademarks, service marks, or product names of the Licensor, except as required for reasonable and customary use in describing the origin of the Work and reproducing the content of the NOTICE file.

7. Disclaimer of Warranty.

Unless required by applicable law or agreed to in writing, Licensor provides the Work (and each Contributor provides its Contributions) on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied, including, without limitation, any warranties or conditions of TITLE, NON-INFRINGEMENT, MERCHANTABILITY, or FITNESS FOR A PARTICULAR PURPOSE. You are solely responsible for determining the appropriateness of using or redistributing the Work and assume any risks associated with Your exercise of permissions under this License.

8. Limitation of Liability.

In no event and under no legal theory, whether in tort (including negligence), contract, or otherwise, unless required by applicable law (such as deliberate and grossly negligent acts) or agreed to in writing, shall any Contributor be liable to You for damages, including any direct, indirect, special, incidental, or consequential damages of any character arising as a result of this License or out of the use or inability to use the Work (including but not limited to damages for loss of goodwill, work stoppage, computer failure or malfunction, or any and all other commercial damages or losses), even if such Contributor has been advised of the possibility of such damages.

9. Accepting Warranty or Additional Liability.

While redistributing the Work or Derivative Works thereof, You may choose to offer, and charge a fee for, acceptance of support, warranty, indemnity, or other liability obligations and/or rights consistent with this License. However, in accepting such obligations, You may act only on Your own behalf and on Your sole responsibility, not on behalf of any other Contributor, and only if You agree to indemnify, defend, and hold each Contributor harmless for any liability incurred by, or claims asserted against, such Contributor by reason of your accepting any such warranty or additional liability.

END OF TERMS AND CONDITIONS

APPENDIX: How to apply the Apache License to your work

To apply the Apache License to your work, attach the following boilerplate notice, with the fields enclosed by brackets "[]" replaced with your own identifying information. (Don't include the brackets!) The text should be enclosed in the appropriate comment syntax for the file format. We also recommend that a file or class name and description of purpose be included on the same "printed page" as the copyright notice for easier identification within third-party archives.

```
Copyright [yyyy] [name of copyright owner]

Licensed under the Apache License, Version 2.0 (the "License");
you may not use this file except in compliance with the License.
You may obtain a copy of the License at

    http://www.apache.org/licenses/LICENSE-2.0

Unless required by applicable law or agreed to in writing, software
distributed under the License is distributed on an "AS IS" BASIS,
WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
See the License for the specific language governing permissions and
limitations under the License.
```