



Informatica® Cloud Data Integration

Databricks Delta Connector

Informatica Cloud Data Integration Databricks Delta Connector  
April 2024

© Copyright Informatica LLC 2018, 2024

This software and documentation are provided only under a separate license agreement containing restrictions on use and disclosure. No part of this document may be reproduced or transmitted in any form, by any means (electronic, photocopying, recording or otherwise) without prior consent of Informatica LLC.

U.S. GOVERNMENT RIGHTS Programs, software, databases, and related documentation and technical data delivered to U.S. Government customers are "commercial computer software" or "commercial technical data" pursuant to the applicable Federal Acquisition Regulation and agency-specific supplemental regulations. As such, the use, duplication, disclosure, modification, and adaptation is subject to the restrictions and license terms set forth in the applicable Government contract, and, to the extent applicable by the terms of the Government contract, the additional rights set forth in FAR 52.227-19, Commercial Computer Software License.

Informatica, the Informatica logo, Informatica Cloud, and PowerCenter are trademarks or registered trademarks of Informatica LLC in the United States and many jurisdictions throughout the world. A current list of Informatica trademarks is available on the web at <https://www.informatica.com/trademarks.html>. Other company and product names may be trade names or trademarks of their respective owners.

Portions of this software and/or documentation are subject to copyright held by third parties. Required third party notices are included with the product.

See patents at <https://www.informatica.com/legal/patents.html>.

DISCLAIMER: Informatica LLC provides this documentation "as is" without warranty of any kind, either express or implied, including, but not limited to, the implied warranties of noninfringement, merchantability, or use for a particular purpose. Informatica LLC does not warrant that this software or documentation is error free. The information provided in this software or documentation may include technical inaccuracies or typographical errors. The information in this software and documentation is subject to change at any time without notice.

#### NOTICES

This Informatica product (the "Software") includes certain drivers (the "DataDirect Drivers") from DataDirect Technologies, an operating company of Progress Software Corporation ("DataDirect") which are subject to the following terms and conditions:

1. THE DATADIRECT DRIVERS ARE PROVIDED "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESSED OR IMPLIED, INCLUDING BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NON-INFRINGEMENT.
2. IN NO EVENT WILL DATADIRECT OR ITS THIRD PARTY SUPPLIERS BE LIABLE TO THE END-USER CUSTOMER FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, CONSEQUENTIAL OR OTHER DAMAGES ARISING OUT OF THE USE OF THE ODBC DRIVERS, WHETHER OR NOT INFORMED OF THE POSSIBILITIES OF DAMAGES IN ADVANCE. THESE LIMITATIONS APPLY TO ALL CAUSES OF ACTION, INCLUDING, WITHOUT LIMITATION, BREACH OF CONTRACT, BREACH OF WARRANTY, NEGLIGENCE, STRICT LIABILITY, MISREPRESENTATION AND OTHER TORTS.

The information in this documentation is subject to change without notice. If you find any problems in this documentation, report them to us at [infa\\_documentation@informatica.com](mailto:infa_documentation@informatica.com).

Informatica products are warranted according to the terms and conditions of the agreements under which they are provided. INFORMATICA PROVIDES THE INFORMATION IN THIS DOCUMENT "AS IS" WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING WITHOUT ANY WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND ANY WARRANTY OR CONDITION OF NON-INFRINGEMENT.

Publication Date: 2024-02-12

# Table of Contents

<b>Preface</b> .....	<b>5</b>
Informatica Resources. . . . .	5
Informatica Documentation. . . . .	5
Informatica Intelligent Cloud Services web site. . . . .	5
Informatica Intelligent Cloud Services Communities. . . . .	5
Informatica Intelligent Cloud Services Marketplace. . . . .	5
Data Integration connector documentation. . . . .	6
Informatica Knowledge Base. . . . .	6
Informatica Intelligent Cloud Services Trust Center. . . . .	6
Informatica Global Customer Support. . . . .	6
<b>Chapter 1: Introduction to Databricks Delta Connector</b> .....	<b>7</b>
Databricks Delta Connector assets. . . . .	7
Databricks compute resources. . . . .	8
External table support in Databricks Delta. . . . .	8
<b>Chapter 2: Connections for Databricks Delta</b> .....	<b>9</b>
Staging prerequisites. . . . .	9
SQL warehouse. . . . .	9
Configure AWS staging. . . . .	9
Configure Azure staging. . . . .	13
Databricks cluster. . . . .	14
Configure Spark parameters. . . . .	14
Configure Secure Agent properties. . . . .	14
Connect to Databricks Delta. . . . .	15
Before you begin. . . . .	15
Connection details. . . . .	16
Advanced settings. . . . .	17
Configure proxy settings. . . . .	21
JDBC URL parameters. . . . .	21
Rules and guidelines for personal staging location. . . . .	22
Private links to access Databricks Delta. . . . .	22
<b>Chapter 3: Mappings for Databricks Delta</b> .....	<b>23</b>
Before you begin. . . . .	23
Verify permissions. . . . .	24
Sources for Databricks Delta. . . . .	25
Source properties for Databricks Delta. . . . .	25
Custom query source type. . . . .	29
Key range partitioning. . . . .	29

Targets for Databricks Delta. . . . .	30
Target properties for Databricks Delta . . . . .	31
Create a target table at runtime. . . . .	33
Override the update operation. . . . .	35
Determine the order of processing for multiple targets. . . . .	36
Lookups for Databricks Delta. . . . .	36
Lookup properties for Databricks Delta. . . . .	37
Enable lookup caching. . . . .	39
SQL transformation. . . . .	39
Dynamic schema handling. . . . .	40
IDENTITY columns . . . . .	40
Mappings in advanced mode example. . . . .	41
Rules and guidelines for mappings. . . . .	42
Rules and guidelines for mappings in advanced mode. . . . .	44
<b>Chapter 4: Databricks Delta SQL ELT optimization. . . . .</b>	<b>46</b>
SQL ELT optimization types. . . . .	46
Previewing SQL ELT optimization. . . . .	47
Configuring SQL ELT optimization . . . . .	47
SQL ELT optimization using a Databricks Delta connection. . . . .	48
Read from and write to Databricks Delta. . . . .	48
Read from Amazon S3 and write to Databricks Delta. . . . .	48
Read from Microsoft Azure Data Lake Storage Gen2 and write to Databricks Delta. . . . .	48
SQL ELT compatibility. . . . .	49
Transformations with Databricks Delta. . . . .	51
Features. . . . .	53
Configuring a custom query for the Databricks Delta source object. . . . .	56
SQL ELT optimization for multiple targets. . . . .	57
Single commit for SQL ELT optimization. . . . .	57
Rules and guidelines for SQL ELT optimization . . . . .	58
Troubleshooting SQL ELT optimization. . . . .	61
<b>Chapter 5: Data type reference. . . . .</b>	<b>63</b>
Databricks Delta and transformation data types. . . . .	63
<b>Index. . . . .</b>	<b>65</b>

# Preface

Use *Databricks Delta Connector* to learn how to read from or write to Databricks Delta by using Cloud Data Integration. Learn to create a connection and develop mappings, mapping tasks, dynamic mapping tasks, and data transfer tasks in Cloud Data Integration.

## Informatica Resources

Informatica provides you with a range of product resources through the Informatica Network and other online portals. Use the resources to get the most from your Informatica products and solutions and to learn from other Informatica users and subject matter experts.

### Informatica Documentation

Use the Informatica Documentation Portal to explore an extensive library of documentation for current and recent product releases. To explore the Documentation Portal, visit <https://docs.informatica.com>.

If you have questions, comments, or ideas about the product documentation, contact the Informatica Documentation team at [infa\\_documentation@informatica.com](mailto:infa_documentation@informatica.com).

### Informatica Intelligent Cloud Services web site

You can access the Informatica Intelligent Cloud Services web site at <http://www.informatica.com/cloud>. This site contains information about Informatica Cloud integration services.

### Informatica Intelligent Cloud Services Communities

Use the Informatica Intelligent Cloud Services Community to discuss and resolve technical issues. You can also find technical tips, documentation updates, and answers to frequently asked questions.

Access the Informatica Intelligent Cloud Services Community at:

<https://network.informatica.com/community/informatica-network/products/cloud-integration>

Developers can learn more and share tips at the Cloud Developer community:

<https://network.informatica.com/community/informatica-network/products/cloud-integration/cloud-developers>

### Informatica Intelligent Cloud Services Marketplace

Visit the Informatica Marketplace to try and buy Data Integration Connectors, templates, and mapplets:

<https://marketplace.informatica.com/>

## Data Integration connector documentation

You can access documentation for Data Integration Connectors at the Documentation Portal. To explore the Documentation Portal, visit <https://docs.informatica.com>.

## Informatica Knowledge Base

Use the Informatica Knowledge Base to find product resources such as how-to articles, best practices, video tutorials, and answers to frequently asked questions.

To search the Knowledge Base, visit <https://search.informatica.com>. If you have questions, comments, or ideas about the Knowledge Base, contact the Informatica Knowledge Base team at [KB\\_Feedback@informatica.com](mailto:KB_Feedback@informatica.com).

## Informatica Intelligent Cloud Services Trust Center

The Informatica Intelligent Cloud Services Trust Center provides information about Informatica security policies and real-time system availability.

You can access the trust center at <https://www.informatica.com/trust-center.html>.

Subscribe to the Informatica Intelligent Cloud Services Trust Center to receive upgrade, maintenance, and incident notifications. The [Informatica Intelligent Cloud Services Status](#) page displays the production status of all the Informatica cloud products. All maintenance updates are posted to this page, and during an outage, it will have the most current information. To ensure you are notified of updates and outages, you can subscribe to receive updates for a single component or all Informatica Intelligent Cloud Services components. Subscribing to all components is the best way to be certain you never miss an update.

To subscribe, on the [Informatica Intelligent Cloud Services Status](#) page, click **SUBSCRIBE TO UPDATES**. You can choose to receive notifications sent as emails, SMS text messages, webhooks, RSS feeds, or any combination of the four.

## Informatica Global Customer Support

You can contact a Global Support Center through the Informatica Network or by telephone.

To find online support resources on the Informatica Network, click **Contact Support** in the Informatica Intelligent Cloud Services Help menu to go to the **Cloud Support** page. The **Cloud Support** page includes system status information and community discussions. Log in to Informatica Network and click **Need Help** to find additional resources and to contact Informatica Global Customer Support through email.

The telephone numbers for Informatica Global Customer Support are available from the Informatica web site at <https://www.informatica.com/services-and-training/support-services/contact-us.html>.

# CHAPTER 1

## Introduction to Databricks Delta Connector

You can use Databricks Delta Connector to securely read data from or write data to Databricks Delta.

When you use Databricks Delta Connector, you can create a Databricks Delta connection to connect to Databricks Delta endpoints hosted on Amazon Web Services, and Microsoft Azure environment.

When you read data from and write data to Databricks Delta, you can stage the data in Amazon Web Services, Microsoft Azure environment, or personal staging location.

You can use Databricks Delta Connector on the Windows and Linux operating systems.

For Linux operating systems, you can switch mappings to advanced mode to include transformations and functions that enable advanced functionality.

You can switch mappings to advanced mode to include transformations and functions that enable advanced functionality. A mapping in advanced mode can run on the advanced cluster hosted on Amazon Web Services, Microsoft Azure environment, or on a self-service cluster.

## Databricks Delta Connector assets

Create assets in Data Integration to integrate data using Databricks Delta Connector.

When you use Databricks Delta Connector, you can include the following Data Integration assets:

- Data transfer task
- Dynamic mapping task
- Mapping
- Mapping task

For more information about configuring assets and transformations, see *Mappings*, *Transformations*, and *Tasks* in the Data Integration documentation.

# Databricks compute resources

You can use a Databricks Delta connection to connect to SQL warehouse or Databricks cluster to read from and write to Databricks Delta tables.

- **SQL warehouse**

The Secure Agent connects to the SQL warehouse at design time and runtime.

You can use a SQL warehouse on the Windows and Linux operating systems.

- **Databricks cluster**

In this document, the term Databricks cluster refers to the all-purpose cluster and job cluster.

The Secure Agent connects to the all-purpose cluster to import the metadata at design time and to the job cluster to run the mappings.

To connect to Databricks cluster, enable the Secure Agent properties for the design time and runtime.

You can use a Databricks cluster only on the Linux operating system.

## External table support in Databricks Delta

External tables store their data in locations outside of the predefined managed storage location associated with the metastore, unity catalog, or schema. An external table references an external storage path by using a `LOCATION` clause. For more information on external tables, see the Databricks Delta documentation.

**Note:** You can read and write data to external tables of file type Delta in Databricks.

The following lists the data types that Databricks Delta supports when you create an external table:

- Array\*
- Binary
- Bigint
- Boolean
- Date
- Decimal
- Double
- Float
- Int
- Map\*
- Smallint
- String
- Struct\*
- Tinyint
- Timestamp

\*Applies only to mappings in advanced mode.

When you create an external table, specify the table location when you create a new target at runtime. For more information, see [“Create a target table at runtime” on page 33](#)



## CHAPTER 2

# Connections for Databricks Delta

Create a Databricks Delta connection to securely read data from or write data to Databricks Delta.

You can use a Databricks Delta connection to specify sources, targets, and lookups in mappings and mapping tasks.

## Staging prerequisites

Before you create a connection, you must perform certain prerequisite tasks to configure the staging environment to connect to SQL warehouse or Databricks cluster.

## SQL warehouse

Configure either the AWS or Azure staging environment for the SQL warehouse based on the deployed environment. You also need to configure the Spark parameters for the SQL warehouse to use Azure and AWS staging.

### Configure AWS staging

Configure IAM AssumeRole authentication to use AWS staging for the SQL warehouse.

#### IAM AssumeRole authentication

You can enable IAM AssumeRole authentication in Databricks Delta for secure and controlled access to the Amazon S3 staging bucket when you run mappings and mapping tasks.

You can configure IAM authentication when the Secure Agent runs on an Amazon Elastic Compute Cloud (EC2) system. When you use a serverless runtime environment, you cannot configure IAM authentication.

Perform the following steps to configure IAM authentication on EC2:

1. Create a minimal Amazon IAM policy.
2. Create the Amazon EC2 role. The Amazon EC2 role is used when you create an EC2 system. For more information about creating the Amazon EC2 role, see the *AWS documentation*.
3. Link the minimal Amazon IAM policy with the Amazon EC2 role.
4. Create an EC2 instance. Assign the Amazon EC2 role that you created to the EC2 instance.
5. Install the Secure Agent on the EC2 system.

## Temporary security credentials using AssumeRole

You can use temporary security credentials using AssumeRole to access AWS resources from same or different AWS accounts.

Ensure that you have the **sts:AssumeRole** permission and a trust relationship established within the AWS accounts to use temporary security credentials. The trust relationship is defined in the trust policy of the IAM role when you create the role. The IAM role adds the IAM user as a trusted entity allowing the IAM users to use temporary security credentials and access AWS accounts.

For more information about how to establish the trust relationship, see the *AWS documentation*.

When the trusted IAM user requests for temporary security credentials, the AWS Security Token Service (AWS STS) dynamically generates the temporary security credentials that are valid for a specified period and provides the credentials to the trusted IAM users. The temporary security credentials consist of access key ID, secret access key, and secret token.

To use the dynamically generated temporary security credentials, provide a value for the **IAM Role ARN** connection property when you create a Databricks Delta connection. The IAM Role ARN uniquely identifies the AWS resources. Then, specify the time duration in seconds during which you can use the temporarily security credentials in the **Temporary Credential Duration** advanced source and target properties.

### External ID

You can specify the external ID for a more secure cross-account access to the Amazon S3 bucket when the Amazon S3 bucket is in a different AWS account.

Optionally, you can specify the external ID in the AssumeRole request to the AWS Security Token Service (STS).

The external ID must be a string.

The following sample shows an external ID condition in the assumed IAM role trust policy:

```
"Statement": [
  {
    "Effect": "Allow",
    "Principal": {
      "AWS": "arn:aws:iam::AWS_Account_ID : user/user_name"
    },
    "Action": "sts:AssumeRole",
    "Condition": {
      "StringEquals": {
        "sts:ExternalId": "dummy_external_id"
      }
    }
  }
]
```

### Temporary security credentials policy

To use temporary security credentials to access AWS resources, both the IAM user and IAM role require policies.

#### Amazon S3 permission policy

Attach the following S3 permission policy to allow access to the Amazon S3 bucket:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Action": [
        "s3:DeleteObject",
        "s3:GetObject",
        "s3:ListBucket",
        "s3:PutObject",
        "s3:PutObjectTagging",

```

```

        "s3:GetBucketAcl"
    ],
    "Effect": "Allow",
    "Resource": "arn:aws:s3:::com.amk"
  },
  {
    "Effect": "Allow",
    "Action": [
      "s3:DeleteObject",
      "s3:GetObject",
      "s3:ListBucket",
      "s3:PutObject",
      "s3:PutObjectTagging",
      "s3:GetBucketAcl"
    ],
    "Resource": "arn:aws:s3:::com.amk/*"
  }
]
}

```

The following section lists the policies required for IAM user and IAM role:

#### **IAM user**

An IAM user must have the `sts:AssumeRole` policy to use temporary security credentials in same or different AWS account.

The following sample policy allows an IAM user to use the temporary security credentials in an AWS account:

```

{
  "Version": "2012-10-17",
  "Statement": {
    "Effect": "Allow",
    "Action": "sts:AssumeRole",
    "Resource": "arn:aws:iam::<ACCOUNT-HYPHENS>:role/<ROLE-NAME>" }
  }
}

```

#### **IAM role**

An IAM role must have the `sts:AssumeRole` policy and a trust policy attached with the IAM role to allow the IAM user to access the AWS resource using temporary security credentials. The policy specifies the AWS resource that the IAM user can access and the actions that the IAM user can perform. The trust policy specifies the IAM user from the AWS account that can access the AWS resource.

The following policy is a sample trust policy:

```

{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": { "AWS": "arn:aws:iam::AWS-account-ID:root" },
      "Action": "sts:AssumeRole"
    }
  ]
}

```

Here, in the `Principal` attribute, you can also provide the ARN of IAM user who can use the dynamically generated temporary security credentials and to restrict further access. For example,

```

"Principal" : { "AWS" : "arn:aws:iam:: AWS-account-ID :user/ user-name " }

```

#### **Temporary security credentials using AssumeRole for EC2**

You can use temporary security credentials using `AssumeRole` for an Amazon EC2 role to access AWS resources from same or different AWS accounts.

The Amazon EC2 role would be able to assume another IAM Role from the same or a different AWS account without requiring the permanent access key and secret key.

Consider the following prerequisites when you use temporary security credentials using AssumeRole for EC2:

- Install the Secure Agent on an AWS service such as Amazon EC2.
- The EC2 role attached to the AWS EC2 service does not need access to Amazon S3 but needs permission to assume another IAM role.
- The IAM role that needs to be assumed by the EC2 role must have a permission policy and a trust policy attached to it.

To configure an EC2 role to assume the IAM role provided in the **IAM Role ARN** connection property, select the **Use EC2 Role to Assume Role** check box in the connection properties.

## Create a minimal Amazon IAM policy

To stage the data in Amazon S3, use the following minimum required permissions:

- PutObject
- GetObject
- DeleteObject
- ListBucket
- ListBucketMultipartUploads. Applicable only for mappings in advanced mode.

You can use the following sample Amazon IAM policy:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "s3:PutObject",
        "s3:GetObject",
        "s3:DeleteObject",
        "s3:ListBucket",

        "s3:ListBucketMultipartUploads"
      ],
      "Resource": [
        "arn:aws:s3:::<bucket_name>/*",
        "arn:aws:s3:::<bucket_name>"
      ]
    }
  ]
}
```

For mappings in advanced mode, you can use different AWS accounts within the same AWS region. Make sure that the Amazon IAM policy confirms access to the AWS accounts used in these mappings.

**Note:** The **Test Connection** does not validate the IAM policy assigned to users. You can specify the Amazon S3 bucket name in the source and target advanced properties.

## Configure Spark parameters for AWS staging

Before you use the Databricks SQL warehouse to run mappings, configure the Spark parameters for SQL warehouse on the Databricks SQL Admin console.

On the Databricks SQL Admin console, navigate to **SQL Warehouse Settings > Data Security**, and then configure the Spark parameters for AWS under **Data access configuration**.

Add the following Spark configuration parameters and restart the SQL warehouse:

- `spark.hadoop.fs.s3a.access.key` <S3 Access Key value>
- `spark.hadoop.fs.s3a.secret.key` <S3 Secret Key value>
- `spark.hadoop.fs.s3a.endpoint` <S3 Staging Bucket endpoint value>

For example, the S3 staging bucket warehouse value is `s3.ap-south-1.amazonaws.com`.

Ensure that the configured access key and secret key have access to the S3 buckets where you store the data for Databricks Delta tables.

## Configure Azure staging

Before you use Microsoft Azure Data Lake Storage Gen2 to stage files, perform the following tasks:

- Create a storage account to use with Microsoft Azure Data Lake Storage Gen2 and enable **Hierarchical namespace** in the Azure portal.  
You can use role-based access control to authorize the users to access the resources in the storage account. Assign the Contributor role or Reader role to the users. The contributor role grants you full access to manage all resources in the storage account, but does not allow you to assign roles. The reader role allows you to view all resources in the storage account, but does not allow you to make any changes.  
**Note:** To add or remove role assignments, you must have write and delete permissions, such as an Owner role.
- Register an application in Azure Active Directory to authenticate users to access the Microsoft Azure Data Lake Storage Gen2 account.  
You can use role-based access control to authorize the application. Assign the Storage Blob Data Contributor or Storage Blob Data Reader role to the application. The Storage Blob Data Contributor role lets you read, write, and delete Azure Storage containers and blobs in the storage account. The Storage Blob Data Reader role lets you only read and list Azure Storage containers and blobs in the storage account.
- Create an Azure Active Directory web application for service-to-service authentication with Microsoft Azure Data Lake Storage Gen2.  
**Note:** Ensure that you have superuser privileges to access the folders or files created in the application using the connector.
- To read and write complex files, set the JVM options for type DTM to increase the -Xms and -Xmx values in the system configuration details of the Secure Agent to avoid java heap space error. The recommended -Xms and -Xmx values are 512 MB and 1024 MB respectively.

## Configure Spark parameters for Azure staging

Before you use the Databricks SQL warehouse to run mappings, configure the Spark parameters for SQL warehouse on the Databricks SQL Admin console.

On the Databricks SQL Admin console, navigate to **SQL Warehouse Settings > Data Security**, and then configure the Spark parameters for Azure under **Data access configuration**.

Add the following Spark configuration parameters and restart the SQL warehouse:

- `spark.hadoop.fs.azure.account.oauth2.client.id.<storage-account-name>.dfs.core.windows.net` <value>
- `spark.hadoop.fs.azure.account.auth.type.<storage-account-name>.dfs.core.windows.net` OAuth
- `spark.hadoop.fs.azure.account.oauth2.client.secret.<storage-account-name>.dfs.core.windows.net` <Value>

- `spark.hadoop.fs.azure.account.oauth.provider.type.<storage-account-name>.dfs.core.windows.net org.apache.hadoop.fs.azurebfs.oauth2.ClientCredsTokenProvider`
- `spark.hadoop.fs.azure.account.oauth2.client.endpoint.<storage-account-name>.dfs.core.windows.net https://login.microsoftonline.com/<Tenant ID>/oauth2/token`

Ensure that the configured client ID and client secret have access to the file systems where you store the data for Databricks Delta tables.

## Databricks cluster

Configure the Spark parameters for the Databricks cluster to use Azure and AWS staging based on where the cluster is deployed.

You also need to enable the Secure Agent properties for runtime and design-time processing on the Databricks cluster.

### Configure Spark parameters

Before you connect to the Databricks cluster, you must configure the Spark parameters on AWS and Azure.

#### Configuration on AWS

Add the following Spark configuration parameters for the Databricks cluster and restart the cluster:

- `spark.hadoop.fs.s3a.access.key <value>`
- `spark.hadoop.fs.s3a.secret.key <value>`
- `spark.hadoop.fs.s3a.endpoint <value>`

Ensure that the access and secret key configured has access to the buckets where you store the data for Databricks Delta tables.

#### Configuration on Azure

Add the following Spark configuration parameters for the Databricks cluster and restart the cluster:

- `fs.azure.account.oauth2.client.id.<storage-account-name>.dfs.core.windows.net <value>`
- `fs.azure.account.auth.type.<storage-account-name>.dfs.core.windows.net <value>`
- `fs.azure.account.oauth2.client.secret.<storage-account-name>.dfs.core.windows.net <Value>`
- `fs.azure.account.oauth.provider.type.<storage-account-name>.dfs.core.windows.net org.apache.hadoop.fs.azurebfs.oauth2.ClientCredsTokenProvider`
- `fs.azure.account.oauth2.client.endpoint.<storage-account-name>.dfs.core.windows.net https://login.microsoftonline.com/<Tenant ID>/oauth2/token`

Ensure that the client ID and client secret configured has access to the file systems where you store the data for Databricks Delta tables.

### Configure Secure Agent properties

When you configure mappings, the SQL warehouse processes the mapping by default. To process the mappings on Databricks cluster, enable the Secure Agent properties.

To connect to all-purpose cluster and job cluster, enable the Secure Agent properties for design time and runtime respectively.

### Setting the property for design time processing

Before you can import metadata and design mappings or mappings in advanced mode, perform the following steps:

1. In **Administrator**, select the Secure Agent listed on the **Runtime Environments** tab.
2. Click **Edit**.
3. In the **System Configuration Details** section, select Data Integration Server as the **Service** and Tomcat JRE as the **Type**.
4. Edit the **JRE\_OPTS** field and set the value to `-DUseDatabricksSql=false`.

Tomcat JRE	JRE_OPTS	'-Xrs -DUseDatabricksSql=false'
------------	----------	---------------------------------

### Setting the property for runtime processing

1. In **Administrator**, select the Secure Agent listed on the **Runtime Environments** tab.
2. Click **Edit**.
3. In the **System Configuration Details** section, select Data Integration Server as the **Service** and DTM as the **Type**.
4. Edit the **JVMOption** field.
  - a. To run mappings, set the value to `-DUseDatabricksSql=false`.

DTM	JVMOption2	'-DUseDatabricksSql=false'
-----	------------	----------------------------

- b. To run mappings enabled with SQL ELT optimization, set the value to `-DUseDatabricksSqlForPdo=false`.

DTM	JVMOption3	'-DUseDatabricksSqlForPdo=false'
-----	------------	----------------------------------

## Connect to Databricks Delta

Let's configure the Databricks Delta connection properties to connect to Databricks Delta.

### Before you begin

Before you get started, you'll need to get information from your Databricks Delta account.

The following video shows you how to get the information you need:



You also need to configure the AWS or Azure staging environment to use the SQL warehouse or the Databricks cluster in the connection.

To learn about the staging prerequisites for the Azure or AWS environment, check out [“SQL warehouse” on page 9](#) or [“Databricks cluster” on page 14](#).

## Connection details

The following table describes the basic connection properties:

Property	Description
Connection Name	Name of the connection. Each connection name must be unique within the organization. Connection names can contain alphanumeric characters, spaces, and the following special characters: _ . + -, Maximum length is 255 characters.
Description	Description of the connection. Maximum length is 4000 characters.
Type	Databricks Delta
Runtime Environment	The name of the runtime environment where you want to run tasks. Select a Secure agent, Hosted Agent, or serverless runtime environment. Hosted Agent is not applicable for mappings in advanced mode.
SQL Warehouse JDBC URL	Databricks SQL Warehouse JDBC connection URL. Required for SQL warehouse. Doesn't apply to Databricks cluster. To get the SQL Warehouse JDBC URL, go to the Databricks console and select the JDBC driver version from the JDBC URL menu. For JDBC URL version 2.6.22 or earlier, use the following syntax: <code>jdbc:spark://&lt;Databricks Host&gt;:443/default;transportMode=http;ssl=1;AuthMech=3;httpPath=/sql/1.0/endpoints/&lt;SQL endpoint cluster ID&gt;;</code> For JDBC URL version 2.6.25 or later, use the following syntax: <code>jdbc:databricks://&lt;Databricks Host&gt;:443/default;transportMode=http;ssl=1;AuthMech=3;httpPath=/sql/1.0/endpoints/&lt;SQL endpoint cluster ID&gt;;</code> Ensure that you set the required environment variables in the Secure Agent. <b>Note:</b> Specify the database name in the Database Name connection property. If you specify the database name in the JDBC URL, it is not considered. The Databricks Host, Organization ID, and Cluster ID properties are not considered if you configure the SQL warehouse JDBC URL property.
Databricks Token	Personal access token to access Databricks. Required for SQL warehouse and Databricks cluster. Ensure that you have permissions to attach to the cluster identified in the <b>Cluster ID</b> property. For mappings, you must have additional permissions to create Databricks clusters.
Catalog Name	The name of an existing catalog in the metastore. Optional for SQL warehouse. Doesn't apply to Databricks cluster. You can also specify the catalog name in the end of the SQL warehouse JDBC URL. <b>Note:</b> The catalog name cannot contain special characters. For more information about unity catalog, see the Databricks Delta documentation.



## Advanced settings

The following table describes the advanced connection properties:

Property	Description
Database	<p>The database name that you want to connect to in Databricks Delta.</p> <p>Optional for SQL warehouse and Databricks cluster.</p> <p>For Data Integration, if you do not provide a database name, all databases available in the workspace are listed. The value you provide here overrides the database name provided in the <b>SQL Warehouse JDBC URL</b> connection property.</p>
JDBC Driver Class Name	<p>The name of the JDBC driver class.</p> <p>Optional for SQL warehouse and Databricks cluster.</p> <p>For JDBC URL versions 2.6.22 or earlier, specify the driver class name as <code>com.simba.spark.jdbc.Driver</code>.</p> <p>For JDBC URL versions 2.6.25 or later, specify the driver class name as <code>com.databricks.client.jdbc.Driver</code>.</p>
Staging Environment	<p>The cloud provider where the Databricks cluster is deployed.</p> <p>Required for SQL warehouse and Databricks cluster.</p> <p>Select one of the following options:</p> <ul style="list-style-type: none"> <li>- AWS</li> <li>- Azure</li> <li>- Personal Staging Location</li> </ul> <p>Default is Personal Staging Location.</p> <p>You can select the Personal Staging Location as the staging environment instead of Azure or AWS staging environments to stage data locally for mappings and tasks.</p> <p>Personal staging location doesn't apply to Databricks cluster.</p> <p>You cannot use personal staging location when you configure mappings in advanced mode.</p> <p><b>Note:</b> You cannot switch between clusters once you establish a connection.</p>
Databricks Host	<p>The host name of the endpoint the Databricks account belongs to.</p> <p>Required for Databricks cluster. Doesn't apply to SQL warehouse.</p> <p>You can get the Databricks Host from the JDBC URL. The URL is available in the Advanced Options of JDBC or ODBC in the Databricks Delta all-purpose cluster.</p> <p>The following example shows the Databricks Host in JDBC URL:</p> <pre>jdbc:spark://&lt;Databricks Host&gt;:443/ default;transportMode=http; ssl=1;httpPath=sq/protocolv1/o/&lt;Org Id&gt;/&lt;Cluster ID&gt;; AuthMech=3; UID=token; PWD=&lt;personal-access-token&gt;</pre> <p>The value of PWD in Databricks Host, Organization Id, and Cluster ID is always <code>&lt;personal-access-token&gt;</code>.</p>
Cluster ID	<p>The ID of the cluster.</p> <p>Required for Databricks cluster. Doesn't apply to SQL warehouse.</p> <p>You can get the cluster ID from the JDBC URL. The URL is available in the Advanced Options of JDBC or ODBC in the Databricks Delta all-purpose cluster</p> <p>The following example shows the Cluster ID in JDBC URL:</p> <pre>jdbc:spark://&lt;Databricks Host&gt;:443/ default;transportMode=http; ssl=1;httpPath=sq/protocolv1/o/&lt;Org Id&gt;/&lt;Cluster ID&gt;; AuthMech=3;UID=token; PWD=&lt;personal-access-token&gt;</pre>

Property	Description
Organization ID	<p>The unique organization ID for the workspace in Databricks. Required for Databricks cluster. Doesn't apply to SQL warehouse.</p> <p>You can get the Organization ID from the JDBC URL. The URL is available in the Advanced Options of JDBC or ODBC in the Databricks Delta all-purpose cluster</p> <p>The following example shows the Organization ID in JDBC URL:  <code>jdbc:spark://&lt;Databricks Host&gt;:443/ default;transportMode=http; ssl=1;httpPath=sql/protocolv1/o/&lt;Organization ID&gt;/ &lt;Cluster ID&gt;;AuthMech=3;UID=token; PWD=&lt;personal-access-token&gt;</code></p>
Min Workers <sup>1</sup>	<p>The minimum number of worker nodes to be used for the Spark job. Minimum value is 1. Required for Databricks cluster. Doesn't apply to SQL warehouse.</p>
Max Workers <sup>1</sup>	<p>The maximum number of worker nodes to be used for the Spark job. If you don't want to autoscale, set Max Workers = Min Workers or don't set Max Workers. Optional for Databricks cluster. Doesn't apply to SQL warehouse.</p>
DB Runtime Version <sup>1</sup>	<p>The version of Databricks cluster to spawn when you connect to Databricks cluster to process mappings. Required for Databricks cluster. Doesn't apply to SQL warehouse. Select the Databricks runtime version 9.1 LTS or 13.3 LTS.</p>
Worker Node Type <sup>1</sup>	<p>The worker node instance type that is used to run the Spark job. Required for Databricks cluster. Doesn't apply to SQL warehouse. For example, the worker node type for AWS can be i3.2xlarge. The worker node type for Azure can be Standard_DS3_v2.</p>
Driver Node Type <sup>1</sup>	<p>The driver node instance type that is used to collect data from the Spark workers. Optional for Databricks cluster. Doesn't apply to SQL warehouse. For example, the driver node type for AWS can be i3.2xlarge. The driver node type for Azure can be Standard_DS3_v2. If you don't specify the driver node type, Databricks uses the value you specify in the worker node type field.</p>
Instance Pool ID <sup>1</sup>	<p>The instance pool ID used for the Spark cluster. Optional for Databricks cluster. Doesn't apply to SQL warehouse. If you specify the Instance Pool ID to run mappings, the following connection properties are ignored:</p> <ul style="list-style-type: none"> <li>- Driver Node Type</li> <li>- EBS Volume Count</li> <li>- EBS Volume Type</li> <li>- EBS Volume Size</li> <li>- Enable Elastic Disk</li> <li>- Worker Node Type</li> <li>- Zone ID</li> </ul>
Elastic Disk <sup>1</sup>	<p>Enables the cluster to get additional disk space. Optional for Databricks cluster. Doesn't apply to SQL warehouse. Enable this option if the Spark workers are running low on disk space.</p>

Property	Description
Spark Configuration <sup>1</sup>	The Spark configuration to use in the Databricks cluster. Optional for Databricks cluster. Doesn't apply to SQL warehouse. The configuration must be in the following format: <code>"key1"="value1";"key2"="value2";...</code> For example, <code>"spark.executor.userClassPathFirst"="False"</code>
Spark Environment Variables <sup>1</sup>	The environment variables to export before launching the Spark driver and workers. Optional for Databricks cluster. Doesn't apply to SQL warehouse. The variables must be in the following format: <code>"key1"="value1";"key2"="value2";...</code> For example, <code>"MY_ENVIRONMENT_VARIABLE"="true"</code>
<sup>1</sup> Doesn't apply to mappings in advanced mode.	

## AWS staging environment

The following table describes the properties for the AWS staging environment:

Property	Description
S3 Access Key	The key to access the Amazon S3 bucket.
S3 Secret Key	The secret key to access the Amazon S3 bucket.
S3 Data Bucket	The existing bucket to store the Databricks Delta data.
S3 Staging Bucket	The existing bucket to store the staging files.
S3 Authentication Mode	The authentication mode to access Amazon S3. Select one of the following authentication modes: <ul style="list-style-type: none"> <li>- Permanent IAM credentials. Uses the S3 access key and S3 secret key to connect to Databricks Delta.</li> <li>- IAM Assume Role<sup>1</sup>. Uses the AssumeRole for IAM authentication to connect to Databricks Delta. Doesn't apply to Databricks cluster.</li> </ul>
IAM Role ARN <sup>1</sup>	The Amazon Resource Number (ARN) of the IAM role assumed by the user to use the dynamically generated temporary security credentials. Set the value of this property if you want to use the temporary security credentials to access the Amazon S3 staging bucket. For more information about how to get the ARN of the IAM role, see the <i>AWS documentation</i> .
Use EC2 Role to Assume Role <sup>1</sup>	Optional. Select the check box to enable the EC2 role to assume another IAM role specified in the IAM Role ARN option. The EC2 role must have a policy attached with a permission to assume an IAM role from the same or different AWS account.
S3 Region Name <sup>1</sup>	The AWS cluster region in which the bucket you want to access resides. Select a cluster region if you choose to provide a custom JDBC URL that does not contain a cluster region name in the JDBC URL connection property.

Property	Description
S3 Service Regional Endpoint	The S3 regional endpoint when the S3 data bucket and the S3 staging bucket need to be accessed through a region-specific S3 regional endpoint. Doesn't apply to Databricks cluster. Default is <code>s3.amazonaws.com</code> .
Zone ID <sup>1</sup>	The zone ID for the Databricks job cluster. Optional for Databricks cluster. Doesn't apply to SQL warehouse. Applies only if you want to create a Databricks job cluster in a particular zone at runtime. For example, <code>us-west-2a</code> . <b>Note:</b> The zone must be in the same region where your Databricks account resides.
EBS Volume Type <sup>1</sup>	The type of EBS volumes launched with the cluster. Optional for Databricks cluster. Doesn't apply to SQL warehouse.
EBS Volume Count <sup>1</sup>	The number of EBS volumes launched for each instance. You can choose up to 10 volumes. Optional for Databricks cluster. Doesn't apply to SQL warehouse. <b>Note:</b> In a Databricks Delta connection, specify at least one EBS volume for node types with no instance store. Otherwise, cluster creation fails.
EBS Volume Size <sup>1</sup>	The size of a single EBS volume in GiB launched for an instance. Optional for Databricks cluster. Doesn't apply to SQL warehouse.
<sup>1</sup> Doesn't apply to mappings in advanced mode.	

## Azure staging environment

The following table describes the properties for the Azure staging environment:

Property	Description
ADLS Storage Account Name	The name of the Microsoft Azure Data Lake Storage account.
ADLS Client ID	The ID of your application to complete the OAuth Authentication in the Active Directory.
ADLS Client Secret	The client secret key to complete the OAuth Authentication in the Active Directory.
ADLS Tenant ID	The ID of the Microsoft Azure Data Lake Storage directory that you use to write data.
ADLS Endpoint	The OAuth 2.0 token endpoint from where authentication based on the client ID and client secret is completed.
ADLS Filesystem Name	The name of an existing file system to store the Databricks Delta data.
ADLS Staging Filesystem Name	The name of an existing file system to store the staging data.

# Configure proxy settings

If your organization uses an outgoing proxy server to connect to the Internet, the agent connects to Informatica Intelligent Cloud Services through the proxy server.

You can configure the Secure Agent and the serverless runtime environment to use the proxy server on Windows and Linux.

You can use only an unauthenticated proxy server to connect to Informatica Intelligent Cloud Services. You can configure proxy settings only when you use the AWS staging environment.

You cannot configure proxy settings for mappings in advanced mode.

To configure the proxy settings for the Secure Agent, perform the following tasks:

- Configure the Secure Agent through the Secure Agent Manager on Windows or shell command on Linux. For instructions, see "Configure the proxy settings on Windows" or "Configure the proxy settings on Linux," in *Getting Started* in the Data Integration documentation.
- Configure the JVM options for the DTM in the Secure Agent properties. For instructions, see the [Proxy server settings](#) Knowledge Base article.

**Note:** If you enable both HTTP and SOCKS proxies, SOCKS proxy is used by default. If you want to use HTTP proxy instead of SOCKS proxy, set the value of the `DisableSocksProxy` property to `true` in the System property.

To configure the proxy settings for the serverless runtime environment, see *Runtime Environments* in the Administrator help.

# JDBC URL parameters

You can utilize the additional JDBC URL parameters field in the Databricks Delta connection to customize and set any additional parameters required to connect to Databricks Delta.

You can configure the following properties as additional JDBC URL parameters in the Databricks Delta connection:

- To pass the unity catalog information to Databricks Delta, specify the catalog name after the SQL warehouse cluster ID in the following format:  
`jdbc:spark://<Databricks Host>:443/  
default;transportMode=http;ssl=1;AuthMech=3;httpPath=/sql/1.0/endpoints/<SQL endpoint  
cluster ID>;ConnCatalog=<catalog_name>;`
- To connect to Databricks Delta using the proxy server, enter the following parameters:  
`jdbc: spark://<Databricks Host>:443/  
default;transportMode=http;ssl=1;AuthMech=3;httpPath=/sql/1.0/warehouses/  
219fe3013963cdce;UseProxy=<Proxy=true>;ProxyHost=<proxy host IPaddress>;ProxyPort=<proxy  
server port number>;ProxyAuth=<Auth_true>;`
- To connect to SSL-enabled Databricks Delta, specify the value in the JDBC URL in the following format:  
`jdbc:spark://<Databricks Host>:443/  
default;transportMode=http;ssl=1;AuthMech=3;httpPath=/sql/1.0/endpoints/<SQL endpoint  
cluster ID>;`

# Rules and guidelines for personal staging location

When you select the personal staging location as a staging environment, the data is first staged in a java temporary location and then copied to a personal staging location of the unity catalog. Both the staged files will be deleted after the mapping runs successfully.

However, to stage the data in a different directory, configure the DTM property `-Djava.io.tmpdir=/my/dir/path` in the JVM options in the system configuration settings of the Administrator service.

To enable data staging in a different directory, you should have read and write permission and enough disk space to stage the data in the directory.

When you specify a personal staging location in the Databricks Delta connection properties for staging, consider the following rules and guidelines:

- You can only specify unity enabled catalog in the SQL warehouse JDBC URL.
- All mappings that are configured run without SQL ELT optimization.
- The data is staged in the folder `stage://tmp/<user_name>` where the `<user_name>` is picked from the Databricks token provided in the connection and requires read and write access to the personal staging location in root location of AWS and Azure.

# Private links to access Databricks Delta

You can access Databricks Delta using Azure Private Link endpoints.

To connect to the Databricks Delta account over the private Azure network, see [Secure connectivity to Azure Data Services](#).

## CHAPTER 3

# Mappings for Databricks Delta

When you configure a mapping, you describe the flow of data from the source to the target.

A mapping defines reusable data flow logic that you can use in mapping tasks.

When you create a mapping, you define the Source, Target, and Lookup transformations to represent a Databricks Delta object. Use the Mapping Designer in Data Integration to add the Source, Target, or Lookup transformations in the mapping canvas and configure the Databricks Delta source, target, and lookup properties.

After you create a mapping, you can run the mapping or you can deploy the mapping in a mapping task. The mapping task allows you to process data based on the data flow logic defined in a mapping.

In advanced mode, the Mapping Designer updates the mapping canvas to include transformations and functions that enable advanced functionality.

You can use Monitor to monitor the jobs.

The following table lists the transformations that are supported by SQL warehouse and Databricks cluster:

Property	SQL warehouse	Databricks cluster
Source	Yes	Yes
Target	Yes	Yes
Filter	Yes	Yes
Lookup	Yes	No
Sorter	Yes	Yes
SQL	Yes	No

## Before you begin

Before you configure and run a mapping, complete the required prerequisites.

## Verify permissions

Permissions define the level of access for the operations that you can perform in Databricks Delta.

The following table lists the permission that you need to run SQL commands to perform specific operations in Databricks Delta:

SQL command	Description
CREATE TABLE	Create a managed table.
CREATE VIEW	Create a view.
CREATE TABLE ... LOCATION <EXTERNAL LOCATION>	Create an external table.
DROP TABLE	Delete a table.
DROP VIEW	Delete a view.
ALTER TABLE	Add a column.
SELECT FROM	Read from a table.
INSERT INTO	Use Insert operation with Write Disposition option set to Append.
INSERT OVERWRITE INTO	Use Insert operation with Write Disposition option set to Truncate.
MERGE INTO	Use Update, Upsert, and Delete operations.
DELETE FROM	Empty a table.
COPY INTO	Load data into target (ecosystem pushdown with CSV as source).
SELECT ON FILE	<p>Retrieve data from a file staged in either Amazon S3, ADLS Gen2, or a personal staging location using the SELECT FROM command.</p> <p>Here are examples of the staging locations where you use the SELECT FROM command to retrieve data:</p> <p><b>Read from a file in Amazon S3:</b></p> <pre>INSERT INTO `tgt_table` (col1, col2) SELECT col1, col2 FROM parquet.`s3a://staging.bucket/infra_tmp-20231120_062134_910212391`</pre> <p><b>Read from a file in ADLS Gen2:</b></p> <pre>INSERT INTO `tgt_table` (col1, col2) SELECT col1, col2 FROM parquet.`stage://tmp/username/infra_tmp-20231120_025852_313310201`</pre> <p><b>Read from a file in personal staging location:</b></p> <pre>INSERT INTO `tgt_table` (col1, col2) SELECT col1, col2 FROM parquet.`abfss://blob@storage_account.dfs.core.windows.net/ infra_tmp-20231120_070825_17050631/`</pre>

### Personal staging location

To use the personal staging location as a staging environment, ensure that you have access to the personal staging location APIs.



# Sources for Databricks Delta

Add a Source transformation to extract data from a source.

When you add a Source transformation to a mapping, you define the source connection, source objects, and source properties related to the Databricks Delta connection type.

The following table lists the source properties that are supported by SQL warehouse and Databricks cluster:

Property	SQL warehouse	Databricks cluster
Source Type - Single Object	Yes	Yes
Source Type - Parameter	Yes	Yes
Source Type - Query	Yes	No
Source Type - Multiple Objects	Yes	No
Parameter	Yes	Yes
Filter	Yes	Yes
Sort	No	No
Database Name	Yes	Yes
Table Name	Yes	Yes
Pre SQL	Yes	No
Post SQL	Yes	No
SQL Override	Yes	No
Staging Location	Yes	Yes
Job Timeout	No	Yes. Applies only to job cluster.
Job Status Poll Interval	No	Yes. Applies only to job cluster.
DB REST API Timeout	No	Yes. Applies only to job cluster.
DB REST API Retry Interval	No	Yes. Applies only to job cluster.
Tracing Level	Yes	Yes
Key Range Partitioning	Yes	No

## Source properties for Databricks Delta

In a mapping, you can configure a Source transformation to represent a Databricks Delta object.

The following table describes the Databricks Delta source properties that you can configure in a Source transformation:

Property	Description
Connection	Name of the source connection. Select a source connection or click <b>New Parameter</b> to define a new parameter for the source connection. <b>Note:</b> You can completely parameterize a parameter file for a source connection only for a single object source type. Parameterization doesn't apply to mappings in advanced mode.
Source Type	Type of the source object. Select any of the following source objects: <ul style="list-style-type: none"> <li>- Single Object</li> <li>- Multiple Objects<sup>1</sup>. You can only use advanced relationships with multiple objects.</li> <li>- Query.</li> <li>- Parameter<sup>1</sup>. Select <b>Parameter</b> to define the source type when you configure the task.</li> </ul> Multiple objects and query source types don't apply to Databricks cluster. <b>Note:</b> Multi-object database override will override the database for all imported objects, while the table override will only override the first table of the multi-object source.
Object	Name of the source object. You cannot use the data preview option if the source fields contain hierarchical data types.
Query <sup>1</sup>	Click on <b>Define Query</b> and enter a valid custom query. The <b>Query</b> property appears only if you select <b>Query</b> as the source type. You can parameterize a custom query object at runtime in a mapping. You can also enable unity catalog settings in a custom query to access a table within a particular catalog.
<sup>1</sup> Doesn't apply to mappings in advanced mode.	

The following table describes the Databricks Delta query options that you can configure in a Source transformation:

Property	Description
Query Options	<p>Filters the source data based on the conditions you specify. Click <b>Configure</b> to configure a filter option.</p> <p>The Filter option filters records and reduces the number of rows that the Secure Agent reads from the source. Add conditions in a read operation to filter records from the source. You can specify the following filter conditions:</p> <ul style="list-style-type: none"> <li>- Not parameterized. Use a basic filter to specify the object, field, operator, and value to select specific records.</li> <li>- Completely parameterized*. Use a parameter to specify the filter query.</li> <li>- Advanced. Use an advanced filter to define a complex filter condition.</li> </ul> <p><b>Note:</b> You can use Contains, Ends With, and Starts With operators to filter records only on SQL endpoints.</p>
Filter	<p>Filters records based on the filter condition.</p> <p>You can specify a simple filter or an advanced filter.</p>

The following table describes the Databricks Delta source advanced properties that you can configure in a Source transformation:

Property	Description
Database Name	<p>Overrides the database name provided in connection and the database name provided during metadata import.</p> <p><b>Note:</b> To read from multiple objects ensure that you have specified the database name in the connection properties.</p>
Table Name	<p>Overrides the table name used in the metadata import with the table name that you specify.</p>
Pre SQL	<p>The pre-SQL command to run on the Databricks Delta source table before the agent reads the data.</p> <p>Doesn't apply to Databricks cluster.</p> <p>For example, if you want to update records in the database before you read the records from the table, specify a pre-SQL statement.</p> <p>The query must include a fully qualified table name. You can specify multiple pre-SQL commands, each separated with a semicolon.</p>
Post SQL	<p>The post-SQL command to run on the Databricks Delta table after the agent completes the read operation.</p> <p>Doesn't apply to Databricks cluster.</p> <p>For example, if you want to delete some records after the latest records are loaded, specify a post-SQL statement.</p> <p>The query must include a fully qualified table name. You can specify multiple post-SQL commands, each separated with a semicolon.</p>

Property	Description
SQL Override	<p>Overrides the default SQL query used to read data from Databricks Delta custom query source.</p> <p>The column names in the SQL override query should match with the column names in the custom query in a SQL transformation.</p> <p><b>Note:</b> The metadata of the source should be the same as SQL override to override the query.</p>
Staging Location	<p>Relative directory path to store the staging files.</p> <ul style="list-style-type: none"> <li>- If the Databricks cluster is deployed on AWS, use the path relative to the Amazon S3 staging bucket.</li> <li>- If the Databricks cluster is deployed on Azure, use the path relative to the Azure Data Lake Store Gen2 staging filesystem name.</li> </ul> <p><b>Note:</b> When you use the unity catalog, a pre-existing location on user's cloud storage must be provided in the Staging Location.</p>
Job Timeout	<p>Maximum time in seconds that is taken by the Spark job to complete processing.</p> <p>Doesn't apply to SQL warehouse.</p> <p>If the job is not completed within the time specified, the Databricks cluster terminates the job and the mapping fails.</p> <p>If the job timeout is not specified, the mapping shows success or failure based on the job completion.</p>
Job Status Poll Interval	<p>Poll interval in seconds at which the Secure Agent checks the status of the job completion.</p> <p>Doesn't apply to SQL warehouse.</p> <p>Default is 30 seconds.</p>
DB REST API Timeout	<p>The Maximum time in seconds for which the Secure Agent retries the REST API calls to Databricks when there is an error due to network connection or if the REST endpoint returns 5xx HTTP error code.</p> <p>Doesn't apply to SQL warehouse.</p> <p>Default is 10 minutes.</p>
DB REST API Retry Interval	<p>The time Interval in seconds at which the Secure Agent must retry the REST API call, when there is an error due to network connection or when the REST endpoint returns 5xx HTTP error code.</p> <p>Doesn't apply to SQL warehouse.</p> <p>This value does not apply to the Job status REST API. Use job status poll interval value for the Job status REST API.</p> <p>Default is 30 seconds.</p>
Tracing Level	<p>Sets the amount of detail that appears in the log file. You can choose terse, normal, verbose initialization, or verbose data.</p> <p>Default is normal.</p>

**Note:** Advanced source properties are not applicable to mappings in advanced mode. Only pre-SQL and post-SQL advanced properties are applicable for custom queries.

## Custom query source type

You can use a custom query as a source object when you use a Databricks Delta connection.

You might want to use a custom query as the source when a source object is large. You can use the custom query to reduce the number of fields that enter the data flow. You can also create a parameter for the source type when you design your mapping so that you can define the query in the Mapping Task wizard.

To use a custom query as a source, select **Query** as the source type when you configure the source transformation and then use valid and supported SQL to define the query.

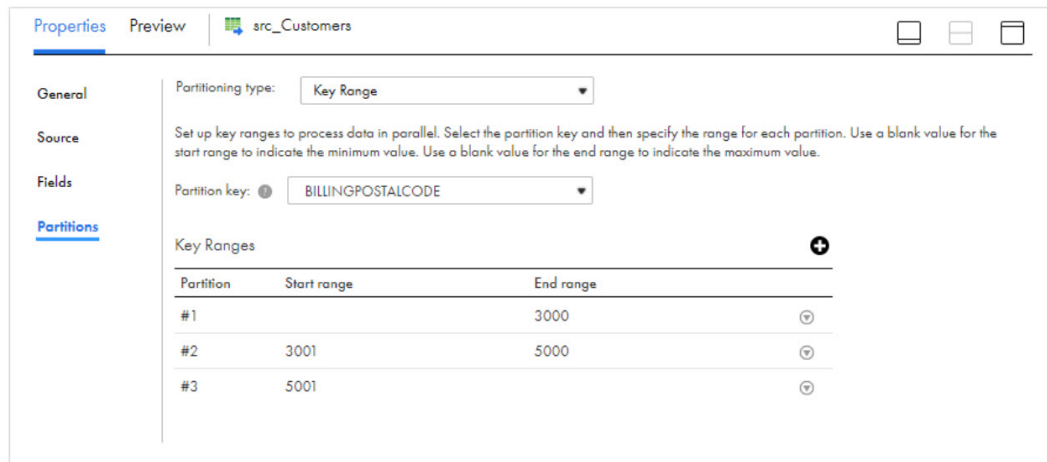
## Key range partitioning

You can configure key range partitioning when you use a mapping to read data from a single Databricks Delta source. The partition type controls how the agent distributes data among partitions at partition points. Partitioning is not applicable for mappings in advanced mode.

Partitioning optimizes the mapping performance at run time. When you run a mapping configured with key range partitioning, the agent distributes rows of source data based on the field that you define as partition keys. The agent compares the field value to the range values for each partition and sends rows to the appropriate partitions.

On the **Partitions** tab for the Source transformation, select key range partitioning and choose the field as the partition key.

You can add key ranges to create partitions, as shown in the following image:



You can configure a partition key for fields of the following data types:

- Int, Bigint, Smallint, and Tiny int
- Timestamp
- Date. Specify the date and time in the following format: YYYY-MM-DD HH24:MI:SS.
- Decimal and Double. Specify the range values as whole numbers.

### Rules and guidelines for key range partitioning

Consider the following rules and guidelines when you use key range partitioning:

- Don't use key range partitioning when you run a mapping in Databricks cluster.
- Don't use key range partitioning when you use Query or Multiple Objects as the source type.
- Don't use key range partitioning in mapping enabled for SQL ELT optimization.

- Don't configure a SQL override with key range partitioning.
- When a column configured as partition key contains null values, the null values are not written to the target.

## Targets for Databricks Delta

Add a Target transformation to write data to a target.

When you add a Target transformation to a mapping, you define the target connection, target objects, and target properties related to the Databricks Delta connection type.

The following table lists the target properties that are supported by SQL warehouse and Databricks cluster:

Property	SQL warehouse	Databricks cluster
Source Type - Single Object	Yes	Yes
Source Type - Parameter	Yes	Yes
Object - Existing	Yes	Yes
Object - Create New at Runtime	Yes	Yes
Create target - Object Name, Table Location, Database Name	Yes	Yes
Create target - Table Properties	Yes	No
Operation - Insert, Update, Upsert, Delete	Yes	Yes
Target Database Name	Yes	Yes
Target Table Name	Yes	Yes
Update Override Query	Yes	No
Write Disposition - Append, Truncate	Yes	Yes
Write Disposition - Truncate Always	Yes	No
Update Mode - Update as update, Update else insert	Yes	Yes
Staging Location	Yes	Yes
Pre SQL	Yes	No
Post SQL	Yes	No
Job Timeout	No	Yes. Applies only to job cluster.
Job Status Poll Interval	No	Yes. Applies only to job cluster.
DB REST API Timeout	No	Yes. Applies only to job cluster.

Property	SQL warehouse	Databricks cluster
DB REST API Retry Interval	No	Yes. Applies only to job cluster.
Forward Rejected Rows	Yes	Yes

## Target properties for Databricks Delta

In a mapping, you can configure a Target transformation to represent a Databricks Delta object.

The following table describes the Databricks Delta properties that you can configure in a Target transformation:

Property	Description
Connection	Name of the target connection. Select a target connection or click <b>New Parameter</b> to define a new parameter for the target connection.
Target Type	Target type. Select one of the following types: <ul style="list-style-type: none"> <li>- Single Object.</li> <li>- Parameter. Select <b>Parameter</b> to define the target type when you configure the task.</li> </ul>
Object	Name of the target object.
Create Target	Creates a target. Enter a name for the target object and select the source fields that you want to use. By default, all source fields are used. You can select an existing target object or create a new target object at runtime. You cannot parameterize the target at runtime.
Operation	Defines the type of operation to be performed on the target table. Select from the following list of operations: <ul style="list-style-type: none"> <li>- Insert (Default)</li> <li>- Update</li> <li>- Upsert</li> <li>- Delete</li> <li>- Data Driven<sup>1</sup></li> </ul> When you use an upsert operation, you must configure the <b>Update Mode</b> in target details as Update else Insert. If the key column gets null value from the source, the following actions take place for different operations: <ul style="list-style-type: none"> <li>- Update. Skips the operation and does not update the row.</li> <li>- Delete. Skips the operation and does not delete the row.</li> <li>- Upsert. Inserts a new row instead of updating the existing row.</li> </ul>
Update Columns	The fields to use as temporary primary key columns when you update, upsert, or delete data on the Databricks Delta target tables. When you select more than one update column, the mapping task uses the AND operator with the update columns to identify matching rows. Applies to update, upsert, delete and data driven operations.

Property	Description
Data Driven Condition <sup>1</sup>	<p>Flags rows for an insert, update, delete, or reject operation based on the expressions that you define.</p> <p>For example, the following IIF statement flags a row for reject if the ID field is null. Otherwise, it flags the row for update:</p> <pre>IIF (ISNULL(ID), DD_REJECT, DD_UPDATE )</pre> <p>Required if you select the data driven operation.</p>
<sup>1</sup> Applies only to mappings in advanced mode.	

The following table describes the Databricks Delta advanced properties that you can configure in a Target transformation:

Advanced Property	Description
Target Database Name <sup>1</sup>	<p>Overrides the database name provided in the connection and the database selected in the metadata browser for existing targets.</p> <p><b>Note:</b> You cannot override the database name when you create a new target at runtime.</p>
Target Table Name <sup>1</sup>	<p>Overrides the table name at runtime for existing targets.</p>
Update Override Query	<p>Overrides the default update query that the agent generates for the update operation specified in this field.</p> <p>Doesn't apply to Databricks cluster.</p> <p>Use the merge command for the update operation.</p>
Write Disposition	<p>Overwrites or adds data to the existing data in a table. You can select from the following options:</p> <ul style="list-style-type: none"> <li>- Append. Appends data to the existing data in the table even if the table is empty.</li> <li>- Truncate. Overwrites the existing data in the table. Only applies to Insert operation and non-empty sources.</li> <li>- Truncate Always. Overwrites the existing data in the table. Applies to insert, update, upsert, and delete target operations for empty and non-empty sources. Doesn't apply to Databricks cluster.</li> </ul> <p><b>Note:</b> You cannot perform <code>Truncate Always</code> when you don't specify the database name in the connection properties and while creating a new target at runtime.</p>
Update Mode <sup>1</sup>	<p>Defines how rows are updated in the target tables. Select from the following options:</p> <ul style="list-style-type: none"> <li>- Update As Update: Rows matching the selected update columns are updated in the target.</li> <li>- Update Else Insert: Rows matching the selected update columns are updated in the target. Rows that don't match are appended to the target.</li> </ul>
Staging Location	<p>Relative directory path to store the staging files.</p> <ul style="list-style-type: none"> <li>- If the Databricks cluster is deployed on AWS, use the path relative to the Amazon S3 staging bucket.</li> <li>- If the Databricks cluster is deployed on Azure, use the path relative to the Azure Data Lake Store Gen2 staging filesystem name.</li> </ul> <p><b>Note:</b> When you use the unity catalog, a pre-existing location on user's cloud storage must be provided in the Staging Location.</p>



Advanced Property	Description
Pre SQL	The pre-SQL command to run before the agent writes to Databricks Delta. For example, if you want to assign sequence object to a primary key field of the target table before you write data to the table, specify a pre-SQL statement. Doesn't apply to Databricks cluster. You can specify multiple pre-SQL commands, each separated with a semicolon.
Post SQL	The post-SQL command to run after the agent completes the write operation. For example, if you want to alter the table created by using create target option and assign constraints to the table before you write data to the table, specify a post-SQL statement. Doesn't apply to Databricks cluster. You can specify multiple post-SQL commands, each separated with a semicolon.
Job Timeout <sup>1</sup>	Maximum time in seconds that is taken by the Spark job to complete processing. Doesn't apply to SQL warehouse. If the job is not completed within the time specified, the Databricks cluster terminates the job and the mapping fails. If the job timeout is not specified, the mapping shows success or failure based on the job completion.
Job Status Poll Interval <sup>1</sup>	Poll interval in seconds at which the Secure Agent checks the status of the job completion. Doesn't apply to SQL warehouse. Default is 30 seconds.
DB REST API Timeout <sup>1</sup>	The Maximum time in seconds for which the Secure Agent retries the REST API calls to Databricks when there is an error due to network connection or if the REST endpoint returns 5xx HTTP error code. Doesn't apply to SQL warehouse. Default is 10 minutes.
DB REST API Retry Interval <sup>1</sup>	The time Interval in seconds at which the Secure Agent must retry the REST API call, when there is an error due to network connection or when the REST endpoint returns 5xx HTTP error code. Doesn't apply to SQL warehouse. This value does not apply to the Job status REST API. Use job status poll interval value for the Job status REST API. Default is 30 seconds.
Forward Rejected Rows	Determines whether the transformation passes rejected rows to the next transformation or drops rejected rows. By default, the agent forwards rejected rows to the next transformation.
<sup>1</sup> Doesn't apply to mappings in advanced mode.	

## Create a target table at runtime

You can use an existing target or create a target to hold the results of a mapping. If you choose to create the target, the agent creates the target if it does not exist already when you run the task.

You can create both managed and external tables for mappings and mappings in advance mode.

To specify the target properties, perform the following tasks:

1. Select the Target transformation in the mapping.

2. To specify the target, click the **Target** tab.
3. Select the target connection.
4. For the target type, choose **Single Object** or **Parameter**.
5. Specify the target object or parameter.
6. To specify a target object, perform the following tasks:
  - a. Click **Select** and choose a target object. You can select an existing target object or create a new target object at runtime.

**Target Object**

Select an existing target object or create a new one. Any new target objects will be created when the mapping task is executed.

Target Object:  Existing  Create New at Runtime

Object Name:

Table Location:

Database Name:

Table Properties:

[?](#)

- b. To create a target object at runtime, select **Create New at Runtime**.
- c. In the **Object Name** field, enter the name of the target table that you want to create. Specify the table name in lowercase.
- d. In the **Table Location** field, enter the location of the target table data.

The table location is relative to the data bucket or data filesystem name specified in the connection. External table is created if you specify the table location.

For unity catalog, specify a pre-existing location to create a new target at runtime.

When you use personal staging location, specify a pre-existing table location with an absolute path to the directory to create a new target at runtime. The following image shows the table location for personal staging location.

Select an existing target object or create a new one. Any new target objects will be created when the mapping task is executed.

Target Object:  Existing  Create New at Runtime

Object Name:

Table Location:

Database Name:

Table Properties:

- e. In the **Database Name** field, specify the Databricks database name.

The database name that you specify in the connection properties takes precedence.

**Note:** If you do not specify the Database Name when you create a new target or in the Database connection attribute, the default database is used to create a new target irrespective of the database name specified in the SQL Warehouse JDBC URL.

- f. Specify the **Table Properties** to optimize the table configuration settings. You can use table properties to tag tables with information that are not tracked by SQL queries. To see the list of table properties and table options, see the Databricks Delta documentation.
- g. Click **OK**.

## Rules and guidelines for create target at runtime

When you configure a mapping with the **Create New at Runtime** option, consider the following rules:

- When a source object consists of Date data type and you use the default create target option in a mapping, the date data gets corrupted. To resolve this issue, navigate to **Edit Metadata** option in the **Target fields** of the target and change **Native Type** to Date.
- When you enable dynamic schema handling in a task and create target at runtime, you must provide the complete path of the target table in the Database Name. Ensure that the table name is in lowercase. For example, database\_name/TABLE/table\_name.
- When you configure a mapping to create a new target at runtime and the source data contains complex fields, the metadata fails to appear on the **Target Fields** tab.
- Hierarchical data types are supported only when you create a new target at runtime and not for existing targets.
- Only **Insert** operation is applicable for hierarchical data types in the Target transformation

## Override the update operation

You can specify an update override to override the update query that the Secure Agent generates for the update operation.

When you configure an update override, the Secure Agent uses the query that you specify, stages the data in files, and then loads that data into a temporary table using the merge command. The data from the temporary table is then loaded to the Databricks target table. The syntax that you specify for the update query must be supported by Databricks Delta:

Specify the update override query in the following format:

```
MERGE INTO <Target table name> AS A USING :TU AS B ON A.<Column1> = B.<Column1> WHEN
MATCHED THEN UPDATE SET A.<Column2> = B.<Column2>, A.<Column3>= B.<Column3> ...
A.<ColumnN>= B.<ColumnN>
```

where, *:TU* represents the incoming data source for the target transformation. The Secure Agent replaces *:TU* with a temporary table name while running the mapping and does not validate the update query.

**Note:** While specifying the name of the target table, you must provide a fully qualified name in the format:

```
<Catalog_Name>.<Database_Name>.<Target_TableName>.
```

When you configure an update override in a mapping to write to Databricks Delta, consider the following rules and guidelines:

- Ensure that the column names for *:TU* matches the target table column names.
- Specify the update query with a valid SQL syntax because Databricks Delta Connector replaces *:TU* with a temporary table name and does not validate the update query.
- Do not change the order of the column in the mappings when you configure the update override option.
- The column names used in the override query must not be unconnected in the target field mapping.
- You can only override a single SQL query at runtime.
- You cannot perform an update override on insert, upsert, delete, data driven, and dynamic schema handling operations.

## Determine the order of processing for multiple targets

When you configure a mapping to write to multiple targets in a single pipeline, you can configure each target for any of the write operations. The order of the target operation is not deterministic.

However, if you want to process the target operations to process in a specific order such as delete, update, and insert, you need to set certain properties either in the Secure Agent or in the task properties.

### Set `-DEnableSingleCommit=true` in the Secure Agent properties

Perform the following tasks to set the property for the Secure Agent:

1. Open Administrator and select **Runtime Environments**.
2. Select the Secure Agent for which you want to set the property.
3. On the upper-right corner of the page, click **Edit**.
4. In the **System Configuration Details** section, select the **Type** as **DTM** for the Data Integration Service.
5. Edit the JVM options and set the property to `-DEnableSingleCommit=true`.

DTM      JVMOption1     

### Set the `EnableSingleCommit` property in the task properties

Perform the following tasks to set the property in the task:

1. On the **Schedule** tab in the mapping task properties, navigate to the Advanced Session Properties section.
2. From the **Session Property Name** list, select **Custom Properties**, and set the Session Property Value to **Yes**.

#### Advanced Session Properties

Session Property Name	Session Property Value
Custom Properties	EnableSingleCommit=Yes

## Lookups for Databricks Delta

Add a Lookup transformation to retrieve data based on a specified lookup condition. When you add a Lookup transformation to a mapping, you define the lookup connection, lookup objects, and lookup properties related to Databricks Delta.

In the Lookup transformation, select the lookup connection and object. Then, define the lookup condition and the outcome for multiple matches.

The mapping queries the lookup source based on the lookup fields and the defined lookup condition. The lookup operation returns the result to the Lookup transformation, which then passes the results downstream.

You can configure the following lookups:

- **Connected.** You can use a cached or uncached connected lookup for mappings. You can also use a dynamic lookup cache to keep the lookup cache synchronized with the target.
- **Unconnected.** You can use a cached lookup. You need to supply input values for an unconnected Lookup transformation from a `:LKP` expression in a transformation that uses an Expression transformation.

Lookup transformation doesn't apply to Databrick cluster.

Lookup objects are not supported for mappings in advanced mode.

For more information about Lookup transformation, see *Transformations* in the Data Integration documentation.

## Lookup properties for Databricks Delta

The following table describes the Databricks Delta lookup object properties that you can configure in a Lookup transformation:

Property	Description
Connection	Name of the lookup connection. You can select an existing connection, create a new connection, or define parameter values for the lookup connection property. If you want to overwrite the lookup connection properties at runtime, select the <b>Allow parameter to be overridden at run time</b> option.
Source Type	Type of the source object. Select Single Object, Query, or Parameter.
Parameter	A parameter file where you define values that you want to update without having to edit the task. Select an existing parameter for the lookup object or click <b>New Parameter</b> to define a new parameter for the lookup object. The <b>Parameter</b> property appears only if you select parameter as the lookup type. If you want to overwrite the parameter at runtime, select the <b>Allow parameter to be overridden at run time</b> option.
Lookup Object	Name of the lookup object for the mapping.
Multiple Matches	Behavior when the lookup condition returns multiple matches. You can return all rows, any row, the first row, the last row, or an error. You can select from the following options in the lookup object properties to determine the behavior: <ul style="list-style-type: none"> <li>- Return first row</li> <li>- Return last row</li> <li>- Return any row</li> <li>- Return all rows</li> <li>- Report error</li> </ul>

The following table describes the Databricks Delta lookup object advanced properties that you can configure in a Lookup transformation:

Advanced Property	Description
Database Name	Overrides the database specified in the connection.
Table Name	Overrides the table specified in the connection.
Staging Location	Relative directory path to store the staging files. <ul style="list-style-type: none"> <li>- If the Databricks cluster is deployed on AWS, use the path relative to the Amazon S3 staging bucket.</li> <li>- If the Databricks cluster is deployed on Azure, use the path relative to the Azure Data Lake Store Gen2 staging filesystem name.</li> </ul>
SQL Qverride	Overrides the default SQL query used to read data from the Databricks Delta custom query source.

Advanced Property	Description
Pre-SQL	SQL statement that you want to run before reading data from the source.
Post-SQL	SQL statement that you want to run after reading data from the source.
Job Timeout	Maximum time in seconds that is taken by the Spark job to complete processing. If the job is not completed within the time specified, the Databricks cluster terminates the job and the mapping fails. If the job timeout is not specified, the mapping shows success or failure based on the job completion.
Job Status Poll Interval	Poll interval in seconds at which the Secure Agent checks the status of the job completion. Default is 30 seconds.
DB REST API Timeout	The Maximum time in seconds for which the Secure Agent retries the REST API calls to Databricks when there is an error due to network connection or if the REST endpoint returns 5xx HTTP error code. Default is 10 minutes.
DB REST API Retry Interval	The time Interval in seconds at which the Secure Agent must retry the REST API call, when there is an error due to network connection or when the REST endpoint returns 5xx HTTP error code. This value does not apply to the Job status REST API. Use job status poll interval value for the Job status REST API. Default is 30 seconds.
Tracing Level	Sets the amount of detail that appears in the log file. You can choose terse, normal, verbose initialization, or verbose data. Default is normal.
Lookup Data Filter <sup>1</sup>	Limits the number of lookups that the mapping performs on the cache of the lookup source table based on the value you specify in the filter condition. This property is applicable when you enable lookup cache on the <b>Advanced</b> tab for the Lookup transformation. Maximum length is 32768 characters. For more information about this property, see <i>Transformations</i> in the Data Integration documentation.
<sup>1</sup> Doesn't apply to mappings in advanced mode.	

## Parameterization

You can parameterize the following properties when you create mappings:

- Source properties. Source type, source connection, query options in source, database name, table name, and advanced properties in the source.
- Target properties. Target type, target connection, target database name, target table name, and target advanced properties.
- Lookup properties. Lookup type, lookup connection, lookup advanced properties.

You can parameterize the following properties when you create mappings in advanced mode:

- Source properties. Source type, source connection, database name, and table name in the source.
- Target properties. Target type, target connection, target database name, and target table name.

## Multiple match options in lookups

When you configure a lookup, you define the behavior when a lookup condition returns more than one match. You can return all rows, any row, the first row, the last row, or an error.

The following configurations have multiple match policy restrictions:

- You cannot parameterize an uncached connected lookup with parameterized prefixes at source.
- When the data is not matched for an uncached connected lookup, incorrect data is logged in the target.

## Enable lookup caching

When you configure a Lookup transformation in a mapping, you can cache the lookup data during the runtime session.

When you select **Lookup Caching Enabled**, Data Integration queries the lookup source once and caches the values for use during the session, which can improve performance. You can configure dynamic lookup caching and can specify the directory to store the cached lookup.

**Note:** You cannot perform insert and update operations to the same target at runtime when you dynamically cache the lookup objects in a mapping.

For information about lookup caching, see the chapter "Lookup transformations" in *Transformations* in the Data Integration documentation.

# SQL transformation

You can configure an SQL transformation in a Databricks Delta mapping to process SQL queries. The SQL transformation needs to be connected to the mapping pipeline.

When you add an SQL transformation to the mapping, on the **SQL** tab, you define the database connection and select the SQL query type for the transformation. The SQL transformation process the entered query that you define in the SQL editor. You can specify functions with a simple SELECT statement. Do not use more than one SQL query in an SQL transformation.

When you configure a SELECT query, specify the column names in the output ports for the functions. For example, when you specify the query `SELECT square(~AGE~), sqrt(~SNAME~)`, specify two output columns for the AGE and SNAME functions.

When you run the mapping, the SQL transformation processes the query and returns the rows. The SQL transformation also returns any errors that occur from the underlying database or from a user syntax that is not valid. When an SQL error occurs, the error is logged to the `SQLException` field, by default.

**Note:** SQL transformation doesn't apply to Databricks cluster.

For more information about SQL queries, see *Transformations* in the Data Integration documentation.

### Rules and guidelines for SQL transformation

- You cannot configure a stored procedure in an SQL transformation.
- The SQL query must be a simple SELECT statement without FROM and WHERE arguments.
- You cannot configure a parameterized query in an SQL transformation.
- You cannot include special characters in the query.

- You need to manually specify the output ports, including the numrowsAffected and passthrough ports that are used downstream in the transformation.

## Dynamic schema handling

You can choose how Data Integration handles changes that you make to the data object schemas. To refresh the schema every time the mapping task runs, you can enable dynamic schema handling in the task.

Configure schema change handling on the **Schedule** page when you configure the task.

The following table describes the schema change handling options:

Option	Description
Asynchronous	Default. Data Integration refreshes the schema when you edit the mapping or mapping task, and when Informatica Intelligent Cloud Services is upgraded.
Dynamic	Data Integration refreshes the schema every time the task runs. You can choose from the following options to refresh the schema: <ul style="list-style-type: none"> <li>- <b>Alter and apply changes.</b> Data Integration applies the following changes from the source schema to the target schema: <ul style="list-style-type: none"> <li>- New fields. Alters the target schema and adds the new fields from the source.</li> </ul> </li> <li>- <b>Don't apply DDL changes.</b> Data Integration does not apply the schema changes to the target.</li> <li>- <b>Drop current and recreate.</b> Drops the existing target table and then recreates the target table at runtime using all the incoming metadata fields from the source.</li> </ul>

**Note:** You cannot enable dynamic schema handling in a task to run a mapping in advanced mode. Dynamic schema handling doesn't apply to Databricks cluster.

For more information, see the "Schema change handling" topic in *Tasks* in the Data Integration help.

## IDENTITY columns

You can use the IDENTITY or generated columns in mappings and mapping tasks.

When you write to a table with IDENTITY or generated columns, the values in the columns are automatically generated based on a user-specified function in the Databricks Delta table.

You can use the IDENTITY column for insert, update, upsert, and delete operations. The query can be used for connected and unconnected ports for matching the columns based on the query in the CREATE TABLE statement.



The following table shows the behavior of the operations for each type of generated column:

Type of query	INSERT [INSERT INTO]	UPDATE [Merge]	UPSERT [Merge]	DELETE [Merge]
GENERATED ALWAYS AS	Port needs to be unconnected	Can be used for matching	Port needs to be unconnected and cannot be used for matching	Can be used for matching
GENERATED ALWAYS AS IDENTITY	Port needs to be unconnected	Can be used for matching	Port needs to be unconnected and cannot be used for matching	Can be used for matching
GENERATED BY DEFAULT AS IDENTITY	Port can be connected or unconnected	Can be used for matching	Port can be connected and can be used for matching	Can be used for matching

**Note:** You cannot define an IDENTITY column when you create a new target at runtime and the source table has generated columns. IDENTITY columns don't apply to Databricks cluster.

## Mappings in advanced mode example

You work for a retail company that offers more than 50,000 products and the stores are distributed across the globe. The company ingests a large amount of customer engagement details from the transactional CRM system into Amazon S3.

The sales team wants to improve customer engagement and satisfaction at every touch point. To create a seamless customer experience and deliver personalized service across the various outlets, the retail company plans to load the data that is stored in the Amazon S3 bucket to Databricks Delta.

You can create a mapping that runs on an advanced cluster to achieve faster performance when you read data from the Amazon S3 bucket and write data to the Databricks Delta target.

You can choose to add transformations to process the raw data that you read from the Amazon S3 bucket and then write the curated data to Databricks Delta.

The following example illustrates how to create a mapping in advanced mode to read from an Amazon S3 source and write to Databricks Delta target:

1. In Data Integration, click **New > Mappings > Mapping**.
2. In the Mapping Designer, click **Switch to Advanced**.  
The Mapping Designer updates the mapping canvas to display the transformations and functions that are available in advanced mode.
3. Enter a name, location, and description for the mapping.
4. Add a Source transformation, and specify a name and description in the general properties.
5. On the **Source** tab, perform the following steps to read data from the Amazon S3 source:
  - a. In the **Connection** field, select the Amazon S3 V2 connection.
  - b. In the **Source Type** field, select single object as the source type.

- c. In the **Object** field, select the parquet file object that contains the customer details.
  - d. In the **Advanced Properties** section, specify the required parameters.
6. On the **Expression** tab, define an expression to change the file name port of the customer parquet file to uppercase based on your business requirement before you write data to the Databricks Delta target.
7. Add a Target transformation, and specify a name and description in the general properties.
8. On the **Target** tab, specify the details to write data to Databricks Delta:
  - a. In the **Connection** field, select the Databricks Delta target connection.
  - b. In the **Target Type** field, select single object.
  - c. In the **Object** field, select the Databricks Delta object to which you want to write the curated customer engagement data.
  - d. In the **Operation** field, select the insert operation.
  - e. In the **Advanced Properties** section, specify the required advanced target properties.
9. Click **Save > Run** to validate the mapping.
 

In Monitor, you can monitor the status of the logs after you run the task.

## Rules and guidelines for mappings

Consider the following rules and guidelines for Databricks Delta objects used as sources, targets, and lookups in mappings:

- You cannot use input or in-out parameters in a parameter file to parameterize multiple objects or parameterize the relationships between the objects.
- A mapping with Null values in an uncached lookup condition generates incorrect results.
- When you do not specify the database name in Databrick Delta connection and read multiple objects with the same table name from different databases, you must append the database name for each object in the advanced relationship.
- When you specify `SESSSTARTTIME` variable in a query in a mapping task to return the Datetime values, specify the query in the following format:
 

```
:select to_timestamp('$$$SESSSTARTTIME', 'MM/dd/yyyy HH:mm:ss.SSSSS') as t;
```
- When you run multiple concurrent mappings to write data to Databricks Delta targets, a transaction commit conflict error might occur and the mappings might fail.
- View objects are displayed in the Table panel instead of the View panel while importing a Databricks Delta object. This issue occurs when the Databricks cluster is deployed on AWS cloud service.
- To avoid java heap space error when you read or write complex files, set the JVM options for type DTM to increase the `-Xms` and `-Xmx` values in the system configuration details of the Secure Agent. The recommended values for `-Xms` is 512 MB and `-Xmx` is 1024 MB.
- When you import views, the Select Source Object dialog box does not display view objects.
- When you test the Databricks Delta connection, the Secure Agent does not validate the values you specify in the Org ID connection parameter.
- You cannot use the Hosted Agent as a runtime environment when you configure a mapping to run on the SQL warehouse to read or write data that contains unicode characters.

- The number of clusters that the Secure Agent creates to run the mapping depends on the number of Databricks Delta connections used in the transformations in a mapping. For example, if multiple transformations use the same Databricks Delta connection, the mapping runs on a single cluster.
- When you keep the mapping designer idle for more than 15 minutes, the metadata fetch throws an exception.
- If you change the database name in the connection and run the existing mappings, the mappings start failing. After you change the database name in the connection, you must reimport the objects in the existing mappings before you run the mappings.
- Use the following formats to run the mapping successfully, when you import a Databricks Delta source object containing Date or Boolean data types with a simple source filter conditions:
  - Boolean = 0 or 1
  - Date = YYYY-MM-DD HH24:MM:SS.US
- When you run a mapping with source column data type as string containing TRUE / FALSE value and write data to target with Boolean data type column of a Databricks Delta table, the Secure Agent writes data as 0 to the target.
- When the Databricks analytics cluster is down and you perform a test connection or import an object, the connection is timed out after 10 minutes.
- When you parameterize the source or target connection in a mapping and you do not specify the database name, ensure that you specify the database name in lowercase when you assign a default value for the parameter.
- When you parameterize the source filter condition or any expressions in a mapping, ensure that you specify the table name in lowercase when you add the source filter condition or the expression in the mapping task. Otherwise, the Secure Agent throws the following exception:
 

```
Invalid expression string for filter condition
```
- When you run a mapping to write data to a Databricks Delta target using create target at runtime and the target table already exists, ensure that the target table schema is same. Otherwise, the mapping fails.
- When you run a mapping to write data to multiple Databricks Delta targets that use the same Databricks Delta connection and the Secure Agent fails to write data to one of targets, the mapping fails and the Secure Agent does not write data to the remaining targets.
- When you use the `Create New at Runtime` option to create a Databricks target, you can parameterize only the target connection and the table name using a parameter file. You cannot parameterize other properties such as `Path` or `DBname`.
- The pre-SQL and post-SQL commands run non-linearly. In the session logs, you will see that the target pre-SQL cases are executed before the source pre-SQL queries.
- When you run pre-SQL and post-SQL commands to read from sources that have semi-colons within the query, the mappings fails. The queries can only have semi-colons at the end.
- When you read or write Unicode data to Databricks on the SQL endpoint, you need to set some properties for the Secure Agent before you run the mapping. You can perform one of the following tasks:
  - Set the following environment variables in the Secure Agent machine, and then restart the Secure Agent.
 

```
export LANGUAGE="en_US.UTF-8"
export LC_ALL="en_US.UTF-8"
```
  - Configure the property `-Dfile.encoding=UTF-8` in the JVM options in the Secure Agent properties.
- Note:** You cannot read or write Unicode data when you use the Hosted Agent in the connection.
- When you configure a mapping where you have staged data in the Personal Staging Location, the temporary data is not deleted when the mapping stops abruptly.

- When you parameterize the source or lookup object, ensure that the column names in the source object and lookup object are not the same. Else, the mapping fails.

## Rules and guidelines for mappings in advanced mode

Consider the following rules and guidelines for Databricks Delta objects used as sources and targets in mappings in advanced mode:

- When you write data to multiple Databricks Delta targets with the same table and configure different target operations for each target, the mapping throws a concurrent append exception.
- When mappings in advanced mode reads NULL values from the source and updates or upserts a column with NOT NULL constraint in the Databricks Delta target table, the mapping fails and the Secure Agent fails to log an appropriate error message.
- When you read data from or write data to Databricks Delta and the source or target object contains 5000 or more columns, the mapping fails.
- A mapping in advanced mode configured to read from or write to Databricks Delta fails in the following cases:

- Data is of the Date data type and the date is less than 1582-10-15.
- Data is of the Timestamp data type and the timestamp is less than 1900-01-01T00:00:00Z.

To resolve this issue, specify the following spark session properties in the mapping task or in the custom properties file for the Secure Agent:

```
- spark.sql.legacy.timeParserPolicy=LEGACY
- spark.sql.parquet.int96RebaseModeInWrite=LEGACY
- spark.sql.parquet.datetimeRebaseModeInWrite=LEGACY
- spark.sql.parquet.int96RebaseModeInRead=LEGACY
- spark.sql.parquet.datetimeRebaseModeInRead=LEGACY
- spark.sql.avro.datetimeRebaseModeInWrite=LEGACY
- spark.sql.avro.datetimeRebaseModeInRead=LEGACY
```

- When you do a data type conversion from Float to Double or use create target at run time, data loss is encountered.
- If a mapping in advanced mode has a source column data type as String containing true or false value and a target data type as Boolean, the Secure Agent writes data as null to the target.
- Use the following formats when you import a Databricks Delta source object containing Boolean, Date, or Timestamp data types with a simple source filter conditions:
  - Boolean = 0 or 1
  - Date = YYYY-MM-DD HH24:MM:SS.US
  - Timestamp = YYYY-MM-DD HH24:MM:SS.US
- You cannot use the following features:
  - View
  - Multipipe

- After you create and run a mapping configuration task, it is recommended to shut down the job cluster. If you modify a mapping task or edit the connection linked to a mapping task, metadata is fetched again and the job cluster restarts.
- When you do a data type conversion from Date or Timestamp to String, the Secure Agent writes the value only in the following default format for both Date and Timestamp:  
MM/DD/YYYY HH24:MI:SS
- When you do a data type conversion from String to Date or Timestamp, the String value must be in the following format:  
MM/DD/YYYY HH24:MI:SS  
To use any other format, you must specify the format in the advanced session property of a mapping task for successful conversion. Null is populated in the target for the unmatched format.
- When you do a data type conversion from Bigint to Double, the target data is written in the exponential format.
- When you perform an update, upsert, or a data driven operation with an IIF condition that includes DD\_DELETE or DD\_UPDATE, ensure that the update column that you specified does not have duplicate rows. Otherwise, the mapping fails with the following error:  
java.lang.UnsupportedOperationException: Cannot perform MERGE as multiple source rows matched and attempted to update the same target row in the Delta table.
- When you perform an insert, update, upsert operation, or DD\_UPDATE and the range of the data in source column is greater than the range of the target column, the mapping does not fail and leads to data truncation.
- When you specify a single constant in the data driven condition, the mapping ignores the data driven condition and the Secure Agent performs insert, update, or delete operation based on the constant. For example, if you specify the data driven condition as DD\_INSERT, the mapping does not consider the update columns and depends on the Write Disposition property.
- When you specify a single constant with the IIF condition in the data driven condition such as IIF(COL\_INT > 20 , DD\_UPDATE), the Secure Agent inserts the data into the target even for those rows that do not satisfy the condition.
- When you specify the DD\_REJECT constant in the data driven condition, the Secure Agent does not log the rejected rows in the error file or the session log.
- You can run mappings with hierarchical data types only on a Linux system.
- When you configure mappings with nested statements that contain hierarchical data types, the mapping fails if the nested field names contain the following special characters:  
- ,  
- (  
- < >
- When you read and write data of hierarchical data type, ensure that the column names do not start with a number.
- When you run a Databricks Delta mapping in advanced mode to write data to multiple Databricks Delta targets and enable SQL ELT optimization, the mapping runs successfully but the data is written to only one target.
- When you read data of hierarchical data type, ensure that the column names do not contain special characters.

## CHAPTER 4

# Databricks Delta SQL ELT optimization

When you run a task configured for SQL ELT optimization, the task converts the transformation logic to an SQL query. The task sends the query to the database, and the database executes the query.

The amount of transformation logic that you can push to the database depends on the database, transformation logic, and task configuration. The Secure Agent processes all transformation logic that it cannot push to the database.

Configure SQL ELT optimization for a mapping in the tasks properties. You can configure SQL ELT optimization in task that references a mapping or a mapping in advanced mode.

**Note:** You can only configure SQL ELT optimization when you use the SQL warehouse to connect to Databricks Delta.

## SQL ELT optimization types

When you apply SQL ELT optimization, the task pushes transformation logic to the source or target database based on the SQL ELT optimization type you specify in the task properties. Data Integration translates the transformation logic into SQL queries or Databricks Delta commands to the Databricks Delta database. The database runs the SQL queries or Databricks Delta commands to process the transformations.

You can configure the following SQL ELT optimization types in a mapping:

### **None**

The task does not push down the transformation logic to the Databricks Delta database.

### **Full**

The task pushes as much of the transformation logic as possible to process in the Databricks Delta target database.

Data Integration analyses all the transformations from the source to the target. If all the transformations are compatible in the target, it pushes the entire mapping logic to the target. If it cannot push the entire mapping logic to the target, Data Integration first pushes as much transformation logic to the source database and then pushes as much transformation logic as possible to the target database.

When a transformation is not supported in the mapping, the task partially pushes down the mapping logic to the point where the transformation is supported for SQL ELT optimization. However, this applies only to SQL warehouse.

When you enable full SQL ELT optimization, you can determine how Data Integration handles the job when SQL ELT optimization does not work in the **Fallback Option** menu.

#### Source

The task pushes as much as the transformation logic as possible to process in the Databricks Delta source database. This applies only to SQL warehouse.

**Note:** You cannot enable source SQL ELT optimization for mappings in advanced mode.

## Previewing SQL ELT optimization

Before you can run a mapping task configured for SQL ELT optimization, you can preview if SQL ELT optimization is possible when you create the mapping. You can preview from the **SQL ELT optimization** panel in the Mapping Designer.

After you select the required SQL ELT optimization options and run the preview, Data Integration creates and runs a temporary SQL ELT preview mapping task. When the job completes, Data Integration displays the SQL queries to be executed and any warnings in the **SQL ELT optimization** panel. The warning messages help you understand which transformations in the configured mapping are not applicable for SQL ELT optimization. If SQL ELT optimization fails, Data Integration lists any queries generated up to the point of failure. You can edit the mapping and fix the required transformations before you run the mapping for SQL ELT optimization.

You can also view the temporary job created under **My Jobs** and download the session log to view the queries generated.

For more information about how to preview SQL ELT optimization, see the topic "SQL ELT optimization preview" in *Mappings* in the Data Integration documentation.

## Configuring SQL ELT optimization

To optimize a mapping, add the mapping to a task, and then configure SQL ELT optimization in the mapping task.

1. Create a mapping task.
2. In the **SQL ELT optimization** section on the **Schedule** tab, set the SQL ELT optimization value to **Full** or **To Source**.
3. If full SQL ELT optimization is not available, select how Data Integration handles SQL ELT optimization in the **Fallback Option** menu:
  - Partial PDO. Default. Data Integration pushes as much transformation logic as possible to the source and target database. The task processes any transformation logic that it can't push to a database. You can use Partial PDO only when you read from and write to Databrick Delta.
  - Non PDO. The task runs without SQL ELT optimization.
  - Fail Task. Data Integration fails the task.

**Note:** The fallback options are not applicable to mappings in advanced mode.

When you run the mapping task, the transformation logic is pushed to the Databricks Delta database.

# SQL ELT optimization using a Databricks Delta connection

You can configure SQL ELT optimization for a mapping that contains a Databricks Delta connection. SQL ELT optimization enhances the mapping performance. You can configure full or source SQL ELT optimization when you read data from a Databricks Delta source and write to a Databricks Delta target.

## Read from and write to Databricks Delta

You can configure SQL ELT optimization in a mapping to read from and write to Databricks Delta using a Databricks Delta connection.

### Example

You work in a motorbike retail company with more than 30,000 dealerships and 2000 inspection centers globally. The company stores millions of records in Databricks Delta hosted on Azure. You want to use Data Integration to perform some transformations on the data before you write back to Databricks Delta.

Use a Databricks Delta connection in the mapping to read from the Databricks Delta source and write the processed data to the Databricks Delta target. Configure full SQL ELT optimization in the mapping to enhance the performance.

## Read from Amazon S3 and write to Databricks Delta

You can configure SQL ELT optimization for a mapping that uses an Amazon S3 V2 connection in the Source transformation to read from Amazon S3 and a Databricks Delta connection in the Target transformation to write to Databricks Delta.

### Example

You work for a healthcare organization. Your organization offers a suite of services to manage electronic medical records, patient engagement, telephonic health services, and care coordination services. The organization uses infrastructure based on Amazon Web Services and stores its data on Amazon S3. The management plans to load data to a data warehouse to perform healthcare analytics and create data points to improve operational efficiency. To load data from an Amazon S3 based storage object to Databricks Delta, you must use ETL and ELT with the required transformations that support the data warehouse model.

Use an Amazon S3 V2 connection to read data from a file object in an Amazon S3 source and a Databricks Delta connection to write to a Databricks Delta target. Configure full SQL ELT optimization in the mapping to optimize the performance.

## Read from Microsoft Azure Data Lake Storage Gen2 and write to Databricks Delta

You can configure SQL ELT optimization for a mapping that uses an Microsoft Azure Data Lake Storage Gen2 connection in the Source transformation to read from Microsoft Azure Data Lake Storage Gen2 and a Databricks Delta connection in the Target transformation to write to Databricks Delta.

### Example

You want to load data from an Microsoft Azure Data Lake Storage Gen2 based storage object to Databricks Delta for analytical purposes. You want to transform the data before it is made available to users. Use an Microsoft Azure Data Lake Storage Gen2 connection to read data from a Microsoft Azure Data Lake Storage Gen2 source and a Databricks Delta connection to write to a Databricks Delta target. Configure full SQL ELT



optimization in the mapping task to optimize the performance of loading data to Databricks Delta. SQL ELT optimization enhances the performance of the task and reduces the cost involved.

## SQL ELT compatibility

You can configure the task to push transformations, functions, and operators to the database.

When you use SQL ELT optimization, the Secure Agent converts the expression in the transformation by determining equivalent operators and functions in the database. If there is no equivalent operator and function, the Secure Agent processes the transformation logic.

### Functions with Databricks Delta

When you use SQL ELT optimization, Data Integration converts the expression in the transformation by determining equivalent functions in the database. If there is no equivalent function, Data Integration processes the transformation logic.

The following table summarizes the availability of SQL ELT functions that you can push to the Databricks Delta database by using full or source SQL ELT optimization:

Function	Function
ABS()	MIN()
ADD_TO_DATE()	MOD()
ASCII()	POWER()
AVG()	RAND()
CEIL()	REG_EXTRACT()
CHR()	REG_MATCH()
CONCAT()	REG_REPLACE()
COS()	REPLACESTR()
COSH()	REPLACECHR()
COUNT()	REVERSE()
CRC32()	ROUND(DATE)
DATE_COMPARE()	ROUND(NUMBER)
DATE_DIFF()	RPAD()
DECODE()	RTRIM()
EXP()	SET_DATE_PART()
FIRST()	SHA256()
FLOOR()	SIN()

Function	Function
GET_DATE_PART()	SINH()
GREATEST()	SQRT()
IIF()	STDDEV()
IN()	SUBSTR()
INDEXOF()	SUM()
INITCAP()	SYSDATE()
INSTR()	SYSTEMSTAMP()
IS_DATE()	TAN()
IS_NULL()	TANH()
IS_NUMBER()	TO_BIGINT
IS_SPACES()	TO_CHAR(DATE)
LAST()	TO_CHAR(NUMBER)
LAST_DAY()	TO_DATE()
LENGTH()	TO_DECIMAL()
LN()	TO_FLOAT()
LOWER()	TO_INTEGER()
LPAD()	TRUNC(DATE)
LTRIM()	TRUNC(NUMBER)
MAKE_DATE_TIME()	UPPER()
MAX()	VARIANCE()
MD5()	

## Operators with Databricks Delta

When you use SQL ELT optimization, the Secure Agent converts the expression in the transformation by determining equivalent operators in the database. If there is no equivalent operator, the Secure Agent processes the transformation logic.

The following table lists the SQL ELT operators that you can push to Databricks Delta:

Operator	Operator
+	=
-	>=
*	<=
/	!=
%	AND
	OR
>	NOT
<	

## Variables with Databricks Delta

You can use full SQL ELT optimization to push the SESSSTARTTIME and SYSDATE variable to the Databricks Delta database.

SYSDATE is stored as a transformation date/time datatype variable. To return a static date and time, use the SESSSTARTTIME variable.

## Transformations with Databricks Delta

When you configure SQL ELT optimization, the Secure Agent tries to push the configured transformation to Databricks Delta.

You can use full or source SQL ELT optimization to push the following transformations to Databricks Delta:

- Aggregator
- Expression
- Filter
- Joiner
- Lookup
- Sorter
- Union
- Router. Doesn't apply to source SQL ELT optimization.
- Rank
- SQL

**Note:** Expression, Filter, Lookup, Rank, and SQL transformations don't apply to mappings in advanced mode.

## Aggregator transformation

You can configure full SQL ELT optimization to push an Aggregator transformation to process in Databricks Delta.

### Aggregate calculations

You can perform the following aggregate calculations:

- AVG
- COUNT
- FIRST
- LAST
- MAX
- MIN
- SUM
- STDDEV
- VARIANCE

### Incoming ports

When you configure an Aggregator transformation and the incoming port is not used in an aggregate function or in a group by field, the output is not deterministic as the ANY\_VALUE() function returns any value from the port.

You can pass only single arguments to the LAST, STDDEV, and VARIANCE functions.

## Lookup transformation

You can configure full SQL ELT optimization to push a Lookup transformation to process in Databricks Delta. This applies to both connected and unconnected lookups.

You can add the following lookups:

- Cached
- Uncached
- Unconnected with cached

When you configure a connected lookup, select the **Multiple Matches** property value as **Return all rows** in the lookup properties for SQL ELT optimization to work.

You can nest the unconnected lookup function with other expression functions.

When you configure an unconnected Lookup transformation, consider the following rules:

- You must select the **Multiple Matches** property value as **Report error** in the unconnected lookup properties for SQL ELT optimization to work.
- You can only configure an Expression transformation for an output received from an unconnected lookup.

## SQL Transformation

You can use an SQL transformation to push supported scalar functions to Databricks Delta.

When you configure SQL ELT optimization for a mapping, you can use scalar functions in a SQL transformation and run queries with the Databricks Delta target endpoint.

You can use a simple SELECT statement without 'FROM' and 'WHERE' arguments. The SQL transformation only supports functions with simple SELECT statement.

The following snippet demonstrates the syntax of a simple SELECT SQL query:

```
SELECT <function_name1>(~Arg~), <function_name2> (~Arg~)...
```

For example, `SELECT SQRT(~AGE~)`

For more information about the supported functions, see the Databricks Delta documentation.

### Rules and guidelines for SQL transformation

Consider the following rules and guidelines when you use SQL transformation:

- You can configure only an SQL query in the SQL transformation. You cannot enable a stored procedure when you push down to Databricks Delta.
- When you enable full SQL ELT optimization, ensure that you use the same connection type for the Source transformation and SQL transformation.
- When you specify a SELECT query, you must also specify the column name and number of columns based on the functions. For example, when you specify the query `select square(~AGE~), sqrt(~SNAME~)`, you must specify two output columns for AGE and SNAME functions each, otherwise the mapping fails.
- If any SQL error occurs, the error is added to the `SQLException` field by default. However, when you run a mapping enabled with SQL ELT optimization, the `SQLException` field remains as Null.
- The `NumRowsAffected` field records the number of rows affected while computing the output buffer. However, for SQL transformation, the `NumRowsAffected` is 0, as the query runs for all the records at the same time.
- You cannot include special characters in the query, as SQL transformation does not support special characters in the arguments.
- You can use an SQL transformation when the SELECT statement is present only in the query property. You cannot configure an SQL transformation with a parameterized query, as dynamic parameter support is limited, and the query fails with a DTM error.

## Features

You can configure SQL ELT optimization for a mapping that reads from the following sources and writes to a Databricks Delta target:

- Databricks Delta source
- Amazon S3 source
- Microsoft Azure Data Lake Storage Gen2 source

When you configure a mapping, some parameters are not supported for a mapping enabled for SQL ELT optimization. You can refer to the list of parameters that each source supports.

## Databricks Delta sources, targets, and lookups

You must configure a Databricks Delta connection with simple or hybrid mode when you enable SQL ELT optimization in a mapping task.

### Source properties

When you configure SQL ELT optimization, the mappings support the following advance properties for a Databricks Delta source:

- Source Object Type
  - Single
  - Multiple
  - Query
  - Parameter

**Note:** When you use the query source type to read from Databricks Delta, you can choose to retain the field metadata and save the mapping. Even if you edit the query and run the mapping, the field metadata specified at design time is retained.

- Query Options
  - Filter. You can use both simple and advanced filter conditions.
- Database Name
- Table Name
- SQL Override

**Note:** Contains, Ends With, and Starts With filter operators are not applicable when you use source filter to filter records.

### Target properties

When you configure SQL ELT optimization, the mappings support the following properties for an Databricks Delta target:

- Target Object Type
  - Single
  - Parameter
  - Create New at Runtime
- Operation
  - Insert
  - Update
  - Upsert
  - Delete
- Create Target
- Target Database Name
- Target Table Name
- Update Mode
- Write Disposition for Insert operation.

**Note:** You cannot run pre-SQL or post-SQL commands in the source and target when you configure mappings for full SQL ELT optimization.

## Lookup properties

When you configure SQL ELT optimization, the mappings support the following advance properties for a Databricks Delta lookup:

- Source Object Type
  - Single
  - Query
  - Parameter
  - Multiple Matches for cached lookup

**Note:** Un-cached unconnected lookup is not supported for SQL ELT optimization. For cached lookup, only **Return all rows** is supported.

- Database Name
- Table Name
- SQL Override

**Note:** If you configure advanced properties that are not supported, the Secure Agent either ignores the properties or logs a SQL ELT optimization validation error in the session logs file. If the Secure Agent logs an error in the session log, the mappings run in the Informatica runtime environment without full SQL ELT optimization.

## Supported features for Amazon S3 V2 source

When you configure SQL ELT optimization, the mappings support the following properties for an Amazon S3 V2 source:

- Source connection parameter
- Source Type - Single, query
- Parameter
- Format - Avro, ORC, Parquet, JSON, and CSV
- Source Type - File and directory. XML source type is not applicable.
- Folder Path
- File Name

When you configure SQL ELT optimization, the mapping supports the following transformations:

- Filter
- Expression
- Aggregator
- Sorter
- Router
- Joiner
- Lookup
- Union
- Rank

For information on how to configure the supported properties, see the Amazon S3 V2 Connector documentation.

## Supported features for Microsoft Azure Data Lake Storage Gen2 source

When you configure SQL ELT optimization, the Microsoft Azure Data Lake Storage Gen2 connection supports the following properties:

- Account Name
- Client ID
- Client Secret
- Tenant ID
- File System Name
- Directory Path
- Adls Gen2 End-point
- Server-side Encryption

When you configure SQL ELT optimization, the mappings support the following properties for a Microsoft Azure Data Lake Storage Gen2 source:

- Source connection, connection parameter
- Source Type - Single, parameter
- Format - Avro, Parquet, JSON, ORC, and CSV.
- Intelligent Structure Model
- Formatting Options
- Filesystem Name Override
- Source Type - File, Directory
- Directory Override - Absolute path; Relative path
- File Name Override - Source object
- Allow Wildcard Characters

When you configure SQL ELT optimization, the mapping supports the following transformations:

- Filter
- Expression
- Aggregator
- Sorter
- Router
- Joiner
- Lookup
- Union
- Rank

For information on how to configure the supported properties, see the Microsoft Azure Data Lake Storage Gen2 Connector documentation.

## Configuring a custom query for the Databricks Delta source object

You can push down a custom query to Databricks Delta.

Before you run a task that contains a custom query as the source object, you must set the **Create Temporary View** session property in the mapping task properties.



**Note:** If you do not set the **Create Temporary View** property, the mapping runs without SQL ELT optimization.

Perform the following task to set the property:

1. In the mapping task, navigate to the **SQL ELT Optimization** section on the **Schedule** tab.
2. Select **Create Temporary View**.
3. Click **Finish**.

## SQL ELT optimization for multiple targets

When you enable full SQL ELT optimization for a mapping to write to multiple Databricks Delta targets, you can further optimize the write operation.

To optimize, you can configure an insert, update, upsert, or delete operation for each target.

**Note:** You cannot configure an insert operation for unconnected target columns for mappings in advanced mode.

You can select the same Databricks Delta target table in multiple Target transformations, configure a different operation for each of the Target transformations independent of each other.

## Single commit for SQL ELT optimization

When you enable full SQL ELT optimization for a mapping to write to multiple Databricks Delta targets, you can configure the mapping to commit the configured operations for all the targets within a connection group together.

You can use single commit to combine the metadata from all the targets and send the metadata for processing in a single execution call. When you use single commit, the Secure Agent separates the targets into connection groups based on equivalent connection attributes and commits the operations together for each connection group. This optimizes the performance of the write operation.

When you run a mapping with multiple targets, the Databricks Delta connections used for these multiple target transformations that have the same connection attribute values are grouped together to form connection groups. As all the targets in a connection group have the same connection attributes, only a single connection is established for each connection group which represents that particular connection group. The transactions on each connection group runs on a single Databricks cluster.

If the Secure Agent fails to write to any of the targets, the task execution stops and the completed transactions for the targets that belong to the same connection group are not rolled back.

To enable single commit to write to multiple targets, set the **EnableSingleCommit=Yes** custom property in the **Advanced Session Properties** section on the **Schedule** tab of the mapping task.

When you run a mapping with single commit enabled, you can view the row statistics details in the session logs.

Single commit is applicable only when you run a mapping on Databricks cluster.

## Rules and guidelines for SQL ELT optimization

Use the following rules and guidelines when you enable a mapping for SQL ELT optimization to a Databricks Delta database:

### Mapping with Databricks Delta source and target

Use the following rules and guidelines when you configure SQL ELT optimization in a mapping that reads from and writes to Databricks Delta:

- When you provide the database name in the connection, the following rules and guidelines apply:
  - For full SQL ELT optimization, temporary staging tables or views are created in the database name provided in the target connection attribute `database`. If the attribute is empty the database name is picked from JDBC URL of the target connection.
  - For Source SQL ELT optimization, temporary staging tables or views are created in the database name provided in the source connection attribute `database`. If the attribute is empty the database name is picked from JDBC URL of source connection.
  - The database name provided in the connection should have read and write permission. If the `database name` field is empty, then the database name provided in the JDBC URL will take precedence and should have read and write permission.

**Note:** The guidelines are only applicable to custom query support, SQL override, and SQL ELT optimization in Databricks using flat file formats.

- LAST function is a non-deterministic function. This function returns different results each time it is called, even when you provide the same input values.
- When you configure a Filter transformation or specify a filter condition, do not specify special characters.
- When you connect to Databricks clusters to process the mapping and define a custom query with multiple tables in the SELECT statement, the mapping displays incorrect data for fields that have the same name. This doesn't apply to Databricks runtime version 9.1 LTA or later.
- When you configure a mapping enabled for full SQL ELT optimization to read from multiple sources and you override the database name and table name from the advanced properties, the mapping fails.
- To configure a Filter transformation or specify a filter condition on columns of date or timestamp in a Databricks Delta table, you must pass the data through the `TO_DATE()` function as an expression in the filter condition.
- When you specify custom query as a source object, ensure that the SQL query does not contain any partitioning hints such as `COALESCE`, `REPARTITION`, or `REPARTITION_BY_RANGE`.
- When you configure a mapping enabled for full SQL ELT optimization on the Databricks Delta SQL engine, you cannot configure single commit to write to multiple targets.
- When you configure a mapping enabled for full SQL ELT optimization on the Databricks Delta SQL engine and push the data to the Databricks Delta target, ensure that you map all the fields in target. Else, the mapping fails.
- When you create a new target at runtime, you must not specify a database name and table name in the **Target Database Name** and **Target Table Name** in the target advanced properties.
- When you read data from a column of Date data type and write data into a column of Date data type, the SQL ELT query pushes the column of Date data type and casts the column to Timestamp data type.
- You cannot completely parameterize a multi-line custom query using a parameter file. If you specify a multi-line custom query in a parameter file, the mapping considers only the first line of the multi-line query.

- When you push the CRC32() function to Databricks Delta, the data type of the return value is either Bigint or String.
- When you push the DATE\_DIFF() function to Databricks Delta, the function returns the integral part of the value and not the fractional part.
- When you push the GREATEST() function to Databricks Delta and configure input value arguments of String data type, you must not specify the caseFlag argument.
- To push the TO\_CHAR(STRING) function to Databricks Delta, use the following string and format arguments:
  - YYYY
  - YY
  - MM
  - MON
  - MONTH
  - DD
  - DDD
  - DY
  - DAY
  - HH12
  - HH24
  - MI
  - Q
  - SS
  - SS.MS
  - SS.US
  - SS.NS
- To push the TO\_DATE(string, format) function to Databricks Delta, you must use the following format arguments:
  - YYYY
  - YY
  - MM
  - MON
  - MONTH
  - DD
  - DDD
  - HH12
  - HH24
  - MI
  - SS
  - SS.MS
  - SS.US
  - SS.NS

- When you enable full SQL ELT optimization in a mapping and use the IFF() condition in an Expression transformation, the mapping fails for the following functions:
  - IS\_SPACES
  - IS\_NUMBER
  - IS\_DATE
- A mapping enabled with full SQL ELT optimization and contains an SQL transformation fails when the column names in the SQL override query don't match with the column names in the custom query.

### Mapping with Amazon S3 source and Databricks Delta target

Use the following rules and guidelines when you configure SQL ELT optimization in a mapping that reads from an Amazon S3 source and writes to a Databricks Delta target:

- When you select the source type as directory in the advanced source properties, ensure that all the files in the directory contain the same schema.
- When you select query as the source type in lookup, you cannot override the database name and table name in the advanced source properties.
- When you include a source transformation in a mapping enabled with SQL ELT optimization, exclude the FileName field from the source. The FileName field is not applicable.
- When you parameterize a lookup object in a mapping enabled with SQL ELT optimization, the mapping fails as you cannot exclude the filename part at runtime.
- When you parameterize the source object in a mapping task, ensure that you pass the source object parameter value with the fully qualified path in the parameter file.
- You cannot use wildcard characters for the source file name and directory name in the source transformation.
- You cannot use wildcard characters for the folder path or file name in the advanced source properties.
- When you read from a partition folder that has a transaction log file, select the source type as Directory in the advanced source properties.
- You cannot configure dynamic lookup cache.
- When you use a Joiner transformation in a mapping enabled with SQL ELT optimization and create a new target at runtime, ensure that the fields do not have a not null constraint.
- Ensure that the field names in Parquet, ORC, AVRO, or JSON files do not contain Unicode characters.

### Mapping with Azure Data Lake Storage Gen2 source and Databricks Delta target

Use the following rules and guidelines when you configure SQL ELT optimization in a mapping that reads from a Azure Data Lake Storage Gen2 source and writes to a Databricks Delta target:

- Mappings fail if the lookup object contains unsupported data types.
- When you select the source type as directory in the advanced source property, ensure that all the files in the directory contain the same schema.
- When you select query as the source type in lookup, you cannot override the database name and table name in the advanced source properties.
- When you include a source transformation in a mapping enabled with SQL ELT optimization, exclude the FileName field from the source. The FileName field is not applicable.
- When you parameterize a lookup object in a mapping enabled with SQL ELT optimization, the mapping fails as you cannot exclude the filename part at runtime.
- When you parameterize the source object in a mapping task, ensure that you pass the source object parameter value with the fully qualified path in the parameter file.

- You cannot use wildcard characters for the source file name and directory name in the source transformation.
- When you read from a partition folder that has a transaction log file, select the source type as Directory in the advanced source properties.
- You cannot configure dynamic lookup cache.
- When you use a Joiner transformation in a mapping enabled with SQL ELT optimization and create a new target at runtime, ensure that the fields do not have a not null constraint.
- Ensure that the field names in Parquet, ORC, AVRO, or JSON files do not contain Unicode characters.

### Cross workspace mappings

When you set up a mapping enabled with full SQL ELT optimization to access data from a Databricks Delta workspace, and the associated metastore resides in a separate workspace, the mapping runs without SQL ELT optimization.

## Troubleshooting SQL ELT optimization

### Mapping fails when configured to read date or timestamp information and write to default date/time format

When you configure a mapping to read date or timestamp information from a string column and process the data with the default date/time format to write to Databricks Delta target, the mapping fails with the following error:

```
[ERROR] The Secure Agent failed to run the full pushdown query due to the following error: [Invalid timestamp: '12/31/1972 00:00:00.000001']
```

To resolve this issue, set the JVM option `-DHonorInfaDateFormat=true` for the Secure Agent.

Perform the following steps to configure the JVM option in Administrator:

1. Select **Administrator > Runtime Environments**.
2. On the Runtime Environments page, select the Secure Agent machine that runs the mapping.
3. Click **Edit**.
4. In the System Configuration Details section, select **Data Integration Server** as the Service and **DTM** as the Type.
5. Edit the JVMOption system property and set the value to `-DHonorInfaDateFormat=true`.
6. Click **Save**.

### IS\_DATE(), IS\_SPACES(), and IS\_NUMBER() functions return 0 or 1 instead of True or False.

When you use IS\_DATE(), IS\_SPACES(), and IS\_NUMBER() functions, the functions return 0 or 1 instead of True or False.

To resolve this issue, set the JVM option `-DDeltaSQLLELTBooleanReturnAsString=true` for the Secure Agent.

Perform the following steps to configure the JVM option in Administrator:

1. Select **Administrator > Runtime Environments**.
2. On the Runtime Environments page, select the Secure Agent machine that runs the mapping.
3. Click **Edit**.
4. In the System Configuration Details section, select **Data Integration Server** as the Service and **DTM** as the Type.

5. Edit the JVMOption system property and set the value to -  
`DDeltaSQLLELTBooleanReturnAsString=true`.
6. Click **Save**.

# CHAPTER 5

## Data type reference

### Databricks Delta native data types

Databricks Delta data types appear in the Fields tab for Source and Target transformations when you choose to edit metadata for the fields.

### Transformation data types

Set of data types that appear in the remaining transformations. They are internal data types based on ANSI SQL-92 generic data types, which Data Integration uses to move data across platforms. Transformation data types appear in all remaining transformations in Data Integration tasks.

When the Data Integration application reads source data, it converts the native data types to the comparable transformation data types before transforming the data. When the Data Integration application writes to a target, it converts the transformation data types to the comparable native data types.

## Databricks Delta and transformation data types

The following table compares the Databricks Delta native data type to the transformation data type:

Databricks Delta Data Type	Transformation Data Type	Range and Description
Array <sup>1</sup>	Array	Unlimited number of characters.
Binary	Binary	1 to 104,857,600 bytes.
Bigint	Bigint	-9,223,372,036,854,775,808 to +9,223,372,036,854,775,807. 8-byte signed integer.
Boolean	Integer	1 or 0.
Date	Date/Time	Date and time values.
Decimal	Decimal	Exact numeric of selectable precision For mappings: Max precision 28, scale 27. For mappings in advanced mode: Max precision 38, scale 37.
Double	Double	Precision 15.

Databricks Delta Data Type	Transformation Data Type	Range and Description
Float	Double	Precision 7.
Int	Integer	-2,147,483,648 to +2,147,483,647.
Map <sup>1</sup>	Map	Unlimited number of characters.
Smallint	Integer	-32,768 to +32,767.
String	String	1 to 104,857,600 characters.
Struct <sup>1</sup>	Struct	Unlimited number of characters.
Tinyint	Integer	-128 to 127
Timestamp	Date/Time	January 1,0001 00:00:00 to December 31,9999 23:59:59.997443. Timestamp values only preserve results up to microsecond precision of six digits. The precision beyond six digits is discarded.
<sup>1</sup> Applies only to mappings in advanced mode. Doesn't apply to Databricks cluster and SQL ELT optimization.		

## Rules and guidelines for data types

Databricks Delta has the following restriction for decimal data type:

The behavior of the decimal data type differs in mappings and mappings in advanced mode. Mappings support decimal max precision 28, while mappings in advanced mode supports max decimal precision of 38.

In mappings, if the decimal data type exceeds 28 precision in the source, the numeric value for the decimal is rounded off after the 18th precision and the remaining digits are replaced with zeroes in the target.

For example, the value 1234567890123456789012345678.9012345678 from the source is rounded off to 1234567890123456900000000000 in the target. However, in mappings in advanced mode, the decimal data from the source remains the same in the target.

To resolve the issue in mappings, specify a precision that is less than or equal to 28 for the Decimal data type in the source table.



# INDEX

## C

Cloud Application Integration community  
URL [5](#)  
Cloud Developer community  
URL [5](#)  
Create target  
rules and guidelines [35](#)  
create target at runtime [33](#)  
custom query [29](#)

## D

Data Integration community  
URL [5](#)  
data types [63](#)  
Databricks Delta  
SQL ELT optimization [47](#)  
SQL ELT optimization overview [46](#)  
Databricks Delta connection  
SQL ELT optimization [48](#)  
Databricks Delta connections  
overview [9](#)  
Databricks Delta connector  
datetime format [44](#)  
rules and guidelines [44](#)  
Databricks Delta Connector  
assets [7](#)  
overview [7](#)  
Databricks Delta connector rules and guidelines [42](#)

## F

field delimiter [25](#)

## I

Informatica Global Customer Support  
contact information [6](#)  
Informatica Intelligent Cloud Services  
web site [5](#)

## M

maintenance outages [6](#)  
mappings  
overview [23](#)

mappings (*continued*)  
source properties [25](#)  
target properties [31](#)  
mappings in advanced mode  
example [41](#)

## N

native data type [63](#)

## P

properties  
in mappings [25, 31](#)

## S

source properties [25](#)  
SQL ELT optimization  
functions [49, 51](#)  
preview [47](#)  
transformations [49, 51](#)  
SQL ELT Optimization  
rules and guidelines [58](#)  
SQL ELT optimization preview [47](#)  
status  
Informatica Intelligent Cloud Services [6](#)  
system status [6](#)

## T

tracing level [25](#)  
transformation data type [63](#)  
transformations  
SQL ELT optimization [51](#)  
trust site  
description [6](#)

## U

upgrade notifications [6](#)

## W

web site [5](#)