

Ingesting Data into CDW using Iceberg Table Format

Date published: 2019-12-16

Date modified: 2024-04-03

The Cloudera logo is displayed in a bold, orange, sans-serif font. The word "CLOUDERA" is written in all caps, with the letter 'E' in "CLouDERA" featuring a unique design where the top bar is a horizontal line that extends to the right and then turns down to form the top of the letter 'R'.

Legal Notice

© Cloudera Inc. 2024. All rights reserved.

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

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

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

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

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

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

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

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

Contents

Ingesting data into CDW using Iceberg table format.....	4
Understand the use case.....	4
Meet the prerequisites.....	4
Create the Iceberg target table in CDW.....	4
Build the data flow.....	4
Create and configure controller services for your data flow.....	6
Configure the processor for your data source.....	7
Configure the processor for your data target.....	8
Start the data flow.....	9

Ingesting data into CDW using Iceberg table format

You can use an Apache NiFi data flow to ingest data in Iceberg table format into Cloudera Data Warehouse (CDW) in CDP Public Cloud by following these steps.

Understand the use case

Learn how to use NiFi to move data to an Iceberg table of a data warehouse in CDP Public Cloud.

Cloudera Data Warehouse (CDW) can be configured to use Iceberg as the table format for querying data with Impala and/or Hive. You can use Apache NiFi to move data from a range of locations into a Data Warehouse cluster running Apache Hive or Impala in CDP Public Cloud.

This use case walks you through creating a data flow that generates FlowFiles with random CSV data and writes this data into an Iceberg table in Hive or Impala. This gets you started with creating an Iceberg ingest data flow. If you want to use a different data source, you can check out the other ingest documents for get and consume data processor options.

For more information on Iceberg table format, see *Apache Iceberg features*.

Related Information

[Apache Iceberg features](#)

Meet the prerequisites

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

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

Create the Iceberg target table in CDW

Learn how to create an Iceberg table for your data flow.

About this task

Before you can ingest data into your Cloudera Data Warehouse (CDW) target in CDP Public Cloud, you need to create the table where you want to send data with NiFi. The following example shows how to create a table with Impala:

```
CREATE TABLE customer (id int, name string, created_at timestamp)
PARTITIONED BY (country_code string)
STORED BY ICEBERG;
```

Build the data flow

Learn how to create an ingest data flow to move data to an Iceberg table in Hive or Impala.

About this task

Building a data flow to move data to an Iceberg table in Hive or Impala involves using NiFi in your Flow Management cluster. The example data flow generates FlowFiles with random CSV data and writes this data into an Iceberg table.

Procedure

1. Open NiFi in Data Hub.
 - a) To access the NiFi service in your Flow Management Data Hub cluster, navigate to Management Console service Data Hub Clusters .
 - b) Click the tile representing the Flow Management Data Hub cluster you want to work with.
 - c) Click the NiFi icon in the Services section of the cluster overview page to access the NiFi UI.
You are logged into NiFi automatically with your CDP credentials.

2. Add the `GenerateFlowFile` processor for data input. This processor creates FlowFiles with random data or custom content.

- a) Drag and drop the processor icon into the canvas. This displays a dialog that allows you to choose the processor you want to add.
- b) Select the `GenerateFlowFile` processor from the list.
- c) Click Add or double-click the required processor type to add it to the canvas.

You will configure the `GenerateFlowFile` processor to define how to create the sample data in *Configure the processor for your data source*.

3. Add the `PutIceberg` processor for data output.

You will configure the `PutIceberg` processor in *Configure the processor for your data target*.

4. Connect the two processors to create a flow.

- a) Drag the connection icon from the first processor, and drop it on the second processor.
A Create Connection dialog appears with two tabs: Details and Settings.

- b) Configure the connection.

You can configure the connection's name, flowfile expiration time period, thresholds for back pressure, load balance strategy and prioritization.

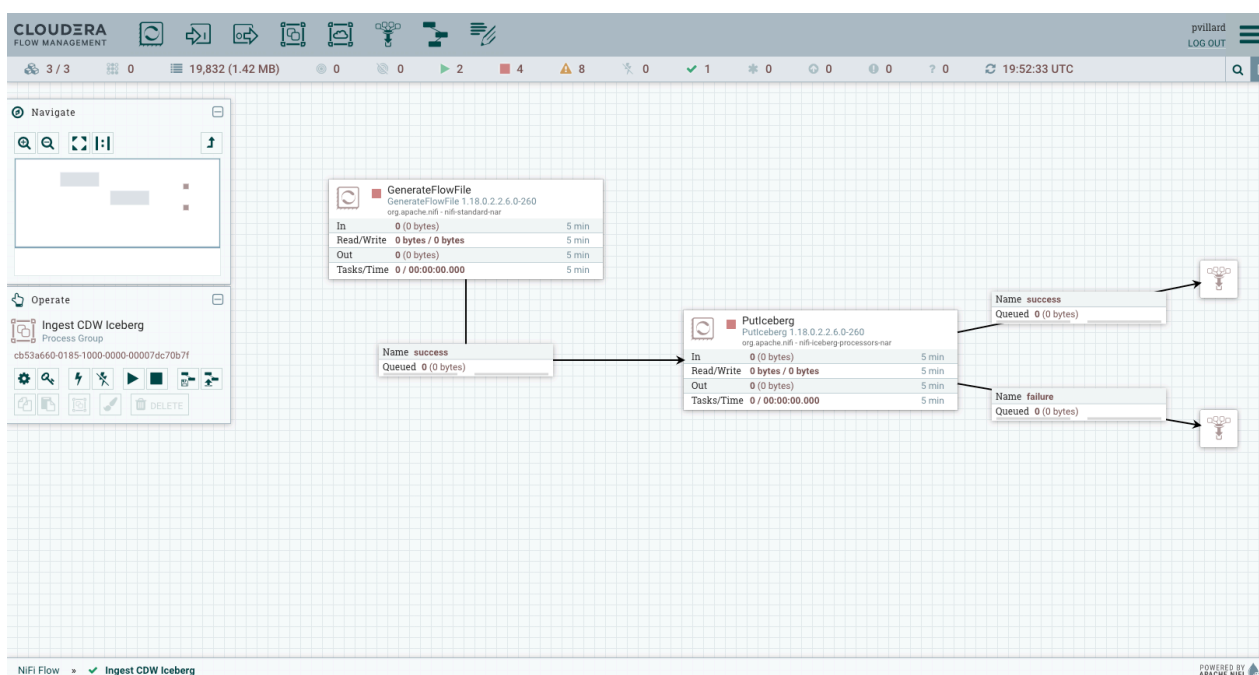
- c) Click Add to close the dialog box and add the connection to your data flow.

5. You can add success and failure funnels to your data flow, which help you see where flow files are being routed when your flow is running. Connect the `PutIceberg` to these funnels with success and failure connections.

If you want to know more about working with funnels, see the *Apache NiFi User Guide*.

Results

Your data flow will look similar to the following:



What to do next

Create controller services for your data flow. You will need these services when configuring the PutIceberg processor.

Create and configure controller services for your data flow

Learn how to create and configure controller services for the CDW Iceberg ingest data flow.

About this task

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

Procedure

1. To add a controller service to your flow, right-click on the canvas and select Configure from the pop-up menu. This displays the Controller Services Configuration window.
2. Select the Controller Services tab.
3. Click the + button to display the Add Controller Service dialog.
4. Select the required controller service and click Add.
5. Click the Configure icon in the right-hand column and configure the options that you need.
6. When you have finished the configuration, click Apply to save the changes.
7. Enable the controller service by clicking the Enable button (flash) in the far-right column of the Controller Services tab.

Example

The following controller services are used in this CDW Iceberg ingest data flow:

- Hive Catalog Controller Service
- Kerberos User Service

Configure the Hive Catalog Controller Service

In the context of CDW, the Hive Metastore URI is not exposed publicly. Instead, you need to provide the `core-site.xml` and `hive-site.xml` files of the CDP environment. In Flow Management DataHub clusters, these files are made available automatically on every node but the paths must be retrieved. If you use the Cloudera DataFlow data service and its Flow Designer, you can use the `#{CDPEnvironment}` parameter.



Note: If you want to push data into Iceberg tables in Data Engineering Data Hub clusters, you can use the Hive Metastore URI instead.

Table 1: Hive Catalog Controller Service properties

Property	Description	Example value for ingest data flow
Hive Metastore URI	Provide the URI of the metastore location.	
Default Warehouse Location	Provide the default warehouse location in the HDFS file system.	
Hadoop Configuration Resources	Add a comma-separated list of Hadoop Configuration files, such as <code>hive-site.xml</code> and <code>core-site.xml</code> for kerberos. Include full paths to the files so that NiFi can refer to those configuration resources from your specified path.	<code>/etc/hive/ conf.cloudera.data_context_connector-975b/ hive-site.xml,/etc/hadoop/ conf.cloudera.stub_dfs/core-site.xml</code>

Configure the Kerberos User Service

Use the Kerberos Password User Service so that you do not need to distribute a keytab file across the NiFi nodes of the cluster.

It is best practice to have a dedicated Machine User created in the control plane for your specific use case so that you can configure specific policies in Ranger and have better control in case of multi-tenancy with many use cases.

Table 2: Kerberos User Service properties

Property	Description	Example value for ingest data flow
Kerberos Principal	Specify the user name that should be used for authenticating with Kerberos. Use your CDP workload username to set this Authentication property.	<code>srv_nifi_to_iceberg</code>
Kerberos Password	Provide the password that should be used for authenticating with Kerberos. Use your CDP workload password to set this Authentication property.	<code>password (sensitive value)</code>

What to do next

Configure the processor for your data source.

Configure the processor for your data source

Learn how to configure the `GenerateFlowFile` data source processor for the CDW Iceberg ingest data flow.

About this task

You can set up a data flow to move data in Iceberg table format into Cloudera Data Warehouse (CDW) from many different locations. This example assumes that you are using sample data generated by the `GenerateFlowFile` processor.

Procedure

1. Launch the Configure Processor window by right clicking the `GenerateFlowFile` processor and selecting Configure. A configuration dialog with the following tabs is displayed: Settings, Scheduling, Properties, and Comments.
2. Configure the processor according to the behavior you expect in your data flow.
The `GenerateFlowFile` processor can create many `FlowFiles` very quickly. Setting the run schedule to a reasonable value is important so that the flow does not overwhelm the system.
3. When you have finished configuring the options you need, save the changes by clicking the Apply button.
Make sure that you set all required properties, because you cannot start the processor until all mandatory properties have been configured.

Example

The following settings and properties are used in this example:

Table 3: GenerateFlowFile processor scheduling

Scheduling	Description	Example value for ingest data flow
Run Schedule	Run schedule dictates how often the processor should be scheduled to run. The valid values for this field depend on the selected Scheduling Strategy.	500 ms

Table 4: GenerateFlowFile processor properties

	Description	Example value for ingest data flow
Custom text	If Data Format is text and if Unique FlowFiles is false, you can provide custom to be used as the content of the generated FlowFiles. The expression statement in the example value generates a random ID between 1 and 10 000, with random last names assigned.	<code>customer_id, customer_name \${random():mod(10000):plus(1)}, Smith \${random():mod(10000):plus(1)}, Johnson \${random():mod(10000):plus(1)}, Steward</code>

What to do next

Configure the processor for your data target.

Configure the processor for your data target

Learn how to configure the target processor of the CDW Iceberg ingest data flow.

About this task

You can set up a data flow to move data to many locations. This example assumes that you are moving data to Cloudera Data Warehouse (CDW). Once your data is collected and available in your flow, you can use the `PutIceberg` processor to send the data in the CDW table.

Procedure

1. Launch the Configure Processor window, by right-clicking the `PutIceberg` processor and selecting Configure. A configuration dialog with the following tabs is displayed: Settings, Scheduling, Properties, and Comments.

- Configure the processor according to the behavior you expect in your data flow.
With the `PutIceberg` processor, you can ingest data using Iceberg table format.
- When you have finished configuring the options you need, click `Apply` to save the changes.
Make sure that you set all required properties, because you cannot start the processor until all mandatory properties are configured.

Example

The following settings and properties are used in this example:

Table 5: PutIceberg processor properties

Property	Description	Example value for ingest data flow
Record Reader	Select the Record Reader of your choice based on the format of your incoming data that should be pushed in the table.	JSON Reader
Catalog Service	Add the Hive Catalog Controller Service you have created.	Hive Catalog Service
Catalog Namespace	Set the Catalog Namespace property to the name of the database where you created the destination table.	default
Table Name	Set the Table Name property to the name of the table previously created.	customer
Kerberos User Service	Add the Kerberos User Service of your choice (based on username/password or based on username/keytab) that you have created.	Kerberos Password User Service

What to do next

Your data flow is ready to ingest data into an Iceberg table in CDW. You can now start the flow.

Start the data flow

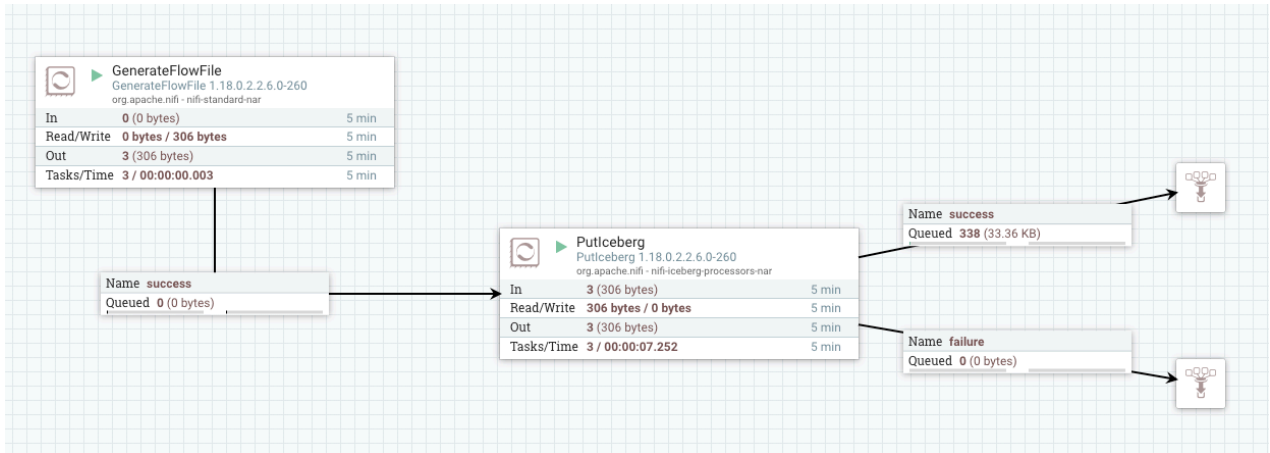
Learn how run the flow and confirm that data is successfully sent to the Iceberg table.

Procedure

- Select all the data flow components you want to start.
- Click the `Start` icon in the `Actions` toolbar.
Alternatively, right-click a single component and choose `Start` from the context menu.

Results

Your flow should be running without any errors. Data should be generated in the form of `FlowFiles`, and the files should be written to your designated Iceberg table.



You can verify the operation of your Iceberg ingest data flow by checking that data is running through the flow that you have built and it actually appears in the target table.

- You can verify that NiFi processors are not producing errors.
- You can look at the processors in the UI to see the amount of data that has gone through them. You can also right-click on the processors, or on connections to view status history.
- You can query the data in Hue:

The screenshot shows the Hue web interface. On the left is a sidebar with navigation options like Editor, Documents, S3, Tables, and Importer. The main area displays an Impala query: `select * from customer limit 100`. Below the query, there are two status messages: `Query 644f0b7f6b5b4d48:57419ace00000000: 0% Complete (0 out of 3168)` and `Query 644f0b7f6b5b4d48:57419ace00000000: 0% Complete (0 out of 3168)`. A 'Query History' section shows a table of results with 100 rows.

id	name	created_at	country_code	
1	697	John Doe	2023-01-25 07:31:20.842000000	FR
2	39	John Doe	2023-01-24 21:29:49.074000000	FR
3	626	John Doe	2023-01-25 06:20:20.966000000	FR
4	1005	John Doe	2023-01-25 12:33:47.506000000	FR
5	767	John Doe	2023-01-25 08:41:21.177000000	FR
6	628	John Doe	2023-01-25 06:22:20.772000000	FR
7	841	John Doe	2023-01-25 09:55:21.294000000	FR
8	993	John Doe	2023-01-25 12:27:21.193000000	FR
9	86	John Doe	2023-01-24 21:30:36.110000000	FR
10	997	John Doe	2023-01-25 12:31:21.198000000	FR
11	659	John Doe	2023-01-25 06:53:20.799000000	FR