

Starburst to Collibra Data Products Integration v1.0.0 (Python)

Name	Starburst to Collibra Integration
Category	Integration
Use Case	Increase the visibility and usability of Starburst Data Products in Collibra.
Target Audience	Data Architect, Data Engineer, Data Analyst, Business Analyst, Data Scientist
Who can set it up	Integration Engineer
Target Business Functions	Any
Target Industries	Any
License Requirements	Collibra Platform Collibra Catalog Starburst Enterprise
Dependencies	Collibra Platform v2021+ Starburst Enterprise 380-e LTS+ Python 3.9+ Jupyter Notebook
Developer	Collibra, Starburst
License	Binary Code License Agreement
Short Description	A Jupyter notebook that automatically brings business metadata for Starburst Data Products and Data Domains into Collibra and then links that metadata to the appropriate data sets (views and/or materialized views) in Collibra that have been ingested through the Starburst driver.

- [Overview](#)
- [Requirements](#)
- [Constraints](#)
- [Functional Design](#)
- [Technical Design](#)
- [Installation](#)
 - [Prerequisites](#)
 - [Install Python on your preferred OS \(Unix/Windows\) or environment](#)
 - [Install PIP - Python Packages Installer](#)
 - [Upload the CMA Migration File to Install the Operating Model](#)
- [Configuration](#)
- [Usage](#)
- [Release History](#)
- [What's Next?](#)

Overview

A [Starburst data product](#) is a schema that contains one or more data sets, which are represented as views and/or materialized views in the data source where they are located. The [Starburst JDBC driver in the Collibra Marketplace](#) will automatically ingest the following metadata for each Starburst data product:

- Schema name
- Data set name
 - Name of each column in the data set
 - Datatype for each column in the data set

The Starburst JDBC driver does not ingest the following metadata for Starburst data products:

- Data product name
- Data product summary
- Data product description
- Data domain
- Catalog (the data source where the data sets are located)
- Data product owners
- Tags
- SQL query used to create each data set

This integration will do the following:

- Create assets in Collibra for each published data product in Starburst
- Create assets for each domain the data products are associated with
- Extract the aforementioned data products metadata from Starburst
- Add the metadata to the corresponding assets in Collibra
- Link the data domain asset to the appropriate data product asset
- Link the data product asset to the appropriate data sets

Requirements

To use this integration, you must meet the following requirements:

- Valid license for Collibra Integration Cloud, Collibra Catalog and Starburst Enterprise
- Collibra Integration Cloud v2021+
- Starburst Enterprise 380-e LTS+
- Python 3.9+
- Jupyter Notebook
- Collibra CMA file containing the required operating model

Constraints

This integration has the following constraints:

Starburst

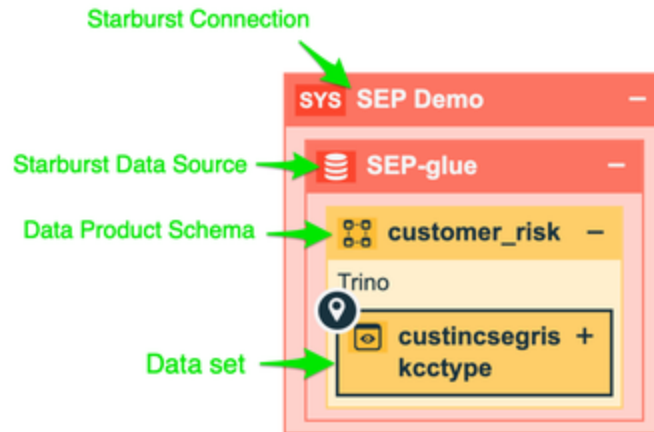
- *The integration is only supported for Starburst Enterprise at this time and is not supported for Starburst Galaxy.*
- Only **published** data products in Starburst Enterprise will be ingested into Collibra. Data products that are in draft status will not be ingested.
- Changes made to data domain and data product assets in Collibra will not be reflected in Starburst Enterprise. Any changes that need to be made on data domains or data products should be made in Starburst Enterprise and then brought into Collibra via this integration.

Collibra

- The integration is only supported for Collibra Integration Cloud.
- Requires Collibra Edge Catalog to ingest metadata using the certified Starburst JDBC driver found on the Collibra Marketplace.
- Create Data Domain and Data Product pulled from Starburst and then enrich the metadata created from the previous bullet.
- Associated Dataset of Starburst Data Products are treated as Views and Materialized Views are represented as Database View and Table, respectively, by Collibra.
- Collibra does not currently support spaces in tags, so any spaces in tags that exist in Starburst will be replaced with underscores ('_') when ingested into Collibra.

Functional Design

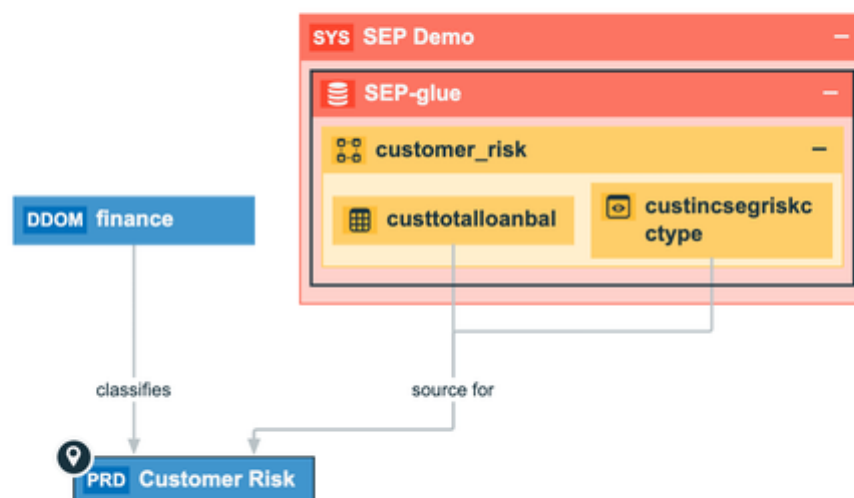
Before you run the integration, you will need to have already ingested metadata into Collibra from Starburst for the data sources ("catalogs") and schemas where your data products are located. For example, the image below depicts a data set (view) and data product schema that were automatically ingested using the certified Starburst JDBC driver.



When you run the integration, you will iterate through each cell in the Jupyter notebook starting at the top of the notebook and working your way to the bottom. Provided below is a high-level overview of the operations that are performed by running the Jupyter notebook.

1. Prompt user for Starburst and Collibra connection information (URL, username, password)
2. Prompt user for the type of pull they want to perform
 - a. *Valid values:* data domains, data products
3. Prompt user for the Collibra community and system where the data products and data domains will be added
4. Retrieve all data domains and their metadata from Starburst and create a domain asset for each one in Collibra
5. Retrieve all "published" data products and their from Starburst and create a data product asset for each one in Collibra
6. Update the views/materialized views associated with each data product in Collibra with the definition SQL queries from Starburst
7. Link the data domain asset and data set(s) (views/materialized views) to the appropriate data product asset in Collibra

After running the integration, you will see a data product asset in Collibra for each published data product that exists in your Starburst Enterprise environment. You will also see a data domain asset in Collibra for each domain that was extracted from the published data products in Starburst Enterprise. Each data product asset in Collibra will be linked to the appropriate data domain asset and the data set(s) that are part of the data product. The example below shows the result of ingesting the **Customer Risk** data product from Starburst Enterprise into Collibra using this integration.



- These assets were created through the Starburst JDBC driver and were not created through this integration:
 - SEP Demo (*system*)
 - SEP-glue (*data source*)
 - customer_risk (*schema*)
 - custtotalloanbal (*materialized view, represented as a table*)
 - custincsegrickctype (*view*)
- These assets were created after running this integration:
 - finance (data domain)
 - Customer Risk (data product)

The table below shows the mapping between the metadata fields for Starburst Enterprise data products and the corresponding fields on the catalog page for the data product in Collibra.

1	Name	Name
2	Catalog	Catalog Name
3	<i>The schema name is automatically generated from the name of the data product.</i>	Schema Name
4	Summary	Definition
5	Description	Description
6	Data product owners	Product Owner
7	Tags	Tags
8	Domain	is classified by Business Dimension
9	Datasets	Is target of Data Structure

Starburst Enterprise

Starburst Enterprise

<> Query editor

Data products

Domain management

INSIGHTS

Overview

Query overview

Cluster history

Usage metrics

Overview

Usage examples

Discussion 0

1 & 3 Customer Risk

BOOKMARKED 0 ★★★★★ 1

Overview

Catalog
glue

Summary
Customer financial profile data to assess risk for inclusion in upsell/cross-sell campaigns.

Number of queries	Past 7 days	Past 30 days	Number of users	Past 7 days	Past 30 days
	1	3		1	1

Description
This data product contains the financial profile data for each customer, which will help identify the degree of risk associated with that customer when considering them as a candidate for upsell and cross-sell campaigns. This data product contains the following data:

- Segmentation
- Income
- Risk appetite
- Mortgage balance(s)
- Auto loan balance(s)
- Credit card type

Datasets
Showing 2 of 2 datasets
custinsegriskcotype
This data set provides the estimated income, customer segmentation, risk app...

Data product owners
Ben Lumbert
ben.lumbert@starburstdata.com
Shaun Van Staden
shaun.vanstaden@starburstdata.com

Tags
customer risk finance PII

Domain
finance

Relevant links

Details
Created on: November 27, 2022
Created by: sa.ben.lumbert
Last updated on: March 30, 2023
Last updated by: sa.ben.lumbert
Last queried: March 30, 2023
Last queried by: sa.ben.lumbert

View in BI tools

Version: 398-e
JDK version: 17.0.3+7-LTS
Environment: official_demo_partner
Uptime: 1d 5h

Collibra Cloud

Starburst Data 2 > finance Data Products

PRD

Customer Risk

1

Data Product ⓘ Candidate ⓘ | 4 0 0%

Add to Data Basket Actions ▾

Add characteristic <

Details

Tags (4)

Comments

Diagram

Pictures

Quality

Responsibilities

References

History

Files

Description ⓘ

This data product contains the financial profile data for each customer, which will help identify the degree of risk associated with that customer when considering them as a candidate for upsell and cross-sell campaigns. This data product contains the following data:

- Segmentation
- Income
- Risk appetite
- Mortgage balance(s)
- Auto loan balance(s)
- Credit card type

Definition ⓘ

Customer financial profile data to assess risk for inclusion in upsell/cross-sell campaigns.

Catalog Name

glue

Schema Name ⓘ

customer_risk

Product Owner ⓘ

Ben Lumbert (ben.lumbert@starburstdata.com)

Shaun Van Staden (shaun.vanstaden@starburstdata.com)

is classified by Business Dimension

Add ⌵

Name ↑	Domain	Description	
finance	Data Domains		

1

is target of Data Structure

Add ⌵

Name ↑	Domain	Description	
custincsegriskcctype	SEP Demo > glue > customer_...		
custtotalloanbal	SEP Demo > glue > customer_...		

2

Data Source ⓘ

Starburst

Tags

customer risk PII finance

Additionally, the definition query for each data set in the data product will be extracted from Starburst Enterprise and added to the **Definition Query** field on the **Details** page for the corresponding view/materialized view in Collibra Cloud.

Starburst Enterprise

☰

Starburst Enterprise

<> Query editor

Data products

Domain management

INSIGHTS

Overview

Query overview

Cluster history

Usage metrics

Define datasets

Your data product is made up of one or more datasets, based on queries you define. Enter the information below to define your dataset.

custincsegriskcctype

custtotalloanbal

Create dataset

Published dataset name *

custincsegriskcctype

i

Reference this as:

glue.customer_risk.custincsegriskcctype

Dataset description

This data set provides the estimated income, customer segmentation, risk appetite, state of residence and credit card type for each customer. The data is sources from PostgreSQL, MySQL and S3. It currently only returns data for customers in the US.

You must show columns before you can save and continue.

Create dataset from query *

1 SELECT c.custkey

2 ,c.state

3 ,c.estimated_income

4 ,cp.customer_segment

5 ,cp.risk_appetite

6 ,a.cc_number

7 ,pp.cc_type

8 FROM postgresql.burst_bank_large.customer_c

Data set definition query

Collibra Cloud

Browse

Search

Business Analysts community

Starburst

SEP Demo > glue > customer_risk

custincsegriskcctype

Database View ⓘ Candidate ⓘ | 0 0 | 5%

Add characteristic

Summary

Details

Columns

Sample data

Diagram

Pictures

Technical Lineage

Quality

Responsibilities

SYS SEP Demo > SEP-glue > customer_risk > custincsegriskcctype

Table Type ⓘ

VIEW

Definition Query ⓘ

```
SELECT c.custkey
,c.state
,c.estimated_income
,cp.customer_segment
,cp.risk_appetite
,a.cc_number
,pp.cc_type
FROM postgresql.burst_bank_large.customer c
JOIN mysql.burst_bank_large.customer_profile cp ON c.custkey = cp.custkey
JOIN glue.burst_bank_large.account a ON c.custkey = a.custkey
JOIN glue.burst_bank_large.product_profile pp ON a.custkey = pp.custkey
WHERE c.country = 'US'
AND c.state NOT IN ('AA', 'AE', 'AP')
```

The table below shows the mapping between the metadata fields for Starburst Enterprise domains and the corresponding fields on the catalog page for the data domain in Collibra.

1	Domain Name	Name
2	Schema location URI	Location

Starburst Enterprise

Query editor

Data products

Domain management

INSIGHTS

Overview

Query overview

Cluster history

Usage metrics

Domain management

Create new domain

8 domains

Domain name ↑	# of data products	Data products assigned to this domain
SK Telecom	0	No data assigned
customer	4	Customers, marketing_metrics, Truist Insurance Customer, gregcusc
finance	4	Customer Risk, Test, rev_unit_margin, Customer Risk2
logistics	0	No data assigned
product	0	No data assigned
risk	1	dp_1
sales	0	No data assigned
transactions	0	No data assigned

Edit: finance

Changes to domain will effect all data products listed below.

Domain name *
finance

Domain description

2048 character remaining

Data Products under this domain will be created as schemas in this location.

Schema location URI
s3://[redacted]/data_products/finance/

Data products assigned to this domain (4)

Customer Risk

Test

rev_unit_margin

Customer Risk2

Cancel Save changes

Collibra Cloud

Browse

Search

Starburst Data 2

Data Domains

DDOM

finance

Type: Data Domain

Status: Candidate

Ask the Expert

Create Issue

Delete Lake Formation Tags With Asset

Add characteristic

Overview

Tags

Comments

Diagram

Description

No value has been given yet. Double click or use the edit button.

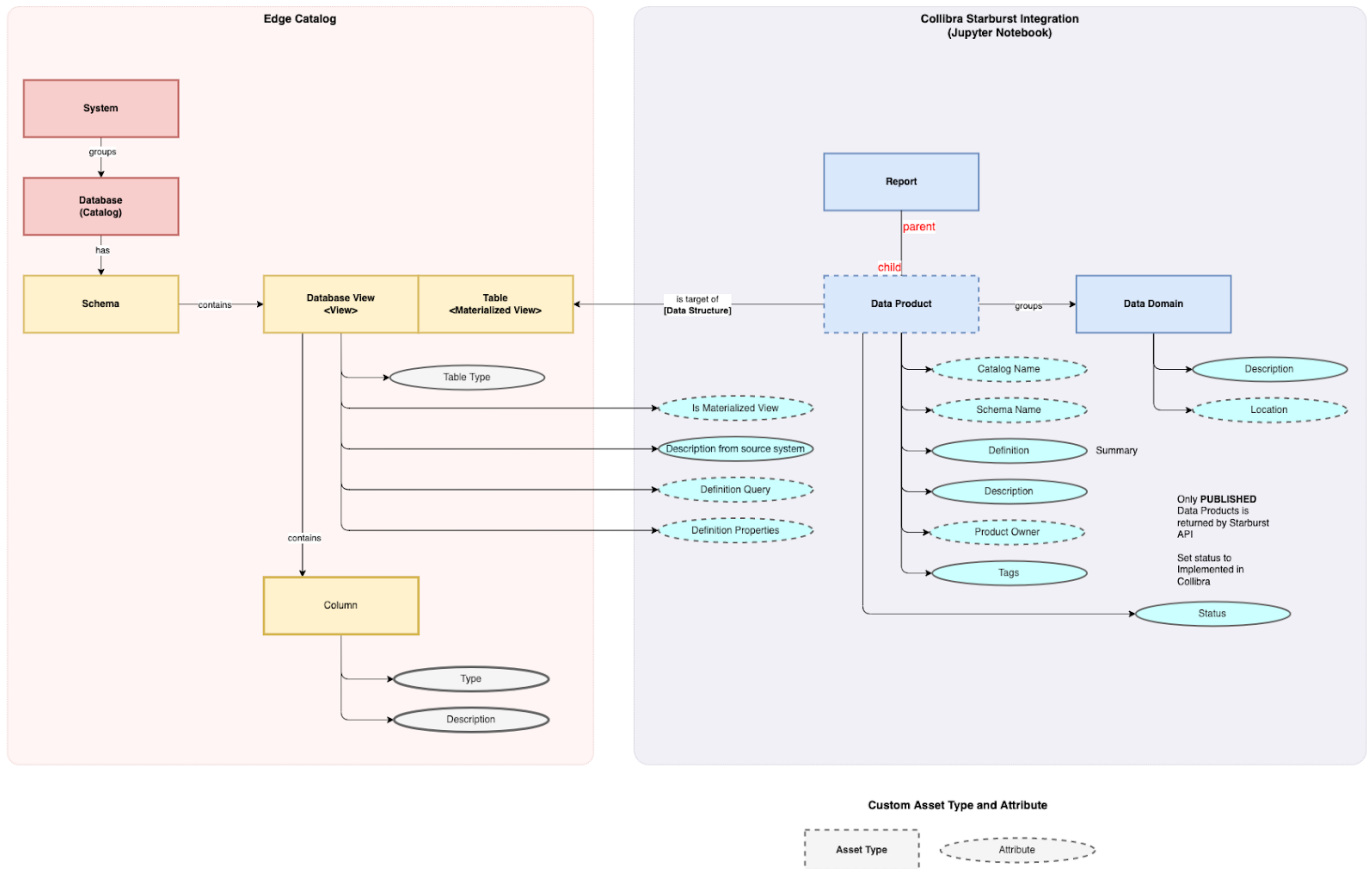
Location

s3://[redacted]/data_products/finance/

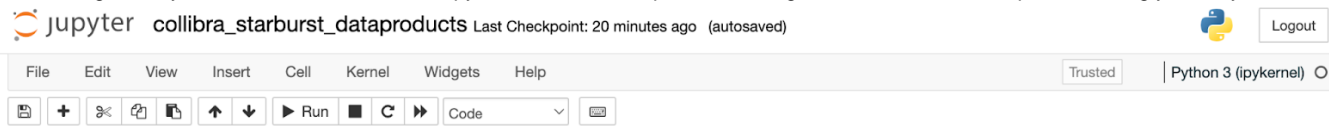
Technical Design

Below is the operating model for this integration. On the left are the asset types and attributes that are ingested by the Edge Catalog. On the right are the metadata that are ingested and linked to the Edge Catalog metadata. To be able to create the custom asset types, attribute types and the Starburst scope on the environment, make sure that you imported the provided CMA file.

Operating Model



To use the integration, you will run each cell in the Jupyter notebook in sequence, starting at the first cell at the top and working your way down.



Colibra Starburst Integration - Data Products

Start

```

In [ ]: ## STEP 1: Install required packages
!pip install -U -r requirements.txt

In [ ]: ## STEP 2: Import common packages, helper classes and helper functions
import getpass
from helper import StarburstColibraFacade

In [ ]: ## STEP 3: Provide Starburst environment details (URL, username, role, password)
sep_url = input('Enter the URL for your Starburst Enterprise instance (example: https://mystarburstcluster.com):\n')
sep_user = input('\n\nEnter your username:\n')
sep_role = input('\n\nIf you are using BIAC, enter your role name (leave blank otherwise):\n')
sep_pwd = getpass.getpass('\n\nEnter your password:\n')

In [ ]: ## STEP 4: Provide Colibra environment details (URL, username, password)
collibra_url = input('Enter the URL for your Colibra Cloud instance (example: https://<domain>.collibra.com): \n')
collibra_user = input('\n\nEnter your username:\n')
collibra_pwd = getpass.getpass('\n\nEnter your password:\n')
  
```

Provided below is an overview of the operations that are performed in each of the cells you see in the notebook.

STEP 1: Install required packages

- Installs common Python packages and Colibra Python packages in your local environment

STEP 2: Import common packages, helper classes and helper functions

- Imports common Python packages
- Imports Colibra Python packages

- Imports helper classes for Starburst and Collibra
- Imports helper functions

STEP 3: Provide Starburst environment details (URL, username, role, password)

- Prompts user for Starburst information
 - Starburst Enterprise URL
 - Starburst username
 - Starburst role (only required if [BIAC](#) is enabled)
 - Starburst password

STEP 4: Provide Collibra environment details (URL, username, password)

- Prompts user for Collibra information
 - Collibra Cloud URL
 - Collibra username
 - Collibra password

STEP 5: Provide temp directory, Collibra community name and system ID

- *Temp directory* - This will be used by the notebook to generate JSON payloads for the Collibra API calls
 - This directory will automatically be created for you in the location where you run the Jupyter notebook from
- *Collibra Community* - The name of the Collibra community where the data products and data domains from Starburst will be added
- *System ID* - The ID of the Starburst Enterprise system in Collibra where the data sets for the data products in Starburst are cataloged

STEP 6: Import Starburst data domains and data products into Collibra

- Instantiates the facade class in order to use the classes and functions defined in helper.py
- Calls the function to extract all data domains from Starburst and import them into Collibra
- Calls the function to extract all “published” data products from Starburst and import them into Collibra

The Jupyter notebook requires 2 additional files for it to operate correctly. These files need to be located in the same directory as the Jupyter notebook.

requirements.txt

This file contains the Python packages that need to be installed in your local environment in order for the Jupyter notebook to operate correctly. These packages will be installed for you when you run the first cell (“##STEP 1: Install required packages”) in the notebook.

helper.py


This file ensures the appropriate packages are imported into your local environment and defines the helper classes and functions required for the Jupyter notebook to operate correctly. The packages will be imported to your local environment and the helper classes and functions will be defined for you when you run the second cell (“## STEP 2: Import common packages, helper classes and helper functions”) in the notebook.

Provided below is an overview of what is contained in this file:

- Import common Python packages
- Import Collibra Python packages
- Define a reusable *CollibraService* class with functions that:
 - Import assets
 - Delete files
 - Read files
 - Write files
- Defines a reusable *ImportCommand* class to generate Import API requests in JSON format
- Defines a reusable *StarburstService* class with functions that:
 - Execute an API call to Starburst to retrieve all data domains
 - Execute an API call to Starburst to retrieve all tags for a data product
 - Execute an API call to Starburst to retrieve all of the metadata for a data product
 - Execute an API call to Starburst to retrieve the list of all ‘published’ data products in Starburst
- Defines a reusable *StarburstCollibraFacade* class
 - This class is used in the Jupyter notebook to call the other classes defined in this file
 - It also defines the following functions:
 - Query and import data domains
 - Pull all domains from Starburst via an API call
 - Create *Starburst Domain* Import Command and then send to Collibra Import API to create as assets
 - Update data product views
 - Link the data products to their respective views and/or materialized views in Collibra
 - Query and import data products
 - Export all “published” data products from Starburst via an API call
 - Create *Starburst Data Product* Import Command and then send to Collibra Import API to create as assets

- Update data product view columns
 - Update descriptions in Colibra for columns in the views/materialized views that are associated with the data products

Installation

 **Python Version:** Use Python 3.9+. Some dependencies in the Jupyter notebook might not be supported if the version is not met.

Prerequisites

- Install Python 3.9+
- Install PIP (Python packages installer)
- Download files from the [starburst-colibra/data_products GitHub repo](#)
 - requirements.txt
 - colibra_starburst_dataproducts.ipynb
 - helper.py
 - cma/starburst_colibra_integration-x.x.x.cma
 - cma/starburst_colibra_integration.xlsx
- You have connected Colibra to your Starburst Enterprise cluster using the Starburst JDBC driver to pull metadata
- You have run metadata extraction for the data sources and schemas where your data products are located
- You have the URL for your Starburst Enterprise instance
- You have the URL for your Colibra Cloud instance
- Your Starburst Enterprise user must have read access to the data sources and schemas where the data products are located
 - If your Starburst Enterprise cluster has [BIAC](#) enabled, be sure to have the name of the [role](#) with the appropriate permissions
- Your Colibra Cloud user must have admin access or at least permission to create and update assets in a Domain
- You have [imported the CMA file](#) for this integration to your Colibra environment


Install Python on your preferred OS (Unix/Windows) or environment

Refer to this [link](#) to download Python for the operating system of your choice.

Install PIP - Python Packages Installer

PIP is a python packages installer that is required in order to install modules specified in the Jupyter notebook. Instructions for installing PIP in your preferred OS can be found [here](#).

Upload the CMA Migration File to Install the Operating Model

 Remember to run the updated CMA file if you are upgrading from a previous version of the code. If there are any specific errors on the attribute and relation types associated with the Starburst asset types, try to delete it manually then import the CMA again.

To create the operating model dependencies, login to the Colibra DGC and then upload the file `cma/starburst_colibra_integration.cma` in the Settings > Migration > Import. More instructions can be found in the [Colibra documentation](#).

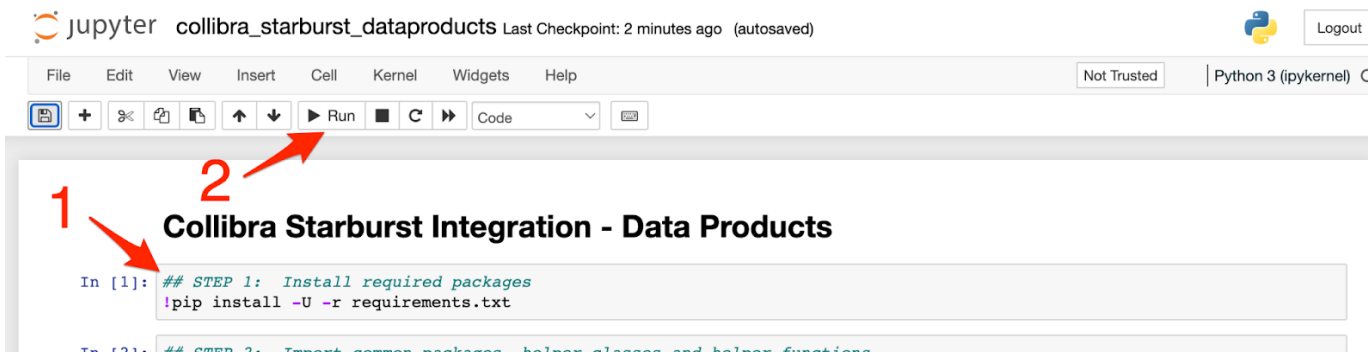
Configuration

After satisfying the prerequisites and completing the installation instructions, you are ready to run the integration. The Jupyter notebook is self contained and imports all of the required packages for you. No additional configuration is required.

Usage

The following instructions will walk you through how to run the integration.

1. Open the Jupyter notebook
2. Select the first cell (1) and click **run** (2).
3. After running the cell, you will automatically be moved to the next cell.
4. Click **run** for each cell until you reach the end.
 - If you encounter an error when running a cell, correct the error and run the cell again.



Below are guidelines for cells that require input from you or that produce output.

- **## STEP 3 - This cell will prompt you for your Starburst connection information:**
 - Enter the full **URL** of your Starburst Enterprise environment, including the HTTP protocol (example: <https://mycluster.mydomain.com/>)
 - Enter the **username** to log into Starburst Enterprise with
 - Enter the name of a role your username has access to
 - This is only required if BIAC is enabled in your Starburst Enterprise cluster
 - Leave this blank if BIAC is not enabled in your cluster
 - Enter the password associated with the username you provided
- **## STEP 4 - This cell will prompt you for your Colibra connection information:**
 - Enter the full **URL** of your Colibra Cloud environment, including the HTTP protocol (example: <https://mycatalog.colibra.com/>)
 - Enter the **username** to log into Colibra Cloud with
 - Enter the **password** associated with the username you provided
- **## STEP 5 - This cell will prompt you for the following information:**
 - A temporary directory where files will be created for the JSON payloads that are needed as input to the Colibra API calls
 - The name of the **Colibra community** where the data domains and data products will be loaded
 - The **system ID** of the data source schema where the data sets for the Starburst data products are located
- **## STEP 6 - This cell is used to run the integration. If everything works correctly, you should see the following output:** About to pull /ingest all the data domains from Starbur to Colibra. About to pull/ingest all the published data products from Starburst to Colibra. Completed.

After running all of the cells in the Jupyter notebook, you can use the [starburst_colibra_integrations.xlsx](#) file in the GitHub repository as a guideline to confirm all of the expected assets have been added to your Colibra environment.

Release History

- Version 1.0.0 - Initial release of this integration.

What's Next?

The initial release of this integration focuses on extracting data domains and data products from Starburst Enterprise, then importing them as assets into Colibra and connecting them to the appropriate data sets (views and/or materialized views). The next phase of this integration will focus on creating a bi-directional sharing of data product and data domain metadata between Starburst Enterprise and Colibra. In that phase, Colibra users will be able to initiate the creation of data products from within Colibra Cloud.

A Colibra user will start the process by providing the following information:

- Select the view(s) and/or materialized view(s) in Colibra that will be used as data sets for the data product.
- Provide metadata (including the data domain) about the data product
- Provide notes regarding any operations (transformations, aggregations, etc) that need to be applied to the data sets

A designated admin in Colibra will be notified when the request has been made. Some of the information above will be used to create a draft data product in Starburst Enterprise. The admin will then use the provided notes to finish creating the data product in Starburst. Once that has been done, the metadata for the newly created data product will be automatically exported from Starburst and imported into Colibra using the functionality provided by the initial release of this integration.