

Cloudera Runtime 7.1.8

Starting and Stopping Apache Impala

Date published: 2020-11-30

Date modified: 2022-08-25

CLOUDERA

<https://docs.cloudera.com/>

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

Modifying Impala Startup Options.....	4
Configuring Client Access to Impala.....	4
Impala Startup Options for Client Connections.....	5
Impala Shell Tool.....	6
Impala Shell Configuration Options.....	6
Impala Shell Configuration File.....	9
Connecting to Impala Daemon in Impala Shell.....	10
Running Commands and SQL Statements in Impala Shell.....	12
Impala Shell Command Reference.....	13
Connecting to a kerberized Impala daemon.....	14
Configuring ODBC for Impala.....	14
Configuring JDBC for Impala.....	15
Configuring Impyla for Impala.....	16
Connecting to DataHub.....	18
Connecting to DataHub Data Mart.....	19
Configuring Delegation for Clients.....	19
Spooling Impala Query Results.....	20
Shut Down Impala.....	21
Setting Timeouts in Impala.....	22
Setting Timeout and Retries for Thrift Connections to Backend Client.....	22
Increasing StateStore Timeout.....	23
Setting the Idle Query and Idle Session Timeouts.....	23
Adjusting Heartbeat TCP Timeout Interval.....	24

Modifying Impala Startup Options

You can view and edit the configuration options for the Impala daemons to customize your Impala environment, such as to specify which hosts and ports to use, to assign directories for logging, and to control resource usage and security.

Configuring Impala Startup Options

Navigate to the following page to configure the settings for all the Impala-related daemons:

ClustersImpalaConfiguration.

If the Cloudera Manager interface does not yet have a form field for a newly added option, or if you need to use special options for debugging and troubleshooting, the Advanced category page for each daemon includes one or more Safety Valve fields where you can enter option names directly.

Checking the Values of Impala Startup Options

You can check the current runtime value of all these settings through the Impala Web UI, available by default at:

- http://impala_hostname:25000/varz for the `impalad` daemon
- http://impala_hostname:25010/varz for the `statestored` daemon
- http://impala_hostname:25020/varz for the `catalogd` daemon

Related Information

["Web User Interface for Debugging"](#)

Configuring Client Access to Impala

Application developers have a number of options to interface with Impala.

The core development language with Impala is SQL. You can also use Java or other languages to interact with Impala through the standard JDBC and ODBC interfaces used by many business intelligence tools. For specialized kinds of analysis, you can supplement the Impala built-in functions by writing user-defined functions in C++ or Java.

You can connect and submit requests to the Impala through:

- The `impala-shell` interactive command interpreter
- The Hue web-based user interface
- JDBC
- ODBC
- Impyla

Impala clients can connect to any Coordinator Impala Daemon (`impalad`) via HiveServer2 over HTTP or over the TCP binary or via Beeswax. All interfaces support Kerberos and LDAP for authentication to Impala. See below for the default ports and the Impala configuration field names to change the ports in Cloudera Manager.

Protocol	Default Port	Cloudera Manager Field to Specify an Alternate Port
HiveServer2 HTTP	28000	Impala Daemon HiveServer2 HTTP Port
HiveServer2 binary TCP	21050	Impala Daemon HiveServer2 Port
Beeswax	21000	Impala Daemon Beeswax Port

Impala Startup Options for Client Connections

Use the following flags to control client connections to Impala when starting Impala Daemon coordinator. If a configuration field exists in Cloudera Manager, the field name is shown in parenthesis next to the flag name.

If the Cloudera Manager interface does not yet have a form field for an option, the Advanced category page for each daemon includes one or more Safety Valve fields where you can enter option names directly.

--accepted_client_cnxn_timeout

Controls how Impala treats new connection requests if it has run out of the number of threads configured by `--fe_service_threads`.

If `--accepted_client_cnxn_timeout > 0`, new connection requests are rejected if Impala can't get a server thread within the specified (in seconds) timeout.

If `--accepted_client_cnxn_timeout=0`, i.e. no timeout, clients wait indefinitely to open the new session until more threads are available.

The default timeout is 5 minutes.

The timeout applies only to client facing thrift servers, i.e., HS2 and Beeswax servers.

--disconnected_session_timeout

When a HiveServer2 session has had no open connections for longer than this value, the session will be closed, and any associated queries will be unregistered.

Specify the value in hours.

The default value is 1 hour.

This flag does not apply to Beeswax clients. When a Beeswax client connection is closed, Impala closes the session associated with that connection.

--fe_service_threads (Impala Daemon Max Client)

Specifies the maximum number of concurrent client connections allowed. The default value is 64 with which 64 queries can run simultaneously.

If you have more clients trying to connect to Impala than the value of this setting, the later arriving clients have to wait for the duration specified by `--accepted_client_cnxn_timeout`. You can increase this value to allow more client connections. However, a large value means more threads to be maintained even if most of the connections are idle, and it could negatively impact query latency. Client applications should use the connection pool to avoid need for large number of sessions.

--idle_client_poll_time_s

The value of this setting specifies how frequently Impala polls to check if a client connection is idle and closes it if the connection is idle. A client connection is idle if all sessions associated with the client connection are idle.

By default, `--idle_client_poll_time_s` is set to 30 seconds.

If `--idle_client_poll_time_s` is set to 0, idle client connections stay open until explicitly closed by the clients.

The connection will only be closed if all the associated sessions are idle or closed. Sessions cannot be idle unless either the flag `--idle_session_timeout` or the `IDLE_SESSION_TIMEOUT` query option is set to greater than 0. If idle session timeout is not configured, a session cannot become idle by definition, and therefore its connection stays open until the client explicitly closes it.

--max_cookie_lifetime_s

Impala uses cookies for authentication when clients connect via HiveServer2 over HTTP. Use the `--max_cookie_lifetime_s` startup flag to control how long generated cookies are valid for.

Specify the value in seconds.

The default value is 1 day.

Setting the flag to 0 disables cookie support.

When an unexpired cookie is successfully verified, the user name contained in the cookie is set on the connection.

Each impalad uses its own key to generate the signature, so clients that reconnect to a different impalad have to re-authenticate.

On a single impalad, cookies are valid across sessions and connections.

--beeswax_port (Impala Daemon Beeswax Port)

Specifies the port for clients to connect to Impala daemon via the Beeswax protocol.

You can disable the Beeswax end point for clients by setting the flag to 0.

--hs2_http_port (Impala Daemon HiveServer2 HTTP Port)

Specifies the port for clients to connect to Impala daemon over HTTP.

You can disable the HTTP end point for clients by setting the flag to 0.

To enable TLS/SSL for HiveServer2 HTTP endpoint, use `--ssl_server_certificate` and `--ssl_private_key`.

--hs2_port (Impala Daemon HiveServer2 Port)

Specifies the port for clients to connect to Impala daemon via the HiveServer2 protocol.

You can disable the binary HiveServer2 end point for clients by setting the flag to 0.

Impala Shell Tool

You can use the Impala shell tool (`impala-shell`) to set up databases and tables, insert data, and issue queries.

For ad-hoc queries and exploration, you can submit SQL statements in an interactive session. To automate your work, you can specify command-line options to process a single statement or a script file. The `impala-shell` accepts all the same SQL statements, plus some shell-only commands that you can use for tuning performance and diagnosing problems.

Cloudera Manager installs `impala-shell` automatically. You might install `impala-shell` manually on other systems not managed by Cloudera Manager, so that you can issue queries from client systems that are not also running the Impala daemon or other Apache Hadoop components.

Impala Shell Configuration Options

You can specify the following options when starting `impala-shell` to control how shell commands are executed. You can specify options on the command line or in the `impala-shell` configuration file.

Command-Line Option	Configuration File Setting	Explanation
<code>-B</code> or <code>--delimited</code>	<code>write_delimited=true</code>	Causes all query results to be printed in plain format as a delimited text file. Useful for producing data files to be used with other Hadoop components. Also useful for avoiding the performance overhead of pretty-printing all output, especially when running benchmark tests using queries returning large result sets. Specify the delimiter character with the <code>--output_delimiter</code> option. Store all query results in a file rather than printing to the screen with the <code>-B</code> option.

Command-Line Option	Configuration File Setting	Explanation
-E or --vertical	vertical=true	Causes all query results to print in vertical format. In this mode, <code>impala-shell</code> will print each row in the following style. <ul style="list-style-type: none"> The first line will contain the line number followed by the row's columns in separate lines. Each column line will have a prefix with its name and a colon. To enable this mode, use the shell option '-E' or '--vertical', or 'set VERTICAL= true' in the interactive mode. To disable it in interactive mode, 'set VERTICAL=false'. NOTE: This vertical option will be disabled if the '-B' option or 'set WRITE_DELIMITED=true' is specified.
--live_progress	live_progress=true	Prints a progress bar showing roughly the percentage complete for each query. Information is updated interactively as the query progresses.
--disable_live_progress	live_progress=false	Disables <code>live_progress</code> in the interactive mode.
-b or --kerberos_host_fqdn	kerberos_host_fqdn= <i>load-balancer-hostname</i>	If set, the setting overrides the expected hostname of the Impala daemon's Kerberos service principal. <code>impala-shell</code> will check that the server's principal matches this hostname. This may be used when <code>impalad</code> is configured to be accessed via a load-balancer, but it is desired for <code>impala-shell</code> to talk to a specific <code>impalad</code> directly.
--print_header	print_header=true	
-o <i>filename</i> or --output_file <i>filename</i>	output_file= <i>filename</i>	Stores all query results in the specified file. Typically used to store the results of a single query issued from the command line with the <code>-q</code> option. Also works for interactive sessions; you see the messages such as number of rows fetched, but not the actual result set. To suppress these incidental messages when combining the <code>-q</code> and <code>-o</code> options, redirect <code>stderr</code> to <code>/dev/null</code> .
--output_delimiter= <i>character</i>	output_delimiter= <i>character</i>	Specifies the character to use as a delimiter between fields when query results are printed in plain format by the <code>-B</code> option. Defaults to tab (<code>'\t'</code>). If an output value contains the delimiter character, that field is quoted, escaped by doubling quotation marks, or both.
-p or --show_profiles	show_profiles=true	Displays the query execution plan (same output as the <code>EXPLAIN</code> statement) and a more detailed low-level breakdown of execution steps, for every query executed by the shell.
--profile_format= text json prettyjson	N/A	json and prettyjson output the JSON representation of each profile in a dense single-line form and in a human-readable multi-line form respectively.
-h or --help	N/A	Displays help information.
N/A	history_max=1000	Sets the maximum number of queries to store in the history file.
-i <i>hostname</i> or --impalad= <i>hostname[:portnum]</i>	impalad= <i>hostname[:portnum]</i>	Connects to the <code>impalad</code> daemon on the specified host. The default port of 21050 is assumed unless you provide another value. You can connect to any host in your cluster that is running <code>impalad</code> . If you connect to an instance of <code>impalad</code> that was started with an alternate port specified by the <code>--fe_port</code> flag, provide that alternative port.
-q <i>query</i> or --query= <i>query</i>	query= <i>query</i>	Passes a query or other <code>impala-shell</code> command from the command line. The <code>impala-shell</code> interpreter immediately exits after processing the statement. It is limited to a single statement, which could be a <code>SELECT</code> , <code>CREATE TABLE</code> , <code>SHOW TABLES</code> , or any other statement recognized in <code>impala-shell</code> . Because you cannot pass a <code>USE</code> statement and another query, fully qualify the names for any tables outside the default database. (Or use the <code>-f</code> option to pass a file with a <code>USE</code> statement followed by other queries.)

Command-Line Option	Configuration File Setting	Explanation
<code>-f query_file</code> or <code>--query_file=query_file</code>	<code>query_file=path_to_query_file</code>	Passes a SQL query from a file. Multiple statements must be semicolon (;) delimited.
<code>-k</code> or <code>--kerberos</code>	<code>use_kerberos=true</code>	Kerberos authentication is used when the shell connects to impalad. If Kerberos is not enabled on the instance of impalad to which you are connecting, errors are displayed.
<code>--query_option=</code> <code>"option=value</code> <code>-Q</code> <code>"option=value"</code>	Header line [<code>impala.query_option</code> s], followed on subsequent lines by <code>option=value</code> , one option per line.	Sets default query options for an invocation of the <code>impala-shell</code> command. To set multiple query options at once, use more than one instance of this command-line option. The query option names are not case-sensitive.
<code>-s kerberos_service_name</code> or <code>--kerberos_service_name=</code> <code>name</code>	<code>kerberos_service_name=name</code>	Instructs <code>impala-shell</code> to authenticate to a particular <code>impalad</code> service principal. If a <code>kerberos_service_name</code> is not specified, <code>impala</code> is used by default. If this option is used in conjunction with a connection in which Kerberos is not supported, errors are returned.
<code>-V</code> or <code>--verbose</code>	<code>verbose=true</code>	Enables verbose output.
<code>--quiet</code>	<code>verbose=false</code>	Disables verbose output.
<code>-v</code> or <code>--version</code>	<code>version=true</code>	Displays version information.
<code>-c</code>	<code>ignore_query_failure=true</code>	Continues on query failure.
<code>-d default_db</code> or <code>--database=default_db</code>	<code>default_db=default_db</code>	Specifies the database to be used on startup. Same as running the <code>USE</code> statement after connecting. If not specified, a database named <code>DEFAULT</code> is used.
<code>--ssl</code>	<code>ssl=true</code>	Enables TLS/SSL for <code>impala-shell</code> .
<code>--http_socket_timeout_s</code>	<code>http_socket_time</code> <code>out_s=HTTP_SOCKET_TIMEOUT</code>	Sets the timeout in seconds after which the socket will time out if the associated operation cannot be completed. Set to <code>None</code> to disable any timeout. This configurable option is only supported for <code>hs2-http</code> mode and the <code>DEFAULT</code> is <code>NONE</code> .
<code>path_to_certificate</code>	<code>ca_cert=path_to_certificate</code>	The local pathname pointing to the third-party CA certificate, or to a copy of the server certificate for self-signed server certificates. If <code>--ca_cert</code> is not set, <code>impala-shell</code> 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).
<code>-l</code>	<code>use_ldap=true</code>	Enables LDAP authentication.
<code>-u</code>	<code>user=user_name</code>	Supplies the username, when LDAP authentication is enabled by the <code>-l</code> option. (Specify the short username, not the full LDAP distinguished name.) The shell then prompts interactively for the password.
<code>--ldap_password_cmd=</code> <code>command</code>	N/A	Specifies a command to run to retrieve the LDAP password, when LDAP authentication is enabled by the <code>-l</code> option. If the command includes space-separated arguments, enclose the command and its arguments in quotation marks.
<code>--config_file=</code> <code>path_to_config_file</code>	N/A	Specifies the path of the file containing <code>impala-shell</code> configuration settings. The default is <code>/etc/impalarc</code> . This setting can only be specified on the command line.
<code>--live_progress</code>	N/A	Prints a progress bar showing roughly the percentage complete for each query. The information is updated interactively as the query progresses.
<code>--live_summary</code>	N/A	Prints a detailed report, similar to the <code>SUMMARY</code> command, showing progress details for each phase of query execution. The information is updated interactively as the query progresses.

Command-Line Option	Configuration File Setting	Explanation
<code>--var= variable_name= value</code>	N/A	Defines a substitution variable that can be used within the <code>impala-shell</code> session. The variable can be substituted into statements processed by the <code>-q</code> or <code>-f</code> options, or in an interactive shell session. Within a SQL statement, you substitute the value by using the notation <code>\${var:variable_name}</code> .
<code>--auth_creds_ok_in_clear</code>	N/A	Allows LDAP authentication to be used with an insecure connection to the shell. WARNING: This will allow authentication credentials to be sent unencrypted, and hence may be vulnerable to an attack.
<code>--protocol= protocol</code>	N/A	Protocol to use for the connection to Impala. Valid <i>protocol</i> values are: <ul style="list-style-type: none"> 'hs2': Impala-shell uses the binary TCP based transport to speak to the Impala Daemon via the HiveServer2 protocol. This is the current default setting. 'hs2-http': Impala-shell uses HTTP transport to speak to the Impala Daemon via the HiveServer2 protocol. 'beeswax': Impala-shell uses the binary TCP based transport to speak to the Impala Daemon via Beeswax. You cannot connect to the 3.2 or earlier versions of Impala using the 'hs2' or 'hs2-http' option. Beeswax support is deprecated and will be removed in the future.

Impala Shell Configuration File

You can store a set of default settings for `impala-shell` in the `impala-shell` configuration file.

The global `impala-shell` configuration file is located in `/etc/impalarc`.

The user-level `impala-shell` configuration file is located in `~/impalarc`.

Note that the global-level file name is different from the user-level file name. The global-level file name does not include a dot (`.`) in the file name.

The default path of the global configuration file can be changed by setting the `$IMPALA_SHELL_GLOBAL_CONFIG_FILE` environment variable.

To specify a different file name or path for the user-level configuration file, start `impala-shell` with the `--config_file impala-shell` option set to the path of the configuration file.

Typically, an administrator creates the global configuration file for the `impala-shell`, and if the user-level configuration file exists, the options set in the user configuration file take precedence over those in the global configuration file.

In turn, any options you specify on the `impala-shell` command line override any corresponding options within the configuration file.

The `impala-shell` configuration file (global or user) must contain a header label `[impala]`, followed by the options specific to `impala-shell`.

The `impala-shell` configuration file consists of key-value pairs, one option per line. Everything after the `#` character on a line is treated as a comment and ignored.

The names of the options in the configuration file are similar (although not necessarily identical) to the long-form command-line arguments to the `impala-shell` command. For the supported options in the configuration file, see [Impala Shell Configuration Options](#) on page 6.

You can specify key-value pair options using `keyval`, similar to the `--var` command-line option. For example, `keyval=variable1=value1`.

The query options specified in the `[impala]` section override the options specified in the `[impala.query_options]` section.

The following example shows a configuration file that you might use during benchmarking tests. It sets verbose mode, so that the output from each SQL query is followed by timing information. `impala-shell` starts inside the database containing the tables with the benchmark data, avoiding the need to issue a `USE` statement or use fully qualified table names.

In this example, the query output is formatted as delimited text rather than enclosed in ASCII art boxes, and is stored in a file rather than printed to the screen. Those options are appropriate for benchmark situations, so that the overhead of `impala-shell` formatting and printing the result set does not factor into the timing measurements. It also enables the `show_profiles` option. That option prints detailed performance information after each query, which might be valuable in understanding the performance of benchmark queries.

```
[impala]
verbose=true
default_db=tpc_benchmarking
write_delimited=true
output_delimiter=,
output_file=/home/tester1/benchmark_results.csv
show_profiles=true
keyval=msg1=hello,keyval=msg2=world
```

The following example shows a configuration file that connects to a specific remote Impala node, runs a single query within a particular database, then exits. Any query options predefined under the `[impala.query_options]` section in the configuration file take effect during the session.

You would typically use this kind of single-purpose configuration setting with the `impala-shell` command-line option `--config_file=path_to_config_file`, to easily select between many predefined queries that could be run against different databases, hosts, or even different clusters. To run a sequence of statements instead of a single query, specify the configuration option `query_file=path_to_query_file` instead.

```
[impala]
impalad=impala-test-nodel.example.com
default_db=site_stats
# Issue a predefined query and immediately exit.
query=select count(*) from web_traffic where event_date = trunc(now(),'dd')

[impala.query_options]
mem_limit=32g
```

Connecting to Impala Daemon in Impala Shell

In an `impala-shell` session, you need to connect to an `impalad` daemon to issue queries. When you connect to an `impalad`, and that daemon coordinates the execution of all queries sent to it.

About this task

Specify the connection information using the following options:

- Through command-line options when you run the `impala-shell` command.
- Through a configuration file that is read when you run the `impala-shell` command.
- During an `impala-shell` session, by issuing a `CONNECT` command.



Note: You cannot connect to the 3.2 or earlier versions of Impala using the 'hs2' or 'hs2-http' protocol (`--protocol` option).

To connect the Impala shell during shell startup:

1. Locate the hostname that is running an instance of the `impalad` daemon. If that `impalad` uses a non-default port (something other than port 21050) for `impala-shell` connections, find out the port number also.

- Use the `-i` option to the `impala-shell` interpreter to specify the connection information for that instance of `impalad`:

```
# When you are connecting to an impalad running on the same machine.
# The prompt will reflect the current hostname.
$ impala-shell
# When you are connecting to an impalad running on a remote machine, and
  impalad is listening
# on a non-default port over the HTTP HiveServer2 protocol.
$ impala-shell -i some.other.hostname:port_number --protocol='hs2-http'
# When you are connecting to an impalad running on a remote machine, and
  impalad is listening
# on a non-default port.
$ impala-shell -i some.other.hostname:port_number
```

To connect to an Impala in the `impala-shell` session:

- Start the Impala shell with no connection:

```
impala-shell
```

- Locate the hostname that is running the `impalad` daemon. If that `impalad` uses a non-default port (something other than port 21050) for `impala-shell` connections, find out the port number also.
- Use the `connect` command to connect to an Impala instance. Enter a command and replace `impalad-host` with the hostname you have configured to run Impala in your environment.

```
[Not connected] > connect impalad-host
[impalad-host:21050] >
```

To start `impala-shell` in a specific database:

You can use all the same connection options as in previous examples. For simplicity, these examples assume that you are logged into one of the Impala daemons.

- Find the name of the database containing the relevant tables, views, and so on that you want to operate on.
- Use the `-d` option to the `impala-shell` interpreter to connect and immediately switch to the specified database, without the need for a `USE` statement or fully qualified names:

```
# Subsequent queries with unqualified names operate on
# tables, views, and so on inside the database named 'staging'.
$ impala-shell -i localhost -d staging

# It is common during development, ETL, benchmarking, and so on
# to have different databases containing the same table names
# but with different contents or layouts.
$ impala-shell -i localhost -d parquet_snappy_compression
$ impala-shell -i localhost -d parquet_gzip_compression
```

To run one or several statements in non-interactive mode:

You can use all the same connection options as in previous examples. For simplicity, these examples assume that you are logged into one of the Impala daemons.

- Construct a statement, or a file containing a sequence of statements, that you want to run in an automated way, without typing or copying and pasting each time.
- Invoke `impala-shell` with the `-q` option to run a single statement, or the `-f` option to run a sequence of statements from a file. The `impala-shell` command returns immediately, without going into the interactive interpreter.

```
# A utility command that you might run while developing shell scripts
# to manipulate HDFS files.
$ impala-shell -i localhost -d database_of_interest -q 'show tables'
```

```
# A sequence of CREATE TABLE, CREATE VIEW, and similar DDL statements
# can go into a file to make the setup process repeatable.
$ impala-shell -i localhost -d database_of_interest -f recreate_tables.sql
```

Running Commands and SQL Statements in Impala Shell

This topic provides the commonly used syntax and shortcut keys in `impala-shell`.

The following are a few of the key syntax and usage rules for running commands and SQL statements in `impala-shell`.

- To see the full set of available commands, press TAB twice.
- To cycle through and edit previous commands, click the up-arrow and down-arrow keys.
- Use the standard set of keyboard shortcuts in GNU Readline library for editing and cursor movement, such as Ctrl-A for the beginning of line and Ctrl-E for the end of line.
- Commands and SQL statements must be terminated by a semi-colon.
- Commands and SQL statements can span multiple lines.
- Use `--` to denote a single-line comment and `/* */` to denote a multi-line comment.

A comment is considered part of the statement it precedes, so when you enter a `--` or `/* */` comment, you get a continuation prompt until you finish entering a statement ending with a semicolon. For example:

```
[impala] > -- This is a test comment
           > SHOW TABLES LIKE 't*';
```

- If a comment contains the `${variable_name}` and it is not for a variable substitution, the `$` character must be escaped, e.g. `-- \${hello}`.

Variable Substitution in impala-shell

You can define substitution variables to be used within SQL statements processed by `impala-shell`.

1. You specify the variable and its value as below.

- On the command line, you specify the option `--var=variable_name=value`
- Within an interactive session or a script file processed by the `-f` option, use the `SET VAR:variable_name=value` command.

2. Use the above variable in SQL statements in the `impala-shell` session using the notation: `${VAR:variable_name}`.

For example, here are some `impala-shell` commands that define substitution variables and then use them in SQL statements executed through the `-q` and `-f` options. Notice how the `-q` argument strings are single-quoted to prevent shell expansion of the `${var:value}` notation, and any string literals within the queries are enclosed by double quotation marks.

```
$ impala-shell --var=tname=table1 --var=colname=x --var=coltype=string -q '
CREATE TABLE ${var:tname} (${var:colname} ${var:coltype}) STORED AS PARQUET'
Query: CREATE TABLE table1 (x STRING) STORED AS PARQUET
```

The below example shows a substitution variable passed in by the `--var` option, and then referenced by statements issued interactively. Then the variable is reset with the `SET` command.

```
$ impala-shell --quiet --var=tname=table1

[impala] > SELECT COUNT(*) FROM ${var:tname};

[impala] > SET VAR:tname=table2;
[impala] > SELECT COUNT(*) FROM ${var:tname};
```

When you run a query, the live progress bar appears in the output of a query. The bar shows roughly the percentage of completed processing. When the query finishes, the live progress bar disappears from the console output.

Impala Shell Command Reference

Use the following commands within `impala-shell` to pass requests to the `impalad` daemon that the shell is connected to. You can enter a command interactively at the prompt or pass it as the argument to the `-q` option of `impala-shell`.

Command	Explanation
Impala SQL statements	You can issue valid SQL statements to be executed.
connect	Connects to the specified instance of <code>impalad</code> . The default port of 21050 is assumed unless you provide another value. You can connect to any host in your cluster that is running <code>impalad</code> . If you connect to an instance of <code>impalad</code> that was started with an alternate port specified by the <code>--fe_port</code> flag, you must provide that alternate port.
help	Help provides a list of all available commands and options.
history	Maintains an enumerated cross-session command history. This history is stored in the <code>~/impalahistory</code> file.
profile	Displays low-level information about the most recent query. Used for performance diagnosis and tuning. The report starts with the same information as produced by the <code>EXPLAIN</code> statement and the <code>SUMMARY</code> command.
quit	Exits the shell. Remember to include the final semicolon so that the shell recognizes the end of the command.
rerun or @	Executes a previous <code>impala-shell</code> command again, from the list of commands displayed by the <code>history</code> command. These could be SQL statements, or commands specific to <code>impala-shell</code> such as <code>quit</code> or <code>profile</code> . Specify an integer argument. A positive integer <code>N</code> represents the command labelled <code>N</code> in the output of the <code>HISTORY</code> command. A negative integer <code>-N</code> represents the <code>N</code> th command from the end of the list, such as <code>-1</code> for the most recent command. Commands that are executed again do not produce new entries in the <code>HISTORY</code> output list.
set	Manages query options for an <code>impala-shell</code> session. These options are used for query tuning and troubleshooting. Issue <code>SET</code> with no arguments to see the current query options, either based on the <code>impalad</code> defaults, as specified by you at <code>impalad</code> startup, or based on earlier <code>SET</code> statements in the same session. To modify option values, issue commands with the syntax <code>set option=value</code> . To restore an option to its default, use the <code>unset</code> command.
shell	Executes the specified command in the operating system shell without exiting <code>impala-shell</code> . You can use the <code>!</code> character as shorthand for the shell command.  Note: Quote any instances of the <code>--</code> or <code>/*</code> tokens to avoid them being interpreted as the start of a comment. To embed comments within <code>source</code> or <code>!</code> commands, use the shell comment character <code>#</code> before the comment portion of the line.
source or src	Executes one or more statements residing in a specified file from the local filesystem. Allows you to perform the same kinds of batch operations as with the <code>-f</code> option, but interactively within the interpreter. The file can contain SQL statements and other <code>impala-shell</code> commands, including additional <code>SOURCE</code> commands to perform a flexible sequence of actions. Each command or statement, except the last one in the file, must end with a semicolon.
summary	Summarizes the work performed in various stages of a query. It provides a higher-level view of the information displayed by the <code>EXPLAIN</code> command. Added in Impala 1.4.0. The time, memory usage, and so on reported by <code>SUMMARY</code> only include the portions of the statement that read data, not when data is written. Therefore, the <code>PROFILE</code> command is better for checking the performance and scalability of <code>INSERT</code> statements. You can see a continuously updated report of the summary information while a query is in progress.
unset	Removes any user-specified value for a query option and returns the option to its default value. You can also use it to remove user-specified substitution variables using the notation <code>UNSET VAR:variable_name</code> .
use	Indicates the database against which to run subsequent commands. Lets you avoid using fully qualified names when referring to tables in databases other than default. Not effective with the <code>-q</code> option, because that option only allows a single statement in the argument.
version	Returns Impala version information.

Connecting to a kerberized Impala daemon

Using an `impala-shell` session you can connect to an `impalad` daemon to issue queries. When you connect to an `impalad`, it coordinates the execution of all queries sent to it. You can run `impala-shell` to connect to a Kerberized Impala instance over HTTP in a cluster.

About this task

Kerberos is an enterprise-grade authentication system Impala supports. 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. Cloudera recommends using `impala-shell` with Kerberos authentication for strong security benefits while accessing an Impala instance.

Before you begin

- Locate the hostname that is running the `impalad` daemon.
- 28000 is the default port `impalad` daemon uses to transmit commands and receive results from client applications over HTTP using the HiveServer2 protocol. Ensure that this port is open.
- Ensure that the host running `impala-shell` has a preexisting kinit-cached Kerberos ticket that `impala-shell` can pass to the `impala` server automatically without the need for the user to reenter the password.
- To override any client connection errors, you should run the Kinit command to retrieve the Ticket Granting Ticket or to extend it if it has already expired.

Procedure

1. To enable Kerberos in the Impala shell, start the `impala-shell` command using the `-k` flag.
2. For `impala-shell` to communicate with the Impala daemon over HTTP through the HiveServer2 protocol, specify `--protocol=hs2-http` as the protocol.

```
impala-shell -i xxxx-cdh-7-2-3.vpc.cloudera.com -k --protocol=hs2-http
Starting Impala Shell with Kerberos authentication using Python 2.7.5
Using service name 'impala'
Warning: --connect_timeout_ms is currently ignored with HTTP transport.
Opened TCP connection to xxxx-cdh-7-2-3.vpc.cloudera.com:28000
Connected to xxxx-cdh-7-2-3.vpc.cloudera.com:28000
Server version: impalad version 4.0.0-SNAPSHOT RELEASE (build d18b1c1d3f
7230d3330b58928513c20e90bab0153)
```

Configuring ODBC for Impala

Download and configure the ODBC driver to integrate your applications with Impala.

About this task

Impala has been tested with the Impala ODBC driver version 2.5.42, and Cloudera recommends that you use this version with the current version of Impala.

Procedure

1. Download and install an ODBC driver.
2. Configure the ODBC port.

Versions 2.5 and 2.0 of the Cloudera ODBC Connector use the HiveServer2 protocol, corresponding to Impala port 21050.

Version 1.x of the Cloudera ODBC Connector uses the original HiveServer1 protocol, corresponding to Impala port 21000.

Configuring JDBC for Impala

Download and configure the JDBC driver to access Impala from a Java program that you write, or a Business Intelligence or similar tool that uses JDBC to communicate with database products.

About this task

The following are the default ports that Impala server accepts JDBC connections through:

Protocol	Default Port	Flag to Specify an Alternate Port
HTTP	28000	##hs2_http_port
Binary TCP	21050	##hs2_port

Make sure the port for the protocol you are using is available for communication with clients, for example, that it is not blocked by firewall software.

If your JDBC client software connects to a different port, specify that alternative port number with the flag in the above table when starting the `impalad`.

Procedure

1. Configure the JDBC port.

Impala server accepts JDBC connections through port 21050 by default. Make sure this port is available for communication with other hosts on your network, for example, that it is not blocked by firewall software. If your JDBC client software connects to a different port, specify that alternative port number in the Impala Daemon HiveServer2 Portfield in Cloudera Manager, in the Configuration tab.

2. Install the JDBC driver.

Impala has been tested using the Impala JDBC driver version 2.5.45 and 2.6.2. Cloudera recommends that you use one of these two versions with Impala.

3. Enable Impala JDBC support on client systems.

4. Establish JDBC connections.

The JDBC driver class depends on which driver you select.

Using the Cloudera JDBC Connector (recommended):

Depending on the level of the JDBC API your application is targeting, you can use the following fully-qualified class names (FQCNs):

- `com.cloudera.impala.jdbc41.Driver`
- `com.cloudera.impala.jdbc41.DataSource`
- `com.cloudera.impala.jdbc4.Driver`
- `com.cloudera.impala.jdbc4.DataSource`
- `com.cloudera.impala.jdbc3.Driver`
- `com.cloudera.impala.jdbc3.DataSource`

The connection string has the following format:

```
jdbc:impala://Host:Port[/Schema];Property1=Value;Property2=Value;...
```

The port value is typically 21050 for Impala.

To connect to an instance of Impala that requires Kerberos authentication, use a connection string of the form `jdbc:impala://host:port/principal=principal_name`. The principal must be the same user principal you used when starting Impala.

To connect to an instance of Impala that requires LDAP authentication, use a connection string of the form `jdbc:impala://host:port/db_name;user=ldap_userid;password=ldap_password`.

To connect to an instance of Impala over HTTP, specify the HTTP port, 28000 by default, and `transportMode=http` in the connection string.



Note: To establish a connection with an Impala instance from your client, using any authentication mode, you must use the connection string `jdbc:impala`. Using the connection string `jdbc:hive2` is not recommended and is not supported.

For updated information on the version of the JDBC driver you are using and for connection string examples for different supported authentications, refer to the link provided under Related Information.

Related Information

["Cloudera Enterprise Connector Documentation"](#)

Configuring Impyla for Impala

Explains how to install Impyla to connect to and submit SQL queries to Impala. Impyla is a Python client wrapper around the HiveServer2 Thrift Service. It connects to Impala and implements Python DB API 2.0.

About this task

Impyla releases are available at pypi.org. To get the available releases, check [Release history](#).



Note: Cloudera will not support versions of Impyla that are built manually from source code.

Key Features of Impyla

- HiveServer2 compliant.
- Works with Impala including nested data.
- [DB API 2.0 \(PEP 249\)](#)-compliant Python client (similar to `sqlite` or `MySQL` clients) supporting Python 2.6+ and Python 3.3+.

- Works with Kerberos, LDAP, SSL.
- [SQLAlchemy](#) connector.
- Converts to pandas DataFrame, allowing easy integration into the Python data stack (including scikit-learn and matplotlib); see the Ibis project for a richer experience.
- For more information, see [here](#).

Before you begin

Different systems require different packages to be installed to enable SASL support in Impyla. The following list shows some examples of how to install the packages on different distributions.

You must have the following installed in your environment before installing impyla. Python 2.6+ or 3.3+ and the pip packages six, bitarray, thrift and thriftpy2 will be automatically installed as dependencies when installing impyla. However if you clone the impyla repo and run their local copy, you must install these pip packages manually.

- Install the latest pip and setuptools:

```
python -m pip install --upgrade pip setuptools
```

Optionally, to install Impyla with Hive and/or GSSAPI (kerberos) support, you will also need to install additional software packages:

- RHEL/CentOS:

```
sudo yum install gcc-c++ cyrus-sasl-md5 cyrus-sasl-plain cyrus-sasl-gssapi
cyrus-sasl-devel
```

- Ubuntu:

```
sudo apt install g++ libsasl2-dev libsasl2-2 libsasl2-modules-gssapi-mit
```

Procedure

1. Using pip you can install the latest release: `pip install impyla`
2. Optionally, to install Impyla with GSSAPI (kerberos) support: `pip install impyla[kerberos]`
3. You also need to pip-install pandas for conversion to DataFrame objects or sqlalchemy for the SQLAlchemy engine.

Example

Sample codes

Impyla implements the [Python DB API v2.0 \(PEP 249\)](#) database interface (refer to it for API details):

```
from impala.dbapi import connect

conn = connect(host = "my.host.com", port = 21050)
cursor = conn.cursor()
cursor.execute("SELECT * FROM mytable LIMIT 100")
print(cursor.description) # prints the result set's schema
results = cursor.fetchall()
cursor.close()
conn.close()
```

The Cursorobject also exposes the iterator interface, which is buffered (controlled by cursor.arraysize):

```
cursor.execute("SELECT * FROM mytable LIMIT 100")
for row in cursor:
    print(row)
```

Furthermore the Cursor object returns you information about the columns returned in the query. This is useful to export your data as a csv file.

```
import csv

cursor.execute("SELECT * FROM mytable LIMIT 100")
columns = [datum[0] for datum in cursor.description]
targetfile = "/tmp/foo.csv"

with open(targetfile, "w", newline = "") as outcsv:
    writer = csv.writer(
        outcsv,
        delimiter = ",",
        quotechar = "'",
        quoting = csv.QUOTE_ALL,
        lineterminator = "\n")
    writer.writerow(columns)
for row in cursor:
    writer.writerow(row)
```

You can also get back a pandas DataFrame object

```
from impala.util import as_pandas

# carry df through scikit-learn, for example
df = as_pandas(cur)
```

Connecting over HTTP/HTTPS

GSSAPI (kerberos) authentication over HTTP has been supported since 0.17a1. Use this example to establish connection over HTTP/HTTPS.

```
from impala.dbapi import connect

conn = connect(
    "impala-coordinator.example.com",
    28000,
    auth_mechanism = "GSSAPI",
    kerberos_service_name = "impala",
    use_http_transport = True,
    http_path = "cliservice",
    auth_cookie_name = "impala.auth")
cursor = conn.cursor()
cursor.execute("SHOW DATABASES")
res = cursor.fetchall()
cursor.close()
conn.close()
```

Connecting to DataHub

Lists an example code to connect to Impala with LDAP over http using LDAP as the authentication mechanism.

Example

Sample code

```
from impala.dbapi import connect

conn = connect(
    host = "aaaaaaa-aaaa-master0.se-sandb.a465-9q4k.cloudera.site",
    port = 443,
    auth_mechanism = "LDAP",
    use_ssl = True,
```

```
use_http_transport = True,
http_path = "aaaaaaa-aaaa/cdp-proxy-api/impala",
user = "aaaaaaa",
password = "xxxxx")
cursor = conn.cursor()
cursor.execute("SELECT * FROM default.emax_temp")

for row in cursor:
    print(row)
cursor.close()
conn.close()
```

Connecting to DataHub Data Mart

Lists an example code to connect to DataHub Data Mart using LDAP as the authentication mechanism.

Before you begin

Must have the latest Impyla release.

Example

Sample code

```
from impala.dbapi import connect

conn = connect(
    "xxxxxxx-data-mart-master0.xxxxxxx.xcu2-8y8x.dev.cldr.work",
    443,
    auth_mechanism = "LDAP",
    user = "XXXXX",
    password = "XXXXX",
    use_ssl = True,
    use_http_transport = True,
    http_path = "xxxxxxx-data-mart/cdp-proxy-api/impala")
cursor = conn.cursor()
cursor.execute("SHOW DATABASES")
res = cursor.fetchall()
print(res)
cursor.close()
conn.close()
```

Configuring Delegation for Clients

Impala supports user and group delegation for client connections.

About this task

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. Impala supports “delegation” where users whose names you specify can delegate the execution of a query to another user. The query runs with the privileges of the delegated user, not the original authenticated user.

You also have an option to delegate using groups. Instead of listing a large number of delegated users, you can create a group of those users and specify the delegated group name in the Impalad startup option. The client sends the delegated user name, and Impala performs an authorization to see if the delegated user belongs to a delegated group.

The name of the delegated user is passed using the HiveServer2 protocol configuration property `impala.doas.user` when the client connects to Impala.

When the client connects over HTTP, the doAs parameter can be specified in the HTTP path. For example:

```
/?doAs=delegated_user
```

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.



Attention:

- When the delegation is enabled in Impala, the Impala clients should take an extra caution to prevent unauthorized access for the delegate-able users.
- Impala requires Apache Ranger on the cluster to enable delegation. Without Ranger installed, the delegation feature will fail with the following error: User user1 is not authorized to delegate to user2. User/group delegation is disabled.

Procedure

To enable delegation:

1. In Cloudera Manager, navigate to ClustersImpala.
2. In the Configuration tab, click Impala-1 (Service-Wide) in the Scope and click Security in the Category.
3. 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.
4. In the Proxy Group Configuration field, type the a semicolon-separated list of key=value pairs of authorized proxy users to the group(s) they can impersonate.
The list of delegated groups are delimited with a comma, e.g. hue=group1, group2.
5. Click Save Changes and restart the Impala service.

Spooling Impala Query Results

In Impala, you can control how query results are materialized and returned to clients, e.g. impala-shell, Hue, JDBC apps.

Result spooling is turned off by default, but can be enabled via the SPOOL_QUERY_RESULTS query option.

- When query result spooling is disabled, Impala relies on clients to fetch results to trigger the generation of more result row batches until all the result rows have been produced. If a client issues a query without fetching all the results then that query continues to be in the running state indefinitely and the query fragments continue to consume the resources until the query is cancelled and unregistered, potentially tying up resources and causing other queries to wait for an extended period of time in admission control.

Impala would materialize rows on-demand where rows are created only when the client requests them.

For example, if a Hue user runs a complex query that returns 1000 rows, but does not scroll through all the returned rows, and then stays idle for a while, the query will remain running and will hold onto all of its resources until it is explicitly closed or the session times out.

- When query result spooling is enabled, result sets of queries are eagerly fetched and spooled in the spooling location, either in memory or on disk.

Once all result rows have been fetched and stored in the spooling location, the resources are freed up. Incoming client fetches can get the data from the spooled results.

Admission Control and Result Spooling

Query results spooling collects and stores query results in memory that is controlled by admission control. Use the following query options to calibrate how much memory to use and when to spill to disk.

MAX_RESULT_SPOOLING_MEM

The maximum amount of memory used when spooling query results. If this value is exceeded when spooling results, all memory will most likely be spilled to disk. Set to 100 MB by default.

MAX_SPILLED_RESULT_SPOOLING_MEM

The maximum amount of memory that can be spilled to disk when spooling query results. Must be greater than or equal to MAX_RESULT_SPOOLING_MEM. If this value is exceeded, the coordinator fragment will block until the client has consumed enough rows to free up more memory. Set to 1 GB by default.

Fetch Timeout

Resources for a query are released when the query completes its execution. To prevent clients from indefinitely waiting for query results, use the FETCH_ROWS_TIMEOUT_MS query option to set the timeout when clients fetch rows. Timeout applies both when query result spooling is enabled and disabled:

- When result spooling is disabled (SPOOL_QUERY_RESULTS = FALSE), the timeout controls how long a client waits for a single row batch to be produced by the coordinator.
- When result spooling is enabled (SPOOL_QUERY_RESULTS = TRUE), a client can fetch multiple row batches at a time, so this timeout controls the total time a client waits for row batches to be produced.

Explain Plans

Below is the part of the EXPLAIN plan output for result spooling.

```
F01:PLAN FRAGMENT [UNPARTITIONED] hosts=1 instances=1
| Per-Host Resources: mem-estimate=4.02MB mem-reservation=4.00MB thread-reservation=1
PLAN-ROOT SINK
| mem-estimate=4.00MB mem-reservation=4.00MB spill-buffer=2.00MB thread-reservation=0
```

- The mem-estimate for the PLAN-ROOT SINK is an estimate of the amount of memory needed to spool all the rows returned by the query.
- The mem-reservation is the number and size of the buffers necessary to spool the query results. By default, the read and write buffers are 2 MB in size each, which is why the default is 4 MB.

PlanRootSink

In Impala, the PlanRootSink class controls the passing of batches of rows to the clients and acts as a queue of rows to be sent to clients.

- When result spooling is disabled, a single batch or rows is sent to the PlanRootSink, and then the client must consume that batch before another one can be sent.
- When result spooling is enabled, multiple batches of rows can be sent to the PlanRootSink, and multiple batches can be consumed by the client.

Related Information

[Impala query options](#)

Shut Down Impala

Explains how to gracefully shut down Impala Daemons by first allowing running queries a specified amount of time to complete the process.

The flow to gracefully shut down an Impala Daemon is as follows:

1. The shutdown is initiated.
2. The grace period starts. The Impala Daemon informs other coordinators not to schedule any new queries on it. This allows queries already scheduled to run on this daemon by other coordinators to start executing.
3. The grace period expires.
4. The Impala Daemon continuously checks if there are no queries or fragments running.
5. If there are no queries or fragments running, it shuts down.
6. Otherwise, when it reaches the `IMPALA_GRACEFUL_SHUTDOWN_DEADLINE` duration, Impala Daemon shuts down.

Once Cloudera Manager initiates the stop/shutdown command, the Impala Daemon starts up the graceful shutdown process, and the process cannot be reverted. However, if you need to change the hard deadline in Cloudera Manager, you can cancel the shutdown command, change the Impala Daemon Graceful Shutdown, and start the shutdown command again.

Procedure

1. Optionally, set the grace period.
 - a) In Cloudera Manager, navigate to Impala ServiceConfigurationScopeImpala Daemon.
 - b) In the Impala Daemon Command Line Argument Advanced Configuration Snippet (Safety Valve) field, specify the grace period:

```
--shutdown_grace_period_s=<new grace period in seconds>
```

The default grace period is 2 minutes.
It is strongly recommended that you use the default value and not change the setting.
2. Optionally, set the hard deadline after which Impala is shut down regardless of whether queries are still running on it.
 - a) In Cloudera Manager, navigate to Impala ServiceConfigurationScopeImpala Daemon.
 - b) In the Impala Graceful Shutdown Deadline field, specify the time to wait for running queries.
The default is 60 minutes.
If you specify 0, Impala will shutdown immediately without waiting for running queries
3. In Cloudera Manager, navigate to Impala ServiceInstances.
4. Click an Impala Daemon role.
5. Click ActionsImpala Daemon Graceful Shutdown.

Related Information

[SHUTDOWN statement](#)

Setting Timeouts in Impala

Depending on how busy your cluster is, you might increase or decrease various timeout values. Increase timeouts if Impala is cancelling operations prematurely, when the system is responding slower than usual but the operations are still successful if given extra time. Decrease timeouts if operations are idle or hanging for long periods, and the idle or hung operations are consuming resources and reducing concurrency.

Setting Timeout and Retries for Thrift Connections to Backend Client

Impala connections to the backend client are subject to failure in cases when the network is momentarily overloaded.

About this task

To avoid failed queries due to transient network problems, you can configure the number of Thrift connection retries using the following option:

Procedure

1. In Cloudera Manager, navigate to Impala serviceConfiguration.
2. In the Impala Daemon Command Line Argument Advanced Configuration Snippet (Safety Valve) field, specify the following.

To avoid failed queries due to transient network problems, you can configure the number of Thrift connection retries using the following option:

- The `--backend_client_connection_num_retries` option specifies the number of times Impala will try connecting to the backend client after the first connection attempt fails. By default, `impalad` will attempt three re-connections before it returns a failure.

You can configure timeouts for sending and receiving data from the backend client. Therefore, if for some reason a query does not respond, instead of waiting indefinitely for a response, Impala will terminate the connection after a configurable timeout.

- The `--backend_client_rpc_timeout_ms` option can be used to specify the number of milliseconds Impala should wait for a response from the backend client before it terminates the connection and signals a failure. The default value for this property is 300000 milliseconds, or 5 minutes.

3. Click Save Changes and restart Impala.

Increasing StateStore Timeout

If you have an extensive Impala schema, for example, with hundreds of databases, tens of thousands of tables, you might encounter timeout errors during startup as the Impala catalog service broadcasts metadata to all the Impala nodes using the StateStore service. To avoid such timeout errors on startup, increase the StateStore timeout value from its default of 10 seconds.

About this task

Increase the timeout value of the StateStore service if you see messages in the `impalad` log such as:

```
Connection with state-store lost
Trying to re-register with state-store
```

Procedure

1. In Cloudera Manager, navigate to Impala serviceConfiguration.
2. In the search field, type `-statestore_subscriber_timeout_seconds`.
3. In the StateStoreSubscriber Timeout field, specify a new timeout value larger than the current value.
4. Click Save Changes and restart Impala.

Setting the Idle Query and Idle Session Timeouts

To keep long-running queries or idle sessions from tying up cluster resources, you can set timeout intervals for both individual queries, and entire sessions.

About this task

Procedure

1. In Cloudera Manager, navigate to Impala serviceConfiguration.
2. In the search field, type `idle`.
3. In the Idle Query Timeout field, specify the time in seconds after which an idle query is cancelled.

This could be a query whose results were all fetched but was never closed, or one whose results were partially fetched and then the client program stopped requesting further results. This condition is most likely to occur in a client program using the JDBC or ODBC interfaces, rather than in the interactive `impala-shell` interpreter. Once a query is cancelled, the client program cannot retrieve any further results from the query.

You can reduce the idle query timeout by using the `QUERY_TIMEOUT_S` query option at the query level. Any non-zero value specified in this field serves as an upper limit for the `QUERY_TIMEOUT_S` query option.

The value of 0 disables query timeouts.

4. In the Idle Session Timeout field, specify the time in seconds after which an idle session expires.

A session is idle when no activity is occurring for any of the queries in that session, and the session has not started any new queries. Once a session is expired, you cannot issue any new query requests to it. The session remains open, but the only operation you can perform is to close it.

The default value of 0 specifies sessions never expire.

You can override this setting with the `IDLE_SESSION_TIMEOUT` query option at the session or query level.

5. Click Save Changes and restart Impala.

Results

Impala checks periodically for idle sessions and queries to cancel. The actual idle time before cancellation might be up to 50% greater than the specified configuration setting. For example, if the timeout setting was 60, the session or query might be cancelled after being idle between 60 and 90 seconds.

Adjusting Heartbeat TCP Timeout Interval

Using the TCP flag, you can prevent the Statestore from waiting indefinitely for a response from the subscribers that fail to respond to the heartbeat RPC within the set period.

About this task

This flag `statestore_heartbeat_tcp_timeout_seconds` defines the time that may elapse before a heartbeat RPC connection request from a Statestore server to an Impalad or a Catalog server (subscribers) should be considered dead.

You can increase the flag value if you see intermittent heartbeat RPC timeouts listed in the statestore's log. You may find the max value of "statestore.priority-topic-update-durations" metric on the statestore to get an idea of a reasonable value to be used in this flag.



Note: The priority topic updates are only small amounts of data that take little time to process, similar to the heartbeat complexity.

Procedure

1. In Cloudera Manager, navigate to Impala serviceConfiguration.
2. In the Impala Daemon Command Line Argument Advanced Configuration Snippet (Safety Valve) field, add the flag `statestore_heartbeat_tcp_timeout_seconds` with an appropriate value.
3. You can also control the maximum number of consecutive heartbeat messages an impalad can miss before being declared failed by the statestore by adding this flag `statestore_max_missed_heartbeats`. Typically, you will not have to change this value.
4. Click Save Changes and restart Impala.