# Integrating CDP Data Catalog with AWS Glue Data Catalog (Preview)

Date published: 2021-08-09
Date modified: 2021-12-08

# Legal Notice

# Contents

# Integrating CDP Data Catalog with AWS Glue Data Catalog

Integrating CDP Data Catalog with AWS Glue Catalog enables the users to browse and discover data as well as register data into SDX (via metadata translation or copy), so that it can be used with Data Hubs and other relevant experiences. While using AWS Glue in Data Catalog, you will be able to experience a complete snapshot metadata view, along with other visible attributes that can power your data governance capabilities.

## How integration works

Assuming that the SDX is running in the users' AWS account (that contains the same AWS account which has GlueDataCatalog and the data that has to be discovered), the credentials with the ExternalDataDiscoveryService (which is hosted in SDX) must be shared, so that these two entities can interact with each other. These credentials are used to launch SDX and other workload clusters on the users' AWS account.

Prerequisites:

- You must have full access to AWS Glue Catalog and also have access to the EMR cluster's Hive Metastore instance.
- You must set up the CDP. For more information, see [Getting Started with CDP](#).
- You must have access to your AWS IT Admin and CDP Admin user credentials, which is required to enable CDP to access AWS/EMR managed data in CDP.

**Note:** AWS policies are managed by AWS IAM and the AWS roles are added to the CDP Management Console. Refer to [Using Instance Profile in AWS](#) and [Using credentials in Management Console](#). For more information about the AWS access, see [AWS Environments](#).

# Setting up AWS Glue Catalog with CDP Data Catalog

Follow these set up tasks to map your CDP Data Catalog instance with AWS Glue Catalog:

1. Enable the entitlement for your Data Catalog instance by running the following command on your CDP environment:

For example: *$ cdp coreadmin grant-entitlement --entitlement-name DATA_CATALOG_ENABLE_AWS_GLUE --account-id {account_id}*

2. You must add relevant permissions in the corresponding AWS account:

   a. Include permission to access Glue Catalog service by editing the policy accordingly.

   Make a note of the **Assumer Instance Profile** role that you intend to use and include full access authorization for AWS Glue.

   Assumer Instance profile role is the one which you select in the available public cloud documentation listed under **step 10**. For more information, see Register a CDP environment.

   Use the following screenshots as a reference to complete the set up.



**Note:** For **Role ARN** and **Instance Profile ARNs**, you must include the appropriate account number and role respectively.

A policy defines the AWS permissions that you can assign to a user, group, or role. You can create and edit a policy in the visual editor and using JSON. Learn more

**Visual editor** | JSON

Import managed policy

Expand all | Collapse all

▶ **STS** (2 actions)                    Clone | Remove

▶ **Glue** (All actions)                 Clone | Remove

⊕ **Add additional permissions**

---

A policy defines the AWS permissions that you can assign to a user, group, or role. You can create and edit a policy in the visual editor and using JSON. Learn more

**Visual editor** | JSON

Import managed policy

Expand all | Collapse all

▶ **STS** (2 actions)                    Clone | Remove

▶ **Glue** (All actions)                 Clone | Remove

▼ Select a service                       Clone | Remove

▼ **Service**  **Select a service below**          Enter service manually
close

🔍 GLuel

Glue ⓘ

**Actions**    Choose a service before defining actions

**Resources**    Choose actions before applying resources

**Request conditions**    Choose actions before specifying conditions

⊕ **Add additional permissions**

b. Search for the role attached to the **Instance Profile** of the CDP environment.

Use the **Instance Profile** that you have configured above with Glue related policy in your AWS Environment creation command.

Use the following examples to setup AWS environment and AWS data lake as part of the Glue setup:

```
cdp environments create-aws-environment --profile default --cli-input-json '
{"environmentName":"ab-ds-cli-7321",
 "credentialName":"cd2d-1234",
 "Region":"us-region-2",
         "securityAccess":{-insert the value--"},
 "Authentication":{---insert the value---"},

"logStorage":{"storageLocationBase":"s3a://demo-e2e-test-state-bucket/ab-ds-cli-7321/logs","instanceProfil
e":"arn:aws:iam::<xxxxxxxxxxx>:instance-profile/<role-name>"},
 "vpcId":"vpc-0123456",
 "subnetIds":["subnet-04fe923b902aa5cf2","subnet-099c7a631f0ebed3c"],
 "s3GuardTableName":"dc-pro-cli-7210",
 "Description":"ab-ds-cli-7321",
```

```
"enableTunnel":false,
 "workloadAnalytics":false,
 "freeIpa":{"instanceCountByGroup":1},
 }'

cdp environments set-id-broker-mappings \
--environment-name "ab-ds-cli-7321" \
--profile default \
--set-empty-mappings \
--data-access-role arn:aws:iam::<xxxxxxxxxxxx>:role/add-role \
--ranger-audit-role arn:aws:iam::<xxxxxxxxxxxx>:role/add-role
```

Similarly, while setting up the data lake use the **Instance Profile** that you configured above with Glue related policy in your data lake creation command:

```
cdp datalake create-aws-datalake --profile default --runtime 7.2.12 --cli-input-json '
{"datalakeName":"ab-ds-cli-7321-sdx",
 "environmentName":"ab-ds-cli-7321",

"cloudProviderConfiguration":{"instanceProfile":"arn:aws:iam::<xxxxxxxxxxxx>:instance-profile/<role-name>","storageBucketLocation":"s3a://demo-e2e-test-state-bucket/ab-ds-cli-7321"},
 "scale":"LIGHT_DUTY",
 }'
```

For more information, see [Creating an AWS environment with a medium duty data lake using the CLI](#).

  c. Navigate to the attached policy for the role.
  d. When you manually create tables in AWS Glue Data Catalog, you must set the **fully qualified path** for the table location.

For example:
```
"s3://manowar-e2e-test-state-bucket.s3-us-west-2.amazonaws.com/dc-pro-721-storage/glue/"
```

3. You must set up the AWS Glue Data Catalog. For more information, see [Populating the Glue Data Catalog](#). You must select only the CSV format which is currently supported for CDP Data Catalog and the delimiter which is used in the data.

4. While creating tables in AWS Glue Data Catalog manually, set the fully qualified path for location. For example: ("s3://manowar-e2e-test-state-bucket.s3-us-west-2.amazonaws.com/dc-pro-721-storage/glue/")

# Working with AWS Glue in CDP Data Catalog

The AWS Glue metadata must be registered with the CDP Data Catalog. The Glue contains the metadata that is synchronized with Data Catalog. The Glue metadata is accessed using the Data lake option in the Data Catalog service. After setting up AWS Glue with CDP Data Catalog and once the Glue synchronization with CDP is complete, the AWS Glue data lake appears in the Data Catalog instance.

**Note:** For each AWS environment, there would be a separate listing under the data lake option. The Glue data lake name / identifier in Data Catalog follows the format: <**glue**: **Data Lake Name**>. The Glue assets are of the type: **GLUE EXTERNAL TABLE**.
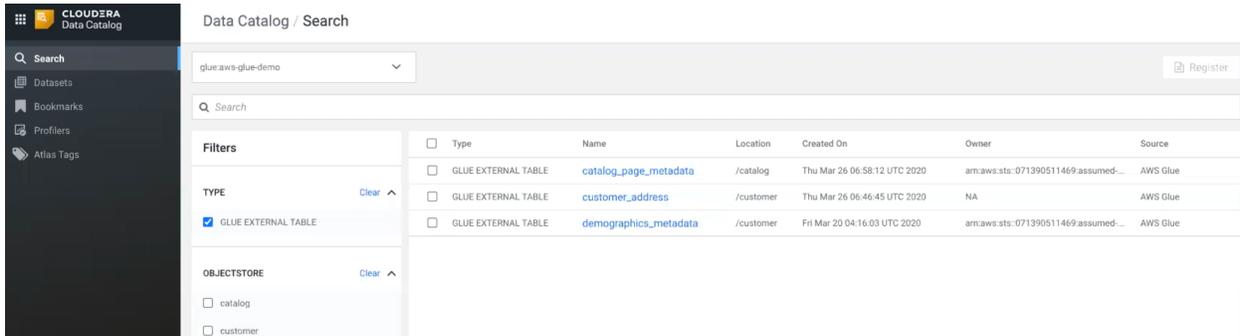
## Accessing AWS Glue in Data Catalog

To access the Glue metadata in Data Catalog, you must note the following in your Data Catalog instance.

- List the Glue Metadata by selecting the Glue data lake.
- Select one or more Glue assets and register the same with CDP
- Verify if the registered Glue assets are listed in the Data Catalog owned data lake
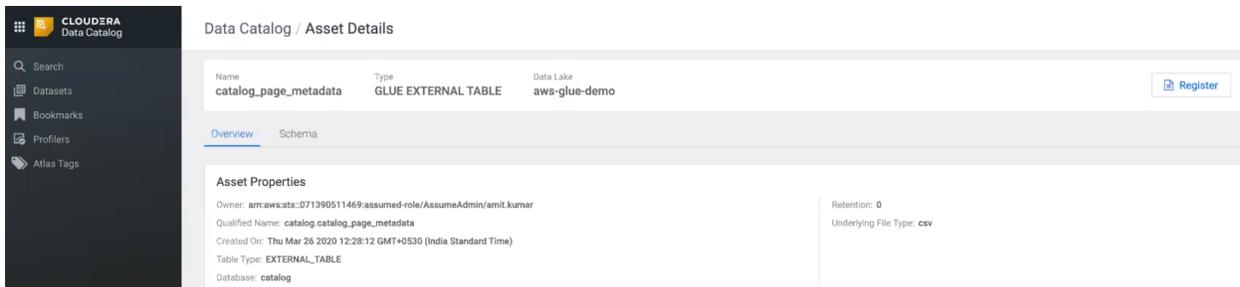- Select the registered Glue asset and click to open the **Asset Details** page

# Listing the Glue assets

In Data Catalog, when you select the AWS Glue data lake, you can view the list of Glue assets. These metadata assets are directly sourced from Glue.



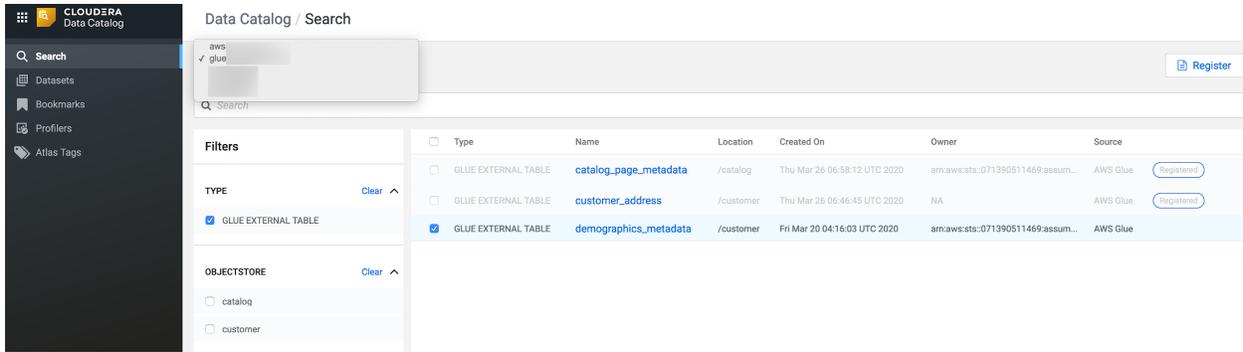When you click on one of the assets, the Asset Details page is displayed.



Next, on the main Data Catalog page, you must select the Glue data lake and select one of the Glue assets and register the asset to CDP. Click **Register**.

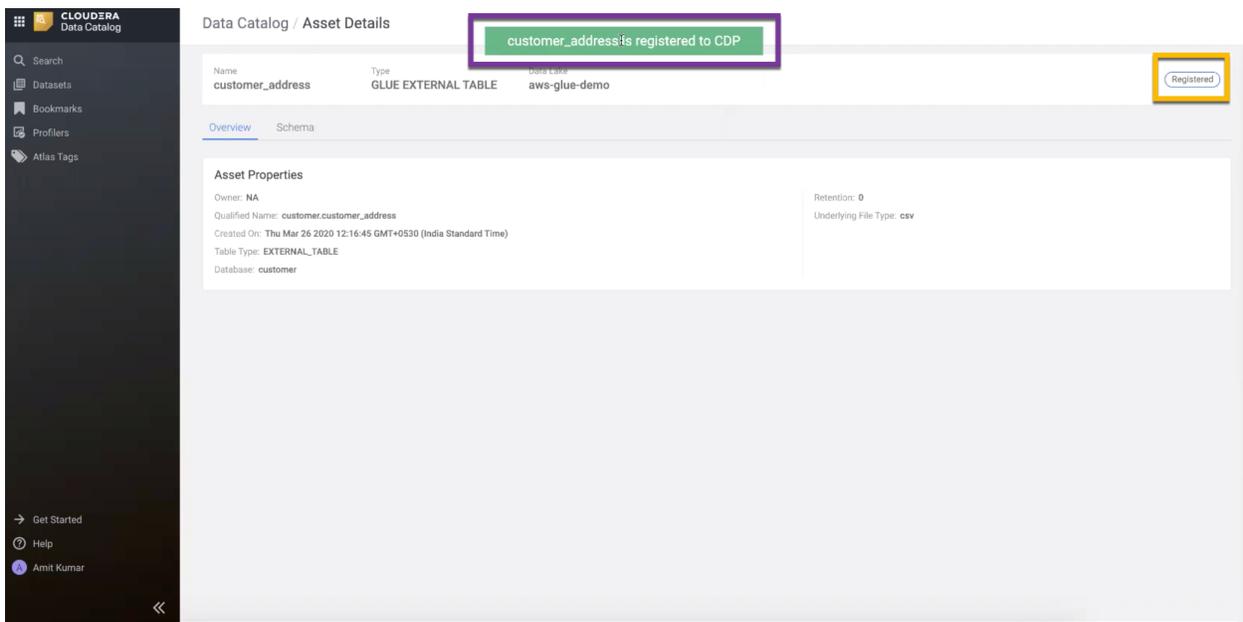Optionally, you directly click on the Glue asset and register the asset on the **Asset Details** page.

**Note:** You can select one or more Glue assets on the asset listing page to register them in CDP.
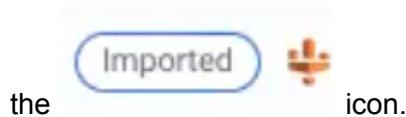
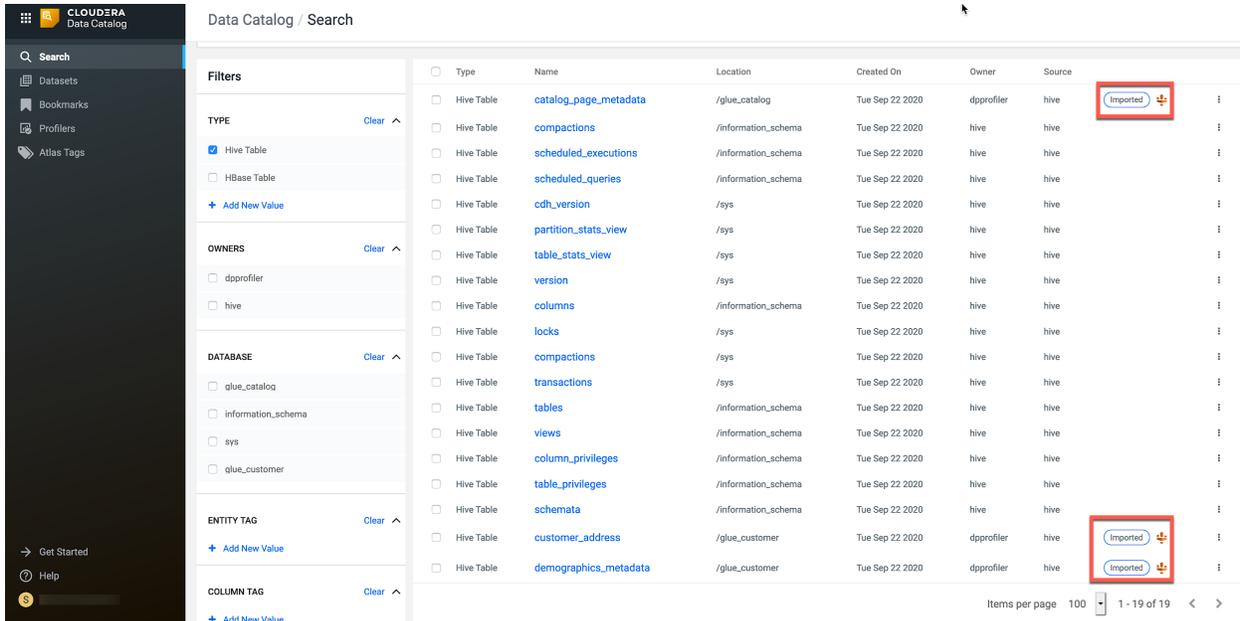Once the Glue asset is registered, the asset is imported into CDP.



Next, navigate back to the Data Catalog main page and select the Data Catalog owned data lake and select the type as **Hive Table**. The search results lists all the Hive table assets and you can view the Glue registered asset(s) as well. The registered Glue asset can be identified using the  icon.
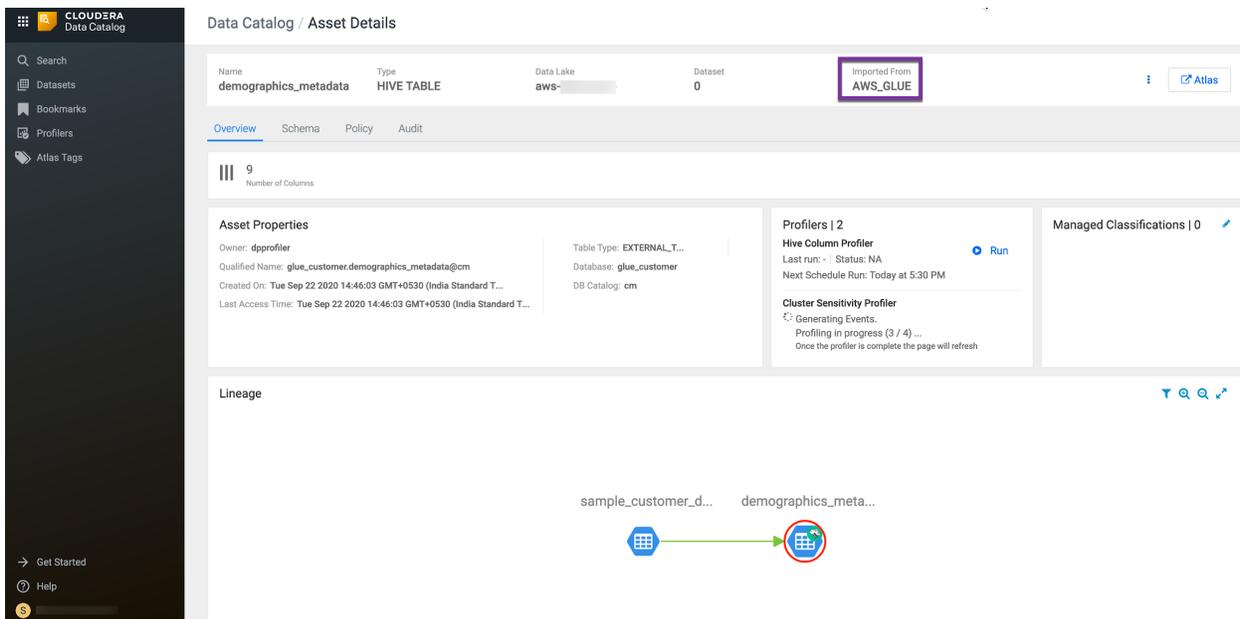
**Note:** The entries in the Data Catalog for Glue assets are created in the Hive Metastore.
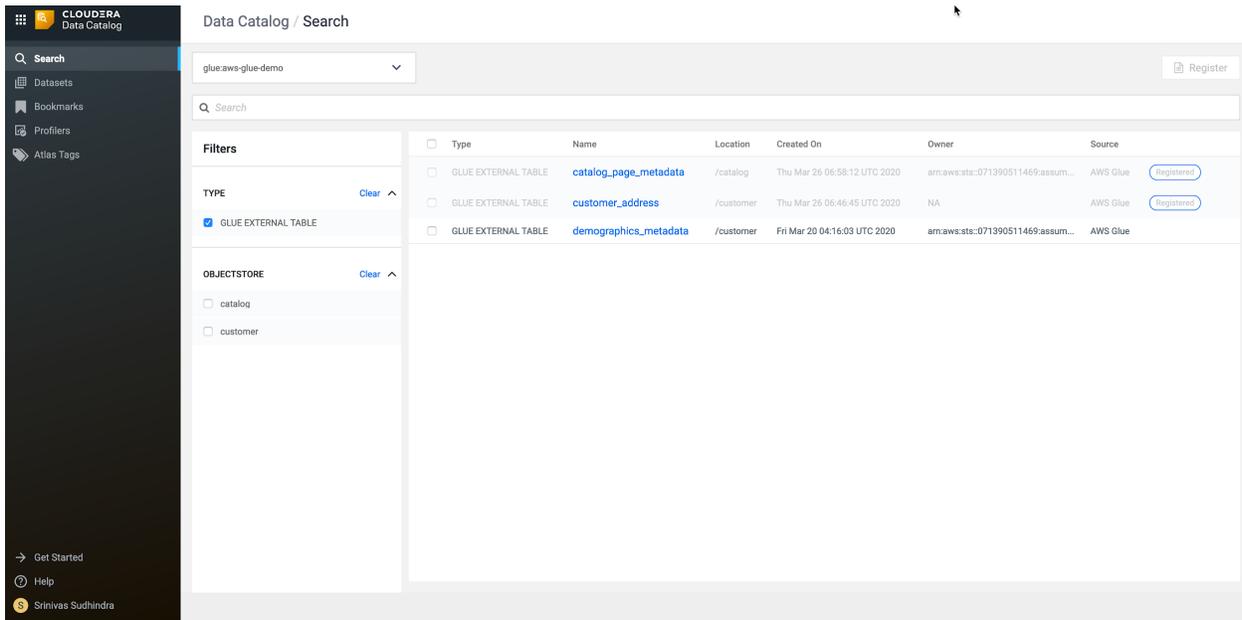
Click on the Glue registered asset and you can view the **Asset Details** page for the selected Glue asset.



The Asset Details page for the Glue asset is populated by Atlas. While registering the Glue data, the data was written to the Hive Metastore and later Atlas synchronised the metadata.

Go back to the main **Data Catalog** page and select the Glue data lake. Note that the registered Glue asset(s) are greyed out or cannot be selected again.



You can still view the registered Glue assets (powered by Atlas) by clicking on the same and it navigates to the **Asset Details** page as seen above in the image.

# Working with Ranger Authorization Service (RAZ) enabled AWS environment

For RAZ enabled AWS environment, you must employ the following permission settings to work with Data Catalog - Glue Integration.

**Note:** By default "**dpprofiler**" user is not included in the allowed users list. You must manually add the "**dpprofiler**" user in the allowed users list.

# Unsupported features

For a Glue data lake, while searching for a specific asset, multiple calls fail with "**400 Bad Request**" which is reflected on the Data Catalog UI in search suggestions. Search for Glue assets in Data Catalog is currently not supported.

# Known Issues

**CDPDSS-1179**: Any Glue registered asset does not fully display the respective Asset Details page when accessed from the Glue Data Lake.

**Problem**: When any asset in the Data Catalog is Glue registered, it does not display any information in the Asset Details page, when viewed from the Glue Data Lake.

**Workaround:** Select the same asset in Data Catalog, which is not Glue Data Lake and click on the same to view the Asset Details information.