## Cloudera Runtime 1.5.5

# **Apache Impala**

Date published: 2020-11-30 Date modified: 2025-06-06



# **Legal Notice**

© Cloudera Inc. 2025. 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**

Apache Impala Overview	4
Components of Impala	4

# **Apache Impala Overview**

The Apache Impala provides high-performance, low-latency SQL queries on data stored in popular Apache Hadoop file formats.

The Impala solution is composed of the following components.

#### Impala

The Impala service coordinates and executes queries received from clients. Queries are distributed among Impala nodes, and these nodes then act as workers, executing parallel query fragments.

#### **Hive Metastore**

Stores information about the data available to Impala. For example, the metastore lets Impala know what databases are available and what the structure of those databases is. As you create, drop, and alter schema objects, load data into tables, and so on through Impala SQL statements, the relevant metadata changes are automatically broadcast to all Impala nodes by the dedicated catalog service.

#### Clients

Entities including Hue, ODBC clients, JDBC clients, Business Intelligence applications, and the Impala Shell can all interact with Impala. These interfaces are typically used to issue queries or complete administrative tasks such as connecting to Impala.

#### Storage for data to be queried

Queries executed using Impala are handled as follows:

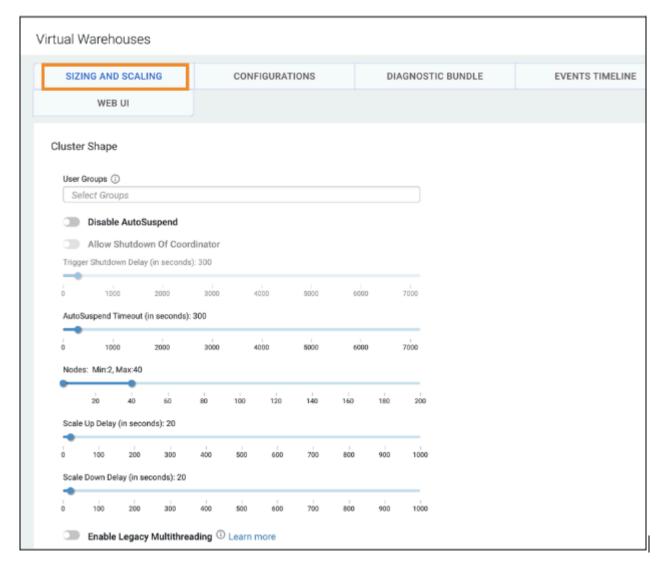
- User applications send SQL queries to Impala through ODBC or JDBC, which provide standardized querying
  interfaces. The user application may connect to any impalad in the cluster. This impalad becomes the coordinator
  for the query.
- **2.** Impala parses the query and analyzes it to determine what tasks need to be performed by impalad instances across the cluster. Execution is planned for optimal efficiency.
- 3. Storage services are accessed by local impalad instances to provide data.
- 4. Each impalad returns data to the coordinating impalad, which sends these results to the client.

# **Components of Impala**

The Impala service is a distributed, massively parallel processing (MPP) database engine. It consists of different daemon processes called as components. This topic describes the different roles these components play for a selected Virtual Warehouse.

When you create an Impala Virtual Warehouse, it is automatically optimized for your workload by the Cloudera Data Warehouse service. Due to the containerized and compute-isolated architecture of Cloudera Data Warehouse, as well as intelligence to assign different default configurations, you do not have to customize your environment to optimize performance or to avoid resource usage spikes and out-of-memory conditions. However, if you must adjust some settings to suit your needs after creating an Impala Virtual Warehouse, you can add particular component configuration details in one of the free-form fields on the SIZING AND SCALING tab or under one of the components available under the CONFIGURATIONS tab for a selected Virtual Warehouse.

Cloudera Runtime Components of Impala

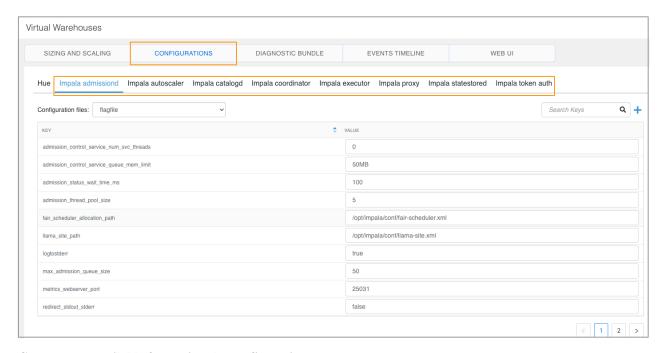


Before making any configuration changes, you can review the default configuration details in the components available under the CONFIGURATIONS tab in a Virtual Warehouse.

Impala service in Cloudera Data Warehouse consists of the following different daemon processes.

- Impala admissiond
- Impala catalogd
- Impala coordinator
- Impala executor
- · Impala statestored
- · Impala autoscaler
- Impala proxy

Cloudera Runtime Components of Impala



### Components available for setting the configuration

This section provides reference information on supported configuration properties under the listed components. For the full list of configuration properties, see Impala Properties in Cloudera Runtime. Based on the information provided under the components, choose the appropriate component/role to tune a configuration if you must.



**Note:** Some of the configuration properties under these components can be tuned to your needs and some of the properties are not editable. In addition, the available properties may differ by Cloudera Data Warehouse Impala Runtime version.

## Impala catalogd

The Catalog Server relays the metadata changes from Impala SQL statements to all the Impala daemons in a cluster. The catalog service avoids the need to issue REFRESH and INVALIDATE METADATA statements when the metadata changes are performed by statements issued through Impala.

### Impala coordinator

A few of the key functions that an Impala coordinator performs are:

- · Reads and writes to data files.
- Accepts queries transmitted from the impala-shell command, Hue, JDBC, or ODBC.
- Parallelizes the queries and distributes work across the cluster.

It is in constant communication with StateStore, to confirm which executors are healthy and can accept new work. Based on the health information it receives it assigns tasks to the executors. It also receives broadcast messages from the Catalog Server daemon whenever a cluster creates, alters, or drops any type of object, or when an INSERT or LOAD DATA statement is processed through Impala. It also communicates to Catalog Server daemon to load table metadata.

#### Impala executor

A few of the key functions that an Impala executor performs are:

- Executes the queries and transmits query results back to the central coordinator.
- Also transmits intermediate query results.

Depending on the size and the complexity of your queries you can select the number of executor nodes that are needed to run a typical query in your workloads. For more information on the

Cloudera Runtime Components of Impala

executor groups and recommendations on sizing your virtual warehouse to handle queries that must be run in your workloads concurrently, see Impala auto-scaling on public clouds.

#### Impala statestored

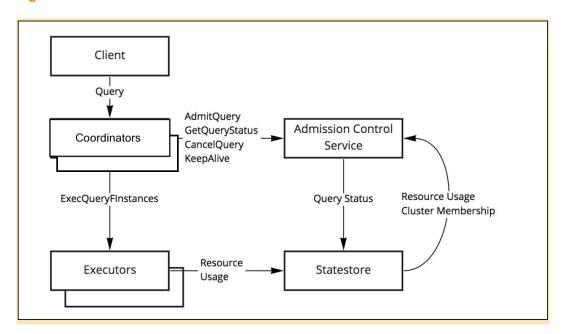
The Impala StateStore checks on the health of all Impala daemons in a cluster, and continuously relays its findings to each of those daemons. If an Impala daemon goes offline due to hardware failure, network error, software issue, or other reason, the StateStore informs all the other Impala daemons so that future queries can avoid making requests to the unreachable Impala daemon.

## Impala admissiond

Before Cloudera Data Warehouse on premises 1.5.5 SP1 release, the admission controller is part of Impala coordinators. Each coordinator runs its local admission controller that uses eventually consistent information about the decisions made by other coordinators. In a multiple-coordinator setup, each coordinator with its local admission controller performs poorly because of the consistent nature of the admission decisions. Also, the traditional Impala multiple coordinator setups result in over admission of cluster resources.

To mitigate the above issue, Cloudera Data Warehouse on premises 1.5.5 SP1 a new service is added for admission control that runs in a separate process. This separation decouples its failure modes from coordinators and executors. However, the functionality is similar to the local admission controller implementation by scheduling queries for all executor groups and then attempting admission using round-robin algorithm. Separating the admission controller from the coordinator and running a single admission controller per cluster allows each coordinator to contact this service for each query to receive an admission decision.

Figure 1: Global admission controller



To impose limits on concurrent SQL queries to avoid resource usage spikes and out-of-memory conditions on busy Cloudera clusters the following configuration flags are set by default:

- The admission\_control\_service\_queue\_mem\_limit flag specifies the limit on RPC payload consumption for the admission control service. The value is specified as number of bytes <int >[bB]?, megabytes <float>[mM], gigabytes <float>[gG], or percentage of the process memory limit <int>%. If no unit is specified, the default unit is bytes.
- The admission\_control\_service\_num\_svc\_threads flag specifies the number of threads for processing the admission control service RPCs. If the default value, 0 is used, the value is set to the number of CPU cores.

- The admission\_thread\_pool\_size flag specifies the size of the thread pool processing AdmitQuery requests.
- The max\_admission\_queue\_size flag specifies the maximum size of the queue for the AdmitQuery thread pool.
- The admission\_status\_wait\_time\_ms flag specifies the time in milliseconds that the GetQueryStatus() RPC in the admission control service waits for the admission to complete before returning.

Adjusting the admission control configuration flags is possible.

#### Components available for diagnostic purpose only

The following components are read-only, and available only to view for diagnostic purposes. The properties listed under them are not tunable.

#### Impala autoscaler

One of the core Impala Virtual Warehouse components is Impala autoscaler. Impala autoscaler is in constant communication with coordinators and executors to determine when more or fewer compute resources are needed, given the incoming query load. When the autoscaler detects an imbalance in resources, it sends a request to the Kubernetes framework to increase or decrease the number of executor groups in the Virtual Warehouse, thereby right-sizing the amount of resources available for queries. This ensures that workload demand is met without wasting cloud resources.

#### Impala proxy

Impala Proxy is a small footprint reverse proxy that will forward every http client request to the Coordinator endpoint.

When you create an Impala warehouse with the 'allow coordinator shutdown' option a proxy 'impala proxy' is created to act as a connection endpoint to the impala clients. This option when enabled allows Impala coordinators to automatically shut down during idle periods. If the coordinator has shut down because of idle period, and if there is a request from a client when the coordinator is not running, impala proxy triggers the coordinator to start. So the proxy acts as a layer between the clients and the coordinator so that the clients connected to the VW do not time out when the coordinator starts up.