

Ingesting data into Apache Hive in CDP Public Cloud

Date published: 2019-12-16

Date modified: 2024-12-10



Legal Notice

© Cloudera Inc. 2024. All rights reserved.

The documentation is and contains Cloudera proprietary information protected by copyright and other intellectual property rights. No license under copyright or any other intellectual property right is granted herein.

Unless otherwise noted, scripts and sample code are licensed under the Apache License, Version 2.0.

Copyright information for Cloudera software may be found within the documentation accompanying each component in a particular release.

Cloudera software includes software from various open source or other third party projects, and may be released under the Apache Software License 2.0 (“ASLv2”), the Affero General Public License version 3 (AGPLv3), or other license terms. Other software included may be released under the terms of alternative open source licenses. Please review the license and notice files accompanying the software for additional licensing information.

Please visit the Cloudera software product page for more information on Cloudera software. For more information on Cloudera support services, please visit either the Support or Sales page. Feel free to contact us directly to discuss your specific needs.

Cloudera reserves the right to change any products at any time, and without notice. Cloudera assumes no responsibility nor liability arising from the use of products, except as expressly agreed to in writing by Cloudera.

Cloudera, Cloudera Altus, HUE, Impala, Cloudera Impala, and other Cloudera marks are registered or unregistered trademarks in the United States and other countries. All other trademarks are the property of their respective owners.

Disclaimer: EXCEPT AS EXPRESSLY PROVIDED IN A WRITTEN AGREEMENT WITH CLOUDERA, CLOUDERA DOES NOT MAKE NOR GIVE ANY REPRESENTATION, WARRANTY, NOR COVENANT OF ANY KIND, WHETHER EXPRESS OR IMPLIED, IN CONNECTION WITH CLOUDERA TECHNOLOGY OR RELATED SUPPORT PROVIDED IN CONNECTION THEREWITH. CLOUDERA DOES NOT WARRANT THAT CLOUDERA PRODUCTS NOR SOFTWARE WILL OPERATE UNINTERRUPTED NOR THAT IT WILL BE FREE FROM DEFECTS NOR ERRORS, THAT IT WILL PROTECT YOUR DATA FROM LOSS, CORRUPTION NOR UNAVAILABILITY, NOR THAT IT WILL MEET ALL OF CUSTOMER’S BUSINESS REQUIREMENTS. WITHOUT LIMITING THE FOREGOING, AND TO THE MAXIMUM EXTENT PERMITTED BY APPLICABLE LAW, CLOUDERA EXPRESSLY DISCLAIMS ANY AND ALL IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO IMPLIED WARRANTIES OF MERCHANTABILITY, QUALITY, NON-INFRINGEMENT, TITLE, AND FITNESS FOR A PARTICULAR PURPOSE AND ANY REPRESENTATION, WARRANTY, OR COVENANT BASED ON COURSE OF DEALING OR USAGE IN TRADE.

Contents

Ingesting data into Hive.....	4
Understand the use case.....	4
Meet the prerequisites.....	4
Configure the service account.....	5
Create IDBroker mapping.....	6
Create the Hive target table.....	8
Add Ranger policies.....	8
Obtain Hive connection details.....	9
Build the data flow.....	10
Configure the controller services.....	12
Configure the processor for your data source.....	14
Configure the processor for your data target.....	16
Start your data flow.....	17
Verify your data flow.....	17
Next steps.....	17

Ingesting data into Hive

Understand the use case

You can use Apache NiFi to move data from a range of locations into a Data Engineering cluster running Apache Hive in CDP Public Cloud.

This use case walks you through the process of creating a data flow focused on ingesting data from Apache Kafka into Apache Hive. The PutHive3Streaming processor is used to ingest data into a Hive table. This processor interacts with the Hive Metastore to determine the storage location of the Hive table and writes data files directly to that location, eliminating the need for a running HiveServer2. After the delta files are written, the NiFi processor updates the table metadata in the Hive Metastore, ensuring that future queries include the newly ingested data.



Note: Although this process does not require a running HiveServer2 or a cluster for executing Hive queries, Cloudera strongly recommends having a Data Engineering cluster available in the same environment so that it can run Hive's compaction background processes which reconcile the delta files created by NiFi into base files. Without compaction, an excessive number of delta files can slow down Hive queries, as the engine has to reconcile all the delta files during query execution.

If you are moving data from a source other than Kafka, review the *Getting Started with Apache NiFi* for information on building data flows and exploring other data ingestion options. You can also find instructions on building a flow for ingesting data into other CDP components.

Related Information

[Getting Started with Apache NiFi](#)

Meet the prerequisites

Use this checklist to make sure that you meet all the requirements before you start building your data flow.

- You have a CDP Public Cloud environment.
- You have a workload username and password set to access Data Hub clusters. The predefined resource role of this user is at least "EnvironmentUser". This resource role provides the ability to view Data Hub clusters and set the FreeIPA password for the environment.
- Your user is synchronized to the CDP Public Cloud environment.
- You have a Flow Management cluster running Apache NiFi.
- You have a running Data Engineering cluster in the same CDP environment.



Note: Data can be ingested into Hive without an active Data Engineering cluster. However, for production workloads, it is strongly recommended to provision a Data Engineering cluster, which ensures that Hive's automated compaction jobs can run concurrently with NiFi data ingestion, optimizing query performance. For more information, see *Understand the use case*.

- You have created a new database and transactional table for Hive data ingest.
- You have a policy in Ranger that allows a user to write data to your Hive table.

Related Information

[Understand the use case](#)

[AWS Onboarding Quickstart](#)

[Azure Onboarding Quickstart](#)

[Understanding roles and resource roles](#)

[CDP workload user](#)

[Setting up your Flow Management cluster](#)

Configure the service account

Configure the Service Account you will use to ingest data into Hive.

Procedure

1. Navigate to the Cloudera Management Console UI.
2. From the User Management tab, create a new machine user.

3. Once the user is created, set a workload password and synchronize the changes to the required environments.

Users / nifi-hive-ingest

- From the Access tab, add the machine user as an Environment User to the environment.

Update Resource Roles for nifi-hive-ingest
×

Resource Roles

<input checked="" type="checkbox"/>	Role	Description
<input type="checkbox"/>	DWAdmin ⓘ	Grants permission to create, delete, and update Cloudera Data Warehouse clusters for a given CDP environment.
<input type="checkbox"/>	DWUser ⓘ	Grants permission to view Cloudera Data Warehouse cluster for a given CDP environment.
<input type="checkbox"/>	EnvironmentAdmin ⓘ	Grants all the rights to an environment.
<input checked="" type="checkbox"/>	EnvironmentUser ⓘ	Grants permission to set the workload password for the environment.
<input type="checkbox"/>	MLAdmin ⓘ	Grants permission to create and delete Cloudera Machine Learning workspaces for a given CDP environment. MLAdmins will also have Site Administrator level access to all the workspaces provisioned using this environment. That is, they can run workloads, monitor, and manage all user activity on these workspaces.
<input type="checkbox"/>	MLBusinessUser ⓘ	Grants permission to list Cloudera Machine Learning workspaces for a given CDP environment. MLBusinessUsers will also be able to view shared machine learning applications
<input type="checkbox"/>	MLUser ⓘ	Grants permission to list Cloudera Machine Learning workspaces for a given CDP environment. MLUsers will also be able to run workloads on all the workspaces provisioned using this environment.

Cancel
Update Roles

- Go back to the environments view and click Actions - Synchronize Users to Free IPA.

Results

Once the synchronization is complete, the required user is present in the environment and has the right permissions to access the object store.

Create IDBroker mapping

To enable your CDP user to utilize the central authentication features CDP provides and to exchange credentials for AWS or Azure access tokens, you have to map this CDP user to the correct IAM role or Azure Managed Service Identity (MSI). The option to add/modify these mappings is available from the Management Console in your CDP environment.

Procedure

- Go to the environment in which your Flow Management and Data Engineering clusters are running.
- To access IDBroker Mappings, click Actions | Manage Access and select the IDBroker Mapping tab in the next screen, where you can provide mappings for users or groups.
- Click Edit.

4. Add a new mapping for your service user, mapping the user to an existing IAM role or Azure Managed Identity Resource ID that has access to the underlying storage which is used by the target Hive table.

For example:

Access
IDBroker Mappings
Workload Password

Last Sync Status

Status
Completed

Sync Needed
No

> Detailed Information

Current Mappings

Data Access Role	arn:aws:iam::381358652250:role/mkohs-dladmin-role
Ranger Audit Role	arn:aws:iam::381358652250:role/mkohs-ranger-audit-role
User or Group	Role
nifi-hive-ingest	arn:aws:iam::381358652250:role/mkohs-customer-ingest-role

5. Add your CDP user and the corresponding AWS or Azure role that provides write access to your folder in your S3 bucket or ADLS folder to the Current Mappings section.



Note:

You can get the AWS IAM role Amazon Resource Name (ARN) from the Roles Summary page in AWS and can copy into the IDBroker role field.

You can get the Azure Managed Identity Resource ID by navigating to Managed Identities->Your Managed Identity->Properties->Resource ID.

6. Click Save and Sync.
7. Ensure that your IDBroker mapping change is synchronized to the environment successfully.

What to do next

Create a Hive table and add Ranger policies that allow your machine user write access to your Hive table.

Related Information

IDBroker

[Onboarding CDP users and groups for AWS cloud storage](#)

[Onboarding CDP users and groups for Azure cloud storage](#)

Create the Hive target table

Before you can ingest data into Apache Hive in CDP Public Cloud, ensure that you have a Hive target table. These steps walk you through creating a simple table. Modify these instructions based on your data ingest target table needs.

About this task

You require a transactional Hive table. This is the type of table created by default.

Procedure

1. From Cloudera Management Console, go to your Data Engineering cluster.
2. From Services, click Data Analytics Studio, and select Compose from the left-navigation pane.
3. Create a new target database.

```
create DATABASE retail
```

4. Create a new target table.

```
create table retail.customer(customer_id int, customer_name string)
```

What to do next

Add policies to Ranger so that the CDP User you are using to create this data flow can ingest data into the Hive target table.

Add Ranger policies

Add Ranger policies to ensure that you have write access to your Hive tables.

Procedure

1. From the Ranger cm_hive service in your Data Lake, create a new policy called Database access, and allow select permissions on the database in which your Hive target table is stored.

2. Create another policy to allow select, update, and Write permissions to ingest into the Hive target table.
For example:

Service Manager

cm_hive Policies

Edit Policy

Policy Details :

Policy Type

Access

Add Validity Period

Policy ID

44

Policy Name *

telco historical data table

enabled

normal

Policy Label

Policy Label

database

*

wwtelco

include

table

*

historical_data

include

column

*

*

include

Description

Audit Logging

YES

Allow Conditions :

hide

Select Role	Select Group	Select User	Permissions	Delegate Admin	
Select Roles	Select Groups	<div><div>nifi_flow_management_ingest</div><div>srv_nifi-ml-ingest</div></div>	<div><div>select</div><div>update</div><div>Write</div><div></div></div>	<input type="checkbox"/>	<div>x</div>

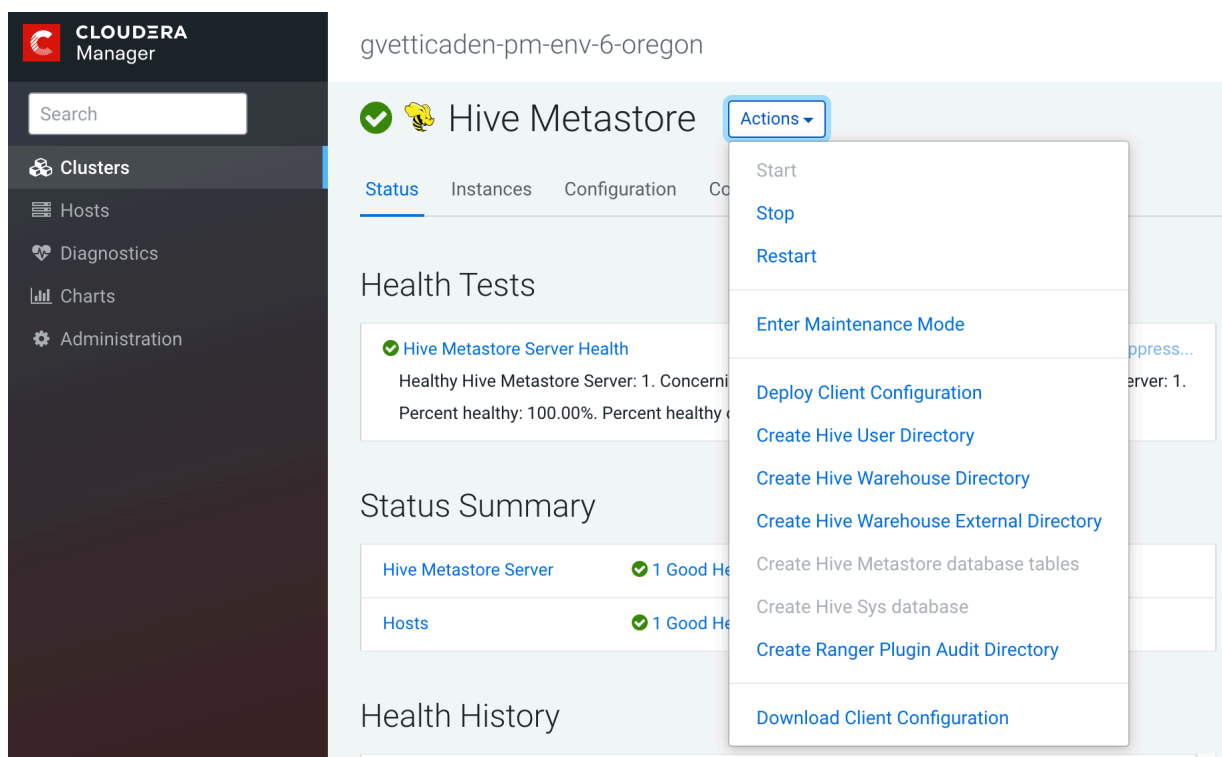
Obtain Hive connection details

To enable an Apache NiFi data flow to communicate with Hive, you must obtain the Hive connection details by downloading several client configuration files. The NiFi processors require these files to parse the configuration values and use those values to communicate with Hive.

Procedure

- From the Services pane in the Data Engineering cluster that contains the Hive instance to which you want to ingest data, click CM-UI.
- Click Hive Metastore, located on the Cloudera Manager Clusters pane.

- From the Hive Metastore view, click **Actions** **Download Client Configuration** to download the client configuration files as a zip file.



- Expand the zip file and upload the following files to each of the NiFi nodes in your Flow Management cluster:

- hdfs-site.xml
- hive-site.xml
- core-site.xml

For example:

```
scp hdfs-site.xml cdp_username@publiclp:/tmp
```

You may also copy these files to a more permanent location later.

- Adjust the file permissions to ensure that the NiFi processors can read them.

For example:

```
chmod 755 /tmp/*-site.xml
```

Build the data flow

From the Apache NiFi canvas, set up the elements of your data flow. This involves opening NiFi in CDP Public Cloud, adding processors to your NiFi canvas, and connecting the processors.

About this task

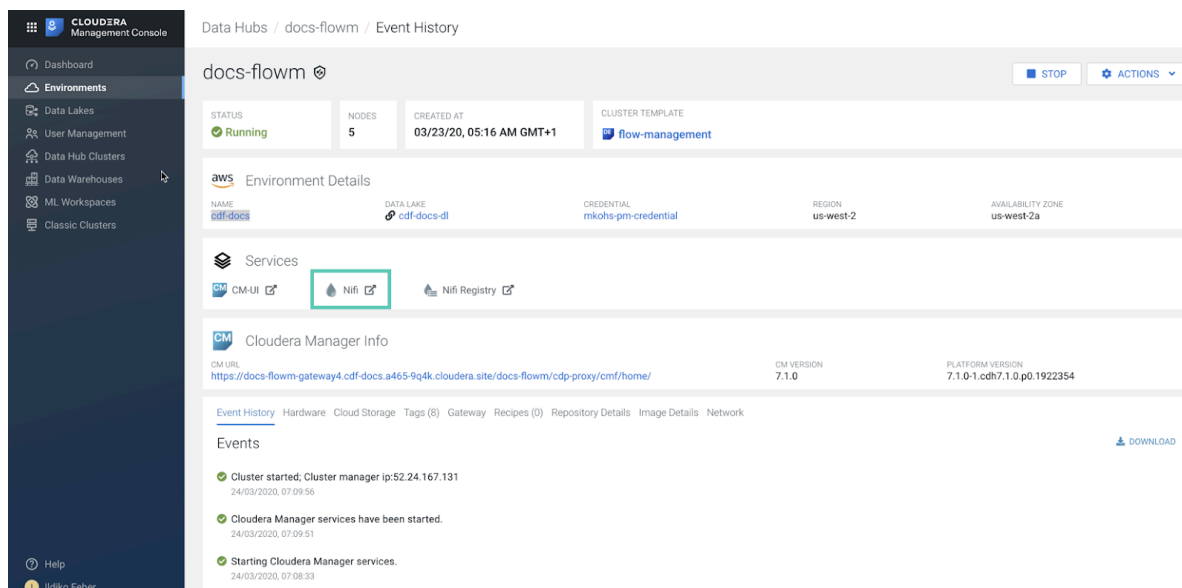
When you are building a data flow to ingest data into Apache Hive following are the preferred Hive processors

- Select_Hive_QL_3 – To get data out of Hive tables.
- PutHive3Streaming – To write data to managed and external Hive tables.

Procedure

1. Open NiFi in CDP Public Cloud.

- To access the NiFi service in your Flow Management cluster, navigate to Management Console service > Data Hub Clusters.
- Click the tile representing the Flow Management cluster with which you want to work.
- Click the NiFi icon in the Services section of the Cluster overview page to access the NiFi UI.



2. Add the NiFi Processors to your canvas.

- Select the Processor icon from the Cloudera Flow Management actions pane, and drag a processor to the Canvas.
- Use the Add Processor filter box to search for the processor you want to add, and then click Add.
- Add each of the processors you want to use for your data flow.

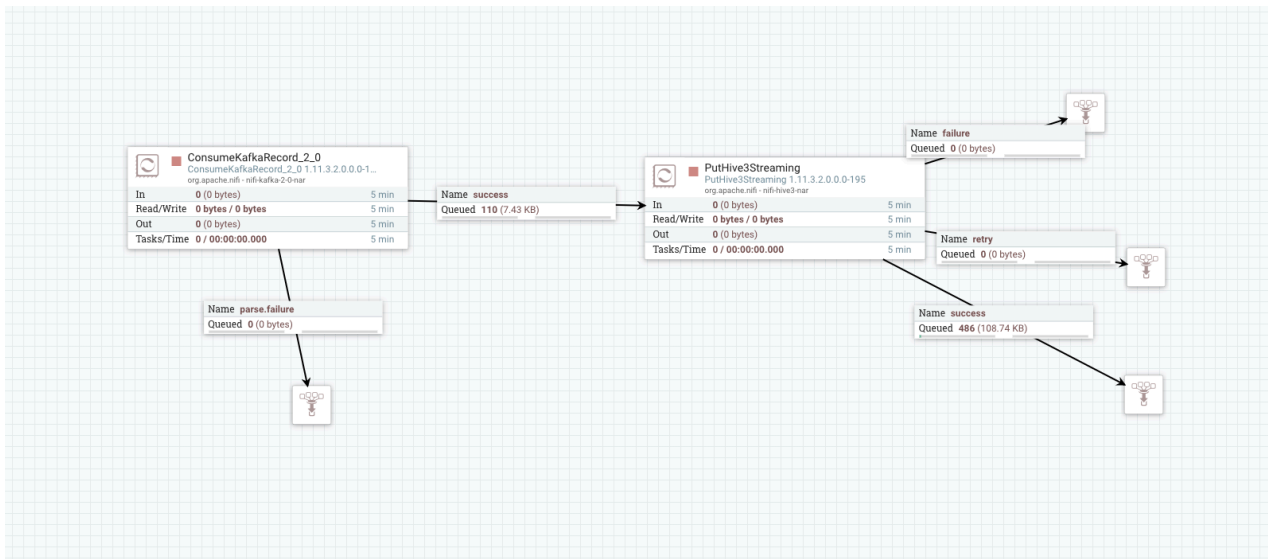
3. Connect the two processors to create a flow.

- Click the connection icon in the first processor, and drag it to the second processor.
- A Create Connection dialog displays. It has Details and Settings tabs. You can configure the connection's name, FlowFile expiration time period, thresholds for back pressure, load balance strategy, and prioritization.
- Click Add to close the dialog box and add the connection to your flow.

Optionally, you can add success and failure funnels to your data flow, which help you see where flow files are routed when your data flow is running.

Results

Your data flow may look similar to the following:



What to do next

Create the Controller service for your data flow. You will need these services later on as you configure your data flow target processor.

Related Information

[Ingesting Data into Apache Kafka in CDP Public Cloud](#)

[Consume KafkaRecord_2_0](#)

[Building a data flow](#)

Configure the controller services

You can add Controller Services to provide shared services to be used by the processors in your data flow. You will use these Controller Services later when you configure your processors.

About this task

You must define the Controller Services for the processors in your data flow in the configuration of the root process group where they will be used.

Procedure

1. To add a Controller Service to your flow, right-click on the Canvas and select Configure from the pop-up menu. This displays the Controller Services Configuration window.
2. Select the Controller Services tab.
3. Click the + button to display the Add Controller Service dialog.



4. Select the required Controller Service and click Add.

Add Controller Service

Source
all groups

Displaying 5 of 67

avro

Type	Version	Tags
AvroReader	1.9.2	comma, reader, record, values, ...
AvroRecordSetWriter	1.9.2	result, set, record, serializer, rec...
AvroSchemaRegistry	1.9.2	schema, registry, csv, json, avro
ConfluentSchemaRegistry	1.9.2	schema, registry, confluent, kaf...
HortonworksSchemaRegistry	1.9.2	schema, registry, hortonworks, ...

avro cache cluster
connection
couchbase csv
database dbcp
distributed enrich
hbase json key
lookup map parse
pooling reader
record recordset
row schema set
value writer

AvroReader 1.9.2 org.apache.nifi - nifi-record-serialization-services-nar

Parses Avro data and returns each Avro record as an separate Record object. The Avro data may contain the schema itself, or the schema can be externalized and accessed by one of the methods offered by the 'Schema Access Strategy' property.

CANCEL
ADD

5. Perform any necessary Controller Service configuration tasks by clicking the Configure icon in the right-hand column.

NIFI Flow Configuration

GENERAL
CONTROLLER SERVICES

Name	Type	Bundle	State	Scope
AvroReader	AvroReader 1.9.2	org.apache.nifi - nifi-record-serialization-services-...	Disabled	NiFi Flow

6. When you have finished configuring the options you need, save the changes by clicking the Apply button.
7. Enable the Controller Service by clicking the Enable button (flash) in the far-right column of the Controller Services tab.

Example

In this example, configure the following Controller Services:

- AvroReader
- AvroWriter
- HortonworksSchemaRegistry

AvroReader Controller Service

Table 1: AvroReader Controller Service properties

Property	Description	Example value for ingest data flow
Schema Access Strategy	Specify how to obtain the schema to be used for interpreting the data.	HWX Content-Encoded Schema Reference
Schema Registry	Specify the Controller Service to use for the Schema Registry.	CDPSchemaRegistry
Schema Name	Specify the name of the schema to look up in the Schema Registry property.	customer

WriteAvro Controller Service

Table 2: WriteAvro Controller Service properties

Property	Description	Example value for ingest data flow
Schema Write Strategy	Specifies how the schema for a Record should be added to the data.	HWX Content Encoded Schema Reference
Schema Access Strategy	Specifies how to obtain the schema that is to be used for interpreting the data.	Use Schema Name Property
Schema Registry	Specifies the Controller Service to use for the Schema Registry.	HortonworksSchemaRegistry
Schema name	Specifies the name of the schema to lookup in the Schema Registry property.	<i>customer</i>

HortonworksSchema Registry

Table 3: HortonworksSchemaRegistry Controller Service properties

Property	Description	Example value for ingest data flow
Schema Registry URL	Provide the URL of the schema registry that this Controller Service should connect to, including version. In the format: https://host:7790/api/v1	https://docs-messaging-master0.cdf-docs.a465-9q4k.cloudera.site:7790/api/v1
SSL Context Service	Specify the SSL Context Service to use for communicating with Schema Registry. Use the pre-configured SSLContextProvider.	Default NiFi SSL Context Service
Kerberos Principal	Specify the user name that should be used for authenticating with Kerberos. Use your CDP workload username to set this Authentication property.	srv_nifi-hive-ingest
Kerberos Password	Provide the password that should be used for authenticating with Kerberos. Use your CDP workload password to set this Authentication property.	password

Configure the processor for your data source

You can set up a data flow to move data from many locations into Apache Hive. This example assumes that you are configuring ConsumeKafkaRecord_2_0. If you are moving data from a location other than Kafka, review

Getting Started with Apache NiFi for information about how to build a data flow, and about other data consumption processor options.

Before you begin

- Ensure that you have the Ranger policies required to access the Kafka consumer group.
- This use case demonstrates a data flow using data in Kafka, and moving it to Hive. To review how to ingest data to Kafka, see *Ingesting Data to Apache Kafka*.

Procedure

1. Launch the Configure Processor window, by right-clicking the processor and selecting Configure.
2. Click the Properties tab.
3. Configure ConsumeKafkaRecord_2_0 with the required values.

The following table includes a description and example values for the properties required to configure an ingest data flow. For a complete list of ConsumeKafkaRecord_2_0, see the *processor documentation*.

Property	Description	Value for ingest to Hive data flow
Kafka Brokers	Provide a comma-separated list of known Kafka Brokers in the format <host>:<port>.	<code>Docs-messaging-broker1.cdf-docs.a465-9q4k.cloudera.site:9093, docs-messaging-broker2.cdf-docs.a465-9q4k.cloudera.site:9093, docs-messaging-broker3.cdf-docs.a465-9q4k.cloudera.site:9093</code>
Topic Name(s)	Enter the name of the Kafka topic from which you want to get data.	<code>customer</code>
Topic Name Format	Specify whether the topics are a comma separated list of topic names, or a single regular expression.	<code>names</code>
Record Reader	Specifies the record reader to use for incoming Flow Files. This can be a range of data formats.	<code>AvroReader</code>
Record Writer	Specifies the format you want to use in order to serialize data before for sending it to the data target. Note: Ensure that the source record writer and the target record reader are using the same Schema.	<code>CSVRecordSetWriter</code>
Honor Transactions	Specifies whether or not NiFi should honor transactional guarantees when communicating with Kafka.	<code>false</code>
Security Protocol	Specifies the protocol used to communicate with brokers. This value corresponds to Kafka's 'security.protocol' property.	<code>SASL_SSL</code>
SASL Mechanism	The SASL mechanism to use for authentication. Corresponds to Kafka's 'sasl.mechanism' property.	<code>plain</code>
Username	Specify the username when the SASL Mechanism is PLAIN or SCRAM-SHA-256.	<code>srv_nifi-hive-ingest</code>

Property	Description	Value for ingest to Hive data flow
Password	Specify the password for the given username when the SASL Mechanism is PLAIN or SCRAM-SHA-256.	<i>Password1!</i>
SSL Context Service	Specifies the SSL Context Service to use for communicating with Kafka.	<i>default</i>
Offset Reset	Allows you to manage the condition when there is no initial offset in Kafka or if the current offset does not exist any more on the server (e.g. because that data has been deleted). Corresponds to Kafka's <code>auto.offset.reset</code> property.	<i>latest</i>
Group ID	A Group ID is used to identify consumers that are within the same consumer group. Corresponds to Kafka's 'group.id' property.	<i>nifi-hive-ingest</i>

What to do next

Configure the processor for your data target.

Related Information

[ConsumeKafka_2_0](#)

[Building a data flow](#)

[Ingesting data into Apache Kafka in CDP Public Cloud](#)

Configure the processor for your data target

You can set up a data flow to move data into many locations. This example assumes that you are moving data into Apache Hive using `PutHive3Streaming`. If you are moving data into another location, review *Getting Started with Apache NiFi* for information about how to build a data flow, and about other data ingest processor options.

Procedure

1. Launch the Configure Processor window, by right-clicking the processor and selecting Configure.
2. Click the Properties tab.
3. Configure `PutHive3Streaming` with the required values.

The following table includes a description and example values for the properties required to configure an ingest data flow into Hive. For a complete list of `PutHive3Streaming`, see the *processor documentation*.

Property	Description	Value for ingest to Hive data flow
Database name	The name of the database in which to put the data.	<i>retail</i>
Table name	The name of the database table in which to put the data.	<i>customer</i>
Record Reader		<i>ReadAvro</i>
Hive Metastore URI	The URI location for the Hive Metastore. Note that this is not the location of the Hive Server. The default port for the Hive metastore is 9083.	<i>thift://<datalake-master-host-name>:9083</i>

Property	Description	Value for ingest to Hive data flow
Hive Configuration Resources	A file or comma separated list of files which contains the Hive configuration.	Replace with the file paths for the three files you downloaded in <i>Obtain the Hive connection details</i> : <ul style="list-style-type: none"> hdfs-site.xml hive-site.xml core-site.xml
Kerberos Principal	Specify your CDP machine user name.	<i>srv_nifi-hive-ingest</i>
Kerberos Password	Specify your CDP machine user password.	<i>Password1!</i>

Start your data flow

Start your data flow to verify that you have created a working dataflow and to begin your data ingest process.

Procedure

1. To initiate your data flow, select all the flow components you want to start.
2. Click the Start icon in the Actions toolbar.

Alternately, right-click a single component and choose Start from the context menu.

Verify your data flow

Learn how you can verify the operation of your Hive ingest data flow.

About this task

There are a number of ways to check that data is running through the flow you have built and it actually appears in Hive.

Procedure

1. Review the Hive ingest data flow on the NiFi canvas and confirm that the NiFi processors are not producing errors.
2. You can look at the processors in the UI, where you can see the amount of data that has gone through them. You can also right click on the processors, or on connections to view status history.
3. From the Data Analytics Studio UI, in the Compose tab, confirm that your Hive table is consuming content:

```
select * from retail.customer
```

Results

You are able to see data flowing into your retail.customer table.

Next steps

Provides information on what to do once you have moved data into Hive in CDP Public Cloud.

You have built a simple data flow for an easy way to move data to Hive. This example data flow enables you to easily design more complex data flows for moving and processing data in Hive. Here are some ideas for next steps:

- Review *NiFi documentation* to learning more about building and managing data flows.
- Review *NiFi Registry documentation* for more information about versioning data flows.

- Review the *Cloudera Runtime data access documentation* for information about working with Hive, Impala, and Hue.

Related Information[NiFi documentation](#)[NiFi Registry documentation](#)[Cloudera Runtime data access information](#)