Collibra Data Intelligence Cloud

# Data Catalog

Collibra Data Intelligence Cloud - Data Catalog

Release date: September 11th, 2022

Revision date: Thu Sep 08, 2022

You can find the most up-to-date technical documentation on our Documentation Center at

https://productresources.collibra.com/docs/collibra/latest/Content/Catalog/to_catalog.htm

# Contents

# Catalog submenu pages

The following table describes each of the submenu items of the Data Catalog application.

| Page | Description |
| --- | --- |
| Data Catalog Home | The landing page when you click the Data Catalog tab. This page is designed to help you quickly and easily find Data Catalog-related assets. |
| Reports | All report assets. |
| Data Sets | All data sets shown as a set of tiles or as a table, with their name, description and, if there are any, connections to existing assets in Collibra. |
| Data Sources | Data sources that are used for data source registrations. |
| Data Dictionary | All data assets in Collibra. |
| Technology Assets | All technology assets in Collibra. |
| Metrics | Contains a variety of statistics related to how the assets of the Catalog are used. |
| Access Requests | The history of your access requests and their status. |
| Advanced Data Types | All advanced data types, which are used during a data source registration. |

# Data Catalog asset pages

The asset pages in Data Catalog provide information about assets. The information depends on the asset type and the asset type's assignment.

## Catalog experience setting

Catalog experience is a setting that improves the user experience of the Data Catalog asset pages. The improvements include:

- Custom tabs that correspond to the page you are working on.
- A streamlined title bar showing general information.
- Quicker and easier navigation that requires less scrolling.

The Catalog experience setting is enabled by default. If required, you can disable it.

> Note   When you use Edge, Catalog experience has to be enabled.

## Page layout

For more information on the Data Catalog asset pages, see the online version of this guide.

## Data Catalog Home

The Collibra Data Catalog Home is the landing page when you click the Data Catalog tab. This page is designed to help you quickly and easily find Data Catalog-related assets.

> Note   You need the Data Catalog global role or Data Catalog Author role to view Data Catalog Home.

The page is organized into five groupings, or sections, of assets and a Data Catalog-specific search field, as described in the following image and table.

> Note   The **Data sets you might like** section is enabled and disabled via Collibra Console. By default, it is enabled (shown) on the page. The other four sections are always shown and cannot be disabled. However, for any of the five sections, if there is no relevant data, nothing is shown on the page, including the section header.

| Element name | Description |
| --- | --- |
| Search field | A Data Catalog-specific search field that you can use to find any asset in CollibraData Catalog, for example assets of asset types **Data Set**, **Schema**, **Table**, **Column**, **Tableau Workbook** and **Tableau View**.<br><br>This search field works in the same manner as does the global search field, but it uses a default 'Data Catalog' filter. |

| Element name | Description |
|---|---|
| Data Catalog Data Sets you might like | Shows up to four data sets you might be interested in, as determined by the recommender, which takes into account your data sets and the data sets of similar users.<br><br>The **Show more** button enables you to view up to eight data sets on this page. |
| Recently viewed | Shows the four most recently viewed Data Catalog-related assets.<br><br>This section uses the Recent widget functionality.<br><br>The **Show more** button enables you to view the eight most recently viewed assets. |
| Reports | Shows the four most recently created assets of asset type **Report** and its child asset types.<br><br>Clicking the asset name takes you to the asset page.<br><br>Clicking **View all reports** takes you to the Catalog reports page. |
| Data sources | Shows the four most recently created assets of asset type **Table**.<br><br>Clicking the asset name takes you to the asset page.<br><br>Clicking **View all data sources** takes you to the Data Sources page. |
| Data sets | Shows the four most recently created assets of asset type **Data Set**.<br><br>Clicking the asset name takes you to the asset page.<br><br>Clicking **View all data sets** takes you to the Data Sets overview page. |

# Recommenders

The recommenders aim to suggest relevant business assets and data sets. Recommenders have to train regularly to update the recommendations. By default, this is

done every night. Recommendations can be calculated on the basis of several algorithms. These algorithms also calculate an error margin for each recommendation, and eventually only the algorithm with the lowest error margin provides the recommendations.

You can edit the settings of the recommenders and matchers to optimize the recommendations.

> Note   The recommender uses statistical information. Therefore, your recommendations will be empty or less useful if your company just started using Collibra Data Intelligence Cloud.

# Recommendation of data sets to users

## Description

The data set recommender recommends data sets to users, based on the data sets of similar users.

If you use some of the same data sets as some other users, you are probably also interested in data sets that they use but you don't. The recommendations are shown on Data Catalog Home.

## Example

## Strategy

The data set recommender compares the data sets used by the users to find relevant data sets. It roughly follows these steps:

1. See which data sets you are currently using.
2. Look for other users that also use your data sets.
3. See which data sets those users use, but you don't.
4. Recommend up to 9 of those data sets to you.

> **Note**
> If the recommender does not have enough data, for example if you just started using Collibra, it only considers 3 parameters:
>
> - Certified
> - Quality
> - Popularity (number of views of the data set asset page)

# Recommendation of business assets to data sets

## Description

The asset recommender recommends business assets to data sets, based on business assets it is related to.

If two data sets have relations to the same business assets, business assets related to only one of the two data sets may be relevant to the other data set as well.

Example

## Strategy

The asset recommender uses the relation **data set related to business asset** set to find relevant assets. It roughly follows these steps:

1. See which business assets are related to the current data set.
2. Look for other data sets related to those business assets.
3. See whether those data sets are also related to other business assets.

4. Recommend those business assets on the data set page and in the **Add related to** dialog box.

> Note  If the recommender does not have enough data, for example if you just started using Collibra, it does not give you any recommendations.

## Recommendation of business assets to column assets

Business assets are recommended to column assets based on the search engine in Collibra. The recommendations are shown in the section of **data asset represented by business asset** relation.

Example

# Recommendation of business assets to Tableau workbook assets and Tableau view assets

Business assets are recommended to Tableau workbook assets and Tableau view assets based on the search engine in Collibra. The recommendations are shown in the section of **report related to business asset** relation.

# Matchers

The matchers aim to suggest assets and data sets that might be interesting for you.

Matchers find similar data sets and schemas based on the name and the attributes.

You can edit the settings of the recommenders and matchers to optimize the recommendations.

> Note   The matcher uses statistical information. Therefore, your recommendations will be empty or less useful if your company just started using Collibra Data Intelligence Cloud.
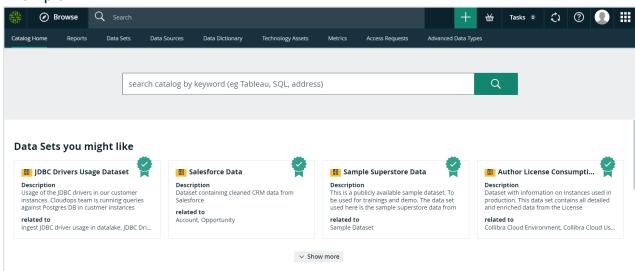
# Data set matcher

The data set matcher looks at the names and attributes of the column assets that a data set contains. It shows similar data sets on the data set asset page.

## Schema matcher

The schema matcher is currently not used in Collibra.

# Data Catalog Search

The Data Catalog Home page has a Data Catalog-specific search field that you can use to find assets in Data Catalog. When you launch a search from Data Catalog, the search page is the regular Collibra search page, but with the **Catalog** search filter applied.

> Note   You need the Data Catalog global role or Data Catalog Author role to view the Data Catalog search page and use the Data Catalog Search.

In the search input field, you can type any text and press `Enter` or click **Search** to launch the search.

The search finds resources that contain a word that begins with your search text. For example, if you type *ca*, the search results could contain 'California' and 'Lewis Carroll', but not 'Meercat'.

You can also use wildcards and symbols to search, see Wildcards and symbols for searching.

# Catalog reports

The **Reports** page is a view that shows:

- All **Report** assets.
- All packaged or manually created child asset types of **Report**, for example BI Report, Tableau View, and Looker Query.

# Report views

You can view the assets in table or tile display mode, and can perform all the same actions you can for any other table or set of tiles.

## Reports in tile display mode

In tile display mode, you can do the following:

- Click an asset name to open the relevant asset page.
- Click anywhere else in the tile to select one or more assets. The list of available actions appears in the action toolbar.

## Reports in table display mode

In table display mode, you can do the following:

- Click an asset name to open the relevant asset page.
- Click anywhere else in the tile to select one or more assets. The list of available actions appears in the action toolbar.
- Edit cells in the table.



## Filters

The default **All reports** view does not contain a filter, so it shows all Report assets. Some of the other packaged views do contain a filter. For example the **Certified reports** view only shows reports that are certified.

You can also create your own filter and, if necessary, save the filtered view as a new view. For example, you can create separate views for Report assets belonging to a specific source, for example Tableau, Looker or Power BI.

# Data Catalog Data Sets

A data set is a logical, handpicked collection of data elements that can come from multiple data sources. For example, Customer Contact information. Data sets allow users to quickly know which data to use for a specific purpose and request access to it.

The Catalog **Data Sets** overview page displays existing data sets in a table or as tiles. The page displays the name of the data set, its description, its certification status, and, if there are any, connections to existing business assets in Collibra Data Intelligence Cloud.

# Data Sets overview page

The Data Sets overview page contains the data sets that are available in Collibra Data Intelligence Cloud. You can view the data sets in table display mode or tile display mode.

## Tile display mode

- Click a data set title to open its details.

- Click anywhere in the tile except for the title to select the data set. The list of actions that you can perform is displayed.



# Table display mode

- Click a name of the data set to open its details.



- Select one or more data sets. The list of actions that you can perform is displayed.



> Note   The Sample Data tab shows the first 100 columns of data. If you have more than 100 columns, they are not shown.

# Data Set asset page

The **Data Set** asset page is basically the same as any asset page in Collibra Data Intelligence Cloud with the following differences:

- The Data Set asset page has a special attribute, namely **Certified**. That attribute indicates whether a data set is certified or not. There are no restrictions for certifying a data set, except the ones your organization chooses. You decide when a data set can or has to be certified. For more information about how to do this, see Certify a data set.
- It contains suggestions for related Business Assets, based on the asset recommender.
- It contains a **Data Profiling** and **Sample data** section which contains respectively a data profile and sample data, if available.

You can perform the following actions on this page:

- Create a view
- Filter data
- Sort Catalog submenu pages
- Request access to data sets and reports
- Delete data sets

# Creating data sets

In this section you can learn how to create a data set and how to add data to it.

## Create a data set

You create data sets to add data to them.

## Steps

1. In the main menu, click the **Create** ($+$) button.
   » The **Create** dialog box appears.
2. In the **Create** dialog box, click the **Asset** tab.
3. Click **Data Set**.

4. In the **Domain** field, select the domain to which you want to add one or more data sets.
5. In the **Name** field, type the name of the data set, press `Enter` to add other data set names.
6. Click **Create**.

# Add data to a data set from an asset page

When you come across an asset that you want to add to a data set, you can add that asset from that asset page.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a resource role with the Attribute > Add resource permission.

## Steps

1. Navigate to an asset page of a schema, table or column asset.
2. In the upper-right corner, click **Add to Data Set**.
3. Enter the required information in the **Add data to data set** dialog box.
   - Existing data set:
     a. Select the data set.
     b. Click **Add to data set**.
   - New data set:
     a. Type a name in the **Data set name** field.
     b. Type a description in the **Data set description** field.
     c. Click **Create & Add data**.

# Add data to a data set from the Data Sources or Data Dictionary page

You can add data to a data set from the Data Sources or Data Dictionary page.

## Prerequisites

- For Data Dictionary: You have a global role with the Data Dictionary global permission, for example Data Dictionary.
- You have a resource role with the Attribute > Add resource permission.

## Steps

1. On the main menu, click ⣿, then ⊟ **Catalog**.

    » The Catalog Home opens.
2. In the submenu, click **Data Sources** or **Data Dictionary**.

    If necessary, filter the list of data assets.
3. Select the check boxes of the data assets you want to add to a specific data set.

    > Note
    > ◦ Some data assets are nested. If you select the top one, all its children are added as well.
    > ◦ Keep in mind that you can only add schemas, tables and columns.

4. Above the table, click **Add to Data Set**.
5. Enter the required information in the **Add data to data set** dialog box.
6. Click **Add to data set**.

    » A notification in the upper-right corner lets you know how many assets you have added to the data set.

# Certify a data set

You can approve, endorse or guarantee the contents of a data set.

## Steps

1. Navigate to the asset page of a data set that you want to certify.
2. Find the **Certified** characteristic and double-click the line of text below it.
3. Click in the field that is displayed.
4. Click **True**.
5. Click **Save**.

> Tip   You can design a workflow to take care of the certification of a data set.

# Delete data sets

If you no longer need a certain data set, you can delete it from the repository.

## Steps

1. On the main menu, click ⦂⦂⦂, then 🗄 **Catalog**.
   » The Catalog Home opens.
2. In the submenu, click **Data Sets**
3. Search for the data sets that you want to delete.
   You can use the Filter pane or sort your data sets.
4. In table mode, select the check boxes of the data sets that you want to delete.
   In tile mode, hold the SHIFT key to select multiple data sets.
5. Click **Delete**.
6. Click **Yes** to confirm.

# Requesting access to data

You can request access to data by adding the relevant data sets or reports to your Data Basket and checking out your Data Basket.

## Adding data sets or reports to the Data Basket

You can add data sets or reports to the Data Basket by clicking **Add to Data Basket**. This button appears:

- When you've selected one or more assets of the type Data Set, Report, or one of their child asset types in Data Catalog.
- On Data Set asset pages, Report asset pages and asset pages of child asset types of Data Set and Report.

When you click **Add to Data Basket**:

- A Data Usage asset is created.
- All of the data sets and reports you selected are shown in the Data Basket.

## Data Usage asset

The Data Usage asset is created in the Data Usages domain. The name of the Data Usage asset is "USER_BASKET_" followed by the UUID of the user.

> Tip   The Data Usages domain is a "hidden" domain in the Business Analysts Community. This means it doesn't appear in the Collibra Browser, which helps to avoid it being inadvertently deleted. To view the Data Usages domain, go to the Access Requests page and click the name of a Data Usage asset. The Data Usages domain appears in the breadcrumb, on the Data Usage asset page.



The Data Usage asset page shows all of the important information related to the access request, including:

- The data sets or reports to which access is requested.
- The purpose for requesting access to the data.
- The access start date and end date.

## Data Basket

The Data Basket is a view that shows all of the data sets or reports you've selected and to which you want access.

To access the Data Basket, click 🛒.

You can remove data sets and reports from the Data Basket by clicking on the relevant tiles, and then clicking **Remove from data basket**.

## Checking out your Data Basket

When you're ready to check out, click **Checkout Data Basket**. This starts the packaged Request Data Sets Access workflow, by which your request is approved or rejected.

When you check out the Data Basket, the Data Usage asset is renamed in the format YYYY-MM-DD #X, where X is a sequential number, for example **2019-12-16 #3**.

All your access requests are shown in Catalog, on the Access Requests page.

## Request access to data sets and reports

To use the data referred to in assets of the type Data Set, Report, or one of their child asset types, you can request access to it via the Data Basket.

## Steps

1. Add assets to your basket by doing one of the following:

| Searching | a. In Collibra Data Intelligence Cloud, search for the data set or report you need and click its name.<br>   » The asset page of the data set or report opens.<br>b. Click **Add to Data Basket**. |
|---|---|
| Navigating | a. On the main menu, click ⠿, then ▤ **Catalog**.<br>   » The Catalog Home opens.<br>b. In the submenu, click **Data Sets** or **Reports**.<br>c. If necessary, search for the data sets or reports that you want to access.<br>d. Select the check boxes of the data sets or reports that you want to access.<br>e. Above the tiles, click **Add to Data Basket** |

   » A message at the upper-right indicates that the assets have been added to your basket.

2. Open your basket by clicking 🛒 on the main menu bar.
3. Review your basket.
To remove unnecessary assets, select them and click **Remove from data basket**.
4. Click **Checkout Data Basket**.
   » A dialog box appears.
5. Enter the required information and click **Submit**.
6. If the **Add Purpose to the Data Usage** dialog box appears, start typing and select the Purpose asset that describes the business use for which you are requesting access to the data sets or reports, and then click **Submit**.

> Note   This dialog box only appears if Collibra Data Privacy is installed.

## What's next?

- A workflow starts to approve the request and to grant you access to the data.
- A Data Usage asset is created in Collibra Data Intelligence Cloud. You can view all your requests and their current status on the Access Requests page. For more information, see Requesting access to data.



# Data Sources page

The **Data Sources** page is page that shows the asset types that are created by Database and S3 registrations. It's a combination of Data and Technology asset types.

You can view the assets in table display mode or tile display mode.

# Data sources in table display mode

With hierarchies enabled, you can expand the assets to consult the structure of the data sources. If needed you can also show other asset types in the lower levels of the hierarchy.

- Click an asset name to open the relevant asset page.



- Select one or more assets. The list of actions that you can perform is then displayed.

# Data sources in tile display mode

- Click an asset name to open the relevant asset page.

- Click anywhere else in the tile to select the asset. The list of actions that you can perform is then displayed.

# Data Dictionary page

The **Data Dictionary** page is a page that shows the assets of asset type Data Asset and its children asset types in Collibra Data Intelligence Cloud.

You can view the assets in table display mode or tile display mode.

On this page, you can perform the following actions:

- Create a view
- Filter assets
- Sort assets by name, description and asset type
- Delete assets
- Move assets
- Add assets to a data set
- Start an asset workflow from an asset table, for assets

# Technology Assets page

The **Technology Assets** page is a view that shows all assets of every technology asset type in Collibra Data Intelligence Cloud.

You can view the assets in table display mode or tile display mode.

On this page, you can perform the following actions:

- Create views.
- Filter assets.
- Sort assets by name, description and asset type.
- Delete assets.
- Move assets to another domain.

# Access Requests page

If you have requested access to one or more data sets, the Access Requests page allows you to view the status of your requests.

When you request access to a data set:

- an asset of the Data Usage type is created in the Data Usages domain in your community.
- the request Assets Access workflow is started.

The names of your requests are automatically generated with the date of your request. You can click the request name to open the asset page, which shows all the information relative to your request.

If you've requested access to many data sets, you can sort on any of the columns on the Access Requests page, to help you find a specific access request.

# Advanced data types

When you profile data when registering a data source, Collibra Data Intelligence Cloud can detect some basic data types, such as numbers and text. Besides these basic data types, you can create your own advanced data types.

> Note  Advanced data types are not taken into account when profiling via Edge. See About profiling and classification via Edge.

In this section, you learn how to work with advanced data types.

# Data type detection

When you run a data profiling when registering a data source, Collibra Data Intelligence Cloud tries to detect the data type of each column.

1. Collibra tries to match the fields of each column with every data type.
2. Collibra remembers the matches for each field, also if a field has multiple matches.
3. Collibra calculates the matching percentage of how many fields of the column match the same data type.

4. Collibra verifies the matching percentage against the data type detection threshold.

> Tip   You can define the data type detection threshold in Collibra Console, see the Collibra Installation and Configuration Guide.

5. Collibra assigns the data type with the highest matching percentage to the source column, provided that the matching percentage exceeds the threshold.

Out of the box, there are several base data types such as integer, text and boolean. With each data profiling, these base data types are evaluated. If your data source contains special data types such as social security numbers or international bank account numbers, you can define them as advanced data types. In the data source registration wizard, you can then choose to also evaluate the data on these advanced data types.

Keep in mind that detecting advanced data types significantly increases the data profiling job execution time.

# Advanced data type management prerequisites

To manage advanced data types, you need the following prerequisites:

- Catalog role
- Advanced Data Type global permission

# Create an advanced data type

If the basic data types, such as numbers and text, are not specific enough, you can create your own advanced data types.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a global role with the Advanced Data Type > Add global permission.

## Steps

1. On the main menu, click ⊞, then 🗄 **Catalog**.

    » The Catalog Home opens.
2. In the submenu, click **Advanced Data Types**.
3. Above the table, to the right, click **Add Advanced Data Type**.
4. In the **Add Advanced Data Type** dialog box, fill in the new data type properties.

| Option | Description |
|---|---|
| Name | The name of the advanced data type. The name has to be unique, including the basic data types. |
| Description | The description of the advanced data type. |

| Option | Description |
|---|---|
| Base data type | The data type used as basis for the advanced data type:<br><br>○ Text<br>○ Geographical<br>○ True/False<br>○ Date<br>○ Time<br>○ Date and Time<br>○ Whole Number<br>○ Decimal Number<br>○ Array<br>○ N/A<br><br>**Examples** |

| Base data type | Field name | Patterns |
|---|---|---|
| Text | Email address | `[a-z0-9]+[_a-z0-9\.-]*[a-z0-9]+@[a-z0-9-]+(\.[a-z0-9-]+)*(\.[a-z]{2,4})` |
| Text | IP address | `\b(?:(?:2(?:[0-4][0-9]|5[0-5])|[0-1]?[0-9]?[0-9])\.){3}(?:(?:2([0-4][0-9]|5[0-5])|[0-1]?[0-9]?[0-9]))\b` |
| Date | Custom Date | `yyyy-MM-dd` |
| Time | Custom Time | `HH mm` |

| Option | Description | | |
|---|---|---|---|
| | **Base data type** | **Field name** | **Patterns** |
| | Date and Time | Custom Date and Time | `MM/dd/yyyy HH:mm:ss` |
| | True/False | Boolean (French) | ○ true: vrai, v<br>○ false: faux, f |

| Option | Description |
|---|---|
| Advanced data type (variable field name) | The field name depends on the selected base data type. |

| Base data type | Field name | Description |
|---|---|---|
| Text | Regular expressions | List of regular expressions.<br><br>For more information about regular expressions, see regular-expressions.info. |
| Geographical | Regular expressions | List of regular expressions.<br><br>For more information about regular expressions, see regular-expressions.info. |
| Date | Date pattern | List of date patterns using the **DateTimeFormatter** format. See the official Java documentation. |
| Time | Time pattern | List of time patterns using the **DateTimeFormatter** format. See the official Java documentation. |

| Option | Description | | |
| --- | --- | --- | --- |
| | **Base data type** | **Field name** | **Description** |
| | Date and Time | Date and Time pattern | List of date and time patterns using the **DateTimeFormatter** format. See the official Java documentation. |
| | Whole Number | Numeric format | Locale for the format of whole numbers. |
| | Decimal Number | Numeric format | Locale for the format of decimal numbers. |
| | True/False | ○ True values<br>○ False values | ○ List of values that are accepted as **True** value.<br>○ List of values that are accepted as **False** value. |
| | Array or N/A | | Not applicable for advanced data type detection. |

5. Click **Save**.

# Edit an advanced data type

If an existing advanced data type is incorrect, you can edit it.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a global role with the Advanced Data Type > Update global permission.

## Steps

1. On the main menu, click ⠿, then 🗄 **Catalog**.

   » The Catalog Home opens.
2. In the submenu, click **Advanced Data Types**.
3. In the row of the data type that you want to edit, click ✎.

   The **Edit Advanced Data Type** dialog box appears.
4. Enter the required information.

   | Option | Description |
   | --- | --- |
   | Name | The name of the advanced data type. The name has to be unique, including the basic data types. |
   | Description | The description of the advanced data type. |

| Option | Description |
|---|---|
| Base data type | The data type used as basis for the advanced data type: <br><br> ◦ Text <br> ◦ Geographical <br> ◦ True/False <br> ◦ Date <br> ◦ Time <br> ◦ Date and Time <br> ◦ Whole Number <br> ◦ Decimal Number <br> ◦ Array <br> ◦ N/A <br><br> **Examples** |

| Base data type | Field name | Patterns |
|---|---|---|
| Text | Email address | `[a-z0-9]+[_a-z0-9\.-]*[a-z0-9]+@[a-z0-9-]+(\.[a-z0-9-]+)*(\.[a-z]{2,4})` |
| Text | IP address | `\b(?:(?:2(?:[0-4][0-9]\|5[0-5])\|[0-1]?[0-9]?[0-9])\.){3}(?:(?:2([0-4][0-9]\|5[0-5])\|[0-1]?[0-9]?[0-9]))\b` |
| Date | Custom Date | `yyyy-MM-dd` |
| Time | Custom Time | `HH mm` |

| Option | Description | | |
|---|---|---|---|
| | **Base data type** | **Field name** | **Patterns** |
| | Date and Time | Custom Date and Time | `MM/dd/yyyy HH:mm:ss` |
| | True/False | Boolean (French) | ○ true: vrai, v<br>○ false: faux, f |

| Option | Description |
|---|---|
| Advanced data type (variable field name) | The field name depends on the selected base data type. |

| Base data type | Field name | Description |
|---|---|---|
| Text | Regular expressions | List of regular expressions.<br><br>For more information about regular expressions, see regular-expressions.info. |
| Geographical | Regular expressions | List of regular expressions.<br><br>For more information about regular expressions, see regular-expressions.info. |
| Date | Date pattern | List of date patterns using the **DateTimeFormatter** format. See the official Java documentation. |
| Time | Time pattern | List of time patterns using the **DateTimeFormatter** format. See the official Java documentation. |

| Option | Description | | |
|---|---|---|---|
| | **Base data type** | **Field name** | **Description** |
| | Date and Time | Date and Time pattern | List of date and time patterns using the **DateTimeFormatter** format. See the official Java documentation. |
| | Whole Number | Numeric format | Locale for the format of whole numbers. |
| | Decimal Number | Numeric format | Locale for the format of decimal numbers. |
| | True/False | ◦ True values<br>◦ False values | ◦ List of values that are accepted as **True** value.<br>◦ List of values that are accepted as **False** value. |
| | Array or N/A | | Not applicable for advanced data type detection. |

You cannot change the base data type.

5. Click **Save**.

# Delete one or more advanced data types

If you no longer use an advanced data type, you can delete it.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a global role with the Advanced Data Type > Remove global permission.

## Steps

1. On the main menu, click ⠿ , then 🗄 **Catalog**.

   » The Catalog Home opens.
2. In the submenu, click **Advanced Data Types**.
3.

| Single advanced data type | a. In the row of the data type that you want to delete, click 🗑 . <br> b. In the **Delete advanced data type** dialog box, click **Delete advanced data type**. |
|---|---|
| Multiple advanced data types | a. Select the check boxes in front of the advanced data types that you want to delete. <br> b. In the action toolbar, click **Delete**. <br><br> Tip  You can select all the visible assets at once by clicking the check box next to the **Name** column header. <br><br> c. In the **Delete (x) advanced data type(s)** dialog box, click **Delete (x) advanced data type(s)**. |

The data type attributes that contain the deleted advanced data type are reset to the base data type that was used for the advanced data type.

# Sort Catalog submenu pages

You can reorder the data on Catalog pages, such as Reports, Data Sets, Data Sources and so on.

# Steps

1. On the main menu, click ⠿, then 🗄 **Catalog**.

   » The Catalog Home opens.
2. Click any of the items in the submenu, for example **Data Sets**.
3. Sort your data:

| Table display mode | Tile display mode (if available) |
| --- | --- |
| Click any column header to sort the data based on that column.<br>Click again to toggle between ascending and descending order.<br> | Click the **Sort by** arrow to sort ascending or descending, and click the drop-down list to select on which field you want to sort.<br> |

# Register a data source

Registering a data source makes metadata from that source available in Collibra to create data sets that can then be used for creating reports and analyzing data. Optionally, Collibra can perform data profilingdata profiling on the registered data and extract sample datasample data from it.

> Note   If you are using a Collibra Data Intelligence Cloud environment with an on-premises Jobserver, they must have the same installer version. You can find the installer version of your Collibra Data Intelligence Cloud environment at the bottom of the sign-in window of its Collibra Console, for example 5.7.13-0

## About registering a data source

By registering a data source, you connect a data source to Collibra. With this, you can make metadata of the data source available in Collibra.

You can register a data source via Jobserver or via Edge.

> **Note**
> When you enable registering a data source via Edge, you can choose to register a data source using Edge or using Jobserver.
>
> 

# Differences between registering a data source via Jobserver or via Edge

The following table shows the differences between registering a data source via Jobserver or via Edge.

| Part of process | Register a data source via Jobserver | Register a data source via Edge |
|---|---|---|
| Permissions | The required permissions to register a data source via Jobserver or via Edge are the same except for the following permission:<br><br>You need a resource role with the following resource permissions on the **Schema** community:<br><br>• Asset > add<br>• Attribute > add<br>• Domain > add<br>• Attachment > add | The required permissions to register a data source via Edge or via Jobserver are the same except for the following permission:<br><br>You need a global role with the View Edge connections and capabilities global permission. |

| Part of process | Register a data source via Jobserver | Register a data source via Edge |
|---|---|---|
| Registering a data source | When you register a data source via Jobserver, you have to enter all database connection properties in the **Register data source** dialog box. | Before you register a data source via Edge, you have to enable data source registration via Edge. You also create JDBC connections to your data source and Edge capabilities with a JDBC Catalog JDBC ingestion capability template. When you register a data source in Data Catalog, you can then select which database you want to add to the JDBC connection. |
| Refreshing or synchronizing | After registering a data source, a Schema asset is created. On the Configuration tab page of the Schema asset page, you can refresh a data source. | After registering a data source, a Database asset is created. The Database asset has a relation of the type "Technology asset groups / is grouped by Technology asset" to the System asset that was selected when registering the data source. On the **Configuration** tab page of the Database asset page, you can synchronize one or many schemas. |
| Profiling options | At the end of the registering process, you can select profiling options to create profiling and sample data. The profiling data is automatically created after the refresh process. | You have to enable profiling and classification via Edge. After registering a data source, you can select profiling options to create profiling data and data classes on a Database asset page. The metadata is profiled and classified automatically or manually after synchronizing a schema. |

# Difference between registering a data source and importing data

When you register a data source, Data Catalog reads and processes metadata of data sources that are not governed in Collibra Data Intelligence Cloud. Collibra will create assets of the relevant types, such as Database, Table and Column.

> Example   You register a data source that contains your financial data in a SAP HANA database. Afterwards, you can use the Collibra to manage the data, for example manage access control through data sets and use traceability to see your data lineage.

When you import data, you create or edit assets or complex relations, with their characteristics, from a view. Collibra will create assets of the type specified in the imported XLSX or CSV file.

> Example   You import an XLSX file containing the most common business terms of your company. You can use Collibra to approve the terms and link them to more technical assets.

# Naming convention

When you register a data source, Collibra follows a strict naming convention for the names of the new assets. Each asset has a display name and full name. You can freely edit the display name. However, you should never edit the full name, because Data Catalog may need it to refresh data sources. Editing the full name may cause unexpected results and break the synchronization process.

> Warning   Editing the full name of the Database and Schema assets may lead to errors during the refresh process.

# Supported data sources for data source registration

Collibra Data Intelligence Cloud supports several databases to register as a data source. Depending on your data source, you can use Collibra-provided Catalog connector, or your own JDBC driver.

# Your own JDBC drivers

For certain data sources, you can use your own JDBC driver. The following table contains the packaged data sources and versions that have been tested.

| Data source | Tested versions | Support for profiling and sample data | JDBC driver version |
|---|---|---|---|
| Amazon Redshift | 1.0 | Yes | v. 1.1.13.1013 |
| Cloudera Hive | 5.10 - 5.14 | No | Apache driver v. 1.2.1 |
| Hortonworks Hive | 2.5, 2.6 | No | Apache driver v. 1.2.1 |
| HP Vertica | 7.0 | Yes | v. 07.01.0200 |
| IBM DB2 | 10.5 | Yes | v. 4.9.78 |
| MySQL | 5.6, 5.7 | Yes | v. 5.1.38 |
| Oracle | 11g, 12c | Yes | v. 12c |
| PostgreSQL | 9.4, 9.5 | Yes | v. 9.4.1207 |
| Microsoft SQL Server | 2014, 2016 | Yes | v. 5.1.38<br><br>Note Only Microsoft drivers and drivers available via Collibra Marketplace are supported. |
| Teradata | 15.0, 16.20.07.01 | Yes | No driver tested |
| Snowflake | | | |

Note We cannot guarantee that other data sources or driver versions work correctly. If you use a generic JDBC driver or an unsupported version, data ingestion, data profiling and sample data may not work as expected.

## Authentication and permissions

Both ingestion andprofiling(including sampling and advanced data type detection) rely on JDBC drivers to operate. Those drivers authenticate to the data sources as a user registered in that data source with specific permissions attached to the user profile in the data source.

To ingest a database without profiling, Data Catalog requires read access to the database metadata: description of schema, tables, columns, including some more complex properties such as the primary and foreign keys.

However, if you enable one or more profiling options, Data Catalog also requires the permission to read the full table. Which permissions are required exactly depends on the data source type, version and configuration. Additionally, they can also differ according to the provider and version of the JDBC driver. Most of the queries required to retrieve the information above are hidden by the driver. As a consequence, Collibra cannot give a exhaustive list.

> Note   Collibra supports several authentication methods, including credentials, NTLM, CyberArk and Kerberos. If you are using a certified Collibra provided driver on the CollibraMarketplace, you can also authenticate using Windows Authentication.

> Tip   If you need more detailed information, we recommend to contact your JDBC driver provider.

> Warning   Currently, if you  register a data source viaEdge, you cannot profile your metadata and you can only use username/password authentication.

# Configuration assets

When you register a database or system as a data source, you enter connection properties and other options. To store the configuration and connection properties, Data Catalog creates a special kind of asset, often called the configuration asset. Some of these assets show parts of the configuration on a dedicated Configuration tab page.

This list contains the most widely used configuration assets:

- Schema assets, if you register a data source using Jobserver
- Database assets, if you register a data source using Edge.
- S3 File System assets
- Tableau Server assets

## Working with configuration assets

Even though you can import or export configuration assets with the import functionality or create them via the global create button, they would not contain any configuration. This means that, if you create a configuration asset in that way, you must also create the configuration and add it to the configuration asset. However, this is not possible for all configuration assets. For example, you cannot configure an S3 File System asset after creation. The only way to configure an S3 File System asset is by connecting to Amazon S3 and synchronizing its content. We highly recommend that you do not create configuration assets by importing them or via the global create button. Instead, use the appropriate procedure, such as registering a data source or registering a system.

> Warning   If you delete a configuration asset, you also delete its configuration. Register your data source or system again to create a new configuration asset or contact support for more information.

# Quartz Cron syntax

Cron is a software utility that specifies commands to run on a given schedule. This schedule is defined by a Cron pattern, which has a specific syntax that will be described in this section.

For example, you can refresh the schema of a data source or synchronize Tableau or Amazon S3 metadata outside office hours to reduce the impact of these actions on the performance of your environment.

> Note   By default, you use Spring Cron expressions to schedule Collibra Console back-ups.

> Warning   If you create an invalid Cron pattern, Collibra Data Intelligence Cloud stops responding.

The Cron pattern consists of six or seven space-separated fields:

```
<second> <minute> <hour> <day of the month> <month> <day of the
week> <year>
```

| Position | Field | Mandatory | Allowed values | Allowed special char- acters | Examples |
|---|---|---|---|---|---|
| 1 | second | Yes | 0-59 | , - * / | • *10*: at the 10th second.<br>• *\*/10*: every 10 seconds. |
| 2 | minute | Yes | 0-59 | , - * / | • *30*: at the 30th minute.<br>• *\*/15*: every 15 minutes.<br>• *5/10*: every 10 minutes start-ing at the 5th minute after the hour |
| 3 | hour | Yes | 0-23 | , - * / | • *10*: at 10 o'clock.<br>• *8-10*: at 8,9 and 10 AM.<br>• *6,18*: at 6 AM and at 6 PM. |

| Position | Field | Mandatory | Allowed values | Allowed special characters | Examples |
|---|---|---|---|---|---|
| 4 | day of the month | Yes | 1-31 | , - * ? / L W | • *3*: on the 3rd day of the month.<br>• *1-4*: every first four days of the month.<br>• *1,15*: the first day of the month and the 15th day of the month.<br>• *L*: on the last day of the month.<br>• *L-3*: on the third-to-last day of the month.<br>• *15W*: on the nearest weekday to the 15th of the month. If the 15th is a Saturday, then the trigger will be on the 14th, if the 15th is a Sunday, then the trigger will be on the 16th.<br><br>Note   If the 1st day of the month is a Saturday, then *1W* corresponds to the 3rd day of the month, since the month is specified in the 5th value of the Cron expression.<br><br>*LW*: on the last weekday of the month. |
| 5 | month | Yes | 1-12 or JAN-DEC | , - * / | • *12*: in December.<br>• *1-3*: every first three months of the year.<br>• *JUL,AUG*: every July and August.<br><br>Tip   The names of the months are not case-sensitive. |

| Position | Field | Mandatory | Allowed values | Allowed special characters | Examples |
|---|---|---|---|---|---|
| 6 | day of the week | Yes | 1-7 or SUN-SAT | , - * ? / L # | • *TUE*: every Tuesday.<br>• *2-6*: every weekday, Monday to Friday.<br>• *MON,WED,FRI*: every Monday, Wednesday and Friday.<br>• *L*: on Saturday, the 7th day of the week.<br>• *2L*: at the last Monday of the month.<br>• *6#3*: on the 3rd Friday of the month.<br><br>Tip   The names of the days are not case-sensitive. |
| 7 | year | No | empty, 1970-2099 | , - * / | • *<empty>*: if your schedule doesn't require a year, you can leave this value empty.<br>• *2021*: in 2021.<br>• *2021-2025*: in the years 2021, 2022, 2023, 2024 and 2025.<br>• *2021,2022,2025*: in the years 2021, 2022 and 2025. |

# Special characters

| Character | Description |
|---|---|
| * | Used to select all values within a field.<br><br>Example   * in the minute field corresponds with every minute. |

| Character | Description |
|---|---|
| ? | Used to specify something in one of the two fields in which the character is allowed, but not the other, mainly used for days of the week.<br><br>Example  If you want your trigger to fire on a particular day of the month, for example the 10th, but don't care what day of the week that happens to be, you could put "10" in the day-of-month field, and "?" in the day of the week field. |
| - | Used to specify ranges.<br><br>Example  `10-12` in the hour field means "the hours 10, 11 and 12". |
| , | Used to specify additional values.<br><br>Example  `MON,WED,FRI` in the day-of-week field means "the days Monday, Wednesday, and Friday". |
| / | Used to specify increments.<br><br>Example  `0/15` in the seconds field means "the seconds 0, 15, 30, and 45". And `5/15` in the seconds field means "the seconds 5, 20, 35, and 50".<br>You can also leave out the number before /, which is equivalent to having 0 before / .<br>`1/3` in the day-of-month field means "fire every 3 days starting on the first day of the month". |

| Character | Description |
|---|---|
| L | Has different meaning in each of the two fields in which it is allowed.<br><br>> Example   The value `L` in the **day-of-month field** means "the last day of the month" - day 31 for January, day 28 for February on non-leap years. You can also specify an offset from the last day of the month, such as "L-3" which would mean the third-to-last day of the calendar month.<br><br>If you use `L` in the **day-of-week field** by itself, it means "7" or "SAT". But if used in the day-of-week field after another value, it means "the last xxx day of the month" - for example "6L" means "the last Friday of the month".<br><br>When using the `L` option, it is important not to specify lists, or ranges of values, because you may get unexpected results. |
| W | Used to specify the weekday (Monday-Friday) nearest the given day.<br><br>> Example   `15W` in the value for the day-of-month field, means the nearest weekday to the 15th of the month:<br><br>• If the 15th is a Saturday, the trigger will fire on Friday the 14th.<br>• If the 15th is a Sunday, the trigger will fire on Monday the 16th.<br>• If the 15th is a Tuesday, then it will fire on Tuesday the 15th.<br><br>However if you specify `1W` as the value for day-of-month, and the 1st is a Saturday, the trigger will fire on Monday the 3rd, as it will not 'jump' over the boundary of a month's days. The 'W' character can only be specified when the value in the day-of-month field specifies a single day, not a range or list of days.<br><br>> Tip   The 'L' and 'W' characters can also be combined in the day-of-month field to yield 'LW', which translates to *"last weekday of the month"*. |
| # | Used to specify "the nth" XXX day of the month.<br><br>> Example   `6#3` in the day-of-week field means "the third Friday of the month" (day 6 = Friday and "#3" = the 3rd one in the month).<br>> Other examples: `2#1` is the first Monday of the month and `4#5` is the fifth Wednesday of the month. Note that if you specify `#5` and there is not 5 of the given day-of-week in the month, then no firing will occur that month. |

> Example
> - `0 0 * ? * * *` = the top of every hour of every day.
> - `*/10 * * * * ?` = every ten seconds.
> - `0 0 8-10 * * ? 2020` = 8, 9 and 10 o'clock of every day during the year 2020.
> - `0 0 6,19 ? * *` = 6:00 AM and 7:00 PM every day.
> - `0 0/30 8-10 ? * *` = 8:00, 8:30, 9:00, 9:30, 10:00 and 10:30 every day.
> - `0 0 9-17 * * MON-FRI` = on the hour nine-to-five weekdays.
> - `0 0 0 25 12 ?` = every Christmas Day at midnight, no matter what day of the week it is.
> - `0 15 10 ? * 6L 2022-2025` = 10:15 AM on every Friday of every month during the years 2022, 2023, 2024 and 2025.
> - `0 30 11 ? * 6#2` = 11:30 AM on the second Friday of every month.

> Warning   Quartz Cron only supports a value in either the 4th or the 6th position, but not in both. At the same time, both positions cannot be empty.

# Foreign key ingestion

A foreign key, in relational databases, is a field in one table that refers to the primary key of another table. A primary key is a column or combination of columns, to uniquely identify table records.

- The table with the primary key is referred to as the referenced table or parent table.
- The table with the foreign key is referred to as the child table.

## Ingesting foreign keys

In Data Catalog, a foreign key is ingested as an asset of the Foreign key type. See Foreign Key asset page.

The Foreign key asset creates relations between columns of different tables. It consists of foreign key mappings between the parent and child table.

In the following example, you see an overview of the tables, columns and a foreign key:

| | |
|---|---|
|  | • **food**: Schema that consists of two tables:<br>    ◦ **food_types**: Table with the columns food_name, food_code and food_type.<br>    ◦ **fastfood**: Table with the columns main_food_code, company_name and main_food_type.<br>• **fastfood_fk**: Foreign key asset consists of one or more foreign key mappings.<br>    ◦ mapping 1 (marked with blue arrows):<br>        ▪ Constrains the column **main_food_code** from the child table.<br>        ▪ References the column **food_code** from the parent table.<br>    ◦ mapping 2:<br>        ▪ Constrains the column **main_food_type** from the child table.<br>        ▪ References the column **food_type** from the parent table. |

# Registering a data source via Jobserver

By registering a data source via Jobserver, you connect a data source to Collibra. With this, you can make metadata of the data source available in Collibra.

During the data source registration process, you create a Schema asset. Via this asset, you can refresh the metadata of the data source.

> Tip   You can also register a data source via Edge.

# Data source ingestion steps

The following table shows the steps required for data source ingestion.

| Step | What? | Description |
|---|---|---|
| 1 | Register a data source | Registering a data source creates a connection between your data source and Collibra. It makes metadata of the data source available in Collibra. <br><br> Note   You can register a data source using a Collibra-provided driver or your own driver. |
| 2 | Ingestion | After registering a data source, Collibra creates a **Physical Data Dictionary** domain and new assets of the type Schema, Table and Column, corresponding to the data in your data source. <br><br> Note   Once you used a connection to successfully register a data source via Jobserver, you cannot change the connection properties. See Error when managing connection properties of a driver for Jobserver. |
| 3 | Refresh a data source | Refreshing the schema of a registered data source updates the metadata of the data source in Collibra. You typically do this when the data in a registered data source has been updated. <br><br> Tip   You can do this manually or automatically at fixed intervals. |

# Profiling data options

When you register your data source, you can choose profiling options for the registered data.

| Option | Description |
|---|---|
| Store Data Profile | Option to perform data profiling on the registered data. |
| Detect advanced data types | Option to detect advanced data types in the data source. |
| Store Sample Data | Option to extract sample data from the registered data. |

| Option | Description |
|---|---|
| Tables excluded from registration | Database tables that will not be ingested.<br><br>**Note**<br>• If required, you can exclude multiple tables. To do this, press *Enter* after typing a value and then type the next.<br>• You can use an asterisk (*) as wildcard to select multiple tables. For example, if you want to exclude the tables that all start with *act_*, you can enter *act_\**.<br>• The table names are case sensitive.<br>• You can add or remove tables from this list by refreshing the schema.<br>• The Table assets that are created after ingestion have an attribute type called Table Type that defines the type of table that is declared in the data source. For example, TABLE, VIEW,... |

# After registering a data source

When the registration is complete:

- A message at the top right tells you that data source registration is complete. A domain and Schema asset are immediately created and an ingestion job is started.
- You can immediately add the registered data source to a data set by clicking the corresponding link in the confirmation message.
- The ingestion job creates assets that represent the metadata of the data source.

  > Note   Table assets that are created after ingestion have an attribute type called Table Type that defines the type of table that is declared in the data source. For example, TABLE, VIEW,...

- A workflow to assign a technical steward to the new domain is started. This is a simple packaged workflow that you can edit to fit your organization's needs. When you have assigned a technical steward, that technical steward has to set the security classification and indicate whether the data elements contain personally identifiable information (PII).

# Register a data source using a Collibra-provided driver

You can register a database as a data source using one of the JDBC drivers provided by Collibra Marketplace.

> Tip   You can also do this with your own JDBC driver.

> Warning
> - This operation should only be executed by your database administrator.
> - The Collibra-provided drivers have been tested with Collibra Data Governance Center version 5.7.5. In older versions, you might encounter unexpected behavior.

## Steps

1. On the main menu, click ⠿, then 🖾 **Catalog**.

   » The Catalog Home opens.
2. In the main menu, click the **Create** (+) button.

   » The **Create** dialog box appears.
3. In the **Create** dialog box, click **Register data source (use a Collibra provided driver)**.
4. If there is no JDBC driver available, add and configure the driver of your preference.
5. In the **Register data source** dialog box, enter the required information.

   | Field | Description |
   | --- | --- |
   | Process on | The jobserver used for ingesting. |
   | Schema name | This name is used in Collibra as schema asset and must therefore be unique. |
   | Schema description | The description of the schema. This is used as description of the schema asset. |
   | Data owner | The owner of the registered data in Collibra. |

6. Click **Next**.
7. Enter the database connection properties.

| Option | Description |
|---|---|
| JDBC driver version | The JDBC driver to connect to your database.<br><br>Note   By default, you see the name of the driver that was used last. |
| Connect via | The jobserver used for ingesting. |
| <Configuration properties> | The connection properties as defined in your JDBC driver.<br><br>Note   For more information on the connection details of supported data sources, see JDBC connection details. |
| Store credentials | Select this option to store the credentials to access the database. With a schema refresh, you can clear this option again. |
| Username | Username to access the database.<br><br>Note   This field is ignored if your data source uses Cyberark, Kerberos or NTLM. |
| Password | Corresponding password to access the database.<br><br>Note   This field is ignored if your data source uses Cyberark, Kerberos or NTLM. |
| Schedule data refresh | Enable or disable a schedule to automatically refresh the data registration. |
| Cron pattern | Schedule of the data refresh as a Cron pattern.<br><br>If you create an invalid Cron pattern, Collibra Data Intelligence Cloud stops responding. |
| Time zone | The time zone of the database. |

> Note   If Collibra DGC cannot connect to the database, you cannot continue the data source registration wizard.

8. Click **Next**.
9. Select the data profiling options.

| Option | Description |
| --- | --- |
| Store Data Profile | Option to perform data profiling on the registered data. |
| Detect advanced data types | Option to detect advanced data types in the data source. |
| Store Sample Data | Option to extract sample data from the registered data. |
| Tables excluded from registration | Database tables that will not be ingested.<br><br>Note<br>○ If required, you can exclude multiple tables. To do this, press *Enter* after typing a value and then type the next.<br>○ You can use an asterisk (*) as wildcard to select multiple tables. For example, if you want to exclude the tables that all start with act_, you can enter *act_\**.<br>○ The table names are case sensitive.<br>○ You can add or remove tables from this list by refreshing the schema.<br>○ The Table assets that are created after ingestion have an attribute type called Table Type that defines the type of table that is declared in the data source. For example, TABLE, VIEW,... |

10. Click **Create**.

# What's next?

The data source is registered and the data is automatically ingested. The ingestion of data is executed in a job. You can see this job in the list of activities.

Click the **Result** button to open the data profiling results.

> **Tip**
> - If the database contains foreign keys, they will be registered as new assets of the Foreign Key asset type. Assets of this type contain the complex relation, which is the link between all column assets that are part of the foreign key definition.
>   However, the complex relation is not created if a column is part of a table that is added to the list of **Tables excluded from registration**.
> - If you exclude a table during the schema refresh, the corresponding table, column assets and foreign key mapping will be deleted.

# Manage Collibra-provided JDBC drivers

To register a database as a data source you need a JDBC driver. You can use one of the JDBC drivers provided by Collibra Marketplace.

This allows you to do the following:

- Edit an existing JDBC driver.
- Install a new JDBC driver for a data source type that has an existing JDBC driver, for example Oracle12c.
- Install a new JDBC driver for a data source type that doesn't have a JDBC driver yet, for example Amazon EMR.

> **Tip**   You can also do this with your own JDBC drivers.

> **Warning**
> - This operation should only be executed by your database administrator.
> - The Collibra-provided drivers have been tested with Collibra Data Governance Center version 5.7.5. In older versions, you might encounter unexpected behavior.

## Steps

1. On the main menu, click ⠿, then ⊟ **Catalog**.

   » The Catalog Home opens.
2. In the main menu, click the **Create** (﹢) button.

   » The **Create** dialog box appears.
3. In the **Create** dialog box, click **Register data source (use a Collibra provided driver)**.
4. If a JDBC driver is already installed for your data source, do the following:
   a. Enter the schema properties.

   | Field | Description |
   |---|---|
   | Schema name | This name is used in Collibra as schema asset and must therefore be unique. |
   | Schema description | The description of the schema. This is used as description of the schema asset. |
   | Data owner | The owner of the registered data in Collibra. |

   b. Click **Next**.

In the **JDBC driver version** field, click **manage drivers...**.

> Note   By default, you see the name of the driver that was used last.

c.

5. Do one of the following:
    - Click **Add JDBC Driver** if you want to create a new JDBC driver.
    - Click ✏ if you want to edit an existing JDBC driver.
6. Enter the required information.

| Field | Description |
| --- | --- |
| JDBC Driver Version Name | The name of the JDBC driver.<br><br>Tip   As a best practice, we recommend you use a strict naming convention which includes the data source and a version number. For example: Google BigQuery 1.5 or MySQL 5.9. |
| ⬆ Upload | Button to upload the relevant files for the data source.<br><br>Note   If you downloaded the JDBC driver from Collibra Marketplace, make sure to unzip the downloaded ZIP file before uploading it to Collibra Data Governance Center.<br><br>Note   The JDBC driver has to be in JAR format. |

| Field | Description |
|---|---|
| Driver files | This table contains a list of uploaded files. You can remove a driver file by clicking 🗑 . |

7. Click **Next**.
8. Configure the JDBC connection.

> Note   For more information on the connection details of supported data sources, see JDBC connection details.

9. Click **Create**.

## What's next?

You can now complete the data source registration wizard for Collibra-provided JDBC drivers.

# Register a data source using your own driver

You can register a database as a data source using one of your own drivers.

> Tip   You can also do this with a Collibra-provided JDBC driver.

 This operation should only be executed by your database administrator.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have set up the JDBC driver of your source data, for example MySQL.
- You have configured one or more Jobservers in Collibra Console. If there is no available Jobserver, the **Register data source** actions will be grayed out in the global create menu of Collibra Data Intelligence Cloud.
- If you are using a Collibra Data Intelligence Cloud environment with an on-premises Jobserver, both must have the same installer version. You can find the installer

version of your Collibra Data Intelligence Cloud environment at the bottom of the sign-in window of its Collibra Console, for example 5.7.13-0

- You have a resource role with the following resource permissions on the **Schema** community:
  - Asset > add
  - Attribute > add
  - Domain > add
  - Attachment > add
- You have the permissions to retrieve the metadata of the following database components through the JDBC Driver Database Metadata methods:
  - Schemas
  - Tables
  - Columns
  - Primary keys
  - Foreign keys

> Note
> - For the list of supported databases and versions, consult the Databases supported versions section.
> - For the JDBC connection details of the various databases, consult the JDBC connection details section.

## Steps

1. On the main menu, click ⠿, then ▣ **Catalog**.

   » The Catalog Home opens.
2. In the main menu, click the **Create** (+) button.

   » The **Create** dialog box appears.
3. In the **Register data source** dialog box, click the type of your data source.
4. If there is no JDBC driver available, add and configure the driver of your preference.
5. In the **Register data source** dialog box, enter the required information.

| Field | Description |
|---|---|
| Process on | The jobserver used for ingesting. |
| Schema name | This name is used in Collibra as schema asset and must therefore be unique. |

| Field | Description |
|---|---|
| Schema description | The description of the schema. This is used as description of the schema asset. |
| Data owner | The owner of the registered data in Collibra. |

6. Click **Next**.
7. Enter the database connection properties.

| Option | Description |
|---|---|
| JDBC driver version | The JDBC driver to connect to your database. |
| Connect via | The jobserver used for ingesting. |
| Database | Name of the database. This field is not available for all data sources. |
| Host | Hostname to access the database. |
| Port | Port to access the database. |

| Option | Description |
|---|---|
| <Configuration properties> | The connection properties as defined in your JDBC driver.

> Note   For more information on the connection details of supported data sources, see JDBC connection details.

If you want to use Kerberos authentication, you also need the following connection properties.

| Label | Description |
|---|---|
| Principal | The Kerberos principal identity. |
| Kerberos realm | The Kerberos realm name. |
| Login context name | The login context name that is used as the index to the configuration. |
| Jaas file name | The name of the Jaas file. |
| Kerberos configuration file | The configuration file containing specific properties for Kerberos authentication. |

If you want to use NTLM authentication, you also need the following connection properties.

| Label | Description |
|---|---|
| Security | The security that enables the authentication |
| Authentication scheme | The used authentication scheme, which is NTLM. |

If you want to use CyberArk authentication, you need the following connection properties. |

| Option | Description |
|---|---|
| | <table><tr><th>Label</th><th>Description</th></tr><tr><td>Keystore file</td><td>The name of the keystore file. The keystore must contain the client key and client certificate or certificate chain.<br><br>If `defaultTruststore` is set to `false`, the keystore has to contain the trusted CA certificate needed to validate the server certificate offered by CyberArk.<br><br>The value must have the following format: `file://<keystore-file name.jks>`.<br><br>**Example**<br>`file://cyberark-keystore.jks`</td></tr><tr><td>Keystore password</td><td>The password required to open the keystore.</td></tr><tr><td>Default truststore</td><td>The indication of the default truststore. The default value is set to `False`.<br><br>○ `False`: The certificate is validated through the keystoreFile property.<br>○ `True`: The certificate is validated through the default truststore from the Java JRE. This is recommended when CyberArk is set up to offer a server certificate that can be validated by a public CA (certification authority).</td></tr></table> |

| Option | Description |
|---|---|
| | <table><tr><th>Label</th><th>Description</th></tr><tr><td>CyberArk address</td><td>The host and port number through which the CyberArk server is accessible. The format of the address is `hostname:port.`<br><br>Example<br>`my.cyberark.com:5502`</td></tr><tr><td>CyberArk application ID</td><td>The application ID as defined in CyberArk.<br><br>This ID should be provided by your network or system administrator.</td></tr><tr><td>CyberArk query</td><td>The CyberArk query.<br><br>This query should be provided by your network or system administrator.</td></tr></table> |
| Store credentials | Select this option to store the credentials to access the database. With a schema refresh, you can clear this option again. |
| Username | Username to access the database.<br><br>Note This field is ignored if your data source uses any authentication method other than credentials. |
| Password | Corresponding password to access the database.<br><br>Note This field is ignored if your data source uses any authentication method other than credentials. |
| Schedule data refresh | Enable or disable a schedule to automatically refresh the data registration. |

| Option | Description |
|---|---|
| Cron pattern | Schedule of the data refresh as a Quartz Cron pattern.<br><br>**Warning**  If you create an invalid Cron pattern, Collibra Data Intelligence Cloud stops responding. |
| Time zone | The time zone of the database. |

**Note**  If Collibra cannot connect to the database, you cannot continue the data source registration wizard.

8. Click **Next**.
9. Select the data profiling options.

| Option | Description |
|---|---|
| Store Data Profile | Option to perform data profiling on the registered data. |
| Detect advanced data types | Option to detect advanced data types in the data source. |
| Store Sample Data | Option to extract sample data from the registered data. |
| Tables excluded from registration | Database tables that will not be ingested.<br><br>**Note**<br>○ If required, you can exclude multiple tables. To do this, press *Enter* after typing a value and then type the next.<br>○ You can use an asterisk (*) as wildcard to select multiple tables. For example, if you want to exclude the tables that all start with act_, you can enter *act_\**.<br>○ The table names are case sensitive.<br>○ You can add or remove tables from this list by refreshing the schema.<br>○ The Table assets that are created after ingestion have an attribute type called Table Type that defines the type of table that is declared in the data source. For example, TABLE, VIEW,... |

10.  Click **Create**.

## What's next?

The data source is registered and the data is automatically ingested. The ingestion of data is executed in a job. You can see this job in the list of activities.

| | Started ▼ | Name | Status | Finished | Results | |
|---|---|---|---|---|---|---|
| Overview | | | | | | |
| Groups | ☐ 12/12/2017 2:04 PM | Export to "Default.csv". | Completed | 12/12/2017 2:04 PM | Result | |
| Responsibilities | ☐ 12/12/2017 1:29 PM | Updating JDBC schema | Completed | 12/12/2017 1:29 PM | Result | |
| History | ☐ 12/12/2017 1:29 PM | Updating JDBC schema | Completed | 12/12/2017 1:29 PM | Result | |
| **Activities** | ☐ 12/12/2017 1:27 PM | Creating schema from JDBC | Completed | 12/12/2017 1:28 PM | Result | ← |
| | ☐ 12/12/2017 1:18 PM | Creating schema from file | Completed | 12/12/2017 1:19 PM | Result | |

Click the **Result** button to open the data profiling results.

> **Tip**
> - If the database contains foreign keys, they will be registered as new assets of the **Foreign Key** asset type. Assets of this type contain the complex relation, which is the link between all column assets that are part of the foreign key definition.
> However, the complex relation is not created if a column is part of a table that is added to the list of **Tables excluded from registration**.
> - If you exclude a table during the schema refresh, the corresponding table, column assets and foreign key mapping will be deleted.

## Register an Excel file as data source

> Note   If you are using a Collibra Data Intelligence Cloud environment with an on-premises Jobserver, they must have the same installer version. You can find the installer version of your Collibra Data Intelligence Cloud environment at the bottom of the sign-in window of its Collibra Console, for example 5.7.13-0

## Prerequisites

- You have downloaded an Excel file.
- You have configured one or more Jobservers in Collibra Console. If there is no available Jobserver, the **Register data source** actions will be grayed out in the global create menu of Collibra Data Intelligence Cloud.
- You have a resource role with the following resource permissions:
    - Asset > add
    - Attribute > add
    - Domain > add
    - Attachment > add

## Steps

1. On the main menu, click ⠿, then ⊜ **Catalog**.

    » The Catalog Home opens.
    Or open any asset of the type Schema, Data Set, Table, Column or Tableau Server.
2. In the main menu, click the **Create** (＋) button.

    » The **Create** dialog box appears.
3. In the **Create** dialog box, click **Register data source (use your own driver)**.

    » The **Register data source (use your own driver)** dialog box appears.
4. In the **Register data source** dialog box, click **Excel**.
5. Enter the data source configuration.

    | Field | Description |
    | --- | --- |
    | Process on | The jobserver used for ingesting. |
    | Schema name | This name is used in Collibra as schema asset and must therefore be unique. |
    | Schema description | The description of the schema. This is used as description of the schema asset. |
    | Data owner | The owner of the registered data in Collibra. |
6. Click **Next**.

7. Select the data profiling options.

| Option | Description |
| --- | --- |
| Store Data Profile | Option to perform data profiling on the registered data. |
| Detect advanced data types | Option to detect advanced data types in the data source. |
| Store Sample Data | Option to extract sample data from the registered data. |
| Tables excluded from registration | Database tables that will not be ingested.<br><br>Note<br>○ If required, you can exclude multiple tables. To do this, press *Enter* after typing a value and then type the next.<br>○ You can use an asterisk (*) as wildcard to select multiple tables. For example, if you want to exclude the tables that all start with act_, you can enter *act_*.<br>○ The table names are case sensitive.<br>○ You can add or remove tables from this list by refreshing the schema.<br>○ The Table assets that are created after ingestion have an attribute type called Table Type that defines the type of table that is declared in the data source. For example, TABLE, VIEW,... |

8. Click **Create**.

# What's next?

The data source is registered and the data is automatically ingested. The ingestion of data is executed in a job. You can see this job in the list of activities.

Click the **Result** button to open the data profiling results.

If you have selected the option to perform data profiling and/or extract sample data, you can go to the schema page to verify if this process has completed in the **Synchronization Status** field. Refresh the schema page until the **Synchronization Status** field has disappeared.

Note that there Collibra may have resolved some small issues:

| Use case | Behavior |
|---|---|
| Missing column name | If the file is missing a column name, a default name will be given, _c + index. The index is the column position in the file starting with 0. For example, **_c4** corresponds with the fifth column in the file. |
| Duplicate column name | If the file has duplicate column names, the column names will be appended with an index. The index is the column position in the file, starting with 0. For example, **mycol1** and **mycol3** are columns 2 and 4 in the file, each with the column name **mycol**. |
| Empty sheet | If the Excel file has empty sheets, they are not registered. |

# Register a CSV file as data source

> Note   If you are using a Collibra Data Intelligence Cloud environment with an on-premises Jobserver, they must have the same installer version. You can find the installer version of your Collibra Data Intelligence Cloud environment at the bottom of the sign-in window of its Collibra Console, for example 5.7.13-0

## Prerequisites

- You have downloaded a CSV file.
- You have configured one or more Jobservers in Collibra Console. If there is no available Jobserver, the **Register data source** actions will be grayed out in the global create menu of Collibra Data Intelligence Cloud.
- You have a resource role with the following resource permissions:
  - Asset > add
  - Attribute > add
  - Domain > add
  - Attachment > add

## Steps

1. On the main menu, click ⠿, then ⊜ **Catalog**.

   » The Catalog Home opens.

   Or open any asset of the type Schema, Data Set, Table, Column or Tableau Server.
2. In the main menu, click the **Create** ( + ) button.

   » The **Create** dialog box appears.
3. In the **Create** dialog box, click **Register data source (use your own driver)**.

   » The **Register data source (use your own driver)** dialog box appears.
4. In the **Register data source** dialog box, click **Csv**.
5. Enter the data source configuration.

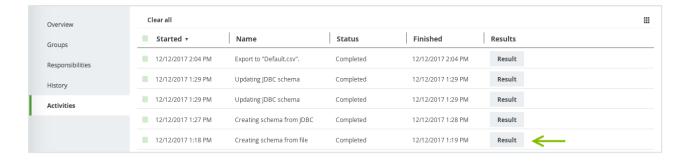   | Field | Description |
   | --- | --- |
   | Process on | The jobserver used for ingesting. |
   | Schema name | This name is used in Collibra as schema asset and must therefore be unique. |
   | Schema description | The description of the schema. This is used as description of the schema asset. |
   | Data owner | The owner of the registered data in Collibra. |
6. Click **Next**.

7. Select the data profiling options.

| Option | Description |
| --- | --- |
| Store Data Profile | Option to perform data profiling on the registered data. |
| Detect advanced data types | Option to detect advanced data types in the data source. |
| Store Sample Data | Option to extract sample data from the registered data. |
| Tables excluded from registration | Database tables that will not be ingested.<br><br>**Note**<br>○ If required, you can exclude multiple tables. To do this, press *Enter* after typing a value and then type the next.<br>○ You can use an asterisk (*) as wildcard to select multiple tables. For example, if you want to exclude the tables that all start with act_, you can enter *act_\**.<br>○ The table names are case sensitive.<br>○ You can add or remove tables from this list by refreshing the schema.<br>○ The Table assets that are created after ingestion have an attribute type called Table Type that defines the type of table that is declared in the data source. For example, TABLE, VIEW,... |

8. Click **Create**.

# What's next?

The data source is registered and the data is automatically ingested. The ingestion of data is executed in a job. You can see this job in the list of activities.

Click the **Result** button to open the data profiling results.

> **Note**
> * Empty rows in the CSV file are ignored. As a consequence, they do not count towards the row count or missing value count.
> * You can define the format of empty values by configuring the data profiling behavior. However, if a field is empty in the CSV file, it will be considered empty even if it does not match the format defined in the configuration.

If you selected the option to perform data profiling and/or extract sample data, you can verify that the process was completed in the Synchronization Status field on the schema asset page. Refresh the schema page until the **Synchronization Status** field disappears.

Note that there Collibra may have resolved some small issues:

| Use case | Behavior |
|---|---|
| Missing column name | If the file is missing a column name, a default name will be given, _c + index.<br><br>The index is the column position in the file starting with 0.<br><br>For example, **_c4** corresponds with the fifth column in the file. |
| Duplicate column name | If the file has duplicate column names, the column names will be appended with an index.<br><br>The index is the column position in the file, starting with 0.<br><br>For example, **mycol1** and **mycol3** are columns 2 and 4 in the file, each with the column name **mycol**. |
| Empty sheet | If the Excel file has empty sheets, they are not registered. |

# Manage your own JDBC drivers

To register a database as a data source you need a JDBC driver. You can use one of your own JDBC drivers.

For more information, see Supported data sources for data source registration.

This allows you to do the following:

- Edit an existing JDBC driver.
- Install a new JDBC driver for a data source type that has an existing JDBC driver, for example Oracle12c.
- Install a new JDBC driver for a data source type that doesn't have a JDBC driver yet, for example Amazon EMR.

> Tip   You can also do this with a Collibra-provided JDBC driver that you download from Collibra Marketplace.

This operation should only be executed by your database administrator.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have downloaded the JDBC driver of your choice as an archive file (for example, ZIP or JAR).
- You have configured one or more Jobservers in Collibra Console. If there is no available Jobserver, the **Register data source** actions will be grayed out in the global create menu of Collibra Data Intelligence Cloud.
- You have a resource role with the following resource permissions on the **Schema** community:
  - Asset > add
  - Attribute > add
  - Domain > add
  - Attachment > add

# Steps

1. On the main menu, click ⠿, then ⊟ **Catalog**.
   - » The Catalog Home opens.
2. In the main menu, click the **Create** (＋) button.
   - » The **Create** dialog box appears.
3. In the **Create** dialog box, click **Register data source (use your own driver)**.
4. In the **Register data source** dialog box, click the type of your data source.



5. If a JDBC driver is already installed for your data source:
   a. Enter the schema properties.

| Field | Description |
|---|---|
| Schema name | This name is used in Collibra as schema asset and must therefore be unique. |
| Schema description | The description of the schema. This is used as description of the schema asset. |
| Data owner | The owner of the registered data in Collibra. |

   b. Click **Next**.

      c.  In the **JDBC driver version** field, click **manage drivers...**.



6.  Do one of the following:

    ◦  Click **Add JDBC Driver** if you want to create a new JDBC driver.

    ◦  Click ✏ if you want to edit an existing JDBC driver.

7.  Enter the required information.

| Field | Description |
|---|---|
| JDBC Driver Version Name | The name of the JDBC driver.<br><br>Tip   As a best practice, we recommend you use a strict naming convention which includes the data source and a version number. For example: Google BigQuery 1.5 or MySQL 5.9. |
| ⬆ Upload | Button to upload the relevant files for the data source.<br><br>The JDBC driver should be in JAR or ZIP format with a valid Java archive structure.<br><br>For authentication with CyberArk, you also need to upload a keystore file in JKS format.<br><br>Note   When you click the button, an **Open** dialog box appears. By default, the dialog box filters on JAR, ZIP and CONF files. However, you can change the filter to show all files.<br><br>For Hortonworks Hive with Kerberos authentication, you need two files: **jaas.conf** and **krb5.conf**. |

| Field | Description |
|---|---|
| Driver files | This table contains a list of uploaded files.<br><br>You can remove a driver file by clicking 🗑 . |

8. Click **Next**.
9. Configure the JDBC connection.

> Note   For more information on the connection details of supported data
> sources, see JDBC connection details of your own drivers.

10. Click **Create**.

## What's next?

You can now complete the data source registration wizard.

## JDBC connection details of your own drivers

In this section, you will see the connection details needed to register a data source or manage your own JDBC driver.

> Note   About the **Connection properties** table:
>
> - The **Label** column is the value that will appear in the connection details dialog box of the **Data Source Registration** wizard.
> - The **Property** column contains the parameters in which the user input will be saved.

## Amazon Redshift

| Label | Property | Mandatory |
|---|---|---|
| Hostname | host | Yes |

| Label | Property | Mandatory |
|-------|----------|-----------|
| Port | port | Yes |
| Database | database | Yes |
| Schema | schema | Yes |

## Cloudera Hive

| Label | Property | Mandatory |
|-------|----------|-----------|
| URL (hostname:port) | host | Yes |
| Principal | principal | Yes |
| Schema | schema | Yes |

## Hortonworks Hive

| Label | Property | Mandatory |
|-------|----------|-----------|
| URL (hostname:port) | host | Yes |
| Schema | schema | Yes |

## HP Vertica

| Label | Property | Mandatory |
|-------|----------|-----------|
| Hostname | host | Yes |
| Port | port | Yes |
| Database | database | Yes |
| Schema | schema | Yes |

## IBM DB2

| Label | Property | Mandatory |
|-------|----------|-----------|
| Hostname | host | Yes |
| Port | port | Yes |
| Database | database | Yes |
| Schema | schema | Yes |

## MapR Hive

| Label | Property | Mandatory |
|-------|----------|-----------|
| URL (hostname:port) | host | Yes |
| Schema | schema | Yes |

## Microsoft SQL Server

| Label | Property | Mandatory |
|-------|----------|-----------|
| Hostname | host | Yes |
| Port | port | Yes |
| Database | databaseName | Yes |
| Schema | schema | Yes |

## MySQL

| Label | Property | Mandatory |
|-------|----------|-----------|
| Hostname | host | Yes |
| Port | port | Yes |
| Database | database | Yes |

## Oracle DB

| Label | Property | Mandatory |
|-------|----------|-----------|
| Hostname | host | Yes |
| Port | port | Yes |
| SID | sid | Yes |
| Schema | schema | Yes |

## PostgreSQL

| Label | Property | Mandatory |
|-------|----------|-----------|
| Hostname | host | Yes |
| Port | port | Yes |
| Database | database | Yes |
| Schema | schema | Yes |

## Teradata

| Label | Property | Mandatory |
|---|---|---|
| Hostname | host | Yes |
| Port | port | Yes |
| Database | database | Yes |
| Schema | schema | Yes |

## Authentication methods

Certain authentication methods require additional connection properties.

# NTLM

If you want to use NTLM authentication, you also need the following connection properties.

| Label | Property | Mandatory |
|---|---|---|
| Security | *integratedSecurity* must be value `True`. | Yes |
| Authentication scheme | *authenticationScheme* must be value *NTLM*. | Yes |

# Kerberos

If you want to use Kerberos authentication, you also need the following connection properties.

| Label | Property | Mandatory |
|---|---|---|
| Principal | principal | Yes |
| Kerberos realm | realm | Yes |

| Label | Property | Mandatory |
|---|---|---|
| Login context name | loginContextName<br><br>You can find the value for this property in the jaas.conf file. | Yes |
| Jaas file name | com.collibra.jobserver.dto.catalog.JdbcConnection.jaasConfig | Yes |
| Kerberos configuration file | com.collibra.jobserver.dto.catalog.JdbcConnection.krbConfig | Yes |

# Cyberark

If you want to use CyberArk authentication, you need the following connection properties. If you use one of the CyberArk connection properties, Data Catalog automatically uses CyberArk authentication.

| Label | Property | Mandatory |
|---|---|---|
| Keystore file | keystoreFile | Yes |
| Keystore password | keystorePass | Yes |
| Default truststore | defaultTruststore | No |
| CyberArk address | cyberarkAddress | Yes |
| CyberArk application ID | cyberarkAppId | Yes |
| CyberArk query | cyberarkQuery | Yes |

# Authentication

If you register a database as data source or manage a JDBC driver, you can use various authentication methods to access your data source.

# CyberArk authentication

CyberArk is middleware to manage authentication and is used to provide access to various data sources. You can use CyberArk to let Data Catalog access and ingest data sources with username and password authentication.

> Note   You can only authenticate to data sources using username and password authentication.

## Setting up CyberArk authentication

You set up CyberArk authentication when you register your data source or manage your JDBC driver. When you register your data source or manage your JDBC driver, you only provide the username, the password you need to authenticate to the data source is stored in CyberArk and is retrieved by the Jobserver. When you ingest a data source using CyberArk authentication, the Jobserver uses certificate-based mutual authentication to authenticate to CyberArk.

> Note   The connection to CyberArk is only supported over HTTPS.

To authenticate via CyberArk, you have to enable CCP WebService in CyberArk and keep the default name AIMWebService unchanged. You also have to provide your own CyberArk certificates via a JKS keystore that you upload to Collibra when you register your data source or manage your JDBC driver. The JKS keystore contains the CyberArk client certificates, the private key and, if required, a server certificate.

# Authentication workflow



| Step | Action |
|---|---|
| 1 | The Jobserver requests credentials from CyberArk through a certificate-based mutual authentication. |
| 2 | CyberArk provides the Jobserver with a username and password. |
| 3 | The Jobserver uses these credentials to authenticate to a data source. |

# Configuration

If you want to use CyberArk authentication, you need the following connection properties. If you use one of the CyberArk connection properties, Data Catalog automatically uses CyberArk authentication.

| Label | Property | Description | Mandatory |
|---|---|---|---|
| Keystore file | keystoreFile | The name of the keystore file. The keystore must contain the client key and client certificate or certificate chain.<br><br>If `defaultTruststore` is set to `false`, the keystore has to contain the trusted CA certificate needed to validate the server certificate offered by CyberArk.<br><br>The value must have the following format: `file://<keystore-file name.jks>`.<br><br>> Example `file://cyberark-keystore.jks` | Yes |
| Keystore password | keystorePass | The password required to open the keystore. | Yes |
| Default truststore | defaultTruststore | The indication of the default truststore. The default value is set to `False`.<br><br>• `False`: The certificate is validated through the keystoreFile property.<br>• `True`: The certificate is validated through the default truststore from the Java JRE. This is recommended when CyberArk is set up to offer a server certificate that can be validated by a public CA (certification authority). | No |
| CyberArk address | cyberarkAddress | The host and port number through which the CyberArk server is accessible. The format of the address is `hostname:port`.<br><br>> Example `my.cyberark.com:5502` | Yes |

| Label | Property | Description | Mandatory |
|-------|----------|-------------|-----------|
| CyberArk application ID | cyberarkAppId | The application ID as defined in CyberArk.<br><br>This ID should be provided by your network or system administrator. | Yes |
| CyberArk query | cyberarkQuery | The CyberArk query.<br><br>This query should be provided by your network or system administrator. | Yes |

# NTLM authentication

NTLM is an authentication protocol used on networks that include systems running the Windows operating system and on stand-alone systems. It uses a challenge-response authentication to connect to the Microsoft SQL Server data source. For more information, see the Microsoft NTLM user guide.

If you have a Microsoft SQL Server data source that uses NTLM authentication, you have to set up specific connection properties when you register the data source or manage the JDBC driver.

## Authentication workflow

When you ingest a Microsoft SQL Server data source using NTLM authentication, the Jobserver connects to the server to request access. The server then sends a challenge for the Jobserver to encrypt and send back. The domain controller validates that response and gives the Jobserver access to the data source.

| Step | Action |
|------|--------|
| 1 | The Jobserver requests access to the Microsoft SQL Server data source. |
| 2 | The server sends a challenge message to the Jobserver to identify the Jobserver. |
| 3 | The Jobserver sends a response back to the server. |
| 4 | The server sends the challenge and response message to the domain controller. |
| 5 | The active directory on the domain controller validates the challenge and response message and sends the result to the server. |
| 6 | The server gives the Jobserver permission to access the data source. |

# Configuration

If you want to use NTLM authentication, you also need the following connection properties.

| Label | Property | Description | Mandatory |
|---|---|---|---|
| Security | *integratedSecurity* must be value `True.` | The security that enables the authentication | Yes |
| Authentication scheme | *authenticationScheme* must be value *NTLM.* | The used authentication scheme, which is NTLM. | Yes |

# Kerberos authentication

You can use Kerberos authentication for registering a Hive data source, for example Cloudera Hive, Hortonworks Hive or MapR Hive.

## Authentication type

We only support Kerberos username and password authentication, not keytab. Ensure that you configure this in the **jaas.conf** file by setting the **useKeyTab** option to *false*.

In the following jaas.conf example, **Client** is the value of the **loginContextName** field when you configure the Kerberos connection configuration.

**Example**

```
Client {
com.sun.security.auth.module.Krb5LoginModule required
useKeyTab=false
useTicketCache=true;
};
```

If there are multiple entries in this configuration file, ask the database administrator or network administrator which one to use. For more information about the Jaas login configuration file, see the Java documentation.

## Example krb5.conf

The following is an example configuration file of Kerberos.

```
[libdefaults]
  renew_lifetime = 7d
  forwardable = true
  default_realm = MY.REALM
  ticket_lifetime = 24h
  dns_lookup_realm = false
  dns_lookup_kdc = false
  default_ccache_name = /tmp/krb5cc_%{uid}

[logging]
  default = FILE:/var/log/krb5kdc.log
  admin_server = FILE:/var/log/kadmind.log
  kdc = FILE:/var/log/krb5kdc.log

[realms]
  MY.REALM = {
    kdc = <kdc.my.realm>
    admin_server = <kadmin.my.realm>
  }
```

# Enable debug for Kerberos authentication issues

If an error occurs during the Kerberos authentication, you can enable debugging to track the root cause of the error.

To enable debugging for the Kerberos authentication:

1. On the server that hosts the Jobserver service, open the file **context_jvm.conf** in **<drive>/collibra/spark-jobserver/conf** for editing.
2. Is the following parameter present in the file: `-Dsun.security.krb5.debug`
    - Yes: Set its value to *true*.
    - No: Add the following line to the file: `-Dsun.security.krb5.debug=true`
3. Save and close the file.
4. Restart the Jobserver service.

The default log file in which to look for Kerberos authentication issues is **<drive>/collibra_data/logs/context_<context-name>/spark-job-server.log**.
In general, you list the **context_<context-name>** directories and pick the most recent one.

> Tip   After resolving the authentication issues, set the parameter to *false*.

# Cancel a data ingestion job

If you are the one that started the data ingestion job, you can cancel it while the data ingestion job is still running.

## Prerequisites

- You have registered a data source.
- You have started the ingestion job.

## Steps

1. On the main menu, click ⟳, then **Show more**.
   » Your profile page opens on the **Activities** tab page.
2. Click ⊗ next to the ingestion job to cancel it.

   > Note   When the job is finished, the ⊗ icon changes into a 🗑 icon. You can't cancel the ingestion job anymore.

» The data ingestion job is canceled.

# About refreshing a schema

Refreshing a schema is the process of updating the metadata of a registered data source in Collibra Data Intelligence Cloud.

You can refresh a schema manually or automatically at fixed intervals. This is particularly useful if the content of the data source changes regularly.

In this section, you can find the relevant actions to successfully refresh a schema.

# Refresh the schema of a registered data source

You can refresh a schema of registered data to update the data and the profiling. It can also be useful to do this to change data types to force the profiling to use the correct type.

> Tip  You can also refresh the schema automatically via a schedule.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have set up the JDBC driver of your source data, for example MySQL.
- You have configured one or more Jobservers in Collibra Console. If there is no available Jobserver, the **Register data source** actions will be grayed out in the global create menu of Collibra Data Intelligence Cloud.
- If you are using a Collibra Data Intelligence Cloud environment with an on-premises Jobserver, both must have the same installer version. You can find the installer version of your Collibra Data Intelligence Cloud environment at the bottom of the sign-in window of its Collibra Console, for example 5.7.13-0
- You have a resource role with the following resource permissions on the **Schema** community:
    - Asset > add
    - Attribute > add
    - Domain > add
    - Attachment > add
- You have the permissions to retrieve the metadata of the following database components through the JDBC Driver Database Metadata methods:
    - Schemas
    - Tables
    - Columns
    - Primary keys
    - Foreign keys

> Note
>   - For the list of supported databases and versions, see Databases supported
>     versions.
>   - For the JDBC connection details of the various databases, see JDBC
>     connection details.

## Steps

1. Open the Schema asset.
   a. On the main menu, click ⊞, then ⊟ **Catalog**.

      » The Catalog Home opens.
   b. In the submenu, click **Data Dictionary** and select the **All Schemas** view.
   c. Click the schema that you want to refresh.

   > Tip   You can also use the Collibra Data Intelligence Cloud search function to
   > look up your schema.

2. In the view bar, to the right, click **Actions** → **Refresh**.
   » The **Refresh Schema** dialog box appears.

   > Tip   If Catalog experience is disabled, the **More** menu is shown instead of
   > **Actions**.

3. Enter the required information.
   This dialog box varies with the data source:
   - Relational database

     > Note
     >   - If you exclude a table during the schema refresh, you will delete
     >     the corresponding table, column assets and the foreign key
     >     mapping (complex relation).
     >   - If you clear the **Store credentials** option, the credentials are no
     >     longer stored.

   - CSV file
   - Excel file

   This step may take some time.

4. Click **Save & Refresh**.
   » The refresh of the schema starts, you can follow the refresh job in the list of activities.

## What's next?

- The representation of the schema is updated: Data Catalog creates, edits and deletes assets as needed.
  - This can lead to refresh conflicts. See Resolve schema refresh conflicts via Jobserver.
  - If you had deleted assets manually, Data Catalog usually doesn't create them again if you refresh the schema. However, if the assets are required to represent the schema structure, Data Catalog can create them again.

> Example
> You ingested a schema that contains a table and three columns. In Data Catalog, this is represented by a Schema asset, a Table asset and three Column assets.
>
> Additionally, the following relations are created between the relevant assets:
>
> - Schema contains/is part of Table
> - Table contains/is part of Column
>
> In the actual data source, the columns are physically inside the table. However, in Data Catalog, they are separate assets linked by relations. As a consequence, you can delete the Table asset without deleting the Column assets. If you did that, Data Catalog creates the Table asset again if you refresh the schema, because the Table asset is needed for the relations to the Column assets.

- If the data source has new values and you selected the checkboxes to store sample data and data profile information, new sample data is generated and all profiling information is updated.
  If you did not select the **Store Sample Data** checkbox, any previously gathered sample data is removed. If you did not select the **Store Data Profile** checkbox, any previously gathered data profiling information is removed.

- Data types or categorical attributes that you changed manually are not updated when you refresh the schema.

    > Note   If you change the data type back to the original value assigned by the profiler, Data Catalog can update it if you refresh the schema.

- If you use this schema of the data source for Tableau stitching, you have to restitch after each schema refresh to make sure that all relations are up to date.

# Schedule a schema refresh

You can refresh a schema manually, but you can also create a schedule to refresh a schema on a regular basis.

You can only create a refresh schedule for schemas of databases that are registered as a data source, not from CSV or Excel files.

> Tip   You can schedule the refresh during the data source registration process or afterwards via the Schema asset.

> Note
> - To enable a scheduled schema refresh, you have to save the credentials in the configuration of a data source registration.
> - The refresh schedule uses Quartz Cron expressions.
> - If you use the schema for Tableau stitching, you have to restitch after each schema refresh to make sure that all relations are up-to-date.

## Prerequisites

- You have registered a data source.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a role with the following resource permissions on the **Schema** community:
    - Asset: add
    - Attribute: add

- ○ Domain: add
- ○ Attachment: add

> Note   These permissions are always necessary when registering a data source.

# Schedule the refresh during the data source registration process

You can create the refresh schedule when you register a data source.

> Example
> When you register a Snowflake data source in Collibra Data Intelligence Cloud, you can create a refresh schedule by selecting **Schedule data refresh**. You can then enter the CRON pattern *0 0 12?\*WED* to refresh every Wednesday at 12:00:00 PM.
>
> 

# Schedule the refresh via the Schema asset

You can create the refresh schedule when you refresh the schema of a registered data source via the Schema asset.

1. Open the Schema asset.
   a. On the main menu, click ⦂⦂⦂, then ▤ **Catalog**.
      » The Catalog Home opens.
   b. In the submenu, click **Data Dictionary** and select the **All Schemas** view.
   c. Click the schema that you want to refresh.

> Tip   You can also use the Collibra Data Intelligence Cloud search function to look up your schema.

2. In the view bar, to the right, click **Actions** → **Refresh**.
   » The **Refresh Schema** dialog box appears.

   > Tip   If Catalog experience is disabled, the **More** menu is shown instead of **Actions**.

3. In the **Login information** section, check **Store credentials** and enter the username and password you use to access your data source.
   » Your credentials are used to automatically connect to your data source and refresh the metadata in Collibra Data Intelligence Cloud.
4. Select **Schedule data refresh**.
5. Enter the required information.

| Option | Description |
| --- | --- |
| Cron pattern | Schedule of the data refresh as a Quartz Cron pattern. <br><br> > Warning   If you create an invalid Cron pattern, Collibra Data Intelligence Cloud stops responding. |
| Time zone | The time zone of the database. |

6. Click **Save**.

> **Example**
> When you refresh a schema of a registered data source, you can create a refresh schedule by selecting **Schedule data refresh**. You can then enter the CRON pattern *0 0 12?\*WED* to refresh every Wednesday at 12:00:00 PM.



# Registering a data source via Edge

Registering a data source via Edge makes metadata from the data source available in Collibra Data Intelligence Cloud.

> **Tip** You can also register a data source via Jobserver.

# Steps

The following table shows the steps required for registering a data source via Edge.

| Step | What? | Description | Results |
|------|-------|-------------|---------|
| 0 | Prerequisites | • Ensure the following settings are enabled:<br>  ◦ Database registration via Edge to allow registering a data source via Edge.<br>  ◦ Database profiling via Edge to allow profiling and classification via Edge.<br>  ◦ Catalog experience to have access to the **Configuration** tab page in the Database asset.<br>• You must have created and installed an Edge site with JDBC connections and capabilities.<br>If the JDBC connections and capabilities are not configured yet, a message on the Register content page shows "No data available" and you cannot continue registering a data source. | When you create a data source, the **Register content** page shows a list of available JDBC connections that you can use to register your database. |
| 1 | Register a data source | Registering a data source creates the structure for the metadata in Collibra. | • A **Physical Data Dictionary** domain containing a Database asset is created.<br>• A list of available schemas is created on the **Configuration** tab page of the Database asset. |

| Step | What? | Description | Results |
|------|-------|-------------|---------|
| 2 | Configure the synchronization of your data source | Making a selection of schemas and tables that you want to ingest.<br><br>When you select a schema to ingest, you can set the table rules to:<br><br>• Include and exclude tables of the schema.<br>• Specify the target domain in which to create assets.<br>• Exclude database views. | The information on the **Configuration** tab page of the Database asset is filled in. |
| 3 | • Synchronize one or more schemas manually<br>• Add a synchronization schedule to synchronize automatically | Synchronizing the schema of a registered data source to make the metadata available in Collibra. | Schema, Table, Column and Foreign Keys assets are created in the specified domain, and registration data becomes available. |
| 4 | If needed, profile and classify the synchronized data. | Data profiling creates a summary of a data source that is registered with Data Catalog and determines the data type of columns in the data source. The summary mainly contains statistics and graphics to give the user an idea what the registered data is about.<br><br>Classification analyzes and predicts the content of registered data sources based on a subset of the data itself, helping you to easily gain insights on what kinds of data you have and where it resides. | The Table and Column assets contain profiling information and the Columns are classified. |

# After registering a data source via Edge

When the registration is complete:

• A message at the top right tells you that the database registration is complete. A domain and a Database asset are immediately created.

- A workflow to assign a technical steward to the new domain is started. This is a simple out-of-the-box workflow that you can edit to fit your organization's needs. When you have assigned a technical steward, that technical steward has to set the security classification and indicate whether the data elements contain personally identifiable information (PII).
- If you registered a database without schemas, a new Schema asset is automatically created with the same name as the database or with a name as defined in the Edge capability.
- You can synchronize schemas in the database, including all tables, columns, views and foreign keys. Collibra creates assets in the selected target domains.
  - The synchronization jobs of all schemas run in parallel.
  - Collibra creates reports:
    - during the synchronization, to show the progress of the synchronization job.
    - after synchronizing, to show the synchronization logs for each synchronized schema.
  - The created assets receive a unique full name based on the following naming convention: [asset parent full name]>[asset name]

| Asset type | Naming convention | Example |
|---|---|---|
| Database | edgeConnectionName>jdbccatalog where jdbccatalog is the name retrieved from the JDBC "catalog" property. | Posgresql xs-gxsQ>-posgresqlsmall |
| Schema | edgeConnectionName>jdbccatalog>schemaName | Posgresql xs-gxsQ>-posgresqlsmall>public |
| Table | edgeConnectionName>jdbccatalog>schemaName>tableName | Posgresql xs-gxsQ>-posgresqlsmall>public>Condition |
| Database view | edgeConnectionName>jdbccatalog>schemaName>viewName | Posgresql xs-gxsQ>-posgresqlsmall>public>PriorConditions |

| Asset type | Naming convention | Example |
|---|---|---|
| Column | edgeConnectionName>jdbccatalog>schemaName>tableName>columnName (column) edgeConnectionName>jdbccatalog>schemaName>viewName>columnName (column) | Posgresql xs-gxsQ>posgresqlsmall>public>Condition>period.end (column) |
| Foreign key | edgeConnectionName>jdbccatalog>schemaName>foreignKeyName (foreign key) | Posgresql xs-gxsQ>posgresqlsmall>public>con.id (foreign key) |

You can view the full name of an asset by editing the asset.

# Register a data source via Edge

Before you can synchronize your data source, you have to register the data source to create an initial structure in a selected community in Data Catalog.

## Prerequisites

- You have enabled data source registration via Edge.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a global role with the View Edge connections and capabilities global permission.

- You have a resource role with the Configure external system resource permission, for example Owner.
- You have set up the required capability and JDBC connections in Edge.
- You made sure that the Catalog experience setting is enabled.
- You have created a System asset.

## Steps

1. On the main menu, click ⠿, then ⊟ **Catalog**.
   » The Catalog Home opens.
2. In the main menu, click the **Create** (＋) button.
3. In the **Register with Edge** section of the **Create** dialog box, click **Register a data source**.
   » The **Register content** page opens.

   > Note   You can register multiple databases via the same JDBC connection. If a database is already registered, the name of the Database asset is shown in the **Data sources** column.

4. In the row of the data source that you want to register, click **Add**.
   » The **Add Database** page opens.
5. Enter the required information.

| Field | Description |
|---|---|
| Community | The name of the parent community in which the assets will be created. |
| System | The name of the System asset.<br><br>After you registered the database, a relation of the type Technology Asset groups / is grouped by technology Asset is created between the System asset and the Database asset. |
| Database name | The name of the database in the data source to which you created a JDBC connection. The drop-down menu automatically shows which databases exist in your data source. You can refresh the list by clicking ↻ next to the drop-down menu.<br><br>The name will also be the name of the Database asset that is created as result of registering a data source via Edge. |

| Field | Description |
| --- | --- |
| Description | A description of the data source. |
| Owner | The owner of the data in the created community.<br><br>By default, the current user is selected. |

6. Click **Register**.
   » A Physical Data Dictionary domain containing a Database asset is created.
   » A relation of the type "Technology asset groups / is grouped by Technology asset
   " is created between the Database asset and the System asset.
   » The Database asset page opens.
   » A list of available schemas in the database is generated.

> Note   We highly recommend to not remove the relation between the Database asset and the System asset.

## What's next?

You can now configure the synchronization of your data source then synchronize the metadata.

# Register content and Add Database pages

When you register a data source via Edge, you first choose the JDBC connection, and then add a database. You perform these steps on different pages:

- a Register content page
- a Add Database page

## Register content page

The **Register content** page shows the data sources that are connected to Data Catalog via Edge. From this page, you can register the data sources to create assets that represent the content of these data sources.

| Column | Description |
|---|---|
| Connection name | The name of the connection to the data source. |
| Description | The description of the connection to the data source. |
| Type | The type of connection. For example, *JDBC* |
| Data sources | The name of the database that is already registered in Data Catalog. If the column is empty, the connection doesn't have a database yet.<br><br>Tip   If you click the name of the asset, you are automatically redirected to the Database asset page. If the name of the asset is grayed out, you don't have permission to go to the asset's page. |
| <Action button> | An **Add** button. Click this button to add a connection to a database and register the data source.<br><br>All added databases are listed in the data sources column. |

# Add Database page

In the **Add Database** page you can select a community, a System asset, the database from your data source that you want to register and the owner of the data.

## Add Database

**Edge data source**
redshift_1

Community *

Select a community from the list ▼

System *

Select a system from the list ▼

Database name *

Select a database name from the list ▼    ↻

Description

Owner *

John Fisher ▼

Cancel    Register

| Field | Description |
|---|---|
| Community | The name of the parent community in which the assets will be created. |
| System | The name of the System asset. After you registered the database, a relation of the type Technology Asset groups / is grouped by technology Asset is created between the System asset and the Database asset. |

| Field | Description |
|---|---|
| Database name | The name of the database in the data source to which you created a JDBC connection. The drop-down menu automatically shows which databases exist in your data source. You can refresh the list by clicking ↻ next to the drop-down menu.<br><br>The name will also be the name of the Database asset that is created as result of registering a data source via Edge. |
| Description | A description of the data source. |
| Owner | The owner of the data in the created community.<br><br>By default, the current user is selected. |

# Synchronizing a schema

Synchronizing schemas is the process of updating the metadata of a registered data source in Collibra Data Intelligence Cloud.

You can synchronize a schema manually or automatically at fixed intervals:

- Synchronize manually if you want to test the synchronization of your data source or if you want to synchronize immediately.
- Synchronize automatically if the content of the data source changes regularly.

In this section, you can find the relevant actions to successfully synchronize a schema in a registered database.

## About synchronizing schemas

Synchronizing schemas is the process of updating the metadata of a registered data source in Collibra Data Intelligence Cloud.

You can synchronize a schema manually or automatically at fixed intervals:

- Synchronize manually if you want to test the synchronization of your data source or if you want to synchronize immediately.
- Synchronize automatically if the content of the data source changes regularly.

# Synchronization process

- After you registered a data source via Edge, Data Catalog connects to your Edge site
  to create a list of schemas from the registered database.
  You can see the schema list on the **Configuration** tab page of the Database asset
  page.
- You can synchronize all schemas that have a table rule.
- During the synchronization process, the Edge site connects to your data source
  again and ingests all schemas, tables and columns according to the table rules. Col-
  libra Data Intelligence Cloud also detects whether there are changes since the last
  synchronization of a schema. Edge resolves the possible conflicts in the following
  way:

| Change in data source | Result in Collibra | Required action |
|---|---|---|
| A table, column or foreign key has been added to the schema. | Collibra creates the assets. | No action is required of you. |
| A table, column or foreign key has been removed from the schema. | The existing asset receives the **Missing from source** status. If it concerns a table, also the related Column assets receive the **Missing from source** status. | If needed, you can manually delete the assets. |
| A column or foreign key has been renamed. | ○ Collibra creates an asset with the new name. ○ The existing asset receives the **Missing from source** status. | If needed, you can apply any manual changes you made to the original asset, to the new asset. And then remove the assets that are no longer applicable. |

| Change in data source | Result in Collibra | Required action |
|---|---|---|
| A table has been renamed. | ○ Collibra creates a Table asset with the new name. Collibra also creates new Column assets for the new Table asset.<br>○ The existing Table and related Column assets receive the **Missing from source** status. | If needed, you can apply any manual changes you made to the original assets, to the new assets. And then remove the assets that are no longer applicable. |

Assets with the **Missing from source** status don't block the synchronization process.

> Note   In the asset diagram, assets with the **Missing from source** status are shown by default. If you don't want to see these assets, apply a filter to the diagram view to only display assets with valid statuses.

- If a schema is synchronized, you can see a check symbol (✓) beside the schema name. If the synchronization of a schema failed, an exclamation mark (❗) is shown instead.
  You can also see the synchronization status in the Activities list.

# Table rules

The table rule determines which tables of a schema you synchronize in Data Catalog. Only schemas that have a table rule can be synchronized. If a schema has a table rule, you can see a table icon (⊞) beside the schema name.



The following table shows fields of table rules:

| Table rule field | Description |
|---|---|
| Include | A comma-separated list of the names of the tables you want to synchronize.<br><br>• In the list, add a space after each comma. For example, `CUSTOMERS, ORDER, SKU`.<br>• You can use `*` as a wildcard.<br>• The default value is `*`, which means all tables are taken into account.<br>• The **Include** field takes priority over the **Exclude** field.<br><br>Example<br>• Out of all tables in a schema, you only want to synchronize the table with name "CUSTOMERS" and the tables with a name that starts with "ORDER".<br>To do this:<br>In the **Include** field, enter: `CUSTOMERS, ORDER*`.<br>• Out of all tables in a schema, you only want to synchronize the tables with a name that contains "SKU".<br>To do this:<br>In the **Include** field, enter: `*SKU*`. |

| Table rule field | Description |
|---|---|
| Exclude | A comma-separated list of the names of the tables you do not want to synchronize.<br><br>• In the list, add a space after each comma. For example, `CUSTOMERS, ORDER, SKU`.<br>• You can use `*` as a wildcard.<br>• The **Include** field takes priority over the **Exclude** field.<br><br>You can use the Exclude table rule to do the following:<br><br>• Synchronize all tables in a schema except the ones defined in the **Exclude** field.<br>• Synchronize only tables as defined in the **Include** field, with the exception of tables that are listed in the **Exclude** field.<br><br>> **Example**<br>> • Out of all tables in a schema, you do not want to synchronize a table with the name "ADDRESS" and tables with a name that ends with "PHONE".<br>> To do this:<br>> In the **Include** field, enter: `*` and in the **Exclude** field, enter: `ADDRESS, *PHONE`.<br>> • Out of all tables in a schema, you want to synchronize the tables with a name that starts with "SKU", but exclude the tables with a name that contains "bkp".<br>> To do this:<br>> In the **Include** field, enter: `SKU*` and in the **Exclude** field, enter: `*bkp*`.<br>> From the following list, only "SKU_1" and "SKU_2" will be synchronized.<br>> SKU_1 , SKU_2, SKU_bkp_1, SKU_bkp_2, New, bkp, bkp_SKU |
| Target domain | The Physical Data Dictionary domain in which the schema is synchronized.<br><br>The default value is **Schema domain**: the metadata is placed in a domain located in the same community as the domain of your Database asset. If that domain doesn't exist yet, Data Catalog creates it.<br><br>You can select any other Physical Data Dictionary domain for which you have a resource role with the Configure external system resource permission. |
| Options | Additional options to specify which type of tables you want to synchronize. |

| Table rule field | Description |
|---|---|
| Skip database views | A checkbox to exclude database views from the synchronization process. If selected, no assets of the type Database view are created.<br><br>Tip  You can also use the table rules to include or exclude specific database views. |

# Configure the synchronization of a data source

After you registered your data source via Edge, you configure the synchronization of your data source by means of table rules to determine which schemas and tables are ingested and how they are ingested. After this, you can synchronize them.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the View Edge connections and capabilities global permission.
- You have registered your database via Edge.

## Steps

1. Open a Database asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Metadata extraction** tab page, select a schema.

   Tip  You can search for a schema in the drop-down list or use the filter to show only schemas with or without table rules.

4.  If required, create or edit the table rule:
    a.  Do one of the following:
        - To create a new table rule, click **Add table rule**.
        - To edit an existing table rule, click **Edit** in the upper right corner.
    b.  Enter the required information.

| Table rule field | Description |
|---|---|
| Include | A comma-separated list of the names of the tables you want to synchronize.<br><br>- In the list, add a space after each comma. For example, `CUSTOMERS, ORDER, SKU`.<br>- You can use `*` as a wildcard.<br>- The default value is `*`, which means all tables are taken into account.<br>- The **Include** field takes priority over the **Exclude** field.<br><br>Example<br>- Out of all tables in a schema, you only want to synchronize the table with name "CUSTOMERS" and the tables with a name that starts with "ORDER".<br>To do this:<br>In the **Include** field, enter: `CUSTOMERS, ORDER*`.<br>- Out of all tables in a schema, you only want to synchronize the tables with a name that contains "SKU".<br>To do this:<br>In the **Include** field, enter: `*SKU*`. |

| Table rule field | Description |
| --- | --- |
| Exclude | A comma-separated list of the names of the tables you do not want to synchronize.<br><br>■ In the list, add a space after each comma. For example, `CUSTOMERS, ORDER, SKU`.<br>■ You can use `*` as a wildcard.<br>■ The **Include** field takes priority over the **Exclude** field.<br>You can use the Exclude table rule to do the following:<br><br>■ Synchronize all tables in a schema except the ones defined in the **Exclude** field.<br>■ Synchronize only tables as defined in the **Include** field, with the exception of tables that are listed in the **Exclude** field.<br><br>Example<br>■ Out of all tables in a schema, you do not want to synchronize a table with the name "ADDRESS" and tables with a name that ends with "PHONE".<br>To do this:<br>In the **Include** field, enter: `*` and in the **Exclude** field, enter: `ADDRESS, *PHONE`.<br>■ Out of all tables in a schema, you want to synchronize the tables with a name that starts with "SKU", but exclude the tables with a name that contains "bkp".<br>To do this:<br>In the **Include** field, enter: `SKU*` and in the **Exclude** field, enter: `*bkp*`.<br>From the following list, only "SKU_1" and "SKU_2" will be synchronized.<br>SKU_1 , SKU_2, SKU_bkp_1, SKU_bkp_2, New, bkp, bkp_SKU |

| Table rule field | Description |
|---|---|
| Target domain | The Physical Data Dictionary domain in which the schema is synchronized. |
| | The default value is **Schema domain**: the metadata is placed in a domain located in the same community as the domain of your Database asset. If that domain doesn't exist yet, Data Catalog creates it. |
| | You can select any other Physical Data Dictionary domain for which you have a resource role with the Configure external system resource permission. |
| Options | Additional options to specify which type of tables you want to synchronize. |
| Skip database views | A checkbox to exclude database views from the synchronization process. If selected, no assets of the type Database view are created. |
| | Tip   You can also use the table rules to include or exclude specific database views. |

  c. Click **Save**.
   » A table icon ( ⊞ ) appears next to the schema name in the schema list.
 5. If required, delete a table rule:
  a. Click **Edit** in the upper right corner.
  b. In the upper-right corner, click ⋮ → **Delete**.

> Note   You can only synchronize schemas that have table rules.

## What's next?

You can now synchronize the schemas to ingest the metadata into Collibra.

## Manually synchronize metadata

To keep the content of Collibra Data Intelligence Cloud synchronized with the schemas in the registered data source, you can synchronize manually or create a synchronization

schedule to automatically do this with a fixed interval. Synchronizing manually can be useful if you want to test your data source, or if you want to synchronize immediately.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the View Edge connections and capabilities global permission.
- You have the permissions to retrieve the metadata of the following database components through the JDBC Driver Database Metadata methods:
  - Schemas
  - Tables
  - Columns
  - Primary keys
  - Foreign keys
- You have registered your data source via Edge.
- You have configured the synchronization of your data source.

## Steps

1. Open a Database asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Metadata extraction** tab page, do one of the following:
   - To synchronize all schemas with a table rule, click **Synchronize metadata**.
   - To synchronize a specific schema in the database:
     i. Select the schema.
     ii. At the right of the schema name, click ⋮ → **Synchronize**.

   » The synchronization job starts.

   » The tables, columns, foreign keys and optionally views from the relevant schemas are ingested based on the configured table rules. See About synchronizing schemas.

## What's next

If required, you can profile and classify the synchronized data..

# Add or edit a synchronization schedule

To keep the content of Collibra Data Intelligence Cloud synchronized with the schemas in the registered data source, you can synchronize manually or add a schedule to automatically do this at fixed intervals.

> Note   You can only add one synchronization schedule for all schemas of a datasource. All schemas with table rules will be synchronized at the specified intervals.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the View Edge connections and capabilities global permission.
- You have the permissions to retrieve the metadata of the following database components through the JDBC Driver Database Metadata methods:
  - Schemas
  - Tables
  - Columns
  - Primary keys
  - Foreign keys
- You have registered your database.
- You have configured the synchronization of your data source.

## Steps

1. Open a Database asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Synchronization schedule** section of the **Metadata extraction** tab page, click **Add schedule** to add a new schedule, or ✏ to edit an existing schedule.
4. Enter the required information.

| Field | Description |
|---|---|
| Repeat | The interval when you want to synchronize the schemas automatically, for example daily, weekly or based on a Cron expression. |
| Cron | The Quartz Cron expression that determines when the synchronization takes place.<br><br>This field is only visible if you select `Cron expression` in the **Repeat** field. |
| Every | The day on which you want to synchronize the schemas, for example Sunday.<br><br>This field is only visible if you select `Weekly` in the **Repeat** field. |
| Every first | The day of the month on which you want to synchronize the schemas , for example Tuesday.<br><br>This field is only visible if you select `Monthly` in the **Repeat** field. |
| At | The time at which you want to synchronize the schemas automatically, for example 14:00.<br><br>This field is only visible if you select `Daily`, `Weekly` or `Monthly` in the **Repeat** field. |
| Time zone | The time zone for the schedule. |

5. Click **Save**.
   » All schemas with table rules are synchronized according to the synchronization schedule.

## Remove a synchronization schedule

You can remove a synchronization schedule from a Database asset to stop automatically synchronizing the schemas.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the View Edge connections and capabilities global permission.
- You have registered your database.
- You have configured the synchronization of your data source.

## Steps

1. Open a Database asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Synchronization schedule** section of the **Metadata extraction** tab page, click 🗑 .
   » The **Remove scheduling** dialog box appears.
4. Click **Remove**.

# Data source registration information

After you have registered a data source via Edge, the following information is added to the Table and Column assets.

| Column attribute | Description | From JDBC property |
|---|---|---|
| Column Name | The column name in the registered table. | COLUMN_NAME |
| Technical Data Type | The data type of the column in the data source. This value can differ from the **Data Type** value. | TYPE_NAME |
| Description from Source | The description of the column in the registered table. | REMARKS |

| Column attribute | Description | From JDBC property |
|---|---|---|
| Char octet Length | Maximum number of bytes in a character type's column. | CHAR_OCTET_LENGTH |
| Column Position | The index of the column in the source table. | ORDINAL_POSITION |
| Is Auto Incremented | Indication whether the data in the column is auto-incremented or not. | IS_AUTOINCREMENT |
| Is Generated | Indication whether the data in the column is generated or not. | IS_ GENERATEDCOLUMN |
| Is Nullable | Indication whether the column can store NULL values or not. | IS_NULLABLE |
| Is Primary Key | Indication whether the column is a primary key or not. | True if the primary keys resultSet contains the COLUMN_NAME |
| Number Of Fractional Digits | The number of decimal digits. | DECIMAL_DIGITS |
| Primary Key Name | The name of the primary key composed by the column. | PK_NAME |
| Size | The size of the column in the table. | COLUMN_SIZE |

| Table attribute | Description | From JDBC property |
|---|---|---|
| Table Name | The table name in the data source. | TABLE_NAME |
| Table Type | The table type in the data source, such as TABLE or VIEW. | TABLE_TYPE |
| Description from Source | The description of the table in the data source. | REMARKS |

For information on profiling, see Data profiling information.

# Sample data

Sample data is a set of randomly collected data from a data source. Sample data can be displayed for Table, Column, or Data Set assets. The purpose of showing sample data is to provide examples of the data so you know what to expect when you use the asset.



You can only view sample data for an asset:

- If the sample data feature is active.
- If you have the required permissions.
- If the asset is a Table, Column, or Data Set asset.
- If sample data is available for the asset.

Sample data is available in:

| Asset type | If Catalog experience is active, you can see the sample data in: | If Catalog experience is not active, you can see the sample data in: |
|---|---|---|
| Table | **Summary** tab pane<br>**Sample data** tab pane | **Details** tab pane<br>**Sample data** tab pane |
| Column | **Summary** tab pane<br>**Data profiling** tab pane | **Details** tab pane<br>**Sample data** tab pane |
| Data Set | **Summary** tab pane<br>**Sample data** tab pane | **Details** tab pane<br>**Sample data** tab pane |

> Tip
> In Table and Data Set assets, you only see sample data for columns for which you
> have the required permissions. If you do not have access, you see the text
> <sensitive> in the column instead of sample data.

The way Collibra handles sample data depends on how the assets are added in Collibra
and how the sample data is collected:

| | Assets are manually added or imported. | Assets are created by registering a data source via Jobserver. | Assets are created by registering a data source via Edge. |
|---|---|---|---|
| Sample data for an asset is uploaded via the Catalog REST API - Profiling. | The sample data is stored in the Collibra cloud repository.<br><br>The sample data is displayed to all users with the required permissions. | The sample data is stored in the Collibra cloud repository.<br><br>This sample data is also used for data classification via the Data Classification Platform.<br><br>The sample data is displayed to all users with the required permissions. | The sample data is stored in the Collibra cloud repository.<br><br>The sample data is displayed to all users with the required permissions. |
| Sample data is collected and stored when the data source is registered via Jobserver.<br><br>See Configure the use of sample data via Jobserver. | Not applicable. | The sample data is stored in the Collibra cloud repository.<br><br>This sample data is also used for data classification via the Data Classification Platform.<br><br>The sample data is displayed to all users with the required permissions. | Not applicable. |

|  | Assets are manually added or imported. | Assets are created by registering a data source via Jobserver. | Assets are created by registering a data source via Edge. |
| --- | --- | --- | --- |
| **Sample data can be manually requested for an asset** that is registered via Edge.<br><br>Note<br>This Edge feature is in Beta testing and must be activated and configured. | Not applicable. | Not applicable. | The requested sample data is cached on the Edge site for 24-48 hours.<br>No sample data is stored in the Collibra cloud repository.<br><br>The sample data is only displayed to users with the required permissions and if the sample data has been requested.<br><br>Note<br>Currently, you can only request sample data via Edge for Table and Column assets. |

For more details on the process, see Understanding the process to display sample data.

For details on the limitations and Beta feature guidelines, see Beta feature: limitations and guidelines.

# Required permissions to view sample data

To view sample data for an asset, you need:

- View permission on the asset. See View permissions.
  View permission is required to access the asset in general.
- Resource permission: **Asset > Data > View Samples**. See Resource permissions.
  View Samples is needed to see the sample data.

> **Tip**
> In Table and Data Set assets, you only see sample data for columns for which you have View permission. If you do not have access, you see the text <sensitive> in the column instead of sample data.

# Configuring the use of sample data

## Configure the use of sample data via Jobserver

You must configure your Collibra environment if you want to display sample data for data sources registered via Jobserver.

|   | Configuration step | More details |
|---|---|---|
| 1 | Ensure the users have the required permissions. | See Required permissions to view sample data. |

| | Configuration step | More details |
|---|---|---|
| 2 | Set the **Data Profiling** setting **Maximum number of samples** to a value higher than 0. | **Warning** For performance reasons, do not set this setting higher than 1,000. The limit of 1,000 will be enforced in a later release. |

**Note** Depending on your environment, you have to follow this procedure either in the Services Configuration section of the Collibra settings or in Collibra Console. This topic shows the procedure in Collibra Console. If your Settings page has a **Configuration** tab, you can follow this procedure in that tab page. For more information, see the online documentation.

## Prerequisites

- You have the ADMIN or SUPER role in Collibra Console.
- You have the SUPER role in Collibra Console.
- You have the ADMIN or SUPER role in Collibra Console.

## Steps

1. Open the DGC service settings for editing:
   a. Open Collibra Console.
      » Collibra Console opens with the **Infrastructure** page.
   b. In the tab pane, expand an environment to show its services.
   c. In the tab pane, click the Data Governance Center service of that environment.
   d. Click **Configuration**.
   e. Click **Edit configuration**.
2. Open the DGC service settings for editing:
   a. Open Collibra Console.
      » Collibra Console opens with the

| Configuration step | More details |
|---|---|
| | **Infrastructure** page.<br>b. In the tab pane, expand an environment to show its services.<br>c. In the tab pane, click the Data Governance Center service of that environment.<br>d. Click **Configuration**.<br>e. Click **Edit configuration**.<br>3. Go to the **Data profiling** section.<br>4. Make sure the setting **Maximum number of samples** is higher than 0.<br>The default value is 100. See also DGC service configuration: options.<br>5. Click the green **Save all** button. |
| 3 | When you register or refresh the data source via Jobserver, select the option **Store Sample Data**. | See Register a data source via Jobserver. |

For detailed information on the sample data process, see Understanding the process to display sample data.

# Configure the use of sample data via Edge

You must configure your Collibra environment if you want to display sample data for data sources registered via Edge.

> Important   This is a Beta feature. See Beta feature: limitations and guidelines.

> Warning   Sample data for data sources registered via Edge is temporarily cached on the Edge site. In the cache, the sample data is not encrypted. This means that the data is available in clear text in the Edge cache for 24-48 hours. Only the key that allows to identify the sample data's origin is encrypted.

| | Configuration step | More details |
|---|---|---|
| 1 | Ensure the users have the required permissions. | See Required permissions to view sample data.<br><br>Important   Several out-of-the-box roles already include those required permissions. Review the permissions assigned to those roles before enabling the public beta feature. |

| | Configuration step | More details |
|---|---|---|
| 2 | • Activate the public beta feature via **Data Profiling** setting **Sample data on Edge**.<br>• Set the **Data Profiling** setting **Maximum number of samples** to a value higher than 0. | **Warning** For performance reasons, do not set this setting higher than 1,000. The limit of 1,000 will be enforced in a later release.<br><br>**Important** The **Maximum number of samples** value applies to both Jobserver and Edge. In mixed environments, increasing the value can result in sample data extraction for data sources registered via Jobserver.<br><br>**Note** Depending on your environment, you have to follow this procedure either in the Services Configuration section of the Collibra settings or in Collibra Console. This topic shows the procedure in Collibra Console. If your Settings page has a **Configuration** tab, you can follow this procedure in that tab page. For more information, see the online documentation.<br><br>## Prerequisites<br><br>• You have the ADMIN or SUPER role in Collibra Console.<br>• You have the SUPER role in Collibra Console.<br>• You have the ADMIN or SUPER role in Collibra Console.<br><br>## Steps<br><br>1. Open the DGC service settings for editing:<br>   a. Open Collibra Console.<br>     » Collibra Console opens with the **Infrastructure** page.<br>   b. In the tab pane, expand an environment to show its services. |

| Configuration step | More details |
|---|---|
| | c. In the tab pane, click the Data Governance Center service of that environment.<br>d. Click **Configuration**.<br>e. Click **Edit configuration**.<br>2. Open the DGC service settings for editing:<br>   a. Open Collibra Console.<br>    » Collibra Console opens with the **Infrastructure** page.<br>   b. In the tab pane, expand an environment to show its services.<br>   c. In the tab pane, click the Data Governance Center service of that environment.<br>   d. Click **Configuration**.<br>   e. Click **Edit configuration**.<br>3. Go to the **Data profiling** section.<br>4. In **Sample data on Edge**, select *True*.<br>5. Make sure the setting **Maximum number of samples** is higher than 0.<br>The default value is 100. See also DGC service configuration: options.<br>6. Click the green **Save all** button. |
| 3   For each data source, add the following Edge capability: **Catalog JDBC Sampling**. | For information on how to add capabilities, see Add an Edge capability to an Edge site.<br><br>• The **Catalog JDBC Sampling** capability allows to collect and cache sample data for a given JDBC data source in the Edge site, and can retrieve sample data from the Edge cache.<br>• Once the capability is selected, define the JDBC connection to which the capability applies. |

For detailed information on the sample data process, see Understanding the process to display sample data.

# Beta feature limitations and guidelines

# Sample data limitations

- Currently, you can only request sample data via Edge for Table and Column assets.

- Sample data in the Edge cache is not encrypted. This means that the data is stored in clear text on the Edge site 24-48 hours. Only the key that allows to identify the sample data origin is encrypted.

- For performance reasons, the number of samples to display should be less than 1,000. This limit is configurable in the **Maximum number of samples** setting, in the Profiling section. The default value is 100.
  See Configure the use of sample data via Edge and Configure the use of sample data via Jobserver.
  The limit of 1,000 will be enforced in a later release.

- For performance reasons, avoid sampling tables with more than 2,000 Column assets.
  This limit is not configurable and will be enforced in a later release.

- The sampling feature always uses push-down sampling if push-down sampling is available for the data source. Push-down sampling increases the sample data extraction speed.
  We advise to only allow sampling on data sources that support push-down sampling. To know if your data source allows for push-down sampling (called partial scan in Edge), see Data sources supported by Edge and Overview of Collibra-provided JDBC drivers (Jobserver).

# Sample data guidelines

During the Beta testing phase, we advise to:

- Keep the maximum number of samples to the default limit of 100 samples or do not exceed 1,000 samples.
- Don't request sample data for tables with more than 2,000 columns.
- Use sampling only on data sources that support push-down sampling.

## Your feedback is welcome

- During the Beta testing phase, we do advise to use sampling only on data sources that support push-down sampling.
  However, we are looking for feedback on data sources that do not allow push-down sampling.

  > Note   If you try sampling on a data source that does not allow push-down sampling, note that the sample data extraction time is proportional to the database table size. The bigger the table, the longer it will take to retrieve the samples.

- We are looking for feedback on large parallel sample data requests. This happens when a lot of users want to see sample data at the same time.

  > Tip   If you experience issues in this situation, you can decrease the number of Edge data sources for which the sampling capability is enabled.

# Delete sample data

The way to remove sample data for a data source depends on how the sample data is made available.

- For Jobserver, perform one of the following:
  - Refresh the related schema and don't select the **Store Sample Data** checkbox.
    As a result, any previously gathered sample data is removed from the Collibra cloud repository.
  - Call the Catalog profiling REST API with an empty array for the samples parameter.
    As a result, any previously gathered sample data is removed from the Collibra cloud repository.
- For Catalog profiling REST API, call the Catalog profiling REST API with an empty array for the samples parameter.
  As a result, any previously gathered sample data is removed from the Collibra cloud repository.

- For Edge, you cannot delete sample data.
  Sample data for data sources registered via Edge is not stored in the Collibra cloud

repository, it is cached on the Edge site. Every day, Edge deletes all sample data that is older than 24 hours from its cache.

> Note If a data source was previously connected to Jobserver or if sample data was pushed using the Catalog profiling REST API, and the data source is now an Edge data source, sample data may still be stored in the Collibra cloud repository for this data source. If you want to remove this sample data, call the Catalog profiling REST API with an empty array for the samples parameter.

# Example of API code that deletes sample data from the Collibra cloud repository

In the example code:

- replace <your_environment> by the name of your environment.
- replace the assetIdentifier section by any combination that uniquely identifies the asset for which you want to delete the sample data.

Example
```
PATCH https://<your_
environment>.collibra.com/rest/catalog/1.0/profiling/columns
{
   "columnProfiles": [
      {
         "assetIdentifier": {
         "assetName": "Catalog postgresql>catalog_
         postg>GDPR>Consumers>Process_id(column)",
         "communityName": "Catalog demo",
         "domainName": "Catalog postgresql > catalog_postg >
         GDPR"
         },
         "samples": {
         "samples": null
         }
      }
   ]
}
```

# Request sample data

> **Note**
> - This Edge feature is in Beta testing and must be activated and configured.
>
> - Currently, you can only request sample data for Table and Column assets.

If a data source is registered via Edge, no sample data is available by default. It means you may need to request the sample data. The requested sample data only remains available via Edge for a limited amount of time.

## Steps

1. Open the Column or Table asset for which the data source is registered via Edge.
2. For a column, click **Data Profiling**.
   For a table, click **Sample data**.
   » Collibra performs checks to determine what happens in the page.
   » If sample data is available in the Edge cache for this data source, the sample data is displayed.

   > **Note** It can take some time for the sample data to be displayed.

   » If no sample data is available in the Edge cache and if you have the required permissions to see sample data, the **Request Sample Data** button is displayed.

3. If displayed, click the **Request Sample Data** button or go to **Actions -> Request Sample Data**.
   » A 'Request sample data' job starts. This job collects and caches the sample data on the Edge site.
   » Collecting and caching the sample data can take some time. You can navigate to other assets while the request is pending.
   » As soon as the request is completed, you receive a notification. When you open the asset again, the sample data will be displayed after some time.

> Note   Requested sample data via Edge only remains available on the Edge
> site for a limited time. Once a day, Edge deletes all the data that is older than
> 24 hours. Once the cached data is deleted, you will need to request the
> sample data again.

For detailed information, see Understanding the process to display sample data.

# Understanding the process to display sample data

If you open a Column, Table or Data Set asset page, Collibra performs a series of checks
to determine if sample data is displayed.

> **Note**
> To use sample data via Edge, you must activate the feature and configure Edge.
>
> - This Edge feature is in Beta testing.
> - Currently, you can only request sample data via Edge for Table and Column assets.

| Check or Action | Description | Positive outcome | Negative outcome |
| --- | --- | --- | --- |
| 1 | Collibra checks if you have the required permissions to view sample data. | You have the required permissions:<br><br>- The process continues with the next check. | You don't have the required permissions:<br><br>- You cannot see the sample data and a message appears on the page.<br>- The process stops. |
| 2 | Collibra checks if sample data is stored in the Collibra cloud repository.<br><br>> **Tip**<br>> This is only possible if:<br>><br>> - Sample data was extracted during the registration of the data source via Jobserver.<br>> - Sample data was uploaded by using the Catalog REST API - Profiling. | Sample data is available in the Collibra cloud repository:<br><br>- The sample data is visible in the page.<br>- The process stops. | No sample data is available in the Collibra cloud repository:<br><br>- The process continues with the next check. |

| Check or Action | Description | Positive outcome | Negative outcome |
|---|---|---|---|
| 3 | Collibra checks if the asset is connected to an Edge site.<br><br>Tip  An asset can be connected to an Edge site via its related Database asset. | The asset is connected to an Edge site:<br><br>• The process continues with the next check. | No sample data is available for the asset:<br><br>• You cannot see the sample data.<br>• The process stops. |
| 4 | Collibra checks if sample data is available in the cache of Edge.<br>This is possible if sample data has been requested before and the cache has not been cleared in the meantime.<br>Once a day, Edge deletes all data that is older than 24 hours from its cache. | Sample data is available in the cache:<br><br>• The sample data is visible in the page.<br><br>Note  It can take some time for the sample data to be displayed.<br><br>• The process stops. | No sample data is available in the cache:<br><br>• The process continues with the next check. |
| 5 | Collibra checks if a sample data request is pending for the data source. | A sample data request is pending for the data source:<br><br>• You need to wait until the sample data has been collected and cached in the Edge site.<br>• The process stops. | No sample data request is pending for the data source:<br><br>• The button **Request Sample Data** appears on the page and in the **Action** drop-down list.<br>• The process stops until you click the button. |

| Check or Action | Description | Positive outcome | Negative outcome |
|---|---|---|---|
| 6 | You click the **Request Sample Data** button. | • A 'Request sample data' job is launched and added to the Activities list.<br>• When Edge is available, the job starts. Sample data is collected and temporarily made available in the cache of the Edge site. The sample data remains available between 24 and 48 hours. Every day, Edge deletes all data that is older than 24 hours from its cache.<br>• While the job is pending, you can navigate away from the page.<br>• Once the job is completed, you receive a notification and the process stops. | |

# Troubleshooting: sample data

# You receive an error code

| Code | Description | Possible causes | Solution |
|------|-------------|-----------------|----------|
| 400 | This message appears if:<br><br>• Something is wrong with the provided asset ID<br>or<br>• The sampling capability is not installed on the Edge site.<br><br>The error message will specify the problem. | • The asset exists but the asset is not a column or table<br>• The table has no columns.<br>• Something is wrong in the relationship of the column, table or database, like a column asset that was not ingested but manually created and no relationship has been defined.<br>• The **Catalog JDBC Sampling** capability has not been defined for the data source Edge connection. | • If it concerns a wrong asset, provide a valid column or table asset id.<br><br>• If the sampling capability is missing, install the **Catalog JDBC Sampling**capability for the data source. |
| 401 | This message appears if you are not authenticated to use the sampling API. | The authentication failed. | Provide valid credentials. |
| 403 | This message appears if you lack permission to any of the columns within the requested asset. | You do not have the required permissions. Both View permission and View Samples permission are needed to see sample data for an asset. | Verify the user has the required permissions. |
| 404 | This message appears if the asset cannot be found. | The asset does not exist. | Provide an existing column or table asset id. |

| Code | Description | Possible causes | Solution |
|------|-------------|-----------------|----------|
| 503 | This message appears if the Edge service gets a timeout or fails. | The Edge service is not available. | Verify that the Edge site is still online and healthy. If not, check the Edge logs to get a better understanding of the issue. If the problem persists, contact Collibra Support for assistance. |

# No sample data is displayed

There are many conditions that can result in no sample data being displayed. Before reporting an issue, check the following:

| Cause | Description | Solution |
|-------|-------------|----------|
| The setting **Maximum number of samples** is set to 0. | The sampling feature is disabled and no samples are displayed. | Set the **Data Profiling** setting **Maximum number of samples** to a value higher than 0. See Configuring the use of sample data |
| The sampling capability is missing for your Edge data source. | Samples can only be extracted if the sampling capability is set for the data source on the corresponding Edge site. | Install the **Catalog JDBC Sampling**capability for the data source. |
| The asset for which you want to collect sample data has no data. | There is no data to show for the asset. | |

| Cause | Description | Solution |
|---|---|---|
| No sample data is stored in the Collibra cloud repository. (not applicable for data sources registered via Edge) | • For Jobserver data sources, sample data is only available in the Collibra cloud repository if the **Store Sample Data** option was selected during the registration of the data source.<br>• For assets created without Jobserver or Edge registration, sample data is only available if they were uploaded to the Collibra cloud repository via the Catalog Profiling REST API. | Configuring the use of sample data |

# You always see old sample data for a data source registered via Edge

Sample data stored in the Collibra cloud repository takes precedence over sample data extraction by Edge. Sample data can be available for an Edge data source in the Collibra cloud repository if this data source was previously connected to Jobserver or if sample data was pushed using the Catalog profiling REST API for the data source.
If you want to remove samples from the Collibra cloud repository, see Delete sample data.
See also Understanding the process to display sample data.

# Collecting the sample data is very slow

- It can take some time to read and display the sample data available in the Edge cache.
- The sample data extraction time via Edge is influenced by multiple factors. For example: table size, number of columns in a table, number of samples to collect, maximum length of samples, and push-down sampling mechanism available for the data source. See also Beta feature limitations and guidelines

# Retrieving sample data log files

For data sources registered via Edge, Edge logs are generated when sample data is extracted from the data source and cached on the Edge site. The logs start with this text: "Writing cache samples with the key...".
Looking at the Edge logs within a 2-day period should give information on the sampling activity.

> Example
> Writing cache samples with the key 'catalog.sample.6385e23cb1ae443a7786c555108d8bb028d23dee39e76ce3169ea a9cdacb1ed3'
> "Cache write sample for table 'Snowflake>SNOWFLAKE_SAMPLE_ DATA>TPCDS_SF100TCL>CALL_CENTER'

# Quality extraction

The quality extraction functionality allows you to ingest Data Quality user-defined rules, metrics, and dimensions into Collibra Data Catalog for registered data sources by using the DQ Connector Edge capability.

> **Warning**   You can use this functionality only if you have the same data sources registered in Collibra Data Quality and you have Edge enabled in your Collibra Data Intelligence Cloud.

## About DQ Connector

The native DQ Connector brings Data Quality into your Collibra Data Intelligence Cloud. The DQ Connector is an Edge capability template that helps you integrate your Data Quality user-defined rules, metrics, and dimensions into Collibra Data Catalog.

> **Note**   To extract data quality statistics from CollibraData Quality both Data Catalog and Data Quality must ingest the same data source.

## DQ Connector requirements

- Collibra Data Intelligence Cloud 2021.07 or newer.
  - Edge
- Collibra Data Quality 2.15 or newer.
  - Existent data quality statistics for the selected data source.

# DQ Connector configuration

1. Connect to a Data Quality source:
    a. Create a Data Quality Edge site.
    b. Connect to your Data Quality source.
    c. Add ingestion capabilities to your Data Quality connection.
    d. Configure destinations for Data Quality assets.
    e. Add Data Quality characteristics to assets.
    f. Add a DQ Connector capability.
2. Register Data Quality Edge connections in Data Catalog:
    a. Create a Data Catalog System Asset.
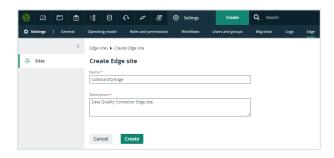    b. Register the Data Quality data source in Data Catalog.

# Connect to a Data Quality source

Because the DQ Connector is an Edge capability, you must be able to ingest data via Edge. For information about enabling and configuring Edge, see the Edge Configuration guide.

## Create a Data Quality Edge site

Create an Edge site with the following properties:

| Field | Description |
| --- | --- |
| Name | The name of the Edge site, for example **Collibra-DQ-Edge**. Do not use spaces or special characters.<br><br>This field is mandatory and the name must be globally unique. |
| Description | The description of the Edge site. We recommend to put at least basic location information of the Edge site.<br><br>This field is mandatory. |

# Install the Data Quality Edge site

Follow the instructions for your environment to Install an Edge site.

> Note   This process automatically creates an Edge user, which you use later in the setup process.

# Connect to your Data Quality source

Create a connection for each CollibraData Quality data source you want to synchronize.

| Section | Property | Value |
| --- | --- | --- |
| Connection set-tings | Name | The same name as the Collibra Data Quality con-nection name. |
| | Description | The description of the JDBC connection. This field is also visible when you register content. |
| | Connection pro-vider | The connection provider, which determines the available connection parameters. Same as Collibra Data Quality. |

| Section | Property | Value |
|---|---|---|
| Connection parameters<br><br>Example for **Username / Password JDBC drive** | Username | The same username as the Collibra Data Quality connection username. |
| | Password | The same password as the Collibra Data Quality connection password. |
| | Driver class name | The same driver name as the Collibra Data Quality connection driver name. |
| | Driver Jar | The same driver JAR file as from Collibra Data Quality. |
| | Connection string | The same URL as the Collibra Data Quality connection URL. |

## Add ingestion capabilities to your Data Quality connection

You must add a Catalog JDBC ingestion Edge capability template for each connection you have created to extract and process data for your data source.

| Field | Description | Requir-ed |
|---|---|---|
| **Capability** | This section contains the general information about the capability. | |
| Name | The name of the Edge capability. | ✓ Yes |
| Description | The description of the Edge capability. | ✗ No |

| Field | Description | Requir-ed |
|---|---|---|
| Capability template | The capability template, which determines the next available sections.<br><br>Select the Catalog JDBC ingestion capability template to register a data source. | ✓ Yes |
| **Connection** | This section contains information to connect to the data source. | |
| JDBC connection | The connection to the data source. | ✓ Yes |
| JDBC data source type | The data source type of the data source that you want to ingest. | ✓ Yes |
| Supports schemas | A text field where you have to enter *True* to enable database registration of data sources that have no schema. If the data source has schemas, you can ignore this field.<br><br>Tip   If the data source does not have a schema, Data Catalog creates a Schema asset with the same name as the full name of the database. | ✕ No |

| Field | Description | Requir-ed |
|---|---|---|
| Others | This section can contain additional capability properties.<br><br>> **Warning**   Adding additional properties can have a significant impact on your Edge site. Only add or update them together with Collibra Support.<br><br>Click **Add property** to add a property.<br>The possible properties are: Show properties | ✕ No |

| Name | Description | Typ-e | Encryp-tion | Value |
|---|---|---|---|---|
| ingest-schem-a-garbag-e-col-lector | The garbage col-lector that is used by the cap-ability.<br>The default garbage col-lector is -XX:+UsePar-allelGC.<br>For information about other pos-sibilities, see AZUL doc-umentation. | Text | Not encryp-ted (plain text) | For example: -XX:+UseSeri-alGC |
| ingest-schem-a-request-s-cpu | The minimum amount of CPU computing power reques-ted by the cap-ability.<br>The amount is expressed in mil-liCPU. | Text | Not encryp-ted (plain text) | For example: 100 |
| ingest-schem-a-limits- | The maximum amount of CPU computing | Text | Not encryp-ted (plain | For example: 950 |

| Field | Description | Required |
|-------|-------------|----------|
| **General** | This section contains general information about logging. |  |
|  | Note This section only applies to JDBC capabilities. |  |
| Debug | An option to automatically send Edge infrastructure log files to Collibra Data Intelligence Cloud. By default, this option is set to *false*. | ✕ No |
|  | Note We highly recommend to only send Edge infrastructure log files to Collibra Data Intelligence Cloud when you have issues with Edge. If you set it to *true*, it will automatically revert to *false* after 24h. |  |
| Log level | An option to determine the verbosity level of Catalog connector log files. By default, this option is set to *No logging*. | ✕ No |

# Configure destinations for Data Quality assets

Data Quality rules, metrics and dimensions require their own domains in Data Catalog. If you don't have existing domains for data quality or wish to use new ones for the quality extraction purpose, create a domain for each type of data quality asset:

- Rules: **Rulebook Domain**
- Metrics: **Business Asset Domain**
- Dimensions: **Business Asset Domain**

# Assign permissions for Data Quality domains

Edge must have the correct resource permissions to manage assets inside the dedicated Data Quality domains. For each dedicated domain, assign the **Technical Steward** role to the Edge user.

> Note   The Edge user is automatically created when you install the Edge site.

## Add Data Quality characteristics to assets

To show Data Quality statistics for your data source, assign the following characteristic types to the **Table** and **Column** asset types:

| Asset type | Characteristic type |
|---|---|
| Table | governed by Governance Asset |
| Column | is governed by Data Quality Rule |

## Add a DQ Connector capability

The DQ Connector facilitates the communication with Collibra Data Quality. Add a DQ Connector capability to your Data Quality Edge site:

| Field | Description | Required |
|---|---|---|
| **Capability** | This section contains the general information about the capability. | |
| Name | The name of the Edge capability. | ✓ Yes |
| Description | The description of the Edge capability. | ✕ No |

| Field | Description | Required |
|---|---|---|
| Capability template | The capability template, which determines the next available sections.<br><br>Select the DQ Connector capability template to ingest Data Quality user-defined rules, metrics, and dimensions into Collibra Data Catalog.<br><br>> Important   Data Quality is only available in beta. Please create a support ticket to get access. | ✓ Yes |
| **DQ** | This section contains information about the Collibra Data Quality connection. | |
| Base URL | Your Collibra Data Quality URL | ✓ Yes |
| Username | The Collibra Data Quality username for this connection. | ✓ Yes |
| Password | The Collibra Data Quality password for this connection. | ✓ Yes |
| Encryption options | Select the type of encryption to use.<br><br>Default: *To be encrypted by Edge management server*. | |
| Issuer of the JWT | If you have selected *Encrypted with public key*, enter your JWT issuer. | ✕ No |
| Collibra metadata model | This section contains information about where to ingest Data Quality assets. | |
| DQ Rules domain id | The UUID of the **Rulebook Domain** for the ingested Data Quality rules. | ✓ Yes |

| Field | Description | Required |
|---|---|---|
| DQ Metrics domain id | The UUID of the **Business Asset Domain** for the ingested Data Quality metrics. | ✓ Yes |
| DQ Dimensions domain id | The UUID of the **Business Asset Domain** for the ingested Data Quality dimensions. | ✓ Yes |
| Default DQ Dimension name | The default **Data Quality Dimension**, for example *Accuracy*, *Completeness*, *Consistency* and so on.<br><br>Default: *Completeness*. | ✓ Yes |
| DQ Metric classified by DQ Dimension relation type id | The UUID of the **Data Quality Metric classified by / classifies Data Quality Dimension** relation. If left unspecified, this relation will not be added. | ✕ No |
| Assets are imported in batches of this size | The batch size of the ingestion.<br><br>Default: *5000*. | ✓ Yes |
| **General** | This section contains general information about logging.<br><br>> Note   This section only applies to JDBC capabilities. | |
| Debug | An option to automatically send Edge infrastructure log files to Collibra Data Intelligence Cloud. By default, this option is set to *false*.<br><br>> Note   We highly recommend to only send Edge infrastructure log files to Collibra Data Intelligence Cloud when you have issues with Edge. If you set it to *true*, it will automatically revert to *false* after 24h. | ✕ No |

| Field | Description | Required |
|---|---|---|
| Log level | An option to determine the verbosity level of Catalog connector log files. By default, this option is set to *No logging*. | ✕ No |

## Next steps

- Register a Data Quality source in Data Catalog.

# Register a Data Quality source in Data Catalog

To make the Data Quality metadata available in Collibra Data Catalog, you must register the data source for each Collibra Data Quality data source you want to synchronize.

## Create a Data Catalog System Asset

As a prerequisite to registering a data source in Data Catalog, you must create a System asset for each connected data source with the following properties:

| Field | Value |
|---|---|
| Type | **System** |
| Domain | The domain to which the new assets will belong. You can only create a asset type in any domain of a domain type that is assigned to a selected asset type. |
| Name | The same name as the CollibraData Quality connection name. |

## Register the Data Quality data source in Data Catalog

Register each Data Quality source in Data Catalog.

# Next steps

- Extract Data Quality metadata.

# Extract Data Quality metadata

After you completed the DQ Connector configuration, you can start ingesting Data Quality metadata.

# Prerequisites

- You have configured the metadata extraction properties for the data source.

# Steps

1. Open a Database asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Quality extraction** section, do one of the following:
   - To select schemas for data quality synchronization:
     i. Click **Edit**.
        » The **Data quality** column becomes editable.
     ii. Select whether to synchronize the available schemas.

     iii. Click **Save**.

   - To synchronize the selected schemas:
     i. Select the schema name to see its configuration.
     ii. Click **Synchronize**.
        » The synchronization job is started for the selected schemas.

Chapter 4

# Data profiling

# About data profiling

Data profiling creates a summary of a data source that is registered with Data Catalog and determines the data type of columns in the data source. The summary mainly contains statistics and graphics to give the user an idea what the registered data is about.

You can create profiling data by:

- Registering a data source via Jobserver or via Edge, and choosing to create profiling data.
- Importing profiling data via the Catalog API.

You can find the profiling information in Table and Column asset pages.
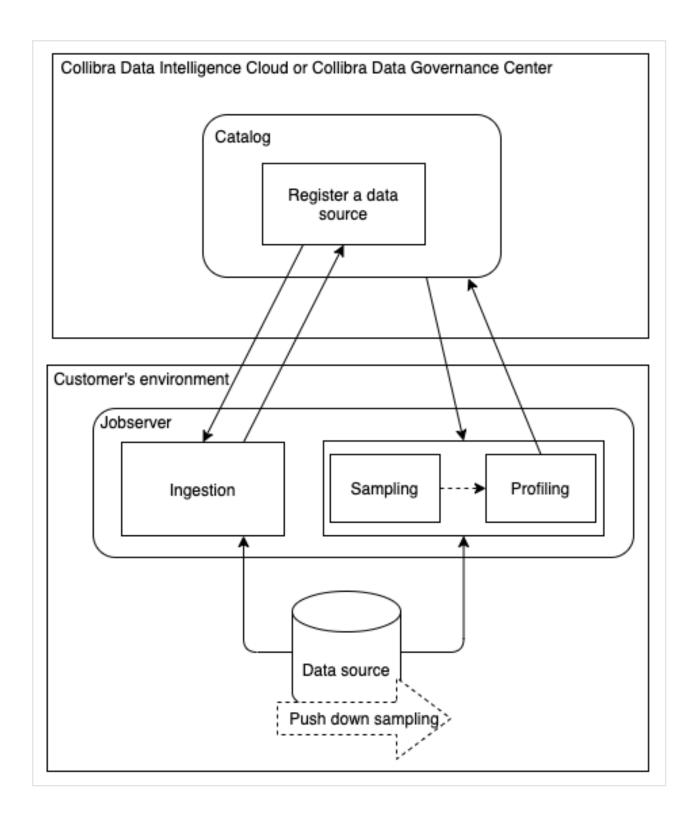
# Profiling process

You can profile data via Edge or via Jobserver.

## Profiling via Jobserver

When you register a data source via Jobserver, Data Catalog triggers the ingestion process. By default, the complete data set is transferred to the Jobserver, which then creates sample data based on your data source. Jobserver then profiles the sample data and sends the result to Data Catalog.
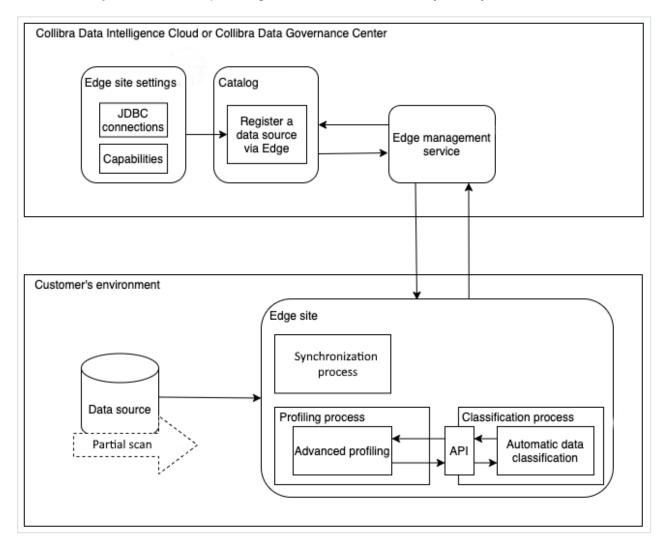
You can enable the **Anonymize data** option to hash or remove profiling information that can be considered sensitive.

# Profiling via Edge

When you registered a data source via Edge and you have created a profiling capability, you can profile and classify the metadata via the Database asset page of the registered data source.

- You install an Edge site close to your data source and create the necessary JDBC connections and capabilities to ingest the data source and to profile and classify its metadata.

- After synchronizing the schemas of the registered database, you can profile and classify the data. The profiling results are automatically anonymized.

## Differences between profiling via Jobserver or via Edge

The following table shows the differences between profiling via Jobserver or via Edge.

| Part of process | Profiling via Jobserver | Profiling via Edge |
|---|---|---|
| Data size | There is a limit on the size of the data that is used to calculate profiling statistics. By default, this is 10 GB. | There is no data size limit. The Edge site calculates the profiling statistics while reading the data. |
| Connectivity | Jobserver requires an HTTP proxy to support reverse connectivity. | Collibra connects to an Edge site. The Edge site is installed in the customer's environment, close to the data source. The Edge site communicates to Collibra Data Intelligence Cloud and other 3rd party systems using an HTTPS connection. |
| Register a data source | When registering a data source via Jobserver, options are available to profile the data and create sample data. | You can only profile the data after you registered a data source and synchronized one or more schemas. You can start the profiling process via the **Configuration** tab page on a Database asset page. |
| Deleting data profiling information | To delete data profiling information for a schema, refresh the schema without storing the data profile. See Refresh the schema of a registered data source. | Once data profiling information is available, you can only delete it by deleting the assets. |

## Profiling sample

To create a data profile, Data Catalog uses a representative sample of the data.

> Note   This profiling sample data is not the same as the sample data available in **Sample data**.

## Creating a profiling sample via Jobserver

If you register a data source via Jobserver, the profiling sample is created when you register the data source.

- If you use Jobserver without push down sampling, the complete data set is transferred to the Jobserver, which then creates the profiling sample based on your data source. The sample size is determined by the **Table profiling data size** setting in Collibra Console or the Services Configuration section of the Collibra settings. By default, the size is 10 GB.
- If you use Jobserver with push down sampling (also called partial scan), the data source itself creates the profiling sample and sends it to Data Catalog.
  The data source creates the sample from randomly selected data and transfers it to the Jobserver. If the cache storage is reached, the process stops. Because the data source already created the sample randomly, the omitted data can be ignored without lowering the representativeness of the sample.

  > Warning   Push down sampling is only available for some data sources. To verify if your data source allows push down sampling, see Collibra-provided JDBC drivers.

## Creating a profiling sample via Edge

Edge profiles and classifies the data on the Edge site itself and only sends the profiling results and classification suggestions to Collibra Data Intelligence Cloud.

- If you use full scan via Edge, all the rows in a table are scanned for profiling, without limit.
- If you use partial scan, the data source itself creates the profiling sample from randomly selected data and sends it to Data Catalog.

> Warning   Partial scan is only available for some data sources. To verify if your data source allows partial scan, see Collibra-provided JDBC drivers.

For more information, see Configure the profiling and classification options via Edge.

# Using push down sampling or partial scan

Push down sampling means that the task of creating the data sample is delegated to the data source itself. In Edge, push down sampling is called partial scan.

- The data source creates the sample from randomly selected data and transfers it to the Jobserver or the Edge site in one fetching process.
If the cache storage is reached nonetheless, the fetching process can be stopped. Because the data source already created the sample randomly, the omitted data can be ignored without lowering the representativeness of the sample.
- Push down sampling can be done using dynamic SQL query, if the data source supports data sampling. For an overview, see Overview of Collibra-provided JDBC drivers.

Push down sampling drastically increases the performance of sampling.

## Enable push down sampling

Push down sampling is not used by default. To use push down sampling, do the following:

| Step | When | Description |
| --- | --- | --- |
| 1 | Manage the driver | Add the **pushDownSampling** connection property. |

| Step | When | Description |
|---|---|---|
| 2 | Register your data source | Follow the usual steps to register a data source, but include the following options:<br><br>1. Enter a value for the **pushDownSampling** connection property.<br><br>Note<br>  ◦ The value must be between *100* and *1 000 000*. Your data source creates the sample of that amount of rows.<br>  ◦ If the size of the amount of rows exceeds the limit of the cache storage (Collibra recommends 10 to 20 GB), the amount of rows is reduced.<br>  ◦ If you typed a value that is bigger than the amount of rows in the data source, the entire data source is used as a sample.<br><br>2. Select **Store Data Profile** and, optionally, **Store Sample Data** to profile via Jobserver. |

Push down sampling is performed via the **Partial scan** option. To use partial scan, do the following:

| Step | When | Description |
|---|---|---|
| 1 | Register and synchronize a data source via Edge. | Follow the usual steps to register and synchronize a data source via Edge. |

| Step | When | Description |
|------|------|-------------|
| 2 | Profile and classify synchronized metadata. | Click the **Profiling and classification** tab in a Database asset's **Configuration** tab page, and do one of the following:<br><br>• To use partial scan for all schemas:<br>  a. In the **Default profiling and classification rule** section, click **Edit**.<br>  b. Select **Partial scan**.<br>  c. Enter the maximum number of rows that you want to use for profiling.<br>  d. Click **Save**.<br>  e. Profile and classify.<br>• To use partial scan for a specific schema only:<br>  a. In the **Schema profiling and classification rules** section, select the schema.<br>  b. Do one of the following:<br>    ▪ To create a new table rule, click **Add table rule**.<br>    ▪ To edit an existing table rule, click **Edit** .<br>  c. Select **Partial scan**.<br>  d. Enter the maximum number of rows that you want to use for profiling.<br>  e. Click **Save**.<br>  f. Profile and classify. |

# Profiling via Jobserver

# Data anonymization via Jobserver

> Tip   If you profile and classify via Edge, data in columns with data type Text or Geo is automatically anonymized before it is sent to Collibra Data Intelligence Cloud.
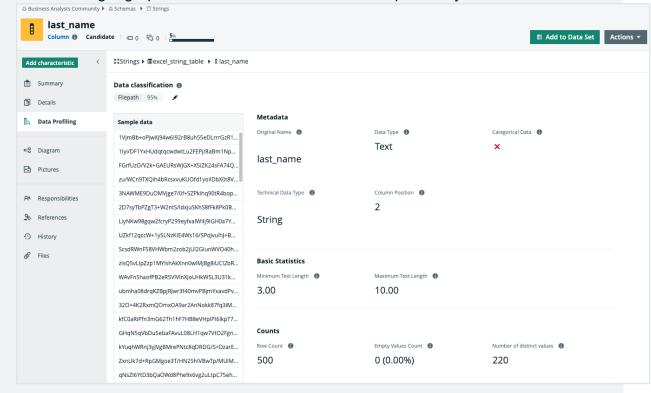
To ensure that sensitive data is not stored in the cloud, you can enable the Anonymize data option in Collibra Console.

With this option enabled, Collibra anonymizes the content of columns with data of the type Text and Geo immediately at the end of the profiling process. As a result, data samples and the values that are shown in the data distribution charts are replaced by a random hash value for columns that contain these data types. Attributes that could contain sensitive data, like attributes of the type Mode or Percentiles, are no longer calculated for columns with data type Text or Geo.

Identical values in a column get the same hash value so that you can still recognize the values as identical.

Collibra detects the data type of a column during profiling and only anonymizes the data if the data type attribute is Text or Geo. However, if Collibra detects a data type that does not correctly correspond with the actual data type, some data may not have been anonymized or has been wrongfully anonymized. To solve this, you can manually modify the column's data type and profile again.

Example   You enabled the Anonymize data option in Collibra Console and profiled a column that has data type Text. If you go to the **Summary** or **Data Profiling** tab, all textual and geographical data has been removed or replaced by hashed values:



Note   Jobserver does not automatically anonymize your data. To ensure that your sensitive data is not stored in the cloud, you must enable the Anonymize data option in Collibra Console. This option is by default disabled.

Warning   Currently, if you enable the data anonymization process you can no longer use automatic data classification via the Data Classification platform. However, you can still classify and anonymize profiling results if you use Edge.

# Data profiling results with Jobserver

When you registered a data source via Jobserver, and you click the **Result** button of a data source registration activity, the **Data Profiling Results** dialog box opens.
A data source registration activity can be:

- Creating schema from JDBC
- Creating schema from file
- Updating JDBC schema
- Updating Excel schema
- Updating CSV schema

The **Data Profiling Results** dialog box contains the following information:

| Item | Description |
|---|---|
| Schema | Name of the schema as added to Collibra Data Intelligence Cloud. |
| Status | Status of the data source registration job. |
| Start time | Date and time when the data source registration job has started. |
| End time | Date and time when the data source registration job has completed. |
| Duration | Elapsed time of the data source registration job. |
| Ingestion Details | Summary of the job, including error messages and the list of tables and columns that have been ingested. |
| Profiling Details | The number of tables that have been correctly profiled. |

# Modify the column data type of registered data

When Collibra Data Intelligence Cloud creates a data profile of registered data, it detects the data type of each column. It's possible that Collibra detects a data type that does not correctly correspond with the actual data type, for example the Text data type is detected for a column, but the actual data in the column are dates.

For more information about the data type detection, see Data type detection.

You can update the data type of each column to ensure that the data is properly managed in Collibra.

> Note   If you use the Jobserver to register a data source and you have enabled the Anonymize data option in Collibra Console, Collibra detects the data type of a column during profiling and only anonymizes the data if the data type attribute is Text or Geo. Other data types are not anonymized. If you use Edge to register a data source, these columns are automatically anonymized.

# Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a resource role with the Attribute > Update resource permission.

# Steps

There are two ways to modify a column's data type:

- In the data sources table.
- On the Column's asset page.

## In the data sources table

1. On the main menu, click ⊞, then ▤ **Catalog**.
   - » The Catalog Home opens.
2. In the submenu, click **Data Sources**.
3. Add the **Data Type** column to the table.
4. Expand the schema and table to see the columns.
5. Double-click in the **Data Type** column and choose the correct data type.
6. Click ✔ to apply the change.

## On the Column asset page

1. Look up the column via the **Search** function.

   > Tip   If you don't know the exact name of the column name, you can find it via **Data Catalog → Data Dictionary** and select the **All Schemas** view. Then click the schema that contains the column and click the column whose data type you want to update.

2. In the tab pane, click **Data Profiling**.
3. In the **Metadata** section, double-click the value of the **Data Type** parameter.
4. Select the correct type from the list.
5. Click **Save**.

When you refresh the schema, this change is not overridden.

# Profiling and classification via Edge

## About profiling and classification via Edge

Profiling and classification via Edge is a functionality offered by Collibra for Collibra Data Intelligence Cloud users. The functionality combines both data profiling and data classification in one process.

- Data profiling creates a summary of a data source that is registered with Data Catalog and determines the data type of columns in the data source. The summary mainly contains statistics and graphics to give the user an idea what the registered data is about.
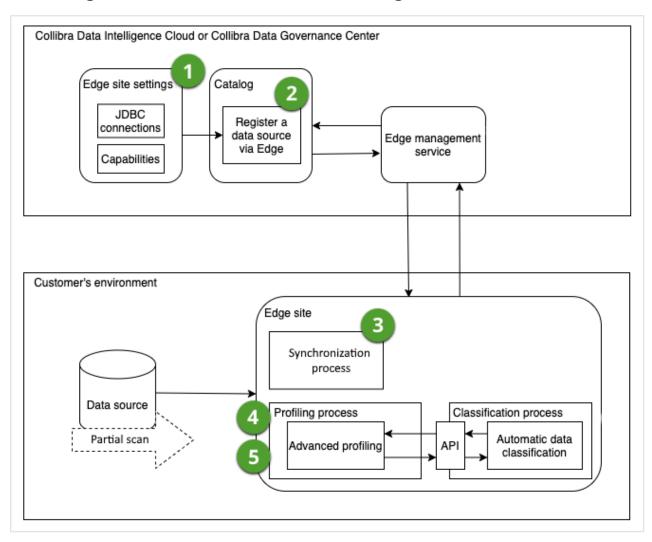
  > Important   Advanced data types are not taken into account when profiling via Edge.

- Automatic Data Classification tries to define the data class of a column. You can accept or reject the suggested data class of each column or add your own new classes.
  Automatic Data Classification can suggest multiple data classes for a column. If the suggestion is accurate, you can accept multiple data classes for the column.

Edge profiles and classifies the data on the Edge site itself and only sends the profiling results and classification suggestions to Collibra Data Intelligence Cloud. As a result, Data Catalog has access to synchronized metadata, profiling results (that are automatically anonymized for columns of type Text and Geo), and classification suggestions, and does not have access to the actual data from your data source.

# Profiling and classification flow in Edge



| Step | Description |
|------|-------------|
| 1 | Create an Edge site with a JDBC connection, a JDBC ingestion capability and a JDBC profiling capability.<br><br>Note   Ensure you have defined the profiling and classification settings. |
| 2 | Register a data source via Edge. |
| 3 | Synchronize one or more schemas. |

| Step | Description |
|---|---|
| 4 | Configure the profiling and classification options for the synchronized schemas. |
| 5 | Profile and classify.<br>The Edge site will initiate the profiling and classification process and send the results to Collibra Data Intelligence Cloud.<br><br>Tip   You can trigger the profiling and classification job manually, set up a schedule or trigger it after synchronizing a schema. |

## Limitations

Profiling via Edge has the following limitations:

- Advanced data types are not supported.
- Not all data sources are certified for Edge.

Automatic Data Classification via Edge has the following limitations:

- Automatic Data Classification via Edge is only available for customers using Collibra Data Intelligence Cloud.
- Currently, data classification on Edge does not retrain the classification model to improve future classification predictions.
- Out-of-the-box, automatic data classification can predict several data classes. You can also create user-defined data classes. Currently, these user-defined data classes are not taken into account by the automatic classification process. You need to assign user-defined data classes manually.
- English is the only supported language, but Automatic Data Classification can run on data in other Latin alphabet-based languages as well.
- Automatic Data Classification needs profiling data to predict the data classes. Data classification is performed automatically after the profiling process on an Edge site. That means that you can only classify columns of data sources registered in Data Catalog via an Edge site that has the JDBC profiling capability.

# Profiling and classification settings

The following settings in the Services Configuration section of the Collibra settings or in Collibra Console are relevant when you want to profile and classify via Edge.

| Setting | Section | Description |
|---|---|---|
| Database registration via Edge | Register data source | An option to enable database registration via Edge.<br><br>• ✔ True: Register a data source via Edge.<br>• ✘ False: Register a data source via Jobserver only.<br><br>Note  Enabling data source registration via Edge does not prevent you from registering a data source via Jobserver as well. |
| Anonymize data | Data profiling | This setting is not relevant. In Edge, data on Edge is automatically anonymized for columns of data type Text and Geo. |
| Database profiling via Edge | Data profiling | An option to enable profiling and classifying synchronized metadata via Edge instead of Jobserver.<br><br>• ✔ True: Profiling and classify via Edge.<br>• ✘ False: Profile via Jobserver and classify via the Data Classification Platform.<br><br>Note  You can only enable Database profiling via Edge if you also enabled Database registration via Edge. |
| Enable Data Classification | Cloud Data Classification configuration | Ensure the **Enable data classification** option in Cloud Data Classification configuration is set to `false`.<br>If the **Enable data classification** option in Cloud Data Classification Configuration is set to `true`, the **Classify** button is available on Column and Table asset pages. This button allows you to classify data via the Data Classification Platform, However, when using profiling and classification via Edge, you no longer need the Data Classification Platform. |

# Configure the profiling and classification options via Edge

Through the profiling and classification options, you can determine:

- whether you want to start the profiling and classification process automatically after each synchronization.
- the default profiling behavior for the schemas, such as whether the profiling is based on all data or on a random subset of the data.
- whether specific schemas do not use the default behavior but instead have their own behavior .
- which schemas you want to profile and classify.

## Prerequisites

- You have created a support ticket to request access to Edge and have created and installed an Edge site.
- Your Edge site has a global role with the following global permissions: Data Catalog and Register Profiling Information.
- Your Edge site has a JDBC profiling capability.
- You have enabled data source registration via Edge.
- You have enabled profiling and classification via Edge.
- You have registered a data source via Edge.
- You have synchronized one or more schemas of the registered database.

## Steps

1. Open a Database asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. Click the **Profiling and classification** tab.
   - » The Profiling and classification options open.

   > Tip   Only the synchronized schemas are available in the list.

4. In the **Default profiling and classification rule** section, click **Edit**.

5. Enter the required information.

| Option | Description |
|---|---|
| Automatically run when a metadata extraction is synchronized | Enable to automatically create a data profile and classify columns every time the synchronization process of one or more schemas finishes.<br><br>This may take a long time. You can also add a schedule to profile and classify at regular intervals. |
| Do not profile unless specified in the schema-specific rules | Select if you don't want to define a default profiling behavior for the schemas.<br><br>Important   Use this option if you only want to profile and classify some of the schemas.<br>If you select this option, Collibra only profiles and classifies the schemas for which a specific profiling and classification rule has been defined. |
| Full scan | Select to, by default, profile the schemas based on all data. |
| Partial scan | Select to, by default, profile schemas based on a subset of the data.<br>If you select this option, the **Maximum number of rows** field becomes available. You can enter the maximum number of rows that you want to use for profiling. By default, the maximum number of rows is 20 000.<br><br>Warning   Partial scan is only available for some data sources. To verify if your data source allows partial scan, see Collibra-provided JDBC drivers. |

For data sources that support partial scan, the **Partial scan** option is selected by default. For data sources that don't support partial scan, the **Do not profile unless specified in the schema-specific rules** option is selected by default.

6. Click **Save**.
7. If you want to define a specific profiling and classification rule for a schema:
   a. In the **Schema profiling and classification rules** section, select the schema.
      » The schema-specific information opens.
   b. Do one of the following:
      - To create a new table rule, click **Add table rule**.
      - To edit an existing table rule, click **Edit** .

c. Enter the required information.

| Option | Description |
| --- | --- |
| Do not profile | Select to indicate you do not want to profile and classify this schema.<br>This option is useful if you want to exclude a schema from the profiling and classification process. |
| Full scan | Select to profile the schema based on all data. |
| Partial scan | Select to profile the schema based on a subset of the data.<br>If you select this option, the **Number of rows scanned (max)** field becomes available. You can enter the maximum number of rows that you want to use for profiling and classification. By default, the maximum number of rows is 20 000.<br><br>Warning  Partial scan is only available for some data sources. To verify if your data source allows partial scan, see Collibra-provided JDBC drivers. |

For data sources that support partial scan, the **Partial scan** option is selected by default. For data sources that don't support partial scan, the **Do not profile** option is selected by default.

d. Click **Save**.

# What's next?

You can now profile and classify the data manually, automatically or add a schedule.

# Profile and classify data via Edge

After you have configured the profiling and classification options, you can start the profiling and classification process for the schemas in the data source.

Tip   Collibra Data Intelligence Cloud only has access to synchronized metadata, profiling results (that are automatically anonymized for columns with data type Text or Geo), and classification suggestions, not to the actual data from your data source.

Important   Advanced data types are not taken into account when profiling via Edge.

## Prerequisites

- You have created a support ticket to request access to Edge and have created and installed an Edge site.
- Your Edge site has a global role with the following global permissions: Data Catalog and Register Profiling Information.
- Your Edge site has a JDBC profiling capability.
- You have enabled data source registration via Edge.
- You have enabled profiling and classification via Edge.
- You have registered a data source via Edge.
- You have synchronized one or more schemas of a registered database.
- You have configured the profiling and classification options.

## Manually profile and classify

1. Open the Database asset page of a registered database.
2. In the tab pane, click ⚙ **Configuration**.
3. Click the **Profiling and classification** tab.
   » The Profiling and classification options open.

   Tip   Only the synchronized schemas are available in the list.

   Important   If you only want to profile and classify one or more schemas, ensure the default profiling and classification option is set to **Don't profile**

> **unless specified in the schema-specific rules**, and that you have only
> defined a specific rule for the relevant schemas.

4. On the **Profiling and classification** tab page, click **Run profiling and classification**.

   » Data Catalog triggers the Edge site to start a profiling and classification job.

   » Depending on your profiling and classification options, the Edge site profiles and
   classifies all or some schemas, based on all synchronized metadata or on a sample.

## Automatically profile and classify after each synchronization

1. Open the Database asset page of a registered database.
2. In the tab pane, click ⚙ **Configuration**.
3. Click the **Profiling and classification** tab.

   » The Profiling and classification options open.

   > Tip   Only the synchronized schemas are available in the list.

4. In the **Default profiling and classification rule** section, click **Edit**.
5. Select **Automatically run when a metadata extraction is synchronized**.
6. Synchronize one or more schemas.

   » When the schemas are synchronized, Data Catalog automatically triggers the
   Edge site to start a profiling and classification job.

   » Depending on your profiling and classification options, the Edge site profiles and
   classifies all or some schemas, based on all synchronized metadata or on a sample.

## Profile and classify based on a schedule

1. Open the Database asset page of a registered database.
2. In the tab pane, click ⚙ **Configuration**.
3. Click the **Profiling and classification** tab.

   » The Profiling and classification options open.

   > Tip   Only the synchronized schemas are available in the list.

4. In **Synchronization schedule**, click **Add Schedule** to add a new schedule, or ✎ to edit an existing schedule.

   » The **Edit schedule** dialog box appears.

5. Enter the required information.

| Field | Description |
|---|---|
| Repeat | The interval when you want to synchronize the schemas automatically, for example daily, weekly or based on a Cron expression. |
| Cron | The Quartz Cron expression that determines when the synchronization takes place.<br><br>This field is only visible if you select `Cron expression` in the **Repeat** field. |
| Every | The day on which you want to synchronize the schemas, for example Sunday.<br><br>This field is only visible if you select `Weekly` in the **Repeat** field. |
| Every first | The day of the month on which you want to synchronize the schemas , for example Tuesday.<br><br>This field is only visible if you select `Monthly` in the **Repeat** field. |
| At | The time at which you want to synchronize the schemas automatically, for example 14:00.<br><br>This field is only visible if you select `Daily`, `Weekly` or `Monthly` in the **Repeat** field. |
| Time zone | The time zone for the schedule. |

6. Click **Save**.

   » The profiling and classification job starts according to the schedule.

   » Depending on your profiling and classification options, the Edge site profiles and classifies all or some schemas, based on all synchronized metadata or on a sample.
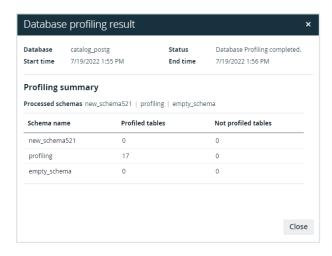
## What's next?

The Edge site completes the profiling and classification process and sends the results to Collibra Data Intelligence Cloud.

- You can see the profiling and classification job in the list of activities. When the activity is completed, the results page gives an overview of the profiled and classified data.
- You can find the profiling information and charts in the Table and Column asset pages.
- You can find the suggested data classes and provide feedback on them via the Table asset page (in the **Columns** tab page) and the Column asset page (in the **Data Profiling** tab page).
- In the **Configuration** of the Database asset page, if a schema is profiled and classified, you see a check symbol (✓) next to the schema name. If the profiling or classification of a schema failed, an exclamation mark (❗) is shown.

| Schema name | ↻ | ▽ |
|---|---|---|
| adt_rename | | ✓ ▦ |
| personal data redshift | | ✓ ▦ |
| pg_catalog | | ❗ ▦ |
| pg_internal | | ✓ ▦ |

# Edge profiling and classification report

When you profiled and classified a data source via Edge, and you click the **Result** button of the profiling and classification activity, the **Database profiling result** dialog box opens.

The following information is available:

| Item | Description |
|---|---|
| Database | Name of the database as added to Collibra Data Intelligence Cloud. |
| Status | Status of the data source profiling and classification job. |
| Start time | Date and time when the profiling and classification job started. |
| End time | Date and time when the profiling and classification job was completed. |
| Duration | Elapsed time of the profiling and classification job. |
| Profiling Summary | The names of the schemas that have been profiled and classified.<br><br>For each schema, you find information on:<br><br>• The number of tables that have been profiled and classified.<br>• The number of tables that have not been profiled and classified.<br>  ◦ If none of the existing tables are profiled and classified, then the required permissions may not be set correctly.<br>  ◦ If only some tables are not profiled and classified, the profiling and classification SQL query for those tables failed. Check the Edge log files for detailed information on why that happened. |

The profiling and classification results are available in the asset pages. See Data profiling information.

# Data profiling information

If you create a data profile of registered data, data profiling information is generated in the table and column assets.

If you use Jobserver to register the data source, the information depends on the profile options that you selected when you registered the data source.
If you use Edge to register the data source, the information is only available after you profiled the data. For an overview of the data that becomes available after the registration of a data source via Edge, see Data source registration information.

| Column attribute | Profiling option (Jobserver) | Statistics | Description | Retrieved from JDBC property |
|---|---|---|---|---|
| Column Name | No option selected | N/A | The column name of the registered table. | COLUMN_NAME |

| Column attribute | Profiling option (Job-server) | Statistics | Description | Retrieved from JDBC property |
|---|---|---|---|---|
| Data Type | Store Data Profile<br><br>If you want to have Advanced Data Type detected, select Detect advanced data types | N/A | The data type of the column. This type is detected by the profiling process. This can differ from the **Technical Data Type** value.<br><br>For example, if a database has a column with text as data type, and the column contains only integer values, the profiling process will set the *Whole Number* data type instead of text.<br><br>If you enable the Anonymize data option in Collibra Console, Collibra anonymizes data in Column assets that have data type Text and Geo.<br><br>If the profiling process has | |

| Column attribute | Profiling option (Job-server) | Statistics | Description | Retrieved from JDBC property |
|---|---|---|---|---|
| | | | detected a wrong data type, you can update it afterwards.

Collibra anonymizes data in Column assets that have data type Text and Geo. | |
| Description from Source | No option selected | N/A | The description of the column in the data source. | REMARKS |
| Row Count | Store Data Profile | Exact | The number of rows in the data source. | |
| Empty Values Count | Store Data Profile | Exact | The number of rows that are empty. | |
| Number of distinct values | Store Data Profile | Exact or approximate depending on column cardinality | The number of unique values in the column. | |

| Column attrib-ute | Profiling option (Job-server) | Statistics | Description | Retrieved from JDBC property |
|---|---|---|---|---|
| Chart | Store Data Profile | Depending on chart type | This column displays whether a chart was generated (⊞) for the column or not (no icon available).<br><br>If you hover over the icon, you see a preview of the chart.<br>The chart type varies per data type. Following charts available:<br><br>• Frequency chart<br>• Histogram that shows dis-tribution<br>• Probability dis-tribution curve | |

| Column attrib-ute | Profiling option (Job-server) | Statistics | Description | Retrieved from JDBC property |
|---|---|---|---|---|
| | | | **Note** Charts are not available for the following data types: <br><br> • Data type = Text and Categorical Data = false <br> • Data type = Array <br> • Data type = N/A | |
| Frequency | Store Data Profile | Exact or approximate depending on column cardinality | A bar chart showing frequency data. | |
| Distribution - Histogram | Store Data Profile | Approximate | A histogram showing the representation of the distribution of numerical data. | |

| Column attribute | Profiling option (Job-server) | Statistics | Description | Retrieved from JDBC property |
|---|---|---|---|---|
| Distribution - Probability distribution curve | Store Data Profile | Approximate | A curve showing the representation of the probability distribution of numerical data. | |
| Technical Data Type | No option selected | N/A | Data type of the column as defined in the source. This value can differ from the **Data Type** value. | TYPE_NAME |
| Descriptive statistics (decile, per-centile, quart-iles) | Store Data Profile | Approximate | The value of the calculated statistic of the registered data. | |
| Categorical Data | Store Data Profile | Exact or approximate depending on column cardinality | Indication whether the data in the column is categorical or not. For example, if 100 000 rows are registered and there are only five distinct values, then the data is considered to be categorical. | |

| Column attrib- ute | Profiling option (Job- server) | Statistics | Description | Retrieved from JDBC property |
|---|---|---|---|---|
| Categories | Store Data Profile | Exact or approximat e depending on column cardinality | List of detected categories. This column has only values if the data is categorical. | |
| Char octet Length | No option selected | N/A | Maximum number of bytes in a char- acter type's column. | CHAR_OCTET_ LENGTH |
| Column Pos- ition | No option selected | N/A | The index of the column in the source table. | ORDINAL_ POSITION |
| Is Auto Incre- mented | No option selected | N/A | Indication whether the data in the column is auto- incremented or not. | IS_ AUTOINCREMENT |
| Is Generated | No option selected | N/A | Indication whether the data in the column is gen- erated or not. | IS_ GENERATEDCOLU- MN |
| Is Nullable | No option selected | N/A | Indication whether the column can store NULL values or not. | IS_NULLABLE |

| Column attribute | Profiling option (Job-server) | Statistics | Description | Retrieved from JDBC property |
|---|---|---|---|---|
| Is Primary Key | No option selected | N/A | Indication whether the column is a primary key or not. | True if the primary keys resultSet contains the COLUMN_NAME |
| Maximum Text Length | Store Data Profile | Exact | The length of the longest text value in the column, including white spaces. | |
| Maximum Value | Store Data Profile | Exact | The maximum value in the column. | |
| Mean | Store Data Profile | Exact | The mean of all the values in the column, excluding empty rows. | |
| Median | Store Data Profile | Exact | The median value of the column. | |
| Minimum Text Length | Store Data Profile | Exact | The length of the shortest text value in the column. | |
| Minimum Value | Store Data Profile | Exact | The minimum value in the column. | |

| Column attrib-ute | Profiling option (Job-server) | Statistics | Description | Retrieved from JDBC property |
|---|---|---|---|---|
| Mode | Store Data Profile | Exact or approximate depending on column cardinality | The value with the highest frequency for categorical data. | |
| Number Of Fractional Digits | No option selected | N/A | The number of fractional digits. | DECIMAL_DIGITS |
| Primary Key Name | No option selected | N/A | The name of the primary key composed by the column. | PK_NAME |
| Size | No option selected | N/A | The size of the column in the table. | COLUMN_SIZE |
| Standard Deviation | Store Data Profile | Exact | The statistical standard deviation of numeric values. | |
| Variance | Store Data Profile | Exact | The statistical variance of numeric values. | |

| Column attribute | Profiling option (Job-server) | Statistics | Description | Retrieved from JDBC property |
|---|---|---|---|---|
| Sample | Store Sample Data | N/A | A random sample of the data set that represents the entire data set.<br><br>Note   In Edge, viewing sample data is not linked to the profiling feature. See sample data. | |

| Table attribute | Profiling option (Job-server) | Statistics | Description | From JDBC property |
|---|---|---|---|---|
| Table Name | No option selected | N/A | The table name in the data source. | TABLE_NAME |
| Table Type | No option selected | N/A | The table type in the data source, such as TABLE or VIEW. | TABLE_TYPE |
| Description from Source | No option selected | N/A | The description of the table in the data source. | REMARKS |

# Data profiling of a table

The location of the data profiling information of a Table asset depends on the Catalog experience setting.

- If the setting is enabled, the information is displayed in **Columns** tab page.
- If Catalog experience is not enabled, the **Data Profiling** tab page displays the information.

The following profiling information is available by default:

- Name
- Data Type
- Row Count
- Empty Values Count
- Number of distinct values
- Chart

For more information about these columns and columns that you can add, see Data profiling information.

You can customize the table by clicking on the Display options icon (⊞). For example, to add more columns, click ⊞ → ✎ **Fields** and then click **Select fields**.

# Data profiling of a column

In the **Data Profiling** tab of a Column asset, you can see the details of the column.

The details are grouped in some fixed sections:

| Section | Content |
|---|---|
| Metadata | Contains the metadata of the column, such as data type, column name and so on. |
| Counts | Contains basic content information, such as number of rows and number of distinct values. |
| Basic Statistics | Contains the basic statistics of the data, such as minimum and maximum value. |

Depending on the column's data type, you can find extra sections:

| Section | Content |
|---|---|
| Quantiles | Contains the descriptive statistics of the data.<br><br>This section is only available if the data type is numerical. |
| Categorical Data | Contains the values of the different categories.<br><br>If there are too many values, only the first 50 and last 50 values are displayed. |
| Chart | Displays the statistics in a graphical way. The chart type varies per data type:<br><br>• bar chart: textual data<br>• data distribution: numerical data and date and time data<br><br>See also Data profiling charts. |

> **Note**
> You can anonymize columns with data type Text or Geo.
>
> • If you use Jobserver to register a data source, you can enable the Anonymize data feature in Collibra Console.
> • If you use Edge to register a data source, these columns are automatically anonymized.
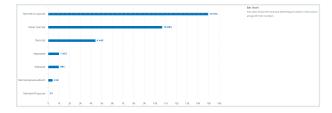
# Data profiling charts

The data profiling process provides a view on the registered data by means of bar charts, distribution data and histograms.

> **Tip**
> • In each chart, you can zoom in by selecting the area of your preference. Click the **Reset zoom** button to return to the original chart view.
> • Charts are not available for every data type. See Data profiling information.
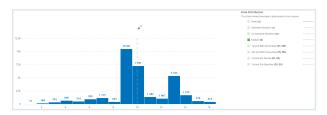
# Bar chart

A bar chart is created when the data type is text. It displays the most and least frequent values of a column along with their number of occurrences.

# Data distribution

The data distribution chart is created when the data type is numerical. It displays how the data is distributed.

In this chart, you can add extra information such as the mean, standard deviation and so on, by selecting the option at the right of the graph.

# Automatic Data Classification

When you register a data source in Collibra Data Intelligence Cloud, the process doesn't stop at ingestion. In order to unlock the full potential of Collibra, the data needs to be contextualized: it needs to be classified and connected to other nodes in the Data Intelligence knowledge graph. Automatic Data Classification adds context to your data.

In the following sections, you will learn more about Collibra's Automatic Data Classification feature.

# About Automatic Data Classification

In Collibra Data Intelligence Cloud, Automatic Data Classification is a feature that analyzes and predicts the content of registered data sources based on a subset of the data itself, helping you to easily gain insights on what kinds of data you have and where it resides. In other words, data classification automatically (with no human input) assigns "class" values to individual columns of data to identify what kind of data is contained in that column. Examples of different data classes are "name", "address", "phone number" and "web browser".

# Why automatic data classification?

When you have ingested data in Data Catalog, the data classification process automatically identifies data structures within the data. As such, it takes less time to learn what kind of data you have ingested.

You can provide feedback by accepting or rejecting the suggested data classes. As a self-learning platform, the Data Classification Platform learns from the feedback, to improve the quality of future predictions.

# Data classification via the Data Classification Platform versus via Edge

The following table shows the differences between classification via the Data Classification Platform or via Edge.

| Part of process | Classification via the Data Classification Platform | Classification via Edge |
| --- | --- | --- |
| Availability | You have enabled Data Classification in Collibra Console. | Data classification is a part of the profiling capability of an Edge site. If you have access to Edge, profiling and classification are available. |
| Sample data | The Data Classification platform requires sample data that needs to be stored in your Collibra environment. | Data classification via Edge classifies data on the Edge site. Sample data is no longer stored in Collibra cloud. |

| Part of process | Classification via the Data Classification Platform | Classification via Edge |
|---|---|---|
| Anonymization | The Data Classification platform uses profiling and sample data to classify. As a result, you cannot classify your data when it is anonymized. | Profiling and classification are performed via an Edge site in the customer's environment. The data is anonymized before it is sent to Collibra Data Intelligence Cloud. |
| Automatic or manual start of the data classification | Data classification must be manually triggered from every table, schema or database. | Data classification is automatically triggered after the profiling process on an Edge site. |
| Retraining | The Data Classification Platform stores your classification selections, along with the associated sample data. This allows to retrain the classification model to improve future classification predictions. | Currently, data classification via Edge does not retrain the classification model.<br>This means that:<br><br>• Your feedback is only stored and is not used for improving classification.<br>• The classification process does not take user-defined classes into account. However, you can create them and assign them manually. |

# Required permissions for Automatic Data Classification

The following table shows the required roles and permissions to use the Automatic Data Classification feature.

| Action | Global Role | Global Permission | Resource Permission (*) | Required for classification via |
|---|---|---|---|---|
| Classify column | Catalog | Catalog | Column asset type's attributes (Asset > Attribute):<br>• Add<br>• Remove<br>• Update<br>Column asset type's data (Asset > Data):<br>• View Samples | Data Classification Platform |
| Classify table | Catalog | Catalog | Table asset type's attributes (Asset > Attribute):<br>• Add<br>• Remove<br>• Update<br>And the resource permissions to classify a column. | Data Classification Platform |
| Accept or reject a classification | Catalog | Catalog | Column asset type's attributes (Asset > Attribute):<br>• Update<br>Column asset type's data (Asset > Data):<br>• View Samples | Data Classification Platform and Edge |

| Action | Global Role | Global Permission | Resource Permission (*) | Required for classification via |
|---|---|---|---|---|
| Add a user-defined clas-sification | Catalog | Catalog > Advanced Data Type:<br><br>• Add | Column asset type's attributes (Asset > Attribute):<br><br>• Add<br>• Update | Data Clas-sification Plat-form<br>and<br>Edge |

(*) As a user, you need a role that has the resource permission.

# Packaged data classes for Automatic Data Classification

The following table shows the data classes that can be detected for columns by the Automatic Data Classification feature.

> Note   This list can evolve over time. When you create a user-defined data class and the number of data samples exceeds a certain threshold, we will add this data class to our system.

| Data Class shown in Collibra Data Intelligence Cloud | Content | Examples |
|---|---|---|
| Canadian Postal Code | Canadian postal codes | • K0C 2K0<br>• G0J 3B0<br>• P3A 1H4<br>• A1N 4E6 |

| Data Class shown in Collibra Data Intelligence Cloud | Content | Examples |
|---|---|---|
| City | Cities | • New York<br>• Los Angeles<br>• Chicago<br>• Houston |
| Country | Countries | • Belgium<br>• Lesotho<br>• Dominica<br>• Nigeria |
| Country code | Countries (short/-code) | • USA<br>• ws<br>• CAF<br>• GIN |
| Credit card number | Credit card number | • 5602223068893246<br>• 1234-1234-1234-1234<br>• 3711-123456-12345<br>• 4123 5123 6123 7123 |
| Currency code | Currency code | • zar<br>• ARS<br>• GBP<br>• kes |
| Date | Date (only) | • 24 January 2004<br>• 11/21/1974<br>• 07-Nov-1982<br>• 11-08-22 |

| Data Class shown in Collibra Data Intelligence Cloud | Content | Examples |
|---|---|---|
| Date time | Datetime | • 2018-08-29 20:25:25.0<br>• 2018-02-05 11:27:10.562<br>• 2017-10-10 05:34:16.216<br>• 2017-07-20 09:03:24.0 |
| Education level | Education (level) | • Doctorate<br>• post-secondary<br>• Doctoral<br>• Upper Secondary School |
| Email | Email | • pdawidas@storify.com<br>• bmcentagartcf@china.com.cn<br>• vgooms6x@barnesandnoble.com<br>• john.smith@[123.123.123.132] |
| Employment status | Employment status | • Freelance<br>• employed part time<br>• office holder<br>• Homemaker |
| Ethnicity | Race | • Hispanic<br>• Latino<br>• White<br>• Asian |
| Filepath | Filepath | • E:\x9xOL\VB2ER_2E\<br>• F:LI\r_dWjux_\<br>• /u_2/tlk4q2/TwaYgn08A/GU/d-fp/z2vHk5iOW/Ael/M_<br>• wUmxr/<br>• BaG_8xxK_m/o1dq4luQ7A/z/kCQXGu.bin |

| Data Class shown in Collibra Data Intelligence Cloud | Content | Examples |
|---|---|---|
| First name | First Name | • Natasha<br>• Manan<br>• Rob<br>• Wojciech |
| Full name | Full name (name + last name) | • lukas yang<br>• Lukas, Yang<br>• Amelia, Dalton<br>• Dickens, Charles |
| Gender | Gender | • M<br>• Male<br>• woman<br>• F |
| IBAN | IBAN - International Bank Account Number | • FO07 4910 6564 9863 03<br>• FR29 5218 3745 58B7 GH7N FYGZ Q50<br>• PS74 TSHR P22C D1DE 5OEB CRUG JRFW W<br>• MK66 115I FYVV SOVS Y00 |
| Internet domain | Web/internet domain | • slashdot.org<br>• usa.gov<br>• time.com<br>• illinois.edu |
| IP address | IP address | • 80.206.17.108<br>• 3a6c:bb28:701a:5aaa:825c:4112:51ea:fadf<br>• 255.139.66.168<br>• 010.010.255.255 |

| Data Class shown in Collibra Data Intelligence Cloud | Content | Examples |
|---|---|---|
| ISBN | ISBN - International Standard Book Number (numeric commercial book identifier) | <ul><li>717393709-4</li><li>106115687-7</li><li>740540459-6</li><li>839089904-3</li></ul> |
| Language | Language | <ul><li>Deccan</li><li>Kazakh</li><li>Zulu</li><li>Greek</li></ul> |
| Language code | Language code | <ul><li>yor</li><li>HAU</li><li>CE</li><li>PS</li></ul> |
| Last name | Surnames / last name | <ul><li>Burke</li><li>Lenaghan</li><li>Balmori</li><li>Balog</li></ul> |
| MAC Address | MAC address | <ul><li>4E-A0-23-78-53-50</li><li>DE:D3:44:A7:7E:13</li><li>a4-53-08-93-70-a4</li><li>83:4f-ca:43:93:32</li></ul> |

| Data Class shown in Collibra Data Intelligence Cloud | Content | Examples |
|---|---|---|
| Marital status | Marital status | <ul><li>unmarried</li><li>Married</li><li>not-in-family</li><li>other-relative</li></ul> |
| Month | Month | <ul><li>Mar</li><li>September</li><li>January</li><li>December</li></ul> |
| NDC Code | FDA NDC code - Food and Drug Administration's National Drug Code | <ul><li>55154-5876</li><li>68927-3491</li><li>58118-0623</li><li>55154-3939</li></ul> |
| Occupation | Occupation | <ul><li>proofer</li><li>transit coach operator</li><li>forging machine tender</li><li>sports worker</li></ul> |
| Personal Email | Email | <ul><li>f0ETKExihcHK@comcast.fr,</li><li>Ffz0Asl0To@comcast.com.br</li><li>jVgNF9v.ranIu@msn.com</li><li>u.L79@verizon.net</li></ul> |

| Data Class shown in Collibra Data Intelligence Cloud | Content | Examples |
|---|---|---|
| Phone number | Phone number | • 532-555-0185<br>• +1 212 555 3000<br>• 829-394-8017<br>• 973-491-8723<br>• 206.782.8410 |
| Religion | Religion | • Buddhist<br>• Confucian<br>• Protestant Anabaptist<br>• Protestant Adventist |
| Routing Number (ABA) | Routing Number | • 058327451<br>• 675702815<br>• 805759224<br>• 305532637 |
| SSN | SSN - Social security number | • 559-03-4491<br>• 284-34-1408<br>• 499-81-8467<br>• 576-17-9443 |
| Street address | Address (first line)<br>Street + number | • 4 Orinda Way<br>• 61 Broadway |
| Time | Time | • 8:52 AM<br>• 7:36 PM<br>• 06:52<br>• 17:08:15 |

| Data Class shown in Collibra Data Intelligence Cloud | Content | Examples |
|---|---|---|
| Title | Honorific | • Honorable<br>• Rev.<br>• Mr<br>• Ms |
| UK Drivers License Number | Drivers License | • ENArq262033Xj32333<br>• ABzPt058106IA18871<br>• wklrS604032zb31785<br>• smeel761300Rc02703 |
| UK National Health Service (NHS) Number | Health Service | • 375 251 3810<br>• 537 649 5407<br>• 784 382 2399<br>• 534 293 9797 |
| URL | URL | • www.sohu.com<br>• http://www.googleweblight.com<br>• https://twitter.com<br>• ftp://mydomena.org/folder1<br>• http://www.-goolgle.com/search?query=my+query |
| US Adoption Taxpayer Identification Number (ATIN) | Tax Identifier | • 944-93-7219<br>• 930-93-3562<br>• 942-93-6471<br>• 932-93-3182 |

| Data Class shown in Collibra Data Intelligence Cloud | Content | Examples |
|---|---|---|
| US Drivers License Number | Drivers License | • QP080580F<br>• W5060999229<br>• Xm939887D<br>• 70kQF62641 |
| US Employer ID | Employer | • 41-0506939<br>• 91-0675223<br>• 43-2942382<br>• 77-4827140 |
| US Individual Taxpayer Identification Number (ITIN) | Tax Identifier | • 915-78-5757<br>• 937-83-1696<br>• 929-75-9337<br>• 966-88-3886 |
| US License Plate Number | License Plate | • 0HB8609<br>• 0qM6428<br>• 0VS0864<br>• 0lq7470 |
| US State | US States | • Illinois<br>• Indiana<br>• Iowa<br>• Kansas |
| US State code | US state code | • il<br>• WI<br>• ut<br>• MT |

| Data Class shown in Collibra Data Intelligence Cloud | Content | Examples |
|---|---|---|
| UUID | GUID/UUID | • 0ee585a5-6bd3-4fde-9383-827095ed08f3<br>• 00000000-0000-0000-0000-000000031108<br>• 0a4281c9-0b6c-4095-b1b6-d8b417cfa952<br>• ffe27556-7c0d-4007-95c4-306633af3f14 |
| Vehicle Identification Number (VIN) | Vehicle | • 4JGDF7DE1EA269698<br>• WDAPF3CC1B9465179<br>• WDBAB33A8EA076439 |
| Web browser | Web browser | • Mozilla<br>• Netscape<br>• Chrome |
| Weekday | Weekday | • Wednesday<br>• Fri<br>• Wed<br>• Tue |

# Calculation components for Automatic Data Classification

The following components are used to calculate data classes via Data Classification Platform or Edge:

| Component | Purpose |
|-----------|---------|
| Neural network | A machine learning tool that is continuously trained to identify linguistic patterns. Training data has been collected to have an initial set of patterns. |
| Regex matcher | A wide range of regular expressions to identify matching patterns. When the matched types in a column exceeds a certain threshold, the result is used in the final calculation of the data class. |
| Dictionary search | The classification is based on a dictionary attack. Multiple data classes only have a limited number of possible values, for example countries, country codes, currencies and days of week. These are all stored in a dictionary.<br>The sample data is matched against these dictionaries. |
| Aggregator | The aggregator gathers the responses from the neural network, regex matcher and dictionary search and creates a final response based on underlying algorithms. |

# How does retraining work?

Currently, data classification on Edge does not retrain the classification model to improve future classification predictions.

In Data Classification Platform, the calculations are all based on the data samples received by the Data Classification Platform. Every time you accept a predicted data class, the sample data used to calculate that data class is added to the Data Classification Platform, to improve future data class predictions. See also Feedback on Automatic Data Classification.

Example
Assume you have a single column, C, containing sample data [a,b,c,d]. You classify this column, and the classification algorithm returns class x with confidence 70%. If you accept this class, then future columns containing the values [a,b,c,d] will be slightly more likely to be classified as x. In the future, a column with the same sample data may be classified as x with confidence 71%. The same can be said for a rejection of the above classification, with future results returning a confidence of, for example, 65%.

Note   In reality, changes will be more discrete and take more than one accepted or rejected data class to become effective.

# Automatic Data Classification via the Data Classification Platform

When you register a data source, you can store a data profile and sample data. This is required if you want to classify columns in the data set. The Data Classification Platform predicts the data classes of selected columns and sends them back to Collibra Data Intelligence Cloud, where you confirm or reject the suggested data classes. The Data Classification Platform uses your feedback to retrain the platform and improve future data classifications.

Warning   If you want to use the Data Classification Platform, request it via your Collibra contact or create a support ticket. See also Data Classification Platform set-up.

# Limitations

- Automatic data classification via the Data Classification Platform is a cloud service. Only if your on-premises environment can reach the cloud service, you can use data classification.
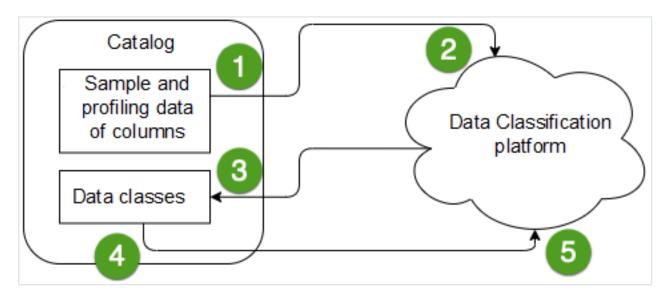
- Out-of-the-box, automatic data classification can predict several data classes. However, you can also create user-defined data classes to increase its prediction quality.
- The only supported language for data classes is English.
- Automatic Data Classification needs sample data and profiling data to be able to predict the data classes.

> Note   You can create sample data and profiling data by registering a data source and choosing to create sample data and profiling data or by importing the data via the Catalog API.

- Automatic Data Classification only works for columns of data sources that are registered in Data Catalog with sample data and profiling data.

# Automatic data classification flow via Data Classification Platform

In the following schema, you can see the different steps of an automatic data classification flow.

| Step | Description |
|------|-------------|
| 1 | You select the columns that you want to classify and send their sample and profiling data to the Data Classification Platform. See Classify columns in a table |
| 2 | The Data Classification Platform predicts the data classes of the columns. |
| 3 | The Data Classification Platform sends the data classes to Collibra. |
| 4 | You provide feedback by accepting or rejecting the predicted data class of each column or by adding your own new classes. The Data Classification Platform can predict multiple data classes for one column. If the prediction is accurate, you can accept multiple data classes for one column. |
| 5 | Your data class selections are sent to the Data Classification Platform. The Data Classification Platform stores your selections, along with the associated sample data, to retrain the classification model and improve future classification predictions. |

# Data Classification Platform set-up

If you want to start using the Data Classification Platform, request it via your Collibra contact or create a support ticket.

## Requirements

You can only use automatic data classification if you comply with the following requirements.

- Data Catalog experience is enabled in the DGC service configuration.
  - » This will give you access to the improved Schema asset page.

- You are using profiling data within Data Catalog.

> Note   Be aware that after you accept the predicted data classes, all sample data and profiling data is stored on the Data Classification Platform.

## Location

We highly recommend to use a Data Classification Platform running in the same region as your Collibra environment.

Currently, Collibra can provide Collibra Data Governance Center environments in Amazon AWS® regions in the following locations:

- United States
- European Union
- United Kingdom
- Canada
- Australia

# Classify columns in a table

By classifying columns in a table, Collibra's Automatic Data Classification predicts their data structures, after which, you can accept or reject the prediction.

> Tip   This information is specific to Data Classification Platform. For information on classifying via Edge, see About profiling and classification via Edge.

You can classify columns via a:

- Database asset page
- Schema asset page
- Table asset page

> Tip   You can also use the physical data connector to manually select a data class for individual columns.
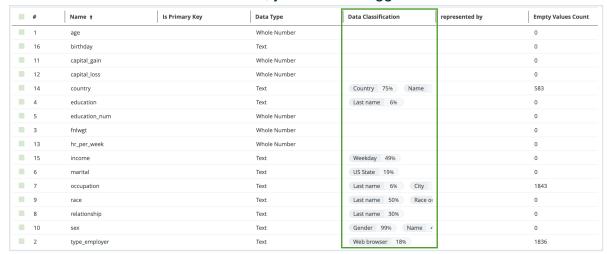
## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have created a support ticket via Zendesk to access to the Automatic Data Classification platform.
- You have configured Automatic Data Classification for the DGC service.
- You have the correct permissions to classify tables and columns.
- You have registered a data source, including these options:
  - Store Data Profile
  - Store Sample Data
- Data Catalog experience is enabled in the DGC service configuration.
  - » This will give you access to the improved Schema asset page.
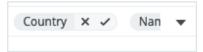- Catalog experience is enabled in the DGC service configuration.

## Via the Database asset page

1. Open the Database asset that contains the tables and columns in the schema you want to classify.
   a. On the main menu, click ⠿, then 🗄 **Catalog**.
      - » The Catalog Home opens.
   b. In the subpages, click **Technology Assets**.
   c. Filter on the Database asset type.
2. Open the relevant database, and then click **Actions** → **Classify**.
   - » You can follow the status of the classification in Activities.
3. Open the database asset with the classified columns.

4. Add the Data Classification column to the table.

   » In the **Data Classification** column, you find the suggested data classes.

| # | Name ↑ | Is Primary Key | Data Type | Data Classification | represented by | Empty Values Count |
|---|--------|----------------|-----------|---------------------|----------------|--------------------|
| 1 | age | | Whole Number | | | 0 |
| 16 | birthday | | Text | | | 0 |
| 11 | capital_gain | | Whole Number | | | 0 |
| 12 | capital_loss | | Whole Number | | | 0 |
| 14 | country | | Text | Country  75%  Name | | 583 |
| 4 | education | | Text | Last name  6% | | 0 |
| 5 | education_num | | Whole Number | | | 0 |
| 3 | fnlwgt | | Whole Number | | | 0 |
| 13 | hr_per_week | | Whole Number | | | 0 |
| 15 | income | | Text | Weekday  49% | | 0 |
| 6 | marital | | Text | US State  19% | | 0 |
| 7 | occupation | | Text | Last name  6%  City | | 1843 |
| 9 | race | | Text | Last name  50%  Race o | | 0 |
| 8 | relationship | | Text | Last name  30% | | 0 |
| 10 | sex | | Text | Gender  99%  Name | | 0 |
| 2 | type_employer | | Text | Web browser  18% | | 1836 |

5. Hover over the classification percentages and accept (✔) or reject (✖) the suggested data class.

   Country  ✖  ✔    Nan  ▼

   • Accepting the classification leaves the classification in the list.
   • Rejecting the classification removes the result from the data classification list.
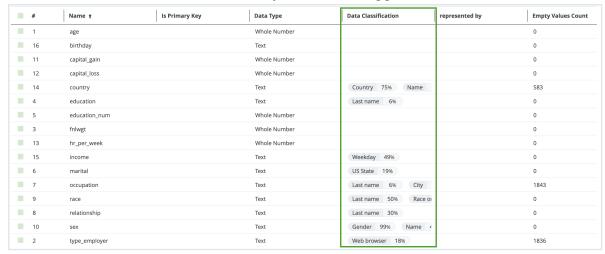
# Via the Schema asset page

1. Open the Schema asset that contains the tables and columns that you want to classify.
   a. On the main menu, click ⠿, then ▤ **Catalog**.
      » The Catalog Home opens.
   b. In the subpages, click **Data Sources**.
   c. Click the relevant schema.
2. Click the Tables tab.
3. Select one or more tables from the schema.
4. To classify all columns in the table, click **Actions → Classify**.

   > Tip   To classify one or more specific columns, select the columns, then click **Actions → Classify**.

   » You can follow the status of the classification job in Activities.

5. Open the Table asset with the classified columns.
6. Add the Data Classification column to the table.
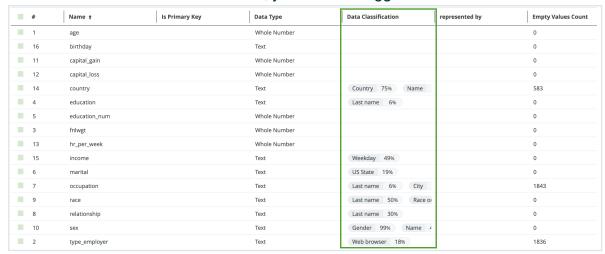   » In the **Data Classification** column, you find the suggested data classes.

| | # | Name ↑ | Is Primary Key | Data Type | Data Classification | represented by | Empty Values Count |
|---|---|---|---|---|---|---|---|
| | 1 | age | | Whole Number | | | 0 |
| | 16 | birthday | | Text | | | 0 |
| | 11 | capital_gain | | Whole Number | | | 0 |
| | 12 | capital_loss | | Whole Number | | | 0 |
| | 14 | country | | Text | Country 75% Name | | 583 |
| | 4 | education | | Text | Last name 6% | | 0 |
| | 5 | education_num | | Whole Number | | | 0 |
| | 3 | fnlwgt | | Whole Number | | | 0 |
| | 13 | hr_per_week | | Whole Number | | | 0 |
| | 15 | income | | Text | Weekday 49% | | 0 |
| | 6 | marital | | Text | US State 19% | | 0 |
| | 7 | occupation | | Text | Last name 6% City | | 1843 |
| | 9 | race | | Text | Last name 50% Race o | | 0 |
| | 8 | relationship | | Text | Last name 30% | | 0 |
| | 10 | sex | | Text | Gender 99% Name | | 0 |
| | 2 | type_employer | | Text | Web browser 18% | | 1836 |

7. Hover over the classification percentages and accept (✔) or reject (✖) the suggested data class.

# Via the Table asset page

1. Open a Table asset that has columns you want to classify.
2. On the Table asset page, do one of the following:
   a. To classify all columns in the table, click **Actions → Classify** in the upper right corner.
   b. To classify specific columns in the table, select the columns and click **Actions → Classify** in the upper right corner.
      » You can follow the status of the classification job in Activities.

3. Open the relevant table, and then add the Data Classification column to the table.
   » In the **Data Classification** column, you find the suggested data classes.

| # | Name ↑ | Is Primary Key | Data Type | Data Classification | represented by | Empty Values Count |
|---|--------|----------------|-----------|---------------------|----------------|--------------------|
| 1 | age | | Whole Number | | | 0 |
| 16 | birthday | | Text | | | 0 |
| 11 | capital_gain | | Whole Number | | | 0 |
| 12 | capital_loss | | Whole Number | | | 0 |
| 14 | country | | Text | Country 75%  Name | | 583 |
| 4 | education | | Text | Last name 6% | | 0 |
| 5 | education_num | | Whole Number | | | 0 |
| 3 | fnlwgt | | Whole Number | | | 0 |
| 13 | hr_per_week | | Whole Number | | | 0 |
| 15 | income | | Text | Weekday 49% | | 0 |
| 6 | marital | | Text | US State 19% | | 0 |
| 7 | occupation | | Text | Last name 6%  City | | 1843 |
| 9 | race | | Text | Last name 50%  Race o | | 0 |
| 8 | relationship | | Text | Last name 30% | | 0 |
| 10 | sex | | Text | Gender 99%  Name | | 0 |
| 2 | type_employer | | Text | Web browser 18% | | 1836 |

4. Hover over the classification percentages and accept (✔) or reject (✖) the suggested data class.

# Feedback on Automatic Data Classification

Each time Collibra predicts data classes for a column, you get the opportunity to send feedback by accepting or rejecting the data class, or by adding a user-defined data class.

To improve future predictions, it is really important to send this feedback.

> Note   When using Edge, the feedback and user-defined classes are only stored, and not used to retrain the classification model.

# Sending feedback

Sending feedback is the act of accepting or rejecting the data classes that are predicted.

- Reject data class: The data class is removed from the column.
  The Data Classification Platform classification model no longer uses the sample data.

- Accept data class: The data class is added to the column.
  The sample data is permanently added to the Data Classification Platform classification model to improve future data class predictions.

To accept or reject a data classification suggestion, hover over the classification suggestion and click the appropriate icon.



For the Data Classification Platform, accepting a data class is more valuable than rejecting, but in general, we recommend that you always send feedback for every prediction. Without your feedback, the classification model cannot be retrained.

# Creating user-defined classes

When columns cannot be classified, you can create your own data classes.

- Avoid duplications. Always check the list of proposed classes before creating a new data class.
- Avoid vague data classes.
- Avoid mixed data classes and accept the best applicable one.

The Data Classification Platform uses this new information to retrain the platform and improve the predictions in the future.

# Create a user-defined data class

If the Automatic Data Classification process cannot detect a data class for a column, you can classify the column yourself.

- If you are using the Data Classification Platform, your new data class will be sent to the Data Classification Platform to improve its future predictive capabilities.
- If you are using Edge, the user-defined classes are only stored. You can use them to manually assign the data class to a column. The user-defined classes are not used to retrain the classification model. See About profiling and classification via Edge.

# Prerequisites

- You have configured Automatic Data Classification for the DGC service.
- You have the correct permissions to classify tables and columns.
- You have registered a data source via Jobserver or via Edge.
  When you use Jobserver, ensure to:
  - Store Data Profile.
  - Store Sample Data. The more sample data, the better the data class predictions.

  > Note   For Jobserver to improve Automatic Data Classification, all sample data and profiling data is stored in the Collibra cloud repository.

- Data Catalog experience is enabled in the DGC service configuration.
  » This will give you access to the improved Schema asset page.

# Create a user-defined data class via the Table asset page

1. Find the table that contains the columns to classify.
2. At the bottom of the **Columns** section, click **See all**.
3. If not yet available, add the Data Classification column to the table.
4. In the **Data Classification** column of the row that you want to classify, click ✏️ .
5. Click the **Select** field.
   The list with existing data classes appears.

6. In the **Select** field, enter the new data class name and press `Enter`.

   > **Note**
   > ○ Data classes are case-sensitive and can contain spaces.
   > ○ You can add more data classes if applicable but avoid it as much as possible.

7. Press `Escape` and click ✔.

   » The new data class is automatically accepted.

# Create a user-defined data class via the Column asset page

1. Find the column you want to classify.
2. In the tab pane, click **Data Profiling**.
3. In the **Data classification** section, click ✏.
4. Enter the new data class name and press `Enter`.

   > **Note**
   > ○ Data classes are case-sensitive and can contain spaces.
   > ○ You can add more data classes if applicable but avoid it as much as possible.

5. Click **Save**.

   » The new data class is automatically accepted.

# Create a user-defined data class via Guided Stewardship

See Add data classes.

# Data Classification Dashboard

The Data Classification Dashboard shows all of the data classes available in your environment.

You can use the Data Classification Dashboard to:

- See information about data classes.
- Add, merge, and delete data classes.
- Link data classes to data concepts and data attributes.

## About the Data Classification dashboard

The Data Classification Dashboard shows the list of data classes in your Collibra environment and gives you more control and visibility on them. When you make changes via the Data Classification Dashboard, feedback is automatically sent to the Data Classification platform.

You access the dashboard via the **Data Classification** subpage in the Stewardship application.

| No. | Name | Description |
|---|---|---|
| 1 | Merge button | A button to merge multiple data classes. |
| 2 | Delete button | A button to delete one or more data classes. |
| 3 | Add button | A button to manually add a new data class. |
| 4 | Table menu (⊞) | The table menu contains buttons to manage the columns shown. |

| No. | Name | Description |
|---|---|---|
| **5** | Table with packaged and manually created data classes | A table that shows all the data classes that exist in your environment. You can also view details about each data class. |
| | Data Classification | The data class name.<br><br>You can manually add, merge, edit or remove the data classes |
| | Column Count | The number of columns classified as the associated data class. |
| | Data Concept | The name of the associated Data Concept assets.<br><br>It connects the data class to your business asset model. |
| | Data Attribute | The name of the associated Data Attribute assets.<br><br>It connects the data class to your logical data model. |
| | Created By | The name of the user who created the class. If the data class is a packaged data class, the user is the *System User*. |
| | Created On | The date the data class was created. |
| | Last Modified By | The name of the user who made the last change. |
| | Last Modified On | The date the data class was last changed. |
| | User Defined | Indicates if the data class was automatically or manually created. |
| **6** | Side pane | A side pane that gives you a clear overview of the data class information of the selected data class. |

# View data class information

You can view data class information on the Classification Dashboard

## Prerequisites

- You have configured Automatic Data Classification in Collibra Console.
- You have the necessary permissions to classify tables and columns.
- You have registered a data source.

## Steps

1. On the main menu, click ⠿, then ↗ Stewardship.
2. In the submenu, click **Data Classification**.
3. Click on a row.
   » The data class information appears in the side pane.

## Data Class side pane

The Data Class side pane gives you a clear overview of related data class information.

When you click the row of a data class in the Data Classification Dashboard, the data class information appears in the side pane.



In the side pane, you find the following information:

| Attribute | Description |
| --- | --- |
| Data class name | The name of the selected data class. You can edit the name by clicking ✎ . |

| Attribute | Description |
|-----------|-------------|
| Data Concepts | The list of data concepts that are associated with the data class.<br><br>This section is only shown if there are associated data concepts. |
| Data Attributes | The list of data attributes that are associated with the data class.<br><br>This section is only shown if there are associated data attributes. |
| <Data class> Columns | The list of columns that are classified with the selected data class. When there are too many columns to show, you can follow a **See all** link. This opens a search results page with all corresponding columns.<br><br>This section is only shown if there are columns with the selected data class. |

# Add data classes

Collibra contains a large number of packaged data classes, but if a certain data class is not available, you can add your own. Data classes that were defined manually are user-defined data classes.

> Tip   You can also create new data classes from a Table or Column asset.

## Prerequisites

- You have configured Automatic Data Classification in Collibra Console.
- You have the necessary permissions to classify tables and columns.
- You have registered a data source.

## Steps

1. On the main menu, click  ⊞ , then ⚲ Stewardship.
2. In the submenu, click **Data Classification**.
   » The table with all data classes is shown.
3. Above the table to the right, click **Add**.
4. Enter the name of a data class and press `Enter`.
   The name of the data class is case-sensitive and it can contain spaces.
   You can enter multiple data classes.
5. Click **Create**.
   » The classes are added.
   If you are using Jobserver, the classes are automatically sent to the Data Classification Platform.
6. Optionally you can link the new classes to a Data Concept or Data Attribute asset.
   a. In the **Data Concept** column, click ✎ .
   b. Click in the **Select** field.
      » The list with existing Data Concept assets appears.
   c. Select one or more Data Concept assets from the drop-down list and click ✔ .
   d. Do the same in the **Data Attribute** column.

## Merge data classes

You can merge two or more data classes via the Data Classification Dashboard. For example, if you have the data classes Email, E-mail and email address, then it is recommended to merge them into the packaged data class Email.

Not only will it keep your data classes list clean, but it will give better results when Collibra performs data classification on ingested data.

> Note   You cannot merge two or more packaged data classes, but you can merge user-defined data classes in a packaged data class. Packaged data classes appear in the **Created By** column as *System User*.

# Prerequisites

- You have configured Automatic Data Classification in Collibra Console.
- You have the necessary permissions to classify tables and columns.
- You have registered a data source.

# Steps

1. On the main menu, click ⊞, then ↗ Stewardship.
2. In the submenu, click **Data Classification**.
3. Select the checkboxes next to the data classes you want to merge.
4. Above the table, click **Merge**.
5. Select the data class you want to merge the selected data classes into.

   > Note
   > - You cannot merge packaged data classes and you can also not merge a packaged data class into a user-defined data class.
   > - The data class attributes Columns Count, Data Concept and Data Attributes are also merged. You can update the list of Data Concepts and Data Attributes after the merge.

6. Click **Merge**.

   

# Edit data classes

You can edit the name of a data class via the Data Classification Dashboard side pane.

# Prerequisites

- You have configured Automatic Data Classification in Collibra Console.
- You have the necessary permissions to classify tables and columns.
- You have registered a data source.

## Steps

1. On the main menu, click ⦂⦂⦂, then ↗ Stewardship.
2. In the submenu, click **Data Classification**.
   » The table with all data classes is shown.
3. Click in the row of the data class that you want to edit.
   » The data class information appears in the side pane.
4. In the side pane, click ✎ next to the data class name.
5. Enter a new name.
6. Click **Save**.
   » The name of the data class is updated.

# Delete a data class

You can delete a data class from the Data Classification Dashboard if it has become obsolete. Note that this is an irreversible action.

# Prerequisites

- You have configured Automatic Data Classification in Collibra Console.
- You have the necessary permissions to classify tables and columns.
- You have registered a data source.

# Steps

1. On the main menu, click ⦂⦂⦂, then ↗ Stewardship.
2. In the submenu, click **Data Classification**.
3. Select the checkboxes next to the data classes you want to delete.
   You cannot delete packaged data classes. These data classes appear in the **Created By** column as *System User* or in the **User Defined** column with ✔.
4. Above the table, click **Delete**.
5. Click **Delete Data Classification**.

# Connect data classes to data layers

You can use the Classification Dashboard to connect data classes to the logical and conceptual data layers.

## Prerequisites

- You have configured Automatic Data Classification in Collibra Console.
- You have the necessary permissions to classify tables and columns.
- You have registered a data source.

## Steps

1. On the main menu, click ⠿, then ↗ Stewardship.
2. In the submenu, click **Data Classification**.
3. In the **Data Concept** or **Data Attribute** column, click ✎ .
4. Click in the **Select** field.
   » The list with existing Data Concept or Data Attribute assets is shown.
5. Click ✔ .

   » The Classification Dashboard creates a relationship between the data class and the logical and conceptual data layers. Column assets that have this data class will be connected to these data layers via their mutual relationship to the data class. Direct relationships between physical and logical information can then be created via Collibra workflows or other methods.

# Guided Stewardship

Guided Stewardship is a set of features designed to help Data Stewards simplify the process of creating connections between physical data assets and their associated logical and conceptual assets. By establishing reliable and fully-connected data structures within your Collibra environment, you can trace relationships across all layers of representation and understand your data in a more complete way.

## Guided Data Stewardship operating model

The Guided Data Stewardship operating model defines the structure of the information in Catalog. For this reason, the Guided Data Stewardship operating model is sometimes also referred to as the Data Catalog operating model.

## Three data layers

The operating model consists of three data layers, representing the three different structural data layers that exist in typical organizations:

- The conceptual data layer represents the overarching structure of objects and elements in your data landscape.
- The logical data layer represents the context-dependent data structures in your organization.
- The physical data layer represents the actual data in your data environment.

The following image shows a complete view of the Data Catalog operating model. It identifies all of the relevant asset types, per data layer, and the relationships that bind them together in the Collibra Data Governance Center.

> Note   Database and System assets are Technology assets that represent the highest level over physical data and logical data organization.

# Conceptual data layer

The conceptual data layer is the highest level of organization in the Data Catalog operating model. It represents the overarching structure of objects and elements within an organization's data landscape. It is where you define concepts, such as Customer and Product and their component fields, without direct reference to system-specific implementations.

Organization of the conceptual data layer is based on many-to-many relationships, which makes the conceptual data layer more concise and flexible than tree-like arrangements that rely strictly on one-to-one and one-to-many relationships.

The conceptual data layer consists of the following asset types:

- Line of Business
- Data Domain
- Data Concept

# Line of Business asset type

The Line of Business asset type is the highest level of abstraction in the conceptual data layer. Also known as business unit or business area, it represents a specific area of business in an organization.

> Example   Finance, Sales, Retail, Investment Management

## Key relation type

Line of Business assets are:

| Related to... | Via the relation type... | Description |
|---|---|---|
| Data Domain assets | Line of Business groups / is grouped by Data Domain | Many-to-many relation, whereby:<br><br>• A Line of Business asset can group many Data Domain assets.<br>• A Data Domain asset can be grouped by many Line of Business assets. |

# Data Domain asset type

Data domains, also known as data categories or subject areas, are high-level, theoretical representations of your data. They represent the structure of concepts in data environments and contain all the different nuances of corresponding business terms.

> Example   Customer, Employee, User, Order, Product

# Key relation types

Data Domain assets are:

| Related to... | Via the relation type... | Description |
|---|---|---|
| Line of Business assets | Business Asset groups / is grouped by Business Asset | Many-to-many relation, whereby:<br><br>• A Line of Business asset can group many Data Domain assets.<br>• A Data Domain asset can be grouped by many Line of Business assets. |
| Data Concept assets | Business Asset groups / is grouped by Business Asset | Many-to-many relation, whereby:<br><br>• A Data Domain asset can group many Data Concept assets.<br>• A Data Concept asset can be grouped by many Data Domain assets. |
| Other Data Domain assets | Data Domain has subtype / is subtype of Data Domain | One-to-many relation, whereby:<br><br>• A Data Domain asset can have many subtype Data Domain assets.<br>• A Data Domain asset can be the subtype of only one Data Domain asset. |

# Data Concept asset type

A Data Concept asset is a high-level theoretical representation of your data and describes one aspect of one or more data domains. These assets represent the most common concepts that are used to organize database content. They allow users to define a context-independent representation of the structure of an organization's data.

They are the most granular level of context-independent structure users can establish within the conceptual data layer, and are comparable to columns in the physical data layer.

> **Example**   Address, Name, ID number, Phone number, Price, Year

For example, if you have a Data Concept asset for Address then this might correlate to a Data Entity asset for Customer Address, Supplier Address and Employee Address.



## Key relation types

Data Concept assets are:

| Related to... | Via the relation type... | Description |
|---|---|---|
| Data Domain assets | Business Asset groups / grouped by Business Asset | Many-to-many relation, whereby:<br><br>• A Data Concept asset can be grouped by many Data Domain assets.<br>• A Data Domain asset can group many Data Concept assets. |
| Other Data Concept assets | Business Asset groups / grouped by Business Asset | Many-to-many relation, whereby:<br><br>• A Data Concept asset can group, and be grouped by, many Data Concept assets. |

| Related to... | Via the relation type... | Description |
|---|---|---|
| Data Attribute assets | Business Dimension classifies / is classified by Asset | Many-to-one relation, whereby:<br>• A Data Concept asset can classify many Data Attribute assets.<br>• A Data Attribute asset can be classified by only one Data Concept asset. |

## Organization based on many-to-many relations

The conceptual data layer is organized such that the relationships between Lines of Business and Data Domain assets, and between Data Domain and Data Concept assets, are many-to-many relationships.

This graph-based approach, based on many-to-many relationships, makes the conceptual data layer more concise and flexible.

Example

In this example, we've identified three lines of business, each of which groups both the Customer data domain and Product data domain. In turn, each data domain groups several data concepts, some of which are grouped by both data domains.

Both data domains group the Name and ID Number data concepts. This is conceivable because Name and ID Number, as Data Concept assets, are abstract representations of these two concepts, rather than specific implementations of them, which are described in the logical data layer and implemented by System assets.

In this way, information stored in the conceptual data layer is kept to a minimum and the Data Domain and Data Concept assets are referred to as often as necessary.



In summary, Line of Business, Data Concept and Data Domain assets are independent assets that do not, by nature, encapsulate or organize the structure of other assets. The Name and ID Number Data Concept assets exist independently of the Data Domain assets that group them. A Customer can have a Name and a Product can have a Name, but you need only one Data Concept asset to encapsulate the idea of "name".

# Conceptual data layer versus the Business Glossary

This section examines the differences and relation between the conceptual data layer and the Collibra Business Glossary.

## Business terms: context-dependent representations of business concepts

In short, the Business Glossary is a system that helps organizations govern their business terms.

> Example   Let's consider the business term Customer, within a multinational consumer goods organization that deals with different consumer groups in different cultural contexts. This organization uses business terms to create a shared understanding of Customer, across different geographical regions. Its offices around the world create their own business terms to encapsulate the specific cultural complexity of a customer, in their own way. Its various business units also have their own definitions, to address different operational, legal and compliance demands.

Business terms are a flexible tool that account for complex business and organizational structures. Anything can be represented by a business term, including the nuanced representations specific to different languages, cultures and branches of business.

Data, on the other hand, can be more explicitly defined and grouped. While there may be several ways to describe Customer, based on cultural and geographic nuance, when we consider data, a customer can be uniquely identified, defined and grouped. This is where the conceptual data layer comes in.

## The conceptual data layer: context-independent representation of the structure of data

A data domain is a container for other data domains and data concepts that encompass associated terminology and definitions that an organization intends to govern.

> Example   Customer Master Data, Product Master Data, Reference Data

While business terms represent Customer in the context of a specific language, culture or branch of business, a customer data domain represents the structure of Customer in a data environment, and encapsulates all of the different nuances of the business term. By abstracting the idea of Customer in a data domain, one can start to consider how customers can be represented by physical data.

The same applies to data concepts, such as Year, Date, Address, and Name. While there may be many business terms that represent Year, across different teams and geographies, the data concept encapsulates all of them and creates a layer of abstraction that allows you to define high-level data structures.

# Logical data layer

The logical data layer defines data structures within an organization's systems, whereas the conceptual data layer represents context-independent data structures within an organization.

The Data Entity-Data Attribute structure is closely related to the Data Domain-Data Concept structure of the conceptual data layer. The main difference between the two is that the conceptual data layer is context-independent, whereas the logical data layer describes the structure in an individual System.

The logical data layer consists of the following asset types:

- Data Model
- Data Entity
- Data Attribute

The logical data layer can be visualized as a tree-like structure, starting with a high-level System and Data Model assets, and branching out with implementation-specific Data Entity and Data Attribute assets.

> Note   Although the System asset type is a Technology Asset, it adds higher-level
> structure to the logical data layer and is considered part of the logical data layer.

## Data Model asset type

The Data Model asset is the highest level of organizational structure in the logical data
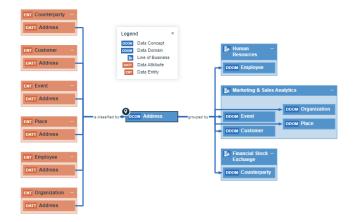layer, and defines the specific structure of data in a System.

## Key relation types

Data Model assets are:

| Related to... | Via the relation type... | Description |
|---|---|---|
| System assets | System implements / is implemented in Data Model | One-to-one relation, whereby:<br><br>• A System asset can implement only one Data Model asset.<br>• A Data Model asset can be implemented in only one System asset.<br><br>> Note   The one-to-one nature of this relationship is what makes Data Models – and, therefore, the entire logical data layer – context-dependent, as opposed to the context-independent conceptual data layer. |
| Data Entity assets | Data Model contains / is contained in Data Entity | One-to-many relation, whereby:<br><br>• A Data Model asset can contain many Data Entity assets.<br>• A Data Entity asset can be contained in only one Data Model asset. |

## Data Entity asset type

Data Entity assets are the logical data layer and correlate to Data Domain assets of the conceptual data layer. Data Entity assets can be thought of as system-specific implementations of Data Domain assets.

For example, if you have a Data Concept asset for Address then this might correlate to a Data Entity asset for Customer Address, Supplier Address and Employee Address.

# Key relation types

Data Entity assets are:

| Related to... | Via the relation type... | Description |
|---|---|---|
| Data Model assets | Data Entity is part of / contains Data Model | One-to-many relation, whereby:<br>• A Data Entity asset can be part of or contained in only one Data Model asset.<br>• A Data Model asset can contain multiple Data Entity assets. |
| Data Domain assets | Data Domain (Business Dimension) classifies / is classified by Data Entity (Asset) | Many-to-many relation, whereby:<br>• A Data Domain asset can classify many Data Entity assets.<br>• A Data Entity asset can be classified by many Data Domain assets. |
| Data Attribute assets | Data Entity contains / is part of Data Attribute | One-to-many relation, whereby:<br>• A Data Entity asset can contain many Data Attribute assets.<br>• A Data Attribute asset can be part of or contained in only one Data Entity asset. |

# Data Attribute asset type

Data Attribute assets are the logical data layer and correlate to Data Concept assets of the conceptual data layer. They can be thought of as system-specific implementations of Data Concept assets.

## Key relation types

Data Attribute assets are:

| Related to... | Via the relation type... | Description |
|---|---|---|
| Data Entity assets | Data Entity contains / is part of Data Attribute | One-to-many relation, whereby:<br><br>• A Data Entity asset can contain many Data Attribute assets.<br>• A Data Attribute asset can be contained by only one Data Entity asset. |
| Data Concept assets | Data Concept classifies / is classified by Data Attribute | One-to-many relation, whereby:<br><br>• A Data Concept asset can classify many Data Attribute assets.<br>• A Data Attribute asset can be classified by only one Data Concept asset. |

# Physical data layer

The physical data layer represents the actual data – the schemas, tables and columns – in an organization's systems.

The physical data layer consists of the following asset types:

- Schema
- Table
- Column

> **Note**
> - Although the Database asset type is a Technology Asset, it is considered part of the physical data layer.
> - The Schema, Table and Column assets in a Collibra Data Intelligence Cloud environment are almost never created manually; rather, they are automatically created via the Data Catalog ingestion process, when registering a data source.

# Schema asset type

A Schema is the highest level of physical structure in a Database. It defines, in a formal language, the structure of the tables and columns in the database.

## Key relation types

Schema assets are:

| Related to... | Via the relation type... | Description |
|---|---|---|
| Database assets | Database has / belongs to Schema | One-to-many relation, whereby:<br><br>• A Database asset can have many Schema assets.<br>• A Schema asset can belong to only one Database asset. |
| Table assets | Schema contains / is part of Table | One-to-many relation, whereby:<br><br>• A Schema assert can contain many Table assets.<br>• A Table asset can be part of only one Schema asset. |

# Table asset type

Table assets represent the physical tables in a data environment.

# Key relation types

Tables assets are:

| Related to... | Via the relation type... | Description |
|---|---|---|
| Schema assets | Table is part of / contains Schema | One-to-many relation, whereby:<br><br>• A Table asset can be a part of only one Schema asset.<br>• A Schema asset can contain many Table assets. |
| Column assets | Table contains / is part of Column | One-to-many relation, whereby:<br><br>• A Table asset can contain many Column assets.<br>• A Column asset can be a part of only one Table asset. |

# Column asset type

Column assets represent the physical columns in a data environment. It is the lowest level of definition in the physical data layer.

# Key relation types

Column assets are:

| Related to... | Via the relation type... | Description |
|---|---|---|
| Table assets | Column is part of / contains Table | One-to-many relation, whereby:<br><br>• A Column asset can be a part of only one Table asset.<br>• A Table asset can contain many Column assets. |

| Related to... | Via the relation type... | Description |
|---|---|---|
| Data Attribute assets | Data Attribute represents / represented by Column | One-to-many relation, whereby:<br><br>• A Data Attribute asset can represent many Column assets.<br>• A Column asset can be represented by only one Data Attribute asset. |

# Technology Assets

Two Technology Assets are included in the Data Catalog operating system:

- System, which is part of the logical data layer.
- Database, which is part of the physical data layer.

## Database asset type

Database assets represent the physical databases in your data environment. They are the highest level of physical data organization in a data environment. Database assets should have specific names, and implement specific technologies, such as PostgreSQL.

## Key relation types

Database assets are:

| Related to... | Via the relation type... | Description |
|---|---|---|
| System assets | System groups / is grouped by Database | One-to-many relation, whereby:<br><br>• A System asset can group many Database assets.<br>• A Database asset can be grouped by only one System asset. |

| Related to... | Via the relation type... | Description |
|---|---|---|
| Schema assets | Database has / belongs to Schema | One-to-many relation, whereby:<br><br>• A Database asset can have many Schema assets.<br>• A Schema asset can belong to only one Database asset. |

# System asset type

System assets represent executable software that an organization uses to automate business functions that help run the business smoothly and efficiently. Systems can be any commercially available or privately developed software that is running in your environment.

> Example   CRM, ERP and EDW software

## Key relation types

System assets are:

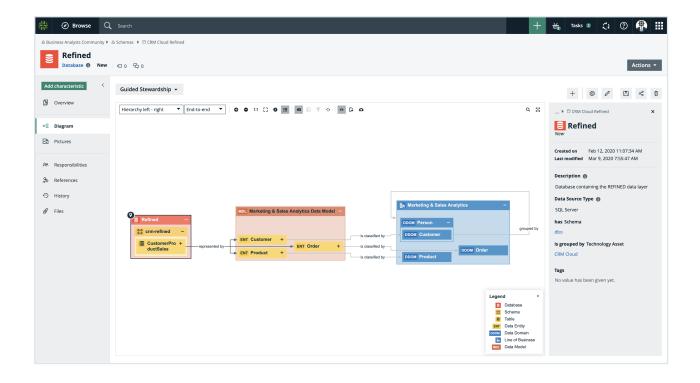| Related to... | Via the relation type... | Description |
|---|---|---|
| Data Model assets | System implements / is implemented in Data Model | One-to-one relation, whereby:<br><br>• A System asset can implement only one Data Model asset.<br>• A Data Model asset can be implemented by only one System asset. |

placeholder

| Related to... | Via the relation type... | Description |
| --- | --- | --- |
| Database assets | System groups / is grouped by Database | One-to-many relation, whereby:<br><br>• A System asset can group many Database assets.<br>• A Database asset can be grouped by only one System asset. |

# Guided Data Stewardship diagram views

For assets in the Guided Data Stewardship operating model, there are two packaged diagram views: Guided Data Stewardship and Guided Data Stewardship - Data Concept. These diagram views show the relation types that bind assets, as established through the Physical Data Connector.
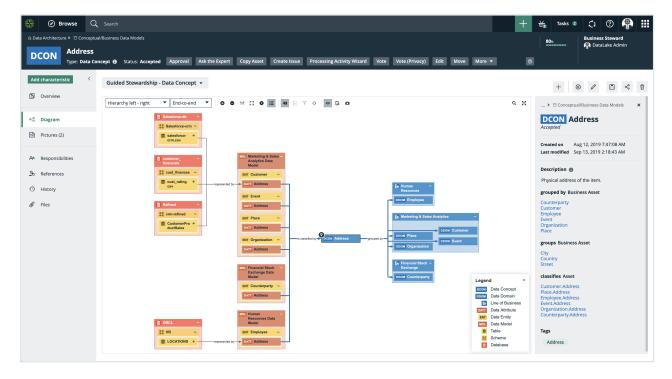
# Guided Data Stewardship view

The Guided Data Stewardship view is the default diagram view designed to help you visualize direct and indirect relations across the entire data environment. For the logical data layer, this view shows the relation types that bind the Data Model, Data Entity, and Data Attribute assets. For the conceptual data layer, it shows the Line of Business and Data Domain assets.

placeholder

placeholder

# Guided Data Stewardship- Data Concept view

The Guided Data Stewardship - Data Concept view is the default diagram view for Data Concept assets only. This diagram view shows the logical and physical data associated with a Data Concept.

For more information, see Diagram views.

# Physical Data Connector

The Physical Data Connector shows a high-level overview of database information on which you can filter.

You can use the Physical Data Connector to:

- Connect the Data Catalog physical data layer to the logical data layer.
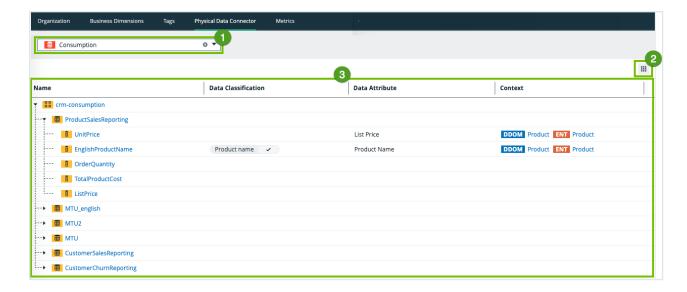- Manually classify columns.

## About the Physical Data Connector

The Physical Data Connector shows a table with a high-level overview of database information. The table has a tree-like structure that enables you to drill down to the column level of a database. It shows the connection between the physical data layer and the logical data layer and enables you to find Data Attribute assets that relate to individual Column assets.

You access the Physical Data Connector via the Physical Data Connector subpage on the Stewardship tab.

| No. | Name | Description |
|---|---|---|
| 1 | Drop-down | A drop-down list to filter on a specific database. |
| 2 | Table menu | The table menu contains buttons for actions you can perform on the table. |
| | ⣿ | A button to manage the columns shown. |

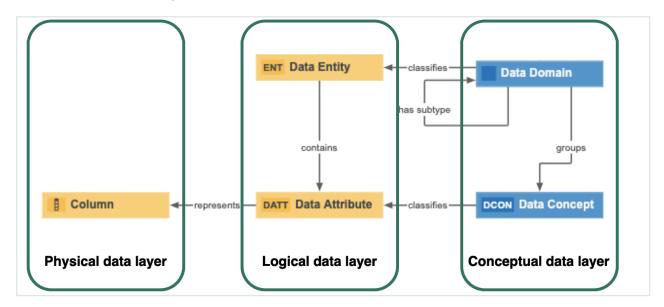| No. | Name | Description |
|-----|------|-------------|
| 3 | Table with database information | A table that shows the content of the registered database and the connections between the physical data layer and logical data layer. |
| | Name | The name of the asset and the icon of the asset type. <br><br> If you click on the asset, the asset page opens. To sort assets alphabetically, click on the column header. |
| | Data Classification | The data class of an asset. <br><br> You can manually add, edit or remove the data class of a Column asset. You can also approve or reject suggested classes |

| No. | Name | Description |
|---|---|---|
| | Data Attribute | The Data Attribute asset linked to the Column asset via relation type "Data Attribute represents / represented by Column". |
| | | When you filter on a Data Domain or Data Entity, the other drop-down lists dynamically update to only show content that relates to your filter. You can select the **Apply filter to all columns in the same table** checkbox to use the same filters to link a Data Attribute to other Column assets in the same table. |
| | |  |
| | | Tip   The physical data connector enables you to quickly connect Data Attribute assets to Column assets. However, you can also connect the physical data layer to the logical data layer via Data Catalog's asset pages by adding a relation of the type Data Attribute represents / represented by Column. |
| | Context | The context of the data. |
| | | This field is read-only and is filled with the Data Domain asset and Data Entity asset related to the Data Attribute asset, if a relation exists. |

# Physical Data Connector relation types

The Physical Data Connector enables you to easily connect the physical data layer to the logical data layer by filtering on the conceptual data layer.

The Physical Data Connector uses the following relation types to connect assets from the different data layers:

- Business Dimension (Data Domain) classifies / is classified by Asset (Data Entity)
- Business Asset (Data Domain) groups / grouped by Business Asset (Data Concept)
- Data Domain has subtype / is subtype of Data Domain
- Business Dimension (Data Concept) classifies Asset (Data Attribute)
- Data Entity contains Data Attribute
- Data Attribute represents Column



# Manually classify columns

The Physical Data Connector enables you to manually add, edit or remove a data class of a Column asset. This is useful, for example, if Automatic Data Classification missed some data classes.

> Tip   You can also automatically classify all columns in a table using Automatic Data Classification.

# Prerequisites

- You have configured Automatic Data Classification for the DGC service.
- You have the correct permissions to classify tables and columns.
- You have registered a data source.
- Data Catalog experience is enabled in the DGC service configuration.

# Steps

1. On the main menu, click ⠿, then ↗ Stewardship.
2. In the submenu, click **Physical Data Connector**.
3. In the drop-down list, filter on a database.
   » The table shows all ingested schemas in the database. You can use the asset tree to drill down to the column level of the database.
4. In the asset tree, find the Column asset that you want to classify.
5. In the Data Classification column, click ✎.
6. Click in the **Select** field.
   » The list with existing data classes appears.
7. In the **Select** field, use the drop-down list to find a data class or enter a new data class name and press `Enter`.

   > Note
   >   ○ Data classes are case-sensitive.
   >   ○ You can add more data classes if applicable, but avoid it as much as possible.
   >   ○ If you created a new data class, it is automatically sent to the Data Classification Platform.
   >   ○ We recommend that you only add one data class to a column.

8. Click ✔.
   » The data class is automatically accepted (✔).

# Connect physical data to logical data

You can use the Physical Data Connector to easily connect a Column asset to a Data Attribute asset via the relation type Data Attribute represents / represented by Column.

A Column asset represents the lowest level of the physical data layer, while a Data Attribute asset represents the lowest level of the logical data layer.

> Tip   You can also add a relation of the type Data Attribute represents / represented by Column via a Data Attribute's or Column's asset page.

# Prerequisites

- You have registered a data source.

# Steps

1. On the main menu, click ⚏, then ⚲ Stewardship.
2. In the submenu, click **Physical Data Connector**.
3. In the drop-down list, filter on a database.
   » The table shows all ingested schemas in the database. You can use the asset tree to drill down to the column level of the database.
4. In the asset tree, find the Column asset that you want to link to a Data Attribute asset.
5. In the **Data Attribute** column, click ✎ .
   » A Data Attribute drop-down list with two filters appears.
6. Link a Data Attribute asset to the Column asset based on the Data Domain and Data Entity filter.
   a. Optionally, select a Data Domain asset and Data Entity asset that are related to the Data Attribute.
      » When you filter on a Data Domain asset or Data Entity asset, the other drop-down lists are dynamically updated to only show content related to your filter.
   b. If you want to use the same filters to find Data Attribute assets for other Column assets in the same table, select the **Apply filter to all columns in the same table** checkbox.
   c. Select the correct Data Attribute asset in the drop-down list.

      > Note   You can only select one Data Attribute asset. The Data Attribute asset must exist in your Collibra environment.

   d. Click ✔ to accept the Data Attribute asset.

» The Data Attribute asset is now linked to the Column asset via the relation type "Data Attribute represents / represented by Column". This relation is also shown on the asset pages of the Column and Data Attribute assets.

» If there is a Data Domain asset and Data Entity asset that is related to the Data Attribute asset, they are shown in the Context column. If you used the filters in the Data Attribute column, the same assets as your filters are shown in the Context column.

Warning   If you click 🗑 to delete a Data Attribute asset in the physical data connector overview, you also delete the relation between the Column asset and the Data Attribute asset from the respective asset pages.

# Working with Amazon S3

Amazon S3 is an online object storage service hosted by Amazon. For more information about Amazon S3, see the Amazon S3 documentation.

In Collibra Data Intelligence Cloud, you can synchronize with Amazon S3 in multiple ways.

| Synchronization method | Advantages and dis-advantages | More information | |
| --- | --- | --- | --- |
| S3 file system integration | The resulting assets represent the folder structure by means of S3 Bucket, Directory, File, Table and Column assets.<br><br>You can't profile and classify columns and tables. | Amazon S3 file system integration | |
| Catalog connector | You can profile and classify the columns and tables in your S3 buckets.<br><br>The folder structure of your S3 bucket isn't represented in Data Catalog. | Jobserver | Edge |
| | | Register an Amazon S3 data source using the AWS Glue Catalog connector | See Register and Synchronze an Amazon S3 file system via Edge |

# About the Amazon S3 file system integration

The Amazon S3 file system integration allows registration of Amazon S3 as a data source in Collibra and synchronization of data in Amazon S3.
After synchronization, the files and directories of Amazon S3 are represented in Collibra by specific asset types, retaining the original names.

> **Note**
> - Only some file types are fully supported.
> - You can restrict the AWS regions to which CollibraData Catalog is allowed to connect. This step is recommended for efficient synchronization.
> - When you use this method, you cannot profile or classify data. See Working with Amazon S3.

## Amazon S3 integration workflow

| Step | What? | Description |
|------|-------|-------------|
| 1 | Register an Amazon S3 file system as a data source | Creates an initial structure of a Storage Catalog domain and S3 File System asset in the selected parent community. |
| 2 | Connect to Amazon S3 | Sets up the connection to Amazon S3. |
| 3 | Create crawlers | Creates crawlers to find and ingest the data of Amazon S3. |
| 4 | Synchronize Amazon S3 | Runs the crawlers to ingest the data of Amazon S3. |

# Password encryption

Collibra's integration of Amazon S3 does not use a separate encryption services, but reuses the Collibra DGC core service encryption method. This method uses the AES/CBC/PKCS5Padding transformation to encrypt your passwords when you connect to Amazon S3.

# Required Amazon Web Services

Collibra relies on AWS Glue and AWS Identity and Access Management to ingest and synchronize data.

## AWS Glue

AWS Glue is an Amazon cloud service to perform extract-transform-load (ETL) processes on data, stored in data sources such as Amazon S3. AWS Glue has the following components:

- Glue crawlers: Glue crawlers analyze and describe a wide range of data sources such as Amazon S3 or MySQL. However, Data Catalog only uses them for the Amazon S3 file system integration.
- Glue database: Glue crawlers store their results in a database in the form of tables and columns. Both the tables and columns in the Glue database contain metadata that describes the content of Amazon S3. Data Catalog reads those databases for data ingestion. The name of the created Glue database is *collibra_catalog_<S3 File System-ID>_<Domain-ID>*.
- ETL processes: The ETL processes can extract data from a data source, process that data, for example, categorize and clean it and produce output. This component is currently not used by Data Catalog.

Though you need an AWS account, you do not have to work in AWS Glue directly, because Collibra does everything for you. For more information about AWS Glue, see the AWS Glue documentation.

> Note   Collibra only uses AWS Glue to ingest data from Amazon S3. All other features, such as crawling other data sources or ETL processes are not integrated.

# AWS Identity and Access Management

Collibra uses the AWS Identity and Access Management (IAM) service to manage access to Amazon S3 and AWS Glue. Similar to AWS Glue, you need an AWS account to use the IAM service, but after setting up the required users and roles, you do not have to work directly with IAM. For more information about IAM, see the IAM documentation.

You need two things in IAM:

- An AWS programmatic user to access Amazon S3 and AWS Glue.
- An IAM role for the crawlers.

## Programmatic user

Collibra needs programmatic access to Amazon S3 and AWS Glue by means of a user. The following policies and permissions are required:

- Policies:

  - AWSGlueServiceRole (AWS managed policy)
  - pass_role (inline policy)
    You can use the following JSON content:

    ```
    {
        "Version": "2012-10-17",
        "Statement":
        [
            {
                "Sid": "VisualEditor0",
                "Effect": "Allow",
                "Action": "iam:PassRole",
                "Resource": "*"
            }
        ]
    }
    ```

- Permissions:
  - In Collibra Data Intelligence Cloud 2020.11 and newer and Collibra Data Governance Center 5.7.7 and newer, the programmatic user needs the following permissions:

```
{
    "Version": "2012-10-17",
    "Statement": [
        {
            "Sid": "VisualEditor0",
            "Effect": "Allow",
            "Action": [
                "glue:GetCrawler",
                "glue:GetCrawlers",
                "glue:DeleteDatabase",
                "glue:GetTables",
                "glue:DeleteCrawler",
                "glue:StopCrawler",
                "s3:ListBucket",
                "glue:GetDatabases",
                "glue:CreateCrawler",
                "glue:GetDatabase",
                "iam:PassRole",
                "glue:StartCrawler",
                "glue:BatchDeleteTable",
                "s3:GetBucketLocation"
            ],
            "Resource": "*"
        }
    ]
}
```

For more information about creating a user with programmatic access, see the IAM documentation.

## IAM role

AWS Glue Crawlers need an IAM role, to allow the crawlers to execute an operation on your behalf. The "pass_role" permission policy of the programmatic user is used to assign this role to the crawler.

You need at least the following parameters:

- Trusted entities: glue.amazonaws.com
- Policies:
    ○ AmazonS3ReadOnlyAccess (AWS managed policy, required when you need to access a private S3 bucket.)
    ○ AWSGlueServiceRole (AWS managed policy)

> Note  You can provide more restrictive permissions to the IAM role, if dictated by your security requirements. Your AWS subject matter expert can create the appropriate permission set using the steps in the IAM documentation. We recommend that you test a crawler with an IAM role that has these permissions in the AWS console, to ensure that it is successful before you use the IAM role in Collibra.

You can also use the IAM role for role-based access control, to authenticate to Amazon AWS without manually entering a user ID and secret access key.

# Amazon S3 asset and domain types

The Amazon S3 file system integration of Collibra Data Intelligence Cloud uses a specific subset of asset types. All of these come out of the box with your software.

| Asset type | Description | Domain type |
|---|---|---|
| Data Asset ▸ Data Element ▸ Column | An atomic unit of data that can be stored in a database table.<br><br>Examples: FST_NM, EMPID | • Physical Data Dictionary<br>• Storage Catalog |
| Data Asset ▸ Data Structure ▸ Table | An implementation of data entities in columns and rows, in a given database system. It is the basic structure of a relational database.<br><br>Examples: Account_tbl, CUST_ADDR | • Physical Data Dictionary<br>• Storage Catalog |
| Data Asset ▸ Data Structure ▸ Table ▸ Database View | A Database View is a virtual table based on the result-set of an SQL statement. | • Physical Data Dictionary<br>• Storage Catalog |

| Asset type | Description | Domain type |
|---|---|---|
| Technology Asset › File Container | An asset type that represents Cloud File Container. | • Storage Catalog<br>• Technology Asset Domain |
| Technology Asset › File Container › Directory | A collection of data that is treated by a computer as a unit, for the purposes of input and output.<br><br>Examples: businessGlossary.xls, dataDictionary05220.csv, datacatalogv25.txt | • Storage Catalog<br>• Technology Asset Domain |
| Technology Asset › File Container › S3 Bucket | An asset type that represents an Amazon S3 Bucket, which is a logical unit of storage containing Amazon S3 Objects. | Storage Catalog |
| Technology Asset › File Group | A collection of physical files which together represent a single logical file. | Storage Catalog |
| Technology Asset › System › File Storage | An asset type that represents a Cloud File Storage bucket. | Storage Catalog |
| Technology Asset › System › File Storage › S3 File System | Amazon S3 (Simple Storage Service) file system abstraction. | Storage Catalog |

# Amazon S3 supported file types

Amazon S3 can contain a wide range of objects in different file types. However, not all file types are fully supported due to limitations of AWS Glue.

The following list shows the file types that are supported by Collibra Data Intelligence Cloud. Note that other file types may work properly as well. For an exhaustive list of supported file types, see the AWS Glue documentation.

- AVRO
- ORC
- PARQUET
- JSON
- BSON
- XML
- ION
- COMBINED_APPACHE
- APACHE
- LINUX_KERNEL
- RUBY_LOGGER
- SQUID
- REDISMONLOG
- REDISLOG
- CSV
- ZIP
- TAR
- RAR
- GZ
- JAR

# Register an Amazon S3 file system

You can register an Amazon S3 file system in Data Catalog.

The newly created S3 file system does not automatically connect to Amazon S3. You connect manually in the S3 File System asset that is created during the registration of the S3 file system.

# Prerequisites

- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a role with the following resource permissions on the S3 community you create when you registered an Amazon S3 file system:
  - Asset: add
  - Attribute: add
  - Domain: add
  - Attachment: add

# Steps

1. On the main menu, click ⠿, then ⊟ **Catalog**.

   » The Catalog Home opens.
2. In the main menu, click the **Create** (＋) button.

   » The **Create** dialog box appears.
3. In the **Create** dialog box, click **Register system**.

   » The **Register system** page appears.
4. In the **Register system** page, click **Amazon S3**.

   » The **Register Amazon S3 file system** dialog box appears.
5. Enter the required information.

   | Field | Description |
   | --- | --- |
   | Community | The parent community in which the initial Amazon S3 structure will be created. |
   | File system name | The name for the S3 file system asset. |

| Field | Description |
| --- | --- |
| Description | The description to provide extra information about the file system.<br><br>This is used as the Description attribute of the S3 File System asset. |
| Owner | The owner name of the data in the created community. |

6. Click **Register**.
   » An S3 File System asset is created.
   » An Storage Catalog domain is created with the same name as the S3 File System asset.
   » The configuration page of the S3 File System asset is automatically opened.

# What's next?

You can now connect to Amazon S3.

# Connect to Amazon S3

To retrieve data from Amazon S3, you have to connect via an S3 File System asset. You always have to do that after registering a new Amazon S3 File System. You can also edit the settings afterwards, for example, if you want to use another Jobserver than the one you originally selected.

# Prerequisites

- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have registered an Amazon S3 file system.
- You have configured one or more Jobservers in Collibra Console. If there is no available Jobserver, the **Register data source** actions will be grayed out in the global cre-

ate menu of Collibra Data Intelligence Cloud.

- You have a programmatic AWS user and IAM role with the required permissions.

## Steps

1. Open an S3 File System asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Connection details** section, click **Edit connection details**.
4. Enter the required information.

| Field | Description |
| --- | --- |
| Connect via | The Jobserver used for synchronizing. |
| Access key ID | The access key ID of the programmatic AWS user. |
| Secret access key | The secret access key of the programmatic AWS user. |
| IAM role | The IAM role to be assigned to the crawlers. |

5. Click **Save**.

## What's next?

You can now create crawlers.

# Connect a file system asset to Amazon S3 via Edge

To retrieve data from Amazon S3, you have to connect via an S3 File System asset. You always have to do that after registering a new Amazon S3 File System. You can also edit the settings, for example, if you want to use another capability than the one you originally selected or if you want to switch to Jobserver.

## Prerequisites

- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have an Edge capability with the S3 synchronization capability template.
- You have registered an Amazon S3 file system.
- You have a global role with the View Edge connections and capabilities global permission.
- You have a programmatic AWS user and IAM role with the required permissions.

## Steps

1. Open an S3 File System asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Connection details** section, click **Edit connection details**.
4. In the right corner, select **Edge**.
5. Select an Edge capability.
6. Click **Save**.

## What's next?

You can now create crawlers.

# Configure role-based Amazon S3 access control

When you register an Amazon S3 file system, you can authenticate to Amazon S3 based on an IAM role. As a result, you can connect to Amazon S3 without an access key ID and secret access key.

## Prerequisites

- You have access to the AWS IAM console.
- You have access to the Amazon EC2 console.
- You have an Amazon EC2 instance.

## Steps

1. In AWS Identity and Access Management, do the following:
    a. Create a new IAM role or select an existing IAM role.
    b. Attach the following policies to the IAM role:
        ◦ AWSGlueServiceRole (AWS managed policy)
        ◦ pass_role (inline policy)

        You can use the following JSON content:

        ```
        {
         "Version": "2012-10-17",
         "Statement":
         [
          {
           "Sid": "VisualEditor0",
           "Effect": "Allow",
           "Action": "iam:PassRole",
           "Resource": "*"
          }
         ]
        }
        ```

2. In the Amazon EC2 console, attach the IAM role to the Amazon EC2 instance.
3. Install the Jobserver service on the Amazon EC2 instance node.

## More information

If the credentials in the Amazon EC2 instance can't be used to authenticate, you can create a credentials file and save it in the **user_home/.aws/** folder. The credentials file should look like this:

```
[default]
aws_access_key_id = <access key ID>
aws_secret_access_key = <secret access key>
```

For more information, see the AWS developer guide.

> Warning   Do not use a credentials file unless absolutely necessary.

## What's next?

You can now connect to Amazon S3 via the jobserver service on the Amazon EC2 instance node.

# Restrict AWS regions

You can restrict the AWS regions to which CollibraData Catalog is allowed to connect to synchronize Amazon S3.

> Note   When there is no restriction, the S3 integration will make requests to all possible AWS regions, which could result in long synchronization times.

> Note   Depending on your environment, you have to follow this procedure either in the Services Configuration section of the Collibra settings or in Collibra Console. This topic shows the procedure in Collibra Console. If your Settings page has a **Configuration** tab, you can follow this procedure in that tab page. For more information, see the online documentation.

## Prerequisites

- You have the ADMIN or SUPER role in Collibra Console.
- You have the SUPER role in Collibra Console.

- You have the ADMIN or SUPER role in Collibra Console.

## Steps

1. Open the DGC service settings for editing:
   a. Open Collibra Console.
      » Collibra Console opens with the **Infrastructure** page.
   b. In the tab pane, expand an environment to show its services.
   c. In the tab pane, click the Data Governance Center service of that environment.
   d. Click **Configuration**.
   e. Click **Edit configuration**.

2. Open the DGC service settings for editing:

    a. Open Collibra Console.

       » Collibra Console opens with the **Infrastructure** page.

    b. In the tab pane, expand an environment to show its services.

    c. In the tab pane, click the Data Governance Center service of that environment.

    d. Click **Configuration**.

    e. Click **Edit configuration**.

3. In the **Register data source** section, enter the required information:

4.

| Setting | Description |
|---------|-------------|
| AWS regions restriction | A list of AWS regions Data Catalog is allowed to connect to. For example, *eu-west-3* and *us-east-2*. For a list of all AWS locations, see the AWS documentation.<br><br>○ If you want to allow Collibra to make a connection to any AWS region, leave the field empty.<br>○ If you remove a region from this list and the region was previously used for an S3 integration, you may want to delete the Glue database from the previously used region manually. By default, Collibra does not remove it. The Glue database has the following naming convention: collibra_catalog_<Asset Id>_<Domain Id><br>For example: collibra_catalog_d3174a88-5ffe-4d50-8fbe-7bf0832ec3af_5d198ce9-4e56-4d0e-a885-58204da50741<br>○ When using Edge, a warning is added to the logs if an invalid region is detected in the restricted regions list. |

5. Click the green **Save all** button.

# Crawlers

A crawler is an automated script that ingests data from Amazon S3 to Data Catalog.

You can create, edit and delete crawlers in Collibra Data Intelligence Cloud. When you synchronize Amazon S3, the crawlers are created in AWS Glue and executed. Each crawler crawls a location in Amazon S3 based on its include path. You can make an S3 bucket accessible for crawlers from the same or other AWS accounts than the account in

which the S3 bucket is located. The results are stored in one AWS Glue database per domain assigned to one or more crawlers. Those databases are ingested in Data Catalog in the form of assets, attributes and relations. The databases are stored in AWS Glue until the next synchronization. At that moment, they are deleted and re-created. The crawlers in AWS Glue are deleted immediately after as the synchronization is finished.

> **Note**
> - By default, AWS Glue allows up to 25 crawlers per account. For more information, see the AWS Glue documentation. This has consequences for Collibra:
>   - If you created crawlers in AWS Glue directly, Collibra can create less crawlers for synchronization.
>   - Because Collibra creates the crawlers in AWS Glue during synchronization, you should avoid having 25 or more crawlers in one S3 File System asset.
>   - You can synchronize several S3 File System assets simultaneously, but if the total number of crawlers exceeds the maximum amount in AWS Glue, synchronization will fail. Since Collibra deletes the crawlers from AWS Glue after synchronization, it is safer to synchronize each S3 File System asset at a unique time.
> - Crawlers in AWS Glue can crawl multiple buckets, but in Collibra, each crawler can only crawl a single bucket.

# Create a crawler

You can create a crawler for an S3 File System asset in Data Catalog.

## Prerequisites

- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have registered an Amazon S3 file system.

- You have configured one or more Jobservers in Collibra Console. If there is no available Jobserver, the **Register data source** actions will be grayed out in the global create menu of Collibra Data Intelligence Cloud.
- You have connected an S3 File System asset to Amazon S3.

## Steps

1. Open an S3 File System asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Crawlers** section, click **Create crawler**.

» The **Create crawler** dialog box appears.

4. Enter the required information.

| Field | Description |
|---|---|
| Domain | The domain in which the assets of the S3 file system are created.<br><br>More information about linking domains to crawlers:<br>○ A specific Storage Catalog domain is created automatically when the S3 File System asset is created. That domain is selected by default. However, you can manually create a new Storage Catalog domain and select it.<br>○ If multiple crawlers point to the same domain, then all assets are created in the same domain.<br>○ If multiple crawlers point to different domains, then all assets are created in their respective domains.<br>○ If multiple crawlers from the same S3 File System asset overlap and point to different domains, then overlapping assets are created in each domain.<br>○ If multiple crawlers from the same S3 File System asset overlap and point to the same domain, then overlapping assets are created once in that domain.<br>○ If crawlers from multiple S3 File System assets overlap and point to different domains, then overlapping assets are created in each domain.<br>○ If crawlers from multiple S3 File System assets overlap and point to the same domain, then overlapping assets are created once in the domain and the S3 Bucket asset has a relation to both S3 File System assets. |

| Field | Description |
|---|---|
| Name | The name of the crawler in Collibra.<br><br>More information about crawler names:<br>◦ You cannot use the same name for two crawlers in the same S3 File System asset.<br>◦ The name of the corresponding crawler in AWS Glue will contain this name. Its name will follow the following convention: `collibra_catalog_<s3fs asset id>_<name_of_the_crawler_in_Collibra>`.<br>◦ The crawler name must be compliant with the AWS Glue limitations:<br>  ■ It has to match the single-line string pattern: `[\u0020-\uD7FF\uE000-\uFFFD\uD800\uDC00-\uDBFF\uDFFF\t]*`.<br>  ■ The length should be between 1 and 255 bytes long, including the fixed prefix that Collibra adds. That means that you can use roughly 65 characters, depending on the characters that were used.<br><br>    Warning   This restriction is imposed by Amazon S3, which allows up to 255 bytes, including the prefix added by Collibra. If you enter too many characters and exceed the byte limit, synchronization fails. |
| Include path | The case-sensitive path to a directory of a bucket in Amazon S3. All objects and subdirectories of this path are crawled.<br><br>For more information and examples, see the AWS Glue documentation. |

| Field | Description |
|---|---|
| Exclude patterns | Glob pattern that represents the objects that are in the include path, but that you want to exclude.<br><br>For more information and examples, see the AWS Glue documentation. |
| Add pattern | Button to add additional exclude patterns. |

5. Click **Create**.

## What's next?

You can now synchronize Amazon S3 manually or define a synchronization schedule.

# Edit a crawler

You can edit a crawler of an S3 File System asset in Data Catalog. For example, you can do this if you want to edit the exclude pattern.

## Prerequisites

- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have registered an Amazon S3 file system.
- You have configured one or more Jobservers in Collibra Console. If there is no available Jobserver, the **Register data source** actions will be grayed out in the global create menu of Collibra Data Intelligence Cloud.
- You have connected an S3 File System asset to Amazon S3.
- You have created one or more crawlers.

# Steps

1. Open an S3 File System asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Crawlers** section, in the row of the crawler that you want to edit, click ✏.
   » The **Edit crawler** window appears.
4. Enter the required information.

| Field | Description |
|---|---|
| Domain | The domain in which the assets of the S3 file system are created.<br><br>More information about linking domains to crawlers:<br>○ A specific Storage Catalog domain is created automatically when the S3 File System asset is created. That domain is selected by default. However, you can manually create a new Storage Catalog domain and select it.<br>○ If multiple crawlers point to the same domain, then all assets are created in the same domain.<br>○ If multiple crawlers point to different domains, then all assets are created in their respective domains.<br>○ If multiple crawlers from the same S3 File System asset over-lap and point to different domains, then overlapping assets are created in each domain.<br>○ If multiple crawlers from the same S3 File System asset over-lap and point to the same domain, then overlapping assets are created once in that domain.<br>○ If crawlers from multiple S3 File System assets overlap and point to different domains, then overlapping assets are created in each domain.<br>○ If crawlers from multiple S3 File System assets overlap and point to the same domain, then overlapping assets are created once in the domain and the S3 Bucket asset has a relation to both S3 File System assets. |

| Field | Description |
|---|---|
| Name | The name of the crawler in Collibra.<br><br>More information about crawler names:<br>○ You cannot use the same name for two crawlers in the same S3 File System asset.<br>○ The name of the corresponding crawler in AWS Glue will contain this name. Its name will follow the following convention: `collibra_catalog_<s3fs asset id>_<name_of_the_crawler_in_Collibra>`.<br>○ The crawler name must be compliant with the AWS Glue limitations:<br>  ■ It has to match the single-line string pattern: `[\u0020-\uD7FF\uE000-\uFFFD\uD800\uDC00-\uDBFF\uDFFF\t]*`.<br>  ■ The length should be between 1 and 255 bytes long, including the fixed prefix that Collibra adds. That means that you can use roughly 65 characters, depending on the characters that were used.<br><br>    Warning   This restriction is imposed by Amazon S3, which allows up to 255 bytes, including the prefix added by Collibra. If you enter too many characters and exceed the byte limit, synchronization fails. |
| Include path | The case-sensitive path to a directory of a bucket in Amazon S3. All objects and subdirectories of this path are crawled.<br><br>For more information and examples, see the AWS Glue documentation. |

| Field | Description |
| --- | --- |
| Exclude patterns | Glob pattern that represents the objects that are in the include path, but that you want to exclude.<br><br>For more information and examples, see the AWS Glue documentation. |
| Add pattern | Button to add additional exclude patterns. |

5. Click **Save**.

# Delete a crawler

You can delete a crawler from an S3 File System asset.

> Note   If you delete an S3 File System asset that contains one or more crawlers, the crawlers are also deleted.

## Prerequisites

- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have registered an Amazon S3 file system.
- You have configured one or more Jobservers in Collibra Console. If there is no available Jobserver, the **Register data source** actions will be grayed out in the global create menu of Collibra Data Intelligence Cloud.
- You have connected an S3 File System asset to Amazon S3.
- You have created one or more crawlers.

## Steps

1. Open an S3 File System asset page.
2. In the tab pane, click ⚙ **Configuration**.

3. In the **Crawlers** section, in the row of the crawler that you want to delete, click 🗑 .

   » The **Delete Crawler** confirmation message appears.
4. Click **Delete crawler**.

# Cross-account crawling

You can make an S3 bucket accessible for crawlers from other AWS accounts than the account in which the S3 bucket is located. To access the external S3 bucket, the programmatic user and the IAM crawling role must be defined in the AWS main account.

## Policy

A policy must be attached to the external S3 bucket to allow:

- the AWS Glue crawler to access and perform S3 actions on an external S3 bucket from another AWS account.
- Data Catalogto execute the S3 GetBucketLocation API on an external S3 bucket via the programmatic user.



You can use the following JSON content:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
        "Sid": "collibra-jobserver-access",
```

```
        "Effect": "Allow",
        "Principal": {
            "AWS": "arn:aws:iam::<enter_id>:role/collibra-job-
server-s3-role"
        },
        "Action": "s3:*",
        "Resource": [
            "arn:aws:s3:::crawler-name",
            "arn:aws:s3:::crawler-name/*"
        ]
    },
    {
        "Sid": "collibra-jobserver-access",
        "Effect": "Allow",
        "Principal": {
            "AWS": "arn:aws:iam::<enter_id>:user/collibra-job-
server"
        },
        "Action": "s3:getBucketLocation",
        "Resource": [
            "arn:aws:s3:::*"
        ]
    }
    ]
}
```

# About synchronizing Amazon S3

Synchronizing Amazon S3 is the process of ingesting metadata from a selected Amazon S3 repository and making the data available in Collibra Data Intelligence Cloud.

When you synchronize Amazon S3, the content of your Amazon S3 repository is analyzed and represented in Collibra by means of assets and their characteristics.
Technically, the synchronization happens in several steps:

1. Collibra creates crawlers in AWS Glue, based on the crawlers defined in Collibra.
2. If AWS Glue contains databases with metadata from a previous synchronization, the databases are deleted.
3. Each AWS Glue crawler crawls a location in Amazon S3 based on its include path. For each domain assigned to one or more crawlers, AWS Glue creates a database with the crawling results.
4. Collibra ingests those databases and creates assets, attributes and relations as

required to match the metadata.
5. The AWS Glue crawlers are deleted.

# Starting the synchronization

You can synchronize manually, or you can automate it by adding a synchronization schedule by means of a cron expression.

You can only synchronize one S3 File System at a time. If a synchronization job is in progress and a second one is triggered, manually or automatically, it will be queued.

If a synchronization job is still running and a new synchronization of the same S3 File System is triggered (manually or automatically), the running synchronization will continue and the new synchronization request is ignored.

# Synchronization results

After synchronization, the resulting assets are in the domain that was specified in the crawler.

> Warning   Do not move the assets to another domain. Doing so may lead to errors during future synchronizations. This is a known limitation.

By default, the assets are shown in a plain list, but you can enable a multi-path hierarchy to show it in a tree structure. For the best result, we recommend that you use the following relations:

1. S3 Bucket contains Directory
2. Directory contains Directory
3. Directory contains File
4. Directory contains File Group
5. File contains Table
6. File Group contains Table
7. Table contains Column

The following images shows the resulting hierarchical table.

Note   In case of a partial synchronization caused by a temporary communication issue, the status of the assets that cannot be synchronized is set to **Missing from source**. During the next fully successful synchronization, the assets are removed or their previous status is restored, depending on their actual status in the source system.

# Naming convention

Synchronizing Amazon S3 relies on a naming convention to match assets during the synchronization process. We highly recommend that you not change the S3 File System asset's full name.

Warning   Editing full name of the S3 File System assets may lead to errors during the synchronization process.

# What if File Groups get the status "Missing from Source" after the S3 synchronization via Edge?

File Group assets can receive the status "missing from source" after the S3 synchronization if the behavior of the AWS crawler is not consistent, meaning AWS classifies files as File Group one day and classifies them as File on another day.
If this happens, File Group assets are created during the first synchronization but no longer exist after the second synchronization, resulting in the status "Missing from source".

Solution:
If you are using Edge, you can add custom parameter **file-group-as-file** to your S3 Edge

capability. By adding the custom parameter, the S3 synchronization will always ingests File groups as File assets. The custom parameter is:

- Name: file-group-as-file
- Value: true

# Synchronize Amazon S3 manually

You can manually start a synchronization job of an S3 File System asset. This can be useful if you want to test your crawlers, or if you want to synchronize immediately.

> Tip   You can also add a synchronization schedule to synchronize automatically.

## Prerequisites

- You have registered an Amazon S3 file system.
- You have configured one or more Jobservers in Collibra Console. If there is no available Jobserver, the **Register data source** actions will be grayed out in the global create menu of Collibra Data Intelligence Cloud.
- You have a programmatic AWS user and IAM role with the required permissions.
- You have connected an S3 File System asset to Amazon S3.
- You have created one or more crawlers.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a resource role with the Configure external system resource permission on the community or domain that contains the S3 File System, for example Owner.
- You have a role with the following resource permissions on the S3 community you created when you registered an Amazon S3 file system:
  - Asset: add
  - Attribute: add
  - Domain: add
  - Attachment: add

## Steps

1. Open an S3 File System asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Crawlers** section, click **Synchronize now**.
   » A notification indicates synchronization has started.

   » The synchronization job appears in the **Activities** list as a bulk synchronization. When the synchronization finishes, the resulting assets, including their attributes and relations, are created, edited or deleted in the selected domain(s) and in the Data Sources page of Data Catalog.
   » The **Synchronization schedule** section displays the time of the last synchronization.

> Note   In case of a partial synchronization caused by a temporary communication issue, the status of the assets that cannot be synchronized is set to **Missing from source**. During the next fully successful synchronization, the assets are removed or their previous status is restored, depending on their actual status in the source system.

## What's next?

You can view a summary of the results from the Activities list.

You can view the assets in their domain.

# Add an S3 synchronization schedule

To keep the content of Collibra Data Intelligence Cloud synchronized with your Amazon S3 File System, you can synchronize manually or create a schedule to automatically do this with a fixed interval.

> Note   You can only create one synchronization schedule.

# Prerequisites

- You have a resource role with the Configure external system resource permission on the community or domain that contains the S3 File System, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have registered an Amazon S3 file system.
- You have a programmatic AWS user and IAM role with the required permissions.
- You have configured one or more Jobservers in Collibra Console. If there is no available Jobserver, the **Register data source** actions will be grayed out in the global create menu of Collibra Data Intelligence Cloud.
- You have connected an S3 File System asset to Amazon S3.
- You have created one or more crawlers.
- You have a role with the following resource permissions on the S3 community you created when you registered an Amazon S3 file system:
    - Asset: add
    - Attribute: add
    - Domain: add
    - Attachment: add

# Steps

1. Open an S3 File System asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Synchronization schedule** section, click **Add Schedule**.

4. Enter the required information.

| Field | Description |
| --- | --- |
| Repeat | The interval when you want to synchronize the schemas auto-matically, for example daily, weekly or based on a Cron expres-sion. |
| Cron | The Quartz Cron expression that determines when the synchronization takes place. <br><br> This field is only visible if you select `Cron expression` in the **Repeat** field. |
| Every | The day on which you want to synchronize the schemas, for example Sunday. <br><br> This field is only visible if you select `Weekly` in the **Repeat** field. |
| Every first | The day of the month on which you want to synchronize the schemas , for example Tuesday. <br><br> This field is only visible if you select `Monthly` in the **Repeat** field. |
| At | The time at which you want to synchronize the schemas automatically, for example 14:00. <br><br> This field is only visible if you select `Daily`, `Weekly` or `Monthly` in the **Repeat** field. |
| Time zone | The time zone for the schedule. |

5. Click **Save**.

# Edit an S3 synchronization schedule

You can edit the synchronization schedule of an Amazon S3 File System asset. For example, you can do this if you think the synchronization job runs too often or not often enough.

## Prerequisites

- You have a resource role with the Configure external system resource permission on the community or domain that contains the S3 File System, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have registered an Amazon S3 file system.
- You have a programmatic AWS user and IAM role with the required permissions.
- You have configured one or more Jobservers in Collibra Console. If there is no available Jobserver, the **Register data source** actions will be grayed out in the global create menu of Collibra Data Intelligence Cloud.
- You have connected an S3 File System asset to Amazon S3.
- You have created one or more crawlers.
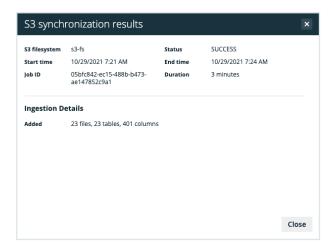- You have added a synchronization schedule.

## Steps

1. Open an S3 File System asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Synchronization schedule** section, click **Edit Schedule**.

4. Enter the required information.

| Field | Description |
|---|---|
| Repeat | The interval when you want to synchronize the schemas automatically, for example daily, weekly or based on a Cron expression. |
| Cron | The Quartz Cron expression that determines when the synchronization takes place. This field is only visible if you select `Cron expression` in the **Repeat** field. |
| Every | The day on which you want to synchronize the schemas, for example Sunday. This field is only visible if you select `Weekly` in the **Repeat** field. |
| Every first | The day of the month on which you want to synchronize the schemas , for example Tuesday. This field is only visible if you select `Monthly` in the **Repeat** field. |
| At | The time at which you want to synchronize the schemas automatically, for example 14:00. This field is only visible if you select `Daily`, `Weekly` or `Monthly` in the **Repeat** field. |
| Time zone | The time zone for the schedule. |

5. Click **Save**.

# Remove an S3 synchronization schedule

You can remove a synchronization schedule from an Amazon S3 File System asset to stop automatically synchronizing Amazon S3.

## Prerequisites

- You have a resource role with the Configure external system resource permission on the community or domain that contains the S3 File System, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have registered an Amazon S3 file system.
- You have a programmatic AWS user and IAM role with the required permissions.
- You have configured one or more Jobservers in Collibra Console. If there is no available Jobserver, the **Register data source** actions will be grayed out in the global create menu of Collibra Data Intelligence Cloud.
- You have connected an S3 File System asset to Amazon S3.
- You have created one or more crawlers.
- You have added a synchronization schedule.

## Steps

1. Open an S3 File System asset page.
2. In the tab pane, click ✿ **Configuration**.
3. In the **Synchronization schedule** section, click **Remove Schedule**.

# View the summary of an Amazon S3 synchronization

After you synchronized Amazon S3, you can view the summary of the results. This shows the impact of the synchronization on the assets in Collibra Data Intelligence Cloud

## Steps

1. Open the Activities list.
2. In the row containing the S3 synchronization job, click **Result**.
   - » The **S3 synchronization results** dialog box appears.

> **Note**
> ○ The **Ingestion Details** section contains information about the total number of resources that were added, modified or removed as a result of the synchronization.
> ○ In case of an error, the **Ingestion Details** section contains additional information about the error.

> Tip   The **Job ID** is useful when troubleshooting your synchronization process with Collibra Support.

# Delete an S3 File System asset from Collibra

You can delete an S3 File System asset from Collibra Data Intelligence Cloud.

> **Note**
> • The crawlers of the S3 File System asset are deleted.
> • The assets that were created by synchronizing are not deleted.

## Prerequisites

- You have registered an Amazon S3 file system.
- You have a global role with the Catalog global permission, for example Catalog

Author.

- You have a resource role with the Asset > Remove resource permission.

## Steps

1. Open an S3 File System asset page.
2. In the view toolbar, click **Actions → Delete**.
   » The **Delete Confirmation** dialog box appears.

   > Tip   If Catalog experience is disabled, the **More** menu is shown instead of **Actions**.

3. Click **Delete S3 File System**.

# Troubleshooting for the S3 file system integration

# Message **Could not add/change/delete crawler '<crawler name>' for S3 File System '<asset name>'.**

You can find more information about the actual problem in the Jobserver logs. The problem is usually described in the AWS SDK error message.

| Cause | Description | Solution |
|---|---|---|
| Incorrect or too limited IAM permissions for the programmatic user defined in the connection details. | While connecting, the verification process only checks that the user can log in, but it doesn't verify permissions. Any further operation may therefore fail if the IAM permissions are wrong or too limited.<br><br>This also applies to the AWS regions. Collibra checks the credentials in the default region, based on the region AWS SDK. Because the IAM service is global, that is sufficient in most cases. However, it is possible to put constraints on specific regions, including the AWS SDK default region. | Edit the IAM permissions or connect to Amazon S3 with another IAM user or role. |

| Cause | Description | Solution |
|---|---|---|
| Maximum number of crawlers in AWS Glue reached. | When you synchronize Amazon S3, Collibra creates crawlers in AWS Glue and executes them. After synchronization, they are deleted.<br><br>By default, each AWS Glue account can only store 25 crawlers. This number can be reached easily, especially if the customer uses AWS Glue apart from Collibra. | • Delete one or more crawlers.<br>• Create an advanced crawler by tweaking the include path and the exclude patterns.<br>• Create additional S3 File System assets and divide the required crawlers between the assets. Then synchronize them at different times.<br>• Synchronize different S3 File Systems at different times.<br>• Ask Amazon support to increase that number.<br><br>For more information, see the AWS Glue documentation. |
| Bucket does not exist | Typo in a bucket name - bucket doesn't exist. | Edit the crawler's include path to correct the bucket name. |
| No permission to access the bucket in Amazon S3. | This includes buckets that exist but belong to different accounts. | Request permission or delete the relevant crawler. |
| Unsupported AWS region. | S3 ingestion in Collibra Data Catalog relies on AWS Glue to analyze S3 buckets. However, AWS Glue is currently not supported in all AWS regions, which may lead to failing crawling creation. The log will display an UnknownHostException. | This is a built-in limitation of AWS Glue. For the list of supported regions for AWS Glue, see the AWS documentation. |

| Cause | Description | Solution |
|---|---|---|
| Incorrect AWS region. | AWS regions can be restricted so that S3 ingestion and synchronization in Collibra Data Catalog can only be performed in the regions your AWS account has access to.<br><br>Example  You will get an error message when:<br><br>• A user with a European account tries to perform S3 ingestion in AWS region Canada.<br>• A user with a European account tries to synchronize S3 buckets for AWS regions Europe and Canada.<br>• A user with a Chinese and Canadian account tries to synchronize buckets for AWS regions Ireland and Canada. | This is a security measure. The AWS regions to which Collibra Data Catalog is allowed to connect can be restricted via Collibra Console. |

Example `[2018-08-03 13:50:38,347] INFO`
`.agent.SprayRoutesProvider [] [] - output: (500 Internal Server`
`Error,{"messageCode":"s3_bucketDoesntExist","messageArguments":`
`["qsdgqsbqfdscs"]})`

# Message **Value not allowed. The connection details of the S3 File System are incorrect.**

| Cause | Description | Solution |
|---|---|---|
| The credentials for the AWS user are incorrect. | This message appears when the credentials for the AWS user are incorrect. The access key ID and/or secret access key are wrong. | Pay attention that they do not contain trailing spaces. |
| Your AWS account doesn't have access to an AWS region where the S3 bucket is located. | This message appears when you add an AWS region in Collibra Console to which your AWS account doesn't have access and then try to ingest an S3 file system. | Make sure that you have access to the AWS region where the S3 bucket is located. |

# Glue Crawler failed and AWS logs show an **Internal server error**

When checking the logs in Jobserver you may notice that one or more crawlers failed in AWS Glue. In that case, you need to open the AWS console and check the crawlers list in AWS Glue. Because crawlers are deleted from AWS Glue after ingestion, you will have to manually re-create the crawlers and run them again before proceeding. The failing crawler has a red exclamation mark and the Failed status. You can check the logs for more information.

Sometimes, the logged message just shows an "Internal server error". The only way to get more information is to contact the Amazon helpdesk. However, we noticed such errors often happen in the following situations

- The number of files to crawl is very large (> 100k)
- There is a series of very small files to crawl (>100).

In both cases, the problem is caused by AWS Glue. All Amazon services are protected against DDoS attacks and they throw throttling exceptions when too many operations are done in a specific time frame. Unfortunately this limit also applies between Amazon services. In this specific case, the AWS Glue database service is denying requests from the AWS Glue crawler service, which causes the crawling process to abort. Because this is an inherent Amazon limitation, Collibra cannot fix this problem. A possible work-around is to use more S3 File System assets with more restricted include paths.

# No assets created after synchronization job is completed

This is usually because AWS Glue didn't find any suitable files to process. A typical problem is a typo in the include path or exclude patterns. AWS Glue does not fail when an include path points to a directory that doesn't exist. Also, always verify there are no leading or trailing spaces in those fields.

# Only part of the expected files or file groups were ingested

Jobs in Collibra can only succeed or fail. It's possible that some of the crawlers are correctly defined while others contain errors, such as a typo in an include path or an unsupported AWS region. In that case, the activity is marked as successful, though part of it didn't succeed. Currently, the only way to confirm this is to read the log files of Collibra and the Jobserver.

> Note   When you start synchronization, the crawlers are created in AWS Glue. Once the crawlers are created, they are executed. If Collibra cannot create one or more crawlers, synchronization fails immediately. If the crawlers are created successfully, but fail later, synchronization only fails if all crawlers fail.

# File size (or other property) is not filled for file xxx.yyy

AWS Glue only provides the file size for known file types (called "classifiers" in the AWS Glue terminology). Files that are classified as Unknown are registered but won't have any property associated. For the list of built-in classifiers, see the AWS Glue documentation.

# A file is wrongly considered as a File Group

AWS Glue preferably considers a directory as a data set when possible. This leads to a File Group being created in Data Catalog. There are multiple cases where it considers (possibly wrongly) one or more files as a File Group. Unfortunately, those rules are not clearly defined in AWS Glue documentation. Collibra noticed that AWS Glue considers a directory as a data set in the following cases:

- A directory only contains one file that belongs to a known classifier (file type).
- All files contained in a directory (including sub-directories) expose a similar schema (for example, all CSV files with columns of text type)

The only work-around that Collibra found, is to experiment with include paths and exclude patterns of the crawlers. For example, if a crawler wrongly takes a directory with subdirectories as a single File Group, the official work-around is to add crawlers with the subdirectories as include paths. Unfortunately, this work-around requires a lot of manual work and is limited by the number of crawlers in AWS Glue (25 by default - can be expanded on request).

# My table name has a strange hash-code at the end

AWS Glue appends a hash code to differentiate two different files of the same name but different directories, for example, csv_boolean_csv_ fe8de80c6f9a2b31463801aa2778a427. This name, including the hash code, is actually transferred to  Data Catalog.

# Synchronizing an S3 File System fails with a **relationMaxLimitReachedTarget** message in logs

This error comes from a broken relation in the assets tree. An asset created by S3 ingestion gets more than one parent asset. For example, a File asset has more thanone parent directory or a Directory asset has more than one parent directory.

This typically happens when a user moves S3 assets to a different domain and then starts a synchronization. In that case, the ingestion jobs try to recreate the missing assets in the original domain while old relations are still present. This can lead to an inconsistency in the relation tree.

We strongly recommend that you never move assets created by S3 ingestion to another domain.

> Example
> You work in domain called Amazon, which contains a Directory asset called Main. The Main Directory asset has a child asset of the File type, called Names.
>
> You move the Main Directory asset to another domain called Local.
>
> When you synchronize again, Data Catalog first recreates the Main Directory asset in the Amazon domain and then it updates the Names File asset.
>
> As a consequence, the Names File has 2 parent directories, which is a relation cardinality error.

# Partial ingestion or update of assets

It is possible to store a very large number of files in S3 buckets, hence leading to a large number of assets, attributes and relations to ingest into Data Catalog. To optimize memory and speed, the ingestion process is not transactional as a whole. It works with small transactional batches. If ingestion fails and aborts after some batches are already executed, it is possible that the ingested data is incomplete (if it is the first synchronization) or only partly updated (if it is not the first synchronization). In this case, it's advised to fix the problem and re-synchronise as soon as possible.

# Synchronization fails when a directory contains a file and a directory with the same name (known issue)

In Amazon S3, you can use periods (.) in the name of a directory. As a consequence, you can give the directory a name that is identical to a file name, for example, Collibra.txt. However, if this happens, ingestion fails. This is a known issue.

# JSON ingestion shows partial value in technical data type attributes (known issue)

For security reasons, all values that contain information between < and > characters are automatically trimmed by Collibra. However, if JSON is ingested by AWS Glue, the technical data type attribute contains those characters to represent the JSON structure. As a consequence, the value is trimmed and thus invalid. In future releases of Collibra, several attribute types will be changed to the plain text kind to avoid this issue.

# Error message **The AWS Access Key Id you provided does not exist in our records** though credentials are accepted

A user may be able to store S3 credentials in the S3 File System asset, though he cannot synchronize Amazon S3, create, edit or delete crawlers. The following message appears:

```
The AWS Access Key Id you provided does not exist in our
records.
(Service: Amazon S3; Status Code: 403; Error Code: Inval-
idAccessKeyId; ...
```

This may be caused by insufficient permissions on AWS Glue services. For more information, see About the Amazon S3 file system integration.

# Some of the folders and files in Amazon S3 are not visible in Collibra

You may notice that the content of your Amazon S3 does not always match the content in Collibra. Some folders from Amazon S3 may not appear in Collibra and some files are merged or split into different assets. This is not a bug in Collibra. When you synchronize Amazon S3, you create and execute crawlers in AWS Glue. Those crawlers create a table with metadata. That table is ingested in Collibra and is the basis for the relevant assets.

However, the crawlers in AWS Glue have some specific behavior to deal with partitioned tables. When the majority of schemas at a folder level are similar, the AWS Glue crawler creates partitions of a table instead of separate tables. Based on that information, the assets in Collibra are created.

See the AWS Glue documentation for more information about folders and tables in Amazon S3 and what happens when a crawler runs.

# Synchronizing Amazon S3 fails because you don't have the necessary permissions

In Collibra Data Intelligence Cloud 2020.11 and newer and Collibra Data Governance Center 5.7.7 and newer, Collibra checks the permissions of the AWS user when you synchronize Amazon S3. Synchronizing Amazon S3 fails if the AWS user does not have the necessary permissions.

A dialog box shows the following message:

```
Could not get/delete Glue database for S3 File System <name-of-
Amazon-S3-file-system>, please make sure you have all the neces-
sary permissions.
```

You must grant the AWS programmatic user the following permissions to synchronize Amazon S3 :

```
{
    "Version": "2012-10-17",
    "Statement": [
        {
            "Sid": "VisualEditor0",
            "Effect": "Allow",
            "Action": [
                "glue:GetCrawler",
                "glue:GetCrawlers",
                "glue:DeleteDatabase",
                "glue:GetTables",
                "glue:DeleteCrawler",
                "glue:StopCrawler",
                "s3:ListBucket",
```

```
                "glue:GetDatabases",
                "glue:CreateCrawler",
                "glue:GetDatabase",
                "iam:PassRole",
                "glue:StartCrawler",
                "glue:BatchDeleteTable",
                "s3:GetBucketLocation"
            ],
            "Resource": "*"
        }
    ]
}
```
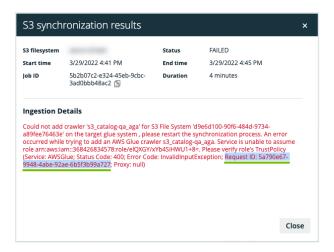
For more information about AWS requirements, see the Amazon S3 file system section.

# Glue Crawler fails with an **Internal Service Exception** error message

This is an AWS Glue crawler error. For possible steps to resolve the issue, see the AWS documentation.

# Where do I find the **Request ID** for AWS troubleshooting?

When an S3 synchronization fails, you can find a detailed error message, including the **Request ID**, in the S3 synchronization results.

Tip  Share the **Request ID** with AWS support to understand why the specific request is failing in AWS. This is typically useful to troubleshoot IAM permission issues in your AWS environment.

# Working with Tableau

Tableau is business intelligence software that helps people see and understand their data. Integrating Tableau in Collibra Data Intelligence Cloud enables you to see metadata from Tableau Server and Tableau Online in CollibraData Catalog.

In this section, we describe how you can ingest Tableau metadata in Collibra Data Catalog and synchronize the metadata via the Data Catalog user interface.

> Note   We have made available a new Tableau integration method that entails use of the lineage harvester, a standalone Java application. The following is true of the new integration method:
>
> - It is a cloud-only feature.
> - It requires a license for Collibra Data Lineage.
> - The new Tableau operating model is only available in Collibra versions 2021.10 and newer.

> Important   The two Tableau integration methods—Tableau integration via the Data Catalog UI and the new integration method via lineage harvester—coexist, and you are free to use the method of your choosing.

The following table shows the features specific to the two integration methods.

| Feature | Integration via Data Catalog UI | Integration via the lineage harvester |
|---|---|---|
| Catalog ingestion | ✓ | ✓ |
| Technical lineage |  | ✓ |
| Automatic stitching |  | ✓ |

| Feature | Integration via Data Catalog UI | Integration via the lineage harvester |
|---|---|---|
| Embedded data source connectivity | | ✓ |
| Custom SQL parsing | | ✓ |
| On-prem credential storage | | ✓ |
| Ingestion via Explorer role (with the Data Management Add-On) | | ✓ |

# About the Tableau integration

Tableau integration means registering Tableau as a system in Collibra Data Intelligence Cloud and ingesting, or synchronizing, the Tableau metadata. After synchronization, metadata from Tableau Server or Tableau Online are represented in Collibra by specific asset and domain types, retaining their original names.

# Tableau ingestion

The table below shows the steps required for ingesting Tableau metadata.

| Step | What? | Description |
|------|-------|-------------|
| 1 | Register Tableau as a system. | Creates an initial structure of a community, BI Catalog domain and Tableau Server asset in the selected parent community. |
| 2 | Connect to Tableau Server or Tableau Online. | Connects to Tableau server or Tableau Online. |
| 3 | Synchronize Tableau Server or Tableau Online. | Ingests the metadata from Tableau. |
| 4 | Stitch Tableau logical data layer and physical data layer. | Optionally, stitch Tableau assets to assets of registered data sources in Data Catalog. |

# Authentication

Data Catalog uses Tableau's REST API to get metadata information and follows Tableau's requirements regarding authentication methods. As a consequence, you need a Tableau user with access to the relevant Tableau sites.

For more information, see the Tableau documentation.

# Password encryption

Collibra's integration of Tableau does not use a separate encryption services, but reuses the Collibra DGC core service encryption method. This method uses the AES/CBC/PKCS5Padding transformation to encrypt your passwords when you connect to Tableau.

# Limitations

Collibra does not support the following Tableau features:

- Gziped encoding in REST results from Tableau.
- Tableau data sources that are created using Custom SQL.
- Tableau data sources that are created using Multiple tables union.

# Supported Tableau Server versions

Collibra Data Intelligence Cloud supports the following Tableau Server versions:

- 10.4
- 10.5
- 2018.x
- 2019.x
- 2020.1
- 2020.2
- 2020.3
- 2020.4
- 2021.1
- 2021.2
- 2021.3
- 2021.4

> Note   Depending on your Tableau version, Data Catalog uses different APIs to integrate Tableau. You need different Tableau permissions according to the Tableau version that you want to integrate.

# Tableau terminology

Before you start using Tableau to ingest data, read more about the Tableau terminology and how it maps with the Collibra Data Intelligence Cloud terminology.

| Tableau term | Description | Collibra equivalent |
|---|---|---|
| Site | A site is a stand-alone collection of content, such as projects, workbooks and users. Each site has its own URL and its own set of users. | Subcommunity and Tableau Site asset |
| Project | A project organizes related content resources. Content resources are workbooks, views and data sources. | Tableau Project asset |
| Workbook | A workbook is a collection of views. | Tableau Workbook asset |
| View | A view is a way to represent data. | Tableau View asset |
| Story | A story contains a sequence of work-sheets or dashboards that work together to convey information. | Tableau Story asset |
| Dashboard | A dashboard is a collection of views from multiple worksheets. | Tableau Dashboard asset |
| Worksheet | A worksheet contains a single view, along with shelves, legends, and the Data pane. | Tableau Worksheet asset |
| Tableau data sources | Tableau Data Sources consist of metadata that describe the connection information, information about how to access or refresh the data and customizations. | Tableau Data Source asset |
| Dimension | Dimensions contain qualitative values (such as names, dates, or geographical data). | Tableau Report Attribute asset |

| Tableau term | Description | Collibra equivalent |
|---|---|---|
| Measure | Measures contain numeric, quantitative values that you can measure. | Tableau Report Attribute asset |
| Tableau data attribute | Tableau Data Attributes define a property of a Tableau data entity. | Column asset |
| Tableau data entity | Tableau Data Entities are an abstraction of the physical implementation of database tables, used for Tableau report creation. | Schema asset and Table asset |
| Tableau data model | Tableau Data Models are an abstraction for the physical implementation of databases, schemas, files, etc., used for Tableau report creation. | Database asset |

# Tableau asset and domain types

The Tableau integration of Collibra Data Intelligence Cloud uses a specific subset of asset types and domain types. All of these come out of the box with your software.

The following table contains the asset an domain types that are used for the Tableau integration. Above each asset type you can see the parent asset types in the breadcrumbs.

| Asset type | Description | Domain type |
|---|---|---|
| Business Asset › Business Dimension › BI Folder › Tableau Project | Collection of Tableau workbooks and data sources. | BI Catalog |
| Business Asset › Business Dimension › BI Folder › Tableau Site | Collection of content (workbooks, data sources, users, …) that's walled off from any other content on that instance of Tableau Server. | BI Catalog |
| Business Asset › Report › BI Report › Tableau View › Tableau Dashboard | A collection of several worksheets and supporting information, shown on a single screen, so that you can simultaneously compare and monitor a variety of data. | BI Catalog |
| Business Asset › Report › BI Report › Tableau View › Tableau Worksheet | A worksheet is a single sheet on which you can build views of your data. | BI Catalog |
| Business Asset › Report › BI Report › Tableau Workbook | Collection of sheets. A sheet can be a worksheet, a dashboard or a story. | BI Catalog |

| Asset type | Description | Domain type |
|---|---|---|
| Data Asset ▸ Data Element ▸ Data Attribute ▸ BI Data Attribute ▸ Tableau Data Attribute | A specification that defines a property of a Tableau data entity.<br><br>Examples: CustomerBirthDate, EmployeeFirstName. | BI Catalog |
| Data Asset ▸ Data Element ▸ Report Attribute ▸ BI Report Attribute ▸ Tableau Report Attribute | An atomic unit of data that represents a Tableau report.<br><br>Examples: ExpenseAmount, RiskAmount | BI Catalog |
| Data Asset ▸ Data Structure ▸ Data Entity ▸ BI Data Entity ▸ Tableau Data Entity | An abstraction from the physical implementation of database tables, used for Tableau report creation. | BI Catalog |
| Data Asset ▸ Data Structure ▸ Data Model ▸ BI Data Model ▸ Tableau Data Model | An abstraction from the physical implementation of database, schema, file, etc., used for Tableau report creation. | BI Catalog |
| Technology Asset ▸ Server ▸ BI Server ▸ Tableau Server | A visual analytics platform for creating interactive dashboards and rich visualisations | BI Catalog |

| Asset type | Description | Domain type |
|---|---|---|
| Technology Asset ▸ System ▸ BI Data Source ▸ Tableau Data Source | The link between Tableau and an external system. A Tableau data source contains the information to connect to external data, table names, the table relationships, and any customizations that you make. | BI Catalog |

Note   The BI Data Catalog domain type was formerly known as the Tableau Data Catalog domain type.

# Tableau business logic

Tableau business users work with Tableau projects, workbooks and worksheets to make business decisions. Collibra's Tableau integration offers business users several advantages:

- Easily find certified Tableau content.
- Shop for Tableau reports.
- Trace Tableau data to its source.
- Find where content is stored in Tableau.

## Tableau asset pages

Tableau metadata is represented by assets of various types. Depending on the Tableau asset type, the asset page shows different information ingested from Tableau. You can find a specific Tableau asset page using Data Catalog search or via the Data Catalog BI domains in which you ingested the Tableau metadata.

## Details

Asset pages show attributes and relations to other assets. This information is synchronized from Tableau. However, you can add additional characteristics, tags or comments.

If you want access to one or more Tableau assets, you can add them to your Data Basket and check out the Data Basket. You can request access to assets of the following types:

- Tableau Workbook
- Tableau Worksheet
- Tableau Dashboard
- Tableau Story

Example The following Tableau Worksheet asset shows in which Tableau Dashboard and Tableau Story it is used and which Tableau Report Attribute it uses. This asset and the related reports are certified, indicating that the data is considered reliable.



# Diagrams

Diagrams is a feature to show and interact with assets based on their relations in an easy-to-read diagram. The diagram helps you to quickly see how assets are related. As such, the diagram can show a high-level presentation of a Tableau Workbook. If the Tableau

assets are stitched to registered assets in Data Catalog, you can also see the stitching results in the diagram. This enables you, for example, to see:

- In which Tableau Project the Tableau Workbook is stored.
- In which Tableau Site the Tableau Project is stored.
- Which Tableau Data Source is the source of the Tableau Report Attributes in the Tableau Workbook.
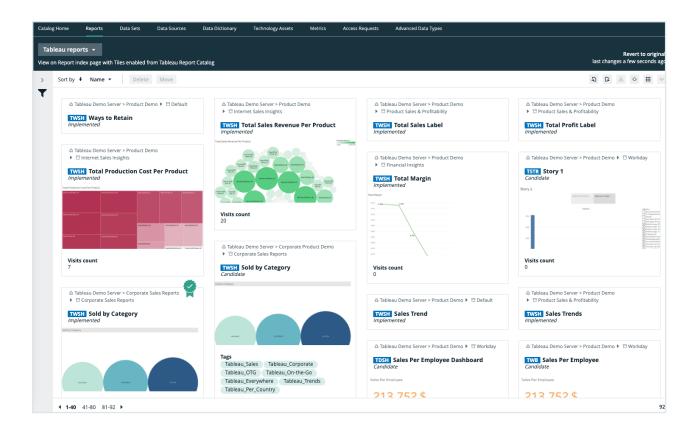- Which Table assets are the source for the Tableau Data Source asset via stitching.

Example   The following diagram shows the *Customer Sales Insights* Tableau Workbook, which is stored in the *Internet Sales Insights* Tableau Project. The Tableau Workbook contains Tableau Report Attributes that have the *CustomerSalesReporting* Tableau Data Source as source. This Tableau Data Source is stitched to the *CustomerSalesReporting* Table asset in the *SQL Server Cloud* data source.



# Report views

The Tableau integration feature enables you to find all ingested Tableau Workbook assets and children of this asset type in a single location.

In the **Reports** tab page in Data Catalog you can see an overview of all BI Report assets and their children. Optionally, you can create a view with a filter to only show Tableau assets. This is useful if you quickly want to find a specific report or if you want to know which reports are certified.

# Register a Tableau server

Before you can synchronize Tableau, you have to register a Tableau server to create an initial structure of a community, meaning a BI Catalog domain and a Tableau Server asset in a selected parent community in Collibra Data Intelligence Cloud.

# Prerequisites

- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a global role that has the Manage all resources global permission.
- You have a role with the following resource permissions on the Tableau community you create when you register a Tableau server:
  - Asset: add
  - Attribute: add

- ◦ Domain: add
- ◦ Attachment: add
- • You have enabled the Tableau metadata API in Collibra Console and in Tableau if you use Tableau 2020.2 or newer.

# Steps

1. On the main menu, click ⠿, then ⬓ **Catalog**.

   » The Catalog Home opens.
2. In the main menu, click the **Create** (＋) button.
3. In the **Create** dialog box, click **Register system**.
4. In the **Register system** dialog box, click **Tableau Server**.
5. In the **Register Tableau server** dialog box, enter the required information.

| Field | Description |
|---|---|
| Community | The name of the parent community in which the initial Tableau structure will be created. |
| Tableau server name | The name of the Tableau server. The name that you fill in here will be the name of the subcommunity, the domain in this subcommunity and the Tableau Server asset. |
| Description | A description to provide extra information about the Tableau server. This content is used as the description of the Tableau Server asset. |
| Owner | The owner of the data in the created community. By default, your user is selected. |

6. Click **Register**.
   » A Tableau Server asset is created.
   » A Tableau Catalog domain is created.
   » The configuration page of the Tableau Server asset is automatically opened.

Chapter 10

# What's next?

You can now connect to Tableau Server or Tableau Online.

# Connect to Tableau

To retrieve data from Tableau, you have to connect to Tableau via a Tableau Server asset in your Collibra Data Intelligence Cloud environment.

> Tip   You have to register the Tableau Server asset before you can connect to it.

You can edit the connection settings at any time, for example, if you want to use another user than the one you originally used.

## Prerequisites

- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have registered Tableau.
- If you connect to Tableau Online, you have a Tableau user with at least Viewer rights.
- If you connect to Tableau Server, you have a Tableau user with access to at least one site.
- You have the necessary Tableau permissions.

## Steps

1. Open a Tableau Server asset page.
2. In the tab pane, click ⚙ **Configuration**.

3. In the **Connection details** section, click **Edit connection details**.
4. Enter the required information.

| Field | Description | Required |
|---|---|---|
| On-premises \| Online | The Tableau product that you use. | ✓ Yes |
| Tableau URL or endpoint | The URL of your Tableau Server or Tableau Online.<br><br>Example   http://my-tableau.collibra.com | ✓ Yes |
| Site ID | The ID of a Tableau site.<br><br>○ If you don't enter a site ID, your Tableau user must have access to the Default site.<br>○ If you enter a site ID, your Tableau user must have access to that site.<br><br>Note   If you connect to Tableau Server, the site ID does not determine which sites you can synchronize from that server. It is used to validate the permissions of the Tableau user. Eeven if you enter one site ID, you can still synchronize the other sites from Tableau Server.<br><br>Tip<br>You can find the site ID in the URL of the Tableau site. The site ID is the string between `/site/` and `/projects/`.<br><br>In the following URL, the site ID is `collibra`.<br><br>`https://example.collibra.online.tableau.com/#/site/collibra/projects` | ✗ No for Tableau Server<br><br>✓ Yes for Tableau Online |

| Field | Description | Requi-red |
|---|---|---|
| Token Name/Use-rname | For Tableau Online with multi-factor authentication, the Per-sonal Access Token (PAT) name of the Tableau user. Other-wise, the username of the Tableau user. | ✓ Yes |
| Token Secret/Pas-sword | For Tableau Online with multi-factor authentication, the Per-sonal Access Token secret of the Tableau user. Otherwise, the password of the Tableau user. | ✓ Yes |

5. Click **Save**.

   » The connection is verified. If successful, you can see the list of available sites in Tableau.

# What's next?

You can now synchronize one or more sites.

# About synchronizing Tableau

Synchronizing Tableau is the process of ingesting metadata from a selected Tableau Server or Tableau Online and making the data available in Collibra Data Intelligence Cloud.
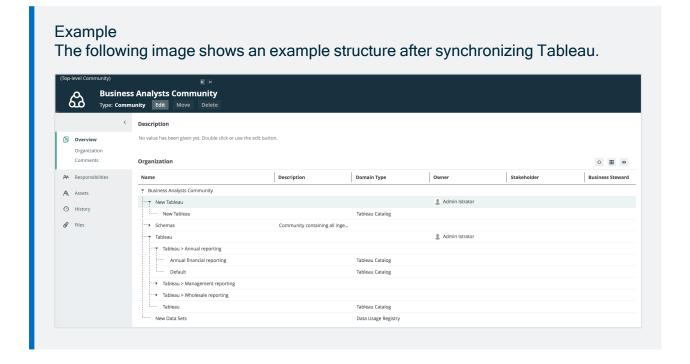
In this section, you can find the relevant actions and permissions to successfully synchronize Tableau.

For complete information on synchronizing Tableau, see the Collibra Data Intelligence Cloud User Guide.

# Synchronizing Tableau

Synchronizing Tableau is the process of ingesting metadata from a selected Tableau Server or Tableau Online and making the data available in Collibra Data Intelligence

Cloud.

Synchronization includes the following actions:

- For each Tableau site, a subcommunity is created in the community that was created during the registration of Tableau Server or Tableau Online.
- For each Tableau project, a Tableau Catalog domain is created in the community.
- In each Catalog BI domain, a Tableau Site asset is created, with the same name as the site.
- In each Catalog BI domain, the relevant assets are created, depending on the Tableau user's permissions.

> **Note**
> - Relations that were created between Tableau assets and other assets via a relation type in the Tableau operating model, are deleted after synchronization.
> - Currently, we only support published Tableau data sources with an extract or a live connection. For more information, see the Tableau documentation.

> **Example**
> The following image shows an example structure after synchronizing Tableau.
>
>

# Starting synchronization

You can synchronize manually, or you can automate the process by adding a synchronization schedule via a cron expression.

You can only synchronize one Tableau Server asset at a time. If a synchronization job is in progress and a second one is triggered (manually or automatically), it will be queued.

If a synchronization job is running and a new synchronization of the same Tableau Server asset is triggered (manually or automatically), the running synchronization continues and the new synchronization request is ignored.

> Note   If you have stitched Tableau's logical data layer to Data Catalog physical data layer, you have to restitch to make sure that all relations are up-to-date.

# Synchronization errors

In the following situations, nothing is synchronized and no subcommunities, domains or assets are created:

- If the job fails to start due to connection problems.
- If the job fails in the middle of the procedure.
- If the job is canceled.

For more information about Tableau synchronization issues, see the troubleshooting section.

> Warning   If you upgrade to Tableau version 2020.2 or newer, but previously synchronized an older Tableau version via the REST API and XML mapping, you have to prepare the migration procedure to prevent losing manually added relations, attributes, tags, comments and stitching results.

# Limitations and considerations

Collibra does not support the following Tableau features:

- Gziped encoding in REST results from Tableau.
- Tableau data sources that are created using Custom SQL.

- Tableau data sources that are created using Multiple tables union.

Collibra does support Tableau data sources that are created using:

- Cross-database joins
- Multiple tables join
- Relationships
- Single table

For more information, see the Tableau documentation.

## Naming convention

When you synchronize Tableau, Collibra follows a strict naming convention for the names of the new assets. Each asset has a display name and full name. The full name represents the asset path from asset to the database in which it is located. You can freely edit the display name. However, you should never edit the full name, because Data Catalog may need it to synchronize and stitch data sources. This can cause unexpected results and break the synchronization process.

> Warning   Editing full name of the Tableau Server or Tableau Online assets may lead to errors during the synchronization process.
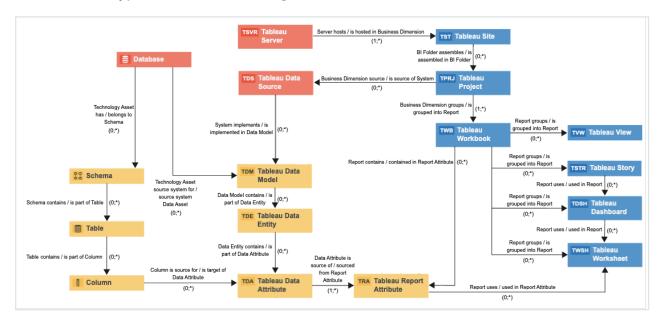
# Synchronized Tableau data

Synchronizing Tableau data means ingesting metadata from Tableau to your Collibra Data Intelligence Cloud environment. The metadata is represented as assets of specific types and their characteristics.

> **Note**
> - The assets have the same names as their counterparts in Tableau.
> - Some asset types are only created if the Tableau user specified in the connection settings has specific permissions.
> - There might be differences between the hierarchy of assets in Data Catalog and in Tableau. For example, Tableau, shows the relation from a parent project to a child project. In Data Catalog, this relation does not exist. Instead, all projects are shown on the Tableau Site asset page and the hierarchy of projects is shown in the Full name of the Tableau Project asset and the name of its domain.
> - If the Tableau data has tags, they are also added to the corresponding assets in Collibra with the prefix *Tableau_*.
> - Relations that were created between Tableau assets and other assets via a relation type in the Tableau operating model, are deleted after synchronization.

## Tableau operating model

The following image shows the relations between Tableau asset types and the cardinality of the relation types in the assets' assignment.



## Synchronized metadata per asset type

This table shows the metadata for each Tableau asset type.

| Asset type | Synchronized metadata |
|---|---|
| Tableau Server | • URL<br>• Server hosts / is hosted in Business Dimension |
| Tableau Site | • URL: The link to the data in Tableau<br>• Original name: The name of the data as used in Tableau<br>• BI Folder assembles / Is assembled in BI Folder<br>• Server hosts / is hosted in Business Dimension |
| Tableau Project | • Description<br>• Original name: The name of the project in Tableau<br>• Business Dimension groups / grouped into Report<br>• Business Dimension source is / source of System<br>• BI Folder assembles / is assembled in BI Folder<br>• Business Asset groups / is grouped by Business Asset |

| Asset type | Synchronized metadata |
|------------|----------------------|
| Tableau Work-book | <ul><li>Certified</li><li>Original name: The name of the workbook in Tableau.</li><li>Report image: The image of the report.</li></ul>Note<ul><li>○ Images are not downloaded or stored in Data Catalog. Instead, Data Catalog stores a link to the image. Every time you open the asset page, the image is fetched from Tableau. If the images do not render correctly, see the Troubleshooting section.</li><li>○ You can also exclude images from synchronization in the **Tableau sites** section on the **Configuration** page of the Tableau Server asset.</li></ul><ul><li>Document size</li><li>Document creation date</li><li>Document modification date</li><li>Report groups / is grouped into Report</li><li>Report is grouped in / groups Business Dimension</li><li>Report related to / is impacted by Business Asset</li><li>Report Attribute contained in / contains in Report</li><li>Technology Asset is source for / sourced from Business Asset</li></ul> |

| Asset type | Synchronized metadata |
|---|---|
| Tableau View | <ul><li>URL: The link to the data in Tableau</li><li>Certified</li><li>Original name: The name of the view in Tableau.</li><li>Report image: The image of the report</li></ul><br>Note<br>  ○ Images are not downloaded or stored in Data Catalog. Instead, Data Catalog stores a link to the image. Every time you open the asset page, the image is fetched from Tableau. If the images do not render correctly, see the Troubleshooting section.<br>  ○ You can also exclude images from synchronization in the **Tableau sites** section on the **Configuration** page of the Tableau Server asset.<br><br><ul><li>Visits count: The number of times that the view has been visited in Tableau</li><li>Document creation date</li><li>Document modification date</li><li>Visible on server</li><li>Tags</li><li>Report groups /is grouped into Report</li><li>Report relates / is impacted by Business Asset</li></ul><br>Note  Assets of this type are only created if the Tableau user does not have the Download/Save As permission on the workbook. |

| Asset type | Synchronized metadata |
|---|---|
| Tableau Story | <ul><li>URL: The link to the data in Tableau</li><li>Certified</li><li>Original name: The name of story in Tableau.</li><li>Report image: The image of the report.</li></ul><blockquote>Note<ul><li>○ Images are not downloaded or stored in Data Catalog. Instead, Data Catalog stores a link to the image. Every time you open the asset page, the image is fetched from Tableau. If the images do not render correctly, see the Troubleshooting section.</li><li>○ You can also exclude images from synchronization in the **Tableau sites** section on the **Configuration** page of the Tableau Server asset.</li></ul></blockquote><ul><li>Visits count: The number of times that the view has been visited in Tableau.</li><li>Document creation date</li><li>Document modification date</li><li>Visible on server</li><li>Tags</li><li>Report groups /is grouped into Report</li><li>Report related to / is impacted by Business Asset</li><li>Report uses / used in Report</li></ul><blockquote>Note   Assets of this type are only created if the Tableau user has the Download/Save As permission on the workbook.</blockquote> |

| Asset type | Synchronized metadata |
|---|---|
| Tableau Dashboard | <ul><li>URL: The link to the data in Tableau</li><li>Certified</li><li>Original name: The name of story in Tableau.</li><li>Report image: The image of the report.</li></ul> |

> **Note**
> ◦ Images are not downloaded or stored in Data Catalog. Instead, Data Catalog stores a link to the image. Every time you open the asset page, the image is fetched from Tableau. If the images do not render correctly, see the Troubleshooting section.
> ◦ You can also exclude images from synchronization in the **Tableau sites** section on the **Configuration** page of the Tableau Server asset.

- Visits count: The number of times that the view has been visited in Tableau.
- Document creation date
- Document modification date
- Visible on server
- Tags
- Report groups /is grouped into Report
- Report related to / is impacted by Business Asset
- Report uses / used in Report

> Note   Assets of this type are only created if the Tableau user has the Download/Save As permission on the workbook.

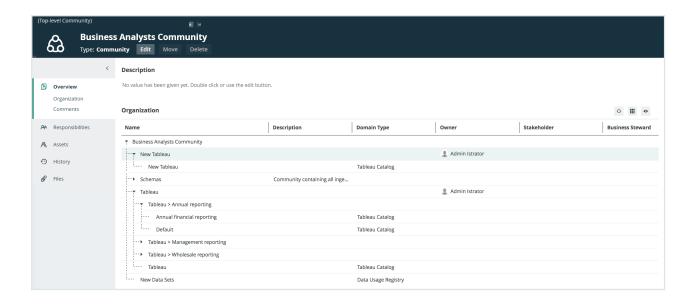| Asset type | Synchronized metadata |
|---|---|
| Tableau Worksheet | <ul><li>URL: The link to the data in Tableau</li><li>Certified</li><li>Original name: The name of the data as used in Tableau</li><li>Report image: The image of the report.</li></ul><br>Note<br>○ Images are not downloaded or stored in Data Catalog. Instead, Data Catalog stores a link to the image. Every time you open the asset page, the image is fetched from Tableau. If the images do not render correctly, see the Troubleshooting section.<br>○ You can also exclude images from synchronization in the **Tableau sites** section on the **Configuration** page of the Tableau Server asset.<br><br><ul><li>Visits count: The number of times that the view has been visited in Tableau.</li><li>Document creation date</li><li>Document modification date</li><li>Visible on server</li><li>Tags</li><li>Report uses / used in Report Attribute</li><li>Report groups / is grouped into Report</li><li>Report related to / impacted by Business Asset</li></ul><br>Note   Assets of this type are only created if the Tableau user has the Download/Save As permission on the workbook. |

| Asset type | Synchronized metadata |
|---|---|
| Tableau Report Attribute | • Description<br>• Original Name: The name of the attribute as used in Tableau<br>• Technical Data Type<br>• Role in Report<br>• Calculation Rule: Formula used in measure<br>• Report Attribute contained in / contains in Report<br>• Report Attribute is source for / is target of Report Attribute<br>• Report Attribute sourced from / is source of Data Attribute<br>• Report uses / used in Report Attribute<br><br>Note<br>• Assets of this type are only created if the Tableau user has the Download/Save As permission on the workbook.<br>• These are only the report attributes that are used in Tableau Worksheet of the Tableau Workbook. |
| Tableau Data Attribute | • Original Name: The name of the attribute as used in Tableau<br>• Technical Data Type: The Data Type of a data asset as it is declared by the data source.<br>• Report Attribute sourced from / is source of Data Attribute<br><br>Note   Assets of this type are only created if the Tableau user has the Download/Save As permission on the data source. |
| Tableau Data Entity | • Data Entity contains / is part of Data Attribute<br>• Data Entity is part of / contains Data Model<br><br>Note   Assets of this type are only created if the Tableau user has the Download/Save As permission on the data source. |

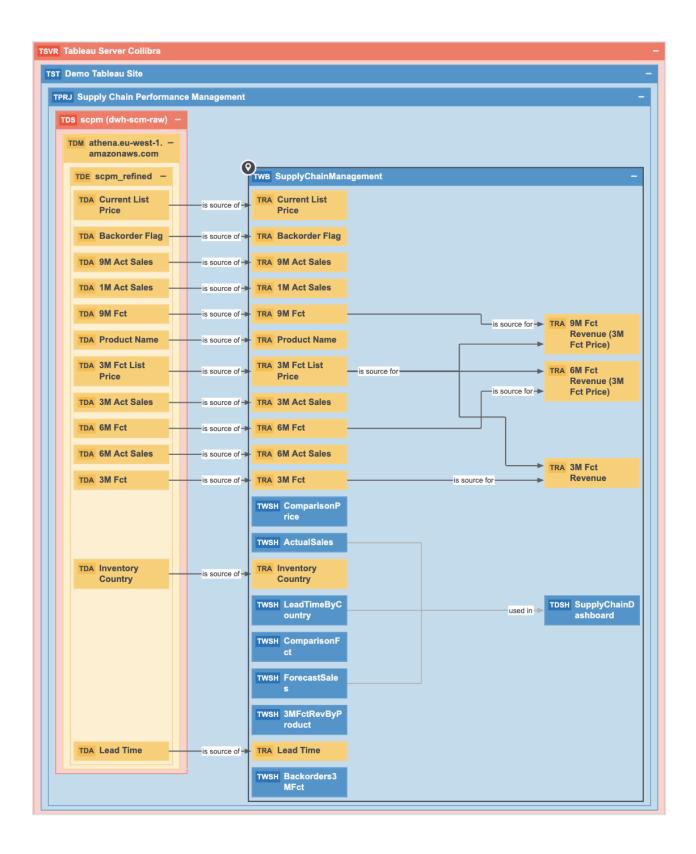| Asset type | Synchronized metadata |
|---|---|
| Tableau Data Model | <ul><li>Data Source Type</li><li>Location</li><li>Data Entity is part of / contains Data Model</li><li>System implements / is implemented in Data Model</li></ul><br>Note   Assets of this type are only created if the Tableau user has the Download/Save As permission on the data source. |
| Tableau Data Source (Published only) | <ul><li>Certified</li><li>Original name: The name of the data as used in Tableau</li><li>Document creation date</li><li>Document modification date</li><li>Business Dimension sources / is source of System</li><li>System implements / is implemented in Data Model</li><li>Technology Asset implements /is implemented in Data Asset</li></ul><br>Note   Currently, we only support published data sources with an extract or a live connection. For more information, see the Tableau documentation. |

# Examples of synchronized metadata

The following image shows an example structure after synchronizing Tableau.

The following image shows an example of a diagram of a Tableau server.

# Tableau permissions and ingestion results

When you synchronize Tableau, you need certain permissions to access the data in Tableau. The extent of your permissions dictates the scope of the ingestion results

The following table shows the minimum role and permissions requirements for successful synchronization and the scope of the ingestion results in Data Catalog.

| Tableau version | Tableau site role | Minimum required permissions | | | Result in Data Catalog |
| --- | --- | --- | --- | --- | --- |
| | | Project | Workbook | Data Source | |
| Older than 2020.2 | Viewer | View | View | View | Tableau Workbooks and Tableau Data Sources are not parsed.<br><br>Resulting asset types:<br><br>• Tableau Server<br>• Tableau Site<br>• Tableau Project<br>• Tableau Data Source<br>• Tableau Workbook<br>• Tableau View |

| Tableau version | Tableau site role | Minimum required permissions | | | Result in Data Catalog |
|---|---|---|---|---|---|
| | | Project | Workbook | Data Source | |
| Older than 2020.2 | Explorer<br><br>Note If your Tableau version is older than 2018.1, the Tableau site role is Interactor. | View | View | View, Download/Save As | Tableau Data Sources are parsed.<br><br>Resulting asset types:<br><br>• Tableau Server<br>• Tableau Site<br>• Tableau Project<br>• Tableau Workbook<br>• Tableau View<br>• Tableau Data Source<br>• Tableau Data Model<br>• Tableau Data Entity<br>• Tableau Data Attribute |

| Tableau version | Tableau site role | Minimum required permissions | | | Result in Data Catalog |
|---|---|---|---|---|---|
| | | Project | Workbook | Data Source | |
| Older than 2020.2 | Explorer<br><br>Note   If your Tableau version is older than 2018.1, the Tableau site role is Interactor. | View | View, Download/Save As | View | Tableau Report Attributes are synchronized and Tableau Workbooks are parsed.<br><br>Resulting asset types:<br><br>• Tableau Server<br>• Tableau Site<br>• Tableau Project<br>• Tableau Data Source<br>• Tableau Workbook<br>• Tableau Story<br>• Tableau Dashboard<br>• Tableau Worksheet |

| Tableau version | Tableau site role | Minimum required permissions | | | Result in Data Catalog |
|---|---|---|---|---|---|
| | | Project | Workbook | Data Source | |
| Older than 2020.2 | Explorer<br><br>Note If your Tableau version is older than 2018.1, the Tableau site role is Interactor. | View | View, Download/Save As | View, Download/Save As | Tableau Report Attributes are synchronized, and Tableau Data Sources and Tableau Workbooks are parsed.<br><br>Resulting asset types:<br><br>• Tableau Server<br>• Tableau Site<br>• Tableau Project<br>• Tableau Data Source<br>• Tableau Data Model<br>• Tableau Data Entity<br>• Tableau Data Attribute<br>• Tableau Workbook<br>• Tableau Story |

| Tableau version | Tableau site role | Minimum required permissions | | | Result in Data Catalog |
|---|---|---|---|---|---|
| | | Project | Workbook | Data Source | |
| | | | | | • Tableau Dashboard<br>• Tableau Worksheet |

| Tableau version | Tableau site role | Minimum required permissions | | | Result in Data Catalog |
|---|---|---|---|---|---|
| | | Project | Workbook | Data Source | |
| 2020.2 and newer | Viewer or Explorer | View | View | View | If you enabled the metadata API, Data Catalog creates new assets according to your content in Tableau without accessing metadata in Tableau databases and tables.<br><br>Resulting asset types:<br><br>• Tableau Server<br>• Tableau Site<br>• Tableau Project<br>• Tableau Data Source<br>• Tableau Workbook<br>• Tableau Story<br>• Tableau Dashboard |

| Tableau version | Tableau site role | Minimum required permissions | | | Result in Data Catalog |
|---|---|---|---|---|---|
| | | Project | Workbook | Data Source | |
| | | | | | • Tableau Worksheet<br><br>If you did not enable the Tableau metadata API, Tableau reports and data sources are ingested in Data Catalog, but with a limited scope.<br><br>Resulting asset types:<br><br>• Tableau Server<br>• Tableau Site<br>• Tableau Project<br>• Tableau Data Source<br>• Tableau Workbook<br>• Tableau View |

| Tableau version | Tableau site role | Minimum required permissions | | | Result in Data Catalog |
|---|---|---|---|---|---|
| | | Project | Workbook | Data Source | |
| 2020.2 and newer | Tableau Server Administrator or Site Administrator | View | View | View | If the metadata API is enabled, Data Catalog creates new assets according to your content in Tableau using metadata in Tableau databases and tables.<br><br>Resulting asset types:<br><br>• Tableau Server<br>• Tableau Site<br>• Tableau Project<br>• Tableau Data Source<br>• Tableau Report Attribute<br>• Tableau Data Model |

| Tableau version | Tableau site role | Minimum required permissions | | | Result in Data Catalog |
|---|---|---|---|---|---|
| | | Project | Workbook | Data Source | |
| | | | | | • Tableau Data Entity<br>• Tableau Data Attribute<br>• Tableau Workbook<br>• Tableau Story<br>• Tableau Dashboard<br>• Tableau Worksheet<br><br>If you did not enable the Tableau metadata API, Tableau reports and data sources are ingested in Data Catalog, but with a limited scope.<br><br>Resulting asset types:<br><br>• Tableau Server<br>• Tableau Site |

| Tableau version | Tableau site role | Minimum required permissions | | | Result in Data Catalog |
|---|---|---|---|---|---|
| | | Project | Workbook | Data Source | |
| | | | | | • Tableau Project<br>• Tableau Data Source<br>• Tableau Workbook<br>• Tableau View |

> **Warning** We do not support a full ingestion of Tableau Server or Tableau Online version 2020.2 or newer if the metadata API is disabled. If you try to synchronize a Tableau Server or Tableau Online asset after a Tableau upgrade to 2020.2 or newer without the metadata API, the synchronization result in Data Catalog will fail. This prevents data loss of manually added relations and attributes.

> **Tip** For more information about Tableau permissions, site roles and licenses, see the Tableau Online Help.

## Tableau data structure

You can only synchronize Tableau elements if the Tableau user specified in the connection settings has permissions to access them. If you have permissions to access a Tableau element, but not its parent elements, the parent elements are skipped when synchronizing Tableau and do not appear in Data Catalog.

This happens in the following situations:

- The Tableau user has permissions to access a Tableau workbook, but not its parent, the Tableau project.
- The Tableau user has permissions to access a Tableau view, but not its parent, the Tableau workbook.

- The Tableau user has permissions to access a Tableau view, but not its parent, the Tableau project.

# Metadata API

If you register a Tableau Server or Tableau Online version 2020.2 or newer, Data Catalog requires the metadata API to synchronize Tableau assets.

Tableau metadata consists of information about Tableau content and assets. Data Catalog creates GraphQL queries to collect metadata from Tableau Online or Tableau Server. If the metadata API is enabled in Tableau and in Collibra Console, Collibra Data Intelligence Cloud uses this metadata to create new assets in Data Catalog.

# Upgrading Tableau to 2020.2 or newer

If you have previously ingested and synchronized a version of Tableau older than 2020.2 and have since upgraded to version 2020.2 or newer, you have to enable the metadata API in Tableau and in Collibra Console. If you synchronize using the metadata API, Data Catalog removes all Tableau assets created via XML mapping and creates new ones using the metadata API. This means that all manually added relations, attributes, tags, comments and stitching results will be lost.

> Tip   We highly recommend to contact your Collibra Customer Success Manager before you synchronize a Tableau Server or Tableau Online asset after upgrading to Tableau version 2020.2 or newer.

# Parsing Tableau metadata

Parsing Tableau metadata is an automated procedure that allows the metadata to be captured and identified in Data Catalog at a more granular level. Typically, the result is that you have more assets of different types in Data Catalog, which leads to more complete information and better lineage diagrams.

Parsing takes place automatically during Tableau synchronization, depending on the Tableau permissions of the Tableau user who launched the synchronization process.

## Parsing Tableau workbooks

Without parsing, Tableau Workbooks contain Tableau Views, without further details. However, if your Tableau user has the Download/Save As permission for the Workbook, the Tableau workbook is parsed. As a consequence, there is no Tableau View asset, but there is at least one Worksheet asset, and, if they exist on Tableau: Tableau Story assets and Tableau Dashboard assets.

| Without Parsing | With Parsing |
|---|---|
| • Tableau Workbook<br>• Tableau View | • Tableau Workbook<br>• Tableau Story<br>• Tableau Dashboard<br>• Tableau Worksheet |

## Parsing Tableau Data Source

Without parsing, Tableau Data Sources do not contain further information about the data source. However, if your Tableau user has the Download/Save As permission for the Data Source , the Tableau Data Source is parsed. As a consequence, there is at least one Tableau Data Model asset and one or more Tableau Data Entity assets and Tableau Data Attribute assets. These assets are required for Tableau stitching.

| Without parsing | With parsing |
|---|---|
| • Tableau Data Source | • Tableau Data Source<br>  ○ Tableau Data Model<br>  ○ Tableau Data Entity<br>  ○ Tableau Data Attribute |

# Working with Tableau APIs

When you register or synchronize a Tableau Server, Data Catalog uses the Tableau APIs to ingest the Tableau metadata. Data Catalog uses different APIs depending on your version of Tableau. This happens automatically and should have little impact on the

resulting assets. However, if you synchronize Tableau 2020.2 or newer, you must perform a few extra actions.

## Tableau versions

The following table shows which APIs Data Catalog uses to register or synchronize a Tableau Server.

| Tableau versions using the REST API and XML parsing | Tableau versions using the REST API in combination with the GraphQL metadata API |
|---|---|
| • 10.4<br>• 10.5<br>• 2018.x<br>• 2019.x<br>• 2020.1 | • 2020.2<br>• 2020. 3<br>• 2020.4 |

> Warning   If you upgrade to Tableau version 2020.2 or newer, but previously synchronized an older Tableau version via the REST API and XML mapping, you have to prepare the migration procedure to prevent losing manually added relations, attributes, tags, comments and stitching results.

## Differences between the metadata API and XML parsing via REST API

The following table shows the differences and similarities between the metadata API and the REST API with XML parsing.

| Part of synchronization process | REST API and XML parsing | Metadata API |
|---|---|---|
| API | Data Catalog connects to Tableau via the REST API and uses custom parsing mechanisms. The result is XML data. | Data Catalog connects to Tableau via the REST API and the metadata API. The result is GraphQL data.<br><br>Note   We highly recommend that you synchronize Tableau after working hours. This is necessary to make sure that no Tableau data is added, changed, renamed or deleted on Tableau's side during the synchronization process. If there are any inconsistencies between the Tableau data collected via the REST API and Tableau data collected via the GraphQL metadata API, the corresponding Tableau assets are not synchronized in Data Catalog. |

| Part of synchronization process | REST API and XML parsing | Metadata API |
|---|---|---|
| Settings | You don't need change any settings to start Tableau synchronization. | You have to enable the Tableau metadata API in Collibra Console before you can ingest or synchronize.<br><br>Note   Also make sure that the Tableau metadata API is enabled in Tableau. |
| Relevant asset types | The resulting Tableau assets that are created after registering or synchronizing a Tableau Server are similar and mainly depend on the permissions of your Tableau user. ||
| Performance | Performance results are similar. ||
| Collibra Data Intelligence Cloudpermissions | The required permissions are the same: a resource role with the Configure external system resource permission. ||
| Stitching | Stitching works the same. ||

## Migration procedure

When you synchronize a Tableau Server for the first time after you upgraded to Tableau 2020.2 or newer, Data Catalog tries to match your Tableau assets that were previously ingested via the Tableau REST API in Data Catalog to their counterparts in Tableau. If the asset names match, Data Catalog changes the full name of the Tableau assets without removing manually added data and stitching results.

To make sure Collibra Data Intelligence Cloud is able to match your Tableau assets in Data Catalog to their counterparts in Tableau, you must prepare the migration procedure.

> Tip   We highly recommend that you create a backup of your Collibra environment before synchronizing a Tableau Server asset after you upgraded to Tableau 2020.2 or newer. We also recommend that you synchronize the first time after working hours.

# Prepare migration after upgrading to Tableau 2020.2 or newer

If you upgraded to Tableau version 2020.2 or newer, but previously synchronized an older version via XML mapping, Data Catalog changes the full names of your Tableau assets to match them to their counterparts in Tableau. This is necessary to prevent losing manually added relations, attributes, tags, comments and stitching results.

You only have to follow these steps once after your upgrade to Tableau 2020.2 or newer. After that, you can follow the default synchronization process.

> Note   Collibra Data Intelligence Cloud can only migrate your assets if:
>
> - All Tableau Report Attribute assets have the same name as their counterparts in Tableau.
> - Each Tableau Report Attribute asset name is unique within the same Tableau workbook.

> Tip   If you never manually changed the name of the assets in Data Catalog, they should automatically be the same as their counterparts in Tableau.

## Prerequisites

- You have registered Tableau.
- You have connected a Tableau Server asset to a Tableau Server or Tableau Online.
- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a resource role with the Asset > Update resource permission.
- You have a global role with the Catalog global permission, for example Catalog Author.

- You have previously ingested Tableau 2020.1 or older and have since upgraded to Tableau 2020.2 or newer.
- Your Tableau user has the right permissions to synchronize Tableau 2020.2 or newer.
- You have enabled the Tableau metadata API in Tableau.

## Steps

1. Match the names of all Tableau Report Attributes assets of a Tableau Workbook with their counterparts in Tableau.
   a. Open a Tableau Report Attribute asset page.
   b. In the resource toolbar, click **Edit**.
      » The **Edit <asset name>** dialog box appears.
   c. Change the name of the asset to the exact name used in Tableau.
   d. Click **Save**.

   > Tip   We highly recommend that you also match the display names of Tableau Data Attribute assets, Tableau Data Entity assets and Tableau Data Model assets. While Data Catalog automatically tries to match these assets to their counterparts in Tableau based on the Tableau Report Attribute asset, making sure the Tableau assets have the same name helps to prevent issues. Unless you manually changed their names in Data Catalog, the names should already be the same as their counterparts in Tableau.

2. Optionally, create a backup of your Collibra environment.

   > Note   We highly recommend that you create a backup before you synchronize a Tableau Server to prevent losing data in Data Catalog if something goes wrong during the migration process.

3. Enable the Tableau metadata API in Collibra Console.
4. Synchronize a Tableau Server asset after working hours.

   > Note   We highly recommend that you synchronize the first time after upgrading to Tableau 2020.2 or newer after working hours. This is necessary to make sure that no Tableau data is added, changed, renamed or deleted on Tableau's side during the synchronization process.

a.  Open a Tableau Server asset page.

b.  In the tab pane, click ⚙ **Configuration**.

c.  In the **Tableau sites** section, do the following:

    i.  Select one or more sites.

   ii.  Enable or disable report images as required for each site.

> **Note** Images are never downloaded or stored in Data Catalog. Depending on the Report image setting, Data Catalog either ignores images completely or stores a link to the image on Tableau and loads that image when you open the relevant asset page.

d.  In the **Tableau sites** section, click **Synchronize now**.

» The synchronization job appears in the **Activities** list as a bulk synchronization.

» The full names of the Tableau assets are updated to include the GraphQL ID.

» The log files show a summary of the migration process.

**Example**

```
"Summary of tableau xml to graphql data migration for
site with id <Tableau-site-ID> and name <Tableau-site-
name> executed on server"
```

The log files also show how many Tableau assets were found in Data Catalog and how many were migrated to match their counterparts in Tableau.

**Example**

```
"Found 50 existing xml assets to migrate.";
"Migrated 48 assets.";
```

If some Tableau assets could not be migrated, Collibra Data Intelligence Cloud recreates the Tableau asset so that it matches in Tableau. The log file shows how many and which assets were recreated in Data Catalog.

Example

```
"It was impossible to migrate 2 assets. These assets
were re-created based on graphql data.";
    List of assets that were not migrated:
        ID: xxxxxxxx-xxxx-xxxx-xxx, Fullname: Tableau-
migration > tableaumigration.xxxxx > [tableau-
migration-asset-name-1] (Tableau Report Attribute).",
        ID: xxxxxxxx-xxxx-xxxx-xxx, Fullname: Tableau-
migration > tableaumigration.xxxxx > [tableau-
migration-asset-name-2] (Tableau Report Attribute)."
```

# Synchronize Tableau site manually

You can manually start a synchronization job of a Tableau Server asset. This can be useful if you don't want to wait for the scheduled job to synchronize your Tableau sites.

Warning   You can choose which sites to synchronize after successfully connecting to Tableau. Select the same or more sites when you synchronize again. If you only synchronize some of the Tableau sites, Data Catalog deletes all other Tableau sites and their content from Collibra Data Intelligence Cloud.

Tip   You can also add a synchronization schedule to synchronize automatically.

## Prerequisites

- You have registered Tableau.
- You have connected a Tableau Server asset to a Tableau Server or Tableau Online.
- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a role with the following resource permissions on the Tableau community you create when you register a Tableau server:

- ○ Asset: add
- ○ Attribute: add
- ○ Domain: add
- ○ Attachment: add

- If you want to stitch Tableau's logical data layer to Data Catalog's physical data layer, the Tableau user must have the Download/Save As permission on the data source.
- You have enabled the Tableau metadata API in Collibra Console and in Tableau if you use Tableau 2020.2 or newer.

> Warning   If you upgrade to Tableau version 2020.2 or newer, but previously synchronized an older Tableau version via the REST API and XML mapping, you have to prepare the migration procedure to prevent losing manually added relations, attributes, tags, comments and stitching results.

## Steps

1. Open a Tableau Server asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Tableau sites** section, do the following:
   a. Select one or more sites.
   b. Enable or disable report images as required for each site.

   > Note   Images are never downloaded or stored in Data Catalog. Depending on the Report image setting, Data Catalog either ignores images completely or stores a link to the image on Tableau and loads that image when you open the relevant asset page.

4. In the **Tableau sites** section, click **Synchronize now**.
   » The synchronization job appears in the **Activities** list as a bulk synchronization.

> Note   We highly recommend that you synchronize a Tableau Server version 2020.02 and newer after working hours. This is necessary to make sure that no Tableau data is added, changed, renamed or deleted on Tableau's side during the synchronization process. If there are any inconsistencies between the Tableau data collected via the REST API and Tableau data collected via the GraphQL metadata API, the corresponding Tableau assets are not synchronized in Data Catalog.

> Tip   If your Tableau synchronization fails, go to the troubleshooting section to find a solution.

## What's next?

When the synchronization finishes, the resulting assets, including their attributes and relations, are created, edited or deleted in the selected domain(s) and in the Data Sources page of Data Catalog.

If you have stitched Tableau's logical data layer to Data Catalog's physical data layer, you have to restitch to make sure that all relations are up to date.

# Add a Tableau synchronization schedule

To keep the content of Collibra Data Intelligence Cloud synchronized with your Tableau Server or Tableau online, you can synchronize manually or create a schedule to automatically do this with a fixed interval.

> Note
> - You can only create one synchronization schedule.
> - If you have stitched Tableau's logical data layer to Data Catalog physical data layer, you have to restitch after each synchronization to make sure that all relations are up to date.
> - We highly recommend that you synchronize a Tableau Server version 2020.02 and newer after working hours. This is necessary to make sure that no Tableau data is added, changed, renamed or deleted on Tableau's side during the synchronization process. If there are any inconsistencies between the Tableau data collected via the REST API and Tableau data collected via the GraphQL metadata API, the corresponding Tableau assets are not synchronized in Data Catalog.

# Prerequisites

- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a role with the following resource permissions on the Tableau community you create when you register a Tableau server:
  - Asset: add
  - Attribute: add
  - Domain: add
  - Attachment: add
- You have registered Tableau.
- You have connected a Tableau Server asset to a Tableau Server or Tableau Online.
- You have enabled the Tableau metadata API in Collibra Console and in Tableau if you use Tableau 2020.2 or newer.

> Warning   If you upgrade to Tableau version 2020.2 or newer, but previously synchronized an older Tableau version via the REST API and XML mapping, you have to prepare the migration procedure to prevent losing manually added relations, attributes, tags, comments and stitching results.

## Steps

1. Open a Tableau Server asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Synchronization schedule** section, click **Add Schedule**.
4. Enter the required information.

| Field | Description |
|---|---|
| Repeat | The interval when you want to synchronize the schemas automatically, for example daily, weekly or based on a Cron expression. |
| Cron | The Quartz Cron expression that determines when the synchronization takes place.<br><br>This field is only visible if you select `Cron expression` in the **Repeat** field. |
| Every | The day on which you want to synchronize the schemas, for example Sunday.<br><br>This field is only visible if you select `Weekly` in the **Repeat** field. |
| Every first | The day of the month on which you want to synchronize the schemas , for example Tuesday.<br><br>This field is only visible if you select `Monthly` in the **Repeat** field. |
| At | The time at which you want to synchronize the schemas automatically, for example 14:00.<br><br>This field is only visible if you select `Daily`, `Weekly` or `Monthly` in the **Repeat** field. |
| Time zone | The time zone for the schedule. |

5. Click **Save**.

> Tip   If your Tableau synchronization fails, go to the troubleshooting section to find a solution.

# Edit a Tableau synchronization schedule

You can edit the synchronization schedule of a Tableau Server asset. For example, you can do this if you think the synchronization job runs too often or not often enough.

## Prerequisites

- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have registered Tableau.
- You have connected a Tableau Server asset to a Tableau Server or Tableau Online.
- You have added a synchronization schedule.

## Steps

1. Open a Tableau Server asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Synchronization schedule** section, click **Edit Schedule**.

4. Enter the required information.

| Field | Description |
|---|---|
| Repeat | The interval when you want to synchronize the schemas automatically, for example daily, weekly or based on a Cron expression. |
| Cron | The Quartz Cron expression that determines when the synchronization takes place. This field is only visible if you select `Cron expression` in the **Repeat** field. |
| Every | The day on which you want to synchronize the schemas, for example Sunday. This field is only visible if you select `Weekly` in the **Repeat** field. |
| Every first | The day of the month on which you want to synchronize the schemas , for example Tuesday. This field is only visible if you select `Monthly` in the **Repeat** field. |
| At | The time at which you want to synchronize the schemas automatically, for example 14:00. This field is only visible if you select `Daily`, `Weekly` or `Monthly` in the **Repeat** field. |
| Time zone | The time zone for the schedule. |

5. Click **Save**.

# Remove a Tableau synchronization schedule

You can remove a synchronization schedule from a Tableau Server asset to stop automatically synchronizing Tableau.

## Prerequisites

- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have registered Tableau.
- You have connected a Tableau Server asset to a Tableau Server or Tableau Online.
- You have added a synchronization schedule.

## Steps

1. Open a Tableau Server asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Synchronization schedule** section, click **Remove Schedule**.

# Delete a Tableau site from Collibra

You can delete a Tableau site and all of its contents from the Tableau site synchronization. Collibra Data Intelligence Cloud then deletes the community related to the Tableau site, including the domains and assets that it contains.

> Note   The **Tableau sites section** on a Tableau Server asset page shows all sites that exist in Tableau. If you want to remove Tableau sites from this list, you must remove them in Tableau.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a resource role with the Configure external system resource permission, for example Owner.
- You have connected a Tableau Server asset to a Tableau Server or Tableau Online.
- You have registered Tableau.
- You have synchronized Tableau at least once.

# Steps

1. Open a Tableau Server asset page.
2. In the tab pane, click ⚙ **Configuration**.
3. In the **Tableau sites** section, clear the sites that you want to delete from Data Catalog.

   > Tip   Only select the Tableau sites that you would like to keep. If you want to delete all Tableau sites from Data Catalog, clear all checkboxes.

4. In the **Tableau sites** section, click **Synchronize now**.
   » The **Synchronize Tableau server** dialog box appears.

   

5. Click **Synchronize and delete**.
   » The synchronization job appears in the **Activities** list as a bulk synchronization.
   After the synchronization, the cleared sites are deleted.

# What's next?

If you deleted the wrong Tableau site or you want to reintroduce it, you can select that Tableau site and synchronize it again.

# Tableau stitching

Stitching is a process that creates relations between assets representing the same data source: the data source of a Tableau report and the Data Catalog database. This allows you to clearly represent the lineage from the data source to the Tableau reports where it is used. As a consequence, you can easily perform impact analyses. For example, you can quickly see which reports will be affected if you refresh a table of your database, or which reports will be impacted if you drop one column from the table.

# About Tableau stitching

Before you can perform stitching, you have to ingest a Tableau report –including its data source– and register that data source separately in Data Catalog. The same data is then represented by Tableau assets as well as by regular Data Catalog assets such as Schema, Table and Column assets. Tableau stitching is based on the matching of the full name of Tableau Data Attribute assets and Column assets of registered data sources in Data Catalog. Follow the steps in the table below to enable Collibra Data Intelligence Cloud to automatically create relations between Tableau assets and assets of a registered data source in Data Catalog.

> Note
> - You can only perform stitching if the Tableau report is based on a database. Stitching Tableau reports based on files such as CSV is not supported.
> - Tableau stitching is based on full names and is case-sensitive. As a consequence, we recommend that you do not manually edit any asset names of data sources or Tableau assets. See the Tableau naming convention for more information.

## Tableau stitching steps

To use Tableau stitching, you have to prepare the assets representing the data source in Tableau's logical data layer and in Data Catalog's physical data layer:

| Step | What | Simplified instructions |
|---|---|---|
| 1 | Prepare the Tableau logical data layer. | 1. Register Tableau Server or Tableau Online.<br>2. Connect to Tableau Server or Tableau Online.<br>3. Synchronize Tableau sites. |
| 2 | Prepare the physical data layer. | 1. Register a database as data source.<br>2. Create a Database asset with the same name as the data source.<br>3. Create a relation between the Database asset and the Schema asset using the Technology Asset has / belongs to Schema relation type. |

| Step | What | Simplified instructions |
|------|------|------------------------|
| 3 | Stitch Tableau logical data layer and physical data layer. | 1. On the Tableau Data Model asset page, click **Stitch with data source**. |
| 4 | View stitching results. | 1. Open the asset page of the Tableau Server asset.<br>2. In the tab pane, click ⊡ **Diagram**.<br>3. In the **Explore** drop-down list, select **Data Catalog Lineage 5.7**. |

Note
- If there were changes in Tableau or the data source, you have to do the following:
  a. Synchronize Tableau. This can be done manually or automatically, by means of a synchronization schedule.
  b. Refresh the schema of your data source. This can be done manually or automatically, by scheduling it during data source registration.
  c. Restitch Tableau's logical data layer or Data Catalog's physical data layer. This has to be done manually.
- You can also remove stitching.

# Data layers

## Tableau's logical data layer

We call the data source in Tableau the logical data layer, because it consists of Tableau metadata, rather than the physical data. It is created when you synchronize a Tableau server. It contains Tableau report metadata, including the data source.

> **Note**
>   - You can combine different data sources in one Tableau data source by using different methods, for example, **Join** or **Union**.
>   - If you combine physical data sources in the Tableau data source with the **Join** method, the Tableau logical data layer is created in Data Catalog. For more information about the **Join** method, see Join Your Data.
>   - If you combine physical data sources in the Tableau data source with other methods, for example, **Union**, the Tableau logical data layer is not created in Data Catalog.

## Data Catalog's physical data layer

We call the data source in Data Catalog the physical data layer, which contains the physical tables and columns. It is created when you register a database as a data source. It contains the physical data of the data source.

## Stitching results

Each element is represented twice in Collibra: once in Tableau's logical data layer and once in Data Catalog's physical data layer.

The corresponding assets are linked by relations:

  - A relation of the type "Technology Asset source system for / source system Data Asset" type between the Database asset and the Tableau Data Model asset.
  - Relations of the type "Data Element targets / sources Data Element" type between the Column assets and the Data Attribute assets, based on the full names of the assets.

| Number | Data Catalog's physical data layer | Tableau's logical data layer | Description |
|---|---|---|---|
| 1 | Database (DB) | Tableau Data Model (TDM) | An abstraction from the physical implementation of database, schema, file, etc., used for Tableau report creation. |
| 2 | Schema (SCM) and Table (TBL) | Tableau Data Entity (TDE) | An abstraction from the physical implementation of database tables, used for Tableau report creation. |
| 3 | Column (COL) | Tableau Data Attribute (TDA) | A specification that defines a property of a Tableau data entity. Examples: CustomerBirthDate, EmployeeFirstName. |

# Naming convention

When you ingest a data source in Tableau, Tableau automatically creates names for the data source, data model, data elements and data attributes. When you create the logical data layer by synchronizing Tableau, Data Catalog uses the names in Tableau to create the corresponding Tableau assets. As a result, in Data Catalog, Tableau assets have as a full name the same name as the original data source names in Tableau.

When you create the physical data layer by registering the data source directly in Data Catalog, you enter the names of the Schema and Database assets manually. To make stitching work, we highly recommend to use the same name as the original data source to which the Tableau assets correspond as well:

- The name of the Schema asset should match a part of the Tableau Data Entity asset's full name. For example, *database-name > schema-name*.
- The name of the Database asset should match a part of the Tableau Data Model asset's full name.

The full name of the asset should match the asset path from the asset to the database it belongs to. For example, the full name of a Column asset would be *database>schema>table>column name*.

> Warning   Editing full name of the Tableau Server or Tableau Online assets may lead to errors during the synchronization process.

# Prepare the Tableau logical data layer

Before you can perform stitching, you have to prepare Tableau's logical data layer and Data Catalog's physical data layer. In this section, we describe how to prepare the logical data layer.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a resource role with the Configure external system resource permission, for example Owner.
- The Tableau user has the Download/Save As permission on the data source.

## Steps

1. Register Tableau Server or Tableau Online.
2. Connect to Tableau Server or Tableau Online.
3. Synchronize Tableau sites.
    » After synchronization, the assets of the following asset types are created in Data

Catalog:
- ◦ Tableau Data Model
- ◦ Tableau Data Entity
- ◦ Tableau Data Attribute

## What's next?

If you haven't done so yet, prepare the Data Catalog physical data layer.

After both the logical data layer and the physical data layer are created, you can stitch them.

# Prepare the Data Catalog physical data layer for Tableau stitching

Before you can perform stitching, you have to prepare Tableau's logical data layer and Data Catalog's physical data layer. In this section, we describe how to prepare the physical data layer.

## Prerequisites

- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a role with the following resource permissions on the **Schema** community:
    - ◦ Asset: add
    - ◦ Attribute: add
    - ◦ Domain: add
    - ◦ Attachment: add

## Steps

1. Register a database as data source.
    » After registration, the assets of the following asset types are created in Data Catalog:
    - ◦ Schema
    - ◦ Table

- ○ Column
2. Create a Database asset.

> Tip   We strongly recommend to use the name as your original data source, so that the name of the Database asset matches Tableau's naming convention.

1. Open Catalog.
2. In the main menu, click the **Create** ( + ) button.
3. Click the **Assets** tab.



4. Click Database.
   » The **Create Asset** dialog box appears.
5. Enter the required information.

| Field | Description |
|---|---|
| Type | The asset type of the asset that you are creating, in this case Database. |
| Domain | The domain to which the new asset will belong. You can only create a asset type in any domain of a domain type that is assigned to a Database asset type. |

| Field | Description |
|---|---|
| Name | The name of the Database asset. This has to match the name of the Tableau Data Model.<br><br>Tip<br>You can create multiple assets in one go.<br>To do this, press `Enter` after typing a value and then type the next. Depending on the settings, asset names may have to be unique in their domain. If you type a name that already exists, it will appear in strike-through style. |

6. Click **Create**.

   » A message at the top-right of your screen confirms that one or more assets are created.

3. Create a relation between the Database asset and the Schema asset using the Technology Asset has / belongs to Schema relation type.

   a. In the tab pane, click **Add Characteristic**.

      » The **Add a characteristic** dialog box appears.

   b. Click **Relations**.

   c. Search for and click  **has schema**.

      » The **Add has schema** dialog box appears.

    d.  Enter the required information.

| Option | Description |
|---|---|
| Assets | The name of the schema. |
| Filter suggested assets by organization | Option to filter the suggestions based on selected communities and domains.<br><br>If this option is selected, the organization tree appears. You can then filter and select domains and communities.<br><br> |
| Start date | Optionally enter the date on which the relation between the assets becomes applicable. Leave this field empty to create a permanent relation. |
| End date | Optionally enter the date on which the relation between the assets is no longer applicable. Leave this field empty to create a permanent relation. |

    e.  Click **Save**.

4. Check that the following relations are created for all Column assets that you want to stitch to Tableau assets:
   - Schema contains / is part of Table
   - Column is part of / contains Table

## What's next?

If you haven't done so yet, prepare the Tableau logical data layer.

After both the logical data layer and the physical data layer are prepared, you can stitch them.

# Supported data sources for Tableau stitching

You can stitch Tableau's logical data layer and Data Catalog's physical data layer for several data sources. The following table contains the packaged data sources and the driver versions that have been tested for Tableau stitching. We cannot guarantee that stitching works as expected for other data sources or versions.

| Data source | Tested versions for Tableau stitching |
| --- | --- |
| Amazon Redshift | 1.0.124969 |
| HP Vertica | 7.1.1-0 |
| IBM DB2 | This data source is not supported by Tableau. |
| MySQL | Tableau stitching is not possible because this data source has no schema. |
| Oracle | 11.2.0.4.0 |
| PostgreSQL | 9.5.1 |
| Microsoft SQL Server | 2014 (12.0.4422.0) |
| Snowflake | Snowflake editions supported by Tableau |

Note   Currently, we only support published Tableau data sources with an extract or a live connection. For more information, see the Tableau documentation.

# Stitch the Tableau logical data layer and the Data Catalog physical data layer

You can stitch Tableau's logical data layer and Data Catalog's physical data layer to represent the lineage from the data source to the Tableau reports.

## Prerequisites

- You have prepared Tableau's logical data layer.
- You have prepared Data Catalog's physical data layer.

## Steps

1. Open the Tableau Data Model asset page.

   > Tip   You can use the Search to quickly find the relevant asset.

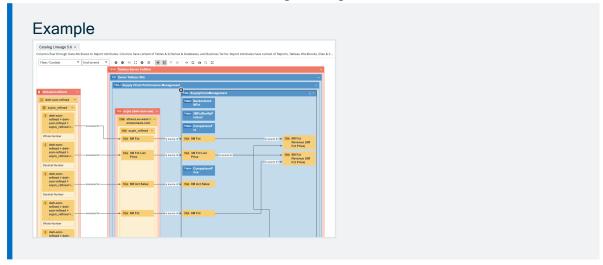2. In the upper-right corner, click **Stitch**.
   - » The **Stitch with data source** dialog box appears.

3. Enter the required information.

| Field | Description |
|---|---|
| Data Source | The Database asset that you want to stitch to this Data Model asset. |
| Filter suggested assets by organization | Option to filter the suggestions based on selected communities and domains.<br><br>If this option is selected, the organization tree appears. You can then filter and select domains and communities. |

4. Click **Stitch**.

> **Note**
> If a relation exists between the Tableau Data Model and the corresponding Database asset, of the "Technology Asset source system for / source system Data Asset" type, stitching happens immediately after clicking **Stitch**, without showing the dialog box.
>
> This occurs if you created the relation manually, or if you restitch.

## What's next?

Stitching is performed, creating relations between assets of Data Catalog's physical data layer and those of Tableau's logical data layer.

More precisely, these relations are created:

- A relation of the type "Technology Asset source system for / source system Data Asset" type between the Database asset and the Tableau Data Model asset.

- Relations of the type "Data Element targets / sources Data Element" type between the Column assets and the Data Attribute assets, based on the full names of the assets.

> Tip   You can view the stitching result as a diagram.

# Restitch the Tableau logical data layer and the Data Catalog physical data layer

After you completed stitching, there might be changes in Tableau or in the data source. For example, Tableau may have a new report and the data source may have a new column. To make sure that the lineage diagrams are also updated, you can restitch the data layers.

## Prerequisites

- You have previously stitched Tableau's logical data layer and Data Catalog's physical data layer.
- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have a resource role with the Attribute > Add resource permission.

## Steps

1. Ensure that Tableau's logical data layer is synchronized.
2. Ensure that Data Catalog's physical data layer is refreshed.
3. Open the Tableau Data Model asset page.

   > Tip   You can use the Search to quickly find the relevant asset.

4. In the upper-right corner, click **Stitch**.

# View stitching results

When stitching is complete, you can view the end-to-end lineage between the database and the Tableau report.

## Prerequisites

- You have prepared Tableau's logical data layer.
- You have prepared Data Catalog's physical data layer.
- You have stitched the logical data layer and the physical data layer.

## Steps

1. Open the Tableau Server asset page.
2. In the tab pane, click ⊡ **Diagram**.
3. In the view selector, select **Data Catalog Lineage 5.7**.



# Remove stitching between the Tableau logical data layer and the Data Catalog physical data layer

You can remove stitching to remove the relations between the logical data layer in Tableau and the physical data layer in Data Catalog.

More precisely, the following relations are removed:

- A relation of the type "Technology Asset source system for / source system Data Asset" type between the Database asset and the Tableau Data Model asset.
- Relations of the type "Data Element targets / sources Data Element" type between the Column assets and the Data Attribute assets, based on the full names of the assets.

## Prerequisites

- You have a resource role with the Configure external system resource permission, for example Owner.
- You have a resource role with the Attribute > Remove resource permission.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have stitched Tableau's logical data layer and Data Catalog's physical data layer.

## Steps

1. Open the Tableau Data Model asset page.

   > Tip  You can use the Search to quickly find the relevant asset.

2. Click **Actions** → **Remove stitching**.

   > Tip  If Catalog experience is disabled, the **More** menu is shown instead of **Actions**.

# Tableau provisioning

With Data Catalog, you can create data sets and convert them to the Tableau format. This enables you to use Collibra-managed data in Tableau.

# The Tableau provisioning file

A Tableau provisioning file is a packaged data source file with the extension TDSX. The packaged data source file is a ZIP file that contains a data source file and any local file data sources. You can import it in Tableau to, for example, analyze the data. It has the extension TDSX.

You can create a Tableau provisioning file from any data set in Data Catalog.

The file contains the following information:

- A TDS file: This is an XML file that contains the data source definition.
- The actual ingested files, if the data set contains data from Excel or CSV data sources.

```
<?xml version="1.0" encoding="UTF-8"?>
<datasource xmlns:user-
r="http://www.tableausoftware.com/xml/user"
formatted-name="<name of your data set>" inline="true" ver-
sion="10.0">
    <connection class="federated">
        <named-connections>
            <named-connection caption="public" name="<connection-ID>">
                <connection authentication="username-password" class-
s="<data-source-type>" dbname="<database-ID>" port="" schem-
a="public" server="<hostname:port>"/>
            </named-connection>
        </named-connections>
        <relation connection="<relation-ID>" name="<name-of-rela-
tion>" table="[public].[<name-of-relation>]" type="table"/>
    </connection>
</datasource>
```

# Required JDBC driver information for Tableau provisioning

To create a Tableau provisioning file from a data set, the JDBC driver of its data source needs the following properties:

| Data source | Required connection properties |
|---|---|
| Amazon Redshift | • host<br>• port<br>• database<br>• schema |
| HP Vertica | • host<br>• port<br>• database<br>• schema |
| MySQL | • host<br>• port<br>• database |
| Oracle | • host<br>• port<br>• database<br>• schema |
| PostgreSQL | • host<br>• port<br>• database<br>• schema |
| SQL Server | • host<br>• port<br>• database<br>• schema |

For more information, see the JDBC configuration details of the various databases.

# Create Tableau provisioning file

In Data Catalog, you can create Tableau provisioning files from data sets.

> **Tip** If your data set's origin is a relational database, you need the credentials to connect to that database. Make sure the JDBC driver has all the required information in the correct format before you create the provisioning file.

## Prerequisites

- You have a resource role with the Access data resource permission, for example Data Analyst Level 2.
- You have a global role with the Catalog global permission, for example Catalog Author.
- You have enabled Tableau provisioning in Collibra Console.

## Steps

1. On the main menu, click ⠿ , then 🗄 **Catalog**.

   » The Catalog Home opens.
2. In the submenu, click **Data Sets**
3. Click the data set that you want to use in Tableau.
4. Above the table, to the right, click **Actions** → **Access Tableau source (beta)**.



   » The Tableau provisioning file in TDSX format is downloaded.

   > **Tip** If Catalog experience is disabled, the **More** menu is shown instead of **Actions**.

## What's next?

You can now import the TDSX file in Tableau.

# Troubleshooting

The following table contains the most common issues that you can encounter while ingesting or synchronizing Tableau.

| Issue | Solution |
| --- | --- |
| Tableau images are not fetched correctly | Synchronizing Tableau data means ingesting metadata from Tableau to your Collibra Data Intelligence Cloud environment. The metadata is represented as assets of specific types and their characteristics. Images such as report thumbnails are not downloaded and stored in Data Catalog. Instead, Data Catalog stores a link to the image. Every time you open the asset page, the image is fetched from Tableau. Images are not fetched correctly if there is a problem with this link. A common issue is caused by the base URL parameter, which is part of the link. If the base URL is not set correctly in Collibra Console, the links to the Tableau images are broken. To fix this issue, edit the base URL in Collibra Console. |

| Issue | Solution |
|---|---|
| When you synchronize a Tableau Server 2020.2, some Tableau data is skipped. | In most cases, this occurs when people are actively using Tableau while Data Catalog is synchronizing the Tableau Server. The technical reason is that the APIs collect Tableau data at different times. If users make changes in Tableau, the data that is collected by the APIs may be inconsistent. When that happens, the corresponding assets are not synchronized in Data Catalog. We highly recommend that you synchronize Tableau after working hours. This reduces the chance that Tableau data is added, changed, renamed or deleted on Tableau's side during the synchronization process. |
| Tableau synchronization fails with error message `Duplicate key`. | The Tableau synchronization fails with the `Duplicate key` error when you have multiple views with the same name in the same workbook. To solve this problem, we highly recommend to give each view in Tableau a unique name before you synchronize the Tableau Server in Data Catalog. |
| When you integrate Tableau, relations between Tableau Data Attributes and Tableau Report Attributes are not being created. | The likely issue is that the report is using embedded data sources or custom SQL. If that is the case, we recommend that you switch to the new Tableau integration method, via the lineage harvester, which supports embedded data sources and custom SQL. For more information, see Features and limitations of Tableau integration via the lineage harvester. |

# Catalog workflows

To keep the information flows that are shipped with the Catalog product configurable, a part of the functionality is achieved through workflows. You can configure the packaged workflows, but they are designed to work together: if you decide to change one of the workflows, verify the other Catalog workflows, since they may depend on one another.

> Tip  For more information about workflows, see the Collibra Developer Portal.

| Name | Description |
|------|-------------|
| Assign Owner To Data Set | This process automates adding owners to data sets. This workflow is automatically triggered when a new Data Set asset is created. |
| Cancel Process | This process notifies the concerned users of a workflow cancellation. |
| Escalation Process | This process is the default mechanism for the escalation of user tasks in workflows. |
| Post Data Ingestion Workflow | This process facilitates assigning the Owner and Technical Steward for newly ingested Schema assets. This workflow is automatically triggered when a new Schema asset is created and after a data source is registered. |
| Propose New Business Asset | This process facilitates the creation of new Business Assets in the **Data Governance Council** community. |
| Propose New Data Asset | This process facilitates the creation of new Data Assets in the **Data Governance Council** community. |

| Name | Description |
|------|-------------|
| Propose New Technology Asset | The Propose New Technology Asset workflow allows you to create a new Technology asset in Collibra Data Intelligence Cloud. By default, the asset is added to the **Data Governance Council** community, in the **New Applications** domain. |
| Request Assets Access | The Request Assets Access workflow allows you to request access to assets that are referenced in your shopping cart. All data owners have to approve the request before you can access the assets.<br><br>More information<br><br>The workflow calculates the name of the asset by combining the creation date with a sequential number for that day, for example 2019-09-30 #1 and sets the asset characteristics according to the data submitted through the start form. The user who started the workflow receives the Requester role. The user with an Owner role approves the request for each data set and the Owner or Technical Steward provides access to the data set elements.<br><br>Note<br>This workflow replaces the Request Data Sets Access workflow.<br><br>If you restore a 5.4.x backup or older, the Requests Data Sets Access will overwrite the packaged Request Assets Access workflow. You have to deploy the Requests Assets Access workflow again and apply all possible customizations to the new workflow.<br><br>You can also manually request access to data sets access to data sets and reports. |
| Simple Approval | The Simple Approval workflow is a single-step process that allows you to approve an asset in Collibra Data Intelligence Cloud. |

| Name | Description |
|------|-------------|
| Voting Sub-Process | The Voting Sub-Process is a workflow that can be called by other workflows when users need to vote. It is used within other packaged workflows such as the Approval Process, the Simple Approval or the Issue Management workflow.<br><br>You can use this sub-process in new custom workflows. The result is a true or false boolean that is provided to the parent workflow. |

# Catalog Troubleshooting

If you are experiencing general issues with the Data Catalog feature, consult the articles in this section.

If you have issues with ingesting a BI source or with Collibra Data Lineage, please visit their individual troubleshooting sections:

- Tableau troubleshooting
- Power BI troubleshooting
- Looker troubleshooting
- Collibra Data Lineage troubleshooting

## What's the difference between Data Catalog and Collibra Connect?

Data Catalog and Collibra Connect have many overlapping features. Which of them is more suited for your situation, depends on a number of factors.

In a nutshell, you use Data Catalog for ingesting metadata from popular database types via a predefined ingestion logic, which is ideal for business users. You can then see the metadata in the form of assets and characteristics. You use Collibra Connect to read and write metadata in any API-supported system and provide the metadata to Collibra Data Intelligence Cloud. Collibra Connect has more flexibility with regard to ingestion, but requires technical skills.

|  | Data Catalog | Collibra Connect |
|---|---|---|
| Definition | The Collibra Data Catalog is an application that helps the business data analyst to discover, describe, assemble and govern data sets, in order to improve trust in analytics based on those data sets. | Collibra Connect is an integration platform that enables integrations between Collibra and other third-party products, such as Informatica, Salesforce.com and JIRA. |
| Purpose | Data Catalog can ingest and represent metadata of specific data sources as assets and characteristics, including diagrams. | Collibra Connect is meant as an advanced interface between Collibra and data sources of any third-party vendors. |
| Processes | <ul><li>Metadata ingestion</li><li>Profiling and data type detection</li><li>Read only</li></ul> | <ul><li>Bidirectional synchronization of metadata</li><li>No profiling</li><li>Read and write</li></ul> |
| Integrations | <ul><li>JDBC-supported databases such as PostgreSQL and IBM DB2.</li><li>File-based databases in Excel and CSV.</li><li>External systems such as Tableau and Amazon S3.</li></ul> | Any system with:<ul><li>API support</li><li>Structured metadata format such as XML and JSON</li></ul> |
| Ingestion | Predefined metamodel and ingestion logic | Flexible and configurable metamodel and ingestion logic |
| Usability | <ul><li>Usable via Collibra</li><li>Business user friendly</li></ul> | <ul><li>Configuration via IDE</li><li>Requires development skills to set up</li></ul> |
| More information | <ul><li>The Data Catalog: What it is, Why you Need it, and How to Make it Successful</li><li>The Data Catalog section of the Collibra user guide.</li></ul> | <ul><li>Introduction to Collibra Connect</li><li>The Collibra Connect user guide.</li></ul> |

# How to enable logging for data ingestion

If you want to troubleshoot issues with data ingestion, you have to enable logging for data ingestion. By default, logging for data ingestion is disabled because your data can be exposed.

For more information, see Environment log settings for DGC services and Environment log settings for Repository services.

> **Warning**   If you have investigated the data ingestion issue, don't forget to revert all the changes from this section.

## Steps

1. Open the Data Governance Center logging settings.
    a. Open Collibra Console.
        » Collibra Console opens with the **Infrastructure** page.
    b. In the tab pane, click the **Data Governance Center** service of the environment whose log settings you need.
    c. Click **Logs**.
    d. Above the table, to the right, click ⚙ **Settings**.
2. Click **Add logger**.
    » The **Add logger** dialog box appears.

3.  Enter the required information.

| Field | Description |
|---|---|
| Logger name | The name of the logger.<br><br>Enter one of the following:<br><br>○ `com.collibra.dgc.catalog.service.schema.impl`<br>○ `com.collibra.dgc.catalog.service.impl`<br>○ `com.collibra.jobserver.client`<br>○ `com.collibra.dgc.catalog.service.datausage.impl`<br>○ `com.collibra.catalog.core.service.datausage.impl`<br>○ `com.collibra.catalog.core.service.schema.impl`<br>○ `com.-`<br>  `collibra.catalog.core.service.schema.impl.ingestion`<br>○ `com.-`<br>  `collibra.catalog.core.service.schema.impl.profiling`<br>○ `com.collibra.catalog.core.service.schema.impl.report`<br>○ `com.collibra.catalog.core.schema.impl`<br>○ `com.collibra.catalog.core.schema.impl.ingestion`<br>○ `com.collibra.catalog.core.schema.impl.profiling`<br>○ `com.collibra.catalog.core.schema.impl.report` |
| Logger level | The amount of log entries you want in the logs.<br><br>Select **DEBUG**. |

4.  Click **Add logger**.
5.  Repeat this until you have added all the loggers.

# The Jobserver logs are out of memory

When the Jobserver log files are out of memory, the logs that are created during ingestion or profiling are deleted immediately after they are created.

# Solution

1. Stop the environment for which you want to update the memory settings.
1. Open a terminal session on the server that hosts the jobserver.
2. Open the file **/opt/collibra/spark-defaults.conf** and do the following.
    a. Add the following line to the configuration file:

    ```
    spark.driver.maxResultSize = 1536m
    ```

    b. Save and close the file.
3. Open the **/opt/collibra/spark-jobserver/conf/log4j-server.properties** file and do the following.
    a. In the `Root logger option` section, update the properties to match this section:

    ```
    # Root logger option
    log4j.rootLogger=INFO,LOGFILE
    log4j.appender.LOGFILE=org.apache.log4j.RollingFileAppender
    log4j.appender.LOGFILE.File=${LOG_DIR}/spark-job-server.log
    log4j.ap-
    pender.LOGFILE.layout=org.apache.log4j.PatternLayout
    log4j.appender.LOGFILE.layout.ConversionPattern=%d{yyyy-MM-
    dd HH:mm:ss.SSS} %-5p [%t] %c{3} - %m%n
    log4j.appender.LOGFILE.maxFileSize=100MB
    log4j.appender.LOGFILE.maxBackupIndex=30
    log4j.logger.org.apache.spark=WARN
    log4j.logger.spark.jobserver.context=WARN
    log4j.logger.akka=WARN
    log4j.logger.com.collibra.jobserver.job=DEBUG
    log4j.logger.com.collibra.catalog.profilers=DEBUG
    log4j.-
    log-
    ger.com.collibra.catalog.profilers.Pass1TableProfiler$=INFO
    log4j.logger.com.collibra.catalog.ingestion=DEBUG
    log4j.logger.com.collibra.jdbc=DEBUG
    ```

    b. Save and close the file.
4. Start the environment again.

# Ingestion out-of-memory error

When you upload a JDBC driver larger than 50 MB or when you have uploaded multiple JDBC drivers, you may encounter an out-of-memory error. Due to this problem, the jobserver does not release the memory needed to store the driver in memory.

# Resolution

To solve this problem, you have to increase the memory of the Jobserver application, for example, increase it to 3 GB.

1. Stop the environment for which you want to update the memory settings.
2. Open a terminal session on the server that hosts the jobserver.
3. Open the file **<drive>/collibra/spark-jobserver/conf/jobserver.conf** for editing.
4. Look up the parameter **driver-memory**.
5. Edit the parameter value, for example, *3G*, corresponding with 3 GB.

   The default value is 2G.
6. Save and close the file.
7. Open the file **<drive>/collibra_data/spark-jobserver/config/server.json** for editing.
8. Look up the parameter **jobserverMemory**.
9. Edit the parameter value, for example, *2048M*, corresponding with 2 GB.

   The default value is 1024M.
10. Save and close the file.
11. Start the environment again.

# Error when managing connection properties of a driver for Jobserver

## Issue

When you want to change the properties of a connection used to register data sources via Jobserver, you receive the following error message:

```
CollibraIllegalStateException: jdbcDriverCannotBeUpdatedWhenLinked
when trying to delete or edit a JDBC driver.
```
or
```
You cannot update the driver because it is linked to a Schema Asset
```

# Reason

Once you have successfully used a connection to register a data source via Jobserver, you cannot update the connection properties anymore.

# Solution

If you want to change the properties for a driver, you need to create a new driver:

1. Open a schema registered via the driver you want to update.
2. Go to **Actions → Refresh**.
3. In **JDBC driver version**, select **Manage drivers...**.
4. Create a new driver for the data source.
   As a best practice for the name of the drivers, use a naming convention which includes the data source and a version number. For example: Google BigQuery 1.5 or MySQL 5.9. If you want to use the same driver version with other properties, add an extra number. For example: Google BigQuery 1.5 v2.
   For details on the properties, see Manage Collibra-provided JDBC drivers.
5. Save the new driver.
   » The new driver is automatically applied to the schema.
6. For each schema that uses the old driver, go to **Actions → Refresh**, and select the new driver.

# Missing schema name during data ingestion

If you ingest a data source with a new JDBC driver, you can receive an error "No schema has been specified".

> Note   In the stacktrace you can see a "`CollibraIllegalArgumentException`" message.

# Solution

Make sure that you defined a schema property for the new JDBC driver.

# Different versions for Collibra and Jobserver

You can install the services of a Collibra Data Intelligence Cloud environment on multiple nodes. If you do so, make sure that you use the same installer on all the nodes. This also applies to upgrading an environment.

If your environment has different versions for the Data Governance Center and Jobserver services, the following errors will occur when you run an ingestion.

- Spark Context's logs

```
[2017-11-07 07:27:15,608] WARN nalRequestDataDeserializer []
[akka://JobServer/user/jobManager-c7-8eec-de0c02029808] - Pack-
age com.collibra.jobserver.dto.catalog.ingestion, different ver-
sion detected: client uses version 1.2.4-SNAPSHOT, server uses
version 1.2.2-SNAPSHOT
```

- Collibra logs

```
20:21:43.407 [Procedure Manager] WARN c.c.j.c.i.s.StateDeseri-
alizer - Package com.collibra.jobserver.dto.catalog.profiling,
different version detected: client uses version 1.1.10, server
uses version 1.1.8
```

# Solution

Install all the Collibra services with the same installer.

# Resolve schema refresh conflicts via Jobserver

> Note   This information only applies to Jobserver. For information on how Edge handles differences between the original schema and the updated schema, see About synchronizing schemas.

If you refresh a schema via Jobserver, the ingestion process detects differences between the original schema, already in Collibra Data Intelligence Cloud, and the updated schema.

If columns or tables have been added to or removed from the schema, the process will create or delete the corresponding Column and Table assets in Collibra. However, the ingestion process results in a refresh conflict if one or more columns or tables were added and others were removed. If that happens, it adds a Refresh conflict attribute to all added and removed columns or tables. You have to resolve these conflicts before you can refresh the schema again. If you do not resolve the refresh conflicts, any future attempts to refresh the data source will fail.

To see if there are any conflicts after a refresh, you have to add the **Refresh Conflict** field to the **Data Sources** view of the schemas.

You may come across the following scenarios:

- A column is deleted from the schema and another one is added to the schema:
    a. You have to manually delete the column asset.
    b. You have to remove the **Refresh conflict** attribute from the added column asset.



- A column is renamed in the schema:
    a. You have to remove the column asset with the updated column name.
    b. You have to rename the original column name to the newly ingested column name and delete the **Refresh Conflict** attribute.

- A column is deleted from the schema: this is automatically detected by the refresh operation. No further action is required of you.
- A column is added to the schema: this is automatically detected by the refresh operation. No further action is required of you.
- A table is renamed in the schema:
  a.  You have to manually delete the renamed new table and all the columns contained in the table.
  b.  You have to manually rename the existing old table and all the columns contained in the table.



- A table is deleted from the schema and another table is added to the schema:
  a.  You have to manually delete the deleted table and all the columns in the table.
  b.  You have to manually delete the Refresh Conflict attribute for the added table.

# Resolve a schema refresh conflict when columns are added and deleted at the same time

If you refresh a schema, the ingestion process will detect conflicts if the data source has the following changes:

- A column has been removed.
- A column has been added.

In the following example, the ingested schema has a column **age** and in the updated schema, the column **age** is removed and a column **birthday** is added.
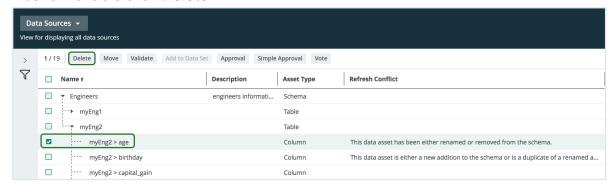
To resolve such a refresh conflict, follow these steps:

1. Look up the data source with the search function or as follows:
   a. On the main menu, click ⠿, then ⊟ **Catalog**.

      » The Catalog Home opens.
   b. In the submenu, click **Data Sources**.
   c. Optionally, add the **Refresh Conflict** column to the table.
   d. In the table, expand the relevant schema and table to find the columns with refresh conflicts.
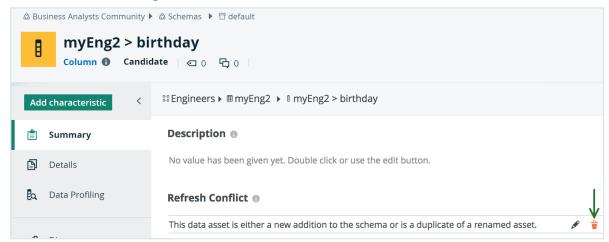
2. Select the column that is removed from the data source. In this example it is the **age** column.

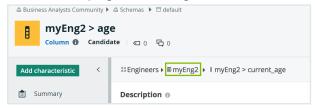   If necessary, select all column assets that are removed from the data source.

3. Above the table click **Delete**.



4. Click **Yes** to confirm the deletion of the column.
5. Click the name of the added column name.
   » The column asset page appear.
6. In the **Refresh Conflict** section of the column asset page, hover over the message and click 🗑 on the right-hand side.



7. Click **Yes** to confirm the deletion of the attribute.
8. Click the browser's **Back** button to return to the **Data Sources** view of the table.

   You can also click on the breadcrumb, as shown in the following image, to open the table asset page of the ingested schema.'
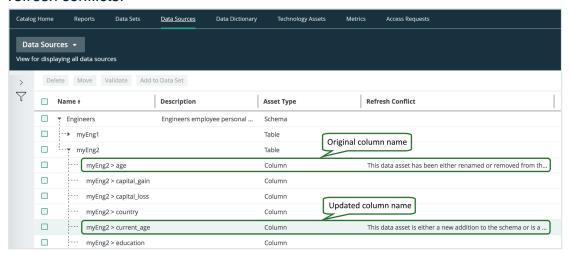
9. Repeat steps 5 to 8 for all other added columns.

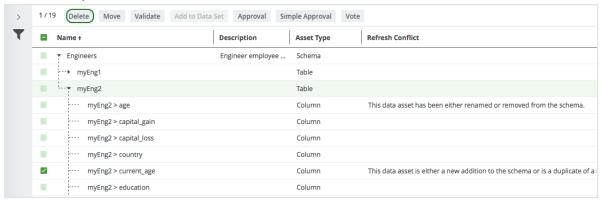# Resolve a schema refresh conflict for a renamed column

If you refresh a schema where the data source contains a column that has been renamed, the ingestion process will detect a conflict. In the following example, the ingested schema contains a column **age**, and in the updated schema, the column name has become **current_age**.

To resolve a refresh conflict due to a column rename, follow these steps:

1. Look up the new column with the search function or as follows:
   a. On the main menu, click ⊞, then 🗄 **Catalog**.

   » The Catalog Home opens.
   b. In the submenu, click **Data Sources**.
   c. Optionally, add the **Refresh Conflict** column to the table.
   d. In the table, expand the relevant schema and table to find the columns with refresh conflicts.

2. Select the updated column name and click **Delete** above the table.

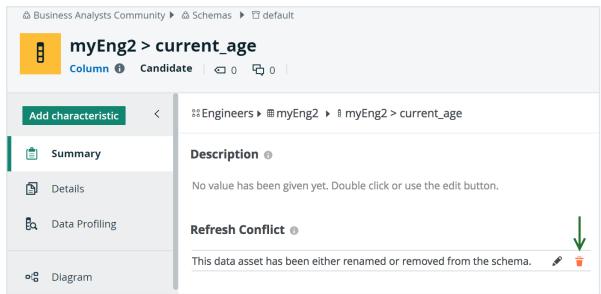   If necessary, select all column assets that are removed from the data source.



3. Click **Yes** to confirm the deletion of the column asset(s).
4. Click the name of the original column name.

   » The column asset page appears.

5. In the resource toolbar, click **Actions** > **Edit**.

   » The **Edit <asset name>** dialog box appears.

6. Change the name to the new ingested name.



7. Click **Save**.
8. Refresh the page.
9. Leave the column asset page open.

10. In the **Refresh Conflict** section of the column asset page, hover over the message and click 🗑 on the right-hand side.



11. Click **Yes** to confirm the deletion of the attribute.
12. Click the browser's **Back** button to return to the **Data Sources** view of the schema. You can also click on the breadcrumb, as shown in the following image, to open the table asset page of the ingested schema.



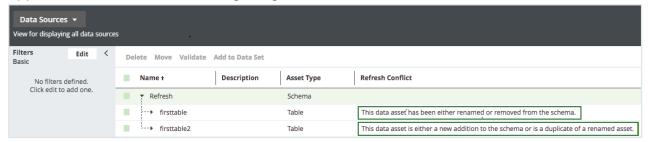13. If necessary, repeat steps 4 to 12 for other renamed column assets.

## What's next?

You can now safely refresh the schema with the new data source; however, keep in mind this may take some time.

# Resolve a schema refresh conflict for a renamed table

If you refresh a schema where the data source contains a table that has been renamed, the ingestion process detects a conflict.

In the following example, the original schema **Refresh** contains the table **firsttable**. This table has been renamed to **firsttable2**. After refreshing the schema, refresh conflicts appear, as shown in the following image:
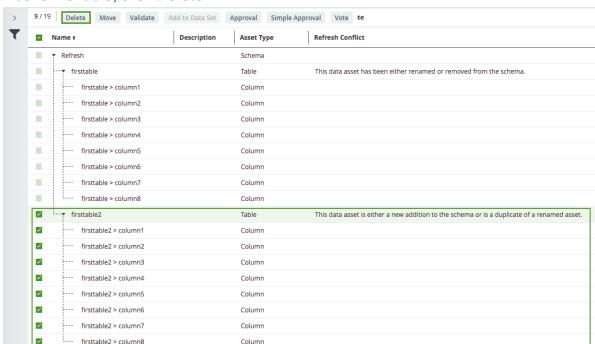


You have to manually resolve the conflicts before you continue. It is not possible to refresh a schema when there are conflicts.

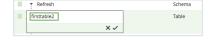> Note You have to add the **Refresh Conflict** column to the table if it is not there already.

## Steps

1. On the main menu, click ⠿, then ⊟ **Catalog**.

    » The Catalog Home opens.

    » The Catalog Home appears

2. In the submenu, click **Data Sources**.

3. Expand the tables to see all the columns that are contained in them.

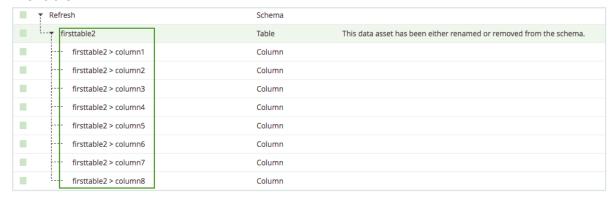4. Select the renamed table and all its contained columns, in this example, **firsttable2**.

5. Above the table, click **Delete**.



6. Click **Yes** to confirm the deletion.
7. Hover over the original table, in this example, **firsttable**, and click ✎ to the right of the table name.
8. Change the name to the new ingested table name, in this example, **firsttable2**, and click ✔ to apply the change.



9. Hover over a column contained in the table you just renamed and click ✎ to the right of the column name.
10. Rename the column by replacing the table part of the name with that of the renamed table and click ✔ to apply the change.

    The column name is a concatenation of the table name and the original column name and so you just have to replace the table part of the name with the new table name. For example, to rename the column name **firsttable > column1** to **firsttable2 > column1**, you just have to change **firsttable** to **firsttable2** so that the column name becomes **firsttable2 > column1**.
11. Repeat this action for all the columns in the renamed table.

    Now, you only see the new ingested table, **firsttable2**, and the columns contained in

the table.



12. Click the name of the renamed table.

&raquo; The table asset page appears.

13. In the **Refresh Conflict** section, hover over the refresh conflict message and click 🗑 on the right-hand side.



14. Click **Yes** to confirm the deletion of the Refresh Conflict attribute.

## What's next?

You can now safely refresh the schema with the data source.

# Resolve a schema refresh conflict when tables are added and deleted at the same time

When you refresh a schema, the ingestion process detects conflicts if the data source has the following changes at the same time:

- A table has been removed.
- A table has been added.

In the following example, the original schema **Postgre** contains the table **Employee** and the table **CompanyList**. A new table **Schools** has been added to the schema and the table **CompanyList** has been deleted. After refreshing the schema, refresh conflicts appear for

the added table and the deleted table, as shown in the following image:



You have to manually resolve the conflicts before you continue. It is not possible to refresh a schema when there are conflicts.

> Note   You have to add the **Refresh Conflict** column to the table if it is not there already.

## Steps

1. On the main menu, click ⠿, then ▤ **Catalog**.
   - » The Catalog Home opens.
   - » The Catalog Home appears.
2. In the submenu, click **Data Sources**.
3. Select the deleted table and all its contained columns, in this example, **CompanyList**.



4. Above the table, click **Delete**.
5. Click **Yes** to confirm the deletion.
6. Click the name of the added table, in this example, **Schools**.
   - » The table asset page appears.
7. In the **Refresh Conflict** section, hover over the refresh conflict message and click 🗑 on the right-hand side.

**Refresh Conflict** ⓘ

This data asset is either a new addition to the schema or is a duplicate of a renamed asset.

8. Click **Yes** to confirm the deletion of the Refresh Conflict attribute.

## What's next?

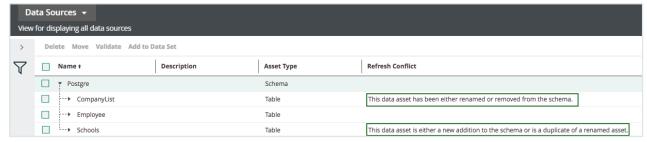You can now safely refresh the schema with the data source.

# Advanced data type detection is slow

Advanced data type (ADT) detection is the process that compares each value in the database with each pattern in the ADT definition list.

The following non-exhaustive list contains the factors that affect the detection time:

- The higher the number of ADTs in Catalog, the longer the detection time.
- The higher the number of patterns in each ADT, the longer the detection time.
  For example, a text ADT can contain one or more regular expressions. The more regular expressions that you add to this ADT, the longer the detection time will take.

> Tip   As a general rule, try to limit both the number of ADTs and the number of patterns per ADT.

# Jobserver troubleshooting

This is a list of known issues in versions older than Collibra 2022.09.

| Problem | Solution |
|---|---|
| One or more of the following error messages appear:<br><br>• `context JS-<context ID> not found` in the Jobserver node in DGC logs.<br>• `manager_start - /opt/collibra/spark-job-server/bin/manager_start.sh: line 73: <process id> killed` in the Jobserver server logs.<br>• Spark context logs are interrupted during Spark processing.<br>• `It is not possible to allocate enough memory` in the Spark process or other process on the same machine. | If the Spark context crashes or is unresponsive, it can be related to a memory shortage. Make sure that you have enough memory.<br><br>• In 5.7, a Jobserver node should have 64GB RAM, 16 CPUs and 500GB SSD.<br>• In 5.7.1, the Spark context process configuration for each Jobserver requires you to change the lower the heap memory to 40GB and replace the -XX:+UseG1GC option by -XX:+UseParallelGC. |
| An ingestion job keeps on running due to lingering Spark Context. | Restart the Jobserver, then restart Collibra. |
| Communication failure occurs between Jobserver and Spark Context when profiling large tables. | The following relevant parameters can be edited in the Jobserver configuration file to decrease the chance that this problem occurs:<br><br>• `acceptable-heartbeat-pause` should be 600s.<br>• `heartbeat-interval` should be 300s.<br>• `threshold` should be 12.0 |

# Jobserver jobs

To ingest data in Collibra Data Intelligence Cloud, you have to register a data source. During the ingestion, you can include to run a data profiling, data sampling and to detect

advanced data types in the data.

The DGC service is responsible for the ingestions, the Jobserver is responsible for the data profiling, data sampling and advanced data type detection.

The following table shows how many jobs it takes to complete a task. The jobs are executed sequentially.

| Task | Number of jobs |
|---|---|
| Data profiling | 4 jobs per table |
| Data sampling | 2 jobs per table |
| Advanced data type detection | 1 job per table |
| Data ingestion | 0 job |

If you have to troubleshoot Jobserver jobs, you need the following log files when you create a diagnostic file.

- Collibra logs
- Jobserver logs: You have to enable the ingestion and profiling logs.
- Spark logs: You have enable to the Spark logs. When you create a diagnostic file, these are included with the Jobserver logs.

# Data Marketplace (Beta)

Data Marketplace is a beta feature that improves the search experience for data catalog in Collibra Data Intelligence Cloud.

Administrators can expose data within an internal marketplace, ensuring only relevant and high quality data is available for organization-wide use.

Data consumers, such as analysts and data engineers, can use a user-friendly search interface to browse, preview, and request the data they need.

## About Data Marketplace

Data Marketplace is a beta feature that improves the search experience for data catalog in Collibra Data Intelligence Cloud.

Administrators can expose data within an internal marketplace, ensuring only relevant and high quality data is available for organization-wide use.

Data consumers, such as analysts and data engineers, can use a user-friendly search interface to browse, preview, and request the data they need.

By default, the feature provides access to specific Data Catalog asset types such as reports, data sets, and tables. The feature also uses some out-of-the box relations for finding and filtering the results. The feature is customizable so that administrators can decide which assets are exposed in the results and which filters and actions are available.

Data Marketplace allows you to find relevant and high quality data more easily thanks to the following features:

### 1. Data Marketplace only searches in data determined by the Data Marketplace scope

Data Marketplace only searches for some Data Catalog-specific asset types, such as data sets, reports and tables. Administrators can extend or refine this scope based on Data Catalog asset types, statuses and organization (communities and domains).
Data Marketplace cannot search Glossary domains and Business Term assets.

### 2. Data Marketplace makes searching for data easier

- Filter tabs allow you to quickly apply a frequently used filter set.
- Filters allow you to refine the search results.
- Data Marketplace can search for data based on specific relations between the data assets. For example, when searching for a column name, you can find the data sets to which the column is added.
- Data Marketplace can guide you to the information you might be looking for via search suggestions.

| 3. Data can be opened in preview |
|---|
| The asset preview shows the asset details, responsibilities, privacy restrictions, related assets, and so on. This allows for investigating the data without losing track of the search results.<br>You have quick and easy access to recent search queries, recently visited assets, and common actions to interact with the data. Administrators can define which actions must be available in various pages in Data Marketplace. |

> Note   Data Marketplace is a beta feature. To test the feature, an administrator must activate the feature and configure Data Marketplace.

# Data Marketplace compatibility

Data Marketplace is only available on Collibra Data Intelligence Cloud.

Administrators can activate and configure Data Marketplace as a beta feature from version 2022.08.

# Data Marketplace permissions

In Data Marketplace, you can be a data consumer or an administrator.

|  | Required global role or permission | Someone with this role or permission can |
|---|---|---|
| Data consumer | You need the Catalog global role. | <ul><li>Access to Data Marketplace from the Applications icon ▦ .</li><li>Search for assets and access the asset preview, if you have view permission on the asset.</li><li>Start the actions that have been assigned to Data Marketplace by the administrator.</li><li>Enable search suggestions.</li></ul> |

| | Required global role or per-mission | Someone with this role or permission can |
|---|---|---|
| Administrator | You need the SysAdmin global role. | • Same as Data consumer.<br>• In Collibra, access Data Marketplace-specific settings to configure the scope, filters, and actions.<br>• In Data Marketplace, access the configuration of the filter sets in the filter tabs. |

# Enable Data Marketplace

Data Marketplace is a beta feature. If you want to test it, you need to activate the feature.

# Prerequisites

You have the ADMIN or SUPER role in Collibra Console.

# Steps

1.  Open the DGC service settings for editing:
    a.  Open Collibra Console.
        »   Collibra Console opens with the **Infrastructure** page.
    b.  In the tab pane, expand an environment to show its services.
    c.  In the tab pane, click the Data Governance Center service of that environment.
    d.  Click **Configuration**.
    e.  Click **Edit configuration**.

2. In the **Beta Features** section, enable the **Data Marketplace** setting:

| Setting | Description |
|---|---|
| Data Marketplace | - ✔ True: Data Marketplace is enabled. Anyone with the required permissions can use or configure the Data Marketplace application from the Applications icon ⊞. <br> - ✘ False: Data Marketplace is not enabled. This is the default value. <br> After you have enabled this setting, reindex Data Marketplace relations or reindex Collibra completely. |

3. Optionally, enable the **Search over relations** setting.

| Setting | Description |
|---|---|
| Search over relations | - ✔ True: The Data Marketplace search also takes specific relations between assets into account. As a result, your search results not only include assets that directly match the search criteria, but also assets that match the criteria through specific relation types. <br><br> > Example   A column with the name "Order" is included in a data set called "Customer". When Search over relations is active and you search in Data Marketplace for "Order", then the data set "Customer" appears in the search results because the data set contains this column. <br><br> For more information on this feature and the default relation types, see Filtering and searching based on relations in Data Marketplace. <br><br> - ✘ False: Data Marketplace search results do not take relations into account. This is the default value. <br> After you have enabled this setting, reindex Data Marketplace relations or reindex Collibra completely. |

4. Optionally, enable the following settings in the **Search index configuration - Search Event Log configuration**:

| Setting | Description |
| --- | --- |
| Asynchronous indexing | ✔ True: The Search service handles changes to your environment in batches. This ensures optimal performance, especially after big changes to your environment, such as importing a lot of assets.<br><br>✖ False: Changes to your environment are processed immediately.<br>This is the default value.<br><br>> **Warning** This setting has an impact on your entire Collibra environment, not only to Data Marketplace.<br><br>For more information, see Asynchronous indexing.<br><br>After you have enabled this setting, stop and restart your environment, and reindex Collibra completely. |

| Setting | Description |
|---------|-------------|
| Consistency for Custom facets and Search over Relations | This setting is specific to Data Marketplace.<br><br>This setting keeps Data Marketplace up to date if you made changes to relations between assets.<br><br>○ ✔ True: The relation information remains consistent between Collibra and Data Marketplace for out-of-the-box filters and relations. If there is a change in the relation between assets , the change will be reflected in the search index after some time.<br><br>> Warning   If you enable this setting, the Asynchronous indexing setting must be enabled too because every relation change results in an event that is processed via asynchronous indexation.<br><br>○ ✘ False: Changes to relations are not automatically indexed. This can cause inconsistencies between Collibra and Data Marketplace. You can, however, reindex Data Marketplace relations manually.<br>This is the default value. |

5. Optionally, enable analytics in the **General settings** to show recommendations in the Data Marketplace landing page.

| Setting | Description |
|---------|-------------|
| Enable Usage analytics(*) | The analytics are used to understand how users interact with Collibra. The information is used to provide historical reporting.<br><br>When you enable the analytics in a Cloud environment, the data collection starts immediately.<br>When you enable it in an on-premises environment, you have to approve the tracking script (pendo.io, app.pendo.io and cdn.pendo.io).<br><br>○ ✔ True (default): Gathering analytics is enabled and sent to Collibra.<br>○ ✘ False: Gathering analytics is disabled.<br>After activating this setting, Data Marketplace needs one day to start providing recommendations. |

6. Click the green **Save all** button.
7. If you have enabled Asynchronous indexing, stop the environment and start the environment.
8. If you have enabled the Data Marketplace, Asynchronous indexing or Search over relations setting, reindex Data Marketplace relations or reindex Collibra completely.

# What's next?

Configuring Data Marketplace

# Configuring Data Marketplace

After Data Marketplace is activated, Data Marketplace is ready for use based on the out-of-the-box experience.
Optionally, Data Marketplace administrators can optimize the Data Marketplace experience.

What can you optimize?

- You can configure the actions available in the landing page.



- You can configure the filter tabs, filters and the scope of Data Marketplace.



- You can configure the actions available in the asset preview.

- You can configure the actions available in the No search results page.



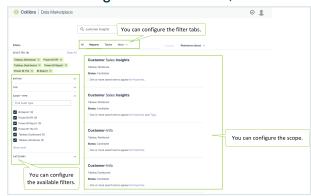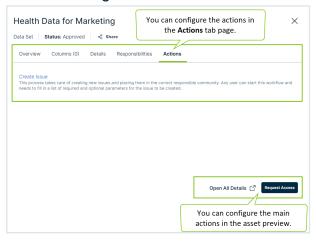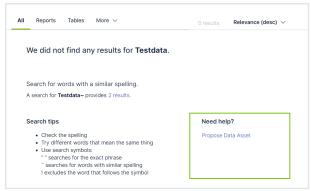| | Step | Description |
|---|---|---|
| 1 | Check that colleagues have the required permissions | Check that colleagues have the required permissions to use or configure Data Marketplace. |
| 2 | Optionally, change the Data Marketplace scope | By default, Data Marketplace finds assets of specific Data Catalog asset types, with any status, anywhere in your organizations.<br>You can change the scope based on asset type, organization and status. You can also include custom asset types and statuses in the scope. |
| 3 | Optionally, make actions available in Data Marketplace | By default, an asset can be shared with a colleague and opened in Collibra from the asset preview.<br>You can make additional actions available. Actions can be available in the Data Marketplace landing page, the No search results page, and the asset preview. |
| 4 | Optionally, configure the filters in Data Marketplace | By default, all out-of-the-box filters that apply to the search results are available to data consumers.<br>You can choose the filters that must be available in the **Filters** pane. |
| 5 | Optionally, configure the filter tabs in Data Marketplace | By default, all out-of-the-box filter tabs are available in the search results.<br>You can configure the filter sets you need. |

| | Step | Description |
|---|------|-------------|
| 6 | Optionally, optimize the search index | The Collibra search index settings have an impact on the Data Marketplace experience. You can use these settings to further optimize the experience. Note that changing these settings has an impact on Collibra as a whole, not only on Data Marketplace.<br><br>For information on the Data Marketplace-specific settings such as Search over relations, see Enable Data Marketplace in your environment. |

# The out-of-the-box Data Marketplace experience

After Data Marketplace is activated, Data Marketplace is ready for use based on the out-of-the-box configuration. This includes the Data Marketplace scope, filter tabs, filters, and actions.

## Scope

By default, Data Marketplace finds assets of the following Data Catalog asset types with any status, anywhere in the organization.

- Data Set
- Report
- Table
- Column
- Database View
- BI Data Set, BI Report
- Power BI Tile, Power BI Data Model, Power BI Dashboard, Power BI Report
- Tableau Dashboard, Tableau View, Tableau Story, Tableau Workbook, Tableau Worksheet

Administrators can configure the scope.

## Filter tabs

Data Marketplace shows the search results in tab pages, which are based on filter sets. Filter sets are typically based on an asset type or on a group of filters that are frequently

used and allow you to quickly apply a filter to the search results.



By default, Data Marketplace provides the following filter sets:

| Filter tab name | Description |
|---|---|
| All | No filter set is applied. All search results in the Data Marketplace scope to which you have View permissions are listed. |
| Tables | Only assets of the type Table or an out-of-the-box child asset type of Table are listed in the search results. |
| Data Sets | Only assets of the type Data Set or an out-of-the-box child asset type of Data Set are listed in the search results. |
| Reports | Only assets of the type Report or an out-of-the-box child asset type of Report are listed in the search results. |

Administrators can configure the filter tabs.

# Filters

By default, Data Marketplace provides the filters in the table below.

> Note   Some of these filters rely on relation path indexing.
> If your Collibra environment does not use one or more of the indexed relation types, the out-of-the-box filters based on relations do not work.

| Filter name | Description | Uses relations |
|---|---|---|
| Business Asset | Business Asset refers to the business assets that are related to the data, such as `Customer Phone Number` or `Customer Order Number`.<br><br>Tip   For a detailed narrative on the differences between business terms, data domains and data concepts, see Conceptual data layer versus the Business Glossary. | Yes |
| Certified | This filter limits your search results to data that has been certified, meaning approved, endorsed or guaranteed.<br><br>Tip   For complete information on certifying reports, see Trusted Business Reporting. | No |

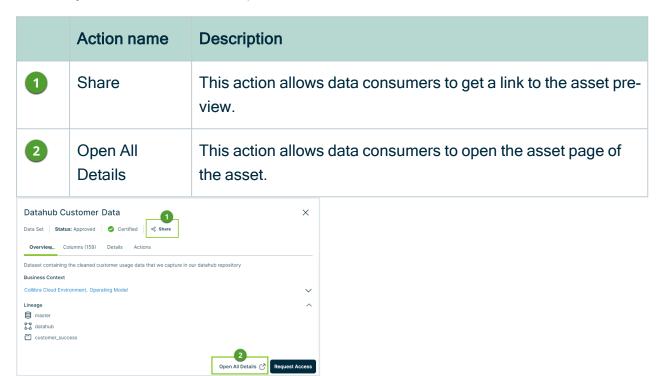| Filter name | Description | Uses relations |
|---|---|---|
| Data Domain | Data Domain refers to a high-level conceptual category to which the data belongs, such as `Customer` or `Order`.<br><br>Assets of this type:<br><br>• Are high-level, theoretical representations of your data.<br>• Represent the structure of concepts in data environments.<br>• Contain all the different nuances of corresponding business terms.<br><br>Important   To benefit from this, your data structure must follow the Data Stewardship metamodel. See Conceptual data layer. | Yes |

| Filter name | Description | Uses relations |
|---|---|---|
| Data Concept | Data Concept refers to the lower-level conceptual category to which the data belongs, such as `Phone Number` or `Country`.<br><br>Assets of this type:<br><br>• Are high-level, theoretical representations of your data.<br>• Represent an aspect of one or more data domains.<br>• Define a context-independent representation of the structure of an organization's data.<br><br>Important   To benefit from this, your data structure must follow the Data Stewardship metamodel. See Conceptual data layer. | Yes |
| Source System | Source System refers to the technology asset that stores the data, like HR Cloud. Via this filter, you can limit the search results to assets that exist in specified systems. | Yes |
| Data Source | Data source limits the search results to assets that exist in specified data sources. | Yes |
| Rating | Rating is a colleague's perception of the data's usefulness and quality. | No |
| Tags | Tag is a label that adds extra context to an asset. | No |

| Filter name | Description | Uses relations |
|---|---|---|
| Status | Status refers to which governance process stage the asset is in, like Approved or Candidate. | No |
| Asset Type | Asset Type refers to the kind of data you are looking for, like data set, report, or table. | No |

Administrators can configure the filters.

# Actions

Data Marketplace provides following actions to the data consumers. The following actions are always available in the asset preview:

| | Action name | Description |
|---|---|---|
| 1 | Share | This action allows data consumers to get a link to the asset preview. |
| 2 | Open All Details | This action allows data consumers to open the asset page of the asset. |

Datahub Customer Data

Data Set | **Status:** Approved | Certified | Share

Overview_ Columns (159) Details Actions

Dataset containing the cleaned customer usage data that we capture in our datahub repository

Business Context

Collibra Cloud Environment, Operating Model

Lineage

master
datahub
customer_success

Open All Details    Request Access

Administrators can make more actions available.

# Filtering and searching based on relations in Data Marketplace

Data Marketplace provides a unique search experience that can be based on relations between assets.

- Data Marketplace can provide filters to refine your search results based on specific relations.
  Next to filters based on asset type, status, tags, creator, modification and creation date, you can also filter based on specific relations. For example, you can filter the results for a specific data domain.
- Data Marketplace can find search results if they have a specific relation with the search term. A data set, for example, can show up in the results if the data set contains a column with a name that matches the searched term.

To make filtering and searching based on relations possible, Data Marketplace indexes specific relation paths, which are chains of relations. The out-of-the-box indexed relation paths are based on the Guided Stewardship metamodel.

## Relations in filters

You can refine the search results with filters that can take relations into account. By following a chain of relations, called a relation path, Data Marketplace can find assets that have a direct or indirect relation to a specific asset. For example, the Source system filter restricts the search results to data sets that are stored in a specific system.

The following table shows the relation paths of the out-of-the-box filters that use relations. If your Collibra environment does not use one or more of these relation types, the out-of-the-box filters based on relations do not work.

| Filter | Relation paths |
|---|---|
| Data Domain | **Relation path for Data Domain - Data Set**:<br><br>1. Data Set contains Column<br>2. Column is represented by Data Attribute<br>3. Data Attribute is part of Data Entity<br>4. Data Entity is classified by Data Domain<br>**Relation path for Data Domain - Table**:<br><br>1. Table contains Column<br>2. Column is represented by Data Attribute<br>3. Data Attribute is part of Data Entity<br>4. Data Entity is classified by Data Domain |
| Data Concept | **Relation path for Data Concept - Data Set**:<br><br>1. Data Set contains Column<br>2. Column is represented by Data Attribute<br>3. Data Attribute is classified by Data Concept<br>**Relation path for Data Concept - Table**:<br><br>1. Table contains Column<br>2. Column is represented by Data Attribute<br>3. Data Attribute is classified by Data Concept |
| Business Asset | **Relation path for Business Asset - Data Set**:<br><br>• Data Set is related to Business Asset<br>**Relation path for Business Asset - Report**:<br><br>• Data Set is related to Business Asset<br>**Relation path for Business Asset - Table**:<br><br>• Table is related to Business Asset<br>**Relation path for Business Asset - Table**:<br><br>1. Table contains Column<br>2. Column is represented by Business Asset |

| Filter | Relation paths |
|---|---|
| Source System | **Relation path for Source System - Data Set**:<br><br>• Data Set is implemented into Source System<br>**Relation path for Source System - Table**:<br><br>1. Table is part of Schema<br>2. Schema belongs to Technology Asset<br>3. Technology Asset is grouped by Source System |
| Data Source | **Relation path for Data Source - Table**:<br><br>• Table is part of Schema<br>  The data source is an attribute of the Schema asset. |

> **Tip**
> If you want to add filters based on the relations in your Collibra environment, ask your Customer Success Manager to contact the Product Management team.
>
> - For each custom filter, you need to provide the filter name and the relation path.
>   - Instead of spaces, use underscores (_) in the name of your custom filters. Data Marketplace replaces the underscores by spaces.
>   - Use meaningful names for each relation in the relation path.
>   - Don't use more than five relations in a relation path.
> - Custom filters and their relation paths are also used in the Search over relation feature.
> - Custom filters are not automatically reindexed, even if Consistency for Custom facets and Search over Relations is enabled. You have to reindex the relations manually on a regular basis.

# Relations in the search

The Search over relations feature is a Data Marketplace-specific beta feature that allows Data Marketplace to find assets based on specific relations. Search results include assets that have a direct match with the search term, but also assets that are related to the search term via specific relation types and based on the name of the asset.

> Example
> The "Customer Sales data" data set contains a column with the name "CustomerId".
> If Search over relations is enabled and you search for "CustomerId", then the
> "Customer Sales data" data set appears in the search results because the column
> name matches the search term and there is a Data set - Column relation between
> "CustomerId' and "Customer Sales data".

By default, the Search over relations feature uses the following relation paths.
If your Collibra environment does not use one or more of these relation types, the out-of-the-box Search over relations feature does not work.

| Relation paths |
| --- |
| Same relation paths as defined for the filters (see above). |
| **Relation path for Data Set - Column**:<br><br>• Data Set contains Column |

For information on how to enable the Search over relations feature, see Enable Data Marketplace.

> Tip
> • If you want to use specific relation paths to search over relations, ask your Customer Success Manager to contact the Product Management team.
>   ○ Complex relations are not supported as custom relation paths.
> • The Search over relation feature uses all indexed relations, meaning the ones linked to filters and the ones that are not linked to filters.
> • Custom relation paths are not automatically reindexed, even if Consistency for Custom facets and Search over Relations is enabled. You have to reindex the relations manually on a regular basis.

# Optimizing Data Marketplace

# Configure the Data Marketplace scope

As an administrator, you can customize the out-of-the-box scope to expose exactly what makes sense for your data consumers. Which data is relevant? Which data meets your quality or governance standards?

You can configure the scope based on asset type, status and organization. Data Marketplace takes all three scope definitions into account. If you choose to only include Reports as asset type, Approved as status and Marketing as organization, you only get results that meet all these conditions.

Example
In your environment, data consumers must request access on data set and report level. As a result, you only want to show Approved data sets and reports in Data Marketplace.
To do this you configure the following scope:

- Asset types: Data Set and Report.
- Status: Approved.
- Organization: all

If you also enable Search over relations, then data consumers can still use column names as search terms. In that case, Data Marketplace shows the related data set and reports in the search results instead of the Column assets.

Example
In your environment, you keep track of approved marketing data in a specific community "Marketing-Approved". As a result, you only want to show data of this community in Data Marketplace.
To do this you configure the following scope:

- Asset types: all
- Status: all
- Community: Marketing-Approved.

Note   Data Marketplace cannot search Glossary domains and Business Term assets.

> Important
> - If you create a new asset type, status, community, or domain in Collibra, and you want it to be found in Data Marketplace, you have to select it in the Data Marketplace scope.
> - If you delete an asset type, status, community, or domain in Collibra, and it was part of the Data Marketplace scope, it will no longer be visible in the Data Marketplace scope list but will still be active in Data Marketplace.
>   Before deleting an asset type, status, community, or domain, remove it from the scope in Data Marketplace.

## Prerequisites

- Data Marketplace is enabled.

- You are an administrator in Data Marketplace.

## Steps

1. On the main menu, click ⠿, then ⚙ **Settings**.

   » The Collibra settings page opens.
2. In the **Data Marketplace (Beta Feature)** section, click **Scope**.
3. If you want to limit the search results to assets of specific asset types, select the required asset types in the **Asset types** tab page.

   > Tip
   > ○ Only Data Catalog asset types can be found via Data Marketplace.
   > ○ You can only select custom asset types in the **Asset Types** tab page if they are a child asset type of a Data Catalog asset type.

4. If you want to limit the search results to assets with specific statuses, select the required statuses in the **Statuses**tab page.
   If you do not define a scope, data consumers can find assets with any statuses via Data Marketplace, including Candidate, Pending, and so on.
5. If you want to limit the search results to assets in specific communities and domains, select the required communities and domains in the **Organization**tab page.
   If you do not define a scope, Data Marketplace searches all communities and

domains .

6. Click **Save**.

## What's next?

Make actions available in Data Marketplace
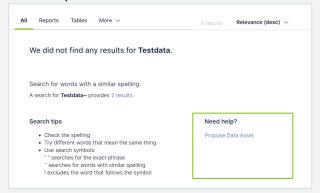
# Make actions available in Data Marketplace

Data Marketplace provides some out-of-the-box actions to the data consumers. As an administrator, you can add actions to and remove them from Data Marketplace. These actions start workflows. Actions can be made available in Data Marketplace:

- At the bottom of the Data Marketplace landing page.
- In the **Need help?** section if the search did not return any results.
- At the bottom of an asset preview.
- In the **Actions** tab page in an asset preview.
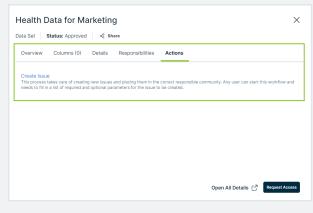
Example
You can configure Data Marketplace to show:

- The Issue creation workflow in the Data Marketplace landing page.
- The Propose new data set workflow in the No search results page.

| All | Reports | Tables | More ∨ | | 0 results | Relevance (desc) ∨ |

We did not find any results for **Testdata**.

Search for words with a similar spelling.
A search for **Testdata~** provides 2 results.

**Search tips**

- Check the spelling
- Try different words that mean the same thing
- Use search symbols:
  " " searches for the exact phrase
  ~ searches for words with similar spelling
  ! excludes the word that follows the symbol

**Need help?**

Propose Data Asset

- The Request Access workflow at the bottom of an asset preview.

Open All Details ⎋    Request Access

- The Create Issue workflow in the **Actions** tab page of an asset preview.

Health Data for Marketing                    ✕
Data Set | **Status:** Approved | ⌁ Share

Overview    Columns (0)    Details    Responsibilities    **Actions**

Create Issue
This process takes care of creating new issues and placing them in the correct responsible community. Any user can start this workflow and needs to fill in a list of required and optional parameters for the issue to be created.

Open All Details ⎋   Request Access

Tip
Data Marketplace shows the **Start Label** of the workflow as the action's button.

## Prerequisites

- Data Marketplace is enabled.

- You are an administrator in Data Marketplace.

## Steps

1. On the main menu, click ⠿, then ⚙ **Settings**.

   » The Collibra settings page opens.
2. In the **Data Marketplace (Beta Feature)** section, click **Actions**.
3. In the **Landing Page Actions** tab page, select the actions that must appear at the bottom of the landing page.

   > **Important**
   > You can only select global workflows for which the option **Show in global create** is selected.

4. In the **No Search Results Actions** tab page, select the actions that must appear in the **Need help?** section if the search did not return any results.

   > **Important**
   > You can only select global workflows for which the option **Show in global create** is selected.

5. In the **Main Asset Preview Actions** tab page, select the actions that must appear in the asset preview.
   a. In **Primary Action**, select the workflow that must start when the primary action button is clicked in the asset preview.
   b. If you want to add a link to an external system from the main actions in the asset preview, enter the **Action Label** and the **System URL** in the **Link Action** section.
      In the below example, **Request Access** is the primary action and **Log IT ticket** is the link action. You cannot customize the **Open All Details** action.

   

   > **Important**
   > ■ You can only select workflows that apply to assets.

> ■ If a data consumer does not have access to the workflow or the workflow does not apply to the asset type, the action will not appear.

6. In the **All Asset Preview Actions** tab page, select the actions that must appear in the **Actions** tab page in the asset preview.

> Important
> You can only select workflows that apply to assets.

7. Click **Save**.

> Tip
> If a data consumer does not have access to a workflow, the action will not appear in Data Marketplace.

## What's next?

Configure the filters in Data Marketplace

## Configure the filters in Data Marketplace

By default, out-of-the-box filters are available in Data Marketplace.
However, administrators can change the filters and the order in which they should appear.

See an example

> Example
> You can ensure the following filters are displayed in this specific order:
>
> - Asset Type
> - Certified
> - Status
> - Data Domain

> Tip
> Only filters that can apply to the search results will be displayed to data consumers.

## Prerequisites

- [Data Marketplace is enabled](#).

- You are an [administrator in Data Marketplace](#).

## Steps

1. On the main menu, click ⊞, then ⚙ **Settings**.
   » The [Collibra settings page](#) opens.
2. In the **Data Marketplace (Beta Feature)** section, click **Filters**.
   » The Management tab page appears and shows the filters that are available in Data Marketplace.

   - The list shows all available filters.

     > Note   If you want to have more filters, also called custom facets, ask your Customer Success Manager to contact the Product Management team.

   - The **Filters** pane in Data Marketplace uses the order of this list.
3. To show a filter, select the **Is Visible** checkbox.
   To hide a filter, clear the **Is Visible** checkbox.
   To move a filter, use the ✛ icon to drag and drop it in the desired location.
4. Click **Save**.

## What's next?

[Configure filter tabs in Data Marketplace](#)

# Configure the filter tabs in Data Marketplace

By default, out-of-the-box filter tabs are available. As an administrator, you can manage the filter tabs. For example, you can create, edit and delete them, make them visible and change their order.

> **Important**
> - Currently, you cannot create a filter set that includes custom filters.
> - Once a filter set is saved, changes to the scope no longer have an impact on the filter set.
>   This means that:
>
>   - Administrators must configure the scope before creating filter tabs.
>   - After editing the scope, administrators should also edit the filter tabs.
>   - Out-of-the-box filter tabs are not impacted by the scope configuration.
>
>   This is a known issue that is planned to be solved in one of the next versions of Data Marketplace.

## Prerequisites

- Data Marketplace is enabled.

- You are an administrator in Data Marketplace.

## Create a new filter tab

1. Select the filters:
   a. On the main menu, click ⠿, then ♛ **Data Marketplace**.
      » The Data Marketplace landing page opens.

   b. Search for * to see all available filters in the search results.
      » The search results appear.
   c. Select the filters you want to save in a filter tab.

      > **Tip**
      > - The filter tab saves the sort order.
      > - The filter tab doesn't save your search terms.

2. Click the arrow next to **More** and then **Save new tab**.

&raquo; The **Save Current Filters** dialog box appears.

3. Enter a name for the new filter tab.
4. Click **Save Changes**.

&raquo; The filter tab is saved and available. Data consumers can select it from the **More** drop-down list.

## Edit or delete filter tabs

1. On the main menu, click ⠿, then 🛒 **Data Marketplace**.

&raquo; The Data Marketplace landing page opens.

2. Search for something, for example * to see all available filters in the search results.

&raquo; The search results appear.

3. Click the arrow next to **More**.
4. Click **Edit filter tabs**.

&raquo; The **Edit Filter Tabs** dialog box appears.

&cir; In the **Arrange Tabs** tab page, you can define the filter tabs that must be visible and the order in which they appear.

&cir; In the **Edit Tabs** tab page, you can change the name of a filter tab or delete it.

> Tip
> - The filter tabs appear in the order of this list.
> - You can't hide the **All** filter tab.
> - You can't rename or delete the out-of-the-box filter tabs, but you can hide them if you don't want to use them.
> - You can hide a filter before you delete it to check the reaction of the data consumers.

5. If required, in the **Arrange Tabs** tab page click the 👁 icon to hide a filter tab.

&raquo; The 🚫 icon indicates the filter does not appear.

6. If required, edit the order of the filter tabs by dragging and dropping the ↕ icon.
7. If required, delete a filter tab by selecting the **Edit Tabs** tab and clicking 🗑 .
8. If required, change the name of a filter tab by selecting the **Edit Tabs** tab and clicking ✏ .
9. Click **Save Changes**.

&raquo; The filter tabs are saved.

## What's next?

Continue the Data Marketplace configuration

# Reindex Data Marketplace relations

Some relations in Collibra are not automatically reindexed for Data Marketplace. This can cause inconsistencies between Collibra and Data Marketplace.
Therefore, you can reindex the relations manually.
The following relations are not indexed automatically:

- Out-of-the-box filters based on relations and search relation paths, if Consistency for Custom facets and Search over Relations setting is not enabled.
- Relations for custom filters and custom search relation paths.

For more information, see Filtering and searching based on relations in Data Marketplace.

## Prerequisites

You have a global role that has the System administration global permission.

## Steps

1. On the main menu, click ⠿, then ⚙ **Settings**.
2. In the **General** section, click **System**.
   - » The **Information** page opens.
3. In the **Relations index** section, click **Rebuild Relations index**.
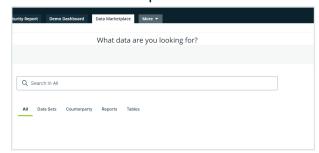
# Add Data Marketplace to a dashboard

You can add the Data Marketplace search field to a dashboard to allow easy access to Data Marketplace.

## Steps

1. On the main menu, click ⚙.
2. If required, switch to the relevant dashboard.
3. On the dashboard bar, click ⋮ → **Edit widgets**.
4. Navigate to the Text widget and drag it onto the dashboard.

   You can browse through the widgets by using the page numbers below them.
5. In the **Content** tab page, click ‹›, to edit the HTML source code.

   » The **Source Code** dialog box opens.
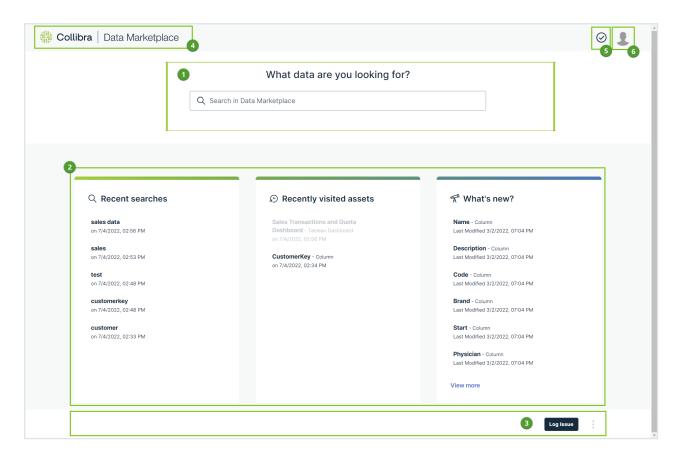6. Add the following code to the dialog box:
   ```
   <div>
   <h1 style="text-align: center; font-size: 24px; font-weight:
   400; font-family: InterVariable,Inter,Open Sans, Helvetica
   Neue,Helvetica,Arial,sans-serif;">What data are you looking
   for?</h1>
   <p><iframe style="width: 100%; height: 60vh; border: 0;"
   src="/apps/shop4data/dashboard-widget"
   allowfullscreen="true"></iframe></p>
   </div>
   ```
7. In the **Widget** tab page, clear the checkboxes **Show title above the widget** and **Show description above the widget**.
8. Click **Save**.
9. Click **Close & save edits**.

   » The Data Marketplace search field and filter tabs are available in the dashboard.

# The Data Marketplace landing page

The Data Marketplace landing page is your starting point for finding the data you need.
The page consists of multiple sections.

| | Section | Description |
|---|---|---|
| **1** | Search | The search field allows you to find any Data Catalog asset, such as Data Set, Schema, Table, Column, Tableau Workbook, BI report, and Power BI Dashboard within the Data Marketplace scope. |
| | | • If search suggestions are enabled, Data Marketplace provides suggestions while you are typing. Click one of the suggestions to use it or press `Enter` to ignore the suggestions.<br>• If search suggestions are not enabled, the search starts automatically after a few seconds. |

| | Section | Description |
|---|---|---|
| **2** | Cards | Cards allow you to navigate quickly to data that might interest you. The available cards depend on whether Analytics are enabled in your environment. |

* If analytics are enabled, you have access to **Popular Data**, **Recent Data Visits**, and **Recommended Data**. View more information.

| | |
|---|---|
| **Popular Data** | This card shows the assets your colleagues visited the most in the last 30 days in Data Marketplace, Collibra, Collibra for Desktop, and Collibra for Mobile and that are in the scope of Data Marketplace. . <br> ○ Data Marketplace checks the analytics once a day. <br> ○ Click **Show more** to view more popular assets. <br> ○ Loading this card can take a few seconds. |
| **My Recent Data Visits** | This card shows the assets for which you opened the preview in Data Marketplace. <br> ○ Click **Show more** to view more recently visited assets. <br> ○ Click ✕ to remove the visited asset. <br><br> ⟳ **My Recent Data Visits** <br> Customer Analytics <br> Data Set     Visited about 1 hour ago ✕ <br><br> ○ An asset is grayed out when the asset is no longer available in Data Marketplace. Usually this is due to a change in the Data Marketplace scope. <br><br> DueDate <br> Column     Visited about 20 hours ago <br> Order Quantity Per Customer <br> Tableau Worksheet     Visited about 20 hours ago |
| **Recommended Data** | This card shows data based on the data you search for and visit in Data Marketplace, Collibra, Collibra for Desktop, and Collibra for Mobile and that are in the scope of Data Marketplace. <br> ○ Data Marketplace checks the analytics once a day. <br> ○ Loading this card can take a few seconds. |

> Tip  No recommendations can be given to new users because no analytics are available yet. Recommendations will be provided after one day.

| | Section | Description |
|---|---|---|
| | | • If analytics are not enabled, you have access to **Recent searches**, **Recent data visits**, and **What's new**. View more information. |

| Recent searches | This card shows the recent searches you performed in Data Marketplace. |
|---|---|
| | ○ Click **Show more** to view more recent search actions. |
| | ○ Click ⊼ to pin the search action to your list. |
| | ○ Click ✕ to remove the recent search action. |
| | 🔍 **Recent searches** <br><br> **customer insights** <br> on 7/5/2022, 09:15 AM     ⊼   ✕ |

| My Recent Data Visits | This card shows the assets for which you opened the preview in Data Marketplace. |
|---|---|
| | ○ Click **Show more** to view more recently visited assets. |
| | ○ Click ✕ to remove the visited asset. |
| | 🕐 **My Recent Data Visits** <br><br> **Customer Analytics** <br> Data Set     Visited about 1 hour ago   ✕ |
| | ○ An asset is grayed out when the asset is no longer available in Data Marketplace. Usually this is due to a change in the Data Marketplace scope. |
| | **DueDate** <br> Column     Visited about 20 hours ago <br><br> **Order Quantity Per Customer** <br> Tableau Worksheet     Visited about 20 hours ago |

| What's new? | This card shows all assets that have recently been updated in Data Catalog and that are in scope of Data Marketplace. Click **Show more** to open the search results page with all updated assets. |
|---|---|

| # | Section | Description |
|---|---|---|
| 3 | Landing page actions | Your Data Marketplace administrator can make buttons available in this section. |
| 4 | Back to landing page | Click this button to bring you back to the Data Marketplace landing page. |
| 5 | Notifications | Click this button to get an overview of the messages you received from Data Marketplace. |

| | Section | Description |
|---|---------|-------------|
| **6** | Avatar | Click your avatar to:<br><br>• Open your profile details.<br>• Enable or disable suggestions.<br>• View recommendations. This option is only available if analytics are enabled in your environment.<br>• Sign off from Data Marketplace.<br>• View the version of Data Marketplace and Collibra. |

# Searching in Data Marketplace

The purpose of Data Marketplace is to search for and interact with the data that is available to you. Searching and refining your search results are therefore the key features. Data Marketplace provides multiple ways for you to get to the Data Catalog information you need.
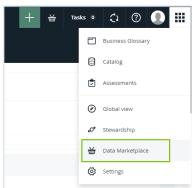
# Search in Data Marketplace

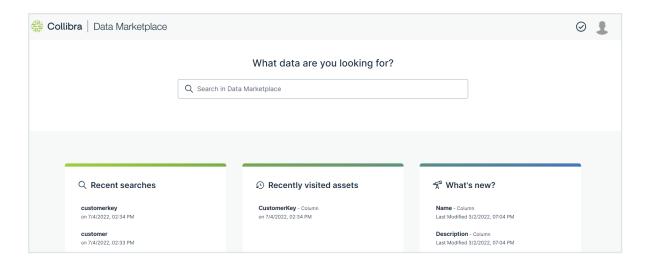1. Open a dashboard that includes the Data Marketplace search field.
   or
   Open Data Marketplace from the ⠿ icon.
   On the main menu, click ⠿, then ⩊ **Data Marketplace**.



   If Data Marketplace does not appear in the application selector after clicking ⠿, the feature is not enabled or you don't have the required permissions.
   » The Data Marketplace landing page opens.

2.  In the search field, enter your search terms.
    ○  Optionally, use wildcards and symbols.
    ○  If search suggestions are not enabled, the search starts automatically after a few seconds.
    ○  If search suggestions are enabled, suggestions appear as you type. Click one of the suggestions to use it or press Enter to ignore the suggestions.



    ○  Click ✕ to remove your search terms.
    »  The search results page opens.

# Wildcards and symbols for searching in Data Marketplace

In the Data Marketplace search field, you can enter your search terms and use wildcards and symbols.

- If you want to search for all assets, type * in the search field.
- If you want to find something, but are not quite sure of the name, you can use the fol-low wildcards in your search terms:

| Wildcard | Description | Example |
|---|---|---|
| ? | A question mark replaces any single character. | Search term: *?ar* Possible results: Car, Bar, Tar. |
| * | An asterisk replaces any string of characters. | Search term: *C\*r* Possible results: Car, Crosshair, and so on. |
| ~ | A tilde represents a fuzzy search, words with spelling similar to the text you searched for.<br><br>Important   Add the tilde sign at the end of your search term. If you add it in front of the term, you receive an error message. | Search term: *Owi~* Possible results: Ozi, Zowie, Bowie, and so on. |
| ! | An exclamation mark excludes words. | Search term: *!David Bowie* Possible results: Lester Bowie, Joe Bowie, and so on. |

- If you want to find a literal piece of text, you can add the terms between double quotation marks.

  Don't forget to close the quotation marks. If you add only one, you receive an error message.

  > Example   By typing *"Tax Statement"* you can find Tax Statement 2017, but not Tax or Statement.

> Important
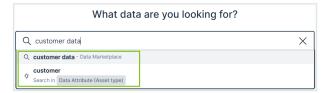> The following characters are not supported in the search field:
> / slash
> : colon
> ^ circumflex
> () parenthesis
> [] square brackets

> {} curly brackets

# About search suggestions in Data Marketplace

Search suggestions guide you to the information you might be looking for.
If you enabled search suggestions, Data Marketplace provides suggestions as you type.
You can click one of the suggestions to use it. The result of clicking a suggestion depends
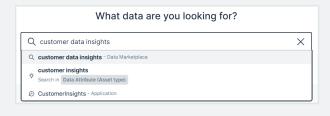on the suggestion type.



## Suggestion types

The icon in front of a suggestion indicates the suggestion type.

| Icon | Meaning | Result when you click the suggestion |
|------|---------|--------------------------------------|
| 🔍 | Recent search | The suggestion shows a previous search action.<br>Clicking the suggestion launches the suggested search action. |
| 🕑 | Recently visited asset | This suggestion shows an asset that you recently previewed.<br>Clicking this suggestion opens the preview of the suggested asset. |
| 💡 | Smart filter | This suggestion provides possible filters you can apply based on your search terms. The smart filter suggestions recognize search terms linked to asset type, system, status, and certification.<br>Clicking this suggestion performs the search and applies the suggested filters. |

> **Example**
> The search terms *customer data insights* result in three suggestions.
>
> - The first suggestion is based on a recent search action. Clicking this suggestion launches the search for *customer data insights* in Data Marketplace again.
> - The second suggestion provides the possibility to search for the terms *customer insights* with the filter *Asset Type* set to *Data Attribute*.
> - The last suggestion points to the Customerinsights application asset that you previewed before. Clicking the suggestion opens the asset preview again.

# Enable or disable search suggestions

Search suggestions guide you to the information you might be looking for. You can enable or disable the search suggestions.
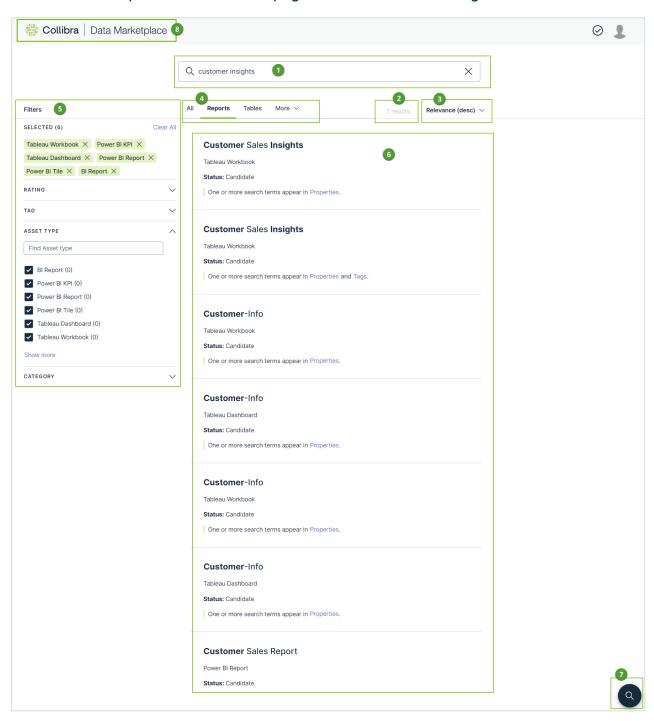
## Steps

1. On the main menu, click ⚏, then 🛒 **Data Marketplace**.
   » The Data Marketplace landing page opens.

2. Click your avatar icon in the upper right corner.
   » The avatar menu opens.

3. Use the **Suggestions** option to enable or disable the search suggestions.

| | |
|---|---|
| Suggestions ⬤◯ | Suggestions are enabled for you. |
| Suggestions ◯ | Suggestions are disabled for you. |

4. Click outside the avatar menu.

# The Data Marketplace search results page

The Data Marketplace search results page consists of the following sections and features.

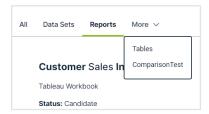| | Section | Description |
|---|---------|-------------|
| **1** | Search | The search field contains the search terms you used. You can perform a new search by changing the search terms. |
| **2** | Number of search res- ults | This section shows the number of results for your search action. |
| **3** | Sort options | The sort options allow you to sort your search results. |
| **4** | Filter tabs | Filter tabs allow you to only show assets of a specific type or apply a group of filters in one action. See also Filter your search results. The available filter tabs are configured by the Data Marketplace administrator. |
| **5** | Filters pane | Filters allow you to refine your search results further. The available filters depend on your search results and on how the Data Marketplace administrator has configured the Filters pane. |
| **6** | Search results | This section shows a list of results that match your search action. For each search result, a summary is available. Click the result to open the details in a preview. The search results depend on multiple factors: <br><br> • The Data Marketplace scope that was configured by the administrator. <br> • Whether you have view permission on the asset. <br> • The filter tab that is selected. <br> • Filters that are applied. |
| **7** | Back to search | This icon guides you to the top of the page. |
| **8** | Back to land- ing page | This button brings you back to the Data Marketplace landing page. |

# Available filter options

You can filter your search results in the following ways:

- By selecting a filter tab
- Via the Filters pane

The selected filters are visible in a specific area in the **Filters** pane.

## Filter tabs

Via a filter tab, you apply a filter set. Filter sets are typically based on an asset type or on a group of filters that are frequently used and allow you to quickly apply a filter to the search results. Via the **More** drop-down list, you have access to additional filter sets.



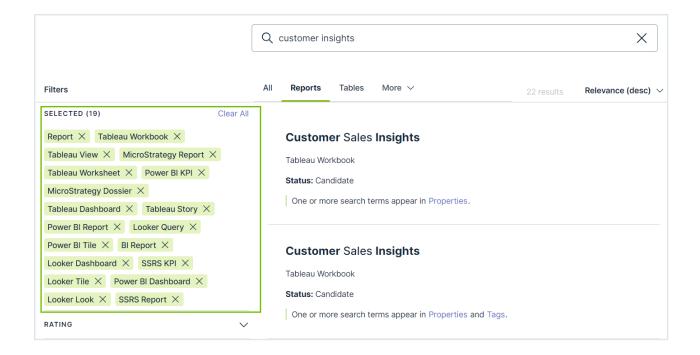Administrators can configure filter tabs in Data Marketplace.

## Filters in the Filters pane

In the **Filters** pane, only filters that are selected by the Data Marketplace administrator and that are relevant to the search results are available.
If you select filters, the **Filters** pane is automatically updated with the filters that are relevant to the refined search results. If the filter does not contain any values, the filter is not visible.

## Selected filters in the Filters pane

In the **Selected** area in the **Filters** pane, you get an overview of all the filters that are applied to the search results. This includes the filters applied via a filter tab as well as the filters applied via the **Filters** pane.
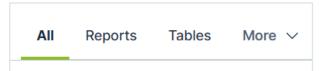
# Filter your search results

You can use filters to refine the search results. You can apply filters via filter tabs and via the **Filters** pane.

Once filters are applied, you can remove a filter or start over by clearing all filters in one action.
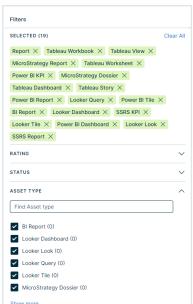
## Steps

1.  Search in Data Marketplace.

    » In the search results page, the **All** filter tab is selected. This tab page shows all search results within the Data Marketplace scope.

2.  If required, select another filter tab to apply a filter set.

    Only a few filter tabs are displayed. Click **More** to see extra filter sets.

    

3.  If required, use the **Filters** pane to refine your search results further.

    The available filters depend on the search results and the configuration of the filters.

» All selected filter criteria are visible in the **Selected** section in the **Filters** pane.



4. If required, use the **Selected** section in the **Filters** pane to remove applied filters.
   - To remove one filter, click ✕ next to the applied filter.
   - To remove all filters in one action, click **Clear All**.
5. Analyze the remaining search results and click the title of a search result to open the asset preview.

# Sort your search results

You can sort search results by the following criteria and in ascending or descending order:

- Relevance.
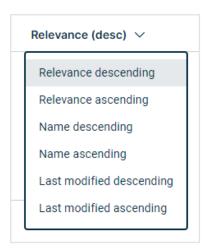
  By default, search results are sorted in order of descending relevance.
- Name.

  The search results are sorted based on the asset name.
- Last modified.

  The search results are sorted based on when the asset was last modified.

# View recommendations

Recommendations are suggestions made by Data Marketplace and based on the data that you and your colleagues are searching for, visiting and updating. The following recommendations are available:

- **Most Visited Assets** shows the data that your colleagues visited the most in the last 30 days in Data Marketplace, Collibra, Collibra for Desktop, and Collibra for Mobile and that is in the scope of Data Marketplace.

- **Last Updated** shows assets that have recently been updated in Data Catalog and that are in the scope of Data Marketplace.

- **Recommendations** shows assets based on the data you searched for and visited in Data Marketplace, Collibra, Collibra for Desktop, and Collibra for Mobile and that are in the scope of Data Marketplace.

## Prerequisites

Recommendations are only available when analytics are enabled in your environment.

## Steps

1. On the main menu, click  ⠿ , then ⛺ **Data Marketplace**.
   » The Data Marketplace landing page opens.

2. Click your avatar and select **Recommendations**.
   » The **Discover** page opens.
3. Click the tab you need.
4. Click the asset you are interested in to open its preview.
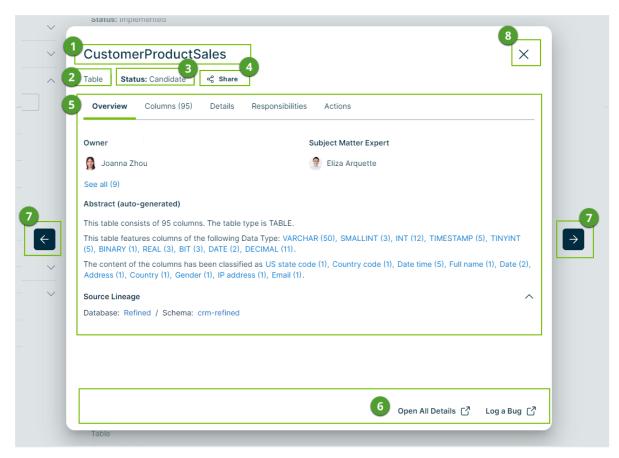
# Previewing assets in Data Marketplace

After searching, you analyze the search results to find what you want. When you click a search result, its details open in an asset preview dialog box.

# The Data Marketplace asset preview

After searching, you analyze the search results to find what you want. When you click a search result, its details open in an asset preview dialog box.

- An asset preview contains multiple sections.



Get more information about the sections

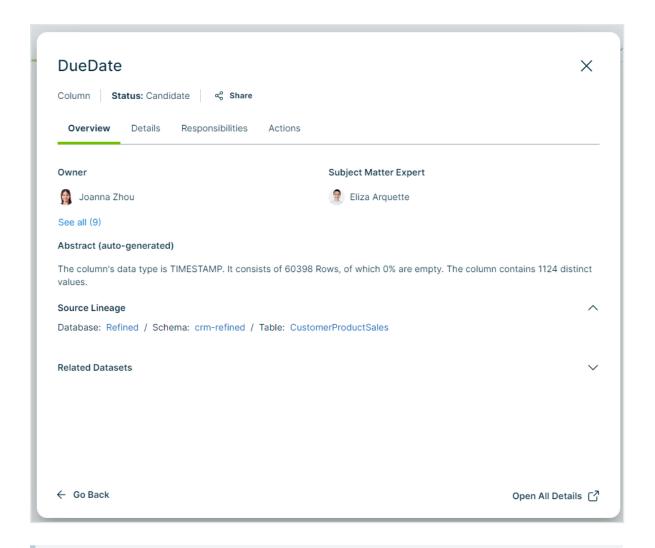| | Section | Description |
|---|---|---|
| 1 | Asset title | The name of the asset. |
| 2 | Asset type | The kind of asset you are previewing. |
| 3 | Asset status | The status of the asset in your governance cycle. |
| 4 | Share | Click **Share** to collect the URL to the asset and share it with colleagues. |

| | Section | Description |
|---|---|---|
| **5** | Tabs and tab content area | The available tabs depend on the asset type. The most common tabs are:<br>◦ **Overview**: this tab page includes, for example:<br>  ▪ The description of the asset.<br>  ▪ All tags applied to the asset.<br>  ▪ All responsibilities configured for the asset.<br>◦ **Details**: this tab page includes, for example ratings and reviews of the asset.<br>◦ **Responsibilities**: this tab page includes the roles assigned to the asset. Inherited responsibilities are not shown in Data Marketplace.<br>◦ **Actions**: this tab pages includes all actions you can perform from this asset. They are the secondary actions. |
| **6** | Main actions | This section shows the actions that administrator configured to be the most important actions.<br>By default, you can always open the asset page in Collibra via **Open All Details**. |
| **7** | Navigation area | Click the arrows to navigate to the previous or to the next search result. |
| **8** | Close | Click **Close** to closes the asset preview. |

- The asset preview content depends on the asset type.
  Most asset types are shown using the same layout, but some asset types follow a specific asset type template.

  Information for a Data Set asset:

Information for a Column asset:

> **Note**
> If, in the asset type template, a field is defined but the information is not available for the selected asset, the field will not be displayed in the preview. This means you never see empty fields in the preview.
> For example: If the data has not been profiled, no profiling details are visible.
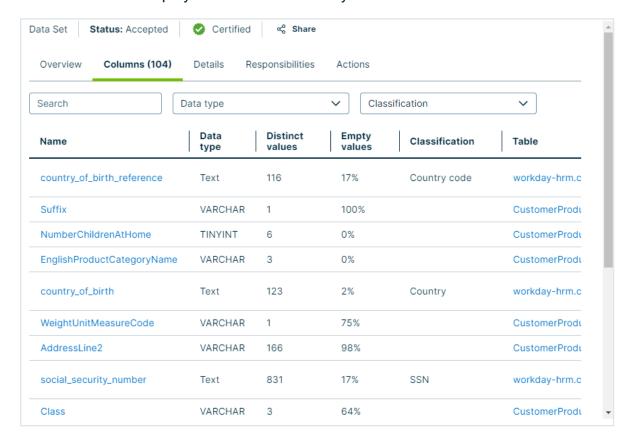
- From an asset preview, you can start multiple actions.

  - Main actions:
    - Main actions are located at the bottom of the asset preview.
    - By default, the **Open All Details** button is available. This opens the asset page in Collibra.
    - The Data Marketplace administrator can add more actions to this section.

- Secondary actions:
  - Secondary actions are all grouped in the **Actions** tab page of the asset preview.
  - The Data Marketplace administrator configures the secondary actions.
- From the title bar section in the preview, you can share an asset link with someone else.

- Column and Table assets show a summary of the asset information that is generated by Data Marketplace based on the asset details.

> **Abstract (auto-generated)**
>
> This table consists of 95 columns. The table type is TABLE.
>
> This table features columns of the following Data Type: VARCHAR (50), SMALLINT (3), INT (12), TIMESTAMP (5), TINYINT (5), BINARY (1), REAL (3), BIT (3), DATE (2), DECIMAL (11).
>
> The content of the columns has been classified as US state code (1), Country code (1), Date time (5), Full name (1), Date (2), Address (1), Country (1), Gender (1), IP address (1), Email (1).

- Some asset previews contain a **Columns** tab. This tab page shows the available columns in the asset.
  The columns are displayed in a table in which you can search.

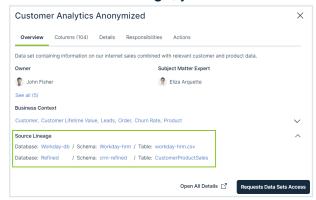| Name | Data type | Distinct values | Empty values | Classification | Table |
|---|---|---|---|---|---|
| country_of_birth_reference | Text | 116 | 17% | Country code | workday-hrm.c |
| Suffix | VARCHAR | 1 | 100% | | CustomerProdu |
| NumberChildrenAtHome | TINYINT | 6 | 0% | | CustomerProdu |
| EnglishProductCategoryName | VARCHAR | 3 | 0% | | CustomerProdu |
| country_of_birth | Text | 123 | 2% | Country | workday-hrm.c |
| WeightUnitMeasureCode | VARCHAR | 1 | 75% | | CustomerProdu |
| AddressLine2 | VARCHAR | 166 | 98% | | CustomerProdu |
| social_security_number | Text | 831 | 17% | SSN | workday-hrm.c |
| Class | VARCHAR | 3 | 64% | | CustomerProdu |

# Find the source of a data set or report

You can find the data that was used to source a data set or report.

## Steps

1. Search for an asset and click it.
   - » Its details appear in the asset preview.
2. In the asset preview, click the **Overview** tab.
3. Under **Source Lineage**, you can see the related data sources.



4. Click the source system name.
   - » The source system details open in the asset preview.

# Data Marketplace troubleshooting

# You see fewer search results in the **All** filter tab compared to other filter tabs

Filter tabs filter the search results based on filter sets that are available out-of-the-box or are created and saved by an administrator. Administrators can also configure the scope of Data Marketplace. The scope determines which asset types, statuses, communities and domains you can search.

Once a filter set is saved, changes to the scope no longer have an impact on the filter set. This means that:

- Administrators must configure the scope before creating filter tabs.
- After editing the scope, administrators should also edit the filter tabs.
- Out-of-the-box filter tabs are not impacted by the scope configuration.

This is a known issue that is planned to be solved in one of the next versions of Data Marketplace.

# You don't see the actions in the asset preview after you enabled them in the Data Marketplace settings

Actions are based on workflows and their configuration. Workflows can be configured to only be available for certain asset types or users, and they can be paused. If you do not see a specific action, the workflow configuration may give you an indication on the reason.

# You receive an error when you want to save a filter tab

When you save a filter tab, you receive the following error: Unexpected error has occurred. The most likely reason for this error is that the filter tab you tried to save includes custom filters. Currently, you cannot create a filter set that includes custom filters.