

Querying Data in Unified Analytics

Date published: 2024-01-01

Date modified: 2024-08-15



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

Unified Analytics basics.....	5
Setting up an Impala Virtual Warehouse.....	5
Setting up a Hive Virtual Warehouse.....	6
Running a query in Cloudera Data Warehouse.....	7
Granting permissions to run SQL queries.....	7
Using DDL statements.....	8
Escaping an invalid identifier.....	8
Client access to Unified Analytics.....	9
Connecting a JDBC client.....	9
Connecting an Impala shell client.....	10
Using the Impala shell client.....	10
Materialized views.....	11
Creating and using a materialized view.....	12
Creating the tables and view.....	12
Verifying the query rewrite.....	13
Query rewrite example.....	14
Managing query rewrites.....	15
Using optimizations from a subquery.....	16
Dropping a materialized view.....	17
Showing materialized views.....	17
Describing a materialized view.....	17
Periodically rebuilding a materialized view.....	19
Incremental rebuilds.....	20
Purposely using a stale materialized view.....	21
Creating and using a partitioned materialized view.....	22
Materialized view commands.....	24
ALTER MATERIALIZED VIEW REBUILD.....	24
ALTER MATERIALIZED VIEW REWRITE.....	24
CREATE MATERIALIZED VIEW.....	24
DESCRIBE EXTENDED and DESCRIBE FORMATTED.....	26
DROP MATERIALIZED VIEW.....	27
SHOW MATERIALIZED VIEWS.....	28
Materialized view recommendations introduction.....	28
Using the recommender in local mode.....	29
Setting up a mini HiveServer.....	29
Getting wi script help.....	31
Running Recommender examples.....	33
Speeding up queries using BI mode.....	42
Get the JDBC driver.....	43
Connect to Unified Analytics with a BI tool.....	43
Enable BI mode to rewrite queries automatically.....	46

Using roll ups and grouping sets.....	46
Using set operations.....	47
Using anti joins.....	48
Creating a table from Parquet data.....	50
Query correlated data.....	50
Subquery restrictions.....	51
Comparing tables using ANY/SOME/ALL.....	51
Handling ETL jobs.....	52
Inserting data into a table.....	52
Updating data in a table.....	53
Merging data in tables.....	53
Deleting data from a table.....	54

Unified Analytics basics

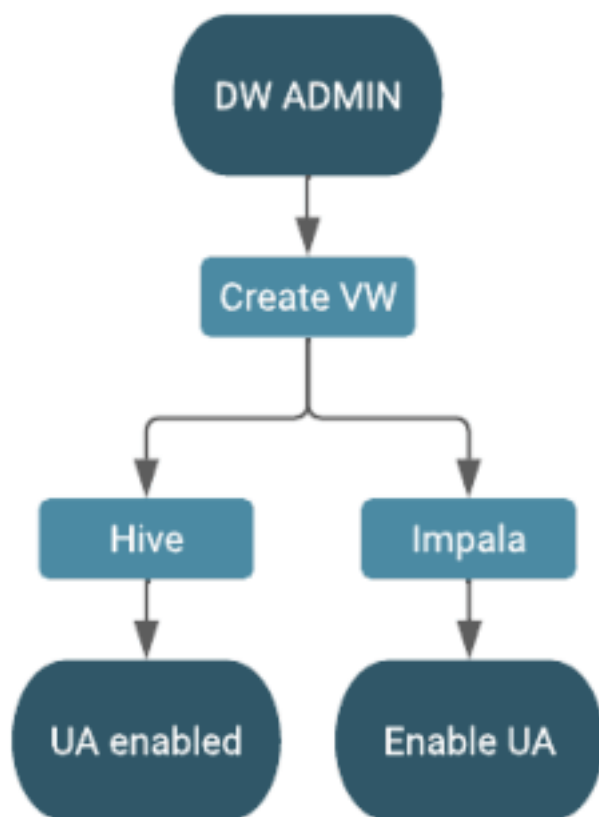
A brief description of Unified Analytics engines and query editors prepares you to run a query. You learn how to select the SQL engine that best suits your needs.

Before using Unified Analytics, [understand the concepts](#) behind it.

In Cloudera Data Warehouse (CDW), Unified Analytics integrates essential warehouse software:

- SQL engines: Hive and Impala
- Query editors: Hue and the Impala shell

As DW Admin, you use a one- or two-step procedure to create a Virtual Warehouse (VW) for querying your data, as shown in the following diagram:



When you create a Hive Virtual Warehouse, Unified Analytics is enabled. When you create an Impala Virtual Warehouse, you need to perform an extra step to enable Unified Analytics.

Related Information

[Activating AWS environments](#)

[Activating Azure environments](#)

[CDP CLI: Create a Database Catalog create-db](#)

Setting up an Impala Virtual Warehouse

You learn how to set up a Unified Analytics Virtual Warehouse for Business Intelligence (BI), or ad-hoc querying. An Impala Virtual Warehouse can read full ACID tables in the Hive metastore, and can write insert-only ACID tables because you select the Impala engine type.

Procedure

1. Log in to the Data Warehouse service as DWAdmin.
2. Go to the **Virtual Warehouses** tab and click New Virtual Warehouse.
The **Create Virtual Warehouse** modal is displayed.
3. Specify a name for your Virtual Warehouse and select an environment, database catalog, Impala, and Hue image versions.
4. Select Impala as the Virtual Warehouse Type.
5. Select a size for the Virtual Warehouse.
6. Select Unified Analytics.
7. Deselect Disable AutoSuspend if you intend to use Impala coordinator shutdown, which you enable in the next step.
The Impala coordinator does not automatically shutdown unless the Impala executors are suspended.
8. Configure Impala coordinator shutdown by turning on Allow Shutdown Of Coordinator and accept the default value Enabled Active-Passive in High Availability (HA).
For the benefits and other information about Coordinator shutdown, see [configure Impala coordinator shutdown](#).
9. Set User Groups that can access endpoints, keys and values for Tagging the Virtual Warehouse.
10. Customize other options as necessary.
11. Click Create Virtual Warehouse.

Related Information

[CDP CLI: Create a Virtual Warehouse create-vw](#)

[Creating a new Virtual Warehouse](#)

Setting up a Hive Virtual Warehouse

You learn how to set up a Hive Unified Analytics Virtual Warehouse. If you have extract, transform, and load (ETL) jobs, you need to learn how to enable query isolation mode.

About this task

When you create a Hive Virtual Warehouse, Unified Analytics is automatically enabled. This type of Virtual Warehouse can read and write full ACID tables in the Hive metastore. ETL jobs, identified by their size, run in query isolation mode if enabled.

Procedure

1. Log in to the Data Warehouse service as DWAdmin.
2. Go to the **Virtual Warehouses** tab and click New Virtual Warehouse.
The **Create Virtual Warehouse** modal is displayed.
3. Specify a name for your Virtual Warehouse and select an environment, database catalog, Impala, and Hue image versions.
4. Select Hive as the Virtual Warehouse Type.
5. Select a size for the Virtual Warehouse.
6. Select the Enable Query Isolation option.
This action sets the internal configuration property `hive.etl.execution.engine` to Tez for query isolation.
7. Set User Groups that can access endpoints, keys and values for Tagging the Virtual Warehouse.
8. Click Create Virtual Warehouse.

Related Information

[CDP CLI: Create a Virtual Warehouse create-vw](#)

[Creating a new Virtual Warehouse](#)

Running a query in Cloudera Data Warehouse

You simply open Hue and submit your query. You do not need to manually start beeline or any other shell.

About this task

As a DWUser, you open Hue from a Virtual Warehouse that you set up, and run the query. The SQL engine reads from and writes to the same metastore, regardless of the type Virtual Warehouse.

Before you begin

- Required role: DW User
- You obtained permission to run SQL queries from the EnvironmentAdmin, who added you to a Hadoop SQL policy.

Procedure

1. Log into the Data Warehouse service as DWUser.
2. Go to the **Virtual Warehouses** tab, locate the Virtual Warehouse using which you want to run queries, and click HUE.


The Hue query editor opens in a new browser tab.

3. To run a query:

- a) Click a database to view the tables it contains.

When you click a database, it sets it as the target of your query in the main query editor panel.

- b)

Type a query in the editor panel and click  to run the query.

Related Information


[Submitting queries with Hue](#)

[Get permission to run SQL queries](#)

Granting permissions to run SQL queries

You must grant Hadoop SQL Ranger permissions to your users so they can query tables.

Procedure

1. Log in to the Data Warehouse service as DWAdmin.
2. Go to the Database Catalog tab, locate your Database Catalog, and click  Open Ranger .
The Ranger application opens in a new browser tab.
3. Click Hadoop SQL on the **Service Manager** page.
4. Click Add New Policy.
5. Specify the databases, tables, and columns on which you want to grant permissions to the users.
6. Under the **Allow Conditions** section, specify users or groups to whom you want to grant permissions.
you can specify {USER} to grant permissions to all users.

7. Click Add Permissions and select permissions such as Create, Alter, Drop, Select, and so on.

8. Click Add to save the settings.

Using DDL statements

Unified Analytics supports DDL statements to create a database schema and to define the type and structure of the data added to a database.

The following list contains the DDL statements that are currently supported through the Unified Analytics engine.

- Create/Drop/Alter/Use Database
- Create/Drop/Truncate Table
- Create Table Like Parquet
- Alter Table/Partition/Column
- Create/Drop/Alter View
- Create/Drop/Alter Materialized View
- Create/Drop/Alter Scheduled Query
- Create/Drop/Grant/Revoke Roles and Privileges
- Show
- Describe
- Analyze Compute / Drop Statistics
- Refresh Table



Note: For more information about these statements, refer to Hive documentation. Be sure to use Hive syntax when you run these statements on either Unified Analytics SQL engines.

Escaping an invalid identifier

When you need to use reserved words, special characters, or a space in a column or partition name, enclose it in backticks (`).

About this task

An identifier in SQL is a sequence of alphanumeric and underscore (_) characters enclosed in backtick characters. In Hive, these identifiers are called quoted identifiers and are case-insensitive. You can use the identifier instead of a column or table partition name.

Before you begin

You have set the following parameter to column in the hive-site.xml file to enable quoted identifiers:

Set the hive.support.quoted.identifiers configuration parameter to column in the hive-site.xml file to enable quoted identifiers in column names. Valid values are none and column. For example, in Hive execute the following command: `SET hive.support.quoted.identifiers = column.`

Procedure

1. Create a table named test that has two columns of strings specified by quoted identifiers:
`CREATE TABLE test (`x+y` String, `a?b` String);`
2. Create a table that defines a partition using a quoted identifier and a region number:
`CREATE TABLE partition_date-1 (key string, value string) PARTITIONED BY (`dt+x` date, region int);`
3. Create a table that defines clustering using a quoted identifier:
`CREATE TABLE bucket_test(`key?1` string, value string) CLUSTERED BY (`key?1`) into 5 buckets;`

Client access to Unified Analytics

You learn which interfaces Unified Analytics supports for connecting a client. Procedures for using the JDBC and Impala shell interfaces include step-by-step instructions.

To connect to Unified Analytics you can download the JDBC driver from the Unified Analytics to give to your client. Also, provide the Unified Analytics endpoint to your client to establish the connection.

You can connect and submit requests using the following interfaces:

- JDBC/ODBC
- Beeline
- Impala shell
- Impyla Python client
- Hue web-based user interface

Connecting a JDBC client


You see how to establish the connection between your JDBC client and Unified Analytics.

About this task

In this task, you connect to a Unified Analytics Virtual Warehouse. First, you download the JDBC driver to give to your client. Next, you obtain the Unified Analytics endpoint for the connection. Finally, you provide necessary information to your client.

Procedure

1. Log in to the Data Warehouse service as DWAdmin.
2. Go to the Resource and Downloads tile and click See More.
3. Click Download JDBC Jar or Download JDBC/ODBC Driver, depending on your Virtual Warehouse SQL engine type.
This action downloads the JAR (driver) file to your local system.

4. Provide the JAR file you downloaded to your JDBC client. For most of the clients, the JAR file should be added under the Libraries folder. Refer to your client documentation for information about the correct location.
5. Go to the Virtual Warehouses tab, locate the Virtual Warehouse you to connect the JDBC client to, and click  Copy JDBC URL .
6. Paste the copied JDBC URL into a text file. The URL is in the following format:

```
jdbc:hive2://<your-virtual-warehouse>.<your-environment>.<your-domain.co
m>/default;transportMode=http;httpPath=cliservice;ssl=true;retries=3
```

7. Provide the Unified Analytics JDBC URL to your client.

```
<your-virtual-warehouse>.<your-environment>.<domain.company.com>
```

8. On the JDBC client end, check that the user can set the authentication credentials:

Authentication: Username and Password

Username: Username you use to connect to the CDP Data Warehouse service.

Password: Password you use to connect to the CDP Data Warehouse service.



Connecting an Impala shell client

You need to provide the commands to your client users for installing and launching the Impala shell to connect to your Unified Analytics. Client users can then query your tables. You learn how to obtain the command for installing the Impala shell on a client and other information to provide to clients.

About this task

To access Unified Analytics from the Impala Shell, clients need to connect to HiveServer (HS2), which is available in the same CDP cluster as Impala. The JDBC URL for the HS2 endpoint uses the strict HS2 protocol option to access Unified Analytics from the Impala Shell. Windows clients are not supported.

Procedure

1. Log in to the Data Warehouse service as DWAdmin.
2. Go to the Virtual Warehouses tab, locate the Unified Analytics Impala Virtual Warehouse to which you want to connect the Impala shell client, and click  Copy Unified Analytics Impala shell Download command
3. Click  again, and select Copy Unified Analytics Impala shell command.

This action copies the command that launches the Impala shell and includes the connection string your client needs to connect to Unified Analytics in your Virtual Warehouse. For example:

```
impala-shell --protocol='hs2-http' --ssl -i 'coordinator-vw-impala.dw-dwx-
rzs556.xcu2-8y8x.dev.cldr.work:443' -u client_max -l
```

4. Provide the commands you copied to your client user.
5. Provide the instructions to your client user to use as described in *Using the Impala shell client*.

Using the Impala shell client

You obtain the commands from the Impala Virtual Warehouse owner for installing and launching the Impala shell to connect to the Impala Virtual Warehouse. You learn how to install and launch the Impala shell that connects to the Impala Virtual Warehouse UI. You then query tables in the Impala Virtual Warehouse, assuming you have authentication credentials. Windows clients are not supported.

Before you begin

- You must have the latest stable version of Python 2.7.
- You must obtain the pip installer compatible with the Python version.

Procedure

1. In a terminal window, run the command provided by the Impala Virtual Warehouse owner to update impyla to the version compatible with the Impala Virtual Warehouse.

For example:

```
pip install impyla==0.18a2
```

Installing/updating to impyla 0.18a2, for example, is required before you install the Impala shell in the next step.

2. Run the Impala shell Download command provided by the Impala Virtual Warehouse owner to install the Impala shell.

The command looks something like this:

```
pip install impala-shell==4.1.0a1
```

The Impala shell is installed.

3. Run the Help command to confirm a successful installation:

```
impala-shell --help
```

4. Run the provided Impala shell command to launch the Impala shell.

The command looks something like this:

```
impala-shell --protocol='hs2-http' --ssl -i 'coordinator-vw-impala.dw-dwx-rzs556.xcu2-8y8x.dev.cldr.work:443' -u client_max -l
```

The Impala shell is launched and you connect to the Impala Virtual Warehouse.

5. Run a SQL command to confirm that you are connected to the Impala Virtual Warehouse instance.

```
SHOW DATABASES;
```

The query returns a list of databases in the Impala Virtual Warehouse.

Materialized views

A materialized view is a Hive-managed database object that can be queried by Hive or Impala users. The materialized view holds a query result you can use to speed up the execution of a query workload. If your queries are repetitive, you can reduce latency and resource consumption by using materialized views. You create materialized views to optimize your queries automatically.

Using a materialized view, the optimizer can compare old and new tables, rewrite queries to accelerate processing, and manage maintenance of the materialized view when data updates occur. The optimizer can use a materialized view to fully or partially rewrite projections, filters, joins, and aggregations.

When you enable BI mode to use Apache DataSketches approximations, a materialized view automatically calls DataSketches for some types of operations.

You can perform the following materialized view operations:

- Create a materialized view of queries or subqueries
- Drop a materialized view

- Show materialized views
- Describe a materialized view
- Enable or disable query rewriting based on a materialized view
- Globally enable or disable rewriting based on any materialized view
- Use partitioning to improve the performance of materialized views

Creating and using a materialized view

You can create a materialized view of a query to calculate and store results of an expensive operation, such as a particular join, on a managed, ACID table that you repeatedly run. When you issue queries specified by that materialized view, the optimizer rewrites the query based on it. This action saves reprocessing. Query performance improves.

About this task

In the tasks that follow, first you create, or use an existing, Hive Virtual Warehouse. You create and populate example tables. The tables are managed tables. You cannot create a materialized view of an external table. You create a materialized view of a join of the tables. Subsequently, you run a query to join the tables, and the query plan takes advantage of the precomputed join to accelerate processing. These over-simplified tasks show the syntax and output of a materialized view, and do not demonstrate accelerated processing that occurs in a real-world task, processing a large amount of data.

Before you begin

- You have access to an existing Hive Virtual Warehouse, or you created a new one.
- You have a client connection, such as Beeline, to Hive in CDP.

Creating the tables and view

You see how to create simple tables, insert the data, and join the tables using a materialized view. You run the query, and the optimizer takes advantage of the precomputation performed by the materialized view to speed response time.

Procedure

1. Log into the Data Warehouse service as DWUser.
2. Go to the **Virtual Warehouses** tab, locate the Virtual Warehouse using which you want to run queries, and click HUE.

The Hue query editor opens in a new browser tab.

3. Select a database, for example default.
4. Create a managed table.

For example, create this emps table in the Hue editor:

```
CREATE TABLE emps (  
  empid INT,  
  deptno INT,  
  name VARCHAR(256),  
  salary FLOAT,  
  hire_date TIMESTAMP);
```

5. Create another table.

For example, create this depts table:

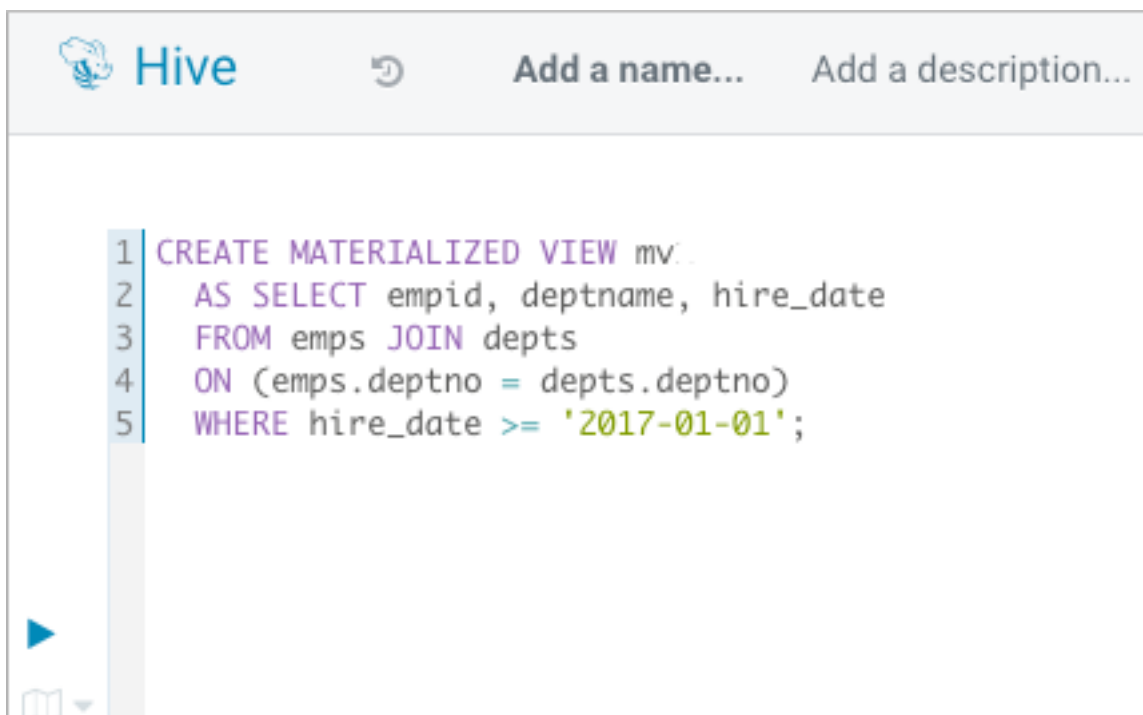
```
CREATE TABLE depts (  
  deptno INT,  
  deptname VARCHAR(256),  
  locationid INT);
```

6. Insert data into the tables by copy/pasting the following statements one-by-one and executing them one-by-one.

```
INSERT INTO TABLE emps VALUES (10001,101,CAST('jane doe' as VARCHAR(256)),250000,'2018-01-10');
INSERT INTO TABLE emps VALUES (10002,100,CAST('somporn klailee' as VARCHAR(256)),210000,'2017-12-25');
INSERT INTO TABLE emps VALUES (10003,200,CAST('jeiranan thongnopneua' as VARCHAR(256)),175000,'2018-05-05');
INSERT INTO TABLE depts VALUES (100,CAST('HR' as VARCHAR(256)),10);
INSERT INTO TABLE depts VALUES (101,CAST('Eng' as VARCHAR(256)),11);
INSERT INTO TABLE depts VALUES (200,CAST('Sup' as VARCHAR(256)),20);
```

7. Create a materialized view to join the tables:

```
CREATE MATERIALIZED VIEW mv
AS SELECT empid, deptname, hire_date
FROM emps JOIN depts
ON (emps.deptno = depts.deptno)
WHERE hire_date >= '2017-01-01';
```



Verifying the query rewrite

You can check that a materialized view is used as a query rewrite.

Procedure

1. Log into the Data Warehouse service as DWUser.
2. Go to the **Virtual Warehouses** tab, locate the Virtual Warehouse using which you want to run queries, and click HUE.

The Hue query editor opens in a new browser tab.

3. In Hue, verify that the query rewrite used the materialized view by running an extended EXPLAIN statement:

```
EXPLAIN EXTENDED SELECT empid, deptname
FROM emps
JOIN depts
```

```
ON (emps.deptno = depts.deptno)
WHERE hire_date >= '2017-01-01'
AND hire_date <= '2019-01-01';
```

4. Export results.

The output shows the alias default.mv1 for the materialized view in the Scan section of the plan.

```
Explain
OPTIMIZED SQL: SELECT `empid`, `deptname`
FROM `default`.`mv`
WHERE '2019-01-01' >= `hire_date`
STAGE DEPENDENCIES:
Stage-1 is a root stage
Stage-0 is a root stage
STAGE PLANS:
Stage: Stage-1
Impala
Impala Plan:
F01:PLAN FRAGMENT [UNPARTITIONED] hosts=1 instances=1
| Per-Instance Resources: mem-estimate=4.02MB mem-reservation=4.00MB thread-reservation=1
PLAN-ROOT SINK
| output exprs: default.mv.empid, default.mv.deptname
| mem-estimate=4.00MB mem-reservation=4.00MB spill-buffer=2.00MB thread-reservation=0
|
01:EXCHANGE [UNPARTITIONED]
| mem-estimate=16.00KB mem-reservation=0B thread-reservation=0
| tuple-ids=0 row-size=32B cardinality=15
| in pipelines: 00(GETNEXT)
|
F00:PLAN FRAGMENT [RANDOM] hosts=1 instances=1
Per-Instance Resources: mem-estimate=16.00MB mem-reservation=24.00KB thread-reservation=1
00:SCAN S3 [default.mv default.mv1, RANDOM]
S3 partitions=1/1 files=1 size=964B
predicates: casttotimestamp('2019-01-01') >= default.mv1.hire_date
stored statistics:
table: rows=unavailable size=964B
columns: unavailable
extrapolated-rows=disabled max-scan-range-rows=unavailable
parquet dictionary predicates: casttotimestamp('2019-01-01') >= default.mv1.hire_date
mem-estimate=16.00MB mem-reservation=24.00KB thread-reservation=0
tuple-ids=0 row-size=32B cardinality=15
in pipelines: 00(GETNEXT)
Stage: Stage-0
Fetch Operator
limit: -1
Processor Tree:
ListSink
```

Query rewrite example

You examine a materialized view definition and see how an example query that matches a materialized view is rewritten transparently.

The following materialized view definition uses dynamic partitioning by year:

```
CREATE MATERIALIZED VIEW mat_view
PARTITIONED ON (d_year)
AS SELECT d_year, d_moy, d_dom, SUM(ss_sales_price) AS sum_sales
FROM store_sales, date_dim
```

```
WHERE ss_sold_date_sk = d_date_sk AND d_year > 2017  
GROUP BY d_year, d_moy, d_dom;
```

With this materialized view in place, when the following query runs, transparent rewriting to speed up the execution occurs.

```
SELECT SUM(ss_sales_price) AS sum_sales  
FROM store_sales, date_dim  
WHERE ss_sold_date_sk = d_date_sk AND d_year = 2018 AND d_mon IN (1,2,3);
```

The query matches most of the constructs, but not the filter, in the materialized view. Calcite recognizes this query as a full, mapped query, minus the predicate. Calcite rewrites the query as a select from the materialized view. Calcite also recognizes that the materialized view is partitioned, and does partition elimination:

```
SELECT SUM(sum_sales) AS sum_sales  
FROM mat_view  
WHERE d_year = 2018 AND d_mon IN (1,2,3);
```

The final rewrite selects from the partition, and applies the predicate to the eliminated partition.

```
SELECT SUM(sum_sales) AS sum_sales  
FROM mat_view(d_year = 2018)  
WHERE d_mon IN (1,2,3);
```

Managing query rewrites

After changes to base tables, the data in a materialized view is stale. You need to know how to prevent the SQL optimizer from rewriting queries in this situation. If you want a query executed without regard to a materialized view, for example to measure the execution time difference, you can disable rewriting and then enable it again.

About this task

As administrator, you can globally enable or disable all query rewrites based on materialized views. By default, the optimizer rewrites a query based on a materialized view.

DWAdmin

Procedure

1. Disable rewriting of a query based on a materialized view named mv1 in the default database.

```
ALTER MATERIALIZED VIEW default.mv1 DISABLE REWRITE;
```

2. Enable rewriting of a query based on materialized view mv1.

```
ALTER MATERIALIZED VIEW default.mv1 ENABLE REWRITE;
```

3. Globally enable rewriting of queries based on materialized views by setting a global property.

```
SET hive.materializedview.rewriting=true;
```

Purposely using a stale materialized view

About this task

When materialized view data becomes stale, or when the optimizer cannot determine the data freshness, which is the case if you use external tables, you can use the materialized view by performing these tasks:

- Schedule the materialized view for rebuilding. For example, schedule a rebuild to occur every x minutes.

- Adjust the rewriting time window to use stale or possibly stale data for a period of time. For example, schedule the window within which to use stale data for $x + y$ minutes.

Procedure

- Create a scheduled query to invoke the rebuild statement every 10 minutes.

```
CREATE SCHEDULED QUERY scheduled_rebuild
EVERY 10 MINUTES AS
ALTER MATERIALIZED VIEW mv_recently_hired REBUILD;
```

- Define the window of time for using stale data.

```
SET hive.materializedview.rewriting.time.window=10min;
```

Using optimizations from a subquery

You can create a query having a subquery that the optimizer rewrites based on a materialized view. You create a materialized view, and then run a query that uses the materialized view.

About this task

In this task, you create a materialized view and use it in a subquery to return the number of destination-origin pairs. Suppose the data resides in a table named `flights_data` that has the following columns:

c_id	dest	origin
1	Chicago	Hyderabad
2	London	Moscow
...		

Procedure

- Create a table schema definition named `flights_data` for destination and origin data.

```
CREATE TABLE flights_data(
  c_id INT,
  dest VARCHAR(256),
  origin VARCHAR(256));
```

- Create a materialized view that counts destinations and origins.

```
CREATE MATERIALIZED VIEW mv1
AS
SELECT dest, origin, count(*)
FROM flights_data
GROUP BY dest, origin;
```

- Take advantage of the materialized view to speed your queries when you have to count destinations and origins again.

For example, use a subquery to select the number of destination-origin pairs like the materialized view.

```
SELECT count(*)/2
FROM(
  SELECT dest, origin, count(*)
  FROM flights_data
  GROUP BY dest, origin
```



```
) AS t;
```

Transparently, the SQL engine uses the work already in place since creation of the materialized view instead of reprocessing.

Dropping a materialized view

You must understand when to drop a materialized view to successfully drop related tables.

About this task

Drop a materialized view before performing a DROP TABLE operation on a related table. You cannot drop a table that has a relationship with a materialized view.

In this task, you drop a materialized view named mv1 from the database named default.

Procedure

Drop a materialized view in my_database named mv1 .
DROP MATERIALIZED VIEW default.mv;

Showing materialized views

You can list all materialized views in the current database or in another database. You can filter a list of materialized views in a specified database using regular expression wildcards.

About this task

You can use regular expression wildcards to filter the list of materialized views you want to see. The following wildcards are supported:

- Asterisk (*)
Represents one or more characters.
- Pipe symbol (|)
Represents a choice.

For example, mv_q* and *mv|q1* match the materialized view mv_q1. Finding no match does not cause an error.

Procedure

1. List materialized views in the current database.
SHOW MATERIALIZED VIEWS;
2. List materialized views in a particular database.
SHOW MATERIALIZED VIEWS IN another_database;

Describing a materialized view

You can get summary, detailed, and formatted information about a materialized view.

About this task

This task builds on the task that creates a materialized view named mv1.

Procedure

1. Get summary information about the materialized view named mv1.

```
DESCRIBE mv1;
```

	Query History	Saved Queries	Results (3)
	col_name	data_type	comment
1	dest	varchar(256)	
2	origin	varchar(256)	
3	_c2	bigint	

2. Get detailed information about the materialized view named mv1.

```
DESCRIBE EXTENDED mv1;
```

	Query History	Saved Queries	Results (5)
	col_name	data_type	
1	dest	varchar(256)	
2	origin	varchar(256)	
3	_c2	bigint	
	Export results	NULL	
5	Detailed Table Information	Table(tableName:mv1, dbName:default	

3. Click Export Results to copy the information.

4. Get formatting details about the materialized view named mv1.

```
DESCRIBE FORMATTED mv1;
```

For example, 14 of the 38 actual results are shown below:

Query History		Saved Queries	Results (38)
	col_name		data_type
1	dest		varchar(256)
2	origin		varchar(256)
3	_c2		bigint
4			NULL
5	# Detailed Table Information		NULL
6	Database:		default
7	OwnerType:		USER
8	Owner:		csso_khahn
9	CreateTime:		Sat May 22 23:57:47 UTC 2021
10	LastAccessTime:		UNKNOWN
11	Retention:		0
12	Location:		s3a://qe-s3-bucket-weekly-ctx4-
13	Table Type:		MATERIALIZED_VIEW
14	Table Parameters:		NULL

Periodically rebuilding a materialized view

Using materialized views can enhance query performance. You need to update materialized view contents when new data is added to the underlying table. Instead of rebuilding the materialized view manually, you can schedule this task. Automatic rebuilding then occurs periodically and transparently to users.

About this task

This task assumes you created the following schemas for storing employee and departmental information:

```
CREATE TABLE emps (
empid INTEGER,
deptno INTEGER,
name VARCHAR(256),
salary FLOAT,
hire_date TIMESTAMP);

CREATE TABLE depts (
deptno INTEGER,
deptname VARCHAR(256),
locationid INTEGER);
```

Imagine that you add data for a number of employees to the table. Assume many users of your database issue queries to access to data about the employees hired during last year including the department they belong to.

You perform the steps below to create a materialized view of the table to address these queries. Imagine new employees are hired and you add their records to the table. These changes render the materialized view contents outdated. You need to update its contents. You create a scheduled query to perform this task. The scheduled rebuilding will not occur unless there are changes to the input tables. You test the scheduled query by bypassing the schedule and executing the schedule immediately. Finally, you change the schedule to rebuild less often.

Procedure

1. To handle many queries to access recently hired employee and departmental data, create a materialized view.

```
CREATE MATERIALIZED VIEW mv_recently_hired AS
  SELECT empid, name, deptname, hire_date FROM emps
  JOIN depts ON (emps.deptno = depts.deptno)
  WHERE hire_date >= '2020-01-01 00:00:00';
```

2. Use the materialized view by querying the employee data.

```
SELECT empid, name FROM emps
JOIN depts ON (emps.deptno = depts.deptno)
WHERE hire_date >= '2020-03-01 00:00:00' AND deptname = 'finance';
```

3. Assuming new hiring occurred and you added new records to the emps table, rebuild the materialized view.

```
ALTER MATERIALIZED VIEW mv_recently_hired REBUILD;
```

The rebuilding updates the contents of the materialized view.

4. Create a scheduled query to invoke the rebuild statement every 10 minutes.

```
CREATE SCHEDULED QUERY scheduled_rebuild
EVERY 10 MINUTES AS
ALTER MATERIALIZED VIEW mv_recently_hired REBUILD;
```

A rebuild executes every 10 minutes, assuming changes to the emp table occur within that period. If a materialized view can be rebuilt incrementally, the scheduled rebuild does not occur unless there are changes to the input tables.

5. To test the schedule, run a scheduled query immediately.

```
ALTER SCHEDULED QUERY scheduled_rebuild EXECUTE;
```

6. Change the frequency of the rebuilding.

```
ALTER SCHEDULED QUERY scheduled_rebuild EVERY 20 MINUTES;
```

Incremental rebuilds

The materialized view is a physical manifestation and summary of the original data in your table. If the source table, or multiple tables in the case of joins, changes, you likely want updates to the materialized view.

An update can occur as an incremental rebuild, or a full rebuild, of the materialized view. A full rebuild can be expensive. An incremental rebuild updates only the affected parts of the materialized view. An incremental rebuild improves the efficiency of maintaining the materialized view by decreasing rebuild step execution time. The incremental rebuild also relies on existing ACID guarantees of view tables.

The following types of incremental rebuilds occur automatically:

- Record-based

The delta is generated based on a snapshot of source tables taken at view creation and at rebuild time. Source tables must be full ACID (transactional) tables, which only Hive can write.

- Partition-based

The partitions that need rebuilding are identified, and then only those partitions are updated by an insert-overwrite. Source tables can be insert-only ACID (transactional) tables, which Hive and Impala can write.

The materialized view must be transactional to qualify for an incremental rebuild. If the materialized view is not transactional, a full rebuild occurs.

Conditions for a record-based incremental rebuild

A record-based incremental rebuild occurs under the following conditions:

- The materialized view includes a count, sum, min, or max aggregation, and one of the source tables has an additional insert after the view is created.
- The materialized view includes a sum aggregation, and one of the source tables has an additional insert or delete after the view is created.
- You are using a Unified Analytics Virtual Warehouse with a Hive engine, which is the only Virtual Warehouse type that supports the incremental, records-based rebuild.
- The materialized view is a full ACID table.

Conditions for a partition-based incremental rebuild

If compaction did not run on the source table after the last rebuild, a partitions-based incremental rebuild occurs under the following conditions:

- The source table is an insert-only transactional tables.
- The materialized view is an insert-only transactional table.

If the materialized view is a full acid, a partition-based incremental rebuild is not attempted. Instead, a record-based incremental rebuild is attempted.

- The materialized view definition includes any supported aggregation.
- You are using a Unified Analytics Virtual Warehouse with the Hive engine.

Purposely using a stale materialized view

When materialized view data becomes stale, or when the optimizer cannot determine the data freshness, which is the case if you use external tables, you might want to use the materialized view anyway:

About this task

You can use the materialized view by performing these tasks:

- Schedule the materialized view for rebuilding. For example, schedule a rebuild to occur every x minutes.
- Adjust the rewriting time window to use stale or possibly stale data for a period of time. For example, schedule the window within which to use stale data for x + y minutes.

Procedure

1. Create a scheduled query to invoke the rebuild statement every 10 minutes.

```
CREATE SCHEDULED QUERY scheduled_rebuild  
EVERY 10 MINUTES AS  
ALTER MATERIALIZED VIEW mv_recently_hired REBUILD;
```

2. Define the window of time for using stale data, 13 minutes.

```
SET hive.materializedview.rewriting.time.window=13min;
```

Creating and using a partitioned materialized view

When creating a materialized view, you can partition selected columns to improve performance. Partitioning separates the view of a table into parts, which often improves query rewrites of partition-wise joins of materialized views with tables or other materialized views.

About this task

This task assumes you created a materialized view of the emps and depts tables and assumes you created these tables. The emps table contains the following data:

empid	deptno	name	salary	hire_date
10001	101	jane doe	250000	2018-01-10
10005	100	somporn klailee	210000	2017-12-25
10006	200	jeiranan thongnopneua	175000	2018-05-05

The depts table contains the following data:

deptno	deptname	locationid
100	HR	10
101	Eng	11
200	Sup	20

In this task, you create two materialized views: one partitions data on department; the other partitions data on hire date. You select data, filtered by department, from the original table, not from either one of the materialized views. The explain plan shows that Hive rewrites your query for efficiency to select data from the materialized view that partitions data by department. In this task, you also see the effects of rebuilding a materialized view.

Procedure

1. Create a materialized view of the emps table that partitions data into departments.

```
CREATE MATERIALIZED VIEW partition_mv_1 PARTITIONED ON (deptno)
AS SELECT hire_date, deptno FROM emps WHERE deptno > 100 AND deptno < 200;
```

2. Create a second materialized view that partitions the data on the hire date instead of the department number.

```
CREATE MATERIALIZED VIEW partition_mv_2 PARTITIONED ON (hire_date)
AS SELECT deptno, hire_date FROM emps where deptno > 100 AND deptno < 200;
```

3. Generate an extended explain plan by selecting data for department 101 directly from the emps table without using the materialized view.

```
EXPLAIN EXTENDED SELECT deptno, hire_date FROM emps_a where deptno = 101;
```

The explain plan shows that Hive rewrites your query for efficiency, using the better of the two materialized views for the job: partition_mv_1.

```
+-----+
| Explain |
+-----+
| OPTIMIZED SQL: SELECT CAST(101 AS INTEGER) AS `deptno`, `hire_date` |
```

```
FROM `default`.`partition_mv_1`
WHERE 101 = `deptno`
STAGE DEPENDENCIES:
  Stage-0 is a root stage
...
```

- Correct Jane Doe's hire date to February 12, 2018, rebuild one of the materialized views, but not the other, and compare contents of both materialized views.

```
INSERT INTO emps VALUES (10001,101,'jane doe',250000,'2018-02-12');
ALTER MATERIALIZED VIEW partition_mv_1 REBUILD;
SELECT * FROM partition_mv_1 where deptno = 101;
SELECT * FROM partition_mv_2 where deptno = 101;
```

The output of selecting the rebuilt partition_mv_1 includes the original row and newly inserted row because INSERT does not perform in-place updates (overwrites).

partition_mv_1.hire_date	partition_mv_1.deptno
2018-01-10 00:00:00.0	101
2018-02-12 00:00:00.0	101

The output from the other partition is stale because you did not rebuild it:

partition_mv_2.deptno	partition_mv_2.hire_date
101	2018-01-10 00:00:00.0

- Create a second employees table and a materialized view of the tables joined on the department number.

```
CREATE TABLE emps2 TBLPROPERTIES AS SELECT * FROM emps;
CREATE MATERIALIZED VIEW partition_mv_3 PARTITIONED ON (deptno) AS
  SELECT emps.hire_date, emps.deptno FROM emps, emps2
  WHERE emps.deptno = emps2.deptno
  AND emps.deptno > 100 AND emps.deptno < 200;
```

- Generate an explain plan that joins tables emps and emps2 on department number using a query that omits the partitioned materialized view.

```
EXPLAIN EXTENDED SELECT emps.hire_date, emps.deptno FROM emps, emps2
  WHERE emps.deptno = emps2.deptno
  AND emps.deptno > 100 AND emps.deptno < 200;
```

The output shows that Hive rewrites the query to use the partitioned materialized view partition_mv_3 even though your query omitted the materialized view.

- Verify that the partition_mv_3 sets up the partition for deptno=101 for partition_mv_3.

```
SHOW PARTITIONS partition_mv_3;
```

Output is:

partition
deptno=101

Materialized view commands

The syntax and examples of materialized view commands provide details for using materialized views in Unified Analytics.

ALTER MATERIALIZED VIEW REBUILD

You must rebuild the materialized view to keep it up-to-date when changes to the data occur.

Syntax

```
ALTER MATERIALIZED VIEW [db_name.]materialized_view_name REBUILD;
```

db_name.materialized_view_name

The database name followed by the name of the materialized view in dot notation.

Description

A rewrite of a query based on a stale materialized view does not occur automatically. If you want a rewrite of a stale or possibly stale materialized view, you can force a rewrite. For example, you might want to use the contents of a materialized view of a non-transactional table because the freshness of such a table is unknown. To enable rewriting of a query based on a stale materialized view, you can run the rebuild operation periodically and set the following property: `hive.materializedview.rewriting.time.window`. For example, `SET hive.materializedview.rewriting.time.window=10min;`

Example

```
ALTER MATERIALIZED VIEW mydb.mv1 REBUILD;
```

ALTER MATERIALIZED VIEW REWRITE

You can change the behavior of Hive to enable or disable the rewriting of queries based on a particular materialized view.

Syntax

```
ALTER MATERIALIZED VIEW [db_name.]materialized_view_name ENABLE|DISABLE REWRITE;
```

db_name.materialized_view_name

The database name followed by the name for the materialized view in dot notation.

Description

To optimize performance, by default, Hive rewrites a query based on materialized views. You can change this behavior to manage query planning and execution manually. By setting the `hive.materializedview.rewriting` global property, you can manage query rewriting based on materialized views for all queries.

Example

```
ALTER MATERIALIZED VIEW mydb.mv1 DISABLE REWRITE;
```

CREATE MATERIALIZED VIEW

If you are familiar with the `CREATE TABLE AS SELECT (CTAS)` statement, you can quickly master how to create a materialized view.

Syntax

```
CREATE MATERIALIZED VIEW [IF NOT EXISTS] [db_name.]materialized_view_name
[DISABLE REWRITE]
[COMMENT materialized_view_comment]
[PARTITIONED ON (column_name, ...)]
[
  [ROW FORMAT row_format]
  [STORED AS file_format]
  | STORED BY 'storage.handler.class.name' [WITH SERDEPROPERTIES (ser
de_property_name=serde_property_value, ...)]
]
[LOCATION hdfs_path]
[TBLPROPERTIES (tbl_property_name=tbl_property_value, ...)]
AS
<query>;
```

Required Parameters

query

The query to execute for results that populate the contents of the materialized view.

Optional Parameters

db_name.materialized_view_name

The database name followed by a name, unique among materialized view names, for the materialized view. The name must conform to Apache Hive specifications for a table name, including case-insensitive alphanumeric and underscore characters.

materialized_view_comment

A string literal enclosed in single quotation marks.

column_name

A key that determines how to do the partitioning, which divides the view of the table into parts.

'storage.handler.class.name'

The name of a storage handler, such as `org.apache.hadoop.hive.druid.DruidStorageHandler`, that conforms to the Apache Hive specifications for storage handlers in a table definition that uses the `STORED BY` clause. When not specified, Hive uses the default `hive.materializedview.fileformat`.

serde_property_name

A property supported by `SERDEPROPERTIES` that you specify as part of the `STORED BY` clause. The property is passed to the `serde` provided by the storage handler. When not specified, Hive uses the default `hive.materializedview.serde`.

serde_property_value

A value of the `SERDEPROPERTIES` property.

hdfs_path

The location on the file system for storing the materialized view.

tbl_property_name

A key that conforms to the Apache Hive specification for `TBLPROPERTIES` keys in a table.

tbl_property_value

The value of a `TBLPROPERTIES` key.

Usage

The materialized view creation statement meets the criteria of being atomic: it does not return incomplete results. By default, the optimizer uses materialized views to rewrite the query. You can store a materialized view in an

external storage system using the `STORED AS` clause followed by a valid storage handler class name. You can set the `DISABLE REWRITE` option to alter automatic rewriting of the query at materialized view creation time.

Examples

```
CREATE MATERIALIZED VIEW druid_t
  STORED AS 'org.apache.hadoop.hive.druid.DruidStorageHandler'
  AS SELECT a, b, c
  FROM src;
```

```
CREATE MATERIALIZED VIEW mv4
  LOCATION '/user/csso_max'
  AS SELECT empid, deptname, hire_date
  FROM emps JOIN depts
  ON (emps.deptno = depts.deptno)
  WHERE hire_date >= '2017-01-01';
```

DESCRIBE EXTENDED and DESCRIBE FORMATTED

You can get extensive formatted and unformatted information about a materialized view.

Syntax

```
DESCRIBE [EXTENDED | FORMATTED] [db_name.]materialized_view_name;
```

db_name

The database name.

materialized_view_name

The name of the materialized view.

Examples

Get summary, details, and formatted information about the materialized view in the default database and its partitions.

```
DESCRIBE FORMATTED default.partition_mv_1;
```

Example output is:

col_name	data_type	comment
# col_name	data_type	comment
name	varchar(256)	
	NULL	NULL
# Partition Information	NULL	NULL
# col_name	data_type	comment
deptno	int	
	NULL	NULL
# Detailed Table Information	NULL	NULL
Database:	default	NULL
OwnerType:	USER	NULL
Owner:	hive	NULL
CreateTime:	Wed Aug 22 19:46:08 UTC 2018	NULL

col_name	data_type	comment
LastAccessTime:	UNKNOWN	NULL
Retention:	0	NULL
Location:	hdfs://myserver:8020/warehouse/tables pace/managed/hive/partition_mv_1	NULL
Table Type:	MATERIALIZED_VIEW	NULL
Table Parameters:	NULL	NULL
	COLUMN_STATS_ACCURATE	{\"BASIC_STATS\":true\"}
	bucketing_version	2
	numFiles	2
	numPartitions	2
	numRows	4
	rawDataSize	380
	totalSize	585
	transient_lastDdlTime	1534967168
	NULL	NULL
# Storage Information	NULL	NULL
SerDe Library:	org.apache.hadoop.hive ql.io.orc.OrcSerde	NULL
InputFormat:	org.apache.hadoop.hive ql.io.orc.OrcInputFormat	NULL
OutputFormat:	org.apache.hadoop.hive ql.io.orc.OrcOutputFormat	NULL
Compressed:	No	NULL
Num Buckets:	-1	NULL
Bucket Columns:	[]	NULL
Sort Columns:	[]	NULL
	NULL	NULL
# Materialized View Information	NULL	NULL
Original Query:	SELECT hire_date, deptno FROM emps WHERE deptno > 100 AND deptno < 200	NULL
Expanded Query:	SELECT `hire_date`, `deptno` FROM (SELECT `emps`.`hire_date`, `emps`.`deptno` FROM `default`.`emps` WHERE `emps`.`deptno` > 100 AND `emps`.`deptno` < 200) `default.partition_mv_1`	NULL
Rewrite Enabled:	Yes	NULL
Outdated for Rewriting:	No	NULL

DROP MATERIALIZED VIEW

You can avoid making a table name unusable by dropping a dependent materialized view before dropping a table.

Syntax

```
DROP MATERIALIZED VIEW [db_name.]materialized_view_name;
```

db_name.materialized_view_name

The database name followed by a name for the materialized view in dot notation.

Description

Dropping a table that is used by a materialized view is not allowed and prevents you from creating another table of the same name. You must drop the materialized view before dropping the tables.

Example

```
DROP MATERIALIZED VIEW mydb.mv1;
```

SHOW MATERIALIZED VIEWS

You can list all materialized views in the current database or in another database. You can filter a list of materialized views in a specified database using regular expression wildcards.

Syntax

```
SHOW MATERIALIZED VIEWS [IN db_name];
```

db_name

The database name.

'identifier_with_wildcards'

The name of the materialized view or a regular expression consisting of part of the name plus wildcards. The asterisk and pipe (* and |) wildcards are supported. Use single quotation marks to enclose the identifier.

Examples

```
SHOW MATERIALIZED VIEWS;
```

mv_name	rewrite_enabled	mode
# MV Name	Rewriting Enabled	Mode
partition_mv_1	Yes	Manual refresh
partition_mv_2	Yes	Manual refresh
partition_mv_3	Yes	Manual refresh

```
SHOW MATERIALIZED VIEWS '*1';
```

mv_name	rewrite_enabled	mode
# MV Name	Rewriting Enabled	Mode
partition_mv_1	Yes	Manual refresh
	NULL	NULL

Materialized view recommendations introduction

Materialized views can improve query performance greatly, but deciding which query results to materialize is difficult. You can use the materialized view recommender command-line interface to help you decide which query results to materialize.

The recommender design is based on *Agrawal et al. - Automated Selection of Materialized Views and Indexes for SQL Databases (VLDB 2000)*. The recommender takes a set of Hive queries as input and creates materialized views as output by exploring many table subsets and syntactical variants for each view candidate.

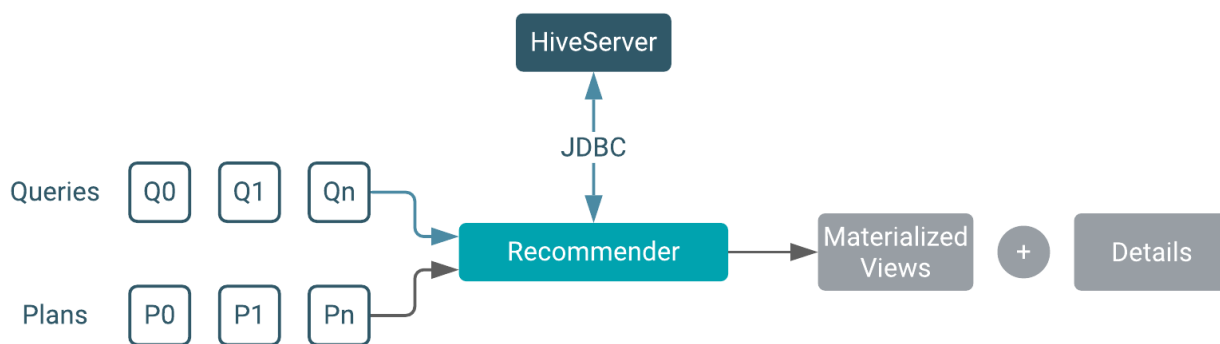
Typically, instead of considering a single, or just a few queries for materializing, you are considering an entire workload. Even the most experienced database experts have difficulty. There are many factors to consider, such as view maintenance and effectiveness. Making a bad decision about what to materialize can make performance worse.

You provide a set of queries or explain plans to the recommender, and then run the Workload Insights (wi) script to get recommendations. You can get recommendations online or offline, depending on the type of input you have:

- Online: input = queries only
- Offline: input = explain plans

If you have explain plans to provide as input to the Workload Insights (wi) script, you can run the script offline; otherwise, you have to provide the HiveServer endpoint for example, `jdbc:hive2://myhost:10000`, to the script. In lieu of explain plans, the recommender needs to gather information about the query, such as the number of rows of each query, how long the query ran, and other statistics, from HiveServer.

The following diagram shows the input and output of the recommender:



How much detail you receive in a recommendation depends on your configuration. Details include the number of recommendations and other information, such as why a particular view is recommended and how to exploit the view.

Workload Insights input requirements

You must have one of the following inputs to the Workload Insights (wi) script:

- One or more directories containing one or more select and aggregation queries in .sql files and multiple queries in files, separated by a semicolon (;).
A .sql file name suffix is required.
- One or more directories containing an explain plan in .json files. A .json file name suffix is required.

Using the recommender in local mode

You can use the recommender on a local computer for learning and testing. You can use a Mac laptop, for example. First, from a Cloudera Data Warehouse (CDW) Virtual Warehouse that enables Unified Analytics, you download software that runs a mini HiveServer in local mode. Next, you provide a set of queries or explain plans to the recommender. Finally, you execute the Workload Insights (wi) script to get pseudo-recommendations.

Setting up a mini HiveServer

You follow a step-by-step procedure to get a docker image of HiveServer (HS2) and example files to simulate working with HiveServer in a cluster. In local mode, you do not get actual recommendations, but you experiment with and understand recommendations using this mode without incurring any cloud costs.


About this task

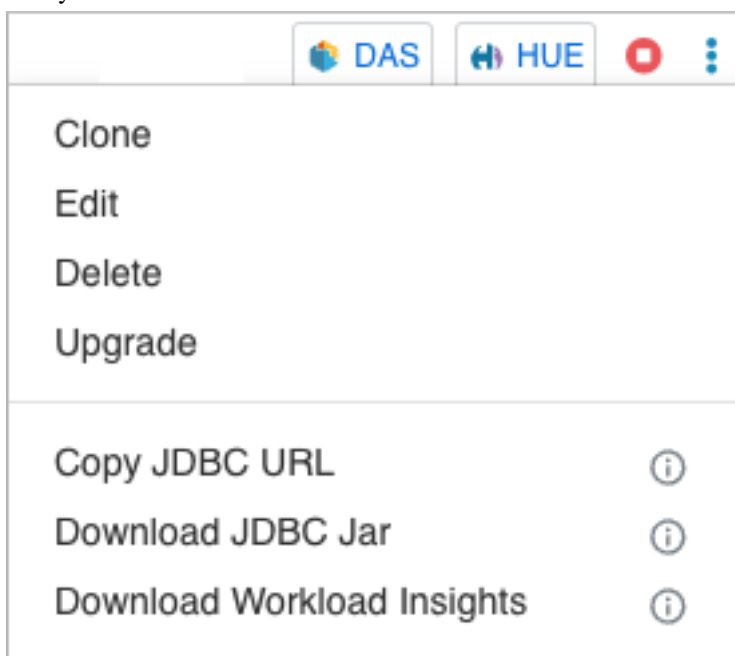
You download the workload-insights archive and unpack it. You run one docker command to start HiveServer and another to load the TPC-DS schema. Example query files, based on the TPC-DS Benchmark DDL queries, are included in the package. No actual data is needed.

Before you begin

- Check that Java 8 or later is installed.
- Install Docker (<https://docs.docker.com/get-docker/>).
- Check that Docker is working: `docker run hello-world`

Procedure

1. On the **Overview** page under Virtual Warehouses, click options  in a Virtual Warehouse that enables Unified Analytics.



2. Click Download Workload Insights to download the workload-insights archive.
3. On your file system, unpack the archive.
4. Set environment variable `WI_HOME="$PWD"`.
For example, on the command line of a terminal:

```
export WI_HOME="$PWD"
```

5. Set environment variable `WI_CLASSPATH`. `WI_CLASSPATH="*:/lib/*:$WI_HOME/target/*:$WI_HOME/target/lib/*"`
6. Start up a HiveServer (HS2) container named mini-hs2 that listens on port 10000.

```
docker run -p "10000:10000" -d --name mini-hs2 zabetak/hs2-embedded:1.0.2.7.2.3.0-220
```

7. Load the TPC-DS schema required for running examples.

```
docker exec mini-hs2 schema-load tpcds
```

Ignore any log4j warning messages.

8. Run the recommender.

Getting wi script help

To get pseudo-recommendations for learning, or testing purposes, you need to know the Workload Insights (wi) script options. You run your first wi script to get online help.

Procedure

- 1. Go to the github directory where you downloaded and built the project.
- 2. Run the wi script with no options to get online help about script usage and options:

```
$ ./wi
```

Output looks something like this:

```
usage: wi -i <file_1,...,file_N> [-o <directory>] [-c <configfile>] [-u
<url>] [-n <username>] [-p
<password>] [-t <threshold>] [-f <format>] [-r <recommendation>] [-l <le
vel>][-d]
```

-i,--input <file_1,...,file_N>	one or more .sql/.json files or a directory (version Y)
-o,--output <directory>	output directory
-c,--config <configfile>	path to configuration .properties) file
-u,--url <url>	JDBC url to connect to
-n,--username <username>	username to connect as
-p,--password <password>	password to connect as
-t,--tables-threshold <threshold>	<div>the threshold controlling above which a table subset is considered interesting if it is estimated to contribute to total workload cost above this threshold (between 0 and 1). The default is 0.1 (10% of total workload cost)</div>

-f,--format <format>	output format TEXT(Default)/JSON
-r,--recommendation <recommendation>	<p>the kind of recommendations to include in the report:</p> <ul style="list-style-type: none"> - SPJ for SELECT-PROJECT-JOIN recommendations - SPJA for SELECT-PROJECT-JOIN-AGG REGATE recommendations - ALL for both SPJ and SPJA recommendations (DEFAULT)
-l,--level <level>	<p>the level of detail to display in the report:</p> <ul style="list-style-type: none"> - MINIMAL displays only the view recommendations - DEFAULT displays recommendations along with basic information about the query workload - FULL displays all information available in the report (useful for debugging purposes)
-d, --delete	<p>suggest which materialized views, not used by a given workload could be deleted. Requires JDBC information to be provided. Available in version Y.</p>

```
usage: wi -i <directory> [-o <directory>] [-c <configfile>] [-u <url>] [-n
<username>] [-p <password>] [-t <threshold>] [-f <format>] [-r <recommend
ation>] [-l <level>]
-i,--input <directory>          input directory with .sql/.json
files
-o,--output <directory>        output directory
-c,--config <configfile>       path to configuration (.properties) file
-u,--url <url>                  JDBC url to connect to
-n,--username <username>       username to connect as
-p,--password <password>       password to connect as
-t,--tables-threshold <threshold> the threshold controlling above wh
ich a table subset is considered interesting if it is estimated to contr
ibute to total workload cost above this threshold (between 0 and 1). The
default is 0.1 (10% of total workload cost)
-f,--format <format>           output format TEXT(Default)/JSON
-r,--recommendation <recommendation> the kind of recommendations to inc
lude in the report:
                                - SPJ for SELECT-PROJECT-JOIN re
                                commendations
                                - SPJA for SELECT-PROJECT-JOIN-AGG
                                REGATE recommendations
                                - ALL for both SPJ and SPJA recom
                                mendations (DEFAULT)
```



```

-l,--level <level>          the level of detail to display in
the report:
                                - MINIMAL displays only the view
recommendations
                                - DEFAULT displays recommendations
                                along with basic information about the query workload
                                - FULL displays all information
                                available in the report (useful for debugging purposes)

```

3. Try the recommender examples.

A text report that has the following information appears:

- View recommendations for the workload.
- Queries in the workload, along with the recommendations for each query (if available), and other information.

4. Stop the mini-hs2.

```
docker stop mini-hs2
```

5. Remove the mini-hs2.

```
docker rm mini-hs2
```

Running Recommender examples

The workload-insights archive includes sample queries and explain plans that demonstrate the recommender. You gain an understanding of the recommender by running examples.

About this task

In local mode, you do not get actual recommendations, but you experiment with and understand recommendations using this mode without incurring any cloud costs. The recommender examples you can run use one or both of the following queries as input:

QueryA

```

SELECT sr_item_sk, sum(sr_net_loss) AS net_loss
FROM store_returns
INNER JOIN reason ON sr_reason_sk = r_reason_sk
WHERE r_reason_desc = 'Found a better price in a store'
GROUP BY sr_item_sk

```

Query A asks these questions:

- Which items were returned because the customer found a better price?
- How much money did the company lose on each returned item?

QueryB

```

SELECT sr_item_sk, sum(sr_net_loss) as net_loss, avg(ss_sales_price) as price
FROM store_sales
INNER JOIN store_returns
  ON ss_item_sk=sr_item_sk AND ss_ticket_number=sr_ticket_number
INNER JOIN reason
  ON sr_reason_sk = r_reason_sk
WHERE r_reason_desc = 'Found a better price in a store'
GROUP BY sr_item_sk
ORDER BY net_loss DESC NULLS LAST, price DESC NULLS LAST
LIMIT 10

```

QueryB asks these questions:

- Which items were returned because the customer found a better price?
- How much money did the company lose on each returned item?
- What was the average selling price of return items?
- What are the top ten items causing the biggest loss?

Before you begin

- The mini-hs2 docker image is running.
- You downloaded and unpacked the workload-insights archive, and it includes the following sample queries that you need to provide as input to the Workload Insights (wi) script:
 - Directories containing select and aggregation queries in .sql files
A .sql file name suffix is required.
 - Directories containing explain plans in .json files
A .json file name suffix is required.

Running the Hello world example

You use the minimum set of required parameters, and get a result that is not interesting. From this example, you learn what to expect in the output when you use the minimum set of required parameters.

About this task

In the input directory workload-insights/demo/queries, a file named queryA.sql contains a single query.

Procedure

1. Go to the workload-insights directory
The Workload Insights (wi) script is located in this directory.
2. Run the script using the -i option to specify file queryA.sql as the input to the recommender.

```
./wi -i queries/queryA.sql -u jdbc:hive2://localhost:10000
```

Output looks something like this:

```
2021-06-10 15:17:37,461 INFO c.c.i.a.m.t.JdbcWorkloadReader Computing plan
  for queryA from JDBC connection
...
Recommended views:
-- materialized_view_1:
CREATE MATERIALIZED VIEW `materialized_view_1` AS
SELECT `store_returns`.`sr_item_sk`, SUM(`store_returns`.`sr_net_loss`) AS
`$f1`
FROM `default`.`reason`, `default`.`store_returns`
WHERE `store_returns`.`sr_reason_sk` = `reason`.`r_reason_sk` AND 'Found a
  better price in a store' = `reason`.`r_reason_desc`
GROUP BY `store_returns`.`sr_item_sk`;
-----
Query details:
-- queryA:SELECT `t0`.`sr_item_sk`, SUM(`t0`.`sr_net_loss`) AS `$f1` FROM
  (SELECT `sr_item_sk`, `sr_reason_sk`, `sr_net_loss` FROM `default`.`sto
re_returns` WHERE `sr_reason_sk` IS NOT NULL) AS `t0` INNER JOIN (SELECT
  `r_reason_sk` FROM `default`.`reason` WHERE `r_reason_desc` = 'Found a
  better price in a store' AND `r_reason_sk` IS NOT NULL) AS `t2` ON `t0`.`
sr_reason_sk` = `t2`.`r_reason_sk` GROUP BY `t0`.`sr_item_sk`
-- RECOMMENDATION : materialized_view_1
```

The result is identical to the input query. The output contains a list of recommendations followed by the queries in the workload. By default, for each query, the materialized views matching the query appear.

Running the exclude aggregations example

If you are knowledgeable about your workload and want to exclude views with aggregations for maintenance reasons, you need to know about the SPJ (select-project-join) recommendation mode.

About this task

In the input directory workload-insights/demo/queries, a file named queryA.sql contains queryA. QueryA contains the GROUP BY aggregation clause, which you want the recommender to exclude from the view..

Procedure

1. Go to the workload-insights directory.

The Workload Insights (wi) script is located in this directory.

2. Run the script using the -r SPJ option to exclude aggregations from recommended views.

```
./wi -i queries/queryA.sql -u jdbc:hive2://localhost:10000 -r SPJ
```

Output looks something like this:

```
...
Recommended views:
-- materialized_view_0:
CREATE MATERIALIZED VIEW `materialized_view_0` AS
SELECT `store_returns`.`sr_item_sk`, `store_returns`.`sr_reason_sk`, `s
tore_returns`.`sr_net_loss`, `store_returns`.`sr_reason_sk` AS `sr_reaso
n_sk0`
FROM `default`.`reason`, `default`.`store_returns`
WHERE `store_returns`.`sr_reason_sk` = `reason`.`r_reason_sk` AND 'Found
a better price in a store' = `reason`.`r_reason_desc`;
-----
Query details:
-- queryA:SELECT `t0`.`sr_item_sk`, SUM(`t0`.`sr_net_loss`) AS `$f1` FROM
(SELECT `sr_item_sk`, `sr_reason_sk`, `sr_net_loss` FROM `default`.`stor
e_returns` WHERE `sr_reason_sk` IS NOT NULL) AS `t0` INNER JOIN (SELECT
`r_reason_sk` FROM `default`.`reason` WHERE `r_reason_desc` = 'Found a b
etter price in a store' AND `r_reason_sk` IS NOT NULL) AS `t2` ON `t0`.`
sr_reason_sk` = `t2`.`r_reason_sk` GROUP BY `t0`.`sr_item_sk`
-- RECOMMENDATION : materialized_view_0
-----
```

The recommender runs faster by excluding aggregations because the recommender explores a smaller search space. Similarly, the user can request to obtain only recommendations with aggregations using the SPJA mode. Using the default ALL option, you would get the same result as above. ALL includes SPJ for SELECT-PROJECT-JOIN recommendations and SPJA for SELECT-PROJECT-JOIN-AGGREGATE.

Running the level of detail examples

You see what the level of detail options, -l MINIMAL, -l DEFAULT, and -l FULL do by running these examples. The recommended materialized view is identical to the input query regardless of the level of detail option you use.

About this task

Changing the level of detail does not affect the number or structure of recommendations but provides more insights for each query/recommendation in the case of FULL, and less in the case of MINIMAL.

In the input directory src/test/resources/demo/queries/scenario02, a file named queryA.sql contains queryA.

Procedure

1. Go to the workload-insights directory.
The Workload Insights (wi) script is located in this directory.
2. Run the script to get a recommendation using the `-l MINIMAL` option.

```
./wi -i queries/queryA.sql -u jdbc:hive2://localhost:10000 -l MINIMAL
```

Only recommendations appear. You can easily copy and paste the CREATE statements for use in another CLI, for example:

```
SUMMARY (END) #####
Recommended views:
CREATE MATERIALIZED VIEW `materialized_view_1` AS
SELECT `store_returns`.`sr_item_sk`, SUM(`store_returns`.`sr_net_loss`) AS
  `$f1`
FROM `default`.`reason`,`default`.`store_returns`
WHERE `store_returns`.`sr_reason_sk` = `reason`.`r_reason_sk` AND 'Found a
  better price in a store' = `reason`.`r_reason_desc`
GROUP BY `store_returns`.`sr_item_sk`;
-----
```

3. Run the script using the `-l FULL` option.

```
./wi -i queries/queryA.sql -u jdbc:hive2://localhost:10000 -l FULL
```

Recommendations are:

```
...
Recommended views:
-- materialized_view_1:
CREATE MATERIALIZED VIEW `materialized_view_1` AS
SELECT `store_returns`.`sr_item_sk`, SUM(`store_returns`.`sr_net_loss`) AS
  `$f1`
FROM `default`.`reason`,`default`.`store_returns`
WHERE `store_returns`.`sr_reason_sk` = `reason`.`r_reason_sk` AND 'Found a
  better price in a store' = `reason`.`r_reason_desc`
GROUP BY `store_returns`.`sr_item_sk`;
-- QUERIES : queryA
-- POTENTIAL : 0.9902040816326532
-----
Total improvement: 0.9902040816326532
-----
```

From the following verbose details in the output, you can understand why a recommendation was selected, or not, for a particular query. You might use these details to troubleshoot problems with the recommender.

```
Query details:
-- queryA:SELECT `t0`.`sr_item_sk`, SUM(`t0`.`sr_net_loss`) AS `$f1` FROM
  (SELECT `sr_item_sk`, `sr_reason_sk`, `sr_net_loss` FROM `default`.`sto
re_returns` WHERE `sr_reason_sk` IS NOT NULL) AS `t0` INNER JOIN (SELECT
  `r_reason_sk` FROM `default`.`reason` WHERE `r_reason_desc` = 'Found a
  better price in a store' AND `r_reason_sk` IS NOT NULL) AS `t2` ON `t0`.`
sr_reason_sk` = `t2`.`r_reason_sk` GROUP BY `t0`.`sr_item_sk`
-- RECOMMENDATION : materialized_view_1
-- PLAN :
  DASAggregate(group=[{0}], agg#0=[sum($2)])
  DASJoin(condition=[($1, $3)], joinType=[inner])
  DASProject(sr_item_sk=$2, sr_reason_sk=$8, sr_net_loss=$19))
  DASFilter(condition=[IS NOT NULL($8)])
```

```

DASTableScan(table=[[default, store_returns]])
DASProject(r_reason_sk=[0])
DASFilter(condition=[AND(=($2, 'Found a better price in a store
'), IS NOT NULL($0))])
DASTableScan(table=[[default, reason]])
-- BT: default.store_returns is a base table
-- Row Count: 1.0|Row size: 208.0
-- BT: default.reason is a base table
-- Row Count: 1.0|Row size: 332.0
-- IMPROVEMENT : 0.9902040816326532
-----
Interesting subsets:
-- [[default, reason], [default, store_returns]]:1.0
  "recommended_views" : [ {
    "name" : "materialized_view_1",
    "sql" : "CREATE MATERIALIZED VIEW `materialized_view_1` AS\nSELECT
`store_returns`.`sr_item_sk`, SUM(`store_returns`.`sr_net_loss`) AS `$f1`\nFROM `default`.`reason`,\n`default`.`store_returns`\nWHERE `store_re
turns`.`sr_reason_sk` = `reason`.`r_reason_sk` AND 'Found a better price in
a store' = `reason`.`r_reason_desc`\nGROUP BY `store_returns`.`sr_item_sk
`;"
  } ],
  "query_details" : [ {
    "name" : "queryA",
    "sql" : "SELECT `t0`.`sr_item_sk`, SUM(`t0`.`sr_net_loss`) AS `$f1`\nF
ROM (SELECT `sr_item_sk`, `sr_reason_sk`, `sr_net_loss`\nFROM `default`.`
store_returns`\nWHERE `sr_reason_sk` IS NOT NULL) AS `t0`\nINNER JOIN (
SELECT `r_reason_sk`\nFROM `default`.`reason`\nWHERE `r_reason_desc` = '
Found a better price in a store' AND `r_reason_sk` IS NOT NULL) AS `t2` ON
`t0`.`sr_reason_sk` = `t2`.`r_reason_sk`\nGROUP BY `t0`.`sr_item_sk`",
    "recommendations" : [ "materialized_view_1" ]
  } ]

```

Running the JSON output example

You see the example JSON output when you run the Workload Insights command with the -f JSON option.

About this task

In the input directory src/test/resources/demo/queries, a file named queryA.sql contains queryA.

Procedure

1. Go to the github directory where you downloaded and built the project.

The Worload Insights (wi) script is located in this directory.

2. Run the script to get a JSON-formatted recommendation.

```
./wi -i queries/queryA.sql -u jdbc:hive2://localhost:10000 -f JSON
```

The output is:

```

{
  "recommended_views" : [ {
    "name" : "materialized_view_1",
    "sql" : "CREATE MATERIALIZED VIEW `materialized_view_1` AS\nSELECT `st
ore_returns`.`sr_item_sk`, SUM(`store_returns`.`sr_net_loss`) AS `$f1`\n
FROM `default`.`reason`,\n`default`.`store_returns`\nWHERE `store_return
s`.`sr_reason_sk` = `reason`.`r_reason_sk` AND 'Found a better price in
a store' = `reason`.`r_reason_desc`\nGROUP BY `store_returns`.`sr_item_s
k`;"
  } ],

```

```

"query_details" : [ {
  "name" : "queryA",
  "sql" : "SELECT `t0`.`sr_item_sk`, SUM(`t0`.`sr_net_loss`) AS `$f1`\n
FROM (SELECT `sr_item_sk`, `sr_reason_sk`, `sr_net_loss`\nFROM `default`\n
.`store_returns`\nWHERE `sr_reason_sk` IS NOT NULL) AS `t0`\nINNER JOIN\n
(SELECT `r_reason_sk`\nFROM `default`.`reason`\nWHERE `r_reason_desc` =\n
'Found a better price in a store' AND `r_reason_sk` IS NOT NULL) AS `t2`\n
ON `t0`.`sr_reason_sk` = `t2`.`r_reason_sk`\nGROUP BY `t0`.`sr_item_sk`",
  "recommendations" : [ "materialized_view_1" ]
} ]

```

Running the multiple query input example

Typically, you want the recommender to base recommendations on multiple queries. You see how to run the wi script on a directory containing multiple queries and how to interact with the recommender.

About this task

Upon completing the procedure below, you see the recommendation you get for queryA is identical to the query itself. The recommendation you get for queryB closely resembles queryA. The recommender does not find any way to consolidate queries into a meaningful recommendation because queries are not similar in any significant way.

Notice that the recommendation for queryB does not contain the ORDER BY and LIMIT clauses. Recommendations cannot contain an ORDER BY and LIMIT. Doing so could significantly reduce the applicability of the view in many queries.

In the input directory demo/queries, one file named queryA.sql contains queryA, and another file named queryB.sql contains queryB.

Procedure

1. Go to the workload-insights directory.

The Workload Insights (wi) script is located in this directory.

2. Run the script on the directory containing multiple .sql files.

```

./wi -i "queries/queryA.sql, queries/queryB.sql" -u jdbc:hive2://localhost:10000

```

Recommended views are:

```

...
Recommended views:
-- materialized_view_1:
CREATE MATERIALIZED VIEW `materialized_view_1` AS
SELECT `store_sales`.`ss_item_sk`, SUM(`store_returns`.`sr_net_loss`) AS `$f1`,
SUM(`store_sales`.`ss_sales_price`) AS `$f2`, COUNT(`store_sales`.`ss_sales_price`) AS `$f3`
FROM `default`.`reason`,
`default`.`store_returns`,
`default`.`store_sales`
WHERE `store_sales`.`ss_ticket_number` = `store_returns`.`sr_ticket_number`
AND `store_returns`.`sr_reason_sk` = `reason`.`r_reason_sk` AND `store_sales`.`ss_item_sk` = `store_returns`.`sr_item_sk`
AND 'Found a better price in a store' = `reason`.`r_reason_desc`
GROUP BY `store_sales`.`ss_item_sk`;
-- materialized_view_4:
CREATE MATERIALIZED VIEW `materialized_view_4` AS
SELECT `store_returns`.`sr_item_sk`, SUM(`store_returns`.`sr_net_loss`) AS `$f1`
FROM `default`.`reason`, `default`.`store_returns`
WHERE `store_returns`.`sr_reason_sk` = `reason`.`r_reason_sk` AND 'Found a better price in a store' = `reason`.`r_reason_desc`

```

```
GROUP BY `store_returns`.`sr_item_sk`;
```

Details in the output are:

Query details:

```
-- queryA:SELECT `t0`.`sr_item_sk`, SUM(`t0`.`sr_net_loss`) AS `$f1` FROM
(SELECT `sr_item_sk`, `sr_reason_sk`, `sr_net_loss` FROM `default`.`sto
re_returns` WHERE `sr_reason_sk` IS NOT NULL) AS `t0` INNER JOIN (SELECT
`r_reason_sk` FROM `default`.`reason` WHERE `r_reason_desc` = 'Found a
better price in a store' AND `r_reason_sk` IS NOT NULL) AS `t2` ON `t0`.`
sr_reason_sk` = `t2`.`r_reason_sk` GROUP BY `t0`.`sr_item_sk`
-- RECOMMENDATION : materialized_view_4
-- queryB:SELECT `t0`.`sr_item_sk` AS `$f0`, SUM(`t0`.`sr_net_loss`) AS `
$f1`, SUM(`t4`.`ss_sales_price`) / COUNT(`t4`.`ss_sales_price`) AS `$f2`
FROM (SELECT `sr_item_sk`, `sr_reason_sk`, `sr_ticket_number`, `sr_net_
loss` FROM `default`.`store_returns` WHERE `sr_item_sk` IS NOT NULL AND
`sr_ticket_number` IS NOT NULL AND `sr_reason_sk` IS NOT NULL) AS `t0` I
NNER JOIN (SELECT `r_reason_sk` FROM `default`.`reason` WHERE `r_reason_
desc` = 'Found a better price in a store' AND `r_reason_sk` IS NOT NULL)
AS `t2` ON `t0`.`sr_reason_sk` = `t2`.`r_reason_sk` INNER JOIN (SELECT
`ss_item_sk`, `ss_ticket_number`, `ss_sales_price` FROM `default`.`store
_sales` WHERE `ss_item_sk` IS NOT NULL AND `ss_ticket_number` IS NOT NULL)
AS `t4` ON `t0`.`sr_item_sk` = `t4`.`ss_item_sk` AND `t0`.`sr_ticket_numbe
r` = `t4`.`ss_ticket_number` GROUP BY `t0`.`sr_item_sk` ORDER BY SUM(`t
0`.`sr_net_loss`) DESC, SUM(`t4`.`ss_sales_price`) / COUNT(`t4`.`ss_sale
s_price`) DESC LIMIT 10
-- RECOMMENDATION : materialized_view_1
```

3. In Hive, run SHOW MATERIALIZED VIEWS to list the recommended materialized views.
4. Run EXPLAIN CBO.

```
beeline -u jdbc:hive2://localhost:10000 -n admin -p mypassword -e "EXPLAIN
CBO `cat <path>/queryB.sql`"
```

5. Take a look at the explain plan for one of the queries.
The output contains the name of the materialized view used in the plan.
6. Run the queryB.sql again and look at the execution time.
Using the materialized view, execution time is likely much shorter.
7. Disable a materialized view n, for example, from being used in rewrites.

```
beeline -u jdbc:hive2://localhost:10000 -n admin -p mypassword -e "ALTER
MATERIALIZED VIEW materialized_view_n DISABLE REWRITE"
```

8. Run EXPLAIN CBO again.
The materialized view does not appear because it has been disabled.

Running the repetitive queries example

You see how the recommender cost model takes repetitive queries into account when generating final recommendations.

About this task

The input directory workload-insights/demo/queries contains 21 .sql files, queryA.sql, queryA1.sql - queryA19.sql, and queryB.sql.

Procedure

1. Go to the workload-insights directory.

The Workload Insights (wi) script is located in this directory.

2. Run the script on the directory that simulates the common, repetitive query situation.

```
./wi -i "queries/queryA.sql, queries/queryAx19.sql, queries/queryB.sql" -u
jdbc:hive2://localhost:10000 -r SPJ
```

Recommendations are:

```
...
Recommended views:
-- merged_59:
CREATE MATERIALIZED VIEW `merged_59` AS
SELECT `store_returns`.`sr_item_sk` AS `$f2`, SUM(`store_returns`.`sr_n
et_loss`) AS `$f1`
FROM `default`.`reason`,`default`.`store_returns`
WHERE `store_returns`.`sr_reason_sk` = `reason`.`r_reason_sk` AND 'Found
a better price in a store' = `reason`.`r_reason_desc`
GROUP BY `store_returns`.`sr_item_sk`;
-- materialized_view_0:
CREATE MATERIALIZED VIEW `materialized_view_0` AS
SELECT `store_returns`.`sr_item_sk`, `store_returns`.`sr_reason_sk`, `s
tore_returns`.`sr_ticket_number`, `store_returns`.`sr_net_loss`, `store_
returns`.`sr_reason_sk` AS `sr_reason_sk0`
FROM `default`.`reason`,`default`.`store_returns`
WHERE `store_returns`.`sr_reason_sk` = `reason`.`r_reason_sk` AND 'Foun
d a better price in a store' = `reason`.`r_reason_desc` AND `store_retur
ns`.`sr_item_sk` IS NOT NULL AND `store_returns`.`sr_ticket_number` IS NOT
NULL;
-----
```

Details in the output look something like this:

```
Query details:
-- queryA:SELECT `t0`.`sr_item_sk`, SUM(`t0`.`sr_net_loss`) AS `$f1` FROM
(SELECT `sr_item_sk`, `sr_reason_sk`, `sr_net_loss` FROM `default`.`sto
re_returns` WHERE `sr_reason_sk` IS NOT NULL) AS `t0` INNER JOIN (SELECT
`r_reason_sk` FROM `default`.`reason` WHERE `r_reason_desc` = 'Found a
better price in a store' AND `r_reason_sk` IS NOT NULL) AS `t2` ON `t0`.`
sr_reason_sk` = `t2`.`r_reason_sk` GROUP BY `t0`.`sr_item_sk`
-- RECOMMENDATION : merged_59
-- queryA1:SELECT `t0`.`sr_item_sk`, SUM(`t0`.`sr_net_loss`) AS `$f1` F
ROM (SELECT `sr_item_sk`, `sr_reason_sk`, `sr_net_loss` FROM `default`.`
store_returns` WHERE `sr_reason_sk` IS NOT NULL) AS `t0` INNER JOIN (SEL
ECT `r_reason_sk` FROM `default`.`reason` WHERE `r_reason_desc` = 'Found a
better price in a store' AND `r_reason_sk` IS NOT NULL) AS `t2` ON `t0`.`
sr_reason_sk` = `t2`.`r_reason_sk` GROUP BY `t0`.`sr_item_sk`

...
-----
```

Cleaning up materialized views

You can delete materialized views, not used by your workload. You provide JDBC information about the queries associated with the materialized views using Workload Insights (wi) script and then use command-line JSON processor jq to delete the unused materialized views.

About this task

You get information about unused materialized views associated with a single or multiple queries, as shown in the following command:

```
./wi -i "demo/queries/all/queryA.sql, demo/queries/all/queryB.sql" -u jdbc:hive2://localhost:10000 -n admin -p mypassword -d
```

In the following procedure, you delete unused materialized views. First, you get and store information about the unused materialized views in a record, next you use jq to read the record of deleted materialized views, and finally you drop the materialized views.

Procedure

1. Delete materialized views associated with query03.sql from being used by your workload and store the deletion record in a file.

```
./wi -i "queries/queryA.sql, queries/queryAx19.sql, queries/queryB.sql" -u jdbc:hive2://localhost:10000 -n admin -p mypassword -f JSON -o out_dir -d
```

2. Read the deletion record.

```
jq '.unused_views.name' < out_dir/DELETE.json
```

Output is:

```
"materialized_view_2"
"materialized_view_9"
```

3. Read the commands for dropping the unused materialized views.

Output is:

```
"DROP MATERIALIZED VIEW materialized_view_2;"
"DROP MATERIALIZED VIEW materialized_view_9;"
```

4. Drop the materialized views using the commands above.

```
beeline -u jdbc:hive2://localhost:10000 -n admin -p mypassword -e "`jq -r '.unused_views.sql' < out_dir/DELETE.json`"
```

Analyzing JSON records

You can use a jq command to extract information about recommendations in JSON. Then, you can rerun the query, take a look at the plan, and see which materialized view was used, and other information.

The jq command that you can use for analyzing JSON records looks something like this:

```
jq -r '.recommended_views | .name, .queries < out_dir/ALL-FULL.json'
```

In this example, output includes the original query submitted by the user. The output also includes the optimized sql that uses the materialized view. Output looks something like this:

```

{
  "name": "query20",
  "sql": "SELECT i_item_id,\n
        sum(cs_ext_sales_price) AS itemre\n
        (PARTITION BY i_class) AS revenueratio\n
        \n AND i_category IN ('Shoes',\n
        ld_date_sk = d_date_sk\n AND d_date BE\n
        GROUP BY i_item_id,\n
        i_item_de\n
        BY i_category,\n
        i_class,\n
  "recommendations": [
    "materialized_view_9"
  ],
  "optimized_sql": "SELECT `i_item_\n
  * 100 / (SUM(`$f5`) OVER (PARTITION BY\n
  \nFROM `materialized_view_9`\nORDER BY\n
  m_desc` NULLS LAST, `$f5` * 100 / (SUM(\n
  UNDEDED FOLLOWING)) NULLS LAST\nLIMIT 100\n
  "plan": "DASSortLimit(sort0=[2],\n
  =[ASC], dir3=[ASC], dir4=[ASC], fetch=[\n
  s=[3], i_current_price=[2], itemreven\n
  Y $3 ROWS BETWEEN UNBOUNDED PRECEDING A\n
  sum($3))\n
  DASJoin(condition=[=(\n

```

This optimized SQL above matches the information in the plan.

Speeding up queries using BI mode

You can use BI mode to automatically rewrite incoming queries to be answered approximately using Apache DataSketches. BI mode can be useful when using data visualization BI tools such as Tableau.

About this task

Queries of very large data sets with a number of distinct values often take too long to return results. If you can accept approximate results, using DataSketch algorithms can save significant time. BI mode integrated with DataSketches accelerates query execution while decreasing resource utilization. BI mode can rewrite a wide range of SQL analytical queries, including COUNT(DISTINCT), CUME_DIST, RANK, and NTILE.

When you enable BI mode to use Apache DataSketches approximations, a materialized view automatically calls DataSketches for some types of operations. You can pre-compute data sketches using materialized views, and then rely on Unified Analytics BI mode to rewrite algorithms to answer queries directly from those materialized views. This accelerates query execution by orders of magnitude without any change to your original queries.

Before you begin

- Choose an authorization model.
- Configure authenticated users for querying Hive through JDBC or ODBC driver. For example, set up a Ranger policy.

Get the JDBC driver

To query, analyze, and visualize data stored in CDW Unified Analytics, you must download the Unified Analytics JDBC driver to your cluster and provide it to your client. You can use the latest version of drivers provided by Cloudera to connect Unified Analytics to Business Intelligence (BI) tools.

Procedure

1. Log in to the Data Warehouse service as DWAdmin.
2. Go to the Resource and Downloads tile, click [See More Unified Analytics JDBC Driver](#). This action downloads the JAR (driver) file to your local system.


Connect to Unified Analytics with a BI tool

After you downloaded the JDBC driver, you must set up a BI tool such as Tableau to connect with the Unified Analytics engine.

Before you begin

Before you can use Tableau with an Impala/Hive Virtual Warehouses, you must populate your Database Catalog with sample data when you create it. You must also create a Impala/Hive Virtual Warehouse and enable Unified Analytics, which is configured to connect to the Database Catalog that is populated with data.

Procedure

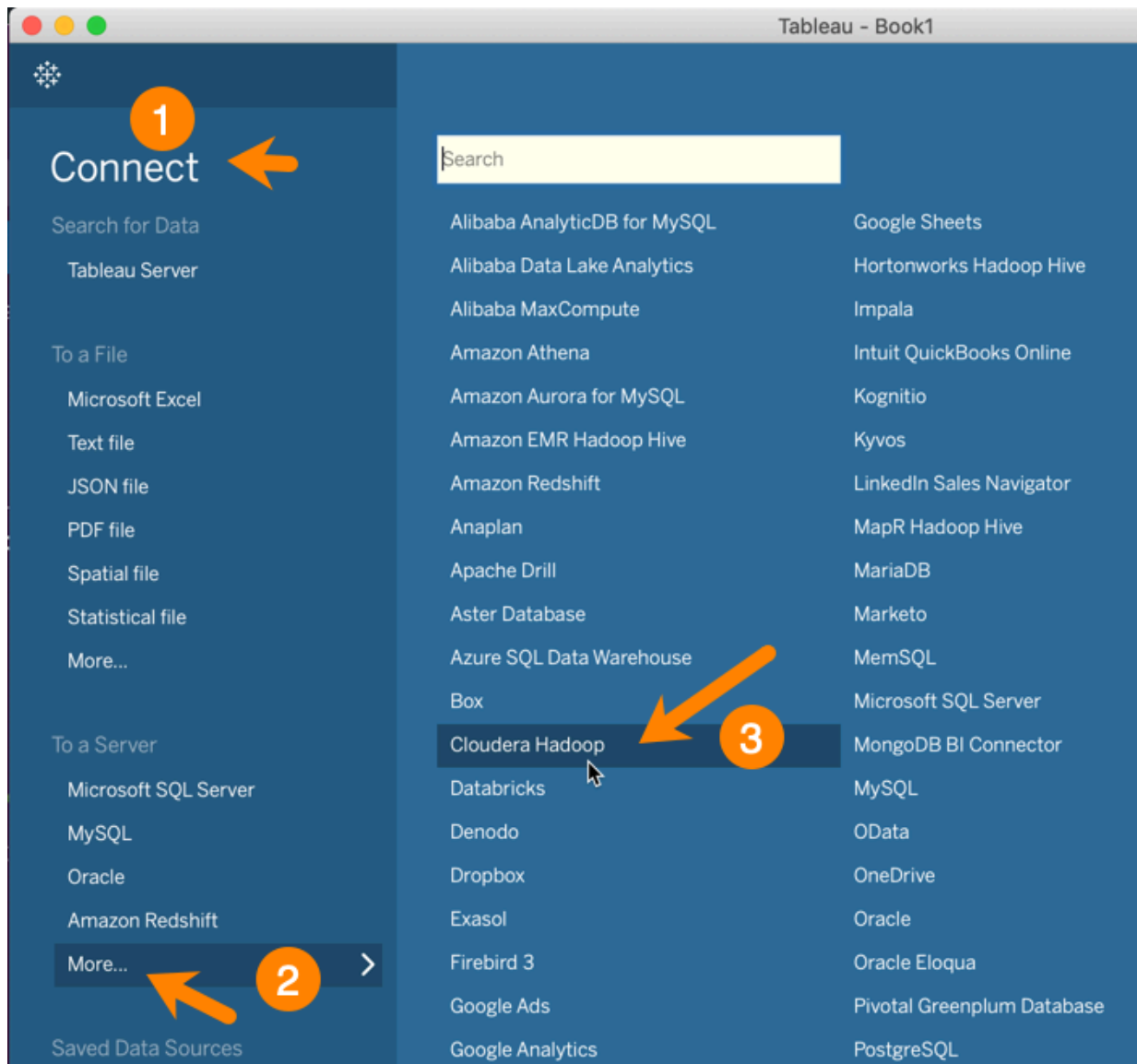
1. Click  on the Impala Virtual Warehouse tile.
2. Select Copy Unified Analytics JDBC URL.
3. Paste the copied JDBC URL into a text file. It should look similar to the following:

```
jdbc:hive2://<your-virtual-warehouse>.<your-environment>.<dwx.company.com>/default;transportMode=http;httpPath=cliservice;ssl=true;retries=3
```

4. From the text file where you just pasted the URL, copy the host name from the JDBC URL to your system's clipboard. For example, in the URL shown in Step 3, the host name is:

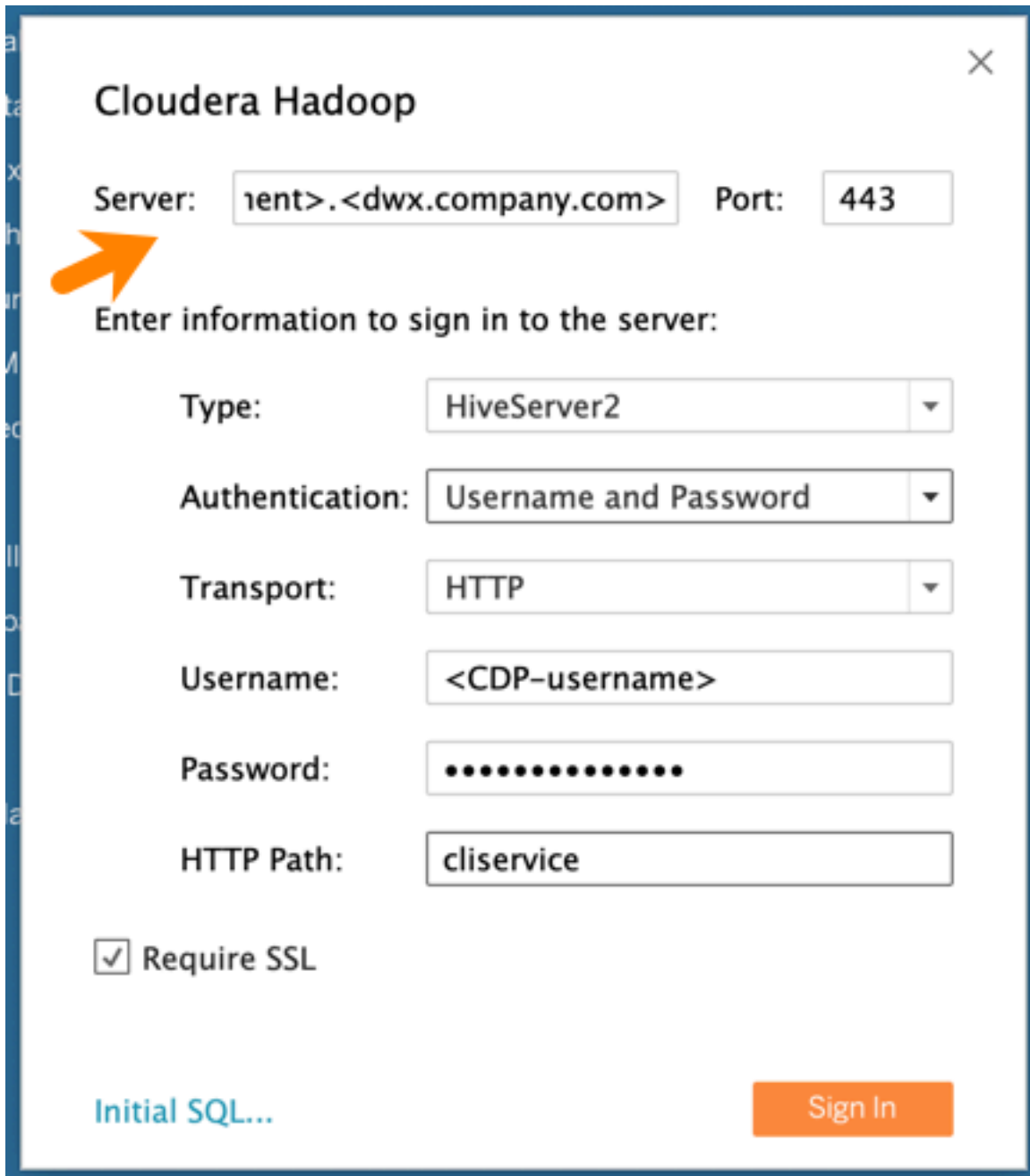
```
<your-virtual-warehouse>.<your-environment>.<yourdomain.com>
```

5. Start Tableau and navigate to ConnectMore...Cloudera Hadoop :



This launches the Cloudera Hadoop dialog box.

6. In the Tableau Cloudera Hadoop dialog box, paste the host name you copied to your clipboard in Step 4 into the Server field:



Cloudera Hadoop

Server: Port:

Enter information to sign in to the server:

Type:

Authentication:

Transport:

Username:

Password:

HTTP Path:

☒ Require SSL

[Initial SQL...](#) [Sign In](#)

7. Then in the Tableau Cloudera Hadoop dialog box, set the following other options:

- Port: 443
- Type: HiveServer2
- Authentication: Username and Password
- Transport: HTTP
- Username: Username you use to connect to the CDP Data Warehouse service.
- Password: Password you use to connect to the CDP Data Warehouse service.
- HTTP Path: cliservice
- Require SSL: Make sure this is checked.

8. Click Sign In.

Enable BI mode to rewrite queries automatically

You must enable BI mode to combine the BI capability with the power of transparent materialized view rewritings.

About this task

If you want your incoming queries, at different granularity levels, to be answered directly from the materialized view rather than being computed from the source tables, you must enable BI mode. BI mode is disabled by default.

Procedure

1. To enable BI mode, execute the following statement from Hue or your preferred client. set hive.optimize.bi.enabled=true;
2. Create a materialized view that stores one of the supported DataSketches (for example, HLL or KLL) on the columns for which the aggregations will be computed during query processing.
The following example shows creating a materialized view for computing the number of unique visitors per region for a given website.

```
CREATE MATERIALIZED VIEW AS
SELECT country, state, city, ds_hll_sketch(userid)
FROM visitors
GROUP BY country, state, city;
```

Results

With the BI mode enabled, any incoming queries will be sliced and diced at different granular levels to be answered directly from the materialized view.

```
SELECT country, count(distinct userid)
FROM visitors
GROUP BY country;

SELECT country, state, count(distinct userid)
FROM visitors
GROUP BY country, state;

SELECT country, state, count(distinct userid)
FROM visitors
WHERE country = 'USA'
GROUP BY country, state;
```

Using roll ups and grouping sets

If your work involves heavy OLAP queries, you need to know about grouping sets and roll ups. Unified Analytics supports grouping sets with rollups and grouping functions.

Before you begin

You created a Unified Analytics Virtual Warehouse based on a Database Catalog that loads the sample data to run queries of airline data.

Procedure

1. Log into the Data Warehouse service as DWUser.

2. Go to the **Virtual Warehouses** tab, locate the Virtual Warehouse using which you want to run queries, and click HUE.

The Hue query editor opens in a new browser tab.

3. Select a database.
4. In Hue, select either the Impala or Hive editor.
5. Run a query that determines the number of flights per country and state.

```
SELECT
  a.country,
  a.state,
  grouping(a.country)
  grouping(a.state)
  SUM(f.Distance),
  COUNT(*)
FROM
  flights f,
  airports a
WHERE
  f.origin = a.iata
GROUP BY
  a.country,
  a.state
WITH ROLLUP;
```

6. Run a query that defines custom grouping sets: country and country/state.

```
SELECT
  a.country,
  a.state,
  grouping(a.country)
  grouping(a.state)
  SUM(f.Distance),
  COUNT(*)
FROM
  flights f,
  airports a
WHERE
  f.origin = a.iata
GROUP BY
  a.country,
  a.state
GROUPING SETS (
  (a.country),
  (a.country, a.state)
);
```

Using set operations

When you want to combine separate query results into a single result, you need to know about the INTERSECT and EXCEPT clauses in addition to UNION.

About this task

In this task, you run a query to find the number of Boeing planes leaving from DFW in 2005. The query filters data using an intersect on the tail number, tailnum, with the unique identification of a plane, planes.manufacturer.

Before you begin

You created a Unified Analytics VW that loads the sample data to run queries of airline data.

Procedure

1. Log into the Data Warehouse service as DWUser.
2. Go to the **Virtual Warehouses** tab, locate the Virtual Warehouse using which you want to run queries, and click HUE.

The Hue query editor opens in a new browser tab.

3. Select a database.
4. In Hue, select either the Impala or Hive editor.
5. Enter a query that intersects the tail number and type of plane.

```
SELECT COUNT(*)
FROM
  SELECT tailnum
  FROM flights
  WHERE
    flights.origin = 'DFW' AND flights.`year`= 2005
  INTERSECT
  SELECT tailnum
  FROM planes
  WHERE
    planes.manufacturer = 'BOEING'
) subq;
```

6. Run a query that returns the planes not manufactured by Boeing leaving from DFW.

```
SELECT COUNT(*)
FROM
  SELECT tailnum
  FROM flights
  WHERE
    flights.origin = 'DFW' AND flights.`year`= 2005
  EXCEPT
  SELECT tailnum
  FROM planes
  WHERE
    planes.manufacturer = 'BOEING'
) subq;
```

Using anti joins


Unified Analytics supports anti join queries to find the rows in a table that do not match rows in another table. The performance of queries involving NULL filters, "not exists" clauses, and "not in" clauses are improved using the anti join optimization.

About this task

In this task, run an anti join query. For example:

```
SELECT * FROM left_table t1 LEFT JOIN right_table t2 ON t1.id = t2.id WHERE
t2.id IS NULL;
```


Procedure

1. Log into the Data Warehouse service as DWUser.
2. Go to the Virtual Warehouses tab, locate the Virtual Warehouse you want to edit, and click  Edit Configurations Hiveserver2 and select hive-site from the Configuration files drop-down menu.
3. Click Add Custom Configuration.
The **Custom Configuration** is displayed.
4. Specify the following values in the Configuration Key: and Configuration Value: fields:
 - Configuration Key: hive.auto.convert.anti.join
 - Configuration Value: true
5. Open Hue, and create two tables, one having more rows than the other.

```
create table tt1 (order_num bigint, price decimal(5, 2));

create table tt2 (order_num bigint);
```

6. Insert several, say five, rows of values (order number, price) in the first table and fewer, say 2, rows of values (order number only) in the second table.

```
insert into tt1 values (42, 300.50), (141, 250.25), (150, 85.00), (150, 500.44), (166, 125.65);
insert into tt2 values (150), (166);
```

7. Take a look at the content of the first table.

```
select * from tt1;
```

```
tt1.order_num  tt1.price
42             300.50
141            250.25
150            85.00
150            500.44
166            125.65
```

8. Take a look at the content of the second table.

```
select * from tt2;
```

```
tt2.order_num
150
166
```

9. Perform an equality left join on the order numbers in the tables.

```
select * from tt1 t1 LEFT JOIN tt2 t2 ON t1.order_num=t2.order_num;
```

```
t1.order_num  t1.price  t2.order_num
42            300.50    NULL
141           250.25    NULL
150           85.00     150
150           500.44     150
166           125.65     166
```

10. Add a null filter to the join query to anti join the tables.

```
select * from tt1 t1 LEFT JOIN tt2 t2 ON t1.order_num=t2.order_num where
t2.order_num is NULL;
```

The anti join returns the rows having NULL in the price column.

```
t1.order_num t1.price t2.order_num
42           300.50  NULL
141          250.25  NULL
```

Creating a table from Parquet data

In Unified Analytics, you can base a new table on a schema in a Parquet file using the CREATE TABLE LIKE FILE PARQUET statement. You are likely familiar with this statement if you have been an Impala user. The Impala counterpart [CREATE TABLE LIKE PARQUET](#) lacks the keyword LIKE, but is functionally similar.

About this task

The column definitions in the new table are inferred from the Parquet data file when you create a table like Parquet in Unified Analytics. Set the following table property for creating the table:

```
hive.parquet.infer.binary.as = <value>
```

Where <value> is binary (the default) or string.

This property determines the interpretation of the unannotated Parquet binary type. Some systems expect binary to be interpreted as string.

Use the following syntax for creating a table like Parquet.

```
CREATE [EXTERNAL] TABLE [IF NOT EXISTS] [db_name.]table_name LIKE PARQUET FI
LE 'object_storage_path_of_parquet_file'
[PARTITIONED BY (col_name data_type [COMMENT 'col_comment'], ...)]
[SORT BY ([column [, column ...])] [COMMENT 'table_comment'] [ROW FORMAT
row_format]
[WITH SERDEPROPERTIES ('key1'='value1', 'key2'='value2', ...)] [STORED AS f
ile_format]
[LOCATION 'object_storage_path']
[TBLPROPERTIES ('key1'='value1', 'key2'='value2', ...)]
```

In this task, you create a table like parquet that uses the default value binary for hive.parquet.infer.binary.as.

Procedure

Create a table based on a parquet table.

```
CREATE TABLE ua_table LIKE FILE PARQUET 's3a://testbucket/files/schema.p
arq';
```

Query correlated data

You can query one table relative to the data in another table.

About this task

A correlated query contains a query predicate with the equals (=) operator. One side of the operator must reference at least one column from the parent query and the other side must reference at least one column from the subquery. An uncorrelated query does not reference any columns in the parent query.

Procedure

Select all state and net_payments values from the transfer_payments table for years during which the value of the state column in the transfer_payments table matches the value of the state column in the us_census table.

```
SELECT state, net_payments
FROM transfer_payments
WHERE EXISTS
  (SELECT year
   FROM us_census
   WHERE transfer_payments.state = us_census.state);
```

This query is correlated because one side of the equals predicate operator in the subquery references the state column in the transfer_payments table in the parent query and the other side of the operator references the state column in the us_census table.

This statement includes a conjunct in the WHERE clause.

A conjunct is equivalent to the AND condition, while a disjunct is the equivalent of the OR condition. The following subquery contains a conjunct:

```
... WHERE transfer_payments.year = "2018" AND us_census.state = "california"
```

The following subquery contains a disjunct:

```
... WHERE transfer_payments.year = "2018" OR us_census.state = "california"
```

Subquery restrictions

To construct queries efficiently, you must understand the restrictions of subqueries in WHERE clauses.

- Subqueries must appear on the right side of an expression.
- Nested subqueries are not supported.
- Subquery predicates must appear as top-level conjuncts.
- Subqueries support four logical operators in query predicates: IN, NOT IN, EXISTS, and NOT EXISTS.
- The IN and NOT IN logical operators may select only one column in a WHERE clause subquery.
- The EXISTS and NOT EXISTS operators must have at least one correlated predicate.
- The left side of a subquery must qualify all references to table columns.
- References to columns in the parent query are allowed only in the WHERE clause of the subquery.
- Subquery predicates that reference a column in a parent query must use the equals (=) predicate operator.
- Subquery predicates may refer only to columns in the parent query.
- Correlated subqueries with an implied GROUP BY statement may return only one row.
- All unqualified references to columns in a subquery must resolve to tables in the subquery.
- Correlated subqueries cannot contain windowing clauses.

Comparing tables using ANY/SOME/ALL

You learn how to use quantified comparison predicates (ANY/SOME/ALL) in non-correlated subqueries according to the SQL standard. SOME is any alias for ANY.

About this task

You can use one of the following operators with a comparison predicate:

- >
- <
- >=
- <=
- <>
- =

ALL:

- If the table is empty, or the comparison is true for every row in subquery table, the predicate is true for that predicand.
- If the comparison is false for at least one row, the predicate is false.

SOME or ANY:

- If the comparison is true for at least one row in the subquery table, the predicate is true for that predicand.
- If the table is empty or the comparison is false for each row in subquery table, the predicate is false.

If the comparison is neither true nor false, the result is undefined.

For example, you run the following query to match any value in c2 of tbl equal to any value in c1 from the same tbl:

```
select c1 from tbl where c1 = ANY (select c2 from tbl);
```

You run the following query to match all values in c1 of tbl not equal to any value in c2 from the same tbl.

```
select c1 from tbl where c1 <> ALL (select c2 from tbl);
```

Handling ETL jobs

Unified Analytics provides DML support for extract, transform, and load tasks. ETL jobs, identified by their size, run in query isolation mode if enabled.

About this task

Unified Analytics supports insert, update, merge, and delete statements for performing ETL tasks.

Procedure

1. Log in to the Data Warehouse service as DWAdmin.
2. Create a Unified Analytics Virtual Warehouse (Hive), selecting the Enable Query Isolation option.
Query isolation is available as a technical preview and requires an entitlement.
3. Click Create Virtual Warehouse.

What to do next

Insert, update, and delete data as described in subsequent topics.

Related Information

[Automating data pipelines with CDE and CDW using Apache Airflow](#)

Inserting data into a table

To insert data into a table you use a familiar ANSI SQL statement. A simple example shows you have to accomplish this basic task.

About this task

To insert data into an ACID table, use the Optimized Row Columnar (ORC) storage format. To insert data into a non-ACID table, you can use other supported formats. You can specify partitioning as shown in the following syntax:

```
INSERT INTO TABLE tablename [PARTITION (partcol1=val1, partcol2=val2 ...)] VALUES values_row [, values_row...]
```

where

values_row is (value [, value]) .

A value can be NULL or any SQL literal.

Procedure

1. Create an ACID table to contain student information.
CREATE TABLE students (name VARCHAR(64), age INT, gpa DECIMAL(3,2));
2. Insert name, age, and gpa values for a few students into the table.
INSERT INTO TABLE students VALUES ('fred flintstone', 35, 1.28), ('barney rubble', 32, 2.32);
3. Create a table called pageviews and assign null values to columns you do not want to assign a value.

```
CREATE TABLE pageviews (userid VARCHAR(64), link STRING, from STRING) PARTITIONED BY (datestamp STRING) CLUSTERED BY (userid) INTO 256 BUCKETS;
INSERT INTO TABLE pageviews PARTITION (datestamp = '2014-09-23') VALUES ('jsmith', 'mail.com', 'sports.com'), ('jdoe', 'mail.com', null);
INSERT INTO TABLE pageviews PARTITION (datestamp) VALUES ('tjohnson', 'sports.com', 'finance.com', '2014-09-23'), ('tlee', 'finance.com', null, '2014-09-21');
```

The ACID data resides in the warehouse.

Updating data in a table

The syntax describes the UPDATE statement you use to modify data already stored in a table. An example shows how to apply the syntax.

About this task

You construct an UPDATE statement using the following syntax:

```
UPDATE tablename SET column = value [, column = value ...] [WHERE expression];
```

Depending on the condition specified in the optional WHERE clause, an UPDATE statement might affect every row in a table. The expression in the WHERE clause must be an expression supported by a SELECT clause. Subqueries are not allowed on the right side of the SET statement. Partition columns cannot be updated.

Before you begin

You must have SELECT and UPDATE privileges to use the UPDATE statement.

Procedure

Create a statement that changes the values in the name column of all rows where the gpa column has the value of 1.0.
UPDATE students SET name = null WHERE gpa <= 1.0;

Merging data in tables

A sample statement shows how you can conditionally insert existing data in Hive tables using the ACID MERGE statement. Additional merge operations are mentioned.

About this task

The MERGE statement is based on ANSI-standard SQL.

Procedure

1. Construct a query to update the customers' names and states in customer table to match the names and states of customers having the same IDs in the new_customer_stage table.
2. Enhance the query to insert data from new_customer_stage table into the customer table if none already exists. Update or delete data using MERGE in a similar manner.

```
MERGE INTO customer USING (SELECT * FROM new_customer_stage) sub ON sub.id
= customer.id
WHEN MATCHED THEN UPDATE SET name = sub.name, state = sub.new_state
WHEN NOT MATCHED THEN INSERT VALUES (sub.id, sub.name, sub.state);
```

Deleting data from a table

You use the DELETE statement to delete data already written to an ACID table.

About this task

Use the following syntax to delete data from a Hive table. DELETE FROM tablename [WHERE expression];

Procedure

Delete any rows of data from the students table if the gpa column has a value of 1 or 0.
DELETE FROM students WHERE gpa <= 1,0;