

Hortonworks Connector for Teradata

(Mar 30, 2015)

Hortonworks Connector for Teradata

Copyright © 2015 Hortonworks, Inc. All rights reserved.

The software documented herein is Copyright (c) 2012-2014 Hortonworks, Inc., all rights reserved. It is distributed under the Hortonworks End User License Agreement (EULA), which you should read and agree to before using the software. You may not use this software except in compliance with the Hortonworks EULA.

Unless required by applicable law or agreed to in writing, software distributed under the EULA is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the EULA for the specific language governing permissions and limitations under the License.

Please visit the [Hortonworks Data Platform](#) page for more information on Hortonworks technology. For more information on Hortonworks services, please visit either the [Support](#) or [Training](#) page. Feel free to [Contact Us](#) directly to discuss your specific needs.

Table of Contents

1. Hortonworks Connector for Teradata	1
Introduction	1
Background	1
Supported Features	1
Software Versions and Installation	2
Connector Version	2
Supported Product Versions	3
Requirements and Dependencies	3
Installation	3
Configuration	4
Database Connection Credentials	4
Configuration Options	4

1. Hortonworks Connector for Teradata

Introduction

Hortonworks Connector for Teradata (powered by the Teradata Connector for Hadoop) is an implementation of a Sqoop connector that enables those conversant with the [Apache Sqoop](#) tool to transfer data between the Teradata MPP DBMS and Apache Hadoop environments.

Background

Sqoop provides facilities for bulk transfer of data between external data stores and the Hadoop environment exploiting the Map Reduce paradigm. Sqoop depends on JDBC interfaces to access the external databases.

Most of the databases also have specialized access methods for high speed bulk data transfers for efficient batch processing needs, such as backups, etc.

To accommodate the varieties of database mechanisms to facilitate bulk transfer, Sqoop provides extensible base implementations of the data transfer functions utilizing the JDBC interface that can optionally be enhanced to suit a database-specific method of data transfer.

Terminology

Sqoop has the notion of *Connectors*, which contain the specialized logic to read and write to external systems.

- The *Hortonworks Connector for Teradata* ("Hortonworks Connector") is a Sqoop Connector implementation for Teradata.
- It is built on the *Teradata Connector for Hadoop*, a Teradata product.

Supported Features

The Hortonworks Connector supports the following features:

- Import/Export tools that run Hadoop MR jobs to transfer data.
- Support for Text, Sequence, ORCFiles, Avro, and RCFiles as the source for export operations and target for import operations.



Note

If you will run Avro jobs, download `avro-mapred-1.7.4-hadoop2.jar` and place it under `$$SQOOP_HOME/lib`.

- Import table or query data from Teradata to:
 - an existing partitioned or non-partitioned Hive table.

- a new partitioned or non-partitioned Hive table created by the connector.
- an HCatalog table.
- Export data from HDFS files, Hive or HCatalog tables to empty or non-empty Teradata tables.
- Facilities for mapping schemas between Teradata and Hive/HCatalog, including necessary data type conversions.

Connector Feature Checklist

Import all tables: Supported.

Incremental import: Sqoop options are not supported but can be emulated, as specified in the sample invocation [Incremental Import](#).

BLOB and CLOB: Limited to 64 KB.

Import data to Sqoop

- **TextFormat, delimited:** Supported.
- **SequenceFile:** Supported.
- **RCFile:** Supported.
- **ORCFile:** Supported
- **Avro file:** Supported

Hive arguments: Support for all standard Hive arguments. All data types except Union are supported.

Export from / import to HCatalog table: Supported.

Automatic schema mapping to/from HCatalog: Supported.

Import using a query: Supported.

Update table: Not supported.

Compression: Not supported.

Software Versions and Installation

Connector Version

This document discusses the Hortonworks Connector for Teradata ("Hortonworks Connector") built on version 1.3.4 of the Teradata Connector for Hadoop.

Supported Product Versions

This section lists the product versions supported in the current release of the Hortonworks Connector.

Teradata Database Versions

The following Teradata database versions are supported:

- Teradata Database 13.00
- Teradata Database 13.10
- Teradata Database 14.00
- Teradata Database 14.10
- Teradata Database 15.00

Hive Version

- Hive 0.14

Hadoop Version

- HDP 2.2 or later (prior versions are untested, but may also work)

Sqoop Versions

- Sqoop 1.4.5

Requirements and Dependencies

System Requirements

The Hortonworks Connector requires JRE/JDK 1.6 or later versions.

Dependencies

1. Teradata GSS Client Driver 14.00 or later versions (tdgssconfig)
2. Teradata JDBC Driver 14.00 or later versions (terajdbc)
3. Teradata Connector for Hadoop 1.3.4

Installation

Installation Dependencies

Sqoop must be installed first.

Installing the Software

1. Download the tarball from the "Add-Ons" for Hortonworks Data Platform 2.2.0 at <http://hortonworks.com/hdp/addons>.
2. Extract the contents of the tar archive to \$SQOOP_HOME/lib. Sqoop will then distribute the contents of the tar to the necessary nodes.

Configuration

This section provides information about connection credentials and configuration options.

Database Connection Credentials

Refer to [Sqoop documentation](#) for the Teradata database connection credentials.

Documentation for Sqoop version 1.4.5 is available here: <http://sqoop.apache.org/docs/1.4.5/index.html>.

Configuration Options

The Hortonworks Connector defines many connector-specific options. A good selection of them is also available as Sqoop options (although not all Sqoop options are directly translatable to Hortonworks Connector options).

Configuration Option Precedence

Options can be specified using any of these techniques:

- a configuration file
- `-D` command line option
- Sqoop options (where applicable); apart from standard Sqoop options, a few connector-specific options are supported

Therefore the following precedence is established:

1. Sqoop connector-specific extra arguments have the highest precedence. (Sqoop command line options must match, or execution will fail.)
2. If `-D` command line options are provided, they override the configuration file values.
3. The value in the configuration file is the default.

As an example, if the configuration file sets the number of input mappers to 4 and the command line option (`-D com.teradata.db.input.num.mappers`) sets it to 5, but the Sqoop option `--num-mappers` is set to 6, then the import job will use 6 mappers.

In some cases, option constraints and the relationships between options affect the configuration value used. For example, import options `job.type` and `file.format` are interrelated. These options are described in [Connector Import Options](#).

Sqoop Options

The Sqoop option `--connection-manager` must be set as follows to use the Hortonworks Connector for Teradata (see the [Sample Invocations](#)):

```
--connection-manager org.apache.sqoop.teradata.TeradataConnManager
```

Some of the Sqoop options are unsupported in the current release of the Hortonworks Connector for Hadoop. See the [Appendix](#) for a list of unsupported Sqoop options.

Hortonworks Connector Options

The [Appendix](#) describes the Hortonworks Connector options, including [Connector Import Options](#), [Connector Export Options](#) and [Connector-specific Extra Arguments](#).

Data Type Support

The Hortonworks Connector data types depend on Teradata database types.

Support for Teradata Data Types

BIGINT	TIME (n)	INTERVAL HOUR (n) TO SECOND (m)
BYTEINT	TIMESTAMP (n)	INTERVAL MINUTE (n)
INTEGER	PERIOD (DATE)	INTERVAL MINUTE (n) TO SECOND (m)
SMALLINT	PERIOD (TIME (n))	INTERVAL SECOND (n)
DOUBLE PRECISION	PERIOD (TIMESTAMP (n))	<i>The following data types are supported with some limitations:</i> <ul style="list-style-type: none"> • BYTE (n) ^{1} • VARBYTE (n) ^{1} • BLOB ^{{1}{2}} • CLOB ^{2} • ARRAY ^{3} <p>^{1} <i>Converted to HEX string</i></p> <p>^{2} <i>BLOB and CLOB datatypes are limited to 64 KB in length</i></p> <p>^{3} <i>Can be used for InputFormat only</i></p>
FLOAT	INTERVAL YEAR (n)	
REAL	INTERVAL YEAR (n) TO MONTH	
DECIMAL (n,m)	INTERVAL MONTH (n)	
NUMERIC (n,m)	INTERVAL DAY (n)	
NUMBER (n,m)	INTERVAL DAY (n) TO HOUR	
CHAR (n)	INTERVAL DAY (n) TO MINUTE	
VARCHAR (n)	INTERVAL DAY (n) TO SECOND (m)	
LONG VARCHAR	INTERVAL HOUR (n)	
DATE	INTERVAL HOUR (n) TO MINUTE	

Support for Hive Data Types

BIGINT	<i>The following Hive types are supported with some limitations:</i> <ul style="list-style-type: none"> • BINARY ^{A} • MAP ^{B} • ARRAY ^{B} • STRUCT ^{B} • TIMESTAMP ^{C}
INT	
SMALLINT	
TINYINT	
DOUBLE	
FLOAT	

STRING	[A] Supported with Hive 0.10.0 or later
BOOLEAN	[B] Converted to/from VARCHAR in JSON format on the Teradata system
	[C] Custom formats are not supported

Unsupported Data Types

<p>These Teradata types are unsupported:</p> <ul style="list-style-type: none"> • GRAPHIC • VARGRAPHIC • LONG VARGRAPHIC 	<p>This Hive type is unsupported:</p> <ul style="list-style-type: none"> • UNION
---	---

Hive and HCatalog Support

Importing from Hive and HCatalog requires that HADOOP_CLASSPATH and LIB_JARS be specified before the **sqoop** command is run. This shows the environment variable setup:

```
export HADOOP_CLASSPATH=$(hcat -classpath)
HIVE_HOME=/usr/lib/hive
HCAT_HOME=/usr/lib/hcatalog

export LIB_JARS=$HCAT_HOME/share/hcatalog/hcatalog-core-0.12.0.2.1.3.0-563.jar,\
                $HIVE_HOME/lib/hive-metastore-0.12.0.2.1.3.0-563.jar,\
                $HIVE_HOME/lib/libthrift-0.9.0.jar,\
                $HIVE_HOME/lib/hive-exec-0.12.0.2.1.3.0-563.jar,\
                $HIVE_HOME/lib/libfb303-0.9.0.jar,\
                $HIVE_HOME/lib/jdo2-api-3.0.1.jar,\
                $HIVE_HOME/lib/slf4j-api-1.7.5.jar,\
                $HIVE_HOME/lib/hive-cli-0.12.0.2.1.3.0-563.jar,\
                $HIVE_HOME/lib/hive-builtins-0.12.0.2.1.3.0-563.jar
```



Note

Change the HIVE_HOME and HCAT_HOME variables as needed and change the versions of the jar to what is available under the directories mentioned.

Hive and HCatalog jobs can be run as shown in the next section.

Sample Invocations

The following examples assume that the SQOOP_HOME environment variable is set to the base directory of the Sqoop installation.

Import Data from Teradata to Hadoop and Hive

```
$SQOOP_HOME/bin/sqoop import \
  -libjars $LIB_JARS \
  --connection-manager org.apache.sqoop.teradata.TeradataConnManager \
  --username tduser \
  --password tduserpass \
  --table tablename \
  --hcatalog-table hcat_table
```

Import Data from Teradata into an HCatalog Table

```
$SQOOP_HOME/bin/sqoop import \  
-libjars $LIB_JARS \  
--connect jdbc:teradata://td-host/Database=dbname \  
--connection-manager org.apache.sqoop.teradata.TeradataConnManager \  
--username tduser \  
--password tduserpass \  
--table tablename \  
--hcatalog-table hcat_table
```

Incremental Import

Teradata incremental import emulates the check-column and last value options. Here is an example for a table which has 'hire_date' as the date column to check against and 'name' as the column that can be used to partition the data.

```
export USER=dbc  
export PASS=dbc  
export HOST=<dbhost>  
export DB=<dbuser>  
export TABLE=<dbtable>  
export JDBCCURL=jdbc:teradata://$HOST/DATABASE=$DB  
export IMPORT_DIR=<hdfs-dir to import>  
export VERBOSE=--verbose  
export MANAGER=org.apache.sqoop.teradata.TeradataConnManager  
export CONN_MANAGER="--connection-manager $MANAGER"  
export CONNECT="--connect $JDBCCURL"  
MAPPERS="--num-mappers 4"  
DATE="'1990-12-31'"  
FORMAT="'yyyy-mm-dd'"  
LASTDATE="cast($DATE as date format $FORMAT)"  
SQOOPQUERY="select * from employees where hire_date < $LASTDATE AND \  
$CONDITIONS"  
$SQOOP_HOME/bin/sqoop import $TDQUERY $TDSPLITBY $INPUTMETHOD $VERBOSE  
$CONN_MANAGER $CONNECT -query "$SQOOPQUERY" --username $USER --password $PASS  
--target-dir $IMPORT_DIR --split-by name
```

Export Data to Teradata

```
$SQOOP_HOME/bin/sqoop export \  
--connect jdbc:teradata://172.16.68.128/Database=employees \  
--connection-manager org.apache.sqoop.teradata.TeradataConnManager \  
--username dbc \  
--password dbc \  
--table employees2 \  
--export-dir /user/hrt_qa/test-sqoop/out \  
--batch
```

Known Issues and Workarounds

Stage Tables

Issue: The export option `--stage-table` does not work.

Cause: The behavior of stage tables is different between Hortonworks Connector and Sqoop, and this causes deadlocks during job cleanup if the Sqoop `-staging-table` option is used.

Workaround: Use the Hortonworks Connector option `teradata.db.output.stage.table.name` for specifying the stage table name.

Fastload

Issue: The export option `'fastload.soclet.host'` does not work.

Cause: The `internal.fastload` method used for Teradata exports can cause resource exhaustion (running out of database AMPs) if the number of reducers exceeds the number of available AMPs.

Workaround: Use the option `teradata.db.output.num.reducers` to restrict the resource usage.

Appendix: Configuration Options

This appendix describes the Hortonworks Connector configuration options and lists the Sqoop options that are currently unsupported.

- [Sqoop Options](#)
- [Hortonworks Connector Options](#)

Sqoop Options

To use the Hortonworks Connector, you must set the Sqoop option `--connection-manager` to `org.apache.sqoop.teradata.TeradataConnManager` as shown in the [Sample Invocations](#).

Some of the Sqoop options are unsupported in the current release of the Hortonworks Connector for Hadoop. The tables below list the unsupported import and export options.



Note

Imports and exports are defined from the Hadoop perspective, that is, an import brings data into Hadoop from the database and an export moves data out of Hadoop into the database.

Unsupported Sqoop Import Options

Import Category	Unsupported Options
Control Options	<code>--append</code> <code>--compression-codec</code> <code>--direct</code> <code>--direct-split-size</code> <code>--where</code> <code>--compress, -z</code>
Incremental Options	<code>--check-column</code> <code>--incremental</code>

Import Category	Unsupported Options
	<code>--last-value</code>
Output Formatting Options	<code>--mysql-delimiters</code> <code>--optionally-enclosed-by</code>
Hive Support Options	<code>--hive-delims-replacement</code> <code>--hive-drop-import-delims</code> <code>--hive-home</code> <code>--hive-overwrite</code> <code>--hive-partition-key</code> <code>--hive-partition-value</code> <code>--map-column-hive</code>
HBase Support Options	<code>--column-family</code> <code>--hbase-create-table</code> <code>--hbase-row-key</code> <code>--hbase-table</code>
Data Mapping Options	<code>--map-column-java</code>

Unsupported Sqoop Export Options

Export Category	Unsupported Options
Control Options	<code>--batch</code> <code>--clear-staging-table</code> <code>--direct</code> <code>--update-key</code> <code>--update-mode</code>
Input Parsing Options	<code>--input-lines-terminated-by</code> <code>--input-optionally-enclosed-by</code>
Data Mapping Options	<code>--map-column-java</code>

Hortonworks Connector Options

This section describes configuration options provided by the Hortonworks Connector.

- [Connector Import Options](#)
- [Connector Export Options](#)

For information about how the options can be specified, see [Configuration Option Precedence](#).



Note

Imports and exports are defined from the Hadoop perspective, that is, an import brings data into Hadoop from the database and an export moves data out of Hadoop into the database.

Connector Import Options

All option names below are prefixed by "teradata.db.input." when specified in the configuration files or in the -D command line option.

For example, the job.type option is specified as teradata.db.input.job.type.



Note

Only options with an overriding Sqoop option are listed.

Connector Import Option (teradata.db.input.*)	Description	Overriding Sqoop Option
job.type	The type of import job. Required: no Supported values: hcat, hive, hdfs Default value: hdfs	None for 'hcat' and 'hive' settings; also none for 'hdfs' when the file format is 'textfile'. But for file formats other than 'textfile' the 'hdfs' job type is reset to 'hive', therefore the following Sqoop option overrides a job.type of 'hdfs': <code>--as-sequencefile</code>
file.format	The format of a to-be-imported data file in HDFS. An 'hcat' or 'hive' job type supports 'orcfile', 'rcfile', 'sequencefile', and 'textfile' file formats; and an 'hdfs' job type supports only 'textfile' format. Required: no Supported values: orcfile, rcfile, sequencefile, textfile Default value: textfile	<code>--as-sequencefile</code> <code>--as-textfile</code>
target.paths	The directory with which to place the imported data. It is required for an 'hdfs' job, optional for a 'hive' job, and not valid for an 'hcat' job. For a 'hive' job, either specify this or the 'target.table' parameter but not both. Required: no Supported values: string Default value: The value of property 'mapred.output.dir'	<code>--target-dir</code> <code>--warehouse-dir</code>
num.mappers	The number of mappers for the import job. It is also the number of splits the Hortonworks Connector will attempt to create. Required: no Supported values: an integer greater than 0 Default value: 2	<code>-m</code> <code>--num-mappers</code>
source.query	The SQL query to select data from a Teradata database; either specify this or the 'source.table' parameter, but not both.	<code>--query</code>

Connector Import Option (<code>teradata.db.input.*</code>)	Description	Overriding Sqoop Option
	<p>Required: no</p> <p>Supported values: The select SQL query (Teradata database supported)</p>	
<code>source.table</code>	<p>The name of the source table in a Teradata system from which the data is imported. Either specify this or the 'source.query' parameter, but not both.</p> <p>Required: no</p> <p>Supported values: string</p>	<code>--table</code>
<code>source.field.names</code>	<p>The names of columns to import from the source table in a Teradata system, in comma-separated format. The order of the source field names must match exactly the order of the target field names for schema mapping. This parameter must be present when the 'target.field.names' parameter is specified. If not specified, then all columns from the source table will be retrieved.</p> <p>Required: no</p> <p>Supported values: string</p>	<code>--columns</code>
<code>target.table</code>	<p>The name of the target table in Hive or HCatalog. It is required with an 'hcat' job, optional with a 'hive' job, and not valid with an 'hdfs' job. For a 'hive' job, either specify this parameter or the 'target.paths' parameter, but not both.</p> <p>Required: no</p> <p>Supported values: string</p>	<code>--hive-table</code>
<code>target.field.names</code>	<p>The names of fields to write to the target file in HDFS, or to the target Hive or HCatalog table, in comma separated format. The order of the target field names must match exactly the order of the source field names for schema mapping. This parameter must be provided when the 'source.field.names' parameter is specified.</p> <p>Required: no</p> <p>Supported values: string</p>	Driven by the imported columns
<code>batch.size</code>	<p>The number of rows a Hortonworks Connector fetches each time from the Teradata system, up to a 1 MB buffer size limit.</p> <p>Required: no</p>	<code>--fetch-size</code>

Connector Import Option (<code>teradata.db.input.*</code>)	Description	Overriding Sqoop Option
	Supported values: an integer greater than 0 Default value: 10000	
separator	The field separator to use with the imported files. This parameter is only applicable with the 'textfile' file format. Required: no Supported values: string Default value: <code>\t</code>	<code>--fields-terminated-by</code>
split.by.column	The name of a table column to be used for splitting import tasks. It is optional with the 'split.by.hash' and 'split.by.value' methods, and not valid with the 'split.by.partition' method. If this parameter is not specified, the first column of the table's primary key or primary index will be used. Required: no Supported values: a valid table column name	<code>--split-by</code>

Connector Export Options

All option names below are prefixed by "teradata.db.output." when specified in the configuration files or in the `-D` command line option.

For example, `target.table` is specified as `teradata.db.output.target.table`.



Note

Only options with an overriding Sqoop option are listed.

Connector Export Option (<code>teradata.db.output.*</code>)	Description	Overriding Sqoop Option
target.table	The name of the target table in a Teradata system. Required: yes Supported values: string	<code>--table</code>
source.paths	The directory of to-be exported source files in HDFS. It is required for an 'hdfs' job, optional with a 'hive' job, and not valid with an 'hcat' job. For a 'hive' job, either specify this or the 'source.table' parameter but not both. Required: no Supported values: string	<code>--export-dir</code>
num.mappers	The maximum number of output mapper tasks. If the value is zero, then	<code>-m</code> <code>--num-mappers</code>

Connector Export Option (teradata.db.output. *)	Description	Overriding Sqoop Option
	<p>the number of mappers will be the same as the number of file blocks in HDFS. Use either this parameter or 'num.reducers', but not both.</p> <p>Required: no</p> <p>Supported values: an integer greater than or equal to zero</p> <p>Default value: 2</p>	
target.field.names	<p>The names of fields to export to the target table in a Teradata system, in comma-separated format. The order of the target field names must match the order of the source field names for schema mapping. This parameter must be provided when the 'source.field.names' parameter is specified.</p> <p>Required: no</p> <p>Supported values: string</p>	--columns
separator	<p>The separator of fields in the source to-be-exported files. This parameter is only valid with 'textfile' file format.</p> <p>Required: no</p> <p>Supported values: string</p> <p>Default value: \t</p>	--input-fields-terminated-by

Connector-specific Extra Arguments

The Hortonworks connector for Teradata has the following connector-specific extra arguments:

Type of Argument	Argument	Description
Common Options		
	jobtype	The job type: hdfs, hive or hcat.
	fileformat	File format: sequencefile, textfile, avrofile, orcfile or refile. Default is textfile.
	usexview	Use X views for metadata queries. (X views take security into consideration.)
	stagedatabase	Database to use for creating stage tables.
	stagetablename	Stage table name to use; if blank, a default name is generated.
	batchsize	Fetch size or insert batch size.
	queryband	Query band for the session.
Import-specific Options		
	numpartitionshating	Number of partitions to be created in the staging table.
	method	One of split.by. {value hash partition amp}
	accesslock	Row lock is used for fetching rows.
	avroschemafile	Avro schema file for Avro imports.

Type of Argument	Argument	Description
	targettableschem	Schema for the partitioning columns. Needed when Hive table is to be created.
	targetpartitionschem	Schema for the partitioning columns. Needed when Hive table is to be created.
	targetfieldnames	Field names for the target fields. Needed when Hive table is to be created.
Export Options		
	sourcetableschem	Schema for the source hive table.
	sourcepartitionschem	Schema for the partitioning columns.
	sourcefieldnames	Field names for the source fields to export.
	fastloadsockethost	Host for Fastload exports.
	fastloadsocketport	Port for the Fastload exports.
	fastloadsockettimeout	Timeout for the Fastload export operation.
	errortablename	Error table name for use with Fastload.
	keepstagetable	Keep stage table after export. (If not present, stage table is dropped after export.)
	forcestage	Force creation of a stage table.