



Informatica® Data Engineering Integration  
10.5.2

# Integration Guide

© Copyright Informatica LLC 2014, 2022

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 **[and any other Informatica-owned trademarks appearing in the document]** 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.

Subject to your opt-out rights, the software will automatically transmit to Informatica in the USA information about the computing and network environment in which the Software is deployed and the data usage and system statistics of the deployment. This transmission is deemed part of the Services under the Informatica privacy policy and Informatica will use and otherwise process this information in accordance with the Informatica privacy policy available at <https://www.informatica.com/in/privacy-policy.html>. You may disable usage collection in Administrator tool.

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

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: 2022-09-29

# Table of Contents

<b>Preface .....</b>	<b>11</b>
Informatica Resources. ....	11
Informatica Network. ....	11
Informatica Knowledge Base. ....	11
Informatica Documentation. ....	12
Informatica Product Availability Matrices. ....	12
Informatica Velocity. ....	12
Informatica Marketplace. ....	12
Informatica Global Customer Support. ....	12
 <b>Part I: Hadoop Integration.....</b>	<b>13</b>
 <b>Chapter 1: Introduction to Hadoop Integration.....</b>	<b>14</b>
Cluster Integration Overview. ....	14
Data Engineering Integration Component Architecture. ....	15
Hadoop Integration. ....	15
Clients and Tools. ....	15
Application Services. ....	15
Repositories. ....	16
Integration with Other Informatica Products. ....	16
 <b>Chapter 2: Before You Begin.....</b>	<b>18</b>
Read the Release Notes. ....	18
Verify System Requirements. ....	18
Verify Product Installations. ....	18
Verify HDFS Disk Space. ....	19
Verify the Hadoop Distribution. ....	19
Verify Port Requirements. ....	19
Uninstall Big Data Management. ....	21
Uninstall for Amazon EMR, Azure HDInsight, and MapR. ....	21
Uninstall for Cloudera CDH. ....	22
Uninstall for Hortonworks HDP. ....	22
Prepare Directories, Users, and Permissions. ....	23
Verify and Create Users. ....	23
Grant Access to Azure ADLS Resources for Informatica Users. ....	25
Create Directories and Set Permissions. ....	27
Configure Access to Secure Hadoop Clusters. ....	30
Configuring Access to an SSL/TLS-Enabled Cluster. ....	31
Generate the OAUTH Token. ....	34
Generating Keytab Files for the SPN User. ....	34

Configure Apache Ranger with HDInsight. . . . .	36
Configure the Metadata Access Service. . . . .	36
Configure the Data Integration Service. . . . .	37
Download the Informatica Server Binaries for the Hadoop Environment. . . . .	37
Edit the etc/hosts File. . . . .	38
Configuring LZ0 Compression Format. . . . .	39
Configuring the Data Integration Service to Use Operating System Profiles . . . . .	40
Configure Data Integration Service Properties. . . . .	40
Prepare a Python Installation. . . . .	41
Install Python for Enterprise Data Preparation. . . . .	42
<b>Chapter 3: Amazon EMR Integration Tasks. . . . .</b>	<b>44</b>
Amazon EMR Task Flows. . . . .	44
Task Flow to Integrate with Amazon EMR. . . . .	45
Task Flow to Upgrade from Version 10.2.1 . . . . .	46
Task Flow to Upgrade from Version 10.2. . . . .	47
Task Flow to Upgrade from a Version Earlier than 10.2. . . . .	48
Prepare for Cluster Import from Amazon EMR. . . . .	49
Configure *-site.xml Files for Amazon EMR. . . . .	49
Prepare the Archive File for Amazon EMR. . . . .	55
Configure Glue as the Hive Metastore . . . . .	55
Create a Cluster Configuration. . . . .	56
Importing a Hadoop Cluster Configuration from a File. . . . .	56
Verify or Refresh the Cluster Configuration . . . . .	57
Verify JDBC Drivers for Sqoop Connectivity. . . . .	58
Verify Design-time Drivers. . . . .	58
Verify Run-time Drivers. . . . .	58
Configure the Files to Use S3 . . . . .	59
Set S3 Access Policies. . . . .	60
Step 1. Identify the S3 Access Policy Elements. . . . .	61
Step 2. Optionally Copy an Existing S3 Access Policy as a Template. . . . .	61
Step 3. Create or Edit an S3 Access Policy. . . . .	62
Configure the Developer Tool. . . . .	62
Configure developerCore.ini. . . . .	63
Configure the Developer Tool for Kerberos. . . . .	63
Complete Upgrade Tasks. . . . .	64
Update Connections . . . . .	64
Update Streaming Objects. . . . .	66
<b>Chapter 4: Azure HDInsight Integration Tasks. . . . .</b>	<b>69</b>
Azure HDInsight Task Flows. . . . .	69
Task Flow to Integrate with Azure HDInsight. . . . .	70
Task Flow to Upgrade from Version 10.2.1 . . . . .	71

Task Flow to Upgrade from Version 10.2. . . . .	72
Task Flow to Upgrade from a Version Earlier than 10.2. . . . .	73
Prepare for Cluster Import from Azure HDInsight. . . . .	74
Configure *-site.xml Files for Azure HDInsight. . . . .	74
Verify HDInsight Cluster Security Settings. . . . .	80
Prepare for Direct Import from Azure HDInsight. . . . .	81
Prepare the Archive File for Import from Azure HDInsight. . . . .	81
Create a Cluster Configuration. . . . .	82
Before You Import. . . . .	82
Importing a Hadoop Cluster Configuration from the Cluster. . . . .	82
Importing a Hadoop Cluster Configuration from a File. . . . .	83
Verify or Refresh the Cluster Configuration . . . . .	84
Configure the Hive Warehouse Connector and Hive LLAP. . . . .	85
Verify JDBC Drivers for Sqoop Connectivity. . . . .	86
Verify Design-time Drivers. . . . .	86
Verify Run-time Drivers. . . . .	87
Configure the Developer Tool. . . . .	87
Configure developerCore.ini. . . . .	87
Configure the Developer Tool for Kerberos. . . . .	88
Complete Upgrade Tasks. . . . .	88
Update Connections . . . . .	88
Update Streaming Objects. . . . .	91
<b>Chapter 5: Cloudera CDH Integration Tasks. . . . .</b>	<b>93</b>
Cloudera CDH Task Flows. . . . .	93
Task Flow to Integrate with Cloudera CDH. . . . .	94
Task Flow to Upgrade from Version 10.2.1 . . . . .	95
Task Flow to Upgrade from Version 10.2. . . . .	96
Task Flow to Upgrade from a Version Earlier than 10.2. . . . .	97
Prepare for Cluster Import from Cloudera CDH. . . . .	98
Configure *-site.xml Files for Cloudera CDH. . . . .	98
Prepare for Direct Import from Cloudera CDH. . . . .	102
Prepare the Archive File for Import from Cloudera CDH. . . . .	102
Create a Cluster Configuration. . . . .	103
Before You Import. . . . .	103
Importing a Hadoop Cluster Configuration from the Cluster. . . . .	103
Importing a Hadoop Cluster Configuration from a File. . . . .	104
Verify or Refresh the Cluster Configuration . . . . .	105
Verify JDBC Drivers for Sqoop Connectivity. . . . .	106
Verify Design-time Drivers. . . . .	106
Verify Run-time Drivers. . . . .	107
Set the Locale for Cloudera CDH 6.x. . . . .	107
Enable Data Preparation of JSON Files on Cloudera CDH. . . . .	108

Complete Upgrade Tasks. . . . .	108
Update Connections . . . . .	108
Update Streaming Objects. . . . .	111
<b>Chapter 6: Cloudera CDP Integration Tasks. . . . .</b>	<b>113</b>
Cloudera CDP Task Flows. . . . .	113
Task Flow to Integrate with Cloudera CDP. . . . .	114
Task Flow to Upgrade from Version 10.4.1. . . . .	115
Prepare for Cluster Import from Cloudera CDP. . . . .	116
Configure *-site.xml Files for Cloudera CDP. . . . .	116
Prepare for Direct Import from Cloudera CDP. . . . .	120
Prepare the Archive File for Import from Cloudera CDP. . . . .	121
Create a Cluster Configuration. . . . .	121
Before You Import. . . . .	121
Importing a Hadoop Cluster Configuration from the Cluster. . . . .	122
Importing a Hadoop Cluster Configuration from a File. . . . .	123
Copy the Truststore File to the Informatica Domain. . . . .	124
Configure the Impersonation User for Operating System Profiles. . . . .	125
Verify JDBC Drivers for Sqoop Connectivity. . . . .	125
Verify Design-time Drivers. . . . .	125
Verify Run-time Drivers. . . . .	125
Set the Locale for Cloudera CDP. . . . .	126
Enable Data Preparation of JSON Files on Cloudera CDP. . . . .	127
Configure the Developer Tool. . . . .	127
Configure developerCore.ini. . . . .	127
Configure the Developer Tool for Kerberos. . . . .	128
<b>Chapter 7: Google Dataproc Integration Tasks. . . . .</b>	<b>129</b>
Google Dataproc Task Flows. . . . .	129
Task Flow to Integrate Versions 10.4 or 10.5. . . . .	130
Task Flow to Upgrade from Version 10.2.2. . . . .	131
Prepare for Cluster Import from Google Dataproc. . . . .	132
Configure *-site.xml Files for Google Dataproc. . . . .	132
Prepare the Archive File for Import from Google Dataproc. . . . .	137
Verify the Distribution Version . . . . .	138
Create a Cluster Configuration . . . . .	138
Importing a Hadoop Cluster Configuration from a File. . . . .	138
Configure the Cluster for the Blaze Engine. . . . .	139
Configure Domain Settings. . . . .	140
Copy the Kerberos Configuration File . . . . .	140
Enable Access to Google Cloud Sources . . . . .	140
Verify JDBC Drivers for Sqoop Connectivity . . . . .	140
Complete Upgrade Tasks. . . . .	141

Complete Connection Upgrade. . . . .	142
Configure the Developer Tool . . . . .	142
Edit the etc/hosts File . . . . .	142
Configure developerCore.ini. . . . .	143
Configure the Developer Tool for Kerberos. . . . .	143
<b>Chapter 8: Hortonworks HDP Integration Tasks. . . . .</b>	<b>144</b>
Hortonworks HDP Task Flows. . . . .	144
Task Flow to Integrate with Hortonworks HDP. . . . .	145
Task Flow to Upgrade from Version 10.2.1 . . . . .	146
Task Flow to Upgrade from Version 10.2. . . . .	147
Task Flow to Upgrade from a Version Earlier than 10.2. . . . .	148
Prepare for Cluster Import from Hortonworks HDP. . . . .	149
Configure *-site.xml Files for Hortonworks HDP. . . . .	149
Prepare for Direct Import from Hortonworks HDP. . . . .	154
Prepare the Archive File for Import from Hortonworks HDP. . . . .	154
Create a Cluster Configuration. . . . .	155
Before You Import. . . . .	155
Importing a Hadoop Cluster Configuration from the Cluster. . . . .	155
Importing a Hadoop Cluster Configuration from a File. . . . .	156
Verify or Refresh the Cluster Configuration . . . . .	157
Configure the Hive Warehouse Connector and Hive LLAP. . . . .	158
Verify JDBC Drivers for Sqoop Connectivity. . . . .	159
Verify Design-time Drivers. . . . .	159
Verify Run-time Drivers. . . . .	160
Configure the Developer Tool. . . . .	160
Configure developerCore.ini. . . . .	160
Configure the Developer Tool for Kerberos. . . . .	161
Complete Upgrade Tasks. . . . .	161
Update Connections . . . . .	161
Update Streaming Objects. . . . .	164
<b>Chapter 9: MapR Integration Tasks. . . . .</b>	<b>166</b>
MapR Task Flows. . . . .	166
Task Flow to Integrate with MapR (copy). . . . .	167
Task Flow to Upgrade from Version 10.2.2 (mapr) (copy). . . . .	168
Task Flow to Upgrade from Version 10.2 (mapr) (copy). . . . .	169
Task Flow to Upgrade from a Version Earlier than 10.2 (mapr) (copy). . . . .	170
Install and Configure the MapR Client . . . . .	171
Prepare for Cluster Import from MapR. . . . .	171
Configure *-site.xml Files for MapR. . . . .	172
Prepare the Archive File for Import from MapR. . . . .	176
Create a Cluster Configuration. . . . .	177

Importing a Hadoop Cluster Configuration from a File. . . . .	177
Verify or Refresh the Cluster Configuration . . . . .	178
Verify JDBC Drivers for Sqoop Connectivity. . . . .	179
Verify Design-time Drivers. . . . .	179
Verify Run-time Drivers. . . . .	179
Generate MapR Tickets. . . . .	180
Generate Tickets. . . . .	180
Configure the Data Integration Service. . . . .	181
Configure the Metadata Access Service. . . . .	182
Configure the Analyst Service. . . . .	182
Complete Upgrade Tasks. . . . .	183
Update Connections . . . . .	183
<b>Part II: Databricks Integration. . . . .</b>	<b>186</b>
<b>Chapter 10: Introduction to Databricks Integration. . . . .</b>	<b>187</b>
Databricks Integration Overview. . . . .	187
Run-time Process on the Databricks Spark Engine. . . . .	188
Native Environment. . . . .	188
Databricks Environment. . . . .	189
Databricks Integration Task Flow. . . . .	189
<b>Chapter 11: Before You Begin Databricks Integration. . . . .</b>	<b>191</b>
Read the Release Notes. . . . .	191
Verify System Requirements. . . . .	191
Configure Preemption for Concurrent Jobs. . . . .	192
Configure Storage Access. . . . .	193
Configure S3 and Redshift Authentication and Encryption on AWS. . . . .	193
Configure AWS Roles and Policies to Access S3 Resources. . . . .	194
Download and Install the JDBC Driver to Enable Delta Lake Access. . . . .	199
Configure ADLS Storage Access. . . . .	200
Configure WASB Storage Access. . . . .	200
Create a Staging Directory for Binary Archive Files. . . . .	200
Create a Staging Directory for Run-time Processing. . . . .	201
Prepare for Token Authentication. . . . .	201
Configure the Data Integration Service. . . . .	201
Configure Data Integration Service Properties. . . . .	201
Install Python Libraries. . . . .	202
<b>Chapter 12: Databricks Integration Tasks. . . . .</b>	<b>203</b>
Create a Databricks Cluster Configuration. . . . .	203
Importing a Databricks Cluster Configuration from the Cluster. . . . .	203
Importing a Databricks Cluster Configuration from a File. . . . .	204



Configure the Databricks Connection. . . . .	206
Complete Upgrade Tasks. . . . .	206
<b>Appendix A: Connections Reference.....</b>	<b>207</b>
Connections Overview. . . . .	208
Cloud Provisioning Configuration. . . . .	208
AWS Cloud Provisioning Configuration Properties. . . . .	209
Azure Cloud Provisioning Configuration Properties. . . . .	212
Databricks Cloud Provisioning Configuration Properties. . . . .	217
Amazon Redshift Connection Properties. . . . .	218
Amazon S3 Connection Properties. . . . .	220
Blockchain Connection Properties. . . . .	222
Cassandra Connection Properties. . . . .	224
Confluent Kafka Connection. . . . .	225
General Properties. . . . .	226
Confluent Kafka Broker Properties. . . . .	226
SSL Properties . . . . .	227
Creating a Confluent Kafka Connection Using infacmd. . . . .	227
Databricks Connection Properties. . . . .	227
Google Analytics Connection Properties. . . . .	229
Google BigQuery Connection Properties. . . . .	230
Google Cloud Spanner Connection Properties. . . . .	232
Google Cloud Storage Connection Properties. . . . .	233
Google PubSub Connection Properties. . . . .	234
Hadoop Connection Properties. . . . .	234
Hadoop Cluster Properties. . . . .	235
Common Properties. . . . .	237
Reject Directory Properties. . . . .	238
Blaze Configuration. . . . .	239
Spark Configuration. . . . .	240
HDFS Connection Properties. . . . .	241
HBase Connection Properties. . . . .	243
HBase Connection Properties for MapR-DB. . . . .	243
Hive Connection Properties. . . . .	244
JDBC Connection Properties. . . . .	247
JDBC Connection String. . . . .	249
Sqoop Connection-Level Arguments. . . . .	251
Delta Lake JDBC Connection Properties. . . . .	254
JDBC V2 Connection Properties. . . . .	254
Kafka Connection Properties. . . . .	256
General Properties. . . . .	257
Kafka Broker Properties. . . . .	258
SSL Properties. . . . .	259

Creating a Kafka Connection Using infacmd. . . . .	259
Kudu Connection Properties . . . . .	259
Microsoft Azure Blob Storage Connection Properties. . . . .	260
Microsoft Azure Cosmos DB SQL API Connection Properties. . . . .	261
Microsoft Azure Data Lake Storage Gen1 Connection Properties. . . . .	262
Microsoft Azure Data Lake Storage Gen2 Connection Properties. . . . .	263
Microsoft Azure SQL Data Warehouse Connection Properties. . . . .	264
Snowflake Connection Properties. . . . .	266
Creating a Connection to Access Sources or Targets. . . . .	267
Creating a Hadoop Connection. . . . .	267
Configuring Hadoop Connection Properties. . . . .	269
Cluster Environment Variables. . . . .	269
Cluster Library Path. . . . .	269
Common Advanced Properties. . . . .	270
Blaze Engine Advanced Properties. . . . .	270
Spark Advanced Properties. . . . .	271
<b>Index. . . . .</b>	<b>277</b>

# Preface

Follow the instructions in the *Informatica Data Engineering Integration Guide* to integrate Informatica with non-native environments.

Integration tasks are required on the Hadoop cluster, the Data Integration Service machine, and the Developer tool machine. As a result, this guide contains tasks for administrators of the non-native environments and Informatica administrators. Tasks required by the Hadoop or Databricks administrator are directed to the administrator.

Use this guide for new integrations and for upgrades. The instructions follow the same task flow. Tasks required for upgrade indicate that they are for upgrade.

## 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 Network

The Informatica Network is the gateway to many resources, including the Informatica Knowledge Base and Informatica Global Customer Support. To enter the Informatica Network, visit <https://network.informatica.com>.

As an Informatica Network member, you have the following options:

- Search the Knowledge Base for product resources.
- View product availability information.
- Create and review your support cases.
- Find your local Informatica User Group Network and collaborate with your peers.

### 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 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 Product Availability Matrices

Product Availability Matrices (PAMs) indicate the versions of the operating systems, databases, and types of data sources and targets that a product release supports. You can browse the Informatica PAMs at <https://network.informatica.com/community/informatica-network/product-availability-matrices>.

## Informatica Velocity

Informatica Velocity is a collection of tips and best practices developed by Informatica Professional Services and based on real-world experiences from hundreds of data management projects. Informatica Velocity represents the collective knowledge of Informatica consultants who work with organizations around the world to plan, develop, deploy, and maintain successful data management solutions.

You can find Informatica Velocity resources at <http://velocity.informatica.com>. If you have questions, comments, or ideas about Informatica Velocity, contact Informatica Professional Services at [ips@informatica.com](mailto:ips@informatica.com).

## Informatica Marketplace

The Informatica Marketplace is a forum where you can find solutions that extend and enhance your Informatica implementations. Leverage any of the hundreds of solutions from Informatica developers and partners on the Marketplace to improve your productivity and speed up time to implementation on your projects. You can find the Informatica Marketplace at <https://marketplace.informatica.com>.

## Informatica Global Customer Support

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

To find your local Informatica Global Customer Support telephone number, visit the Informatica website at the following link:

<https://www.informatica.com/services-and-training/customer-success-services/contact-us.html>.

To find online support resources on the Informatica Network, visit <https://network.informatica.com> and select the eSupport option.

# Part I: Hadoop Integration

This part contains the following chapters:

- [Introduction to Hadoop Integration, 14](#)
- [Before You Begin, 18](#)
- [Amazon EMR Integration Tasks, 44](#)
- [Azure HDInsight Integration Tasks, 69](#)
- [Cloudera CDH Integration Tasks, 93](#)
- [Cloudera CDP Integration Tasks, 113](#)
- [Google Dataproc Integration Tasks, 129](#)
- [Hortonworks HDP Integration Tasks, 144](#)
- [MapR Integration Tasks, 166](#)

# CHAPTER 1

## Introduction to Hadoop Integration

This chapter includes the following topics:

- [Cluster Integration Overview, 14](#)
- [Data Engineering Integration Component Architecture, 15](#)
- [Integration with Other Informatica Products, 16](#)

### Cluster Integration Overview

You can integrate the Informatica domain with Hadoop clusters through Data Engineering Integration.

The Data Integration Service automatically installs Hadoop binaries to integrate the Informatica domain with the Hadoop environment.

The integration requires Informatica connection objects and cluster configurations. A cluster configuration is a domain object that contains configuration parameters that you import from the Hadoop cluster. You then associate the cluster configuration with connections to access the Hadoop environment.

Perform the following tasks to integrate the Informatica domain with the Hadoop environment:

1. Install or upgrade to the current Informatica version.
2. Perform pre-import tasks, such as verifying system requirements and user permissions.
3. Import the cluster configuration into the domain. The cluster configuration contains properties from the \*-site.xml files on the cluster.
4. Create a Hadoop connection and other connections to run mappings within the Hadoop environment.
5. Perform post-import tasks specific to the Hadoop distribution that you integrate with.

When you run a mapping, the Data Integration Service checks for the binary files on the cluster. If they do not exist or if they are not synchronized, the Data Integration Service prepares the files for transfer. It transfers the files to the distributed cache through the Informatica Hadoop staging directory on HDFS. By default, the staging directory is /tmp. This transfer process replaces the requirement to install distribution packages on the Hadoop cluster.

# Data Engineering Integration Component Architecture

The Data Engineering Integration components include client tools, application services, repositories, and third-party tools that Data Engineering Integration uses for a data engineering project. The specific components involved depend on the task you perform.

## Hadoop Integration

The Informatica domain can connect to clusters that run different Hadoop distributions. Hadoop is an open-source software framework that enables distributed processing of large data sets across clusters of machines. You might also need to use third-party software clients to set up and manage your Hadoop cluster.

The domain can connect to the supported data source in the Hadoop environment, such as HDFS, HBase, or Hive, and push job processing to the Hadoop cluster. To enable high performance access to files across the cluster, you can connect to an HDFS source. You can also connect to a Hive source, which is a data warehouse that connects to HDFS.

It can also connect to NoSQL databases such as HBase, which is a database comprising key-value pairs on Hadoop that performs operations in real-time. The Data Integration Service can push mapping jobs to the Spark or Blaze engine, and it can push profile jobs to the Blaze engine in the Hadoop environment.

Data Engineering Integration supports more than one version of some Hadoop distributions. By default, the cluster configuration wizard populates the latest supported version.

## Clients and Tools

You can use several Informatica tools and clients to manage data engineering projects.

Based on your product license and use case, you can use the Administrator tool, Developer tool, Analyst tool, and command line interface.

## Application Services

Informatica uses application services in the Informatica domain to process data.

Use the Administrator tool to create connections, monitor jobs, and manage application services.

Data Engineering Integration uses the following application services:

### **Analyst Service**

The Analyst Service runs the Analyst tool in the Informatica domain. The Analyst Service manages the connections between service components and the users that have access to the Analyst tool.

### **Data Integration Service**

The Data Integration Service can process mappings in the native environment or push the mapping for processing to a compute cluster in a non-native environment. The Data Integration Service also retrieves metadata from the Model repository when you run a Developer tool mapping or workflow. The Analyst tool and Developer tool connect to the Data Integration Service to run profile jobs and store profile results in the profiling warehouse.

### **Mass Ingestion Service**

The Mass Ingestion Service manages and validates mass ingestion specifications that you create in the Mass Ingestion tool. The Mass Ingestion Service deploys specifications to the Data Integration Service. When a specification runs, the Mass Ingestion Service generates ingestion statistics.

### **Metadata Access Service**

The Metadata Access Service allows the Developer tool to import and preview metadata from a Hadoop cluster.

The Metadata Access Service contains information about the Service Principal Name (SPN) and keytab information if the Hadoop cluster uses Kerberos authentication. You can create one or more Metadata Access Services on a node. Based on your license, the Metadata Access Service can be highly available.

HBase, HDFS, Hive, and MapR-DB connections use the Metadata Access Service when you import an object from a Hadoop cluster. Create and configure a Metadata Access Service before you create HBase, HDFS, Hive, and MapR-DB connections.

### **Model Repository Service**

The Model Repository Service manages the Model repository. The Model Repository Service connects to the Model repository when you run a mapping, mapping specification, profile, or workflow.

### **REST Operations Hub**

The REST Operations Hub Service is an application service in the Informatica domain that exposes Informatica product functionality to external clients through REST APIs.

## **Repositories**

Data Engineering Integration uses repositories and other databases to store data related to connections, source metadata, data domains, data profiling, data masking, and data lineage. Data Engineering Integration uses application services in the Informatica domain to access data in repositories.

Data Engineering Integration uses the following databases:

### **Model repository**

The Model repository stores profiles, data domains, mapping, and workflows that you manage in the Developer tool. The Model repository also stores profiles, data domains, and mapping specifications that you manage in the Analyst tool.

### **Profiling warehouse**

The Data Integration Service runs profiles and stores profile results in the profiling warehouse.

## **Integration with Other Informatica Products**

To expand functionality and to process data more efficiently, you can use Data Engineering Integration in conjunction with other Informatica products.

Data Engineering Integration integrates with the following Informatica products:

- PowerExchange adapters. Connect to data sources through adapters.
- Enterprise Data Catalog. Perform data lineage analysis for sources and targets.
- Enterprise Data Preparation. Discover raw data and publish it in a lake as a Hive table.



- Data Quality. Perform address validation and data discovery.
- Data Replication. Replicate change data to a Hadoop Distributed File System (HDFS).
- Data Transformation. Process complex file sources from the Hadoop environment.
- Data Engineering Streaming. Stream data as messages, and process it as it becomes available.
- Edge Data Streaming. Collect and ingest data in real time to a Kafka queue.
- Dynamic Data Masking. Mask or prevent access to sensitive data.

## CHAPTER 2

# Before You Begin

This chapter includes the following topics:

- [Read the Release Notes, 18](#)
- [Verify System Requirements, 18](#)
- [Uninstall Big Data Management, 21](#)
- [Prepare Directories, Users, and Permissions, 23](#)
- [Configure Access to Secure Hadoop Clusters, 30](#)
- [Configure the Metadata Access Service, 36](#)
- [Configure the Data Integration Service, 37](#)
- [Install Python for Enterprise Data Preparation, 42](#)

## Read the Release Notes

Read the Release Notes for updates to the installation and upgrade process. You can also find information about known and fixed issues for the release.

Find the Release Notes on the Informatica [documentation portal](#).

## Verify System Requirements

Verify that your environment meets the minimum system requirements for the installation process, disk space requirements, port availability, and third-party software.

For more information about product requirements and supported platforms, see the Product Availability Matrix on Informatica Network:

<https://network.informatica.com/community/informatica-network/product-availability-matrices>

## Verify Product Installations

Before you begin Data Engineering Integration between the domain and Hadoop environments, verify that Informatica and third-party products are installed.

You must install the following products:

### Informatica domain and clients

Install and configure the Informatica domain and the Developer tool. The Informatica domain must have a Model Repository Service, a Data Integration Service, and a Metadata Access Service.

### Hadoop File System and MapReduce

The Hadoop installation must include a Hive data warehouse with a non-embedded database for the Hive metastore. Verify that Hadoop is installed with Hadoop File System (HDFS) and MapReduce on each node. Install Hadoop in a single node environment or in a cluster. For more information, see the Apache website: <http://hadoop.apache.org>.

### Database client software

To access relational databases in the Hadoop environment, install database client software and drivers on each node in the cluster.

## Verify HDFS Disk Space

When the Data Integration Service integrates the domain with the Hadoop cluster, it uploads the Informatica binaries to HDFS.

Verify with the Hadoop administrator that the distributed cache has at least 1.5 GB of free disk space.

## Verify the Hadoop Distribution

Verify the version of the Hadoop distribution in the Hadoop environment.

Data Engineering Integration supports the following Hadoop distributions:

- Amazon EMR
- Azure HDInsight

**Note:** Verify that the HDInsight cluster is of "Hadoop" or "Spark" type.

- Cloudera CDH
- Cloudera CDP
- Hortonworks HDP
- MapR

In each release, Informatica can add, defer, and drop support for the non-native distributions and distribution versions. Informatica might reinstate support for deferred versions in a future release. To see a list of the latest supported versions, see the Product Availability Matrix on the Informatica Customer Portal:

<https://network.informatica.com/community/informatica-network/product-availability-matrices>.

## Verify Port Requirements

Open a range of ports to enable the Informatica domain to communicate with the Hadoop cluster and the distribution engine.

To ensure access to ports, the network administrator needs to complete additional tasks in the following situations:

- The Hadoop cluster is behind a firewall. Work with the network administrator to open a range of ports that a distribution engine uses.

- The Hadoop environment uses Azure HDInsight.
  - When the domain is deployed on the Azure platform in the same VNet as the HDInsight cluster, you do not have to open ports.
  - When the domain is on the Azure platform in a separate VNet from the HDInsight cluster, open ports on both VNets to enable communication between the domain and the cluster and its storage resources.
  - When the domain is on-premises, open ports on the Azure VNet to accept communication from the domain.
  - Work with the network administrator to enable VPN between the Informatica domain and the Azure cloud network.

The following table lists the ports to open:

Port	Description
7180	Cluster management web app for Cloudera. Required for Cloudera only.
8020	NameNode RPC. Required for all supported distributions except MapR.
8032	ResourceManager. Required for all distributions.
8080	Cluster management web app. Used by distributions that use Ambari to manage the cluster: HDInsight, Hortonworks.
8088	Resource Manager web app. Required for all distributions.
8443	MapR control system. Required for MapR only.
9080	Blaze monitoring console. Required for all distributions if you run mappings using Blaze.
9083	Hive metastore. Required for all distributions.
12300 to 12600	Default port range for the Blaze distribution engine. A port range is required for all distributions if you run mappings using Blaze.
19888	YARN JobHistory server webapp. Optional for all distributions.
50070	HDFS Namenode HTTP. Required for all distributions.

**Note:** The network administrators must ensure that the port used by the Metadata Access Service is accessible from the cluster nodes.

### Spark Engine Monitoring Port

Spark engine monitoring requires the cluster nodes to communicate with the Data Integration Service over a socket. The Data Integration Service picks the socket port randomly from the port range configured for the domain. You can view the port range in the advanced properties of the primary node. By default, the minimum port number is 12000 and the maximum port number is 13000. The network administrators must ensure that the port range is accessible from the cluster nodes to the Data Integration Service. If the administrators cannot provide a port range access, you can configure the Data Integration Service to use a fixed port with the SparkMonitoringPort custom property. The network administrator must ensure that the configured port is accessible from the cluster nodes to the Data Integration Service.

# Uninstall Big Data Management

If you are upgrading from version 10.1.1 or earlier and have a previous version installed on the Hadoop environment, Informatica recommends that you uninstall the previous version.

<b>Perform this task in the following situation:</b>
- You upgraded from version 10.1.1 or earlier.



## Uninstall for Amazon EMR, Azure HDInsight, and MapR

Complete the following prerequisite tasks before you uninstall:

1. Verify that the Big Data Management administrator can run `sudo` commands.
2. If you are uninstalling in a cluster environment, configure the root user to use a passwordless Secure Shell (SSH) connection between the machine where you want to run the Big Data Management uninstall and all of the nodes where Big Data Management is installed.
3. If you are uninstalling in a cluster environment using the `HadoopDataNodes` file, verify that the `HadoopDataNodes` file contains the IP addresses or machine host names of each of the nodes in the Hadoop cluster from which you want to uninstall Big Data Management. The `HadoopDataNodes` file is located on the node from where you want to launch the installation. You must add one IP address or machine host name of the nodes in the Hadoop cluster for each line in the file.

Complete the following tasks to perform the uninstallation:

1. Log in to the machine as root user. The machine you log in to depends on the environment and uninstallation method.
  - To uninstall in a single node environment, log in to the machine on which Big Data Management is installed.
  - To uninstall in a cluster environment using the `HADOOP_HOME` environment variable, log in to the primary name node.
  - To uninstall in a cluster environment using the `HadoopDataNodes` file, log in to any node.
2. Run the following command to start the uninstallation in console mode:

```
bash InformaticaHadoopInstall.sh
sh InformaticaHadoopInstall.sh
./InformaticaHadoopInstall.sh
```
3. Press **y** to accept the terms of agreement.
4. Press **Enter**.
5. Select **3** to uninstall.
6. Press **Enter**.
7. Select the uninstallation option, depending on the environment:
  - Select **1** to uninstall from a single node environment.
  - Select **2** to uninstall from a cluster environment.
8. Press **Enter**.

9. If you are uninstalling in a cluster environment, select the uninstallation option, depending on the uninstallation method:
  - Select **1** to uninstall from the primary name node.
  - Select **2** to uninstall using the `HadoopDataNodes` file.
10. Press **Enter**.
11. If you are uninstalling from a cluster environment from the primary name node, type the absolute path for the Hadoop installation directory. Start the path with a slash.

The uninstaller deletes all of the binary files from the following directory: `/<Big Data Management installation directory>/Informatica`

In a cluster environment, the uninstaller deletes the binary files from all nodes within the Hadoop cluster.

## Uninstall for Cloudera CDH

Uninstall Big Data Management on Cloudera from the Cloudera Manager.

1. In Cloudera Manager, browse to **Hosts > Parcels > Informatica**.
2. Select **Deactivate**.  
Cloudera Manager stops the Informatica Big Data Management instance.
3. Select **Remove**.  
The cluster uninstalls Informatica Big Data Management.

## Uninstall for Hortonworks HDP

To uninstall the stack deployment of Big Data Management, you use the Ambari configuration manager to stop and deregister the Big Data Management service, and then perform manual removal of Informatica files from the cluster.

1. In the Ambari configuration manager, select **INFORMATICA BDM** from the list of services.
2. Click the **Service Actions** dropdown menu and select **Delete Service**.
3. To confirm that you want to delete Informatica Big Data Management, perform the following steps:
  - a. In the **Delete Service** dialog box, click **Delete**.
  - b. In the Confirm Delete dialog box, type `delete` and then click **Delete**.
  - c. When the deletion process is complete, click **OK**.

Ambari stops the Big Data Management service and deletes it from the listing of available services. To fully delete Big Data Management from the cluster, continue with the next steps.
4. In a command window, delete the `INFORMATICABDM` folder from the following directory on the name node of the cluster: `/var/lib/ambari-server/resources/stacks/<Hadoop distribution>/<Hadoop version>/services/`
5. Delete the `INFORMATICABDM` folder from the following location on all cluster nodes where it was installed: `/var/lib/ambari-agent/cache/stacks/<Hadoop distribution>/<Hadoop version>/services`
6. Perform the following steps to remove RPM binary files:
  - a. Run the following command to determine the name of the RPM binary archive:

```
rpm -qa |grep Informatica
```

- b. Run the following command to remove RPM binary files:

```
rpm -ev <output_from_above_command>
```

For example:

```
rpm -ev InformaticaHadoop-10.1.1-1.x86_64
```

7. Repeat the previous step to remove RPM binary files from each cluster node.
8. Delete the following directory, if it exists, from the name node and each client node: `/opt/Informatica/`.
9. Repeat the last step on each cluster node where Big Data Management was installed.
10. On the name node, restart the Ambari server.

## Prepare Directories, Users, and Permissions

The Data Integration Service needs access to the Hadoop environment for integration and staging.

### Informatica Users

Verify and create the following Informatica users and grant them permission to access the cluster and its storage resources and staging directories:

- Hadoop impersonation user
- Service Principal Name (SPN) for the Data Integration Service
- Hadoop staging user
- Blaze user
- Operating system profile user
- Mapping impersonation user

### Directories

Prepare the following directories and permissions:

- Informatica cluster staging directory
- Hive warehouse directory
- Hive staging directory
- Blaze engine directories
- Spark engine staging directory
- Reject file directory

## Verify and Create Users

The Data Integration Service requires different users to access the Hadoop environment. Any user that you create for an Azure HDInsight distribution must be an Azure Active Directory user. For other distributions, use Linux users.

One or more of the following users with access to the cluster is known as the "Informatica user" on the cluster.

## Hadoop impersonation user

Verify that every node on the cluster has an impersonation user that can be used in a Hadoop connection. Create one if it does not exist. The Data Integration Service impersonates this user to run jobs in the Hadoop environment.

The following distributions use a Hadoop impersonation user:

### Azure HDInsight

To run Sqoop mappings on the Spark engine, add the Hadoop impersonation user as a Linux user on the machine that hosts the Data Integration Service.

**Note:** If the impersonation user contains mixed case characters, add the realm name along with the impersonation user.

Also see [“Verify and Create Users for HDInsight” on page 25](#).

### Cloudera CDP Public Cloud

The Hadoop impersonation user must have access to the Hive warehouse directory.

### MapR

If the MapR distribution uses Ticket or Kerberos authentication, the name must match the system user that starts the Informatica daemon and the gid of the user must match the gid of the MapR user.

## Service principal name (SPN) for the Data Integration Service

If the cluster uses Kerberos authentication, verify that the SPN corresponding to the cluster keytab file matches the name of the system user that starts the Informatica daemon.

## Hadoop staging user

Optionally, create an HDFS user that performs operations on the cluster staging directory. If you do not create a staging user, the Data Integration Service uses the operating system user that starts the Informatica daemon.

## Blaze user

Optionally, create an operating system user account that the Blaze engine uses to write to staging and log directories. If you do not create a Blaze user, the Data Integration Service uses the Hadoop impersonation user.

## Operating system profile user

If operating system profiles are configured for the Data Integration Service, the Data Integration Service runs jobs with permissions of the operating system user that you define in the profile. You can choose to use the operating system profile user instead of the Hadoop impersonation users to run jobs in a Hadoop environment.

To use operating system profile users with Cloudera CDP Public Cloud, configure an impersonation user, add the impersonation user to FreeIPA, and map the user to a cloud role using Knox IDBroker.

## Mapping impersonation user

A mapping impersonation user is valid for the native run time environment. Use mapping impersonation to impersonate the Data Integration Service user that connects to Hive, HBase, or HDFS sources and targets that use Kerberos authentication. Configure functionality in the Data Integration Service and the mapping properties. The mapping impersonation user uses the following format:

```
<Hadoop service name>/<host name>@<Kerberos realm>
```



## Verify and Create Users for HDInsight

Any user that you create for an Azure HDInsight distribution must be an Azure Active Directory user. The Azure platform uses Azure Active Directory and Azure AD Connect to synchronize Active Directory users to the HDInsight cluster.

Verify that Informatica users exist in Azure Active Directory (AAD), and synchronize users from AAD to the cluster. See the [Azure documentation](#).

**Note:** If synchronization results in duplication of a user in Azure Active Directory, the duplicated user cannot access cluster storage resources or be able to run a mapping on the cluster.

For information about how the Informatica domain interacts with Azure Active Directory, see the *Informatica Security Guide*.

If an Azure HDInsight cluster uses Enterprise Security Package, and mappings access ADLS storage, grant the required permissions. For the permissions, see [“Grant Access to Azure ADLS Resources for Informatica Users” on page 25](#).

## Grant Access to Azure ADLS Resources for Informatica Users

To run mappings that read from or write to ADLS storage resources, grant access permissions to Informatica users on the resources.

Depending on your Data Engineering Integration implementation, the Informatica user may be one or more of the following:

- Hadoop impersonation user
- Service principal name (SPN)
- Hadoop staging user
- Blaze user
- Operating system (OS) profile user

For more information about these users, see [“Verify and Create Users” on page 23](#).

## Grant Access Permissions to ADLS Gen1 Storage

Depending on how the HDInsight cluster is secured, perform the tasks in this section to grant access to ADLS Gen1 storage.

### Enterprise Security Package-enabled Clusters

If an Azure HDInsight cluster is enabled with the Enterprise Security Package, grant the following permissions to all the users:

- Execute permission on the root folder and its subfolders of the ADLS Gen1 account.
- Read and execute permissions on the following directory and its contents: `/hdp/apps/<version>`
- Read, write, and execute permissions on the following directories:

- `/tmp`
- `/app-logs`
- `/hive/warehouse`
- `/blaze/workdir`
- `/user`

```
/var/log/hadoop-yarn/apps  
/mr-history  
/tezstaging  
/mapreducestaging
```

**Note:** If the directories are not available, create the directories and grant the required permissions.

- Assign the Owner role to the Azure HDInsight cluster service principal display name. See [“Assigning the Owner Role to the Service Principal User” on page 26](#).

## Non-Enterprise Security Package-enabled Clusters

If the cluster is not enabled with the Enterprise Security Package:

- Populate the following property in core-site.xml:
  - dfs.adls.oauth2.client.id

For more information, see [“Configure \\*-site.xml Files for Azure HDInsight” on page 74](#).

## Assigning the Owner Role to the Service Principal User

To use ADLS Gen1 storage with an HDInsight cluster, assign the Owner role to the Azure HDInsight cluster service principal user.

1. In the Azure portal web interface, browse to **Home > Storage accounts**
2. Select the ADLS Gen1 storage resource associated with the HDInsight cluster.
3. Select **Access control (IAM)**.
4. Select the **Role** assignments tab.
5. Click **+ Add role assignment**.
6. In the **Role** dropdown, select "Owner."
7. In the **Assign Access to** dropdown, select "Azure AD user, group, or service principal."

## Grant Access Permissions to ADLS Gen2 Storage

When the HDInsight cluster uses an ADLS Gen2 storage resource, grant permission to the storage for Informatica users.

### Run Scripts to Automate Permissions Setting

ADLS Gen2 does not support recursive permissions, so you must set permissions on directories and files individually through the command line. To make permission setting easier, Informatica provides scripts that you can download, customize and run. You can also set permissions manually.

To use the scripts, download and uncompress the [script .zip archive](#) from the Informatica documentation portal. Follow the instructions in the archive Readme file to grant access to Informatica users, and then run the scripts on a Windows machine.

### Perform Manual Steps to Set Permissions

In addition to the readme file, the script archive contains a file titled Folders\_provided\_permission.docx. Page 1 of this file contains a list of twenty-one directories from (a.) to (v.) and the permissions that the script grants to each directory.

Before running the script, manually create the following directories listed under (e.) and (f.) on the list:

- /blaze
- /blaze/workdir

After running the script, perform the following additional steps, which are also listed at the bottom of page 1 of the `Folders_provided_permission.docx` file:

1. Grant read and execute permissions on the blob container root directory where the ADLS Gen2 account resides.
2. Add read, write and execute permissions for Access and Default permissions on the Spark staging directory specified in the Hadoop connection for the Spark impersonation user.

## Create Directories and Set Permissions

Prepare the following directories and permissions:

- Informatica cluster staging directory
- Hive warehouse directory
- Hive staging directory
- Blaze engine directories
- Spark engine staging directory
- Reject file directory

### Create a Cluster Staging Directory

Optionally, create a directory on HDFS that the Data Integration Service uses to stage the Informatica binary archive files.

By default, the Data Integration Service writes the files to the HDFS directory `/tmp`.

Grant permission to the Hadoop staging user and all blaze users. If you did not create a Hadoop staging user, the Data Integration Services uses the operating system user that starts the Informatica daemon.

### Grant Permissions on the Hive Warehouse Directory

Grant access to the absolute HDFS file path of the default database for the hive warehouse.

Grant read and write permissions on the Hive warehouse directory. You can find the location of the warehouse directory in the `hive.metastore.warehouse.dir` property of the `hive-site.xml` file. For example, the default might be `/user/hive/warehouse` or `/apps/hive/warehouse`.

Grant permission to the Hadoop impersonation user. Optionally, you can assign `-777` permissions on the directory.

### Create a Hive Staging Directory

The Blaze and Spark engines require access to the Hive staging directory. You can use the default directory, or you can create a directory on HDFS. For example, if you create a directory, you might run the following command:

```
hadoop fs -mkdir /staging
```

If you use the default directory or create a directory, you must grant execute permission to the Hadoop impersonation user and the mapping impersonation users.

## Create a Spark Staging Directory

When the Spark engine runs a job, it stores temporary files in a staging directory.

Optionally, create a staging directory on HDFS for the Spark engine. For example:

```
hadoop fs -mkdir -p /spark/staging
```

If you want to write the logs to the Informatica Hadoop staging directory, you do not need to create a Spark staging directory. By default, the Data Integration Service uses the HDFS directory `/tmp/SPARK_<user name>`.

Grant permission to the following users:

- Hadoop impersonation user
- SPN of the Data Integration Service
- Mapping impersonation users

Optionally, you can assign -777 permissions on the directory.

If you create a staging directory on a CDP Data Hub cluster, grant Access Control List (ACL) permissions for the staging directory to the Hive user and the impersonation user. To grant ACL permissions, run the following command on the CDP Data Hub cluster:

```
hadoop fs -setfacl -m user:user:rwX <staging directory>
```

## Create a Sqoop Staging Directory

When you run Sqoop jobs on the Spark engine, the Data Integration Service creates a Sqoop staging directory named `sqoop_staging` within the Spark staging directory by default. You can configure the Spark staging directory that you want to use in the Hadoop connection.

However, based on your processing requirements, you might need to create the directory manually and give write permissions to the Hive super user. When you create the `sqoop_staging` directory manually, the Data Integration Service uses this directory instead of creating another one.

Create a Sqoop staging directory named `sqoop_staging` manually in the following situations:

- You run a Sqoop pass-through mapping on the Spark engine to read data from a Sqoop source and write data to a Hive target that uses the Text format.
- You use a Cloudera CDH cluster with Sentry authorization, a Cloudera CDP cluster with Ranger authorization, or a Hortonworks HDP cluster with Ranger authorization.

After you create the `sqoop_staging` directory, you must add an Access Control List (ACL) for the `sqoop_staging` directory and grant write permissions to the Hive super user. Run the following command on the Cloudera CDH cluster or the Hortonworks HDP cluster to add an ACL for the `sqoop_staging` directory and grant write permissions to the Hive super user:

```
hdfs dfs -setfacl -m default:user:hive:rwX /<Spark staging directory>/sqoop_staging/
```

For information about Sentry authorization, see the Cloudera documentation. For information about Ranger authorization, see the Hortonworks documentation.

## Create Blaze Engine Directories

Create a blaze user account and directories required by the Blaze engine.

Complete the following tasks to prepare the Hadoop cluster for the Blaze engine:

### Create a home directory for the blaze user.

If you created a blaze user, create home directory for the blaze user. For example,

```
hdfs hadoop fs -mkdir /user/blaze
hdfs hadoop fs -chown blaze:blaze /user/blaze
```

If you did not create a blaze user, the Hadoop impersonation user is the default user.

### Optionally, create a local services log directory.

By default, the Blaze engine writes the service logs to the YARN distributed cache. For example, run the following command:

```
mkdir -p /opt/informatica/blazeLogs
```

\$HADOOP\_NODE\_INFA\_HOME gets set to the YARN distributed cache. If you create a directory, you must update the value of the advanced property in the Hadoop connection.

### Create an aggregated HDFS log directory.

Create a log directory on HDFS to contain aggregated logs for local services. For example:

```
hadoop fs -mkdir -p /var/log/hadoop-yarn/apps/informatica
```

Ensure that value of the advanced property in the Hadoop connection matches the directory that you created.

### Optionally, create a Blaze staging directory.

You can write the logs to the Informatica Hadoop staging directory, or you can create a Blaze staging directory. If you do not want to use the default location, create a staging directory on the HDFS. For example:

```
hadoop fs -mkdir -p /blaze/workdir
```

**Note:** If you do not create a staging directory, clear the Blaze staging directory property value in the Hadoop connection and the Data Integration Service uses the HDFS directory `/tmp/blaze_<user name>`.

If you create a staging directory on a CDP Data Hub cluster, grant Access Control List (ACL) permissions for the staging directory to the Hive user and the impersonation user. To grant ACL permissions, run the following command on the CDP Data Hub cluster:

```
hadoop fs -setfacl -m user:user:rwX <staging directory>
```

### Grant permissions on the local services log directory, aggregated HDFS log directory, and the staging directory.

Grant permission to the following users:

- Blaze user
- Hadoop impersonation user
- Mapping impersonation users

If the blaze user does not have permission, the Blaze engine uses a different user, based on the cluster security and the mapping impersonation configuration. Blaze users must also have write permissions on `/tmp`.

## Create a Reject File Directory

You can choose to store reject files on HDFS for the Blaze and Spark engines.

Reject files can be very large, and you can choose to write them to HDFS instead of the Data Integration Service machine. You can configure the Hadoop connection object to write to the reject file directory.

Grant permission to the following users:

- Blaze user

- Hadoop impersonation user
- Mapping impersonation users

If the blaze user does not have permission, the Blaze engine uses a different user, based on the cluster security and the mapping impersonation configuration.

## Create a Proxy Directory for MapR

If the Hadoop cluster runs on MapR, you must create a proxy directory for the user who will impersonate other users.

Verify the following requirements for the proxy user:

- Create a user or verify that a user exists on every Data Integration Service machine and on every node in the Hadoop cluster.
- Verify that the uid and the gid of the user match in both environments.
- Verify that a directory exists for the user on the cluster. For example, `/opt/mapr/conf/proxy/<user name>`

# Configure Access to Secure Hadoop Clusters

If the Hadoop cluster uses Kerberos authentication or SSL/TLS, you must configure the Informatica domain to access the cluster. If the cluster uses transparent encryption, you must configure the Key Management Server (KMS) for Informatica user access.

Depending on the security implementation on the cluster, you must perform the following tasks:

### Cluster uses Kerberos authentication.

You must configure the Kerberos configuration file on the Data Integration Service machine to match the Kerberos realm properties of the Hadoop cluster. Verify that the Hadoop Kerberos properties are configured in the Data Integration Service and the Metadata Access Service.

If the cluster uses Kerberos and does not use the Enterprise Security Package, you can configure Ranger authorization separately to grant permissions to Informatica users.

### Cluster uses SSL/TLS.

You must import security certificates to the Data Integration Service and the Metadata Access Service machines. See [“Configuring Access to an SSL/TLS-Enabled Cluster” on page 31](#).

### Cluster uses transparent encryption.

If the transparent encryption uses Cloudera Java KMS, Cloudera Navigator KMS, or Apache Ranger KMS, you must configure the KMS for Informatica user access.

### Cluster uses Enterprise Security Package.

If the cluster uses Enterprise Security Package and ADLS Gen1 or Gen2 storage, perform the following tasks:

- Create a keytab file on any one of the cluster nodes for the specific user. To create a keytab file, see [“Using the ktutil Utility to Create a Keytab File” on page 35](#).
- Configure proxy users in `core-site.xml`. Populate the following properties:
  - `hadoop.proxyuser.<SPN username>.groups=user1,user2,user3`

- `hadoop.proxyuser.<SPN username>.hosts=user1,user2,user3`
- `hadoop.proxyuser.<SPN username>.users=user1,user2,user3`

... where <SPN username> is the Kerberos Service Principal User. For more information, see [“Configure \\*-site.xml Files for Azure HDInsight” on page 74](#).

- If the cluster uses ADLS Gen1 storage, assign the Owner role to the Azure HDInsight cluster service principal display name. See [“Assigning the Owner Role to the Service Principal User” on page 26](#).

#### Cluster uses WASB Storage

If the cluster uses WASB storage, configure the storage account access key in `core-site.xml`. Populate the following property:

- `fs.azure.account.key.<your account>.blob.core.windows.net`

For more information, see [“Configure \\*-site.xml Files for Azure HDInsight” on page 74](#).

## Configuring Access to an SSL/TLS-Enabled Cluster

When you use an SSL-enabled or TLS-enabled cluster, you must configure the Informatica domain to communicate with the secure cluster.

Based on the cluster distribution that uses SSL, you perform the following tasks:

#### Amazon EMR cluster uses SSL/TLS

Import security certificates from the cluster to the Informatica domain. If you created a Hive or S3 connection object manually, configure the connection string properties to access the SSL-enabled cluster.

#### Cloudera CDH, Cloudera CDP, Dataproc, or Hortonworks HDP cluster uses SSL

Import security certificates from the cluster to the Informatica domain. If you created a Hive connection manually, configure the connection string properties to access the SSL-enabled cluster.

#### MapR cluster uses SSL

Make sure that the MapR client is configured to communicate with a secure cluster. If you created a Hive connection object manually, configure the connection string properties to access the SSL-enabled cluster.

## Configure the Hive Connection for SSL-Enabled Clusters

If you created the Hive connection when you created cluster configurations, the cluster configuration creation wizard enables access to a cluster that uses SSL. If you manually created a Hive connection, you must configure the connection string properties to enable access to a cluster that uses SSL.

If you manually created a Hive connection, add the following property-value pair to the metadata connection string and data access connection string properties:

```
ssl=true
```

For example:

```
jdbc:hive2://<hostname>:<port>/<db>;ssl=true
```

**Note:** Insert the `ssl=true` flag before the `kerberos principal` element when you create the Hive connection manually.

## Import Security Certificates from an SSL-Enabled Cluster

When you use custom, special, or self-signed security certificates to secure the Hadoop cluster, Informatica services that connect to the cluster require these certificates to be present on the machines that run the application services. Use the keytool utility to import certificates from the cluster.

For more information about the keytool utility, refer to the Oracle documentation.

If a cluster uses SSL, perform the following steps to import security certificates from the cluster to the Data Integration Service and Metadata Access Service machines:

1. Run the following keytool -exportcert command on the cluster to export the certificates:

```
keytool -exportcert
  -alias <alias name>
  -keystore <custom.truststore file location>
  -file <exported certificate file location>
  -storepass <password>
```

Where:

- -alias specifies the alias name associated with the truststore file.
- -keystore specifies the location of the truststore file on the cluster.
- -file specifies the file name and location for the exported certificate file.
- -storepass specifies the password for the keystore on the cluster.

The keytool -exportcert command produces a certificate file associated with the alias.

2. Run the following keytool -importcert command on one Data Integration Service machine to import the security certificates:

```
keytool -importcert -trustcacerts
  -alias <alias name>
  -file <exported certificate file location>
  -keystore <java cacerts location>
  -storepass <password>
```

Where:

- -alias specifies the alias name associated with the certificate file.
- -file specifies the file name and location of the exported certificate file.
- -keystore specifies the location of the truststore file on the domain.
- -storepass specifies the password for the keystore on the domain.

**Important:** Import the certificate files one time and then copy them to all machines that host the Data Integration Service and Metadata Access Service. If the Data Integration Service runs on a grid, mappings that you push to the Hadoop environment can fail with initialization errors due to inconsistent binary hex values.

Depending on whether the Informatica domain uses SSL, you specify the keystore location as follows:

- If the domain is SSL-enabled, import the certificate file to the following location:  
<Informatica installation directory>/services/shared/security/infa\_truststore.jks
- If the domain is not SSL-enabled, import the certificate file to the following location:  
<Informatica installation directory>/java/jre/lib/security/cacerts

The keytool -importcert command imports the security certificates to the keystore location you specify.

### Example. Import Security Certificates

The environment includes a Cloudera CDH cluster that uses SSL and an Informatica domain that does not use SSL. You export the security certificate for the user bigdata\_user1 from the custom.keystore on the



Cloudera CDH cluster to the file `exported.cer`. Then, you import the `export.cer` certificate file to the Informatica domain location.

1. Run the following export command:

```
keytool -exportcert -alias bigdata_user1 -keystore ~/custom.truststore -file ~/exported.cer
```

2. Run the following import command on the Data Integration Service machine:

```
keytool -importcert -alias bigdata_user1 -file ~/exported.cer -keystore <Informatica installation directory>/java/jre/lib/security/cacerts
```

3. Copy the certificate file to all other machines that host the Data Integration Service and the Metadata Access Service.

## Rules and Guidelines for Importing Security Certificates from an SSL-Enabled Cluster

Consider the following rules and guidelines when you import security certificates from an SSL-enabled cluster:

- If a MapR cluster is SSL-enabled, you do not have to import the security certificates. Make sure that the MapR client on the Data Integration Service and Metadata Access Service machines is configured to access an SSL-enabled cluster.
- If a Cloudera CDP cluster is Auto-TLS enabled, import the security certificates before you import the cluster configuration into the domain.
- After you import certificates from a Cloudera CDP cluster, configure the `LD_LIBRARY_PATH` environment variable for the Data Integration Service. Add the following path to the beginning of the value:  
`<Informatica installation directory>/CDH_7.1/lib/native`

## Import Security Certificates from a TLS-Enabled Domain

When the domain is configured to use TLS, you must import the certificates to the default or custom truststore file that the Informatica domain uses.

### Default truststore file

If the domain is TLS-enabled and the cluster uses server managed keys, you must import the Baltimore CyberTrust Root certificate to the default truststore file.

Use the `keytool` utility to import the security certificate.

The default truststore file is located in the following directory: `<Informatica installation home>/services/shared/security/infa_truststore.jks`

### Custom truststore file

If the domain is TLS-enabled with a custom truststore file, and the cluster uses server managed keys, get the custom truststore file location from Informatica Administrator, and then import the Baltimore CyberTrust Root certificate to the custom truststore file.

Use the `keytool` utility to import the security certificate.

To get the custom truststore file location, perform the following steps:

1. In the Administrator tool, click the Manage tab.
2. Click the Services and Nodes view.
3. In the Domain Navigator, select the domain.
4. Get the custom truststore file location from the domain properties.

You can download the Baltimore CyberTrust Root certificates from <https://www.digicert.com/digicert-root-certificates.htm>.

For more information about downloading the certificates, see <https://docs.microsoft.com/en-us/azure/java-add-certificate-ca-store>.

## Generate the OAUTH Token

The cluster uses an OAUTH token for authentication for the impersonation user, the Hadoop staging user, and any other user who runs jobs on the cluster from the Data Integration Service.

Log in to Ambari Web UI with the Azure Active Directory user credentials to generate the OAUTH token for authentication for the Informatica users. For definitions of these users, see [“Verify and Create Users” on page 23](#).

**Note:** The HDInsight cluster requires an OAUTH token only if the mapping accesses ADLS Gen1 resources.

## Generating Keytab Files for the SPN User

Use the Microsoft Windows Server ktpass utility to generate a keytab file for the SPN user.

You must generate the keytab files on a member server or on a domain controller within the Active Directory domain. You cannot generate keytab files on a workstation operating system such as Microsoft Windows 10.

To use ktpass to generate a keytab file, run the following command:

```
ktpass.exe -out <keytab filename> -princ <service principal name> -mapuser <user account> [-pass <user account password>] -crypto <keys> -ptype <principal type>
```

The following table describes the command options:

Option	Description
-out	The file name of the Kerberos keytab file to generate as shown under the <code>KEY_TAB_NAME</code> column in the <code>SPNKeytabFormat.txt</code> file.
-princ	The service principal name displayed under the <code>SPN</code> column in the <code>SPNKeytabFormat.txt</code> file.
-mapuser	The Active Directory user account to associate with the SPN.
-pass	The password set in Active Directory for the Active Directory user account, if applicable.
-crypto	Specifies the key types generated in the keytab file. Set to all to use all supported cryptographic types.
-ptype	The principal type. Set to <code>KRB5_NT_PRINCIPAL</code> .

When you run ktpass, you associate each node account and HTTP process account with the corresponding SPN in Active Directory. The following table shows the association between the accounts and the SPNs described in this article:

User Account	Keytab Type	SPN
nodeuser01	NODE_SPN	isp/node01/InfraDomain/COMPANY.COM
httpuser01	NODE_HTTP_SPN	HTTP/US001DEV.company.com@COMPANY.COM

User Account	Keytab Type	SPN
nodeuser02	NODE_SPN	isp/node02/Infadomain/COMPANY.COM
httpuser02	NODE_HTTP_SPN	HTTP/US005DEV.company.com@COMPANY.COM
nodeuser03	NODE_SPN	isp/node03/Infadomain/COMPANY.COM

## Generate the Keytab File

Generate a keytab file for the SPN user.

1. Create a keytab file for the Kerberos principal user account that you created for each node in Active Directory.

Copy the file name from the `KEY_TAB_NAME` column in the `SPNKeytabFormat.txt` file.

The following example creates a keytab file for the `nodeuser01` user:

```
ktpass.exe -out node01.keytab -princ isp/node01/Infadomain/COMPANY.COM -mapuser
nodeuser01 -pass password -crypto all -ptype KRB5_NT_PRINCIPAL
```

2. Create a keytab file for each HTTP process Kerberos principal user account that you created in Active Directory.

If the domain uses Kerberos cross realm authentication, the principal user account can exist in any Kerberos realm the domain uses.

Copy the keytab file name from the `KEY_TAB_NAME` column in the `SPNKeytabFormat.txt` file. Copy the service principal name from the `SPN` column in the `SPNKeytabFormat.txt` file.

The following example creates a keytab file for a Kerberos principal user account named `httpuser01`:

```
ktpass.exe -out webapp_http.keytab -princ HTTP/US001DEV.company.com@COMPANY.COM -mapuser
httpuser01 -crypto all -ptype KRB5_NT_PRINCIPAL
```

3. Create a keytab for the LDAP bind user account that is used to access and search Active Directory during LDAP synchronization.

Structure the value for the `-princ` option as `<principal name>@<KERBEROS REALM>`. Include the name of the LDAP configuration for the Active Directory server in the keytab file name. Structure the keytab file name as follows: `<Active Directory LDAP configuration_name>.keytab`.

The following example creates a keytab file for a service principal user account named `ldapuser`:

```
ktpass.exe -out ActiveDirectoryServer1.keytab -princ ldapuser@COMPANY.COM -mapuser
ldapuser -crypto all -ptype KRB5_NT_PRINCIPAL
```

## Using the ktutil Utility to Create a Keytab File

Use the `ktutil` utility to create a keytab file.

Before you begin, get the Kerberos principal user name from the cluster administrator.

1. Log in to any cluster VM.
2. From the command line, type `ktutil` to launch the utility.
3. Type the following command:

```
addent -password -p <user name> -k 1 -e RC4-HMAC
```

where `<user name>` is the Kerberos principal user. For example:

```
addent -password -p myname -k 1 -e RC4-HMAC
```

4. When prompted, enter the password for the Kerberos principal user.
5. Type the following command to create a keytab:  

```
wkt /tmp/keytabs/<user name>.keytab
```

where `/tmp/keytabs/` is the path where you want to store keytabs. For example:  

```
wkt /tmp/keytabs/myname.keytab
```
6. Type `q` to quit the `ktutil` utility.

## Configure Apache Ranger with HDInsight

The HDInsight administrator can configure Apache Ranger authorization on the cluster.

Ranger is also configured when the cluster uses the Enterprise Security Package, which is a bundle of Ranger and Kerberos.

To configure Apache Ranger authorization with an HDInsight cluster, see the [Azure documentation](#).

## Configure the Metadata Access Service

Configure the Metadata Access Service to integrate with the Hadoop environment.

Perform this task in the following situations:
<ul style="list-style-type: none"><li>- You are integrating for the first time.</li><li>- You upgraded from version 10.2 or earlier.</li></ul>



The following table describes the Metadata Access Service properties that you need to configure:

Property	Description
Use Operating System Profiles and Impersonation	If enabled, the Metadata Access Service uses the operating system profiles to access the Hadoop cluster.
Hadoop Kerberos Service Principal Name	Service Principal Name (SPN) of the Metadata Access Service to connect to a Hadoop cluster that uses Kerberos authentication. Not applicable for the MapR distribution.
Use logged in user as impersonation user	Required if the Hadoop cluster uses Kerberos authentication. If enabled, the Metadata Access Service uses the impersonation user to access the Hadoop environment. Default is false.

# Configure the Data Integration Service

Configure the Data Integration Service to integrate with the Hadoop environment.

Perform the following pre-integration tasks:

1. Download Informatica Hadoop binaries to the Data Integration Service machine if the operating systems of the Hadoop environment and the Data Integration Service are different.
2. Configure the Data Integration Service properties, such as the cluster staging directory, Hadoop Kerberos service principal name, and the path to the Kerberos keytab file.
3. Prepare an installation of Python on the Data Integration Service machine or on the Hadoop cluster if you plan to run the Python transformation.
4. Copy the krb5.conf file to the following location on the machine that hosts the Data Integration Service:
  - <Informatica installation directory>/java/jre/lib/security
  - <Informatica installation directory>/services/shared/security
5. Copy the keytab file to the following directory: <Informatica installation directory>/isp/config/keys

## Download the Informatica Server Binaries for the Hadoop Environment

If the domain and the Hadoop environments use different supported operating systems, you must configure the Data Integration Service to be compatible with the Hadoop environment. To run a mapping, the local path to the Informatica server binaries must be compatible with the Hadoop operating system.

The Data Integration Service can synchronize the following operating systems: SUSE and Redhat

The Data Integration Service machine must include the Informatica server binaries that are compatible with the Hadoop cluster operating system. The Data Integration Service uses the operating system binaries to integrate the domain with the Hadoop cluster.

You must run the installer to extract the installation binaries into custom Hadoop OS path and then exit the installer.

1. Create a directory on the Data Integration Service host machine to store the Informatica server binaries associated with the Hadoop operating system.

If the Data Integration Service runs on a grid, Informatica recommends extracting the files to a location that is shared by all services on the grid. If the location is not shared, you must extract the files to all Data Integration Service machines that run on the grid.

The directory names in the path must not contain spaces or the following special characters: @ | \* \$ # ! % ( ) { } [ ]
2. Download and extract the Informatica server binaries from the Informatica download site. For example,

```
tar -xvf <Informatica server binary tar file>
```
3. Run the installer to extract the installation binaries into the custom OS path.

Perform the following steps to run the installer:

  - Run the `sh Server/install.bin -DINSTALL_MODE=CONSOLE -DINSTALL_TYPE=0` file.
  - Press **Y** to continue the installation.
  - Press **1** to install Informatica Data Engineering products.
  - Press **3** to run the installer.

- Press **2** to accept the terms and conditions.
  - Press **2** to continue the installation for Data Engineering products only.
  - Press **2** to configure the Informatica domain to run on a network with Kerberos authentication.
  - Enter the path and file name of the Informatica license key and press an option to tune the services.
  - Enter the custom Hadoop OS path.
  - Type **Quit** to quit the installation.
4. Set the custom Hadoop OS path in the Data Integration Service and then restart the service.
  5. Optionally, you can delete files that are not required. For example, run the following command:
 

```
rm -Rf <Informatica server binary file> ./source/*.7z
```

**Note:** If you subsequently install an Informatica EBF, you must also install it in the path of the Informatica server binaries associated with the Hadoop environment.

## Edit the etc/hosts File

Edit the etc/hosts file on the domain hosts and on client machines if you want to run mappings with the Blaze engine, or if you want to process mappings on an Azure HDInsight or Google Dataproc cluster.

### Edit the hosts File for the Blaze Engine

To run the Blaze engine on every node in the cluster, verify that the /etc/hosts file on every node has entries for all other nodes.

Each node in the cluster requires an entry for the IP address and the fully qualified domain name (FQDN) of all other nodes. For example,

```
127.0.0.1 localhost node1.node.com
208.164.186.1 node1.node.com node1
208.164.186.2 node2.node.com node2
208.164.186.3 node3.node.com node3
```

Changes take effect after you restart the network.

### Edit the Hosts File for Access to Azure HDInsight

Ensure that Informatica can access the HDInsight cluster by updating the /etc/hosts file on all machines that host the Data Integration Service.

#### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

### Configure Dynamic Updates for the Cluster Headnode Host

Identify the headnode host, and schedule a script to regularly update the entry.

The HDInsight cluster designates one cluster node as the headnode. If the node fails or stops for maintenance, the cluster designates another node as the headnode. When this happens, the headnodehost

entry in `/etc/hosts` files on Informatica domain nodes requires updating. You can schedule a script to perform this update.

1. In the `/etc/hosts` file on each machine that hosts the Data Integration Service, enter the IP address, DNS name, and DNS short name for each data node on the cluster. Use `headnodehost` to identify the host as the cluster headnode host.

For example:

```
10.20.30.40 hn0-rndhdi.abcdefghaouniivfp3bet13d.ix.internal.cloudapp.net
headnodehost
```

2. Download the [headnode\\_update\\_script.zip](#) file from the Informatica documentation portal and uncompress it to get `headnode_update_script.sh`.  
The script gets the the IP address, DNS name, and DNS short name for the head node host from the `x-ms-hdi-active` property on the cluster. The script then replaces the value of `headnodehost` in the hosts file with the `x-ms-hdi-active` value.
3. Schedule the script to run regularly on each machine that hosts the Data Integration Service. Set up a schedule based on your requirements, such as daily updates.

### Optional: Configure IP Addresses for ADLS Storage

If the HDInsight cluster is integrated with ADLS storage, you also need to enter the IP addresses and DNS names for the hosts listed in the cluster property `fs.azure.datalake.token.provider.service.urls`.

For example:

```
1.2.3.67 gw1-ltsa.1320suh5abcdefghgaz0izgnhe.gx.internal.cloudapp.net
1.2.3.68 gw0-ltsa.1320suh5abcdefghgaz0izgnhe.gx.internal.cloudapp.net
```

**Note:** To get the IP addresses, run a telnet command from the cluster host using each host name found in the `fs.azure.datalake.token.provider.service.urls` property.

### Edit the hosts File for Google Dataproc

To enable Informatica to access the Dataproc cluster, edit the `/etc/hosts` file on the domain host.

Enter the IP address, DNS name, and DNS short name for each data node on the cluster.

For example:

```
10.20.30.40 dataprocABC-m.c.MyUsernameINFA54321.internal dataprocABC-m
10.20.30.44 dataprocABC-w-0.c.MyUsernameINFA54321.internal dataprocABC-w-0
10.20.30.43 dataprocABC-w-1.c.MyUsernameINFA54321.internal dataprocABC-w-1
10.20.30.42 dataprocABC-w-2.c.MyUsernameINFA54321.internal dataprocABC-w-2
```

Edit the hosts file with cluster node information on each machine that hosts the Data Integration Service.

## Configuring LZO Compression Format

To write `.jar` files in the LZO compression format, you must copy the files for LZO compression to the machine where the Data Integration Service runs.

Perform the following steps to configure the Data Integration Service for LZO compression:

1. Copy the `lzo.jar` file from the cluster to the following directory on the machine on which the Data Integration Service runs: `<Informatica installation directory>/<distribution>/infaLib`
2. Copy the LZO native binaries from the cluster to one of the following directories on the machine on which the Data Integration Service runs:
  - `<Informatica installation directory>/<distribution>/lib/native`
  - `<Informatica installation directory>/<distribution>/lib/native/Linux-amd64-64` for MapR clusters

3. On the Data Integration Service **Processes** tab, add or update the LD\_LIBRARY\_PATH environment variable to include the path the to LZO native binaries on the Data Integration Service machine.
4. Restart the Data Integration Service.

## Configuring the Data Integration Service to Use Operating System Profiles

Configure the Data Integration Service to run mappings, workflows, and profiling jobs with operating system profiles.

The operating system user you define in the operating system profile must have access to the directories you configure in the operating system profile and to the directories the Data Integration Service accesses at run time. For example, pmsuid is a tool that the DTM process, command tasks, and parameter files use to switch between operating system users. You must provide permissions to operating system users to run pmsuid with the permissions of the Data Integration Service administrator user.

The operating system profile user must be a member of the Data Integration Service user group.

**Note:** If you enable the Data Integration Service to use operating system profiles, you cannot enable cache connection, the SQL Service Module, and the Web Service Module.

Complete the following steps to configure the Data Integration Service to use operating system profiles:

1. Configure system permissions on the files and directories that the operating system profile user needs access at run time.
2. In the Administrator tool, enable the Data Integration Service to use operating system profiles.
3. On the Security page of the Administrator tool, create operating system profiles.

## Configure Data Integration Service Properties

The Data Integration Service contains properties that integrate the domain with the Hadoop cluster.

The following table describes the Data Integration Service properties that you need to configure:

Property	Description
Cluster Staging Directory	The directory on the cluster where the Data Integration Service pushes the binaries to integrate the native and non-native environments and to store temporary files during processing. Default is /tmp.
Hadoop Staging User	The HDFS user that performs operations on the Hadoop staging directory. The user requires write permissions on Hadoop staging directory. Default is the operating system user that starts the Informatica daemon.
Custom Hadoop OS Path	<p>The local path to the Informatica server binaries compatible with the Hadoop operating system. Required when the Hadoop cluster and the Data Integration Service are on different supported operating systems. The Data Integration Service uses the binaries in this directory to integrate the domain with the Hadoop cluster. The Data Integration Service can synchronize the following operating systems:</p> <ul style="list-style-type: none"> <li>- SUSE and Redhat</li> </ul> <p>Include the source directory in the path. For example, &lt;Informatica server binaries&gt;/source.</p> <p>Changes take effect after you recycle the Data Integration Service.</p> <p><b>Note:</b> When you install an Informatica EBF, you must also install it in this directory.</p>



Property	Description
Hadoop Kerberos Service Principal Name	Service Principal Name (SPN) of the Data Integration Service to connect to a Hadoop cluster that uses Kerberos authentication. Not required for the MapR distribution.
Hadoop Kerberos Keytab	The file path to the Kerberos keytab file on the machine on which the Data Integration Service runs. Not required for the MapR distribution.
Custom Properties	<p>Properties that are unique to specific environments.</p> <p>You can configure run-time properties for the Hadoop environment in the Data Integration Service, the Hadoop connection, and in the mapping. You can override a property configured at a high level by setting the value at a lower level. For example, if you configure a property in the Data Integration Service custom properties, you can override it in the Hadoop connection or in the mapping. The Data Integration Service processes property overrides based on the following priorities:</p> <ol style="list-style-type: none"> <li>1. Mapping custom properties set using <code>infacmd ms runMapping</code> with the <code>-cp</code> option</li> <li>2. Mapping run-time properties for the Hadoop environment</li> <li>3. Hadoop connection advanced properties for run-time engines</li> <li>4. Hadoop connection advanced general properties, environment variables, and classpaths</li> <li>5. Data Integration Service custom properties</li> </ol> <p><b>Note:</b> When a mapping uses Hive Server 2 to run a job or parts of a job, you cannot override properties that are configured on the cluster level in preSQL or post-SQL queries or SQL override statements. Workaround: Instead of attempting to use the cluster configuration on the domain to override cluster properties, pass the override settings to the JDBC URL. For example: <code>beeline -u "jdbc:hive2://&lt;domain host&gt;:&lt;port_number&gt;/tpch_text_100" --hiveconf hive.execution.engine=tez</code></p>

## Prepare a Python Installation

If you want to use the Python transformation, you must ensure that the worker nodes on the Hadoop cluster contain an installation of Python. You must complete different tasks depending on the product that you use.

### Installing Python for Data Engineering Integration

To use the Python transformation in a mapping, the worker nodes on the cluster must contain a uniform installation of Python. You can ensure that the installation is uniform in one of the following ways:

#### Verify that the Python installation exists.

Verify that all worker nodes on the cluster contain an installation of Python in the same directory, such as `/usr/lib/python`, and that each Python installation contains all required modules. You do not reinstall Python, but you must reconfigure the following Spark advanced property in the Hadoop connection:

```
infaspark.pythontx.executorEnv.PYTHONHOME
```

#### Install Python.

Install Python on every Data Integration Service machine. You can create a custom installation of Python that contains specific modules that you can reference in the Python code. When you run mappings, the Python installation is propagated to the worker nodes on the cluster.

If you choose to install Python on the Data Integration Service machines, complete the following tasks:

1. Install Python.

2. Optionally, install any third-party libraries such as numpy, scikit-learn, and cv2. You can access the third-party libraries in the Python transformation.
3. Copy the Python installation folder to the following location on the Data Integration Service machine:

```
<Informatica installation directory>/services/shared/spark/python
```

**Note:** If the Data Integration Service machine already contains an installation of Python, you can copy the existing Python installation to the above location.

Changes take effect after you recycle the Data Integration Service.

## Installing Python for Data Engineering Streaming

To use the Python transformation in Data Engineering Streaming, you must install Python and the Jep package. Because you must install Jep, the Python version that you use must be compatible with Jep. You can use one of the following versions of Python:

2.7  
3.3  
3.4  
3.5  
3.6

To install Python and Jep, complete the following tasks:

1. Install Python with the **--enable-shared** option to ensure that shared libraries are accessible by Jep.
2. Install Jep. To install Jep, consider the following installation options:
  - Run `pip install jep`. Use this option if Python is installed with the pip package.
  - Configure the Jep binaries. Ensure that `jep.jar` can be accessed by Java classloaders, the shared Jep library can be accessed by Java, and Jep Python files can be accessed by Python.
3. Optionally, install any third-party libraries such as numpy, scikit-learn, and cv2. You can access the third-party libraries in the Python transformation.
4. Copy the Python installation folder to the following location on the Data Integration Service machine:

```
<Informatica installation directory>/services/shared/spark/python
```

**Note:** If the Data Integration Service machine already contains an installation of Python, you can copy the existing Python installation to the above location.

Changes take effect after you recycle the Data Integration Service.

# Install Python for Enterprise Data Preparation

Enterprise Data Preparation uses the Apache Solr indexing capabilities to provide recommendations of related data assets. Apache Solr requires Python modules.

You must install Python with the following modules on every node that hosts the Interactive Data Preparation Service associated with Enterprise Data Preparation:

argparse  
sys  
getopt

os  
urllib  
httplib2  
ConfigParser

## CHAPTER 3

# Amazon EMR Integration Tasks

This chapter includes the following topics:

- [Amazon EMR Task Flows, 44](#)
- [Prepare for Cluster Import from Amazon EMR, 49](#)
- [Configure Glue as the Hive Metastore , 55](#)
- [Create a Cluster Configuration, 56](#)
- [Verify or Refresh the Cluster Configuration , 57](#)
- [Verify JDBC Drivers for Sqoop Connectivity, 58](#)
- [Configure the Files to Use S3 , 59](#)
- [Set S3 Access Policies, 60](#)
- [Configure the Developer Tool, 62](#)
- [Complete Upgrade Tasks, 64](#)

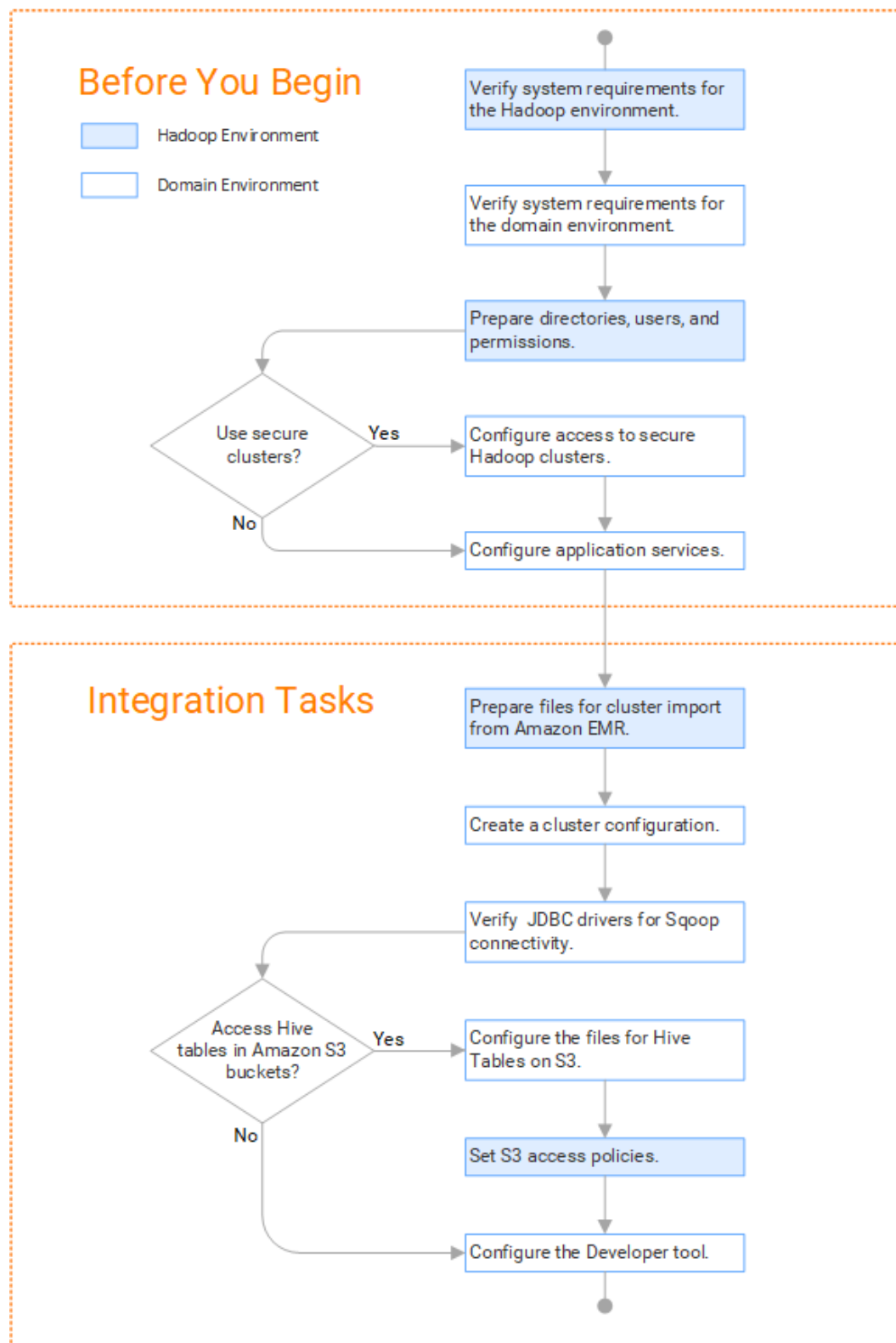
## Amazon EMR Task Flows

Depending on whether you want to integrate or upgrade Data Engineering Integration in an Amazon EMR environment, you can use the flow charts to perform the following tasks:

- Integrate the Informatica domain with Amazon EMR for the first time.
- Upgrade from version 10.2.1.
- Upgrade from version 10.2.
- Upgrade from a version earlier than 10.2.

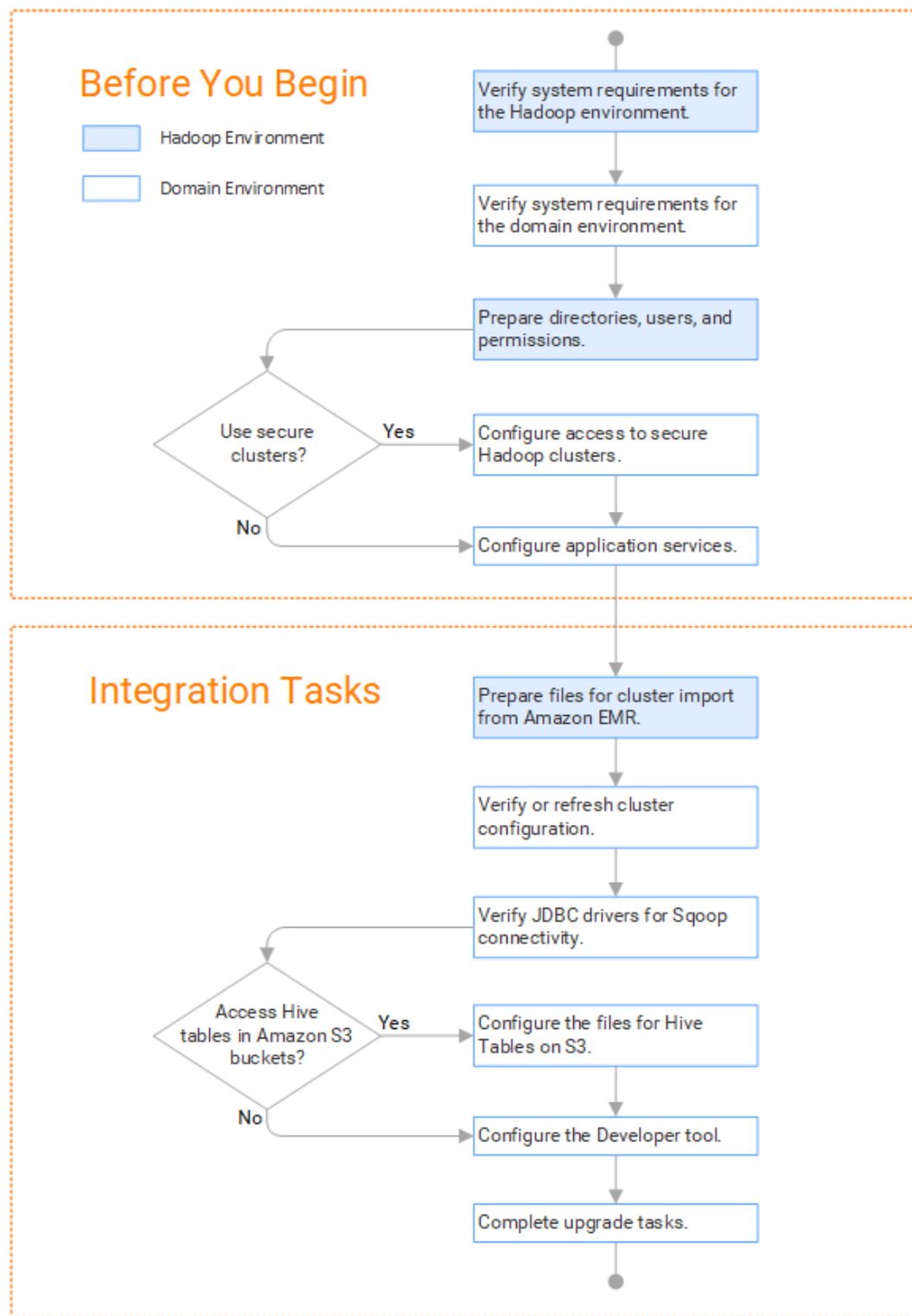
# Task Flow to Integrate with Amazon EMR

The following diagram shows the task flow to integrate the Informatica domain with Amazon EMR:



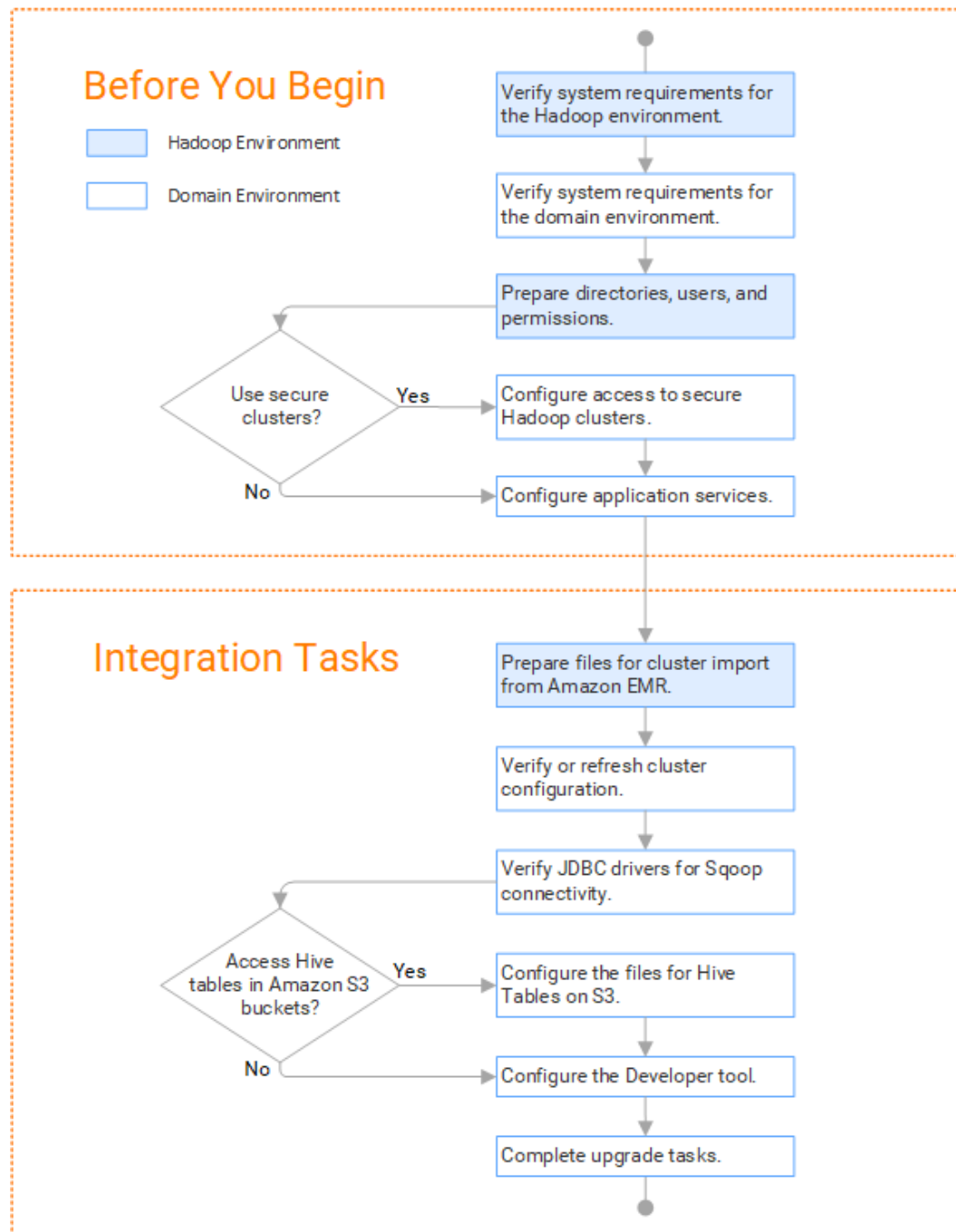
## Task Flow to Upgrade from Version 10.2.1

The following diagram shows the task flow to upgrade version 10.2.1 for Amazon EMR:



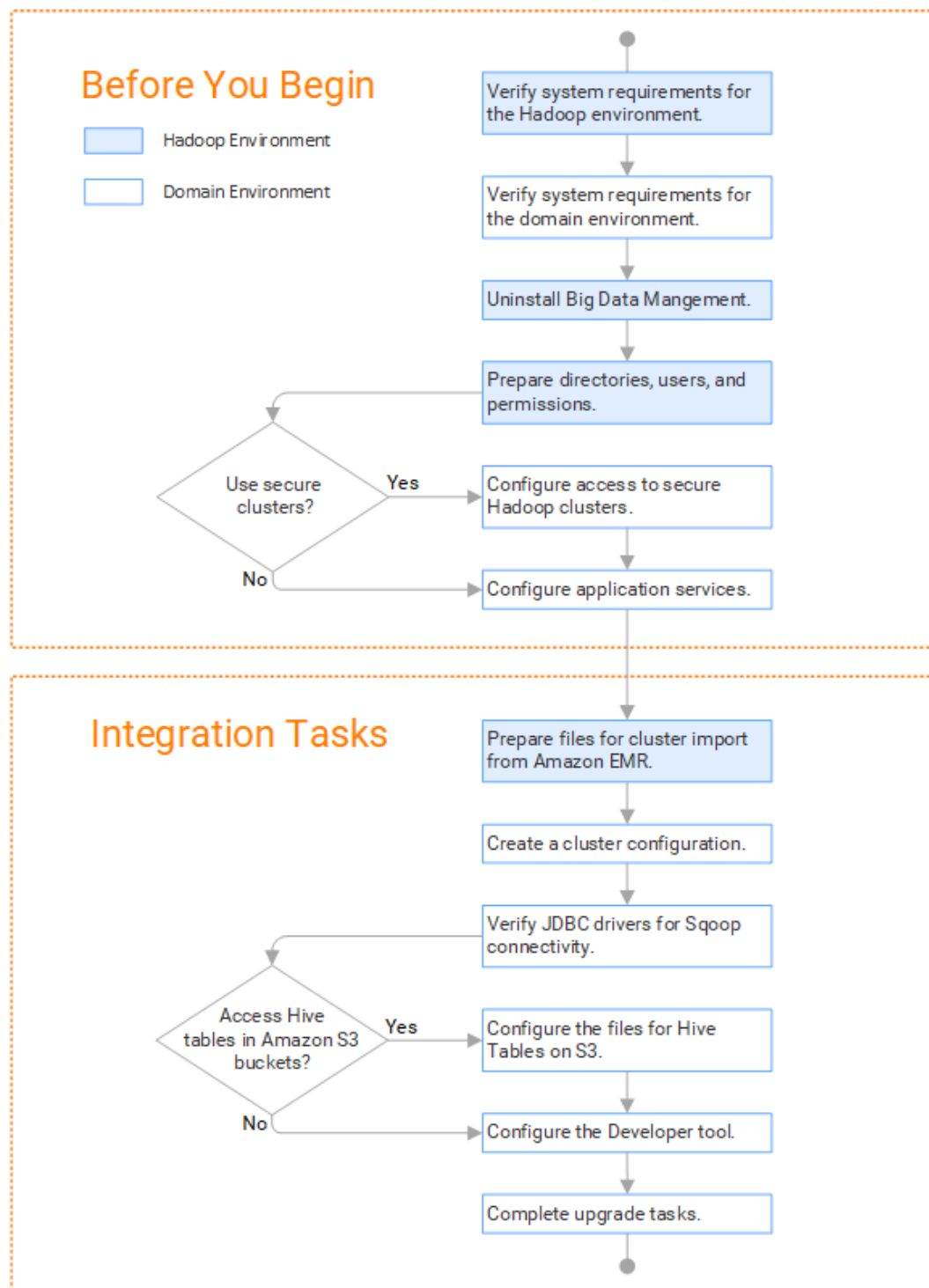
## Task Flow to Upgrade from Version 10.2

The following diagram shows the task flow to upgrade version 10.2 for Amazon EMR:



## Task Flow to Upgrade from a Version Earlier than 10.2

The following diagram shows the task flow to upgrade from a version earlier than 10.2 for Amazon EMR:





# Prepare for Cluster Import from Amazon EMR

Before the Informatica administrator can import cluster information to create a cluster configuration in the Informatica domain, the Hadoop administrator must perform some preliminary tasks.

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

**Note:** If you are upgrading from a previous version, verify the properties and suggested values, as Data Engineering Integration might require additional properties or different values for existing properties.

Complete the following tasks to prepare the cluster before the Informatica administrator creates the cluster configuration:

1. Verify property values in \*-site.xml files that the Data Integration Service needs to run mappings in the Hadoop environment.
2. Prepare the archive file to import into the domain.

**Note:** You cannot import cluster information directly from the Amazon EMR cluster into the Informatica domain.

## Configure \*-site.xml Files for Amazon EMR

The Hadoop administrator needs to configure \*-site.xml file properties and restart impacted services before the Informatica administrator imports cluster information into the domain.

### capacity-scheduler.xml

Configure the following properties in the capacity-scheduler.xml file:

**yarn.scheduler.capacity.<queue path>.disable\_preemption**

Disables preemption for the Capacity Scheduler. The Blaze engine does not support preemption. If YARN preemption is enabled for the cluster, you need to disable it for the queue allocated to the Blaze engine.

Set to TRUE for queue allocated to the Blaze engine.

### core-site.xml

Configure the following properties in the core-site.xml file:

**fs.s3.awsAccessKeyId**

The ID for the run-time engine to connect to the Amazon S3 file system. Required for the Blaze engine and for the Spark engine if the Data Integration if S3 policy does not allow EMR access, or if you use EMRFS access and the Informatica domain does not reside on an EC2 instance.

**Note:** If the Data Integration Service is deployed on an EC2 instance and the IAM roles and policies allow access to S3 and other resources, this property is not required. If the Data Integration Service is deployed on-premises, then you can choose to configure the value for this property in the cluster configuration on the Data Integration Service after you import the cluster configuration. Configuring the AccessKeyId value on the cluster configuration is more secure than configuring it in core-site.xml on the cluster.

Set to your access ID.

**fs.s3.awsSecretAccessKey**

The access key for the Blaze and Spark engines to connect to the Amazon S3 file system. Required for the Blaze engine and for the Spark engine if the Data Integration if S3 policy does not allow EMR access, or if you use EMRFS access and the Informatica domain does not reside on an EC2 instance.

**Note:** If the Data Integration Service is deployed on an EC2 instance and the IAM roles and policies allow access to S3 and other resources, this property is not required. If the Data Integration Service is deployed on-premises, then you can choose to configure the value for this property in the cluster configuration on the Data Integration Service after you import the cluster configuration. Configuring the AccessKeyID value on the cluster configuration is more secure than configuring it in core-site.xml on the cluster.

Set to your access key.

**fs.s3.enableServerSideEncryption**

Enables server side encryption for S3 buckets. Required for SSE and SSE-KMS encryption.

Set to: TRUE

**fs.s3a.server-side-encryption-algorithm**

The server-side encryption algorithm for S3. Required for SSE and SSE-KMS encryption. Set to the encryption algorithm used.

**fs.s3a.endpoint**

URL of the entry point for the web service.

For example:

```
<property>
  <name>fs.s3a.endpoint</name>
  <value>s3-us-west-1.amazonaws.com</value>
</property>
```

**fs.s3a.bucket.BUCKET\_NAME.server-side-encryption.key**

Server-side encryption key for the S3 bucket. Required if the S3 bucket is encrypted with SSE-KMS.

For example:

```
<property>
  <name>fs.s3a.bucket.BUCKET_NAME.server-side-encryption.key</name>
  <value>arn:aws:kms:us-west-1:*****</value>
  <source>core-site.xml</source>
</property>
```

where BUCKET\_NAME is the name of the S3 bucket.

**hadoop.proxyuser.<proxy user>.groups**

Defines the groups that the proxy user account can impersonate. On a secure cluster the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to group names of impersonation users separated by commas. If less security is preferred, use the wildcard " \* " to allow impersonation from any group.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

**hadoop.proxyuser.<proxy user>.hosts**

Defines the host machines that a user account can impersonate. On a secure cluster the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to a single host name or IP address, or set to a comma-separated list. If less security is preferred, use the wildcard "\*" to allow impersonation from any host.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.proxyuser.hive.hosts**

Comma-separated list of hosts that you want to allow the Hive user to impersonate on a non-secure cluster.

When `hive.server2.enable.doAs` is false, append a comma-separated list of Informatica server host names or IP address where the Data Integration Service is running. If less security is preferred, use the wildcard "\*" to allow impersonation from any host.

**Note:** After you make changes to this property, restart the cluster services that use core-site configuration values.

#### **hadoop.proxyuser.yarn.groups**

Comma-separated list of groups that you want to allow the YARN user to impersonate on a non-secure cluster.

Set to group names of impersonation users separated by commas. If less security is preferred, use the wildcard "\*" to allow impersonation from any group.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.proxyuser.yarn.hosts**

Comma-separated list of hosts that you want to allow the YARN user to impersonate on a non-secure cluster.

Set to a single host name or IP address, or set to a comma-separated list. If less security is preferred, use the wildcard "\*" to allow impersonation from any host.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.security.auth\_to\_local**

Translates the principal names from the Active Directory and MIT realm into local names within the Hadoop cluster. Based on the Hadoop cluster used, you can set multiple rules.

Set to: `RULE:[1:$1@$0](^.*@YOUR.REALM)s/^.*(.)@YOUR.REALM\.COM$/\1/g`

Set to: `RULE:[2:$1@$0](^.*@YOUR.REALM\.)s/^.*(.)@YOUR.REALM\.COM$/\1/g`

#### **io.compression.codecs**

Enables compression on temporary staging tables.

Set to a comma-separated list of compression codec classes on the cluster.

#### [fair-scheduler.xml](#)

Configure the following properties in the fair-scheduler.xml file:

##### **allowPreemptionFrom**

Enables preemption for the Fair Scheduler. The Blaze engine does not support preemption. If YARN preemption is enabled for the cluster, you need to disable it for the queue allocated to the Blaze engine.

Set to FALSE for the queue allocated to the Blaze engine.

For example:

```
<queue name="Blaze">
  <weight>1.0</weight>
  <allowPreemptionFrom>false</allowPreemptionFrom>
  <schedulingPolicy>fsp</schedulingPolicy>
  <aclSubmitApps>*</aclSubmitApps>
  <aclAdministerApps>*</aclAdministerApps>
</queue>
```

### hbase-site.xml

Configure the following properties in the hbase-site.xml file:

#### **hbase.use.dynamic.jars**

Enables metadata import and test connection from the Developer tool. Required for an HDInsight cluster that uses ADLS storage or an Amazon EMR cluster that uses HBase resources in S3 storage.

Set to: false

#### **zookeeper.znode.parent**

Identifies HBase master and region servers.

Set to the relative path to the znode directory of HBase.

### hive-site.xml

Configure the following properties in the hive-site.xml file:

#### **hive.compactor.initiator.on**

Runs the initiator and cleaner threads on metastore instance. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: TRUE

#### **hive.compactor.worker.threads**

The number of worker threads to run in a metastore instance. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: 1

#### **hive.conf.hidden.list**

Comma-separated list of hidden configuration properties.

Set to:

javax.jdo.option.ConnectionPassword,hive.server2.keystore.password,fs.s3n.awsAccessKeyId,fs.s3n.awsSecretAccessKey,fs.s3a.access.key,fs.s3a.secret.key,fs.s3a.proxy.password

#### **hive.enforce.bucketing**

Enables dynamic bucketing while loading to Hive. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: TRUE

#### **hive.exec.dynamic.partition**

Enables dynamic partitioned tables for Hive tables. Applicable for Hive versions 0.9 and earlier.

Set to: TRUE

#### **hive.exec.dynamic.partition.mode**

Allows all partitions to be dynamic. Required for the Update Strategy transformation in a mapping that writes to a Hive target. Also required if you use Sqoop and define a DDL query to create or replace a partitioned Hive target at run time.

Set to: nonstrict

**hive.support.concurrency**

Enables table locking in Hive. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: TRUE

**hive.txn.manager**

Turns on transaction support. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: org.apache.hadoop.hive.ql.lockmgr.DbTxnManager

**Note:** The following properties enable pre-task and post-task monitoring statistics for Amazon EMR jobs in the Developer tool:

**hive.async.log.enabled**

Enables asynchronous logging. Required when you enable pre-task and post-task monitoring statistics on a Dataproc cluster.

Set to: FALSE

**hive.server2.in.place.progress**

Allows HiveServer2 to send progress bar update information. Takes effect only when you enable Tez. Required when you enable pre-task and post-task monitoring statistics on a Dataproc cluster.

Set to: TRUE

**hive.server2.logging.operation.enabled**

Enables logs to be saved. Required when you enable pre-task and post-task monitoring statistics on a Dataproc cluster.

Set to: TRUE

**hive.server2.logging.operation.level**

Hive Server2 logging level at the session level. Requires hive.server2.logging.operation.enabled to be set to TRUE. Required when you enable pre-task and post-task monitoring statistics on a Dataproc cluster.

Set to: EXECUTION

[kms-site.xml](#)

Configure the following properties in the kms-site.xml file:

**hadoop.kms.authentication.kerberos.name.rules**

Translates the principal names from the Active Directory and MIT realm into local names within the Hadoop cluster. Based on the Hadoop cluster used, you can set multiple rules.

Set to: RULE:[1:\$1@\$0](^.\*@YOUR.REALM\.COM\$)s/^.\*(.)@YOUR.REALM\.COM\$/1/g

Set to: RULE:[2:\$1@\$0](^.\*@YOUR.REALM\.COM\$)s/^.\*(.)@YOUR.REALM\.COM\$/1/g

[mapred-site.xml](#)

Configure the following properties in the mapred-site.xml file:

**mapreduce.framework.name**

The run-time framework to run MapReduce jobs. Values can be local, classic, or yarn. Required for Sqoop.

Set to: yarn

**yarn.app.mapreduce.am.staging-dir**

The HDFS staging directory used while submitting jobs.

Set to the staging directory path.

**tez-site.xml**

Configure the following properties in the tez-site.xml file:

**tez.am.tez-ui.history-url.template**

Tez UI URL template for the application. The application manager uses this URL to redirect the user to the Tez UI. Required when you enable pre-task and post-task monitoring statistics on a Dataproc cluster.

Set value to:

```
_HISTORY_URL_BASE?%2F%23%2Ftez-app%2FAPPLICATION_ID
```

**Note:** The values of `_HISTORY_URL_BASE_` and `_APPLICATION_ID` are resolved at runtime. Do not edit the string to supply values.

**tez.task.generate.counters.per.io**

Enables pre-task and post-task monitoring statistics on an Amazon EMR or Dataproc cluster.

Set to: TRUE

**yarn-site.xml**

Configure the following properties in the yarn-site.xml file:

**yarn.application.classpath**

Required for dynamic resource allocation.

Add `spark_shuffle.jar` to the class path. The `.jar` file must contain the class `"org.apache.spark.network.yarn.YarnShuffleService."`

**yarn.nodemanager.resource.memory-mb**

The maximum RAM available for each container. Set the maximum memory on the cluster to increase resource memory available to the Blaze engine.

Set to 16 GB if value is less than 16 GB.

**yarn.nodemanager.resource.cpu-vcores**

The number of virtual cores for each container. Required for Blaze engine resource allocation.

Set to 10 if the value is less than 10.

**yarn.scheduler.minimum-allocation-mb**

The minimum RAM available for each container. Required for Blaze engine resource allocation.

Set to 6 GB if the value is less than 6 GB.

**yarn.nodemanager.vmem-check-enabled**

Disables virtual memory limits for containers. Required for the Blaze and Spark engines.

Set to: false

**yarn.nodemanager.aux-services**

Required for dynamic resource allocation for the Spark engine.

Add an entry for `"spark_shuffle."`

**yarn.nodemanager.aux-services.spark\_shuffle.class**

Required for dynamic resource allocation for the Spark engine.

Set to: `org.apache.spark.network.yarn.YarnShuffleService`

**yarn.resourcemanager.scheduler.class**

Defines the YARN scheduler that the Data Integration Service uses to assign resources.

Set to: org.apache.hadoop.yarn.server.resourcemanager.scheduler

**yarn.node-labels.enabled**

Enables node labeling.

Set to: TRUE

**yarn.node-labels.fs-store.root-dir**

The HDFS location to update node label dynamically.

Set to: <hdfs://[Node name]:[Port]/[Path to store]/[Node labels]/>

## Prepare the Archive File for Amazon EMR

After you verify property values in the \*-site.xml files, create a .zip or a .tar file that the Informatica administrator can use to import the cluster configuration into the domain.

Create an archive file that contains the following files from the cluster:

- core-site.xml
- hbase-site.xml. Required only if you access HBase sources and targets.
- hdfs-site.xml
- hive-site.xml
- mapred-site.xml or tez-site.xml. Include the mapred-site.xml file or the tez-site.xml file based on the Hive execution type used on the Hadoop cluster.
- yarn-site.xml

**Note:** To import from Amazon EMR, the Informatica administrator must use an archive file.

## Configure Glue as the Hive Metastore

You can configure Amazon Glue as the Hive metastore with an Amazon EMR 5.29 or 6.1 cluster.

To enable integration with an EMR cluster with Glue, copy .jar files from the cluster to the domain, and then enable the Hive metastore setting in the hive-site.xml configuration before you create the cluster configuration.

Consider the following rules and guidelines:

- Glue does not support Hive transactions.
- Amazon supports Glue only when the EMR cluster is not Kerberos enabled.

# Create a Cluster Configuration

After the Hadoop administrator prepares the cluster for import, the Informatica administrator must create a cluster configuration.

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from version 10.1.1 or earlier.

A cluster configuration is an object in the domain that contains configuration information about the Hadoop cluster. The cluster configuration enables the Data Integration Service to push mapping logic to the Hadoop environment. Import configuration properties from the Hadoop cluster to create a cluster configuration.

The import process imports values from \*-site.xml files into configuration sets based on the individual \*-site.xml files. When you perform the import, the cluster configuration wizard can create Hadoop, HBase, HDFS, and Hive connection to access the Hadoop environment. If you choose to create the connections, the wizard also associates the cluster configuration with the connections.

**Note:** If you are integrating for the first time and you imported the cluster configuration when you ran the installer, you *must* re-create or refresh the cluster configuration.

## Importing a Hadoop Cluster Configuration from a File

You can import properties from an archive file to create a cluster configuration.

Before you import from the cluster, you must get the archive file from the Hadoop administrator.

1. From the **Connections** tab, click the **ClusterConfigurations** node in the Domain Navigator.
2. From the Actions menu, select **New > Cluster Configuration**.

The **Cluster Configuration** wizard opens.

3. Configure the following properties:

Property	Description
Cluster configuration name	Name of the cluster configuration.
Description	Optional description of the cluster configuration.
Distribution type	The cluster Hadoop distribution type.



Property	Description
Method to import the cluster configuration	Choose <b>Import from file</b> to import properties from an archive file.
Create connections	<p>Choose to create Hadoop, HDFS, Hive, and HBase connections.</p> <p>If you choose to create connections, the Cluster Configuration wizard associates the cluster configuration with each connection that it creates.</p> <p>The Hadoop connection contains default values for properties such as cluster environment variables, cluster path variables, and advanced properties. Based on the cluster environment and the functionality that you use, you can add to the default values or change the default values of these properties. For a list of Hadoop connection properties to configure, see <a href="#">"Configuring Hadoop Connection Properties" on page 269</a>.</p> <p>If you do not choose to create connections, you must manually create them and associate the cluster configuration with them.</p> <p><b>Important:</b> When the wizard creates the Hive connection, it populates the Metadata Connection String and the Data Access Connection String properties with the value from the hive.metastore.uris property. If the Hive metastore and HiveServer2 are running on different nodes, you must update the Metadata Connection String to point to the HiveServer2 host.</p>

- Click **Browse** to select a file. Select the file and click **Open**.
- Click **Next** and verify the cluster configuration information on the summary page.

## Verify or Refresh the Cluster Configuration

You might need to refresh the cluster configuration or update the distribution version in the cluster configuration when you upgrade.

### Perform this task in the following situation:

- You upgraded from version 10.2 or later.

### Verify the Cluster Configuration

The cluster configuration contains a property for the distribution version. The verification task depends on the version you upgraded:

#### Upgrade from 10.2

If you upgraded from 10.2 and you changed the distribution version, you need to verify the distribution version in the General properties of the cluster configuration.

#### Upgrade from 10.2.1

Effective in version 10.2.1, Informatica assigns a default version to each Hadoop distribution type. If you configure the cluster configuration to use the default version, the upgrade process upgrades to the assigned default version if the version changes. If you have not upgraded your Hadoop distribution to Informatica's default version, you need to update the distribution version property.

For example, suppose the assigned default Hadoop distribution version for 10.2.1 is  $n$ , and for 10.2.2 is  $n+1$ . If the cluster configuration uses the default supported Hadoop version of  $n$ , the upgraded cluster

configuration uses the default version of  $n+1$ . If you have not upgraded the distribution in the Hadoop environment you need to change the cluster configuration to use version  $n$ .

If you configure the cluster configuration to use a distribution version that is not the default version, you need to update the distribution version property in the following circumstances:

- Informatica dropped support for the distribution version.
- You changed the distribution version.

### Refresh the Cluster Configuration

If you updated any of the \*-site.xml files noted in the topic to prepare for cluster import, you need to refresh the cluster configuration in the Administrator tool.

## Verify JDBC Drivers for Sqoop Connectivity

Verify that you have the JDBC drivers to access JDBC-compliant databases in the Hadoop environment. You might need separate drivers for metadata import and for run-time processing.

#### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from version 10.2.1 or earlier.

You download drivers based on design-time and run-time requirements:

- **Design-time.** To import metadata, you can use the DataDirect drivers packaged with the Informatica installer if they are available. If they are not available, use any Type 4 JDBC driver that the database vendor recommends.
- **Run-time.** To run mappings, use any Type 4 JDBC driver that the database vendor recommends. Some distributions support other drivers to use Sqoop connectors. You cannot use the DataDirect drivers for run-time processing.

## Verify Design-time Drivers

Use the DataDirect JDBC drivers packaged with the Informatica installer to import metadata from JDBC-compliant databases. If the DataDirect JDBC drivers are not available for a specific JDBC-compliant database, download the Type 4 JDBC driver associated with that database.

Copy the JDBC driver .jar files to the following location on the Developer tool machine:

```
<Informatica installation directory>\clients\externaljdbcjars
```

## Verify Run-time Drivers

Verify run-time drivers for mappings that access JDBC-compliant databases in the Hadoop environment. Use any Type 4 JDBC driver that the database vendor recommends.

1. Download Type 4 JDBC drivers associated with the JCBC-compliant databases that you want to access.

2. To optimize the Sqoop mapping performance on the Spark engine while writing data to an HDFS complex file target of the Parquet format, download the following .jar files:
  - `parquet-hadoop-bundle-1.6.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-hadoop-bundle/1.6.0/>
  - `parquet-avro-1.6.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-avro/1.6.0/>
  - `parquet-column-1.5.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-column/1.5.0/>
3. Copy all of the .jar files to the following directory on the machine where the Data Integration Service runs:

```
<Informatica installation directory>\externaljdbcjars
```

Changes take effect after you recycle the Data Integration Service. At run time, the Data Integration Service copies the .jar files to the Hadoop distribution cache so that the .jar files are accessible to all nodes in the cluster.

## Configure the Files to Use S3

To run mappings using S3 resources, you must configure the files from the master node to the Data Integration Service machine.

### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any Informatica version and changed the distribution version.

Perform the following steps to configure master node files to integrate with EMR:

1. Copy the `hadoop-assembly` .JAR file.
2. Create an AWS configuration file on the Data Integration Service machine.
3. Create and configure an environment variable.
4. Copy and replace the `hadoop-common` .JAR file.
5. Recycle the Data Integration Service.

Perform these steps regardless of the version of EMR you are integrating with, but note that the version number parts of file names may vary depending on the version.

### Example

This example contains file names that support integration with EMR 5.20.

#### Copy the .jar file

To integrate with EMR 5.20, get `emrfs-hadoop-assembly-2.29.0.jar` from the Hadoop administrator. Copy the file to the following locations on each Data Integration Service machine:

```
/<Informatica installation directory>/services/shared/hadoop/EMR_5.16/lib
```

Required when you run mapping on the Spark engine.

### Create a file

Create a `~/.aws/config` on the Data Integration Service machine. The file must contain the AWS location.

For example,

```
[default] region=us-west-2
```

Required when you run mapping on the Spark engine.

### Create an environment variable

Create the `AWS_CONFIG_FILE` environment variable on the Data Integration Service machine. Set the value to `<EMR_5.16>/conf/aws.default`

Required when you run mapping on the Spark and Blaze engines.

### Copy and replace a file

Copy `hadoop-common-2.8.5-amzn-1.jar` from the following location in the EMR 5.20 cluster:

```
/usr/lib/hadoop
```

Replace the file in the following location:

```
<Informatica installation directory>/services/shared/hadoop/EMR_5.16/lib
```

Required when you run mapping on the Spark engine.

### Recycle the Data Integration Service

You must recycle the Data Integration Service to reflect the changes.

## Set S3 Access Policies

The AWS administrator must set S3 access policies to grant users the required access to S3 resources.

#### Perform this task in the following situations:

- You are integrating for the first time.

S3 access policies allow control of user access to S3 resources and the actions that users can perform. The AWS administrator uses policies to control access and actions for specific users and resources, depending on the use case that mappings and workflows require.

AWS uses a JSON statement for S3 access policies. To set the S3 access policy, determine the principal, actions, and resources to define, then create or edit an existing S3 access policy JSON statement.

For more information about Amazon S3 access policies, see [AWS documentation](#).

## Step 1. Identify the S3 Access Policy Elements

Identify the principal, actions, and resources to insert in the access policy.

The following table describes the tags to set in the access policy:

Tag	Description
Principal	The user, service, or account that receives permissions that are defined in a policy. Assign the owner of the S3 bucket resources as the principal. <b>Note:</b> The S3 bucket owner and the owner of resources within the bucket can be different.
Action	The activity that the principal has permission to perform. In the sample, the Action tag lists two put actions and one get action. You must specify both get and put actions to grant read and write access to the S3 resource.
Resource	The S3 bucket, or folder within a bucket. Include only resources in the same bucket.

### Sample S3 Policy JSON Statement

The following JSON statement contains the basic elements of an S3 bucket access policy:

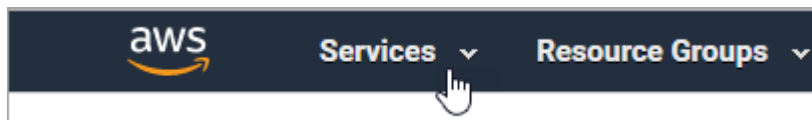
```
{
  "Version": "<date>",
  "Id": "Allow", "Statement": [
    { "Sid": "<Statement ID>", "Effect": "Allow",
      "Principal": {
        "AWS": "arn:aws:iam::<account_2_ID>:<user>"
      },
      "Action": [
        "s3:PutObject", "s3:PutObjectAcl",
        "s3:GetObject"
      ],
      "Resource": [
        "arn:aws:s3:::<bucket_1_name>/foldername/*"
      ]
    }
  ]
}
```

## Step 2. Optionally Copy an Existing S3 Access Policy as a Template

When the AWS administrator selects a role for cluster users, the AWS console generates a default access policy. After the AWS console generates the default policy, you can copy it and customize it to grant access to specific resources to specific users.

Complete the following steps to copy an existing S3 access policy:





1. In the AWS console, click the **Services** menu.  
The image below shows the **Services** menu in the menu bar:



2. Type "IAM" in the search bar and press Enter.  
The **Welcome to Identity and Access Management** screen opens.

3. In the menu on the left, select **Policies**.  
The console displays a list of existing policies.
4. Type "S3" in the search bar and press Enter.  
The console displays a list of existing S3 access policies.

The image below shows an example of a list of S3 access policies:

Filter: Policy type ▾		Q S3	Showing 4 results	
	Policy name ▾	Type	Description	
<input type="radio"/>	▶  AmazonDMSRedshiftS3Role	AWS managed	Provides access to manage S3 settings for Redshift endpoints for DMS.	
<input type="radio"/>	▶  AmazonS3FullAccess	AWS managed	Provides full access to all buckets via the AWS Management Console.	
<input type="radio"/>	▶  AmazonS3ReadOnlyAccess	AWS managed	Provides read only access to all buckets via the AWS Management Console.	
<input type="radio"/>	▶  QuickSightAccessForS3StorageMan...	AWS managed	Policy used by QuickSight team to access customer data produced by S3 Storage ...	

5. Click the name of the policy that you want to copy.  
The policy opens in a read-only window.
6. Highlight and copy the policy statement.

After you copy the JSON statement, you can edit it in a text editor or in the bucket policy editor.

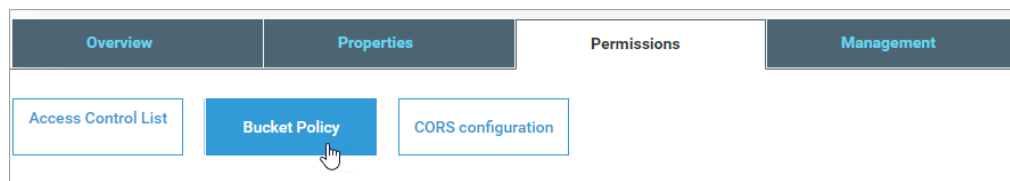
## Step 3. Create or Edit an S3 Access Policy

Create an S3 access policy or edit an existing policy. The AWS administrator can enter a JSON statement, based on a template. The administrator can copy and customize the S3 policy from another bucket.

1. In the AWS console, click the **Services** menu.
2. In the **Storage** section, choose **S3**.  
The AWS console displays a list of existing buckets.
3. Use the **search box** to find the bucket you want to set a policy for, and select the bucket from the results.
4. Click the **Permissions** tab, then click **Bucket Policy**.

The **Bucket Policy Editor** opens.

The image below shows the **Bucket Policy** button:



5. Type the bucket access policy, or edit the existing policy, and click **Save**.  
AWS applies the access policy to the bucket.

## Configure the Developer Tool

You can configure the Developer tool to enable you to import complex files or import metadata when the domain is Kerberos-enabled.

Edit the developerCore.ini file to import complex files. Edit the file on each Developer tool machine.

## Configure developerCore.ini

Edit the developerCore.ini file to import complex files.

Edit the developerCore.ini file on each machine that hosts the Developer tool.

When you import a complex file, such as Avro or Parquet, the imported object includes metadata associated with the distribution in the Hadoop environment. If the file resides on the Developer tool machine, the import process picks up the distribution information from the developerCore.ini file. You must edit the developerCore.ini file to point to the distribution directory on the Developer tool machine.

You can find the developerCore.ini file in the following directory:

```
<Informatica installation directory>\clients\DeveloperClient
```

Add the following property:

```
-DINFA_HADOOP_DIST_DIR=hadoop<distribution>_<version>
```

The change takes effect when you restart the Developer tool.

## Configure the Developer Tool for Kerberos

To import metadata from Hive, HBase, and complex file sources, import configuration files from the Kerberos-enabled cluster, and generate the Kerberos credentials file on the Developer tool machine.

### Import configuration files

The Hadoop cluster uses a set of XML files named \*-site.xml to store configuration settings. The domain uses the same set of files to create the cluster configuration object.

To enable you to import metadata from the cluster, import the \*-site.xml files to each Developer tool machine:

1. Log in to the Administrator tool and navigate to **Connections > Cluster Configuration > CCO**. Locate the cluster configuration associated with the Hadoop cluster.
2. Extract the \*-site.xml files in the cluster configuration, including sensitive properties, to the following directory on the Developer tool machine: `<Informatica installation directory>\clients\DeveloperClient\hadoop\<Hadoop distribution>\conf`  
For more information about sensitive properties, see the "Active Properties View" topic in the *Data Engineering Administrator Guide*.

**Note:** If you refresh the cluster configuration, repeat these steps.

### Generate the Kerberos credentials file

1. Copy the krb5.conf file from `<Developer tool installation directory>/services/shared/security` to `C:/Windows`.
2. Rename krb5.conf to krb5.ini.
3. In the krb5.ini file, verify the value of the forwardable option to determine how to use the kinit command. If `forwardable=true`, run the command with the `-f` option. Otherwise, run the command without the `-f` option.
4. To generate the Kerberos credentials file, run the kinit command from the following location: `<Developer tool installation directory>/clients/java/bin/kinit.exe`  
For example, you might run the following command: `kinit joe/domain12345@MY-REALM`

# Complete Upgrade Tasks

If you upgraded the Informatica platform from version 10.2.2, update connections.

## Update Connections

You might need to update connections based on the version you are upgrading from.

If you did not create connections when you created the cluster configuration, you need to update the connections.

### Configure the Hadoop Connection

To use properties that you customized in the `hadoopEnv.properties` file, you must configure the Hadoop connection properties such as cluster environment variables, cluster path variables, and advanced properties.

<b>Perform this task in the following situation:</b>
- You upgraded from version 10.1.1 or earlier.

When you run the Informatica upgrade, the installer backs up the existing `hadoopEnv.properties` file. You can find the backup `hadoopEnv.properties` file in the following location:

```
<Previous Informatica installation directory>/services/shared/hadoop/<Hadoop distribution name>_<version>/infaConf
```

Edit the Hadoop connection in the Administrator tool or the Developer tool to include any properties that you manually configured in the `hadoopEnv.properties` file. The Hadoop connection contains default values for properties such as cluster environment and path variables and advanced properties. You can update the default values to match the properties in the `hadoopEnv.properties` file.

### Replace the Connections with New Connections

If you created connections when you imported the cluster configuration, you need to replace connections in mappings with the new connections.

<b>Perform this task in the following situation:</b>
- You upgraded from version 10.1.1 or earlier.

The method that you use to replace connections in mappings depends on the type of connection.

#### **Hadoop connection**

Run the following commands to replace the connections:

- `infacmd dis replaceMappingHadoopRuntimeConnections`. Replaces connections associated with mappings that are deployed in applications.
- `infacmd mrs replaceMappingHadoopRuntimeConnections`. Replaces connections associated with mappings that you run from the Developer tool.

For information about the `infacmd` commands, see the *Informatica Command Reference*.



## Hive, HDFS, and HBase connections

You must replace the connections manually.

## Complete Connection Upgrade

If you did not create connections when you imported the cluster configuration, you must update connection properties for Hadoop, Hive, HDFS, and HBase connections.

### Perform this task in the following situations:

- You upgraded from version 10.2.2 or earlier.

Perform the following tasks to update the connections:

### Update changed properties

Review connections that you created in a previous release to update the values for connection properties. For example, if you added nodes to the cluster or if you updated the distribution version, you might need to verify host names, URIs, or port numbers for some of the properties.

### Associate the cluster configuration

The Hadoop, Hive, HDFS, and HBase connections must be associated with a cluster configuration. Complete the following tasks:

1. Run `infacmd isp listConnections` to identify the connections that you need to upgrade. Use `-ct` to list connections of a particular type.
2. Run `infacmd isp UpdateConnection` to associate the cluster configuration with the connection. Use `-cn` to name the connection and `-o clusterConfigID` to associate the cluster configuration with the connection.

For more information about `infacmd`, see the *Informatica Command Reference*.

## Replace Hive Run-time Connections with Hadoop Connections

Effective in version 10.2.2, Big Data Management dropped support for the Hive engine and Hive run-time connections. If you used Hive connections to run mappings on the Hadoop cluster, you must generate Hadoop connections from the Hive connections.

### Perform this task in the following situations:

- You upgraded from version 10.1.1 or earlier.
- The Hive connections are configured to run mappings in the Hadoop environment.

The upgrade process generates a connection name for the Hadoop connection and replaces the connection name in the mappings. It does not create the physical connection object. When the upgrade is complete, you must run a command to generate the connection. You generate Hadoop connections from Hive connections that are configured to run mappings in the Hadoop environment.

You must generate Hadoop connections from Hive connections that are configured to run mappings in the Hadoop environment.

1. Run `infacmd isp generateHadoopConnectionFromHiveConnection` to generate a Hadoop connection from a Hive connection that is configured to run in the Hadoop environment.  
The command names the connection as follows: "Autogen\_<Hive connection name>." If the connection name exceeds the 128 character limit, the command fails.
2. If the command fails, perform the following tasks:
  - a. Rename the connection to meet the character limit and run the command again.
  - b. Run `infacmd dis replaceMappingHadoopRuntimeConnections` to replace connections associated with mappings that are deployed in applications.
  - c. Run `infacmd mrs replaceMappingHadoopRuntimeConnections` to replace connections associated with mappings that you run from the Developer tool.
3. If the Hive connection was parameterized, you must update the connection names in the parameter file. Verify that the Hive sources, Hive targets, and the Hive engine parameters are updated with the correct connection name.
4. If any properties changed in the cluster, such as host names, URIs, or port numbers, you must update the properties in the connections.

For information about the `infacmd` commands, see the *Informatica® Command Reference*.

## Update Streaming Objects

Data Engineering Streaming uses Spark Structured Streaming to process data instead of Spark Streaming. To support Spark Structured Streaming, some header ports are added to the data objects, and support to some of the data objects and transformations are deferred to a future release. The behavior of some of the data objects is also updated.

After you upgrade, the existing streaming mappings become invalid because of the unavailable header ports, the unsupported transformations or data objects, and the behavior change of some data objects.

<b>Perform this task in the following situations:</b>
- You upgraded from version 10.1.1, 10.2.0, or 10.2.1.



To use an existing streaming mapping, perform the following tasks:

- Re-create the physical data objects. After you re-create the physical data objects, the data objects get the required header ports, such as timestamp, partitionID, or key based on the data object.
- In a Normalizer transformation, if the **Occurs** column is set to Auto, re-create the Normalizer transformation. You must re-create the Normalizer transformation because the type configuration property of the complex port refers to the physical data object that you plan to replace.
- Update the streaming mapping. If the mapping contains Kafka target, Aggregator transformation, Joiner transformation, or Normalizer transformation, replace the data object or transformation, and then update the mapping because of the changed behavior of these transformations and data objects.
- Verify the deferred data object types. If the streaming mapping contains unsupported transformations or data objects, contact Informatica Global Customer Support.

## Re-create the Physical Data Objects

When you re-create the physical data objects, the physical data objects get the header ports and some properties are not available for some data objects. Update the existing mapping with the newly created physical data objects.

1. Go to the existing mapping, select the data object from the mapping.
2. Click the **Properties** tab. On the **Column Projection** tab, click **Edit Schema**.
3. Note the schema information from the **Edit Schema** dialog box.
4. Note the parameters information from the **Parameters** tab.
5. Create new physical data objects.

After you re-create the data objects, the physical data objects get the required header ports. The Microsoft Azure does not support the following properties and are not available for Azure Event Hubs data objects:

- Consumer Properties
- Partition Count

## Re-create the Normalizer Transformation

If the mapping contains a Normalizer transformation with the **Occurs** column set to Auto, re-create the Normalizer transformation. When you re-create the Normalizer transformation, the type configuration property of the complex port refers to the re-created physical data object.

## Update the Streaming Mappings

After you re-create the data object, replace the existing data objects with the re-created data objects. If the mapping contains Normaliser Transformation, Aggregator transformation, or Joiner transformation, update the mapping because of the changed behavior of these transformations and data objects.

### Transformation Updates

If a transformation uses a complex port, configure the type configuration property of the port because the property refers to the physical data object that you replaced.

### Aggregator and Joiner Transformation Updates

An Aggregator transformation must be downstream from a Joiner transformation. A Window transformation must be directly upstream from both Aggregator and Joiner transformations. Previously, you could use an Aggregator transformation anywhere in the streaming mapping.

If a mapping contains an Aggregator transformation upstream from a Joiner transformation, move the Aggregator transformation downstream from a Joiner transformation. Add a Window transformation directly upstream from both Aggregator and Joiner transformations.

## Verify the Deferred Data Object Types

After you upgrade, the streaming mappings might contain some transformations and data objects that are deferred.

The following table lists the data object types to which the support is deferred to a future release:

Object Type	Object
Transformation	Data Masking

If you want to continue using the mappings that contain deferred data objects or transformations, you must contact Informatica Global Customer Support.

## CHAPTER 4

# Azure HDInsight Integration Tasks

This chapter includes the following topics:

- [Azure HDInsight Task Flows, 69](#)
- [Prepare for Cluster Import from Azure HDInsight, 74](#)
- [Create a Cluster Configuration, 82](#)
- [Verify or Refresh the Cluster Configuration , 84](#)
- [Configure the Hive Warehouse Connector and Hive LLAP, 85](#)
- [Verify JDBC Drivers for Sqoop Connectivity, 86](#)
- [Configure the Developer Tool, 87](#)
- [Complete Upgrade Tasks, 88](#)

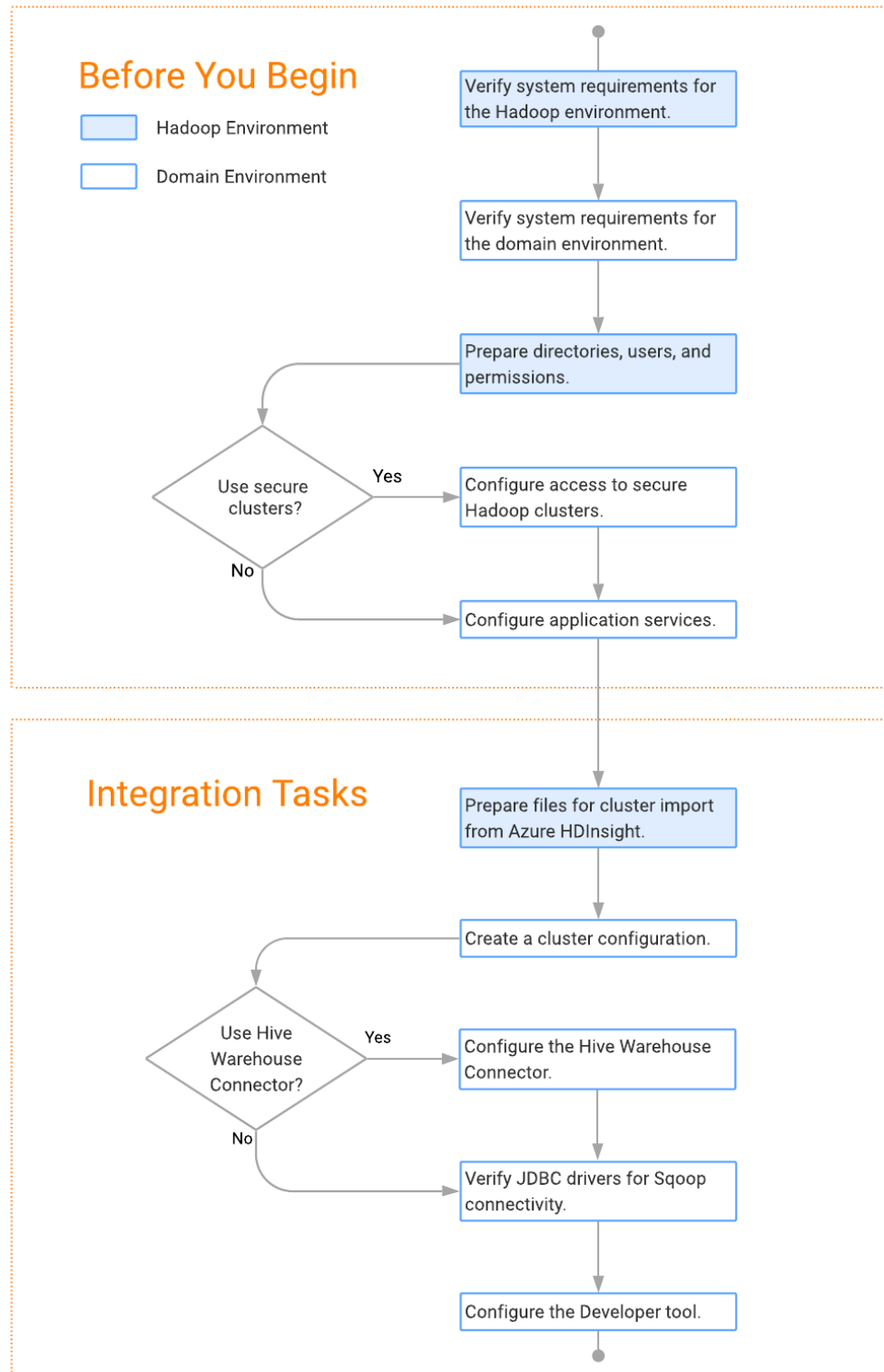
## Azure HDInsight Task Flows

Depending on whether you want to integrate or upgrade Data Engineering Integration in an Azure HDInsight environment, you can use the flow charts to perform the following tasks:

- Integrate the Informatica domain with Azure HDInsight for the first time.
- Upgrade from version 10.2.1.
- Upgrade from version 10.2.
- Upgrade from a version earlier than 10.2.

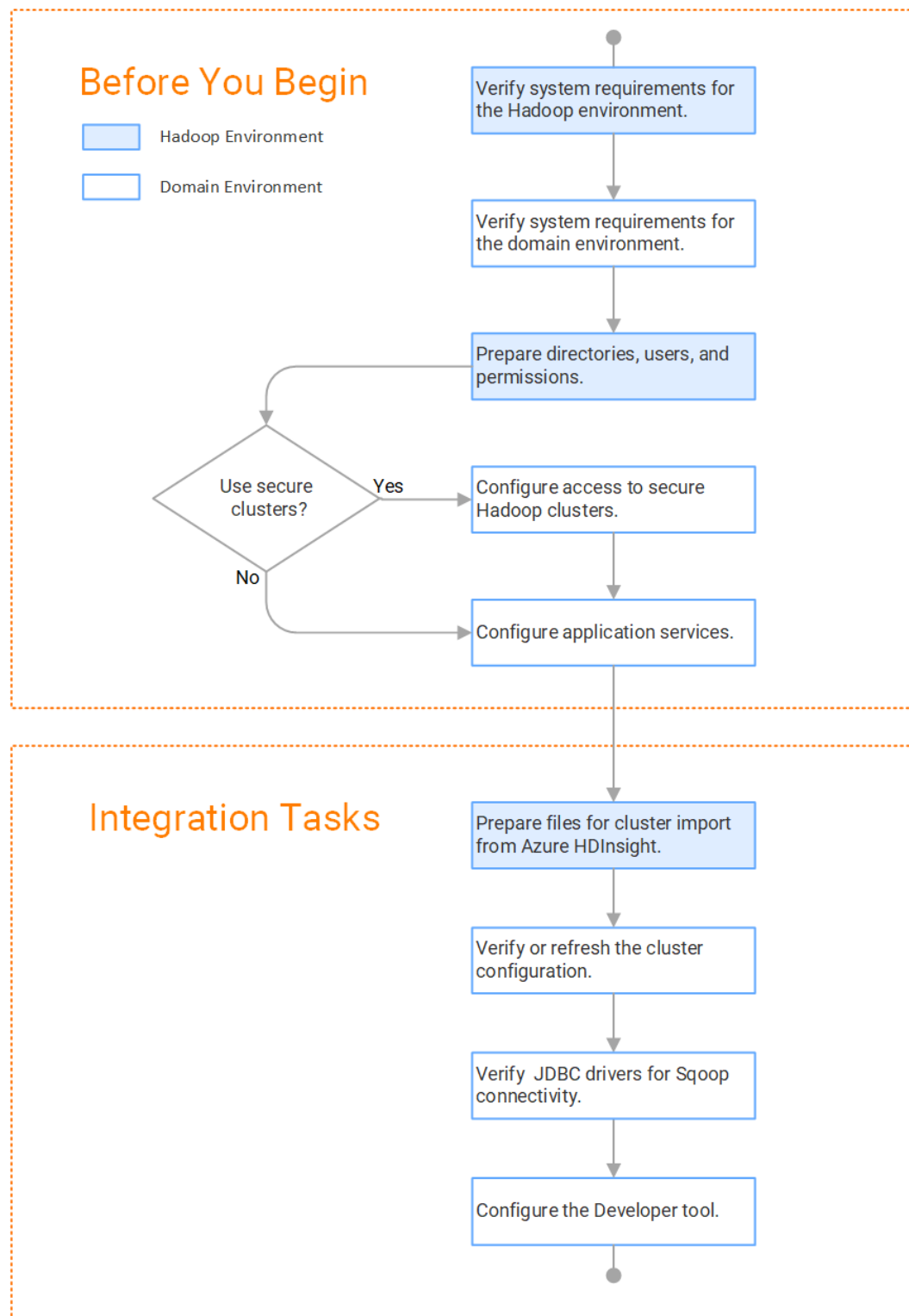
## Task Flow to Integrate with Azure HDInsight

The following diagram shows the task flow to integrate the Informatica domain with Azure HDInsight:



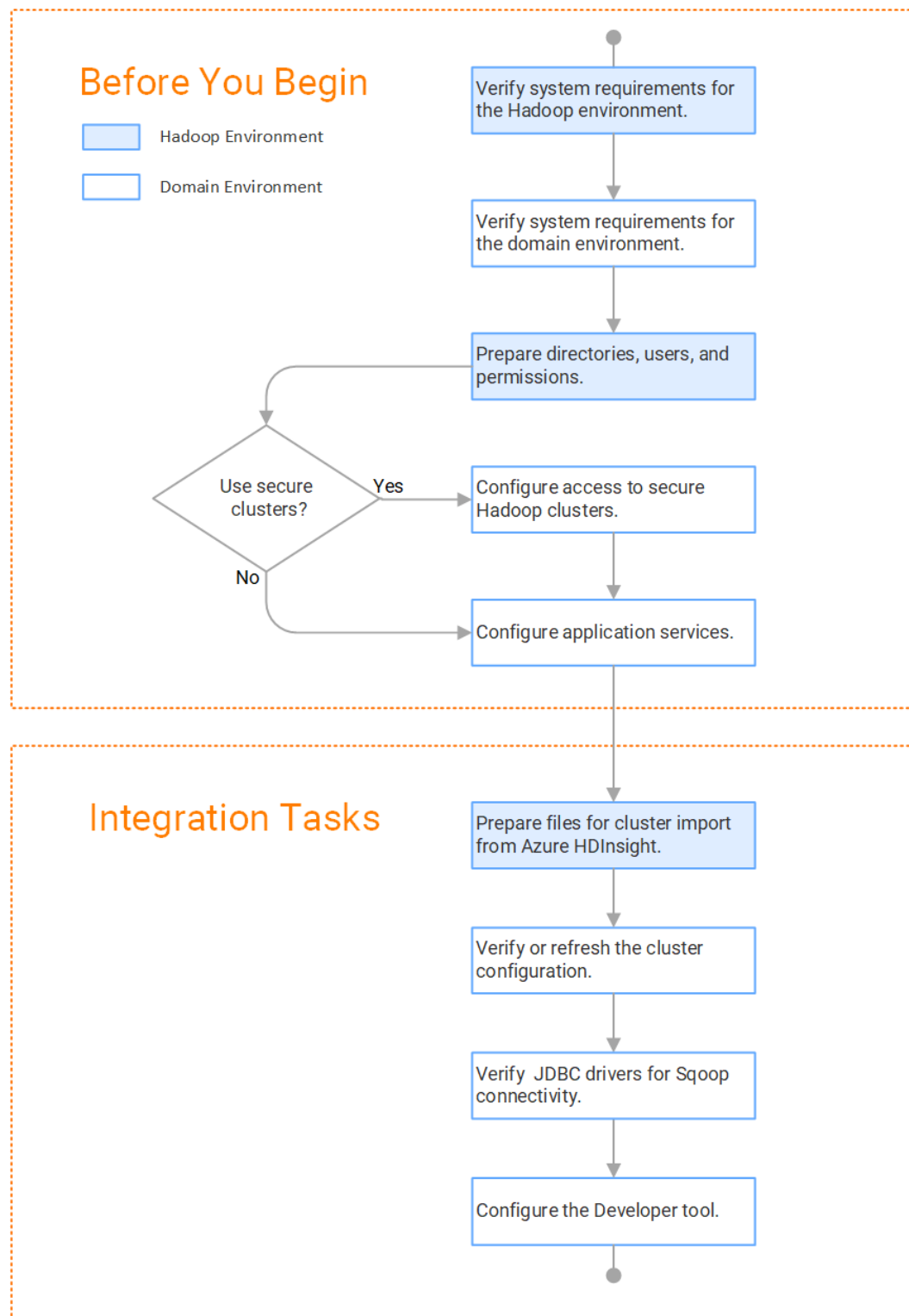
## Task Flow to Upgrade from Version 10.2.1

The following diagram shows the task flow to upgrade version 10.2.1 for Azure HDI:



## Task Flow to Upgrade from Version 10.2

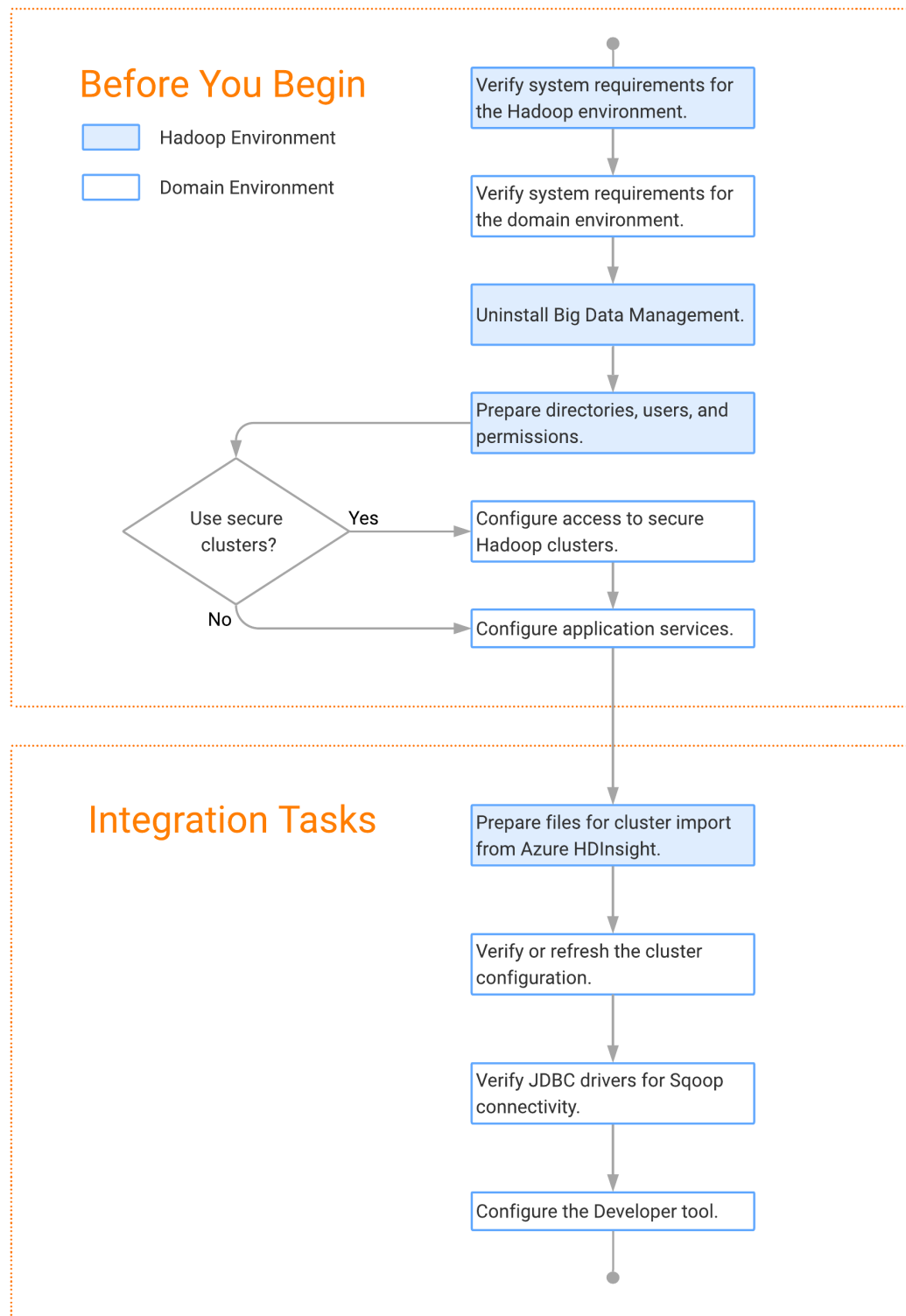
The following diagram shows the task flow to upgrade version 10.2 for Azure HDInsight:





## Task Flow to Upgrade from a Version Earlier than 10.2

The following diagram shows the task flow to upgrade from a version earlier than 10.2 for Azure HDInsight:



# Prepare for Cluster Import from Azure HDInsight

Before the Informatica administrator can import cluster information to create a cluster configuration in the Informatica domain, the Hadoop administrator must perform some preliminary tasks.

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

**Note:** If you are upgrading from a previous version, verify the properties and suggested values, as Data Engineering Integration might require additional properties or different values for existing properties.

Complete the following tasks to prepare the cluster before the Informatica administrator creates the cluster configuration:

1. When the Informatica domain is on-premises, verify that the VPN is enabled between the domain and the Azure HDInsight cloud network.
2. When the Informatica domain is deployed on the Azure cloud, verify the following requirements:
  - Verify that the domain can access the private or internal IP addresses of all HDInsight cluster nodes and can connect to the required ports. For a list of ports, see the HDInsight ports listed in the article [Configuring Ports for Big Data Products](#).
  - When the domain and the HDInsight cluster reside in different virtual networks, known as "Vnets," see the [Azure documentation](#) to enable peering between virtual networks.
3. Verify property values in \*-site.xml files that Data Engineering Integration needs to run mappings in the Hadoop environment.
4. Provide information to the Informatica administrator that is required to import cluster information into the domain. Depending on the method of import, perform one of the following tasks:
  - To import directly from the cluster, give the Informatica administrator cluster authentication information to connect to the cluster.
  - To import from an archive file, export cluster information and provide an archive file to the Informatica administrator.

## Configure \*-site.xml Files for Azure HDInsight

The Hadoop administrator must configure \*-site.xml file properties on the cluster and restart the credential service and other impacted services before the Informatica administrator imports cluster information into the domain.

### capacity-scheduler.xml

Configure the following properties in the capacity-scheduler.xml file:

**yarn.scheduler.capacity.<queue path>.disable\_preemption**

Disables preemption for the Capacity Scheduler. The Blaze engine does not support preemption. If YARN preemption is enabled for the cluster, you need to disable it for the queue allocated to the Blaze engine.

Set to TRUE for queue allocated to the Blaze engine.

### core-site.xml

Configure the following properties in the core-site.xml file:

#### **fs.azure.account.key.<your account>.blob.core.windows.net**

Required for Azure HDInsight cluster that uses WASB storage. The storage account access key required to access the storage.

You can contact the HDInsight cluster administrator to get the storage account key associated with the HDInsight cluster. If you are unable to contact the administrator, perform the following steps to decrypt the encrypted storage account key:

- Copy the value of the `fs.azure.account.key.<your account>.blob.core.windows.net` property.

```
<property>
<name>fs.azure.account.key.<youraccount>.blob.core.windows.net</name>
<value>STORAGE ACCOUNT KEY</value>
</property>
```

- Decrypt the storage account key. Run the `decrypt.sh` specified in the `fs.azure.shellkeyprovider.script` property along with the encrypted value you copied in the previous step.

```
<property>

<name>fs.azure.shellkeyprovider.script</name>
<value>/usr/lib/hdinsight-common/scripts/decrypt.sh</value>
</property>
```

- Copy the decrypted value and update the value of `fs.azure.account.key.youraccount.blob.core.windows.net` property in the cluster configuration `core-site.xml`.

#### **dfs.adls.oauth2.client.id**

Required for Azure HDInsight cluster that uses ADLS Gen1 storage without Enterprise Security Package. The application ID associated with the Service Principal required to authorize the service principal and access the storage.

To find the application ID for a service principal, in the Azure Portal, click **Azure Active Directory > App registrations > Service Principal Display Name**.

#### **dfs.adls.oauth2.refresh.url**

Required for Azure HDInsight cluster that uses ADLS Gen1 storage without Enterprise Security Package. The OAuth 2.0 token endpoint required to authorize the service principal and access the storage.

To find the refresh URL OAuth 2.0 endpoint, in the Azure portal, click **Azure Active Directory > App registrations > Endpoints**.

#### **dfs.adls.oauth2.credential**

Required for Azure HDInsight cluster that uses ADLS Gen1 storage without Enterprise Security Package. The password required to authorize the service principal and access the storage.

To find the password for a service principal, in the Azure portal, click **Azure Active Directory > App registrations > Service Principal Display Name > Settings > Keys**.

#### **fs.azure.account.key.<your account>.dfs.core.windows.net**

Required for Azure HDInsight cluster that uses ADLS Gen2 storage without Enterprise Security Package. The decrypted account key for the storage account.

You can contact the HDInsight cluster administrator to get the storage account key associated with the HDInsight cluster.

#### **hadoop.proxyuser.<proxy user>.groups**

Defines the groups that the proxy user account can impersonate. On a secure cluster the `<proxy user>` is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the `<proxy user>` is the system user that runs the Informatica daemon.

Set to group names of impersonation users separated by commas. If less security is preferred, use the wildcard " \* " to allow impersonation from any group.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.proxyuser.<proxy user>.users**

Required for all cluster types. Defines the user account that the proxy user account can impersonate. On a secure cluster, the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to a single user account or set to a comma-separated list. If less security is preferred, use the wildcard " \* " to allow impersonation from any group.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.proxyuser.<proxy user>.hosts**

Defines the host machines that a user account can impersonate. On a secure cluster the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to a single host name or IP address, or set to a comma-separated list. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.proxyuser.hive.hosts**

Comma-separated list of hosts that you want to allow the Hive user to impersonate on a non-secure cluster.

When `hive.server2.enable.doAs` is false, append a comma-separated list of Informatica server host names or IP address where the Data Integration Service is running. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to this property, restart the cluster services that use core-site configuration values.

#### **hadoop.proxyuser.yarn.groups**

Comma-separated list of groups that you want to allow the YARN user to impersonate on a non-secure cluster.

Set to group names of impersonation users separated by commas. If less security is preferred, use the wildcard " \* " to allow impersonation from any group.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.proxyuser.yarn.hosts**

Comma-separated list of hosts that you want to allow the YARN user to impersonate on a non-secure cluster.

Set to a single host name or IP address, or set to a comma-separated list. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

**io.compression.codecs**

Enables compression on temporary staging tables.

Set to a comma-separated list of compression codec classes on the cluster.

**hadoop.security.auth\_to\_local**

Translates the principal names from the Active Directory and MIT realm into local names within the Hadoop cluster. Based on the Hadoop cluster used, you can set multiple rules.

Set to: RULE:[1:\$1@\$0](^.\*@YOUR.REALM)s/^.\*(.)@YOUR.REALM\.COM\$/ \$1/g

Set to: RULE:[2:\$1@\$0](^.\*@YOUR.REALM\.)s/^.\*(.)@YOUR.REALM\.COM\$/ \$1/g

**hbase-site.xml**

Configure the following properties in the hbase-site.xml file:

**hbase.use.dynamic.jars**

Enables metadata import and test connection from the Developer tool. Required for an HDInsight cluster that uses ADLS storage or an Amazon EMR cluster that uses HBase resources in S3 storage.

Set to: false

**zookeeper.znode.parent**

Identifies HBase master and region servers.

Set to the relative path to the znode directory of HBase.

**hive-site.xml**

Configure the following properties in the hive-site.xml file:

**hive.cluster.delegation.token.store.class**

The token store implementation. Required for HiveServer2 high availability and load balancing.

Set to: org.apache.hadoop.hive.thrift.DBTokenStore

**hive.compactor.initiator.on**

Runs the initiator and cleaner threads on metastore instance. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: TRUE

**hive.compactor.worker.threads**

The number of worker threads to run in a metastore instance. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: 1

**hive.enforce.bucketing**

Enables dynamic bucketing while loading to Hive. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: TRUE

**hive.exec.dynamic.partition**

Enables dynamic partitioned tables for Hive tables. Applicable for Hive versions 0.9 and earlier.

Set to: TRUE

**hive.exec.dynamic.partition.mode**

Allows all partitions to be dynamic. Required for the Update Strategy transformation in a mapping that writes to a Hive target. Also required if you use Sqoop and define a DDL query to create or replace a partitioned Hive target at run time.

Set to: nonstrict

**hive.support.concurrency**

Enables table locking in Hive. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: TRUE

**hive.server2.support.dynamic.service.discovery**

Enables HiveServer2 dynamic service discovery. Required for HiveServer2 high availability.

Set to: TRUE

**hive.server2.zookeeper.namespace**

The value of the ZooKeeper namespace in the JDBC connection string. Required for HiveServer2 high availability.

Set to: jdbc:hive2://<zookeeper\_ensemble>/  
default;serviceDiscoveryMode=zooKeeper;zooKeeperNamespace=hiveserver2

**hive.txn.manager**

Turns on transaction support. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: org.apache.hadoop.hive.ql.lockmgr.DbTxnManager

**hive.warehouse.subdir.inherit.perms**

Enables permissions inheritance in Hive warehouse subdirectories. When set to TRUE, subdirectories inherit permissions, groups, and access control lists (ACLs).

Set to FALSE to enable mappings that access Hive tables to run on an ESP-enabled HDInsight 4.x using WASBS.

**hive.zookeeper.quorum**

Comma-separated list of ZooKeeper server host:ports in a cluster. The value of the ZooKeeper ensemble in the JDBC connection string. Required for HiveServer2 high availability.

Set to: jdbc:hive2://<zookeeper\_ensemble>/default;serviceDiscoveryMode=zooKeeper;

**[mapred-site.xml](#)**

Configure the following properties in the mapred-site.xml file:

**mapreduce.framework.name**

The run-time framework to run MapReduce jobs. Values can be local, classic, or yarn. Required for Sqoop.

Set to: yarn

**yarn.app.mapreduce.am.staging-dir**

The HDFS staging directory used while submitting jobs.

Set to the staging directory path.

## yarn-site.xml

Configure the following properties in the yarn-site.xml file:

### **yarn.application.classpath**

Required for dynamic resource allocation.

Add spark\_shuffle.jar to the class path. The .jar file must contain the class "org.apache.spark.network.yarn.YarnShuffleService."

### **yarn.nodemanager.resource.memory-mb**

The maximum RAM available for each container. Set the maximum memory on the cluster to increase resource memory available to the Blaze engine.

Set to 16 GB if value is less than 16 GB.

### **yarn.nodemanager.resource.cpu-vcores**

The number of virtual cores for each container. Required for Blaze engine resource allocation.

Set to 10 if the value is less than 10.

### **yarn.scheduler.minimum-allocation-mb**

The minimum RAM available for each container. Required for Blaze engine resource allocation.

Set to 6 GB if the value is less than 6 GB.

### **yarn.nodemanager.vmem-check-enabled**

Disables virtual memory limits for containers. Required for the Blaze and Spark engines.

Set to: false

### **yarn.nodemanager.aux-services**

Required for dynamic resource allocation for the Spark engine.

Add an entry for "spark\_shuffle."

### **yarn.nodemanager.aux-services.spark\_shuffle.class**

Required for dynamic resource allocation for the Spark engine.

Set to: org.apache.spark.network.yarn.YarnShuffleService

### **yarn.resourcemanager.scheduler.class**

Defines the YARN scheduler that the Data Integration Service uses to assign resources.

Set to: org.apache.hadoop.yarn.server.resourcemanager.scheduler

### **yarn.node-labels.enabled**

Enables node labeling.

Set to: TRUE

### **yarn.node-labels.fs-store.root-dir**

The HDFS location to update node label dynamically.

Set to: <hdfs://[Node name]:[Port]/[Path to store]/[Node labels]/>

## tez-site.xml

Configure the following properties in the tez-site.xml file:

### **tez.runtime.io.sort.mb**

The sort buffer memory. Required when the output needs to be sorted for Blaze and Spark engines.

Set value to 270 MB.

## Verify HDInsight Cluster Security Settings

Verify cluster security settings depending on whether or not the cluster is configured with the Enterprise Security Package (ESP).

### Clusters with ESP Enabled

When the HDInsight cluster is configured with the Enterprise Security Package (ESP), verify security settings to access Azure Data Lake Storage (ADLS) Gen2 resources.

The Enterprise Security Package uses Kerberos authentication and Apache Ranger authorization to enable Active Directory (AD) based authentication, multi-user support, and role-based access control.

To verify that you have the required environment, see the Azure documentation.

- [Configure a HDInsight cluster with Enterprise Security Package by using Azure Active Directory Domain Services](#)
- [Use Azure Data Lake Storage Gen2 with Azure HDInsight Clusters](#)

In addition, set the following properties on the cluster when you integrate an HDInsight 4.x cluster with WASBS storage:

#### **fs.permissions.umask-mode**

The umask used to set default permissions on created files and directories.

Set to 000.

#### **Run as end user instead of Hive user (doAs)**

Enables the cluster to run jobs as the impersonation user and not the Hive user.

Set to FALSE.

### Clusters without ESP Enabled

When the HDInsight 4.x cluster does not have the Enterprise Security Package enabled, verify the following properties on the cluster when it uses WASBS storage:

#### **Hive Authorization Manager**

Authorization provider for the cluster.

Set to

`org.apache.hadoop.hive.q1.security.authorization.MetaStoreAuthzAPIAuthorizerEmbedOnly`.

For more information, see the [Azure HDInsight documentation](#).

#### **Run as end user instead of Hive user (doAs)**

Enables the cluster to run jobs as the impersonation user and not the Hive user.

Set to TRUE.



## Prepare for Direct Import from Azure HDInsight

If you plan to provide direct access to the Informatica administrator to import cluster information, provide the required connection information.

The following table describes the information that you need to provide to the Informatica administrator to create the cluster configuration directly from the cluster:

Property	Description
Host	Host name or IP address of the cluster manager.
Port	Port of the cluster manager.
User ID	Cluster user ID.
Password	Password for the user.
Cluster name	Name of the cluster. Use the display name if the cluster manager manages multiple clusters. If you do not provide a cluster name, the wizard imports information based on the default cluster.

## Prepare the Archive File for Import from Azure HDInsight

When you prepare the archive file for cluster configuration import from HDInsight, include all required \*-site.xml files and edit the file manually after you create it.

Create a .zip or .tar file that contains the following \*-site.xml files:

- core-site.xml
- hbase-site.xml. Required only to access HBase sources and targets.
- hdfs-site.xml
- hive-site.xml
- mapred-site.xml or tez-site.xml. Include the mapred-site.xml file or the tez-site.xml file based on the Hive execution type used on the Hadoop cluster.
- yarn-site.xml

### Update the Archive File

After you create the archive file, edit the Hortonworks Data Platform (HDP) version string wherever it appears in the archive file. Search for the string `${hdp.version}` and replace all instances with the HDP version that HDInsight includes in the Hadoop distribution.

For example, the edited `tez.task.launch.cluster-default.cmd-opts` property value looks similar to the following:

```
<property>
<name>tez.task.launch.cluster-default.cmd-opts</name>
<value>-server -Djava.net.preferIPv4Stack=true -Dhdp.version=2.6.0.2-76</value>
</property>
```

# Create a Cluster Configuration

After the Hadoop administrator prepares the cluster for import, the Informatica administrator must create a cluster configuration.

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from version 10.1.1 or earlier.

A cluster configuration is an object in the domain that contains configuration information about the Hadoop cluster. The cluster configuration enables the Data Integration Service to push mapping logic to the Hadoop environment. Import configuration properties from the Hadoop cluster to create a cluster configuration.

The import process imports values from \*-site.xml files into configuration sets based on the individual \*-site.xml files. When you perform the import, the cluster configuration wizard can create Hadoop, HBase, HDFS, and Hive connection to access the Hadoop environment. If you choose to create the connections, the wizard also associates the cluster configuration with the connections.

**Note:** If you are integrating for the first time and you imported the cluster configuration when you ran the installer, you *must* re-create or refresh the cluster configuration.

## Before You Import

Before you can import the cluster configuration, you must get information from the Hadoop administrator based on the method of import.

If you import directly from the cluster, contact the Hadoop administrator to get cluster connection information. If you import from a file, get an archive file of exported cluster information.

## Importing a Hadoop Cluster Configuration from the Cluster

When you import the Hadoop cluster configuration directly from the cluster, you provide information to connect to the cluster.

Get cluster connection information from the Hadoop administrator.

1. From the **Connections** tab, click the **ClusterConfigurations** node in the Domain Navigator.
2. From the Actions menu, select **New > Cluster Configuration**.

The **Cluster Configuration** wizard opens.

3. Configure the following General properties:

Property	Description
Cluster configuration name	Name of the cluster configuration.
Description	Optional description of the cluster configuration.
Distribution type	The cluster Hadoop distribution type.

Property	Description
Method to import the cluster configuration	Choose <b>Import from cluster</b> .
Create connections	<p>Choose to create Hadoop, HDFS, Hive, and HBase connections.</p> <p>If you choose to create connections, the Cluster Configuration wizard associates the cluster configuration with each connection that it creates.</p> <p>The Hadoop connection contains default values for properties such as cluster environment variables, cluster path variables, and advanced properties. Based on the cluster environment and the functionality that you use, you can add to the default values or change the default values of these properties. For a list of Hadoop connection properties to configure, see Hadoop Cluster Connection Properties.</p> <p>If you do not choose to create connections, you must manually create them and associate the cluster configuration with them.</p> <p><b>Important:</b> When the wizard creates the Hive connection, it populates the Metadata Connection String and the Data Access Connection String properties with the value from the hive.metastore.uris property. If the Hive metastore and HiveServer2 are running on different nodes, you must update the Metadata Connection String to point to the HiveServer2 host.</p>

The cluster properties appear.

- Configure the following properties:

Property	Description
Host	Host name or IP address of the cluster manager.
Port	Port of the cluster manager.
User ID	Cluster user ID.
Password	Password for the user.
Cluster name	Name of the cluster. Use the display name if the cluster manager manages multiple clusters. If you do not provide a cluster name, the wizard imports information based on the default cluster.

- Click **Next** and verify the cluster configuration information on the summary page.

## Importing a Hadoop Cluster Configuration from a File

You can import properties from an archive file to create a cluster configuration.

Before you import from the cluster, you must get the archive file from the Hadoop administrator.

- From the **Connections** tab, click the **ClusterConfigurations** node in the Domain Navigator.
- From the Actions menu, select **New > Cluster Configuration**.

The **Cluster Configuration** wizard opens.

3. Configure the following properties:

Property	Description
Cluster configuration name	Name of the cluster configuration.
Description	Optional description of the cluster configuration.
Distribution type	The cluster Hadoop distribution type.
Method to import the cluster configuration	Choose <b>Import from file</b> to import properties from an archive file.
Create connections	<p>Choose to create Hadoop, HDFS, Hive, and HBase connections.</p> <p>If you choose to create connections, the Cluster Configuration wizard associates the cluster configuration with each connection that it creates.</p> <p>The Hadoop connection contains default values for properties such as cluster environment variables, cluster path variables, and advanced properties. Based on the cluster environment and the functionality that you use, you can add to the default values or change the default values of these properties. For a list of Hadoop connection properties to configure, see <a href="#">“Configuring Hadoop Connection Properties” on page 269</a>.</p> <p>If you do not choose to create connections, you must manually create them and associate the cluster configuration with them.</p> <p><b>Important:</b> When the wizard creates the Hive connection, it populates the Metadata Connection String and the Data Access Connection String properties with the value from the hive.metastore.uris property. If the Hive metastore and HiveServer2 are running on different nodes, you must update the Metadata Connection String to point to the HiveServer2 host.</p>

4. Click **Browse** to select a file. Select the file and click **Open**.
5. Click **Next** and verify the cluster configuration information on the summary page.

## Verify or Refresh the Cluster Configuration

You might need to refresh the cluster configuration or update the distribution version in the cluster configuration when you upgrade.

### Perform this task in the following situation:

- You upgraded from version 10.2 or later.

### Verify the Cluster Configuration

The cluster configuration contains a property for the distribution version. The verification task depends on the version you upgraded:

#### Upgrade from 10.2

If you upgraded from 10.2 and you changed the distribution version, you need to verify the distribution version in the General properties of the cluster configuration.

### Upgrade from 10.2.1

Effective in version 10.2.1, Informatica assigns a default version to each Hadoop distribution type. If you configure the cluster configuration to use the default version, the upgrade process upgrades to the assigned default version if the version changes. If you have not upgraded your Hadoop distribution to Informatica's default version, you need to update the distribution version property.

For example, suppose the assigned default Hadoop distribution version for 10.2.1 is  $n$ , and for 10.2.2 is  $n+1$ . If the cluster configuration uses the default supported Hadoop version of  $n$ , the upgraded cluster configuration uses the default version of  $n+1$ . If you have not upgraded the distribution in the Hadoop environment you need to change the cluster configuration to use version  $n$ .

If you configure the cluster configuration to use a distribution version that is not the default version, you need to update the distribution version property in the following circumstances:

- Informatica dropped support for the distribution version.
- You changed the distribution version.

### Refresh the Cluster Configuration

If you updated any of the \*-site.xml files noted in the topic to prepare for cluster import, you need to refresh the cluster configuration in the Administrator tool.

## Configure the Hive Warehouse Connector and Hive LLAP

Optionally, you can configure the Hive Warehouse Connector and Hive LLAP to improve performance when you read from and write to Hive targets.

#### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

The Hive Warehouse Connector reads from and writes to Hive tables without using temporary staging tables that require additional storage overhead. Use the Hive Warehouse Connector on the Spark engine to allow Spark code to interact with Hive targets and to use ACID-enabled Hive tables. When you enable the Hive Warehouse Connector, mappings use Hive LLAP to run Hive queries rather than HiveServer2.

Before you enable the Hive Warehouse Connector, enable Hive LLAP on the Hadoop cluster. To enable the connector, configure the following properties in the Spark advanced properties for the Hadoop connection:

#### **infaspark.useHiveWarehouseAPI**

Enables the Hive Warehouse Connector. Set to TRUE.

For example, `infaspark.useHiveWarehouseAPI=true`.

#### **spark.datasource.hive.warehouse.load.staging.dir**

Directory for the temporary HDFS files used for batch writes to Hive. Required when you enable the Hive Warehouse Connector.

For example, set to `/tmp`

**spark.datasource.hive.warehouse.metastoreUri**

URI for the Hive metastore. Required when you enable the Hive Warehouse Connector. Use the value for *hive.metastore.uris* from the *hive\_site.xml* cluster configuration properties.

For example, set the value to `thrift://mycluster-1.com:9083` .

**spark.hadoop.hive.llap.daemon.service.hosts**

Application name for the LLAP service. Required when you enable the Hive Warehouse Connector. Use the value for *hive.llap.daemon.service.hosts* from the *hive\_site.xml* cluster configuration properties.

**spark.hadoop.hive.zookeeper.quorum**

Zookeeper hosts used by Hive LLAP. Required when you enable the Hive Warehouse Connector. Use the value for *hive.zookeeper.quorum* from the *hive\_site.xml* cluster configuration properties.

**spark.sql.hive.hiveserver2.jdbc.url**

URL for HiveServer2 Interactive. Required to use the Hive Warehouse Connector. Use the value in Ambari for HiveServer2 JDBC URL.

## Verify JDBC Drivers for Sqoop Connectivity

Verify that you have the JDBC drivers to access JDBC-compliant databases in the Hadoop environment. You might need separate drivers for metadata import and for run-time processing.

**Perform this task in the following situations:**

- You are integrating for the first time.
- You upgraded from version 10.2.1 or earlier.

You download drivers based on design-time and run-time requirements:

- **Design-time.** To import metadata, you can use the DataDirect drivers packaged with the Informatica installer if they are available. If they are not available, use any Type 4 JDBC driver that the database vendor recommends.
- **Run-time.** To run mappings, use any Type 4 JDBC driver that the database vendor recommends. Some distributions support other drivers to use Sqoop connectors. You cannot use the DataDirect drivers for run-time processing.

## Verify Design-time Drivers

Use the DataDirect JDBC drivers packaged with the Informatica installer to import metadata from JDBC-compliant databases. If the DataDirect JDBC drivers are not available for a specific JDBC-compliant database, download the Type 4 JDBC driver associated with that database.

Copy the JDBC driver .jar files to the following location on the Developer tool machine:

```
<Informatica installation directory>\clients\externaljdbcjars
```

## Verify Run-time Drivers

Verify run-time drivers for mappings that access JDBC-compliant databases in the Hadoop environment. Use any Type 4 JDBC driver that the database vendor recommends.

1. Download Type 4 JDBC drivers associated with the JCBC-compliant databases that you want to access.
2. To optimize the Sqoop mapping performance on the Spark engine while writing data to an HDFS complex file target of the Parquet format, download the following .jar files:

- `parquet-hadoop-bundle-1.6.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-hadoop-bundle/1.6.0/>
- `parquet-avro-1.6.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-avro/1.6.0/>
- `parquet-column-1.5.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-column/1.5.0/>

3. Copy all of the .jar files to the following directory on the machine where the Data Integration Service runs:

```
<Informatica installation directory>\externaljdbcjars
```

Changes take effect after you recycle the Data Integration Service. At run time, the Data Integration Service copies the .jar files to the Hadoop distribution cache so that the .jar files are accessible to all nodes in the cluster.

## Configure the Developer Tool

You can configure the Developer tool to enable you to import complex files or import metadata when the domain is Kerberos-enabled.

Edit the `developerCore.ini` file to import complex files. Edit the file on each Developer tool machine.

### Configure developerCore.ini

Edit the `developerCore.ini` file to import complex files.

Edit the `developerCore.ini` file on each machine that hosts the Developer tool.

When you import a complex file, such as Avro or Parquet, the imported object includes metadata associated with the distribution in the Hadoop environment. If the file resides on the Developer tool machine, the import process picks up the distribution information from the `developerCore.ini` file. You must edit the `developerCore.ini` file to point to the distribution directory on the Developer tool machine.

You can find the `developerCore.ini` file in the following directory:

```
<Informatica installation directory>\clients\DeveloperClient
```

Add the following property:

```
-DINFA_HADOOP_DIST_DIR=hadoop\<distribution>_<version>
```

The change takes effect when you restart the Developer tool.

## Configure the Developer Tool for Kerberos

To import metadata from Hive, HBase, and complex file sources, import configuration files from the Kerberos-enabled cluster, and generate the Kerberos credentials file on the Developer tool machine.

### Import configuration files

The Hadoop cluster uses a set of XML files named `*-site.xml` to store configuration settings. The domain uses the same set of files to create the cluster configuration object.

To enable you to import metadata from the cluster, import the `*-site.xml` files to each Developer tool machine:

1. Log in to the Administrator tool and navigate to **Connections > Cluster Configuration > CCO**. Locate the cluster configuration associated with the Hadoop cluster.
2. Extract the `*-site.xml` files in the cluster configuration, including sensitive properties, to the following directory on the Developer tool machine: `<Informatica installation directory>\clients\DeveloperClient\hadoop\<Hadoop distribution>\conf`  
For more information about sensitive properties, see the "Active Properties View" topic in the *Data Engineering Administrator Guide*.

**Note:** If you refresh the cluster configuration, repeat these steps.

### Generate the Kerberos credentials file

1. Copy the `krb5.conf` file from `<Developer tool installation directory>/services/shared/security` to `C:/Windows`.
2. Rename `krb5.conf` to `krb5.ini`.
3. In the `krb5.ini` file, verify the value of the `forwardable` option to determine how to use the `kinit` command. If `forwardable=true`, run the command with the `-f` option. Otherwise, run the command without the `-f` option.
4. To generate the Kerberos credentials file, run the `kinit` command from the following location: `<Developer tool installation directory>/clients/java/bin/kinit.exe`  
For example, you might run the following command: `kinit joe/domain12345@MY-REALM`

## Complete Upgrade Tasks

If you upgraded the Informatica platform from version 10.2.2, update connections.

### Update Connections

You might need to update connections based on the version you are upgrading from.

If you did not create connections when you created the cluster configuration, you need to update the connections.



## Configure the Hadoop Connection

To use properties that you customized in the `hadoopEnv.properties` file, you must configure the Hadoop connection properties such as cluster environment variables, cluster path variables, and advanced properties.

### Perform this task in the following situation:

- You upgraded from version 10.1.1 or earlier.

When you run the Informatica upgrade, the installer backs up the existing `hadoopEnv.properties` file. You can find the backup `hadoopEnv.properties` file in the following location:

```
<Previous Informatica installation directory>/services/shared/hadoop/<Hadoop distribution name>_<version>/infaConf
```

Edit the Hadoop connection in the Administrator tool or the Developer tool to include any properties that you manually configured in the `hadoopEnv.properties` file. The Hadoop connection contains default values for properties such as cluster environment and path variables and advanced properties. You can update the default values to match the properties in the `hadoopEnv.properties` file.

## Replace the Connections with New Connections

If you created connections when you imported the cluster configuration, you need to replace connections in mappings with the new connections.

### Perform this task in the following situation:

- You upgraded from version 10.1.1 or earlier.

The method that you use to replace connections in mappings depends on the type of connection.

#### Hadoop connection

Run the following commands to replace the connections:

- `infacmd dis replaceMappingHadoopRuntimeConnections`. Replaces connections associated with mappings that are deployed in applications.
- `infacmd mrs replaceMappingHadoopRuntimeConnections`. Replaces connections associated with mappings that you run from the Developer tool.

For information about the `infacmd` commands, see the *Informatica Command Reference*.

#### Hive, HDFS, and HBase connections

You must replace the connections manually.

## Complete Connection Upgrade

If you *did not* create connections when you imported the cluster configuration, you must update connection properties for Hadoop, Hive, HDFS, and HBase connections.

### Perform this task in the following situations:

- You upgraded from version 10.2.2 or earlier.

Perform the following tasks to update the connections:

### Update changed properties

Review connections that you created in a previous release to update the values for connection properties. For example, if you added nodes to the cluster or if you updated the distribution version, you might need to verify host names, URIs, or port numbers for some of the properties.

### Associate the cluster configuration

The Hadoop, Hive, HDFS, and HBase connections must be associated with a cluster configuration. Complete the following tasks:

1. Run `infacmd isp listConnections` to identify the connections that you need to upgrade. Use `-ct` to list connections of a particular type.
2. Run `infacmd isp UpdateConnection` to associate the cluster configuration with the connection. Use `-cn` to name the connection and `-o clusterConfigID` to associate the cluster configuration with the connection.

For more information about `infacmd`, see the *Informatica Command Reference*.

## Replace Hive Run-time Connections with Hadoop Connections

Effective in version 10.2.2, Big Data Management dropped support for the Hive engine and Hive run-time connections. If you used Hive connections to run mappings on the Hadoop cluster, you must generate Hadoop connections from the Hive connections.

#### Perform this task in the following situations:

- You upgraded from version 10.1.1 or earlier.
- The Hive connections are configured to run mappings in the Hadoop environment.

The upgrade process generates a connection name for the Hadoop connection and replaces the connection name in the mappings. It does not create the physical connection object. When the upgrade is complete, you must run a command to generate the connection. You generate Hadoop connections from Hive connections that are configured to run mappings in the Hadoop environment.

You must generate Hadoop connections from Hive connections that are configured to run mappings in the Hadoop environment.

1. Run `infacmd isp generateHadoopConnectionFromHiveConnection` to generate a Hadoop connection from a Hive connection that is configured to run in the Hadoop environment.  
The command names the connection as follows: "Autogen\_<Hive connection name>." If the connection name exceeds the 128 character limit, the command fails.
2. If the command fails, perform the following tasks:
  - a. Rename the connection to meet the character limit and run the command again.
  - b. Run `infacmd dis replaceMappingHadoopRuntimeConnections` to replace connections associated with mappings that are deployed in applications.
  - c. Run `infacmd mrs replaceMappingHadoopRuntimeConnections` to replace connections associated with mappings that you run from the Developer tool.
3. If the Hive connection was parameterized, you must update the connection names in the parameter file. Verify that the Hive sources, Hive targets, and the Hive engine parameters are updated with the correct connection name.
4. If any properties changed in the cluster, such as host names, URIs, or port numbers, you must update the properties in the connections.

For information about the `infacmd` commands, see the *Informatica® Command Reference*.

## Update Streaming Objects

Data Engineering Streaming uses Spark Structured Streaming to process data instead of Spark Streaming. To support Spark Structured Streaming, some header ports are added to the data objects, and support to some of the data objects and transformations are deferred to a future release. The behavior of some of the data objects is also updated.

After you upgrade, the existing streaming mappings become invalid because of the unavailable header ports, the unsupported transformations or data objects, and the behavior change of some data objects.

<b>Perform this task in the following situations:</b>
<ul style="list-style-type: none"><li>- You upgraded from version 10.1.1, 10.2.0, or 10.2.1.</li></ul>



To use an existing streaming mapping, perform the following tasks:

- Re-create the physical data objects. After you re-create the physical data objects, the data objects get the required header ports, such as timestamp, partitionID, or key based on the data object.
- In a Normalizer transformation, if the **Occurs** column is set to Auto, re-create the Normalizer transformation. You must re-create the Normalizer transformation because the type configuration property of the complex port refers to the physical data object that you plan to replace.
- Update the streaming mapping. If the mapping contains Kafka target, Aggregator transformation, Joiner transformation, or Normalizer transformation, replace the data object or transformation, and then update the mapping because of the changed behavior of these transformations and data objects.
- Verify the deferred data object types. If the streaming mapping contains unsupported transformations or data objects, contact Informatica Global Customer Support.

### Re-create the Physical Data Objects

When you re-create the physical data objects, the physical data objects get the header ports and some properties are not available for some data objects. Update the existing mapping with the newly created physical data objects.

1. Go to the existing mapping, select the data object from the mapping.
2. Click the **Properties** tab. On the **Column Projection** tab, click **Edit Schema**.
3. Note the schema information from the **Edit Schema** dialog box.
4. Note the parameters information from the **Parameters** tab.
5. Create new physical data objects.

After you re-create the data objects, the physical data objects get the required header ports. The Microsoft Azure does not support the following properties and are not available for Azure Event Hubs data objects:

- Consumer Properties
- Partition Count

## Update the Streaming Mappings

After you re-create the data object, replace the existing data objects with the re-created data objects. If the mapping contains Normaliser Transformation, Aggregator transformation, or Joiner transformation, update the mapping because of the changed behavior of these transformations and data objects.

### Transformation Updates

If a transformation uses a complex port, configure the type configuration property of the port because the property refers to the physical data object that you replaced.

### Aggregator and Joiner Transformation Updates

An Aggregator transformation must be downstream from a Joiner transformation. A Window transformation must be directly upstream from both Aggregator and Joiner transformations. Previously, you could use an Aggregator transformation anywhere in the streaming mapping.

If a mapping contains an Aggregator transformation upstream from a Joiner transformation, move the Aggregator transformation downstream from a Joiner transformation. Add a Window transformation directly upstream from both Aggregator and Joiner transformations.

## Verify the Deferred Data Object Types

After you upgrade, the streaming mappings might contain some transformations and data objects that are deferred.

The following table lists the data object types to which the support is deferred to a future release:

Object Type	Object
Transformation	Data Masking

If you want to continue using the mappings that contain deferred data objects or transformations, you must contact Informatica Global Customer Support.

## CHAPTER 5

# Cloudera CDH Integration Tasks

This chapter includes the following topics:

- [Cloudera CDH Task Flows, 93](#)
- [Prepare for Cluster Import from Cloudera CDH, 98](#)
- [Create a Cluster Configuration, 103](#)
- [Verify or Refresh the Cluster Configuration , 105](#)
- [Verify JDBC Drivers for Sqoop Connectivity, 106](#)
- [Set the Locale for Cloudera CDH 6.x, 107](#)
- [Enable Data Preparation of JSON Files on Cloudera CDH, 108](#)
- [Complete Upgrade Tasks, 108](#)

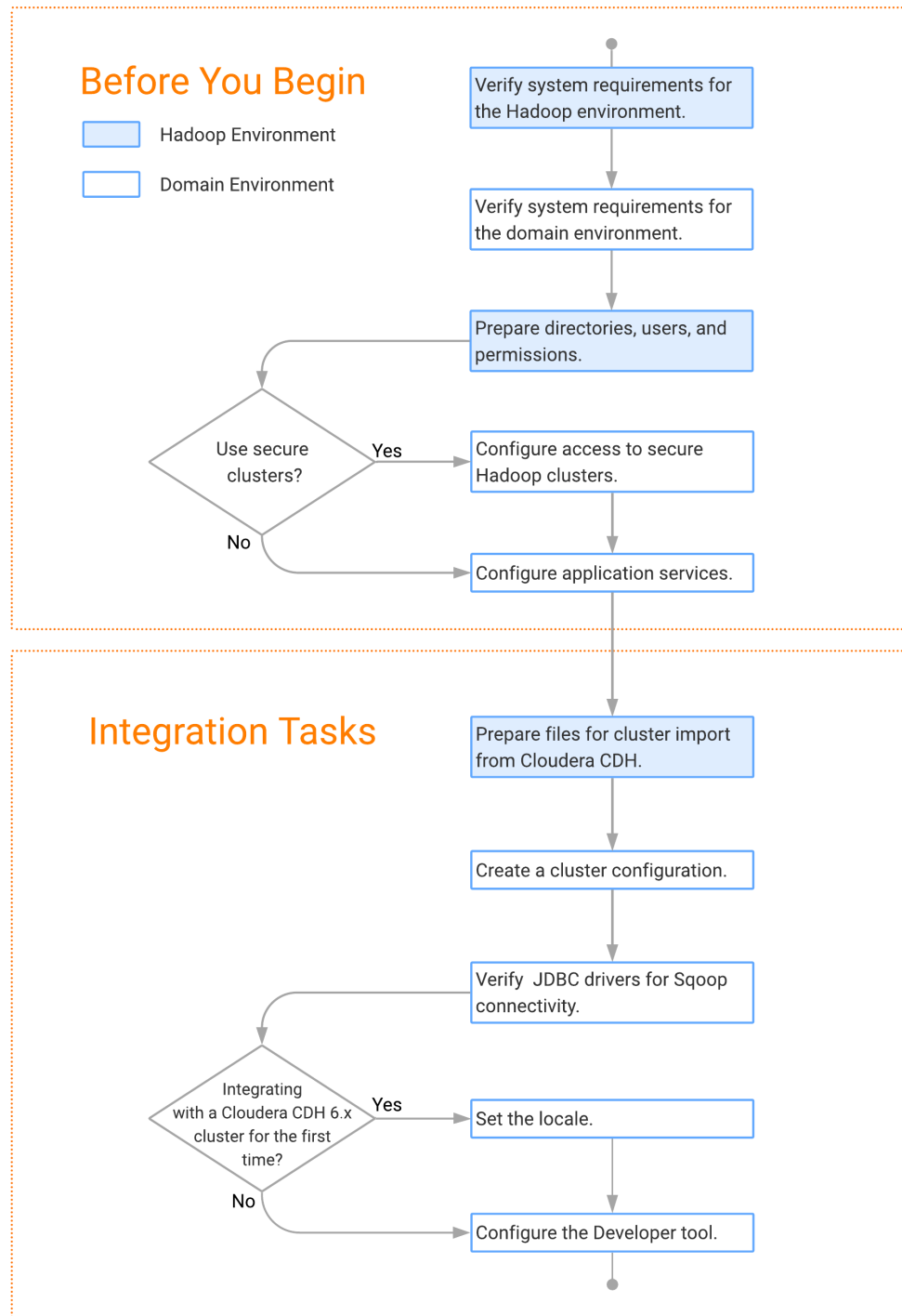
## Cloudera CDH Task Flows

Depending on whether you want to integrate or upgrade Data Engineering Integration in a Cloudera CDH environment, you can use the flow charts to perform the following tasks:

- Integrate the Informatica domain with Cloudera CDH for the first time.
- Upgrade from version 10.2.1.
- Upgrade from version 10.2.
- Upgrade from a version earlier than 10.2.

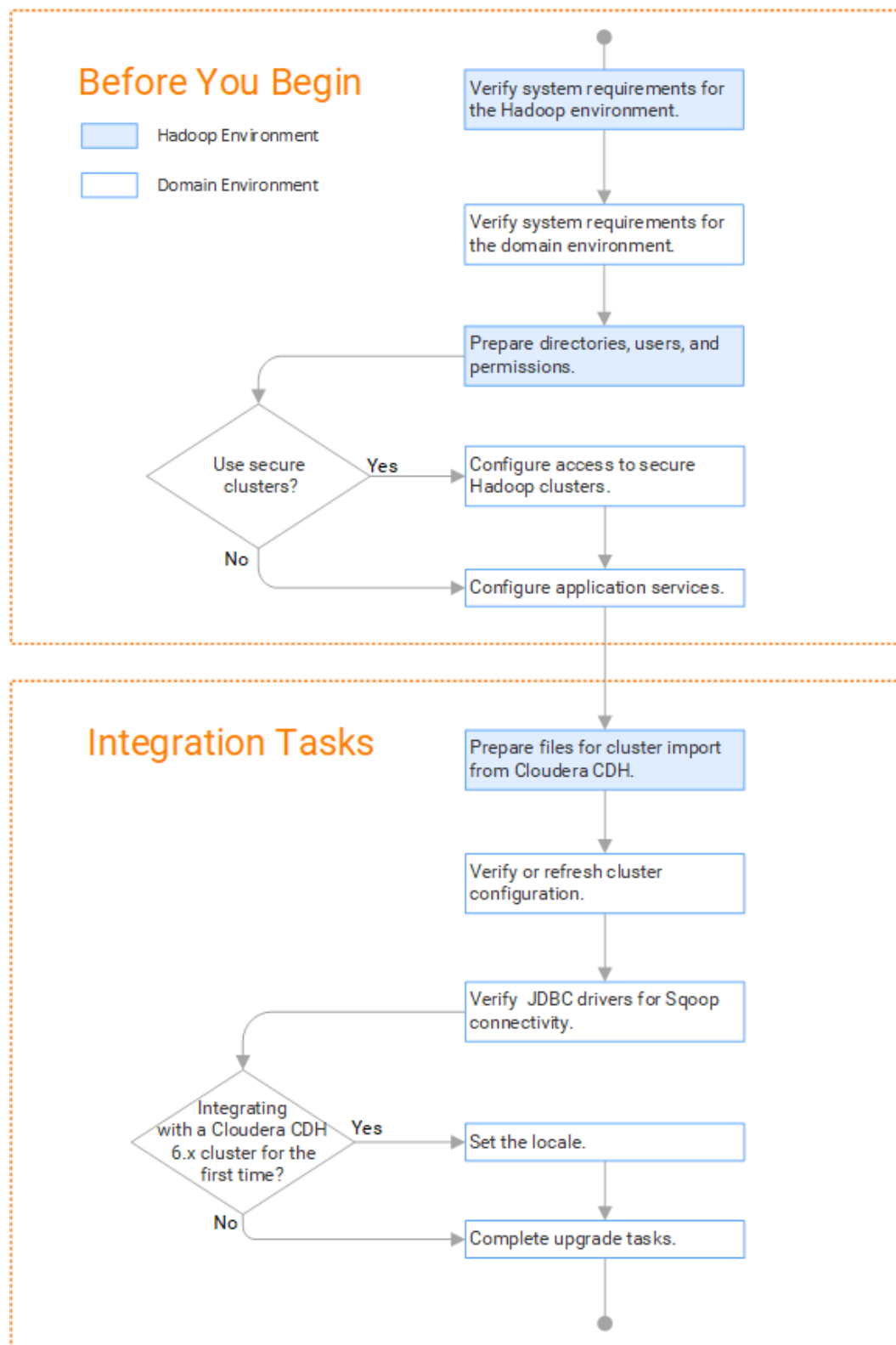
## Task Flow to Integrate with Cloudera CDH

The following diagram shows the task flow to integrate the Informatica domain with Cloudera CDH:



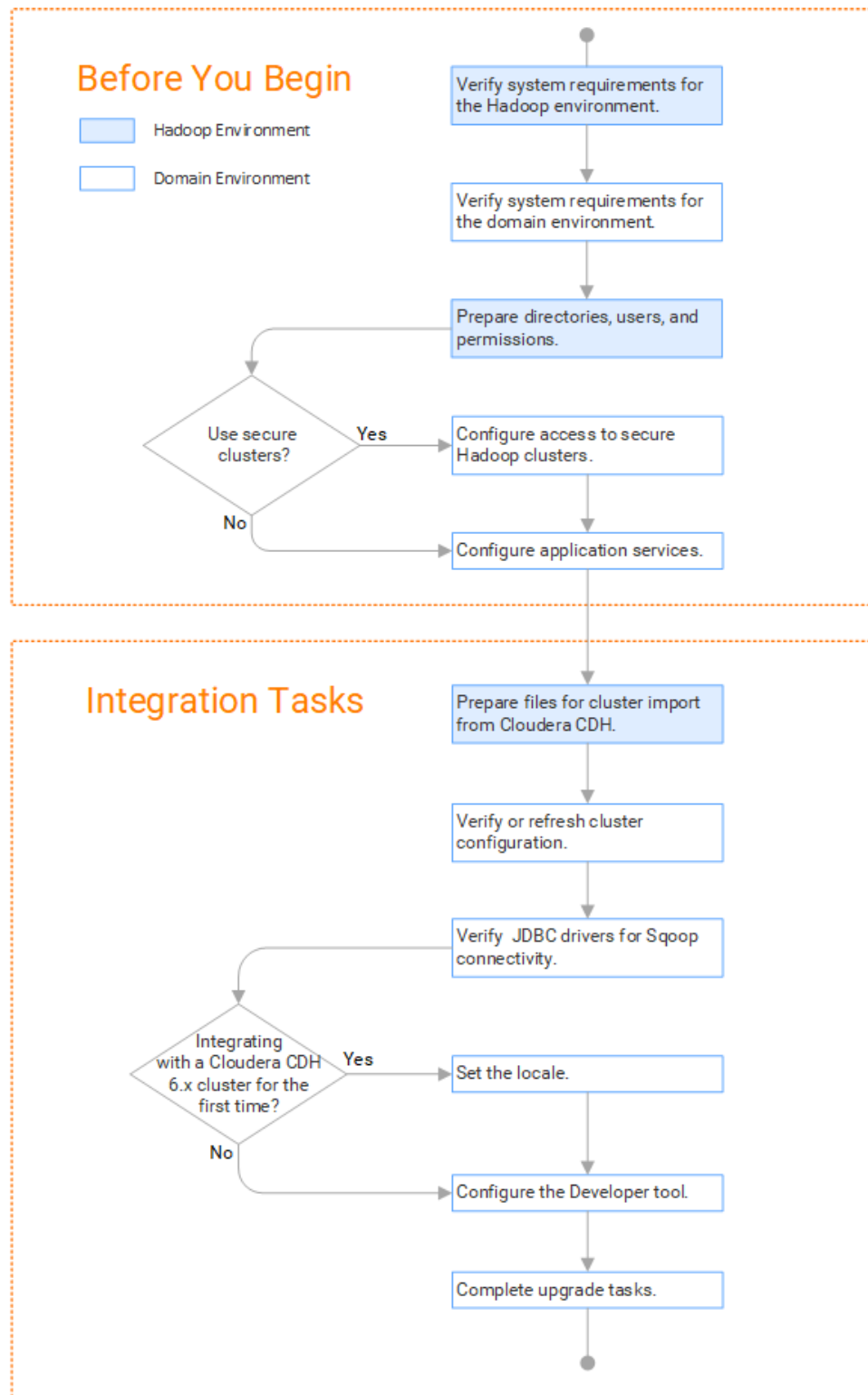
## Task Flow to Upgrade from Version 10.2.1

The following diagram shows the task flow to upgrade version 10.2.1 for Cloudera CDH:



## Task Flow to Upgrade from Version 10.2

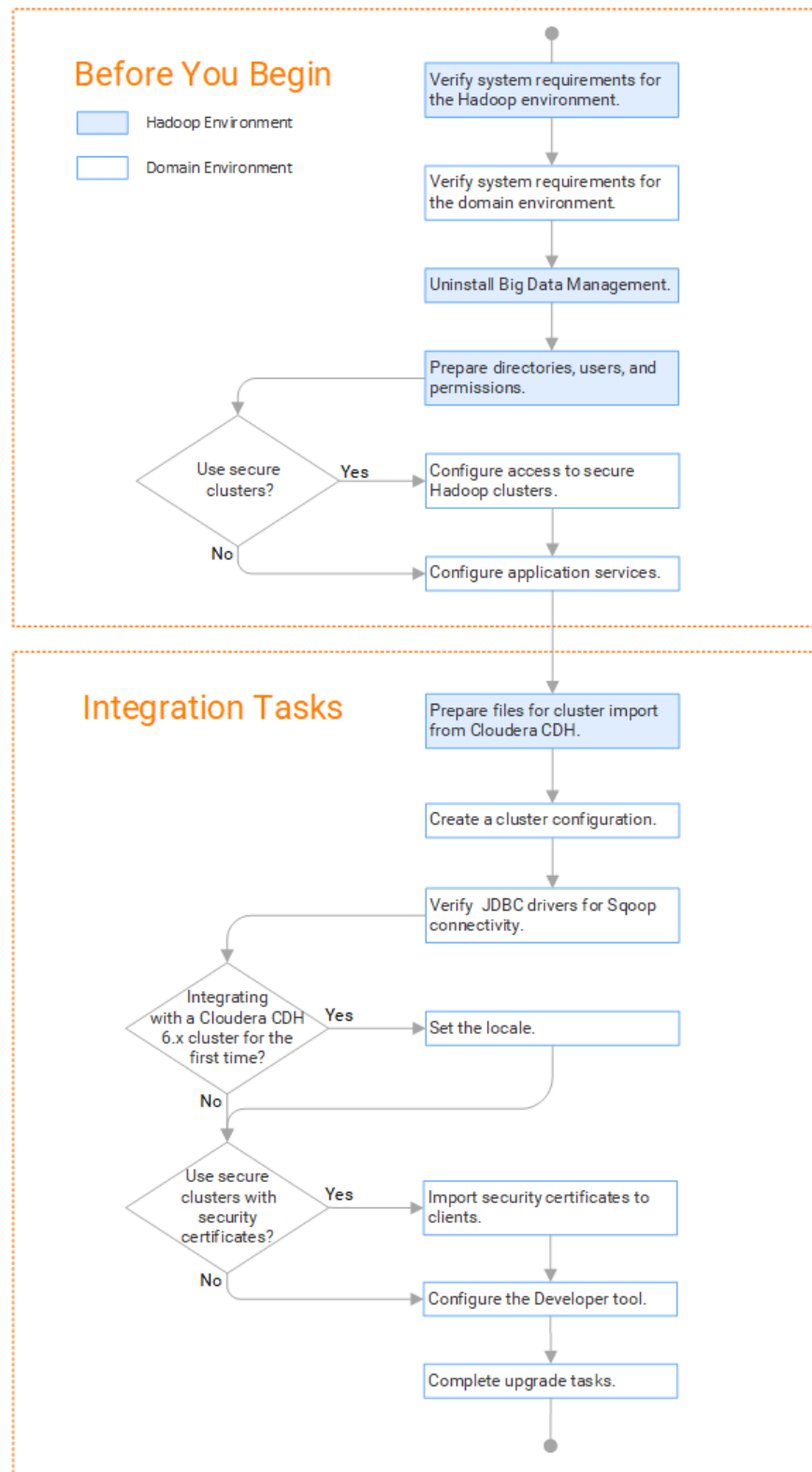
The following diagram shows the task flow to upgrade version 10.2 for Cloudera CDH:





## Task Flow to Upgrade from a Version Earlier than 10.2

The following diagram shows the task flow to upgrade from a version earlier than 10.2 for Cloudera CDH:



# Prepare for Cluster Import from Cloudera CDH

Before the Informatica administrator can import cluster information to create a cluster configuration in the Informatica domain, the Hadoop administrator must perform some preliminary tasks.

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

**Note:** If you are upgrading from a previous version, verify the properties and suggested values, as Data Engineering Integration might require additional properties or different values for existing properties.

Complete the following tasks to prepare the cluster before the Informatica administrator creates the cluster configuration:

1. Verify property values in \*-site.xml files that Data Engineering Integration needs to run mappings in the Hadoop environment.
2. Provide information to the Informatica administrator that is required to import cluster information into the domain. Depending on the method of import, perform one of the following tasks:
  - To import directly from the cluster, give the Informatica administrator cluster authentication information to connect to the cluster.
  - To import from an archive file, export cluster information and provide an archive file to the Data Engineering Integration administrator.

## Configure \*-site.xml Files for Cloudera CDH

The Hadoop administrator needs to configure \*-site.xml file properties and restart impacted services before the Informatica administrator imports cluster information into the domain.

### core-site.xml

Configure the following properties in the core-site.xml file:

#### **fs.s3.enableServerSideEncryption**

Enables server side encryption for S3 buckets. Required for SSE and SSE-KMS encryption.

Set to: TRUE

#### **fs.s3a.access.key**

The ID for the Blaze and Spark engines to connect to the Amazon S3 file system.

Set to your access key.

#### **fs.s3a.secret.key**

The password for the Blaze and Spark engines to connect to the Amazon S3 file system

Set to your access ID.

#### **fs.s3a.server-side-encryption-algorithm**

The server-side encryption algorithm for S3. Required for SSE and SSE-KMS encryption. Set to the encryption algorithm used.

#### **hadoop.proxyuser.<proxy user>.groups**

Defines the groups that the proxy user account can impersonate. On a secure cluster the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to group names of impersonation users separated by commas. If less security is preferred, use the wildcard " \* " to allow impersonation from any group.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.proxyuser.<proxy user>.hosts**

Defines the host machines that a user account can impersonate. On a secure cluster the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to a single host name or IP address, or set to a comma-separated list. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.proxyuser.hive.hosts**

Comma-separated list of hosts that you want to allow the Hive user to impersonate on a non-secure cluster.

When `hive.server2.enable.doAs` is false, append a comma-separated list of Informatica server host names or IP address where the Data Integration Service is running. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to this property, restart the cluster services that use core-site configuration values.

#### **io.compression.codecs**

Enables compression on temporary staging tables.

Set to a comma-separated list of compression codec classes on the cluster.

#### **hadoop.security.auth\_to\_local**

Translates the principal names from the Active Directory and MIT realm into local names within the Hadoop cluster. Based on the Hadoop cluster used, you can set multiple rules.

Set to: `RULE:[1:$1@$0](^.*@YOUR.REALM)s/^.*(.)@YOUR.REALM\.COM$/ $1/g`

Set to: `RULE:[2:$1@$0](^.*@YOUR.REALM\.$)s/^.*(.)@YOUR.REALM\.COM$/ $1/g`

#### [fair-scheduler.xml](#)

Configure the following properties in the fair-scheduler.xml file:

##### **allowPreemptionFrom**

Enables preemption for the Fair Scheduler. The Blaze engine does not support preemption. If YARN preemption is enabled for the cluster, you need to disable it for the queue allocated to the Blaze engine.

Set to FALSE for the queue allocated to the Blaze engine.

For example:

```
<queue name="Blaze">
  <weight>1.0</weight>
  <allowPreemptionFrom>>false</allowPreemptionFrom>
  <schedulingPolicy>fsp</schedulingPolicy>
```

```

        <aclSubmitApps>*</aclSubmitApps>
        <aclAdministerApps>*</aclAdministerApps>
    </queue>

```

## hbase-site.xml

Configure the following properties in the hbase-site.xml file:

### **zookeeper.znode.parent**

Identifies HBase master and region servers.

Set to the relative path to the znode directory of HBase.

## hdfs-site.xml

Configure the following properties in the hdfs-site.xml file:

### **dfs.encryption.key.provider.uri**

The KeyProvider used to interact with encryption keys when reading and writing to an encryption zone. Required if sources or targets reside in the HDFS encrypted zone on Java KeyStore KMS-enabled Cloudera CDH cluster or a Ranger KMS-enabled Hortonworks HDP cluster.

Set to: kmf://http@xx11.xyz.com:16000/kms

## hive-site.xml

Configure the following properties in the hive-site.xml file:

### **hive.cluster.delegation.token.store.class**

The token store implementation. Required for HiveServer2 high availability and load balancing.

Set to: org.apache.hadoop.hive.thrift.DBTokenStore

### **hive.exec.dynamic.partition**

Enables dynamic partitioned tables for Hive tables. Applicable for Hive versions 0.9 and earlier.

Set to: TRUE

### **hive.exec.dynamic.partition.mode**

Allows all partitions to be dynamic. Required if you use Sqoop and define a DDL query to create or replace a partitioned Hive target at run time.

Set to: nonstrict

### **hiveserver2\_load\_balancer**

Enables high availability for multiple HiveServer2 hosts.

Set to: jdbc:hive2://<HiveServer2 Load Balancer>:<HiveServer2 Port>/default;principal=hive/<HiveServer2 load Balancer>@<REALM>

## mapred-site.xml

Configure the following properties in the mapred-site.xml file:

### **mapreduce.framework.name**

The run-time framework to run MapReduce jobs. Values can be local, classic, or yarn. Required for Sqoop.

Set to: yarn

### **mapreduce.jobhistory.address**

Location of the MapReduce JobHistory Server. The default port is 10020. Required for Sqoop.

Set to: <MapReduce JobHistory Server>:<port>

**mapreduce.jobhistory.intermediate-done-dir**

Directory where MapReduce jobs write history files. Required for Sqoop.

Set to: /mr-history/tmp

**mapreduce.jobhistory.done-dir**

Directory where the MapReduce JobHistory Server manages history files. Required for Sqoop.

Set to: /mr-history/done

**mapreduce.jobhistory.principal**

The Service Principal Name for the MapReduce JobHistory Server. Required for Sqoop.

Set to: mapred/\_HOST@YOUR-REALM

**mapreduce.jobhistory.webapp.address**

Web address of the MapReduce JobHistory Server. The default value is 19888. Required for Sqoop.

Set to: <host>:<port>

**yarn.app.mapreduce.am.staging-dir**

The HDFS staging directory used while submitting jobs.

Set to the staging directory path.

[yarn-site.xml](#)

Configure the following properties in the yarn-site.xml file:

**yarn.application.classpath**

Required for dynamic resource allocation.

Add spark\_shuffle.jar to the class path. The .jar file must contain the class "org.apache.spark.network.yarn.YarnShuffleService."

**yarn.nodemanager.resource.memory-mb**

The maximum RAM available for each container. Set the maximum memory on the cluster to increase resource memory available to the Blaze engine.

Set to 16 GB if value is less than 16 GB.

**yarn.nodemanager.resource.cpu-vcores**

The number of virtual cores for each container. Required for Blaze engine resource allocation.

Set to 10 if the value is less than 10.

**yarn.scheduler.minimum-allocation-mb**

The minimum RAM available for each container. Required for Blaze engine resource allocation.

Set to 6 GB if the value is less than 6 GB.

**yarn.nodemanager.vmem-check-enabled**

Disables virtual memory limits for containers. Required for the Blaze and Spark engines.

Set to: false

**yarn.nodemanager.aux-services**

Required for dynamic resource allocation for the Spark engine.

Add an entry for "spark\_shuffle."

**yarn.nodemanager.aux-services.spark\_shuffle.class**

Required for dynamic resource allocation for the Spark engine.

Set to: org.apache.spark.network.yarn.YarnShuffleService

**yarn.resourcemanager.scheduler.class**

Defines the YARN scheduler that the Data Integration Service uses to assign resources.

Set to: org.apache.hadoop.yarn.server.resourcemanager.scheduler

## Prepare for Direct Import from Cloudera CDH

If you plan to provide direct access to the Informatica administrator to import cluster information, provide the required connection information.

The following table describes the information that you need to provide to the Informatica administrator to create the cluster configuration directly from the cluster:

Property	Description
Host	IP address of the cluster manager.
Port	Port of the cluster manager.
User ID	Cluster user ID.
Password	Password for the user.
Cluster name	Name of the cluster. Use the display name if the cluster manager manages multiple clusters. If you do not provide a cluster name, the wizard imports information based on the default cluster. To find the correct Cloudera cluster name when you have multiple clusters, perform the following steps: 1. Log in to Cloudera Manager adding the following string to the URL: /api/v8/clusters 2. Provide the Informatica Administrator the cluster property name that appears in the browser tab.

## Prepare the Archive File for Import from Cloudera CDH

If you plan to provide an archive file for the Informatica administrator, ensure that you include all required site-\*.xml files.

Create a .zip or .tar file that contains the following \*-site.xml files:

- core-site.xml
- hbase-site.xml. Required only for access to HBase sources and targets.
- hdfs-site.xml
- hive-site.xml
- mapred-site.xml
- yarn-site.xml

Give the Informatica administrator access to the archive file to import the cluster information into the domain.

# Create a Cluster Configuration

After the Hadoop administrator prepares the cluster for import, the Informatica administrator must create a cluster configuration.

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from version 10.1.1 or earlier.

A cluster configuration is an object in the domain that contains configuration information about the Hadoop cluster. The cluster configuration enables the Data Integration Service to push mapping logic to the Hadoop environment. Import configuration properties from the Hadoop cluster to create a cluster configuration.

The import process imports values from \*-site.xml files into configuration sets based on the individual \*-site.xml files. When you perform the import, the cluster configuration wizard can create Hadoop, HBase, HDFS, and Hive connection to access the Hadoop environment. If you choose to create the connections, the wizard also associates the cluster configuration with the connections.

**Note:** If you are integrating for the first time and you imported the cluster configuration when you ran the installer, you *must* re-create or refresh the cluster configuration.

## Before You Import

Before you can import the cluster configuration, you must get information from the Hadoop administrator based on the method of import.

If you import directly from the cluster, contact the Hadoop administrator to get cluster connection information. If you import from a file, get an archive file of exported cluster information.

## Importing a Hadoop Cluster Configuration from the Cluster

When you import the Hadoop cluster configuration directly from the cluster, you provide information to connect to the cluster.

Get cluster connection information from the Hadoop administrator.

1. From the **Connections** tab, click the **ClusterConfigurations** node in the Domain Navigator.
2. From the Actions menu, select **New > Cluster Configuration**.

The **Cluster Configuration** wizard opens.

3. Configure the following General properties:

Property	Description
Cluster configuration name	Name of the cluster configuration.
Description	Optional description of the cluster configuration.
Distribution type	The cluster Hadoop distribution type.

Property	Description
Method to import the cluster configuration	Choose <b>Import from cluster</b> .
Create connections	<p>Choose to create Hadoop, HDFS, Hive, and HBase connections.</p> <p>If you choose to create connections, the Cluster Configuration wizard associates the cluster configuration with each connection that it creates.</p> <p>The Hadoop connection contains default values for properties such as cluster environment variables, cluster path variables, and advanced properties. Based on the cluster environment and the functionality that you use, you can add to the default values or change the default values of these properties. For a list of Hadoop connection properties to configure, see Hadoop Cluster Connection Properties.</p> <p>If you do not choose to create connections, you must manually create them and associate the cluster configuration with them.</p> <p><b>Important:</b> When the wizard creates the Hive connection, it populates the Metadata Connection String and the Data Access Connection String properties with the value from the hive.metastore.uris property. If the Hive metastore and HiveServer2 are running on different nodes, you must update the Metadata Connection String to point to the HiveServer2 host.</p>

The cluster properties appear.

4. Configure the following properties:

Property	Description
Host	Host name or IP address of the cluster manager.
Port	Port of the cluster manager.
User ID	Cluster user ID.
Password	Password for the user.
Cluster name	Name of the cluster. Use the display name if the cluster manager manages multiple clusters. If you do not provide a cluster name, the wizard imports information based on the default cluster.

5. Click **Next** and verify the cluster configuration information on the summary page.

## Importing a Hadoop Cluster Configuration from a File

You can import properties from an archive file to create a cluster configuration.

Before you import from the cluster, you must get the archive file from the Hadoop administrator.

1. From the **Connections** tab, click the **ClusterConfigurations** node in the Domain Navigator.
2. From the Actions menu, select **New > Cluster Configuration**.

The **Cluster Configuration** wizard opens.



3. Configure the following properties:

Property	Description
Cluster configuration name	Name of the cluster configuration.
Description	Optional description of the cluster configuration.
Distribution type	The cluster Hadoop distribution type.
Method to import the cluster configuration	Choose <b>Import from file</b> to import properties from an archive file.
Create connections	<p>Choose to create Hadoop, HDFS, Hive, and HBase connections.</p> <p>If you choose to create connections, the Cluster Configuration wizard associates the cluster configuration with each connection that it creates.</p> <p>The Hadoop connection contains default values for properties such as cluster environment variables, cluster path variables, and advanced properties. Based on the cluster environment and the functionality that you use, you can add to the default values or change the default values of these properties. For a list of Hadoop connection properties to configure, see <a href="#">“Configuring Hadoop Connection Properties” on page 269</a>.</p> <p>If you do not choose to create connections, you must manually create them and associate the cluster configuration with them.</p> <p><b>Important:</b> When the wizard creates the Hive connection, it populates the Metadata Connection String and the Data Access Connection String properties with the value from the hive.metastore.uris property. If the Hive metastore and HiveServer2 are running on different nodes, you must update the Metadata Connection String to point to the HiveServer2 host.</p>

4. Click **Browse** to select a file. Select the file and click **Open**.
5. Click **Next** and verify the cluster configuration information on the summary page.

## Verify or Refresh the Cluster Configuration

You might need to refresh the cluster configuration or update the distribution version in the cluster configuration when you upgrade.

### Perform this task in the following situation:

- You upgraded from version 10.2 or later.

### Verify the Cluster Configuration

The cluster configuration contains a property for the distribution version. The verification task depends on the version you upgraded:

#### Upgrade from 10.2

If you upgraded from 10.2 and you changed the distribution version, you need to verify the distribution version in the General properties of the cluster configuration.

### Upgrade from 10.2.1

Effective in version 10.2.1, Informatica assigns a default version to each Hadoop distribution type. If you configure the cluster configuration to use the default version, the upgrade process upgrades to the assigned default version if the version changes. If you have not upgraded your Hadoop distribution to Informatica's default version, you need to update the distribution version property.

For example, suppose the assigned default Hadoop distribution version for 10.2.1 is  $n$ , and for 10.2.2 is  $n+1$ . If the cluster configuration uses the default supported Hadoop version of  $n$ , the upgraded cluster configuration uses the default version of  $n+1$ . If you have not upgraded the distribution in the Hadoop environment you need to change the cluster configuration to use version  $n$ .

If you configure the cluster configuration to use a distribution version that is not the default version, you need to update the distribution version property in the following circumstances:

- Informatica dropped support for the distribution version.
- You changed the distribution version.

### Refresh the Cluster Configuration

If you updated any of the \*-site.xml files noted in the topic to prepare for cluster import, you need to refresh the cluster configuration in the Administrator tool.

## Verify JDBC Drivers for Sqoop Connectivity

Verify that you have the JDBC drivers to access JDBC-compliant databases in the Hadoop environment. You might need separate drivers for metadata import and for run-time processing.

#### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from version 10.2.1 or earlier.

You download drivers based on design-time and run-time requirements:

- **Design-time.** To import metadata, you can use the DataDirect drivers packaged with the Informatica installer if they are available. If they are not available, use any Type 4 JDBC driver that the database vendor recommends.
- **Run-time.** To run mappings, use any Type 4 JDBC driver that the database vendor recommends. Some distributions support other drivers to use Sqoop connectors. You cannot use the DataDirect drivers for run-time processing.

## Verify Design-time Drivers

Use the DataDirect JDBC drivers packaged with the Informatica installer to import metadata from JDBC-compliant databases. If the DataDirect JDBC drivers are not available for a specific JDBC-compliant database, download the Type 4 JDBC driver associated with that database.

Copy the JDBC driver .jar files to the following location on the Developer tool machine:

```
<Informatica installation directory>\clients\externaljdbcjars
```

## Verify Run-time Drivers

Verify run-time drivers for mappings that access JDBC-compliant databases in the Hadoop environment. Use any Type 4 JDBC driver that the database vendor recommends.

1. Download Type 4 JDBC drivers associated with the JCBC-compliant databases that you want to access.
2. To use Sqoop TDCH Cloudera Connector Powered by Teradata, perform the following tasks:
  - Download all .jar files in the Cloudera Connector Powered by Teradata package from the following location: <http://www.cloudera.com/downloads.html>. The package has the following naming convention: `sqoop-connector-teradata-<version>.tar`  
If you use Cloudera CDH version 6.x, you must download the `sqoop-connector-teradata-1.7c6.jar` file.
  - Download `terajdbc4.jar` and `tdgssconfig.jar` from the following location: <http://downloads.teradata.com/download/connectivity/jdbc-driver>  
If you use Cloudera CDH version 6.x, you must also download the `junit-4.11.jar` file.
3. To optimize the Sqoop mapping performance on the Spark engine while writing data to an HDFS complex file target of the Parquet format, download the following .jar files:
  - `parquet-hadoop-bundle-1.6.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-hadoop-bundle/1.6.0/>
  - `parquet-avro-1.6.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-avro/1.6.0/>
  - `parquet-column-1.5.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-column/1.5.0/>
4. Copy all of the .jar files to the following directory on the machine where the Data Integration Service runs:  

```
<Informatica installation directory>\externaljdbcjars
```

Changes take effect after you recycle the Data Integration Service. At run time, the Data Integration Service copies the .jar files to the Hadoop distribution cache so that the .jar files are accessible to all nodes in the cluster.

## Set the Locale for Cloudera CDH 6.x

If you want to process data that contains non-ASCII characters, you must integrate the locale setting on Data Engineering Integration with the locale setting on the cluster.

### Perform this task in the following situations:

- You are integrating with a Cloudera CDH 6.x cluster for the first time.

To integrate the locale setting, complete the following tasks:

1. In the Hadoop connection, navigate to **Hadoop Cluster Properties**. As the value for the property **Cluster Environment Variables**, configure the locale environment variables, such as the `LANG` or `LC_ALL` environment variable.  
  
The locale setting in the Hadoop connection must match the locale setting that is configured in the domain.

2. In Cloudera Manager, add the environment variables to the following YARN property:

```
yarn.nodemanager.env-whitelist
```

## Enable Data Preparation of JSON Files on Cloudera CDH

If you integrate Enterprise Data Preparation with a Cloudera CDH Hadoop cluster, you must specify the location of a Hive .jar file in the Hive Auxiliary JARs Directory in CDH to enable data preparation of JSON files.

Perform this task in the following situations:
<ul style="list-style-type: none"><li>- You are integrating for the first time.</li></ul>



1. Search for hive-hcatalog-core.jar in the Cloudera CDH installation directory.

You can generally find the .jar file in the following directory:

```
/opt/cloudera/parcels/CDH/lib/hive-hcatalog/share/hcatalog/
```

2. Log in to Cloudera Manager.
3. Select **Hive** in the cluster.
4. Click the **Configuration** tab, then select the **Advanced** category.
5. Enter the path to the directory containing the .jar file in the Hive Auxiliary JARs Directory property.

If the file is in the location noted, the path to the directory is:

```
/opt/cloudera/parcels/CDH/lib/hive-hcatalog/share/hcatalog/
```

6. Restart the Hive server.

## Complete Upgrade Tasks

If you upgraded the Informatica platform from version 10.2.2, update connections.

### Update Connections

You might need to update connections based on the version you are upgrading from.

If you did not create connections when you created the cluster configuration, you need to update the connections.

## Configure the Hadoop Connection

To use properties that you customized in the `hadoopEnv.properties` file, you must configure the Hadoop connection properties such as cluster environment variables, cluster path variables, and advanced properties.

### Perform this task in the following situation:

- You upgraded from version 10.1.1 or earlier.

When you run the Informatica upgrade, the installer backs up the existing `hadoopEnv.properties` file. You can find the backup `hadoopEnv.properties` file in the following location:

```
<Previous Informatica installation directory>/services/shared/hadoop/<Hadoop distribution name>_<version>/infaConf
```

Edit the Hadoop connection in the Administrator tool or the Developer tool to include any properties that you manually configured in the `hadoopEnv.properties` file. The Hadoop connection contains default values for properties such as cluster environment and path variables and advanced properties. You can update the default values to match the properties in the `hadoopEnv.properties` file.

## Replace the Connections with New Connections

If you created connections when you imported the cluster configuration, you need to replace connections in mappings with the new connections.

### Perform this task in the following situation:

- You upgraded from version 10.1.1 or earlier.

The method that you use to replace connections in mappings depends on the type of connection.

#### Hadoop connection

Run the following commands to replace the connections:

- `infacmd dis replaceMappingHadoopRuntimeConnections`. Replaces connections associated with mappings that are deployed in applications.
- `infacmd mrs replaceMappingHadoopRuntimeConnections`. Replaces connections associated with mappings that you run from the Developer tool.

For information about the `infacmd` commands, see the *Informatica Command Reference*.

#### Hive, HDFS, and HBase connections

You must replace the connections manually.

## Complete Connection Upgrade

If you *did not* create connections when you imported the cluster configuration, you must update connection properties for Hadoop, Hive, HDFS, and HBase connections.

### Perform this task in the following situations:

- You upgraded from version 10.2.2 or earlier.

Perform the following tasks to update the connections:

### Update changed properties

Review connections that you created in a previous release to update the values for connection properties. For example, if you added nodes to the cluster or if you updated the distribution version, you might need to verify host names, URIs, or port numbers for some of the properties.

### Associate the cluster configuration

The Hadoop, Hive, HDFS, and HBase connections must be associated with a cluster configuration. Complete the following tasks:

1. Run `infacmd isp listConnections` to identify the connections that you need to upgrade. Use `-ct` to list connections of a particular type.
2. Run `infacmd isp UpdateConnection` to associate the cluster configuration with the connection. Use `-cn` to name the connection and `-o clusterConfigID` to associate the cluster configuration with the connection.

For more information about `infacmd`, see the *Informatica Command Reference*.

## Replace Hive Run-time Connections with Hadoop Connections

Effective in version 10.2.2, Big Data Management dropped support for the Hive engine and Hive run-time connections. If you used Hive connections to run mappings on the Hadoop cluster, you must generate Hadoop connections from the Hive connections.

#### Perform this task in the following situations:

- You upgraded from version 10.1.1 or earlier.
- The Hive connections are configured to run mappings in the Hadoop environment.

The upgrade process generates a connection name for the Hadoop connection and replaces the connection name in the mappings. It does not create the physical connection object. When the upgrade is complete, you must run a command to generate the connection. You generate Hadoop connections from Hive connections that are configured to run mappings in the Hadoop environment.

You must generate Hadoop connections from Hive connections that are configured to run mappings in the Hadoop environment.

1. Run `infacmd isp generateHadoopConnectionFromHiveConnection` to generate a Hadoop connection from a Hive connection that is configured to run in the Hadoop environment.  
The command names the connection as follows: "Autogen\_<Hive connection name>." If the connection name exceeds the 128 character limit, the command fails.
2. If the command fails, perform the following tasks:
  - a. Rename the connection to meet the character limit and run the command again.
  - b. Run `infacmd dis replaceMappingHadoopRuntimeConnections` to replace connections associated with mappings that are deployed in applications.
  - c. Run `infacmd mrs replaceMappingHadoopRuntimeConnections` to replace connections associated with mappings that you run from the Developer tool.
3. If the Hive connection was parameterized, you must update the connection names in the parameter file. Verify that the Hive sources, Hive targets, and the Hive engine parameters are updated with the correct connection name.
4. If any properties changed in the cluster, such as host names, URIs, or port numbers, you must update the properties in the connections.

For information about the `infacmd` commands, see the *Informatica® Command Reference*.

## Update Streaming Objects

Data Engineering Streaming uses Spark Structured Streaming to process data instead of Spark Streaming. To support Spark Structured Streaming, some header ports are added to the data objects, and support to some of the data objects and transformations are deferred to a future release. The behavior of some of the data objects is also updated.

After you upgrade, the existing streaming mappings become invalid because of the unavailable header ports, the unsupported transformations or data objects, and the behavior change of some data objects.

<b>Perform this task in the following situations:</b>
<ul style="list-style-type: none"><li>- You upgraded from version 10.1.1, 10.2.0, or 10.2.1.</li></ul>



To use an existing streaming mapping, perform the following tasks:

- Re-create the physical data objects. After you re-create the physical data objects, the data objects get the required header ports, such as timestamp, partitionID, or key based on the data object.
- In a Normalizer transformation, if the **Occurs** column is set to Auto, re-create the Normalizer transformation. You must re-create the Normalizer transformation because the type configuration property of the complex port refers to the physical data object that you plan to replace.
- Update the streaming mapping. If the mapping contains Kafka target, Aggregator transformation, Joiner transformation, or Normalizer transformation, replace the data object or transformation, and then update the mapping because of the changed behavior of these transformations and data objects.
- Verify the deferred data object types. If the streaming mapping contains unsupported transformations or data objects, contact Informatica Global Customer Support.

### Re-create the Physical Data Objects

When you re-create the physical data objects, the physical data objects get the header ports and some properties are not available for some data objects. Update the existing mapping with the newly created physical data objects.

1. Go to the existing mapping, select the data object from the mapping.
2. Click the **Properties** tab. On the **Column Projection** tab, click **Edit Schema**.
3. Note the schema information from the **Edit Schema** dialog box.
4. Note the parameters information from the **Parameters** tab.
5. Create new physical data objects.

After you re-create the data objects, the physical data objects get the required header ports. The Microsoft Azure does not support the following properties and are not available for Azure Event Hubs data objects:

- Consumer Properties
- Partition Count

## Update the Streaming Mappings

After you re-create the data object, replace the existing data objects with the re-created data objects. If the mapping contains Normaliser Transformation, Aggregator transformation, or Joiner transformation, update the mapping because of the changed behavior of these transformations and data objects.

### Transformation Updates

If a transformation uses a complex port, configure the type configuration property of the port because the property refers to the physical data object that you replaced.

### Aggregator and Joiner Transformation Updates

An Aggregator transformation must be downstream from a Joiner transformation. A Window transformation must be directly upstream from both Aggregator and Joiner transformations. Previously, you could use an Aggregator transformation anywhere in the streaming mapping.

If a mapping contains an Aggregator transformation upstream from a Joiner transformation, move the Aggregator transformation downstream from a Joiner transformation. Add a Window transformation directly upstream from both Aggregator and Joiner transformations.

## Verify the Deferred Data Object Types

After you upgrade, the streaming mappings might contain some transformations and data objects that are deferred.

The following table lists the data object types to which the support is deferred to a future release:

Object Type	Object
Transformation	Data Masking

If you want to continue using the mappings that contain deferred data objects or transformations, you must contact Informatica Global Customer Support.



## CHAPTER 6

# Cloudera CDP Integration Tasks

This chapter includes the following topics:

- [Cloudera CDP Task Flows, 113](#)
- [Prepare for Cluster Import from Cloudera CDP, 116](#)
- [Create a Cluster Configuration, 121](#)
- [Copy the Truststore File to the Informatica Domain, 124](#)
- [Configure the Impersonation User for Operating System Profiles, 125](#)
- [Verify JDBC Drivers for Sqoop Connectivity, 125](#)
- [Set the Locale for Cloudera CDP, 126](#)
- [Enable Data Preparation of JSON Files on Cloudera CDP, 127](#)
- [Configure the Developer Tool, 127](#)

## Cloudera CDP Task Flows

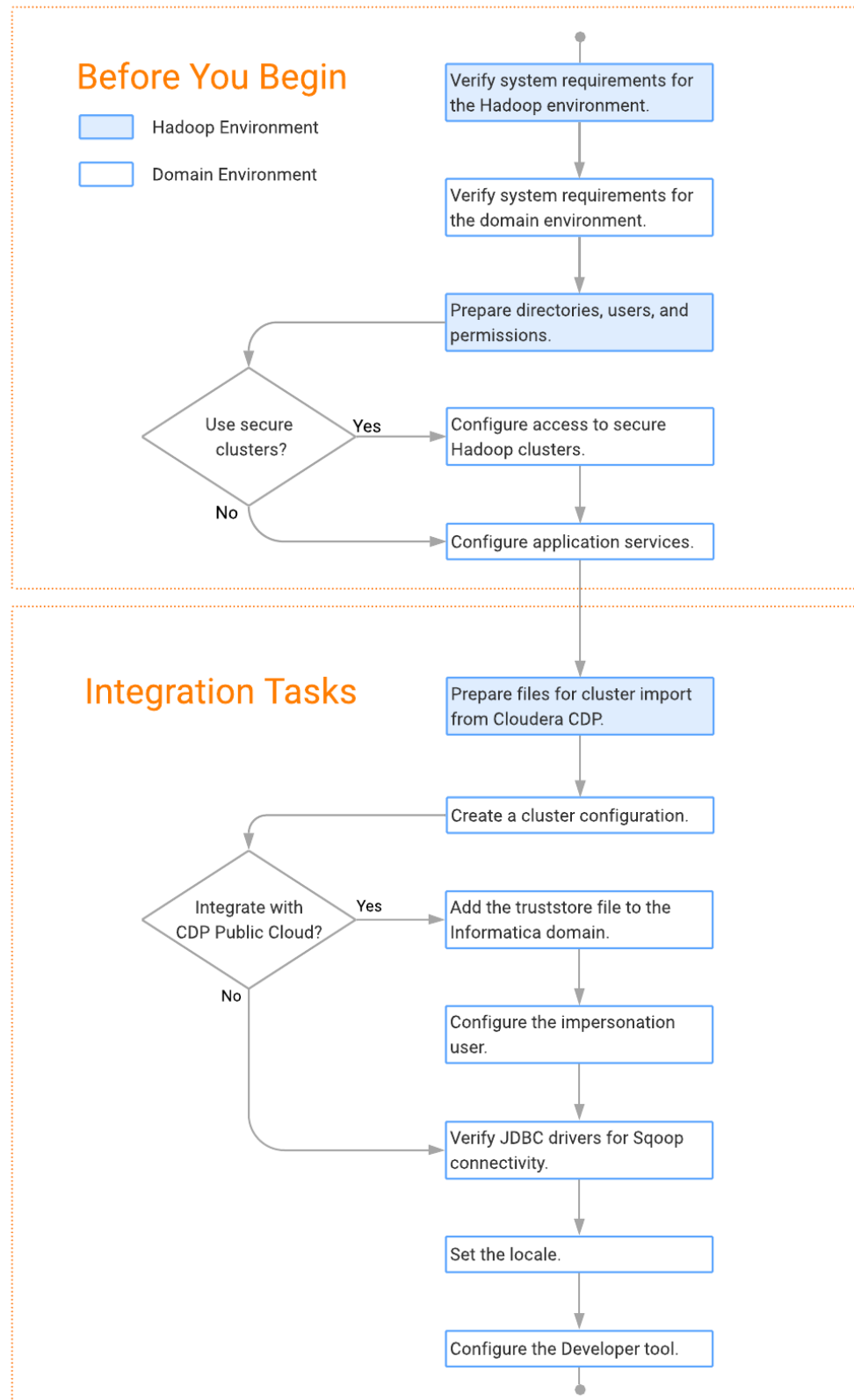
Depending on whether you want to integrate or upgrade Data Engineering Integration in a Cloudera CDP environment, you can use the flow charts to perform the following tasks:

- Integrate the Informatica domain with Cloudera CDP for the first time.
- Upgrade from version 10.4.1.

All tasks apply to both CDP Private Cloud and CDP Public Cloud unless otherwise noted.

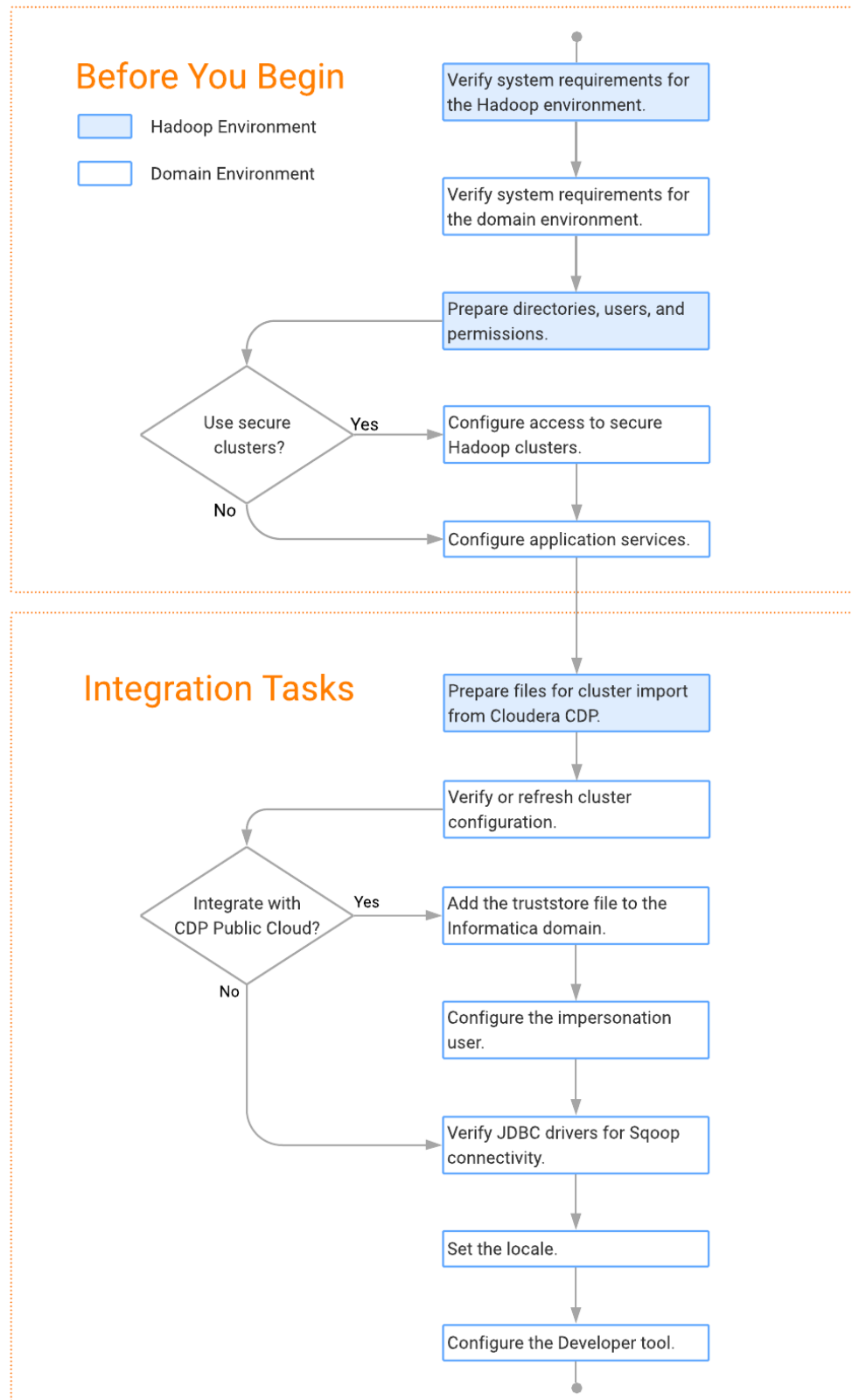
# Task Flow to Integrate with Cloudera CDP

The following diagram shows the task flow to integrate the Informatica domain with Cloudera CDP:



## Task Flow to Upgrade from Version 10.4.1

The following diagram shows the task flow to upgrade version 10.4.1 for Cloudera CDP:



# Prepare for Cluster Import from Cloudera CDP

Before the Informatica administrator can import cluster information to create a cluster configuration in the Informatica domain, the Hadoop administrator must perform some preliminary tasks.

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

**Note:** If you are upgrading from a previous version, verify the properties and suggested values, as Data Engineering Integration might require additional properties or different values for existing properties.

Complete the following tasks to prepare the cluster before the Informatica administrator creates the cluster configuration:

1. Verify property values in \*-site.xml files that Data Engineering Integration needs to run mappings in the Hadoop environment.
2. Provide information to the Informatica administrator that is required to import cluster information into the domain. Depending on the method of import, perform one of the following tasks:
  - To import directly from the cluster, give the Informatica administrator cluster authentication information to connect to the cluster. To integrate with a CDP Public Cloud cluster, provide information for the Data Hub cluster.
  - To import from an archive file, export cluster information and provide an archive file to the Data Engineering Integration administrator.
3. If you plan to use an Auto-TLS enabled CDP cluster, import security certificates to the Informatica domain before you import cluster information into the domain.  
Refer to [“Import Security Certificates from an SSL-Enabled Cluster” on page 32](#).

## Configure \*-site.xml Files for Cloudera CDP

The Hadoop administrator needs to configure \*-site.xml file properties and restart impacted services before the Informatica administrator imports cluster information into the domain.

### core-site.xml

Configure the following properties in the core-site.xml file:

#### **fs.s3.enableServerSideEncryption**

Enables server side encryption for S3 buckets. Required for SSE and SSE-KMS encryption.

Set to: TRUE

#### **fs.s3a.access.key**

The ID for the Blaze and Spark engines to connect to the Amazon S3 file system.

Set to your access key.

#### **fs.s3a.secret.key**

The password for the Blaze and Spark engines to connect to the Amazon S3 file system

Set to your access ID.

### **fs.s3a.server-side-encryption-algorithm**

The server-side encryption algorithm for S3. Required for SSE and SSE-KMS encryption. Set to the encryption algorithm used.

### **hadoop.proxyuser.<proxy user>.groups**

Defines the groups that the proxy user account can impersonate. On a secure cluster the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to group names of impersonation users separated by commas. If less security is preferred, use the wildcard " \* " to allow impersonation from any group.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

### **hadoop.proxyuser.<proxy user>.hosts**

Defines the host machines that a user account can impersonate. On a secure cluster the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to a single host name or IP address, or set to a comma-separated list. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

### **hadoop.proxyuser.hive.hosts**

Comma-separated list of hosts that you want to allow the Hive user to impersonate on a non-secure cluster.

When `hive.server2.enable.doAs` is false, append a comma-separated list of Informatica server host names or IP address where the Data Integration Service is running. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to this property, restart the cluster services that use core-site configuration values.

### **io.compression.codecs**

Enables compression on temporary staging tables.

Set to a comma-separated list of compression codec classes on the cluster.

### **hadoop.security.auth\_to\_local**

Translates the principal names from the Active Directory and MIT realm into local names within the Hadoop cluster. Based on the Hadoop cluster used, you can set multiple rules.

Set to: `RULE:[1:$1@$0](^.*@YOUR.REALM)s/^(.*)@YOUR.REALM\$.COM$/1/g`

Set to: `RULE:[2:$1@$0](^.*@YOUR.REALM\.$)s/^(.*)@YOUR.REALM\$.COM$/1/g`

### **fair-scheduler.xml**

Configure the following properties in the fair-scheduler.xml file:

#### **allowPreemptionFrom**

Enables preemption for the Fair Scheduler. The Blaze engine does not support preemption. If YARN preemption is enabled for the cluster, you need to disable it for the queue allocated to the Blaze engine.

Set to FALSE for the queue allocated to the Blaze engine.

For example:

```
<queue name="Blaze">
  <weight>1.0</weight>
  <allowPreemptionFrom>false</allowPreemptionFrom>
  <schedulingPolicy>fsp</schedulingPolicy>
  <aclSubmitApps>*</aclSubmitApps>
  <aclAdministerApps>*</aclAdministerApps>
</queue>
```

### hbase-site.xml

Configure the following properties in the hbase-site.xml file:

#### **zookeeper.znode.parent**

Identifies HBase master and region servers.

Set to the relative path to the znode directory of HBase.

### hdfs-site.xml

Configure the following properties in the hdfs-site.xml file:

#### **dfs.encryption.key.provider.uri**

The KeyProvider used to interact with encryption keys when reading and writing to an encryption zone. Required if sources or targets reside in the HDFS encrypted zone on Java KeyStore KMS-enabled Cloudera CDH cluster or a Ranger KMS-enabled Hortonworks HDP cluster.

Set to: kmf://http@xx11.xyz.com:16000/kms

### hive-site.xml

Configure the following properties in the hive-site.xml file:

#### **hive.cluster.delegation.token.store.class**

The token store implementation. Required for HiveServer2 high availability and load balancing.

Set to: org.apache.hadoop.hive.thrift.DBTokenStore

#### **hive.exec.dynamic.partition**

Enables dynamic partitioned tables for Hive tables. Applicable for Hive versions 0.9 and earlier.

Set to: TRUE

#### **hive.exec.dynamic.partition.mode**

Allows all partitions to be dynamic. Required if you use Sqoop and define a DDL query to create or replace a partitioned Hive target at run time.

Set to: nonstrict

#### **hive.server2.zookeeper.namespace**

The value of the ZooKeeper namespace in the JDBC connection string. Required for HiveServer2 high availability.

Set to: jdbc:hive2://<zookeeper\_ensemble>/  
default;serviceDiscoveryMode=zooKeeper;zooKeeperNamespace=hiveserver2

#### **hive.zookeeper.quorum**

Comma-separated list of ZooKeeper server host:ports in a cluster. The value of the ZooKeeper ensemble in the JDBC connection string. Required for HiveServer2 high availability.

Set to: jdbc:hive2://<zookeeper\_ensemble>/default;serviceDiscoveryMode=zooKeeper;

### mapred-site.xml

Configure the following properties in the mapred-site.xml file:

**mapreduce.framework.name**

The run-time framework to run MapReduce jobs. Values can be local, classic, or yarn. Required for Sqoop.

Set to: yarn

**mapreduce.jobhistory.address**

Location of the MapReduce JobHistory Server. The default port is 10020. Required for Sqoop.

Set to: <MapReduce JobHistory Server>:<port>

**mapreduce.jobhistory.intermediate-done-dir**

Directory where MapReduce jobs write history files. Required for Sqoop.

Set to: /mr-history/tmp

**mapreduce.jobhistory.done-dir**

Directory where the MapReduce JobHistory Server manages history files. Required for Sqoop.

Set to: /mr-history/done

**mapreduce.jobhistory.principal**

The Service Principal Name for the MapReduce JobHistory Server. Required for Sqoop.

Set to: mapred/\_HOST@YOUR-REALM

**mapreduce.jobhistory.webapp.address**

Web address of the MapReduce JobHistory Server. The default value is 19888. Required for Sqoop.

Set to: <host>:<port>

**yarn.app.mapreduce.am.staging-dir**

The HDFS staging directory used while submitting jobs.

Set to the staging directory path.

[tez-site.xml](#)

Configure the following properties in the tez-site.xml file:

**tez.runtime.io.sort.mb**

The sort buffer memory. Required when the output needs to be sorted for Blaze and Spark engines.

Set value to 270 MB.

[yarn-site.xml](#)

Configure the following properties in the yarn-site.xml file:

**yarn.application.classpath**

Required for dynamic resource allocation.

Add spark\_shuffle.jar to the class path. The .jar file must contain the class "org.apache.spark.network.yarn.YarnShuffleService."

**yarn.nodemanager.resource.memory-mb**

The maximum RAM available for each container. Set the maximum memory on the cluster to increase resource memory available to the Blaze engine.

Set to 16 GB if value is less than 16 GB.

**yarn.nodemanager.resource.cpu-vcores**

The number of virtual cores for each container. Required for Blaze engine resource allocation.

Set to 10 if the value is less than 10.

**yarn.scheduler.minimum-allocation-mb**

The minimum RAM available for each container. Required for Blaze engine resource allocation.

Set to 6 GB if the value is less than 6 GB.

**yarn.nodemanager.vmem-check-enabled**

Disables virtual memory limits for containers. Required for the Blaze and Spark engines.

Set to: false

**yarn.nodemanager.aux-services**

Required for dynamic resource allocation for the Spark engine.

Add an entry for "spark\_shuffle."

**yarn.nodemanager.aux-services.spark\_shuffle.class**

Required for dynamic resource allocation for the Spark engine.

Set to: org.apache.spark.network.yarn.YarnShuffleService

**yarn.resourcemanager.scheduler.class**

Defines the YARN scheduler that the Data Integration Service uses to assign resources.

Set to: org.apache.hadoop.yarn.server.resourcemanager.scheduler

## Prepare for Direct Import from Cloudera CDP

If you plan to provide direct access to the Informatica administrator to import cluster information, provide the required connection information.

The following table describes the information that you need to provide to the Informatica administrator to create the cluster configuration directly from the cluster:

Property	Description
Host	Host of the cluster manager.
Port	Port of the cluster manager.
User ID	Cluster user ID.
Password	Password for the user.
Cluster name	Name of the cluster. Use the display name if the cluster manager manages multiple clusters. If you do not provide a cluster name, the wizard imports information based on the default cluster. To find the correct Cloudera cluster name when you have multiple clusters, perform the following steps: 1. Log in to Cloudera Manager adding the following string to the URL: /api/v8/clusters 2. Provide the Informatica Administrator the cluster property name that appears in the browser tab.

**Note:** For CDP Public Cloud, use the host and port information of the Data Hub cluster.



## Prepare the Archive File for Import from Cloudera CDP

If you plan to provide an archive file for the Informatica administrator, ensure that you include all required site-\*.xml files.

Create a .zip or .tar file that contains the following \*-site.xml files:

- core-site.xml
- hbase-site.xml. Required only for access to HBase sources and targets.
- hdfs-site.xml
- hive-site.xml
- mapred-site.xml
- tez-site.xml
- yarn-site.xml

**Note:** When you configure a CDP Public Cloud cluster, the hbase-site.xml file is on the Data Lake cluster. The other files are on the Data Hub cluster.

Give the Informatica administrator access to the archive file to import the cluster information into the domain.

## Create a Cluster Configuration

After the Hadoop administrator prepares the cluster for import, the Informatica administrator must create a cluster configuration.

### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

A cluster configuration is an object in the domain that contains configuration information about the Hadoop cluster. The cluster configuration enables the Data Integration Service to push mapping logic to the Hadoop environment. Import configuration properties from the Hadoop cluster to create a cluster configuration.

The import process imports values from \*-site.xml files into configuration sets based on the individual \*-site.xml files. When you perform the import, the cluster configuration wizard can create Hadoop, HBase, HDFS, and Hive connection to access the Hadoop environment. If you choose to create the connections, the wizard also associates the cluster configuration with the connections.

**Note:** If you are integrating for the first time and you imported the cluster configuration when you ran the installer, you *must* re-create or refresh the cluster configuration.

## Before You Import

Before you can import the cluster configuration, you must get information from the Hadoop administrator based on the method of import.

If you import directly from the cluster, contact the Hadoop administrator to get cluster connection information. If you import from a file, get an archive file of exported cluster information.

## Importing a Hadoop Cluster Configuration from the Cluster

When you import the Hadoop cluster configuration directly from the cluster, you provide information to connect to the cluster.

Get cluster connection information from the Hadoop administrator.

1. From the **Connections** tab, click the **ClusterConfigurations** node in the Domain Navigator.
2. From the Actions menu, select **New > Cluster Configuration**.

The **Cluster Configuration** wizard opens.

3. Configure the following General properties:

Property	Description
Cluster configuration name	Name of the cluster configuration.
Description	Optional description of the cluster configuration.
Distribution type	The cluster Hadoop distribution type.
Distribution version	<p>Version of the Hadoop distribution.</p> <p>Each distribution type has a default version. The default version is the latest version of the Hadoop distribution that Data Engineering Integration supports.</p> <p><b>Note:</b> When the cluster version differs from the default version and Informatica supports more than one version, the cluster configuration import process populates the property with the most recent supported version. For example, consider the case where Informatica supports versions 5.10 and 5.13, and the cluster version is 5.12. In this case, the cluster configuration import process populates this property with 5.10, because 5.10 is the most recent supported version before 5.12.</p> <p>You can edit the property to choose any supported version. Restart the Data Integration Service for the changes to take effect.</p>
Method to import the cluster configuration	Choose <b>Import from cluster</b> .
Create connections	<p>Choose to create Hadoop, HDFS, Hive, and HBase connections.</p> <p>If you choose to create connections, the Cluster Configuration wizard associates the cluster configuration with each connection that it creates.</p> <p>The Hadoop connection contains default values for properties such as cluster environment variables, cluster path variables, and advanced properties. Based on the cluster environment and the functionality that you use, you can add to the default values or change the default values of these properties. For a list of Hadoop connection properties to configure, see <a href="#">"Configuring Hadoop Connection Properties" on page 269</a>.</p> <p>If you do not choose to create connections, you must manually create them and associate the cluster configuration with them.</p> <p><b>Important:</b> When the wizard creates the Hive connection, it populates the Metadata Connection String and the Data Access Connection String properties with the value from the hive.metastore.uris property. If the Hive metastore and HiveServer2 are running on different nodes, you must update the Metadata Connection String to point to the HiveServer2 host.</p>

The cluster properties appear.

4. Configure the following properties. For CDP Public Cloud, use the host and port information of the Data Hub cluster.

Property	Description
Host	Host name or IP address of the cluster manager.
Port	Port of the cluster manager.
User ID	Cluster user ID.
Password	Password for the user.
Cluster name	Name of the cluster. Use the display name if the cluster manager manages multiple clusters. If you do not provide a cluster name, the wizard imports information based on the default cluster.

5. Click **Next** and verify the cluster configuration information on the summary page.

## Importing a Hadoop Cluster Configuration from a File

You can import properties from an archive file to create a cluster configuration.

Before you import from the cluster, you must get the archive file from the Hadoop administrator.

1. From the **Connections** tab, click the **ClusterConfigurations** node in the Domain Navigator.
2. From the Actions menu, select **New > Cluster Configuration**.

The **Cluster Configuration** wizard opens.

3. Configure the following properties:

Property	Description
Cluster configuration name	Name of the cluster configuration.
Description	Optional description of the cluster configuration.
Distribution type	The cluster Hadoop distribution type.

Property	Description
Method to import the cluster configuration	Choose <b>Import from file</b> to import properties from an archive file.
Create connections	<p>Choose to create Hadoop, HDFS, Hive, and HBase connections.</p> <p>If you choose to create connections, the Cluster Configuration wizard associates the cluster configuration with each connection that it creates.</p> <p>The Hadoop connection contains default values for properties such as cluster environment variables, cluster path variables, and advanced properties. Based on the cluster environment and the functionality that you use, you can add to the default values or change the default values of these properties. For a list of Hadoop connection properties to configure, see <a href="#">“Configuring Hadoop Connection Properties” on page 269</a>.</p> <p>If you do not choose to create connections, you must manually create them and associate the cluster configuration with them.</p> <p><b>Important:</b> When the wizard creates the Hive connection, it populates the Metadata Connection String and the Data Access Connection String properties with the value from the <code>hive.metastore.uris</code> property. If the Hive metastore and HiveServer2 are running on different nodes, you must update the Metadata Connection String to point to the HiveServer2 host.</p>

- Click **Browse** to select a file. Select the file and click **Open**.
- Click **Next** and verify the cluster configuration information on the summary page.

## Copy the Truststore File to the Informatica Domain

After you import a CDP Public Cloud cluster configuration, copy the truststore file `cm-auto-global_truststore.jks` from the cluster to the Informatica domain.

### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

- Find the value for the property `ssl.client.truststore.location` in the following file on the cluster: `/etc/hadoop/conf/ssl-client.xml`  
 The value of this property is the file path for the file `cm-auto-global-truststore.jks`. For example, `/var/lib.cloudera-scm-agent/agent-cert/cm-auto-global_truststore.jks`
- Create the same directory structure on the Informatica domain.  
 For example, `<Informatica installation directory>/var/lib.cloudera-scm-agent/agent-cert/`
- Find the `.jks` file on the cluster in the file path you found in step 1 and copy the file to the directory you created on the Informatica domain in step 2.

# Configure the Impersonation User for Operating System Profiles

For Cloudera CDP Public Cloud, if you configure the Data Integration Service to use operating system profiles, you must configure the Hadoop impersonation user before the Spark engine can run mappings.

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

For information about configuring the impersonation user, see ["Verify and Create Users" on page 23](#).

## Verify JDBC Drivers for Sqoop Connectivity

Verify that you have the JDBC drivers to access JDBC-compliant databases in the Hadoop environment. You might need separate drivers for metadata import and for run-time processing.

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

You download drivers based on design-time and run-time requirements:

- **Design-time.** To import metadata, you can use the DataDirect drivers packaged with the Informatica installer if they are available. If they are not available, use any Type 4 JDBC driver that the database vendor recommends.
- **Run-time.** To run mappings, use any Type 4 JDBC driver that the database vendor recommends. Some distributions support other drivers to use Sqoop connectors. You cannot use the DataDirect drivers for run-time processing.

## Verify Design-time Drivers

Use the DataDirect JDBC drivers packaged with the Informatica installer to import metadata from JDBC-compliant databases. If the DataDirect JDBC drivers are not available for a specific JDBC-compliant database, download the Type 4 JDBC driver associated with that database.

Copy the JDBC driver .jar files to the following location on the Developer tool machine:

```
<Informatica installation directory>\clients\externaljdbcjars
```

## Verify Run-time Drivers

Verify run-time drivers for mappings that access JDBC-compliant databases in the Hadoop environment. Use any Type 4 JDBC driver that the database vendor recommends.

1. Download Type 4 JDBC drivers associated with the JCBC-compliant databases that you want to access.

2. To use Sqoop TDCH Cloudera Connector Powered by Teradata, download all the .jar files in the Cloudera Connector Powered by Teradata package from the following location:  
<http://www.cloudera.com/downloads.html>.

The package has the following naming convention: sqoop-connector-teradata-<version>.tar

To use the Cloudera CDP version 7.x, you must download the sqoop-connector-teradata-1.8.0c7.jar file.

3. Append /opt/cloudera/parcels/CDH/lib/hive/lib/hive-exec-{version}.jar to the mapreduce.application.classpath property in the mapred-site.xml file.
4. To optimize the Sqoop mapping performance on the Spark engine while writing data to an HDFS complex file target of the Parquet format, download the following .jar files:

File name	Location
parquet-hadoop-bundle-1.6.0.jar	<a href="https://repo1.maven.org/maven2/com/twitter/parquet-hadoop-bundle/1.6.0/">https://repo1.maven.org/maven2/com/twitter/parquet-hadoop-bundle/1.6.0/</a>
parquet-avro-1.6.0.jar	<a href="https://repo1.maven.org/maven2/com/twitter/parquet-avro/1.6.0/">https://repo1.maven.org/maven2/com/twitter/parquet-avro/1.6.0/</a>
parquet-column-1.5.0.jar	<a href="https://repo1.maven.org/maven2/com/twitter/parquet-column/1.5.0/">https://repo1.maven.org/maven2/com/twitter/parquet-column/1.5.0/</a>

5. Copy all of the .jar files to the following directory on the machine where the Data Integration Service runs:

```
<Informatica installation directory>\externaljdbcjars
```

Changes take effect after you recycle the Data Integration Service. At run time, the Data Integration Service copies the .jar files to the Hadoop distribution cache so that the .jar files are accessible to all nodes in the cluster.

## Set the Locale for Cloudera CDP

If you want to process data that contains non-ASCII characters, you must integrate the locale setting on Data Engineering Integration with the locale setting on the cluster.

### Perform this task in the following situations:

- You are integrating for the first time.

To integrate the locale setting, complete the following tasks:

1. In the Hadoop connection, navigate to **Hadoop Cluster Properties**. As the value for the property **Cluster Environment Variables**, configure the locale environment variables, such as the LANG or LC\_ALL environment variable.  
  
The locale setting in the Hadoop connection must match the locale setting that is configured in the domain. To view the locale environment variable values set in the domain, run the following command on any node in the cluster: `locale`
2. In Cloudera Manager, add the environment variables to the following YARN property:  

```
yarn.nodemanager.env-whitelist
```

# Enable Data Preparation of JSON Files on Cloudera CDP

If you integrate Enterprise Data Preparation with a Cloudera CDP Hadoop cluster, you must specify the location of a Hive .jar file in the Hive Auxiliary JARs Directory in CDP to enable data preparation of JSON files.

## Perform this task in the following situations:

- You are integrating for the first time.

1. Search for hive-hcatalog-core.jar in the Cloudera CDP installation directory.  
You can generally find the .jar file in the following directory:  
`/opt/cloudera/parcels/CDH/jars/`
2. Log in to Cloudera Manager.
3. Select **Hive** in the cluster.
4. Click the **Configuration** tab, then select the **Advanced** category.
5. Enter the path to the directory containing the .jar file in the Hive Auxiliary JARs Directory property.  
If the file is in the location noted, the path to the directory is:  
`/opt/cloudera/parcels/CDH/jars/`
6. Restart the Hive server.

## Configure the Developer Tool

You can configure the Developer tool to enable you to import complex files or import metadata when the domain is Kerberos-enabled.

Edit the developerCore.ini file to import complex files. Edit the file on each Developer tool machine.

### Configure developerCore.ini

Edit the developerCore.ini file to import complex files.

Edit the developerCore.ini file on each machine that hosts the Developer tool.

When you import a complex file, such as Avro or Parquet, the imported object includes metadata associated with the distribution in the Hadoop environment. If the file resides on the Developer tool machine, the import process picks up the distribution information from the developerCore.ini file. You must edit the developerCore.ini file to point to the distribution directory on the Developer tool machine.

You can find the developerCore.ini file in the following directory:

```
<Informatica installation directory>\clients\DeveloperClient
```

Add the following property:

```
-DINFA_HADOOP_DIST_DIR=hadoop<distribution>_<version>
```

The change takes effect when you restart the Developer tool.

## Configure the Developer Tool for Kerberos

To import metadata from Hive, HBase, and complex file sources, import configuration files from the Kerberos-enabled cluster, and generate the Kerberos credentials file on the Developer tool machine.

### Import configuration files

The Hadoop cluster uses a set of XML files named `*-site.xml` to store configuration settings. The domain uses the same set of files to create the cluster configuration object.

To enable you to import metadata from the cluster, import the `*-site.xml` files to each Developer tool machine:

1. Log in to the Administrator tool and navigate to **Connections > Cluster Configuration > CCO**. Locate the cluster configuration associated with the Hadoop cluster.
2. Extract the `*-site.xml` files in the cluster configuration, including sensitive properties, to the following directory on the Developer tool machine: `<Informatica installation directory>\clients\DeveloperClient\hadoop\<Hadoop distribution>\conf`  
For more information about sensitive properties, see the "Active Properties View" topic in the *Data Engineering Administrator Guide*.

**Note:** If you refresh the cluster configuration, repeat these steps.

### Generate the Kerberos credentials file

1. Copy the `krb5.conf` file from `<Developer tool installation directory>/services/shared/security` to `C:/Windows`.
2. Rename `krb5.conf` to `krb5.ini`.
3. In the `krb5.ini` file, verify the value of the `forwardable` option to determine how to use the `kinit` command. If `forwardable=true`, run the command with the `-f` option. Otherwise, run the command without the `-f` option.
4. To generate the Kerberos credentials file, run the `kinit` command from the following location: `<Developer tool installation directory>/clients/java/bin/kinit.exe`  
For example, you might run the following command: `kinit joe/domain12345@MY-REALM`



## CHAPTER 7

# Google Dataproc Integration Tasks

This chapter includes the following topics:

- [Google Dataproc Task Flows, 129](#)
- [Prepare for Cluster Import from Google Dataproc, 132](#)
- [Create a Cluster Configuration , 138](#)
- [Configure the Cluster for the Blaze Engine, 139](#)
- [Configure Domain Settings, 140](#)
- [Complete Upgrade Tasks, 141](#)
- [Configure the Developer Tool , 142](#)

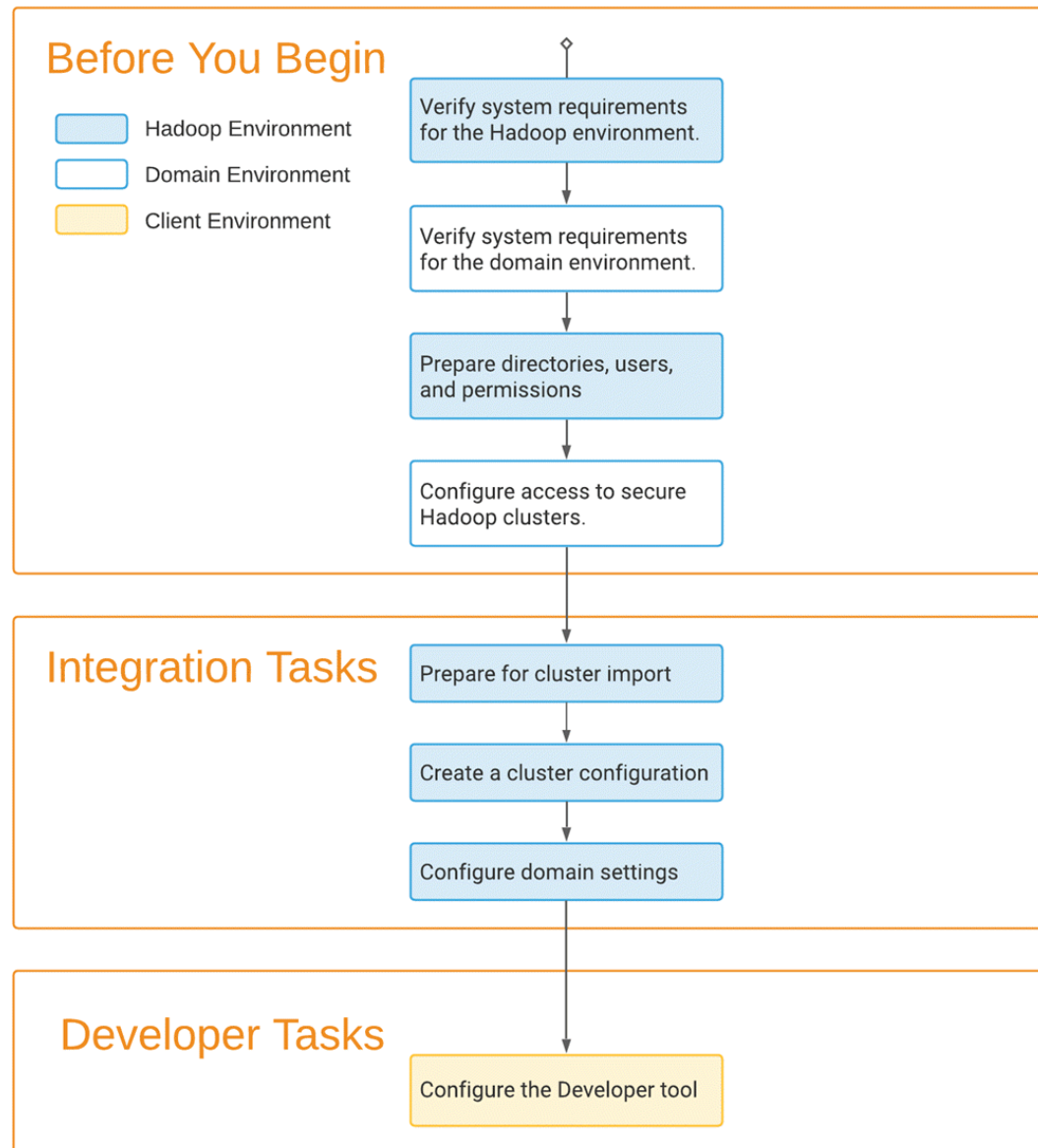
## Google Dataproc Task Flows

Depending on whether you want to integrate or upgrade Data Engineering Integration in a Google Dataproc environment, you can use the flow charts to perform the following tasks:

- Integrate the Informatica domain with Google Dataproc for the first time.
- Upgrade from version 10.2.2.

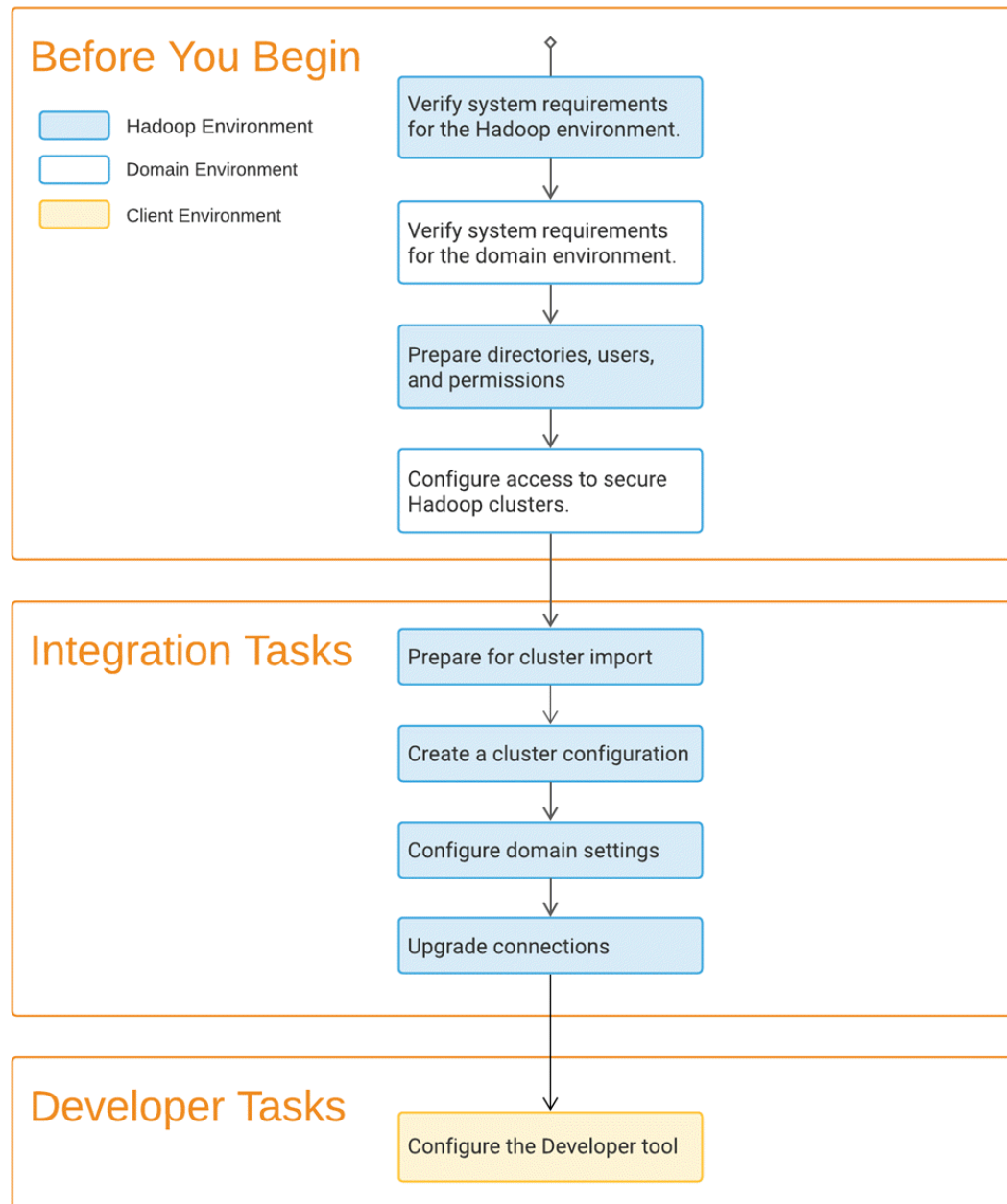
## Task Flow to Integrate Versions 10.4 or 10.5

The following image shows the task flow to integrate Google Dataproc with Informatica 10.4 or 10.5:



## Task Flow to Upgrade from Version 10.2.2

The following image shows the task flow to upgrade version 10.2.2 for Google Dataproc:



# Prepare for Cluster Import from Google Dataproc

Before the Informatica administrator can import cluster information to create a cluster configuration in the Informatica domain, the Hadoop administrator must perform some preliminary tasks.

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

Complete the following tasks to prepare the cluster before the Informatica administrator creates the cluster configuration:

1. Configure \*-site.xml files on the cluster.
2. Prepare an archive file of cluster configuration settings.
3. Generate keytab files for the Service Principal Name user.

## Configure \*-site.xml Files for Google Dataproc

The Hadoop administrator needs to configure \*-site.xml file properties and restart impacted services before the Informatica administrator imports cluster information into the domain.

### core-site.xml

Configure the following properties in the core-site.xml file:

#### **fs.s3.enableServerSideEncryption**

Enables server side encryption for S3 buckets. Required for SSE and SSE-KMS encryption.

Set to: TRUE

#### **fs.s3a.access.key**

The ID for the Blaze and Spark engines to connect to the Amazon S3 file system.

Set to your access key.

#### **fs.s3a.secret.key**

The password for the Blaze and Spark engines to connect to the Amazon S3 file system

Set to your access ID.

#### **fs.s3a.server-side-encryption-algorithm**

The server-side encryption algorithm for S3. Required for SSE and SSE-KMS encryption. Set to the encryption algorithm used.

#### **hadoop.proxyuser.<proxy user>.groups**

Defines the groups that the proxy user account can impersonate. On a secure cluster the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to group names of impersonation users separated by commas. If less security is preferred, use the wildcard " \* " to allow impersonation from any group.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.proxyuser.<proxy user>.hosts**

Defines the host machines that a user account can impersonate. On a secure cluster the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to a single host name or IP address, or set to a comma-separated list. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.proxyuser.hive.hosts**

Comma-separated list of hosts that you want to allow the Hive user to impersonate on a non-secure cluster.

When `hive.server2.enable.doAs` is false, append a comma-separated list of Informatica server host names or IP address where the Data Integration Service is running. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to this property, restart the cluster services that use core-site configuration values.

#### **io.compression.codecs**

Enables compression on temporary staging tables.

Set to a comma-separated list of compression codec classes on the cluster.

#### **hadoop.security.auth\_to\_local**

Translates the principal names from the Active Directory and MIT realm into local names within the Hadoop cluster. Based on the Hadoop cluster used, you can set multiple rules.

Set to: `RULE:[1:$1@$0](^.*@YOUR.REALM)s/^(.*)@YOUR.REALM\.COM$/1/g`

Set to: `RULE:[2:$1@$0](^.*@YOUR.REALM\.)s/^(.*)@YOUR.REALM\.COM$/1/g`

#### [fair-scheduler.xml](#)

Configure the following properties in the fair-scheduler.xml file:

##### **allowPreemptionFrom**

Enables preemption for the Fair Scheduler. The Blaze engine does not support preemption. If YARN preemption is enabled for the cluster, you need to disable it for the queue allocated to the Blaze engine.

Set to FALSE for the queue allocated to the Blaze engine.

For example:

```
<queue name="Blaze">
  <weight>1.0</weight>
  <allowPreemptionFrom>false</allowPreemptionFrom>
  <schedulingPolicy>fsp</schedulingPolicy>
  <aclSubmitApps>*</aclSubmitApps>
  <aclAdministerApps>*</aclAdministerApps>
</queue>
```

#### [hbase-site.xml](#)

Configure the following properties in the hbase-site.xml file:

##### **zookeeper.znode.parent**

Identifies HBase master and region servers.

Set to the relative path to the znode directory of HBase.

## hdfs-site.xml

Configure the following properties in the hdfs-site.xml file:

### **dfs.encrypted.key.provider.uri**

The KeyProvider used to interact with encryption keys when reading and writing to an encryption zone. Required if sources or targets reside in the HDFS encrypted zone on Java KeyStore KMS-enabled Cloudera CDH cluster or a Ranger KMS-enabled Hortonworks HDP cluster.

Set to: kmf://http@xx11.xyz.com:16000/kms

### **dfs.namenode.rpc-bind-host**

The actual address the Remote Procedure Call (RPC) server will bind to. If this optional address is set, it overrides the hostname portion of dfs.namenode.rpc-address. Enables the cluster to listen on private and public network interfaces, allowing remote access and datanode access. Required when you run mappings on a non-VPN Dataproc cluster.

Set to: 0.0.0.0 to enable the cluster to listen on private and public ports, allowing remote access and datanode access.

### **dfs.namenode.servicerpc-bind-host**

The actual address the Remote Procedure Call (RPC) server will bind to. If this optional address is set, it overrides the hostname portion of dfs.namenode.rpc-address. Enables the cluster to listen on private and public network interfaces, allowing remote access and datanode access. Required when you run mappings on a non-VPN Dataproc cluster.

Set to: 0.0.0.0 to enable the cluster to listen on private and public ports, allowing remote access and datanode access.

### **dfs.namenode.http-bind-host**

The actual address the Remote Procedure Call (RPC) server will bind to. If this optional address is set, it overrides the hostname portion of dfs.namenode.rpc-address. Enables the cluster to listen on private and public network interfaces, allowing remote access and datanode access. Required when you run mappings on a non-VPN Dataproc cluster.

Set to: 0.0.0.0 to enable the cluster to listen on private and public ports, allowing remote access and datanode access.

### **dfs.namenode.https-bind-host**

The actual address the Remote Procedure Call (RPC) server will bind to. If this optional address is set, it overrides the hostname portion of dfs.namenode.rpc-address. Enables the cluster to listen on private and public network interfaces, allowing remote access and datanode access. Required when you run mappings on a non-VPN Dataproc cluster.

Set to: 0.0.0.0 to enable the cluster to listen on private and public ports, allowing remote access and datanode access.

## hive-site.xml

Configure the following properties in the hive-site.xml file:

### **hive.async.log.enabled**

Enables asynchronous logging. Required when you enable pre-task and post-task monitoring statistics on a Dataproc cluster.

Set to: FALSE

### **hive.cluster.delegation.token.store.class**

The token store implementation. Required for HiveServer2 high availability and load balancing.

Set to: org.apache.hadoop.hive.thrift.DBTokenStore

**hive.exec.dynamic.partition**

Enables dynamic partitioned tables for Hive tables. Applicable for Hive versions 0.9 and earlier.

Set to: TRUE

**hive.exec.dynamic.partition.mode**

Allows all partitions to be dynamic. Required if you use Sqoop and define a DDL query to create or replace a partitioned Hive target at run time.

Set to: nonstrict

**hive.server2.in.place.progress**

Allows HiveServer2 to send progress bar update information. Takes effect only when you enable Tez. Required when you enable pre-task and post-task monitoring statistics on a Dataproc cluster.

Set to: TRUE

**hive.server2.logging.operation.level**

Hive Server2 logging level at the session level. Requires hive.server2.logging.operation.enabled to be set to TRUE. Required when you enable pre-task and post-task monitoring statistics on a Dataproc cluster.

Set to: EXECUTION

**hive.server2.logging.operation.enabled**

Enables logs to be saved. Required when you enable pre-task and post-task monitoring statistics on a Dataproc cluster.

Set to: TRUE

**hive.server2.zookeeper.namespace**

The value of the ZooKeeper namespace in the JDBC connection string. Required for HiveServer2 high availability.

Set to: jdbc:hive2://<zookeeper\_ensemble>/  
default;serviceDiscoveryMode=zooKeeper;zooKeeperNamespace=hiveserver2

**hive.zookeeper.quorum**

Comma-separated list of ZooKeeper server host:ports in a cluster. The value of the ZooKeeper ensemble in the JDBC connection string. Required for HiveServer2 high availability.

Set to: jdbc:hive2://<zookeeper\_ensemble>/default;serviceDiscoveryMode=zooKeeper;

[mapred-site.xml](#)

Configure the following properties in the mapred-site.xml file:

**mapreduce.framework.name**

The run-time framework to run MapReduce jobs. Values can be local, classic, or yarn. Required for Sqoop.

Set to: yarn

**mapreduce.jobhistory.address**

Location of the MapReduce JobHistory Server. The default port is 10020. Required for Sqoop.

Set to: <MapReduce JobHistory Server>:<port>

**mapreduce.jobhistory.intermediate-done-dir**

Directory where MapReduce jobs write history files. Required for Sqoop.

Set to: /mr-history/tmp

**mapreduce.jobhistory.done-dir**

Directory where the MapReduce JobHistory Server manages history files. Required for Sqoop.

Set to: /mr-history/done

**mapreduce.jobhistory.principal**

The Service Principal Name for the MapReduce JobHistory Server. Required for Sqoop.

Set to: mapred/\_HOST@YOUR-REALM

**mapreduce.jobhistory.webapp.address**

Web address of the MapReduce JobHistory Server. The default value is 19888. Required for Sqoop.

Set to: <host>:<port>

**yarn.app.mapreduce.am.staging-dir**

The HDFS staging directory used while submitting jobs.

Set to the staging directory path.

[tez-site.xml](#)

Configure the following properties in the tez-site.xml file:

**tez.am.tez-ui.history-url.template**

Tez UI URL template for the application. The application manager uses this URL to redirect the user to the Tez UI. Required when you enable pre-task and post-task monitoring statistics on a Dataproc cluster.

Set value to:

`_HISTORY_URL_BASE?%2F%23%2Ftez-app%2FAPPLICATION_ID`

**Note:** The values of `_HISTORY_URL_BASE_` and `_APPLICATION_ID` are resolved at runtime. Do not edit the string to supply values.

**tez.runtime.io.sort.mb**

The sort buffer memory. Required when the output needs to be sorted for Blaze and Spark engines.

Set value to 270 MB.

**tez.task.generate.counters.per.io**

Enables pre-task and post-task monitoring statistics on an Amazon EMR or Dataproc cluster.

Set to: TRUE

[yarn-site.xml](#)

Configure the following properties in the yarn-site.xml file:

**yarn.application.classpath**

Required for dynamic resource allocation.

Add `spark_shuffle.jar` to the class path. The `.jar` file must contain the class `"org.apache.spark.network.yarn.YarnShuffleService."`

**yarn.nodemanager.resource.memory-mb**

The maximum RAM available for each container. Set the maximum memory on the cluster to increase resource memory available to the Blaze engine.

Set to 16 GB if value is less than 16 GB.

**yarn.nodemanager.resource.cpu-vcores**

The number of virtual cores for each container. Required for Blaze engine resource allocation.



Set to 10 if the value is less than 10.

**yarn.scheduler.minimum-allocation-mb**

The minimum RAM available for each container. Required for Blaze engine resource allocation.

Set to 6 GB if the value is less than 6 GB.

**yarn.nodemanager.vmem-check-enabled**

Disables virtual memory limits for containers. Required for the Blaze and Spark engines.

Set to: false

**yarn.nodemanager.aux-services**

Required for dynamic resource allocation for the Spark engine.

Add an entry for "spark\_shuffle."

**yarn.nodemanager.aux-services.spark\_shuffle.class**

Required for dynamic resource allocation for the Spark engine.

Set to: org.apache.spark.network.yarn.YarnShuffleService

**yarn.resourcemanager.scheduler.class**

Defines the YARN scheduler that the Data Integration Service uses to assign resources.

Set to: org.apache.hadoop.yarn.server.resourcemanager.scheduler

**yarn.node-labels.enabled**

Enables node labeling.

Set to: TRUE

**yarn.node-labels.fs-store.root-dir**

The HDFS location to update node label dynamically.

Set to: <hdfs://[Node name]:[Port]/[Path to store]/[Node labels]/>

## Prepare the Archive File for Import from Google Dataproc

To provide an archive file for the Informatica administrator, ensure that you include all required site-\*.xml files.

Create a .zip or .tar file that contains the following \*-site.xml files:

- core-site.xml
- fair-scheduler.xml
- hbase-site.xml. Required only for access to HBase sources and targets.
- hdfs-site.xml
- hive-site.xml
- mapred-site.xml
- tez-site.xml
- yarn-site.xml

Give the Informatica administrator access to the archive file to import the cluster information into the domain.

## Verify the Distribution Version

Verify that the Distribution version property is set to 2.0.

In Informatica version 10.5.1x and lower, Informatica creates the cluster configuration with the Distribution version property set to "1.4 (default)" by default. To integrate with Dataproc 2.x clusters, manually update the Distribution version property to 2.0.

1. In the Administrator tool, click **Connections**.
2. Expand the **Cluster Configurations** node in the Domain Navigator and select the Dataproc cluster configuration.
3. Set the Distribution Version property to 2.0.
4. Save the changes and restart the domain.

## Create a Cluster Configuration

### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

A cluster configuration is an object in the domain that contains configuration information about the Hadoop cluster. The cluster configuration enables the Data Integration Service to push mapping logic to the Hadoop environment. Import configuration properties from the Hadoop cluster to create a cluster configuration.

The import process imports values from \*-site.xml files into configuration sets based on the individual \*-site.xml files. When you perform the import, the cluster configuration wizard can create Hadoop, HBase, HDFS, and Hive connection to access the Hadoop environment. If you choose to create the connections, the wizard also associates the cluster configuration with the connections.

**Note:** If you are integrating for the first time and you imported the cluster configuration when you ran the installer, you *must* re-create or refresh the cluster configuration.

## Importing a Hadoop Cluster Configuration from a File

You can import properties from an archive file to create a cluster configuration.

Before you import from the cluster, you must get the archive file from the Hadoop administrator.

1. From the **Connections** tab, click the **ClusterConfigurations** node in the Domain Navigator.
2. From the Actions menu, select **New > Cluster Configuration**.

The **Cluster Configuration** wizard opens.

3. Configure the following properties:

Property	Description
Cluster configuration name	Name of the cluster configuration.
Description	Optional description of the cluster configuration.
Distribution type	The cluster Hadoop distribution type.
Method to import the cluster configuration	Choose <b>Import from file</b> to import properties from an archive file.
Create connections	<p>Choose to create Hadoop, HDFS, Hive, and HBase connections.</p> <p>If you choose to create connections, the Cluster Configuration wizard associates the cluster configuration with each connection that it creates.</p> <p>The Hadoop connection contains default values for properties such as cluster environment variables, cluster path variables, and advanced properties. Based on the cluster environment and the functionality that you use, you can add to the default values or change the default values of these properties. For a list of Hadoop connection properties to configure, see <a href="#">“Configuring Hadoop Connection Properties” on page 269</a>.</p> <p>If you do not choose to create connections, you must manually create them and associate the cluster configuration with them.</p> <p><b>Important:</b> When the wizard creates the Hive connection, it populates the Metadata Connection String and the Data Access Connection String properties with the value from the hive.metastore.uris property. If the Hive metastore and HiveServer2 are running on different nodes, you must update the Metadata Connection String to point to the HiveServer2 host.</p>

4. Click **Browse** to select a file. Select the file and click **Open**.
5. Click **Next** and verify the cluster configuration information on the summary page.

## Configure the Cluster for the Blaze Engine

To run mappings on the Blaze engine, the cluster requires the libncurses.so.5 library.

### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

You can use the following command to install the 64-bit libncurses.so.5 library on all cluster nodes:

```
sudo apt-get install libncurses5
```

# Configure Domain Settings

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

Perform the following steps to configure the Informatica domain to run mappings on the Dataproc cluster:

1. Copy the krb5.conf file from the cluster to the domain.
2. Enable access to Google Cloud Services.
3. Verify JDBC drivers to enable Sqoop connectivity.

## Copy the Kerberos Configuration File

Copy the Kerberos configuration file Krb5.conf from the cluster to the following directories:

- <Informatica installation directory>/services/shared/security
- <Informatica installation directory>/java/jre/lib/security

## Enable Access to Google Cloud Sources

You can use Google Cloud resources as mapping sources. You can enable access to Google Cloud storage by saving account credentials in a JSON file on each Data Integration Service machine.

1. Create service account credentials in JSON file format.  
Set permissions for the account to access the Google Cloud resources that act as data sources for mappings. To create the JSON credentials file, see the [Google documentation](#).
2. Save the JSON file on the machine that hosts the Data Integration Service. Set full permissions to enable the Data Integration Service to access the file.
3. In the Administrator tool, select the Data Integration Service and edit the Environment Variables properties to point to the JSON file.

Add the following key-value pair to the environment variables:

Key: GOOGLE\_APPLICATION\_CREDENTIALS

Value: <path to the JSON file>

4. Recycle the Data Integration Service.

**Note:** Perform these steps on each machine that hosts the Data Integration Service.

## Verify JDBC Drivers for Sqoop Connectivity

Verify that you have the JDBC drivers to access JDBC-compliant databases in the Hadoop environment. You might need separate drivers for metadata import and for run-time processing.

## Perform this task in the following situations:

- You are integrating for the first time.

You download drivers based on design-time and run-time requirements:

- **Design-time.** To import metadata, you can use the DataDirect drivers packaged with the Informatica installer if they are available. If they are not available, use any Type 4 JDBC driver that the database vendor recommends.
- **Run-time.** To run mappings, use any Type 4 JDBC driver that the database vendor recommends. Some distributions support other drivers to use Sqoop connectors. You cannot use the DataDirect drivers for run-time processing.

## Verify Design-time Drivers

Use the DataDirect JDBC drivers packaged with the Informatica installer to import metadata from JDBC-compliant databases. If the DataDirect JDBC drivers are not available for a specific JDBC-compliant database, download the Type 4 JDBC driver associated with that database.

Copy the JDBC driver .jar files to the following location on the Developer tool machine:

```
<Informatica installation directory>\clients\externaljdbcjars
```

## Verify Run-time Drivers

Verify run-time drivers for mappings that access JDBC-compliant databases in the Hadoop environment. Use any Type 4 JDBC driver that the database vendor recommends.

1. Download Type 4 JDBC drivers associated with the JCBC-compliant databases that you want to access.
2. To optimize the Sqoop mapping performance on the Spark engine while writing data to an HDFS complex file target of the Parquet format, download the following .jar files:

- `parquet-hadoop-bundle-1.6.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-hadoop-bundle/1.6.0/>
- `parquet-avro-1.6.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-avro/1.6.0/>
- `parquet-column-1.5.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-column/1.5.0/>

3. Copy all of the .jar files to the following directory on the machine where the Data Integration Service runs:

```
<Informatica installation directory>\externaljdbcjars
```

Changes take effect after you recycle the Data Integration Service. At run time, the Data Integration Service copies the .jar files to the Hadoop distribution cache so that the .jar files are accessible to all nodes in the cluster.

# Complete Upgrade Tasks

If you upgraded the Informatica platform from version 10.2.2, update connections.

## Complete Connection Upgrade

If you *did not* create connections when you imported the cluster configuration, you must update connection properties for Hadoop, Hive, HDFS, and HBase connections.

### Perform this task in the following situations:

- You upgraded from version 10.2.2 or earlier.

Perform the following tasks to update the connections:

#### Update changed properties

Review connections that you created in a previous release to update the values for connection properties. For example, if you added nodes to the cluster or if you updated the distribution version, you might need to verify host names, URIs, or port numbers for some of the properties.

#### Associate the cluster configuration

The Hadoop, Hive, HDFS, and HBase connections must be associated with a cluster configuration. Complete the following tasks:

1. Run `infacmd isp listConnections` to identify the connections that you need to upgrade. Use `-ct` to list connections of a particular type.
2. Run `infacmd isp UpdateConnection` to associate the cluster configuration with the connection. Use `-cn` to name the connection and `-o clusterConfigID` to associate the cluster configuration with the connection.

For more information about `infacmd`, see the *Informatica Command Reference*.

## Configure the Developer Tool

Perform these steps on each machine that hosts the Developer tool.

### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

## Edit the `etc/hosts` File

To enable Informatica to access the Dataproc cluster, edit the `/etc/hosts` file on client machines.

Enter the IP address, DNS name, and DNS short name for each data node on the cluster.

For example:

```
10.20.30.40 dataprocABC-m.c.MyUsernameINF54321.internal dataprocABC-m
10.20.30.44 dataprocABC-w-0.c.MyUsernameINF54321.internal dataprocABC-w-0
10.20.30.43 dataprocABC-w-1.c.MyUsernameINF54321.internal dataprocABC-w-1
10.20.30.42 dataprocABC-w-2.c.MyUsernameINF54321.internal dataprocABC-w-2
```

## Configure developerCore.ini

Edit the developerCore.ini file to import complex files.

Edit the developerCore.ini file on each machine that hosts the Developer tool.

When you import a complex file, such as Avro or Parquet, the imported object includes metadata associated with the distribution in the Hadoop environment. If the file resides on the Developer tool machine, the import process picks up the distribution information from the developerCore.ini file. You must edit the developerCore.ini file to point to the distribution directory on the Developer tool machine.

You can find the developerCore.ini file in the following directory:

```
<Informatica installation directory>\clients\DeveloperClient
```

Add the following property:

```
-DINFA_HADOOP_DIST_DIR=hadoop<distribution>_<version>
```

The change takes effect when you restart the Developer tool.

## Configure the Developer Tool for Kerberos

To import metadata from Hive, HBase, and complex file sources, import configuration files from the Kerberos-enabled cluster, and generate the Kerberos credentials file on the Developer tool machine.

### Import configuration files

The Hadoop cluster uses a set of XML files named \*-site.xml to store configuration settings. The domain uses the same set of files to create the cluster configuration object.

To enable you to import metadata from the cluster, import the \*-site.xml files to each Developer tool machine:

1. Log in to the Administrator tool and navigate to **Connections > Cluster Configuration > CCO**. Locate the cluster configuration associated with the Hadoop cluster.
2. Extract the \*-site.xml files in the cluster configuration, including sensitive properties, to the following directory on the Developer tool machine: `<Informatica installation directory>\clients\DeveloperClient\hadoop\<Hadoop distribution>\conf`  
For more information about sensitive properties, see the "Active Properties View" topic in the *Data Engineering Administrator Guide*.

**Note:** If you refresh the cluster configuration, repeat these steps.

### Generate the Kerberos credentials file

1. Copy the krb5.conf file from `<Developer tool installation directory>/services/shared/security` to `C:/Windows`.
2. Rename krb5.conf to krb5.ini.
3. In the krb5.ini file, verify the value of the forwardable option to determine how to use the kinit command. If `forwardable=true`, run the command with the `-f` option. Otherwise, run the command without the `-f` option.
4. To generate the Kerberos credentials file, run the kinit command from the following location: `<Developer tool installation directory>/clients/java/bin/kinit.exe`  
For example, you might run the following command: `kinit joe/domain12345@MY-REALM`

## CHAPTER 8

# Hortonworks HDP Integration Tasks

This chapter includes the following topics:

- [Hortonworks HDP Task Flows, 144](#)
- [Prepare for Cluster Import from Hortonworks HDP, 149](#)
- [Create a Cluster Configuration, 155](#)
- [Verify or Refresh the Cluster Configuration , 157](#)
- [Configure the Hive Warehouse Connector and Hive LLAP, 158](#)
- [Verify JDBC Drivers for Sqoop Connectivity, 159](#)
- [Configure the Developer Tool, 160](#)
- [Complete Upgrade Tasks, 161](#)

## Hortonworks HDP Task Flows

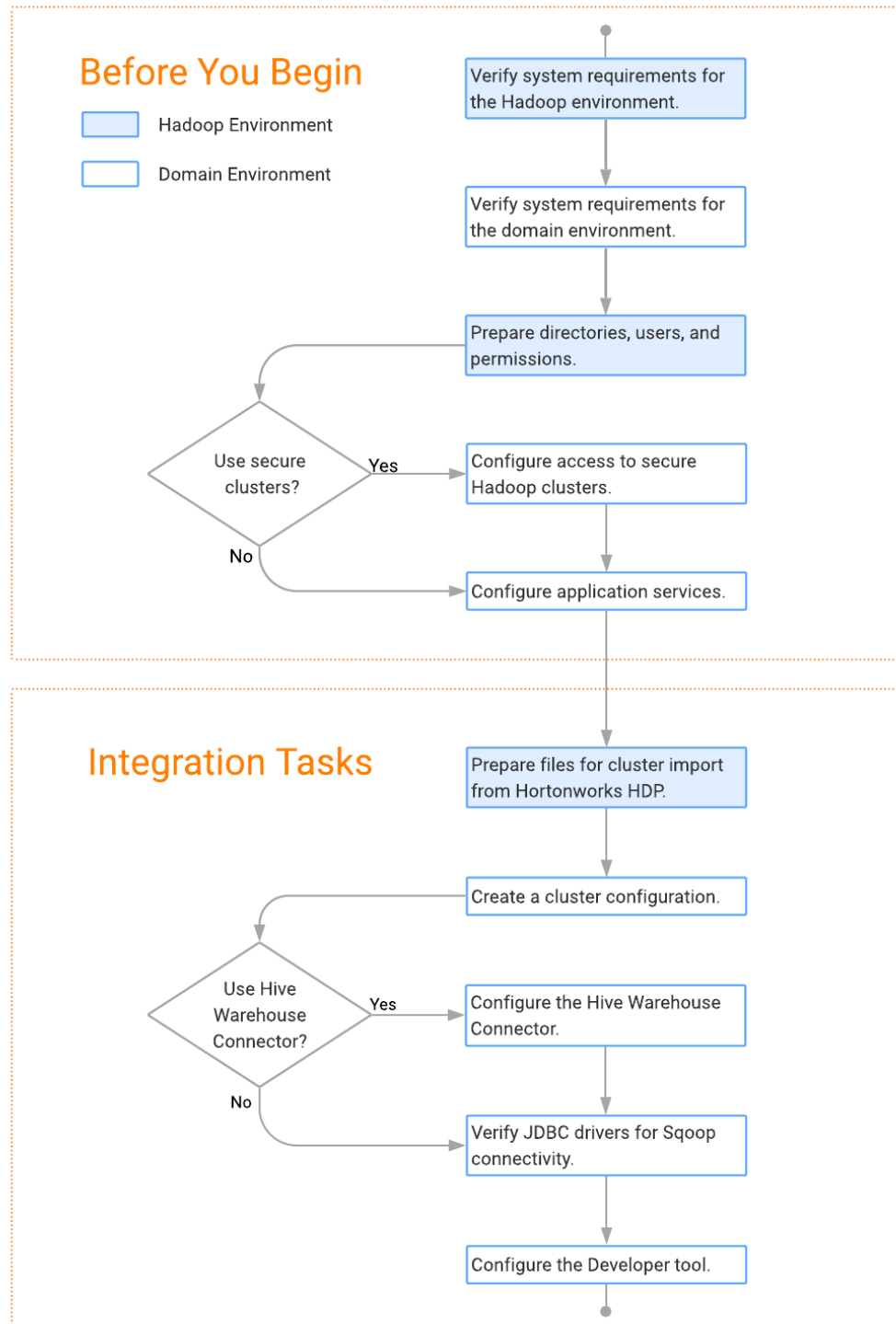
Depending on whether you want to integrate or upgrade Data Engineering Integration in a Hortonworks HDP environment, you can use the flow charts to perform the following tasks:

- Integrate the Informatica domain with Hortonworks HDP for the first time.
- Upgrade from version 10.2.1.
- Upgrade from version 10.2.
- Upgrade from a version earlier than 10.2.



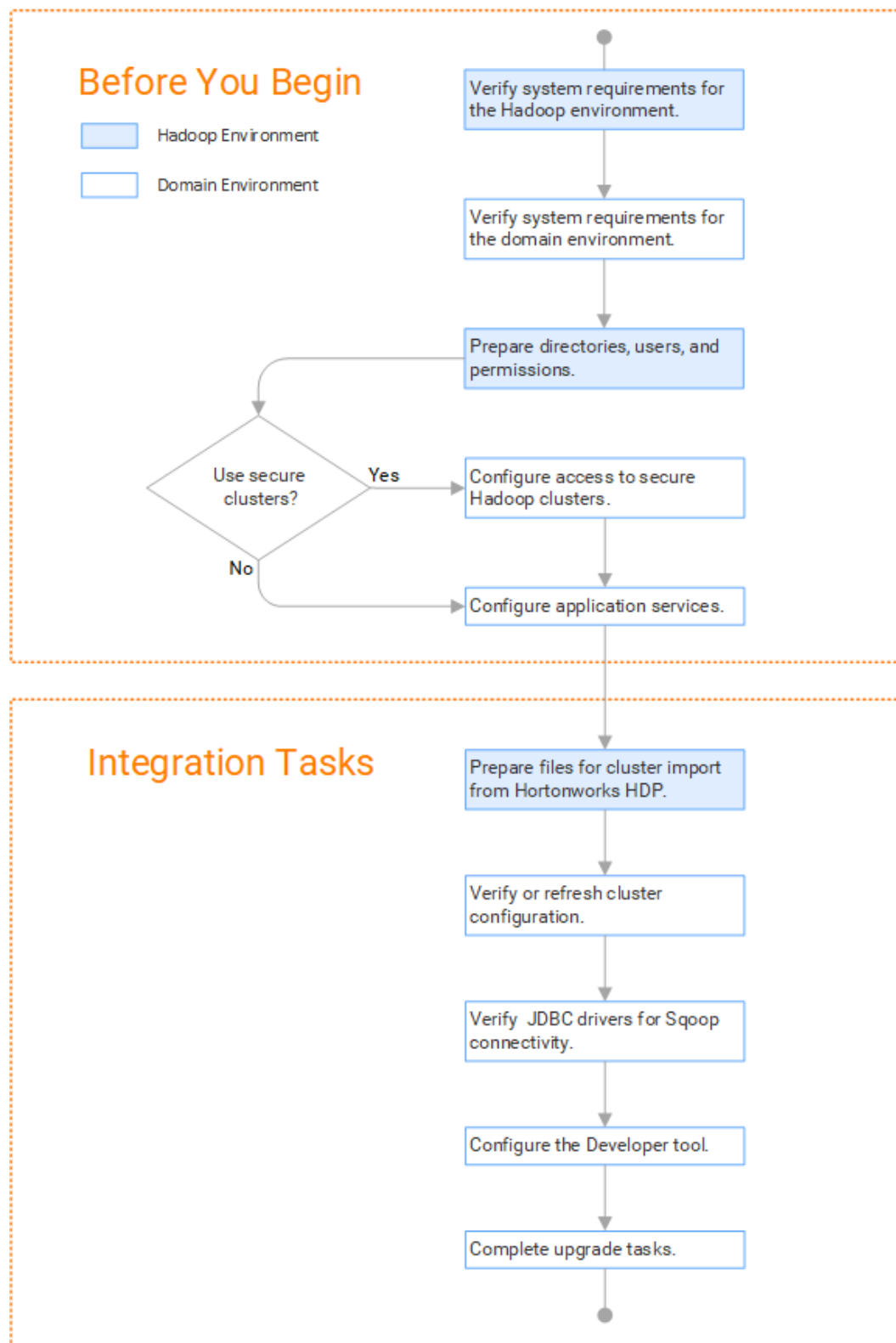
## Task Flow to Integrate with Hortonworks HDP

The following diagram shows the task flow to integrate the Informatica domain with Hortonworks HDP:



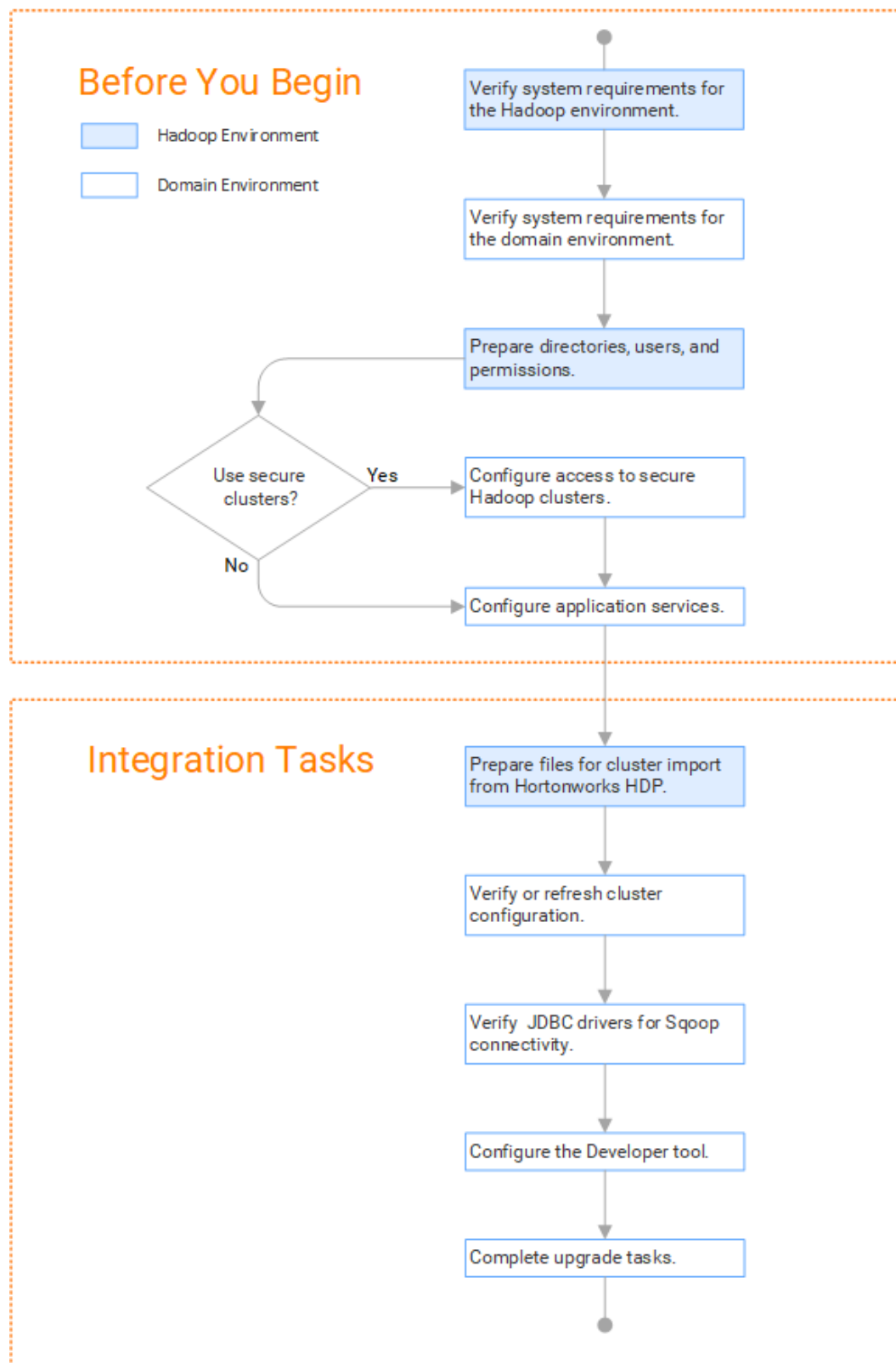
## Task Flow to Upgrade from Version 10.2.1

The following diagram shows the task flow to upgrade version 10.2.1 for Hortonworks HDP:



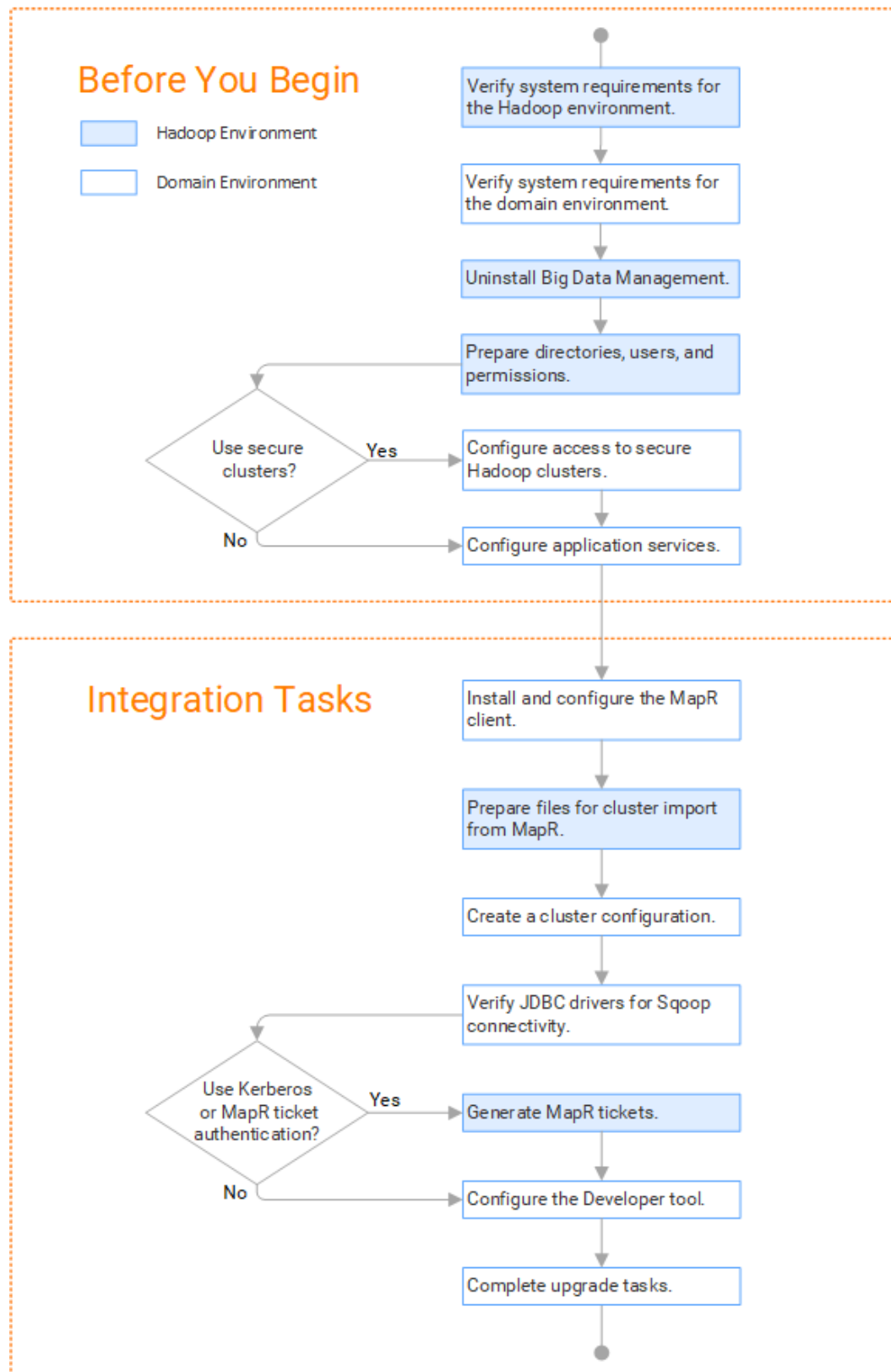
## Task Flow to Upgrade from Version 10.2

The following diagram shows the task flow to upgrade version 10.2 for Hortonworks HDP:



## Task Flow to Upgrade from a Version Earlier than 10.2

The following diagram shows the task flow to upgrade from a version earlier than 10.2 for Hortonworks HDP:



# Prepare for Cluster Import from Hortonworks HDP

Before the Informatica administrator can import cluster information to create a cluster configuration in the Informatica domain, the Hadoop administrator must perform some preliminary tasks.

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

**Note:** If you are upgrading from a previous version, verify the properties and suggested values, as Data Engineering Integration might require additional properties or different values for existing properties.

Complete the following tasks to prepare the cluster before the Informatica administrator creates the cluster configuration:

1. Verify property values in \*-site.xml files that Data Engineering Integration needs to run mappings in the Hadoop environment.
2. Provide information to the Informatica administrator that is required to import cluster information into the domain. Depending on the method of import, perform one of the following tasks:
  - To import directly from the cluster, give the Informatica administrator cluster authentication information to connect to the cluster.
  - To import from an archive file, export cluster information and provide an archive file to the Data Engineering Integration administrator.

## Configure \*-site.xml Files for Hortonworks HDP

The Hadoop administrator needs to configure \*-site.xml file properties and restart impacted services before the Informatica administrator imports cluster information into the domain.

### capacity-scheduler.xml

Configure the following properties in the capacity-scheduler.xml file:

#### **yarn.scheduler.capacity.<queue path>.disable\_preemption**

Disables preemption for the Capacity Scheduler. The Blaze engine does not support preemption. If YARN preemption is enabled for the cluster, you need to disable it for the queue allocated to the Blaze engine.

Set to TRUE for queue allocated to the Blaze engine.

### core-site.xml

Configure the following properties in the core-site.xml file:

#### **fs.s3.enableServerSideEncryption**

Enables server side encryption for S3 buckets. Required for SSE and SSE-KMS encryption.

Set to: TRUE

#### **fs.s3a.access.key**

The ID for the Blaze and Spark engines to connect to the Amazon S3 file system.

Set to your access key.

**fs.s3a.secret.key**

The password for the Blaze and Spark engines to connect to the Amazon S3 file system

Set to your access ID.

**fs.s3a.server-side-encryption-algorithm**

The server-side encryption algorithm for S3. Required for SSE and SSE-KMS encryption. Set to the encryption algorithm used.

**hadoop.proxyuser.<proxy user>.groups**

Defines the groups that the proxy user account can impersonate. On a secure cluster the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to group names of impersonation users separated by commas. If less security is preferred, use the wildcard " \* " to allow impersonation from any group.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

**hadoop.proxyuser.<proxy user>.hosts**

Defines the host machines that a user account can impersonate. On a secure cluster the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to a single host name or IP address, or set to a comma-separated list. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

**hadoop.proxyuser.hive.hosts**

Comma-separated list of hosts that you want to allow the Hive user to impersonate on a non-secure cluster.

When `hive.server2.enable.doAs` is false, append a comma-separated list of Informatica server host names or IP address where the Data Integration Service is running. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to this property, restart the cluster services that use core-site configuration values.

**hadoop.proxyuser.yarn.groups**

Comma-separated list of groups that you want to allow the YARN user to impersonate on a non-secure cluster.

Set to group names of impersonation users separated by commas. If less security is preferred, use the wildcard " \* " to allow impersonation from any group.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

**hadoop.proxyuser.yarn.hosts**

Comma-separated list of hosts that you want to allow the YARN user to impersonate on a non-secure cluster.

Set to a single host name or IP address, or set to a comma-separated list. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.security.auth\_to\_local**

Translates the principal names from the Active Directory and MIT realm into local names within the Hadoop cluster. Based on the Hadoop cluster used, you can set multiple rules.

Set to: RULE:[1:\$1@\$0](^.\*@YOUR.REALM)s/^.\*(.)@YOUR.REALM\.COM\$/ \$1/g

Set to: RULE:[2:\$1@\$0](^.\*@YOUR.REALM\.)s/^.\*(.)@YOUR.REALM\.COM\$/ \$1/g

#### **hbase-site.xml**

Configure the following properties in the hbase-site.xml file:

##### **zookeeper.znode.parent**

Identifies HBase master and region servers.

Set to the relative path to the znode directory of HBase.

#### **hdfs-site.xml**

Configure the following properties in the hdfs-site.xml file:

##### **dfs.encryption.key.provider.uri**

The KeyProvider used to interact with encryption keys when reading and writing to an encryption zone. Required if sources or targets reside in the HDFS encrypted zone on Java KeyStore KMS-enabled Cloudera CDH cluster or a Ranger KMS-enabled Hortonworks HDP cluster.

Set to: kmf://http@xx11.xyz.com:16000/kms

#### **hive-site.xml**

Configure the following properties in the hive-site.xml file:

##### **hive.cluster.delegation.token.store.class**

The token store implementation. Required for HiveServer2 high availability and load balancing.

Set to: org.apache.hadoop.hive.thrift.DBTokenStore

##### **hive.compactor.initiator.on**

Runs the initiator and cleaner threads on metastore instance. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: TRUE

##### **hive.compactor.worker.threads**

The number of worker threads to run in a metastore instance. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: 1

##### **hive.enforce.bucketing**

Enables dynamic bucketing while loading to Hive. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: TRUE

##### **io.compression.codecs**

Enables compression on temporary staging tables.

Set to a comma-separated list of compression codec classes on the cluster.

**hive.exec.dynamic.partition.mode**

Allows all partitions to be dynamic. Required for the Update Strategy transformation in a mapping that writes to a Hive target. Also required if you use Sqoop and define a DDL query to create or replace a partitioned Hive target at run time.

Set to: nonstrict

**hive.support.concurrency**

Enables table locking in Hive. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: TRUE

**hive.server2.support.dynamic.service.discovery**

Enables HiveServer2 dynamic service discovery. Required for HiveServer2 high availability.

Set to: TRUE

**hive.server2.zookeeper.namespace**

The value of the ZooKeeper namespace in the JDBC connection string. Required for HiveServer2 high availability.

Set to: jdbc:hive2://<zookeeper\_ensemble>/  
default;serviceDiscoveryMode=zooKeeper;zooKeeperNamespace=hiveserver2

**hive.txn.manager**

Turns on transaction support. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: org.apache.hadoop.hive.ql.lockmgr.DbTxnManager

**hive.zookeeper.quorum**

Comma-separated list of ZooKeeper server host:ports in a cluster. The value of the ZooKeeper ensemble in the JDBC connection string. Required for HiveServer2 high availability.

Set to: jdbc:hive2://<zookeeper\_ensemble>/default;serviceDiscoveryMode=zooKeeper;

[mapred-site.xml](#)

Configure the following properties in the mapred-site.xml file:

**mapreduce.framework.name**

The run-time framework to run MapReduce jobs. Values can be local, classic, or yarn. Required for Sqoop.

Set to: yarn

**yarn.app.mapreduce.am.staging-dir**

The HDFS staging directory used while submitting jobs.

Set to the staging directory path.

[yarn-site.xml](#)

Configure the following properties in the yarn-site.xml file:

**yarn.application.classpath**

Required for dynamic resource allocation.



Add spark\_shuffle.jar to the class path. The .jar file must contain the class "org.apache.spark.network.yarn.YarnShuffleService."

**yarn.nodemanager.resource.memory-mb**

The maximum RAM available for each container. Set the maximum memory on the cluster to increase resource memory available to the Blaze engine.

Set to 16 GB if value is less than 16 GB.

**yarn.nodemanager.resource.cpu-vcores**

The number of virtual cores for each container. Required for Blaze engine resource allocation.

Set to 10 if the value is less than 10.

**yarn.scheduler.minimum-allocation-mb**

The minimum RAM available for each container. Required for Blaze engine resource allocation.

Set to 6 GB if the value is less than 6 GB.

**yarn.nodemanager.vmem-check-enabled**

Disables virtual memory limits for containers. Required for the Blaze and Spark engines.

Set to: false

**yarn.nodemanager.aux-services**

Required for dynamic resource allocation for the Spark engine.

Add an entry for "spark\_shuffle."

**yarn.nodemanager.aux-services.spark\_shuffle.class**

Required for dynamic resource allocation for the Spark engine.

Set to: org.apache.spark.network.yarn.YarnShuffleService

**yarn.resourcemanager.scheduler.class**

Defines the YARN scheduler that the Data Integration Service uses to assign resources.

Set to: org.apache.hadoop.yarn.server.resourcemanager.scheduler

**yarn.node-labels.enabled**

Enables node labeling.

Set to: TRUE

**yarn.node-labels.fs-store.root-dir**

The HDFS location to update node label dynamically.

Set to: <hdfs://[Node name]:[Port]/[Path to store]/[Node labels]/>

[tez-site.xml](#)

Configure the following properties in the tez-site.xml file:

**tez.runtime.io.sort.mb**

The sort buffer memory. Required when the output needs to be sorted for Blaze and Spark engines.

Set value to 270 MB.

## Prepare for Direct Import from Hortonworks HDP

If you plan to provide direct access to the Informatica administrator to import cluster information, provide the required connection information.

The following table describes the information that you need to provide to the Informatica administrator to create the cluster configuration directly from the cluster:

Property	Description
Host	Host name or IP address of the cluster manager.
Port	Port of the cluster manager.
User ID	Cluster user ID.
Password	Password for the user.
Cluster name	Name of the cluster. Use the display name if the cluster manager manages multiple clusters. If you do not provide a cluster name, the wizard imports information based on the default cluster.

## Prepare the Archive File for Import from Hortonworks HDP

When you prepare the archive file for cluster configuration import from Hortonworks, include all required \*-site.xml files and edit the file manually after you create it.

The Hortonworks cluster configuration archive file must have the following contents:

- core-site.xml
- hbase-site.xml. hbase-site.xml is required only if you access HBase sources and targets.
- hdfs-site.xml
- hive-site.xml
- mapred-site.xml or tez-site.xml. Include the mapred-site.xml file or the tez-site.xml file based on the Hive execution type used on the Hadoop cluster.
- yarn-site.xml

### Update the Archive File

After you create the archive file, edit the Hortonworks Data Platform (HDP) version string wherever it appears in the archive file. Search for the string `${hdp.version}` and replace all instances with the HDP version that Hortonworks includes in the Hadoop distribution.

For example, the edited `tez.lib.uris` property looks similar to the following:

```
<property>
<name>tez.lib.uris</name>
<value>/hdp/apps/2.5.0.0-1245/tez/tez.tar.gz</value>
</property>
```

# Create a Cluster Configuration

After the Hadoop administrator prepares the cluster for import, the Informatica administrator must create a cluster configuration.

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from version 10.1.1 or earlier.

A cluster configuration is an object in the domain that contains configuration information about the Hadoop cluster. The cluster configuration enables the Data Integration Service to push mapping logic to the Hadoop environment. Import configuration properties from the Hadoop cluster to create a cluster configuration.

The import process imports values from \*-site.xml files into configuration sets based on the individual \*-site.xml files. When you perform the import, the cluster configuration wizard can create Hadoop, HBase, HDFS, and Hive connection to access the Hadoop environment. If you choose to create the connections, the wizard also associates the cluster configuration with the connections.

**Note:** If you are integrating for the first time and you imported the cluster configuration when you ran the installer, you *must* re-create or refresh the cluster configuration.

## Before You Import

Before you can import the cluster configuration, you must get information from the Hadoop administrator based on the method of import.

If you import directly from the cluster, contact the Hadoop administrator to get cluster connection information. If you import from a file, get an archive file of exported cluster information.

## Importing a Hadoop Cluster Configuration from the Cluster

When you import the Hadoop cluster configuration directly from the cluster, you provide information to connect to the cluster.

Get cluster connection information from the Hadoop administrator.

1. From the **Connections** tab, click the **ClusterConfigurations** node in the Domain Navigator.
2. From the Actions menu, select **New > Cluster Configuration**.

The **Cluster Configuration** wizard opens.

3. Configure the following General properties:

Property	Description
Cluster configuration name	Name of the cluster configuration.
Description	Optional description of the cluster configuration.
Distribution type	The cluster Hadoop distribution type.

Property	Description
Method to import the cluster configuration	Choose <b>Import from cluster</b> .
Create connections	<p>Choose to create Hadoop, HDFS, Hive, and HBase connections.</p> <p>If you choose to create connections, the Cluster Configuration wizard associates the cluster configuration with each connection that it creates.</p> <p>The Hadoop connection contains default values for properties such as cluster environment variables, cluster path variables, and advanced properties. Based on the cluster environment and the functionality that you use, you can add to the default values or change the default values of these properties. For a list of Hadoop connection properties to configure, see Hadoop Cluster Connection Properties.</p> <p>If you do not choose to create connections, you must manually create them and associate the cluster configuration with them.</p> <p><b>Important:</b> When the wizard creates the Hive connection, it populates the Metadata Connection String and the Data Access Connection String properties with the value from the hive.metastore.uris property. If the Hive metastore and HiveServer2 are running on different nodes, you must update the Metadata Connection String to point to the HiveServer2 host.</p>

The cluster properties appear.

- Configure the following properties:

Property	Description
Host	Host name or IP address of the cluster manager.
Port	Port of the cluster manager.
User ID	Cluster user ID.
Password	Password for the user.
Cluster name	Name of the cluster. Use the display name if the cluster manager manages multiple clusters. If you do not provide a cluster name, the wizard imports information based on the default cluster.

- Click **Next** and verify the cluster configuration information on the summary page.

## Importing a Hadoop Cluster Configuration from a File

You can import properties from an archive file to create a cluster configuration.

Before you import from the cluster, you must get the archive file from the Hadoop administrator.

- From the **Connections** tab, click the **ClusterConfigurations** node in the Domain Navigator.
- From the Actions menu, select **New > Cluster Configuration**.

The **Cluster Configuration** wizard opens.

3. Configure the following properties:

Property	Description
Cluster configuration name	Name of the cluster configuration.
Description	Optional description of the cluster configuration.
Distribution type	The cluster Hadoop distribution type.
Method to import the cluster configuration	Choose <b>Import from file</b> to import properties from an archive file.
Create connections	<p>Choose to create Hadoop, HDFS, Hive, and HBase connections.</p> <p>If you choose to create connections, the Cluster Configuration wizard associates the cluster configuration with each connection that it creates.</p> <p>The Hadoop connection contains default values for properties such as cluster environment variables, cluster path variables, and advanced properties. Based on the cluster environment and the functionality that you use, you can add to the default values or change the default values of these properties. For a list of Hadoop connection properties to configure, see <a href="#">“Configuring Hadoop Connection Properties” on page 269</a>.</p> <p>If you do not choose to create connections, you must manually create them and associate the cluster configuration with them.</p> <p><b>Important:</b> When the wizard creates the Hive connection, it populates the Metadata Connection String and the Data Access Connection String properties with the value from the hive.metastore.uris property. If the Hive metastore and HiveServer2 are running on different nodes, you must update the Metadata Connection String to point to the HiveServer2 host.</p>

4. Click **Browse** to select a file. Select the file and click **Open**.
5. Click **Next** and verify the cluster configuration information on the summary page.

## Verify or Refresh the Cluster Configuration

You might need to refresh the cluster configuration or update the distribution version in the cluster configuration when you upgrade.

### Perform this task in the following situation:

- You upgraded from version 10.2 or later.

### Verify the Cluster Configuration

The cluster configuration contains a property for the distribution version. The verification task depends on the version you upgraded:

#### Upgrade from 10.2

If you upgraded from 10.2 and you changed the distribution version, you need to verify the distribution version in the General properties of the cluster configuration.

### Upgrade from 10.2.1

Effective in version 10.2.1, Informatica assigns a default version to each Hadoop distribution type. If you configure the cluster configuration to use the default version, the upgrade process upgrades to the assigned default version if the version changes. If you have not upgraded your Hadoop distribution to Informatica's default version, you need to update the distribution version property.

For example, suppose the assigned default Hadoop distribution version for 10.2.1 is  $n$ , and for 10.2.2 is  $n+1$ . If the cluster configuration uses the default supported Hadoop version of  $n$ , the upgraded cluster configuration uses the default version of  $n+1$ . If you have not upgraded the distribution in the Hadoop environment you need to change the cluster configuration to use version  $n$ .

If you configure the cluster configuration to use a distribution version that is not the default version, you need to update the distribution version property in the following circumstances:

- Informatica dropped support for the distribution version.
- You changed the distribution version.

### Refresh the Cluster Configuration

If you updated any of the \*-site.xml files noted in the topic to prepare for cluster import, you need to refresh the cluster configuration in the Administrator tool.

## Configure the Hive Warehouse Connector and Hive LLAP

Optionally, you can configure the Hive Warehouse Connector and Hive LLAP to improve performance when you read from and write to Hive targets.

#### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from any previous version.

The Hive Warehouse Connector reads from and writes to Hive tables without using temporary staging tables that require additional storage overhead. Use the Hive Warehouse Connector on the Spark engine to allow Spark code to interact with Hive targets and to use ACID-enabled Hive tables. When you enable the Hive Warehouse Connector, mappings use Hive LLAP to run Hive queries rather than HiveServer2.

Before you enable the Hive Warehouse Connector, enable Hive LLAP on the Hadoop cluster. To enable the connector, configure the following properties in the Spark advanced properties for the Hadoop connection:

#### **infaspark.useHiveWarehouseAPI**

Enables the Hive Warehouse Connector. Set to TRUE.

For example, `infaspark.useHiveWarehouseAPI=true`.

#### **spark.datasource.hive.warehouse.load.staging.dir**

Directory for the temporary HDFS files used for batch writes to Hive. Required when you enable the Hive Warehouse Connector.

For example, set to `/tmp`

**spark.datasource.hive.warehouse.metastoreUri**

URI for the Hive metastore. Required when you enable the Hive Warehouse Connector. Use the value for *hive.metastore.uris* from the *hive\_site.xml* cluster configuration properties.

For example, set the value to `thrift://mycluster-1.com:9083`.

**spark.hadoop.hive.llap.daemon.service.hosts**

Application name for the LLAP service. Required when you enable the Hive Warehouse Connector. Use the value for *hive.llap.daemon.service.hosts* from the *hive\_site.xml* cluster configuration properties.

**spark.hadoop.hive.zookeeper.quorum**

Zookeeper hosts used by Hive LLAP. Required when you enable the Hive Warehouse Connector. Use the value for *hive.zookeeper.quorum* from the *hive\_site.xml* cluster configuration properties.

**spark.sql.hive.hiveserver2.jdbc.url**

URL for HiveServer2 Interactive. Required to use the Hive Warehouse Connector. Use the value in Ambari for HiveServer2 JDBC URL.

## Verify JDBC Drivers for Sqoop Connectivity

Verify that you have the JDBC drivers to access JDBC-compliant databases in the Hadoop environment. You might need separate drivers for metadata import and for run-time processing.

**Perform this task in the following situations:**

- You are integrating for the first time.
- You upgraded from version 10.2.1 or earlier.

You download drivers based on design-time and run-time requirements:

- **Design-time.** To import metadata, you can use the DataDirect drivers packaged with the Informatica installer if they are available. If they are not available, use any Type 4 JDBC driver that the database vendor recommends.
- **Run-time.** To run mappings, use any Type 4 JDBC driver that the database vendor recommends. Some distributions support other drivers to use Sqoop connectors. You cannot use the DataDirect drivers for run-time processing.

## Verify Design-time Drivers

Use the DataDirect JDBC drivers packaged with the Informatica installer to import metadata from JDBC-compliant databases. If the DataDirect JDBC drivers are not available for a specific JDBC-compliant database, download the Type 4 JDBC driver associated with that database.

Copy the JDBC driver .jar files to the following location on the Developer tool machine:

```
<Informatica installation directory>\clients\externaljdbcjars
```

## Verify Run-time Drivers

Verify run-time drivers for mappings that access JDBC-compliant databases in the Hadoop environment. Use any Type 4 JDBC driver that the database vendor recommends.

1. Download Type 4 JDBC drivers associated with the JCBC-compliant databases that you want to access.
2. To use Sqoop TDCH Hortonworks Connector for Teradata, perform the following task:

Download all .jar files in the Hortonworks Connector for Teradata package from the following location : <http://hortonworks.com/downloads/#addons>

The package has the following naming convention: `hdp-connector-for-teradata-<version>-distro.tar.gz`

3. To optimize the Sqoop mapping performance on the Spark engine while writing data to an HDFS complex file target of the Parquet format, download the following .jar files:
  - `parquet-hadoop-bundle-1.6.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-hadoop-bundle/1.6.0/>
  - `parquet-avro-1.6.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-avro/1.6.0/>
  - `parquet-column-1.5.0.jar` from <https://repo1.maven.org/maven2/com/twitter/parquet-column/1.5.0/>
4. Copy all of the .jar files to the following directory on the machine where the Data Integration Service runs:

```
<Informatica installation directory>\externaljdbcjars
```

Changes take effect after you recycle the Data Integration Service. At run time, the Data Integration Service copies the .jar files to the Hadoop distribution cache so that the .jar files are accessible to all nodes in the cluster.

## Configure the Developer Tool

You can configure the Developer tool to enable you to import complex files or import metadata when the domain is Kerberos-enabled.

Edit the `developerCore.ini` file to import complex files. Edit the file on each Developer tool machine.

### Configure `developerCore.ini`

Edit the `developerCore.ini` file to import complex files.

Edit the `developerCore.ini` file on each machine that hosts the Developer tool.

When you import a complex file, such as Avro or Parquet, the imported object includes metadata associated with the distribution in the Hadoop environment. If the file resides on the Developer tool machine, the import process picks up the distribution information from the `developerCore.ini` file. You must edit the `developerCore.ini` file to point to the distribution directory on the Developer tool machine.

You can find the `developerCore.ini` file in the following directory:

```
<Informatica installation directory>\clients\DeveloperClient
```

Add the following property:

```
-DINFA_HADOOP_DIST_DIR=hadoop\<distribution>_<version>
```



The change takes effect when you restart the Developer tool.

## Configure the Developer Tool for Kerberos

To import metadata from Hive, HBase, and complex file sources, import configuration files from the Kerberos-enabled cluster, and generate the Kerberos credentials file on the Developer tool machine.

### Import configuration files

The Hadoop cluster uses a set of XML files named \*-site.xml to store configuration settings. The domain uses the same set of files to create the cluster configuration object.

To enable you to import metadata from the cluster, import the \*-site.xml files to each Developer tool machine:

1. Log in to the Administrator tool and navigate to **Connections > Cluster Configuration > CCO**. Locate the cluster configuration associated with the Hadoop cluster.
2. Extract the \*-site.xml files in the cluster configuration, including sensitive properties, to the following directory on the Developer tool machine: `<Informatica installation directory>\clients\DeveloperClient\hadoop\<Hadoop distribution>\conf`  
For more information about sensitive properties, see the "Active Properties View" topic in the *Data Engineering Administrator Guide*.

**Note:** If you refresh the cluster configuration, repeat these steps.

### Generate the Kerberos credentials file

1. Copy the krb5.conf file from `<Developer tool installation directory>/services/shared/security` to `C:/Windows`.
2. Rename krb5.conf to krb5.ini.
3. In the krb5.ini file, verify the value of the forwardable option to determine how to use the kinit command. If `forwardable=true`, run the command with the `-f` option. Otherwise, run the command without the `-f` option.
4. To generate the Kerberos credentials file, run the kinit command from the following location: `<Developer tool installation directory>/clients/java/bin/kinit.exe`  
For example, you might run the following command: `kinit joe/domain12345@MY-REALM`

## Complete Upgrade Tasks

If you upgraded the Informatica platform from version 10.2.2, update connections.

### Update Connections

You might need to update connections based on the version you are upgrading from.

If you did not create connections when you created the cluster configuration, you need to update the connections.

## Configure the Hadoop Connection

To use properties that you customized in the `hadoopEnv.properties` file, you must configure the Hadoop connection properties such as cluster environment variables, cluster path variables, and advanced properties.

### Perform this task in the following situation:

- You upgraded from version 10.1.1 or earlier.

When you run the Informatica upgrade, the installer backs up the existing `hadoopEnv.properties` file. You can find the backup `hadoopEnv.properties` file in the following location:

```
<Previous Informatica installation directory>/services/shared/hadoop/<Hadoop distribution name>_<version>/infaConf
```

Edit the Hadoop connection in the Administrator tool or the Developer tool to include any properties that you manually configured in the `hadoopEnv.properties` file. The Hadoop connection contains default values for properties such as cluster environment and path variables and advanced properties. You can update the default values to match the properties in the `hadoopEnv.properties` file.

## Replace the Connections with New Connections

If you created connections when you imported the cluster configuration, you need to replace connections in mappings with the new connections.

### Perform this task in the following situation:

- You upgraded from version 10.1.1 or earlier.

The method that you use to replace connections in mappings depends on the type of connection.

#### Hadoop connection

Run the following commands to replace the connections:

- `infacmd dis replaceMappingHadoopRuntimeConnections`. Replaces connections associated with mappings that are deployed in applications.
- `infacmd mrs replaceMappingHadoopRuntimeConnections`. Replaces connections associated with mappings that you run from the Developer tool.

For information about the `infacmd` commands, see the *Informatica Command Reference*.

#### Hive, HDFS, and HBase connections

You must replace the connections manually.

## Complete Connection Upgrade

If you *did not* create connections when you imported the cluster configuration, you must update connection properties for Hadoop, Hive, HDFS, and HBase connections.

### Perform this task in the following situations:

- You upgraded from version 10.2.2 or earlier.

Perform the following tasks to update the connections:

### Update changed properties

Review connections that you created in a previous release to update the values for connection properties. For example, if you added nodes to the cluster or if you updated the distribution version, you might need to verify host names, URIs, or port numbers for some of the properties.

### Associate the cluster configuration

The Hadoop, Hive, HDFS, and HBase connections must be associated with a cluster configuration. Complete the following tasks:

1. Run `infacmd isp listConnections` to identify the connections that you need to upgrade. Use `-ct` to list connections of a particular type.
2. Run `infacmd isp UpdateConnection` to associate the cluster configuration with the connection. Use `-cn` to name the connection and `-o clusterConfigID` to associate the cluster configuration with the connection.

For more information about `infacmd`, see the *Informatica Command Reference*.

## Replace Hive Run-time Connections with Hadoop Connections

Effective in version 10.2.2, Big Data Management dropped support for the Hive engine and Hive run-time connections. If you used Hive connections to run mappings on the Hadoop cluster, you must generate Hadoop connections from the Hive connections.

#### Perform this task in the following situations:

- You upgraded from version 10.1.1 or earlier.
- The Hive connections are configured to run mappings in the Hadoop environment.

The upgrade process generates a connection name for the Hadoop connection and replaces the connection name in the mappings. It does not create the physical connection object. When the upgrade is complete, you must run a command to generate the connection. You generate Hadoop connections from Hive connections that are configured to run mappings in the Hadoop environment.

You must generate Hadoop connections from Hive connections that are configured to run mappings in the Hadoop environment.

1. Run `infacmd isp generateHadoopConnectionFromHiveConnection` to generate a Hadoop connection from a Hive connection that is configured to run in the Hadoop environment.  
The command names the connection as follows: "Autogen\_<Hive connection name>." If the connection name exceeds the 128 character limit, the command fails.
2. If the command fails, perform the following tasks:
  - a. Rename the connection to meet the character limit and run the command again.
  - b. Run `infacmd dis replaceMappingHadoopRuntimeConnections` to replace connections associated with mappings that are deployed in applications.
  - c. Run `infacmd mrs replaceMappingHadoopRuntimeConnections` to replace connections associated with mappings that you run from the Developer tool.
3. If the Hive connection was parameterized, you must update the connection names in the parameter file. Verify that the Hive sources, Hive targets, and the Hive engine parameters are updated with the correct connection name.
4. If any properties changed in the cluster, such as host names, URIs, or port numbers, you must update the properties in the connections.

For information about the `infacmd` commands, see the *Informatica® Command Reference*.

## Update Streaming Objects

Data Engineering Streaming uses Spark Structured Streaming to process data instead of Spark Streaming. To support Spark Structured Streaming, some header ports are added to the data objects, and support to some of the data objects and transformations are deferred to a future release. The behavior of some of the data objects is also updated.

After you upgrade, the existing streaming mappings become invalid because of the unavailable header ports, the unsupported transformations or data objects, and the behavior change of some data objects.

<b>Perform this task in the following situations:</b>
- You upgraded from version 10.1.1, 10.2.0, or 10.2.1.



To use an existing streaming mapping, perform the following tasks:

- Re-create the physical data objects. After you re-create the physical data objects, the data objects get the required header ports, such as timestamp, partitionID, or key based on the data object.
- In a Normalizer transformation, if the **Occurs** column is set to Auto, re-create the Normalizer transformation. You must re-create the Normalizer transformation because the type configuration property of the complex port refers to the physical data object that you plan to replace.
- Update the streaming mapping. If the mapping contains Kafka target, Aggregator transformation, Joiner transformation, or Normalizer transformation, replace the data object or transformation, and then update the mapping because of the changed behavior of these transformations and data objects.
- Verify the deferred data object types. If the streaming mapping contains unsupported transformations or data objects, contact Informatica Global Customer Support.

### Re-create the Physical Data Objects

When you re-create the physical data objects, the physical data objects get the header ports and some properties are not available for some data objects. Update the existing mapping with the newly created physical data objects.

1. Go to the existing mapping, select the data object from the mapping.
2. Click the **Properties** tab. On the **Column Projection** tab, click **Edit Schema**.
3. Note the schema information from the **Edit Schema** dialog box.
4. Note the parameters information from the **Parameters** tab.
5. Create new physical data objects.

After you re-create the data objects, the physical data objects get the required header ports. The Microsoft Azure does not support the following properties and are not available for Azure Event Hubs data objects:

- Consumer Properties
- Partition Count

## Update the Streaming Mappings

After you re-create the data object, replace the existing data objects with the re-created data objects. If the mapping contains Normaliser Transformation, Aggregator transformation, or Joiner transformation, update the mapping because of the changed behavior of these transformations and data objects.

### Transformation Updates

If a transformation uses a complex port, configure the type configuration property of the port because the property refers to the physical data object that you replaced.

### Aggregator and Joiner Transformation Updates

An Aggregator transformation must be downstream from a Joiner transformation. A Window transformation must be directly upstream from both Aggregator and Joiner transformations. Previously, you could use an Aggregator transformation anywhere in the streaming mapping.

If a mapping contains an Aggregator transformation upstream from a Joiner transformation, move the Aggregator transformation downstream from a Joiner transformation. Add a Window transformation directly upstream from both Aggregator and Joiner transformations.

## Verify the Deferred Data Object Types

After you upgrade, the streaming mappings might contain some transformations and data objects that are deferred.

The following table lists the data object types to which the support is deferred to a future release:

Object Type	Object
Transformation	Data Masking

If you want to continue using the mappings that contain deferred data objects or transformations, you must contact Informatica Global Customer Support.

## CHAPTER 9

# MapR Integration Tasks

This chapter includes the following topics:

- [MapR Task Flows, 166](#)
- [Install and Configure the MapR Client , 171](#)
- [Prepare for Cluster Import from MapR, 171](#)
- [Create a Cluster Configuration, 177](#)
- [Verify or Refresh the Cluster Configuration , 178](#)
- [Verify JDBC Drivers for Sqoop Connectivity, 179](#)
- [Generate MapR Tickets, 180](#)
- [Complete Upgrade Tasks, 183](#)

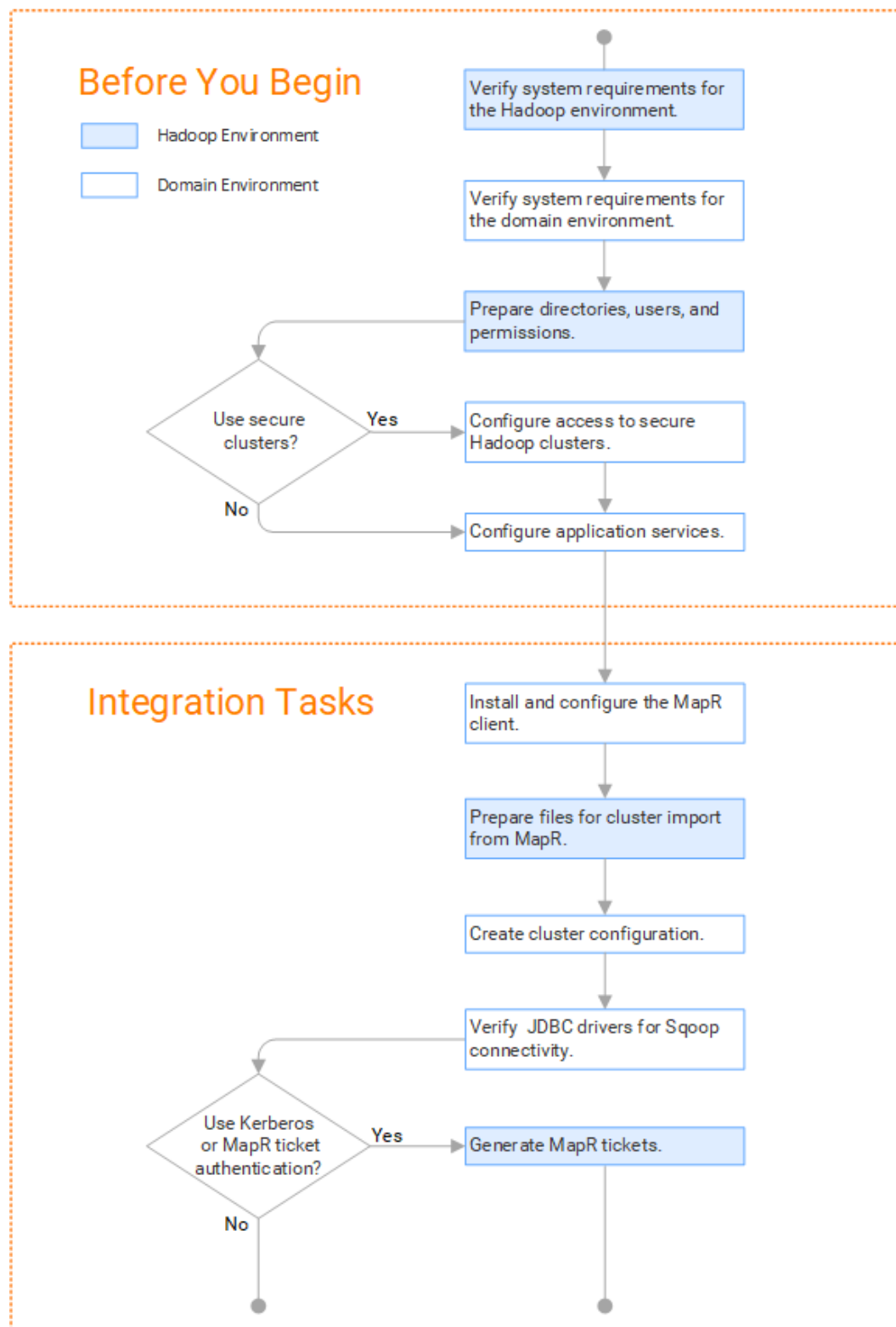
## MapR Task Flows

Depending on whether you want to integrate or upgrade Data Engineering Integration in a MapR environment, you can use the flow charts to perform the following tasks:

- Integrate the Informatica domain with MapR for the first time.
- Upgrade from version 10.2.1.
- Upgrade from version 10.2.
- Upgrade from a version earlier than 10.2.

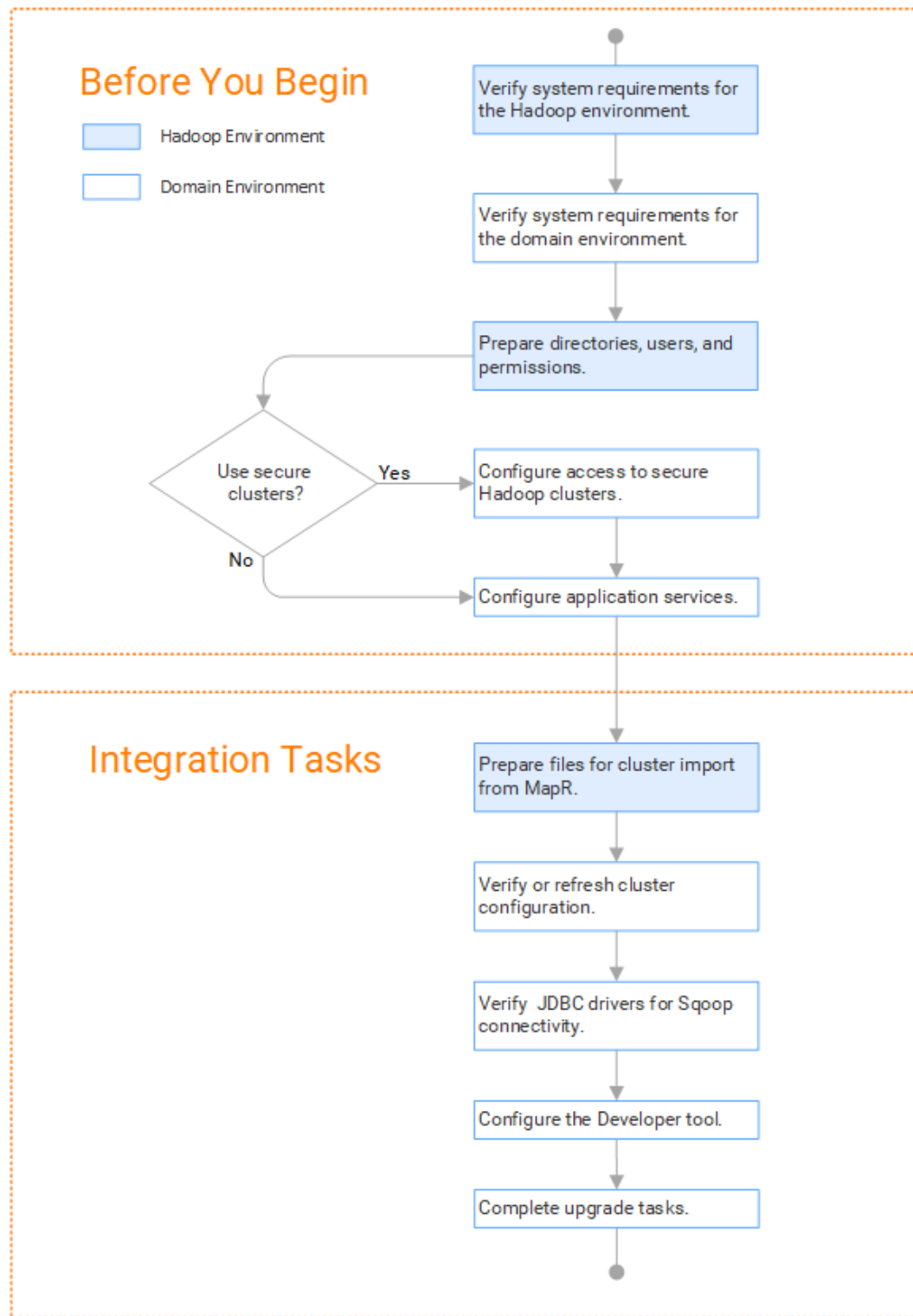
## Task Flow to Integrate with MapR (copy)

The following diagram shows the task flow to integrate the Informatica domain with MapR:



## Task Flow to Upgrade from Version 10.2.2 (mapr) (copy)

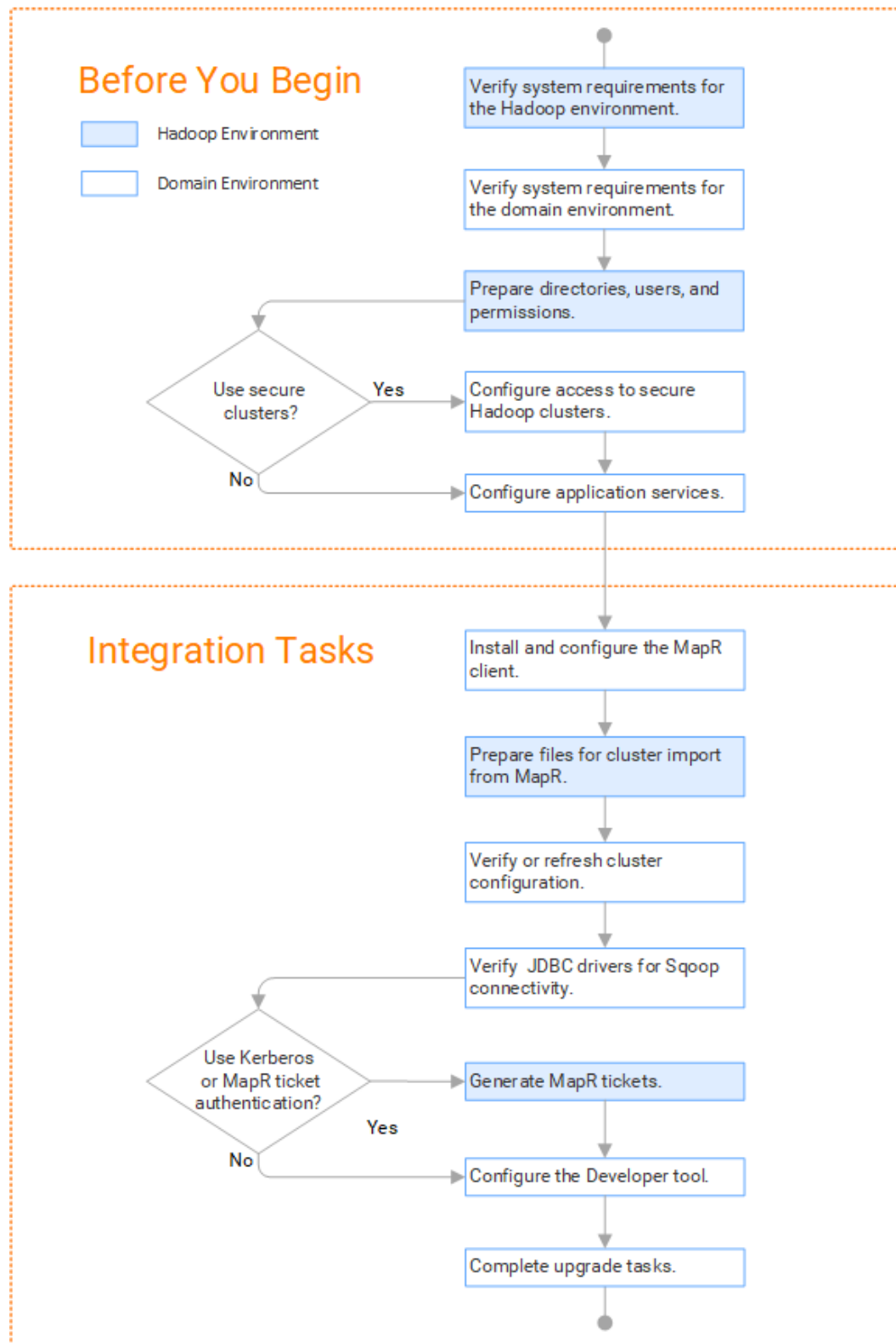
The following diagram shows the task flow to upgrade version 10.2.1 for MapR:





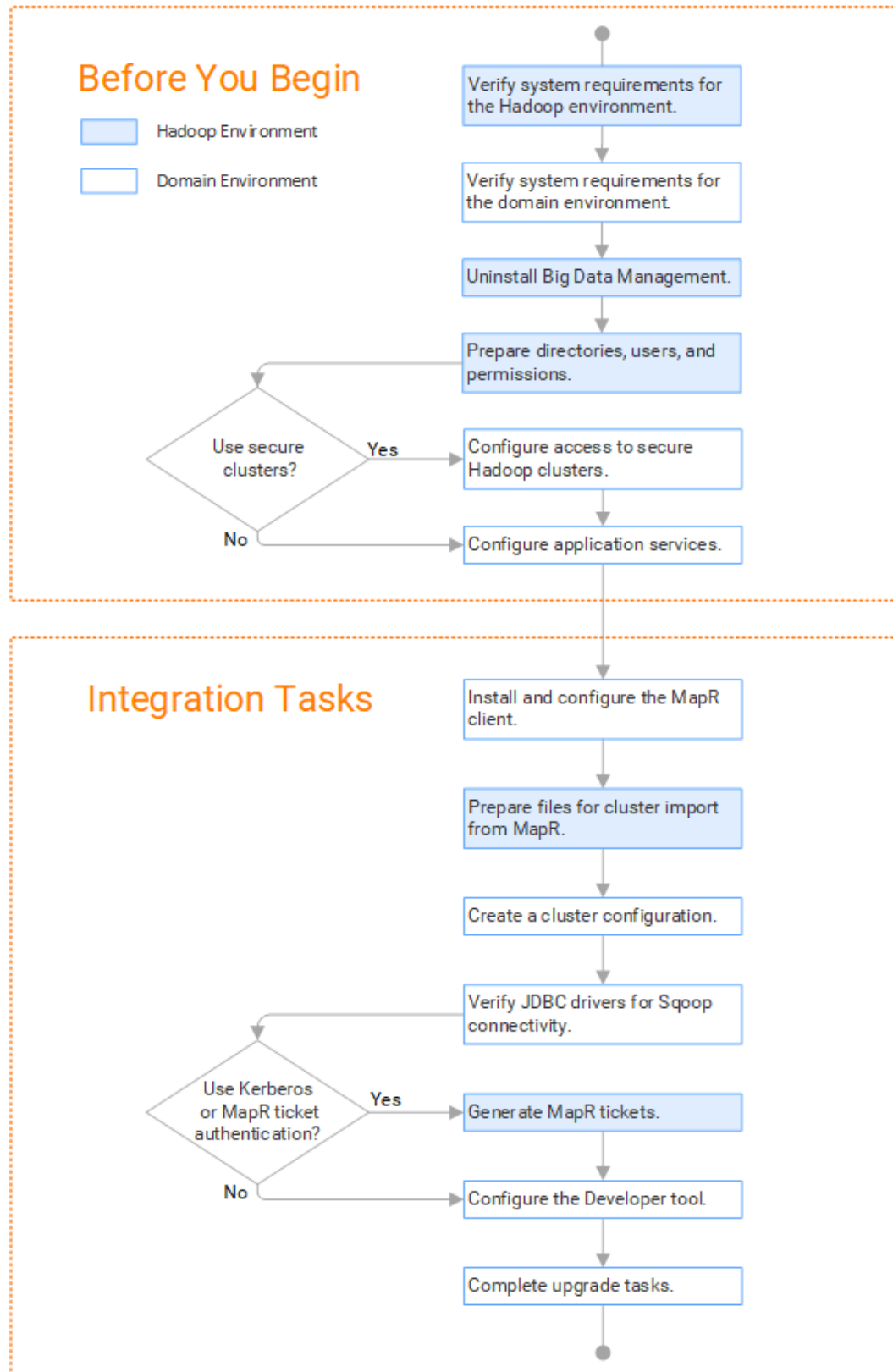
## Task Flow to Upgrade from Version 10.2 (mapr) (copy)

The following diagram shows the task flow to upgrade version 10.2 for MapR:



## Task Flow to Upgrade from a Version Earlier than 10.2 (mapr) (copy)

The following diagram shows the task flow to upgrade from a version earlier than 10.2 for MapR:



# Install and Configure the MapR Client

To enable communication between the Informatica domain and the MapR cluster, install and configure the MapR client on the application service machines. The MapR client version on the MapR cluster and the application service machines must match.

<b>Perform this task in the following situations:</b>
<ul style="list-style-type: none"><li>- You are integrating for the first time.</li><li>- You upgraded from version 10.2 or earlier.</li></ul>



You install the MapR client on the Data Integration Service, Metadata Access Service, and Analyst Service machines in the following directory:

`/opt/mapr`

For instructions about installing and configuring the MapR client, refer to the MapR documentation at <https://mapr.com/docs/60/AdvancedInstallation/SettingUptheClient-install-mapr-client.html>.

## Prepare for Cluster Import from MapR

Before the Informatica administrator can import cluster information to create a cluster configuration in the Informatica domain, the Hadoop administrator must perform some preliminary tasks.

<b>Perform this task in the following situations:</b>
<ul style="list-style-type: none"><li>- You are integrating for the first time.</li><li>- You upgraded from any previous version.</li></ul>



**Note:** If you are upgrading from a previous version, verify the properties and suggested values, as Data Engineering Integration might require additional properties or different values for existing properties.

Complete the following tasks to prepare the cluster before the Informatica administrator creates the cluster configuration:

1. Verify property values in \*-site.xml files that Data Engineering Integration needs to run mappings in the Hadoop environment.
2. Prepare the archive file to import into the domain.

**Note:** You cannot import cluster information directly from the MapR cluster into the Informatica domain.

## Configure \*-site.xml Files for MapR

The Hadoop administrator needs to configure \*-site.xml file properties and restart impacted services before the Informatica administrator imports cluster information into the domain.

### capacity-scheduler.xml

Configure the following properties in the capacity-scheduler.xml file:

#### **yarn.scheduler.capacity.<queue path>.disable\_preemption**

Disables preemption for the Capacity Scheduler. The Blaze engine does not support preemption. If YARN preemption is enabled for the cluster, you need to disable it for the queue allocated to the Blaze engine.

Set to TRUE for queue allocated to the Blaze engine.

### core-site.xml

Configure the following properties in the core-site.xml file:

#### **fs.s3.enableServerSideEncryption**

Enables server side encryption for S3 buckets. Required for SSE and SSE-KMS encryption.

Set to: TRUE

#### **fs.s3a.access.key**

The ID for the Blaze and Spark engines to connect to the Amazon S3 file system.

Set to your access key.

#### **fs.s3a.secret.key**

The password for the Blaze and Spark engines to connect to the Amazon S3 file system

Set to your access ID.

#### **fs.s3a.server-side-encryption-algorithm**

The server-side encryption algorithm for S3. Required for SSE and SSE-KMS encryption. Set to the encryption algorithm used.

#### **hadoop.proxyuser.<proxy user>.groups**

Defines the groups that the proxy user account can impersonate. On a secure cluster the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to group names of impersonation users separated by commas. If less security is preferred, use the wildcard " \* " to allow impersonation from any group.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.proxyuser.<proxy user>.hosts**

Defines the host machines that a user account can impersonate. On a secure cluster the <proxy user> is the Service Principal Name that corresponds to the cluster keytab file. On a non-secure cluster, the <proxy user> is the system user that runs the Informatica daemon.

Set to a single host name or IP address, or set to a comma-separated list. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.proxyuser.hive.hosts**

Comma-separated list of hosts that you want to allow the Hive user to impersonate on a non-secure cluster.

When `hive.server2.enable.doAs` is false, append a comma-separated list of Informatica server host names or IP address where the Data Integration Service is running. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to this property, restart the cluster services that use core-site configuration values.

#### **hadoop.proxyuser.yarn.groups**

Comma-separated list of groups that you want to allow the YARN user to impersonate on a non-secure cluster.

Set to group names of impersonation users separated by commas. If less security is preferred, use the wildcard " \* " to allow impersonation from any group.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **hadoop.proxyuser.yarn.hosts**

Comma-separated list of hosts that you want to allow the YARN user to impersonate on a non-secure cluster.

Set to a single host name or IP address, or set to a comma-separated list. If less security is preferred, use the wildcard " \* " to allow impersonation from any host.

**Note:** After you make changes to proxyuser properties, restart the credential service and other cluster services that use core-site configuration values.

#### **io.compression.codecs**

Enables compression on temporary staging tables.

Set to a comma-separated list of compression codec classes on the cluster.

#### **hadoop.security.auth\_to\_local**

Translates the principal names from the Active Directory and MIT realm into local names within the Hadoop cluster. Based on the Hadoop cluster used, you can set multiple rules.

Set to: `RULE:[1:$1@$0](^.*@YOUR.REALM)s/^(.*)@YOUR.REALM\.COM$/S1/g`

Set to: `RULE:[2:$1@$0](^.*@YOUR.REALM\.$)s/^(.*)@YOUR.REALM\.COM$/S1/g`

#### **fair-scheduler.xml**

Configure the following properties in the fair-scheduler.xml file:

##### **allowPreemptionFrom**

Enables preemption for the Fair Scheduler. The Blaze engine does not support preemption. If YARN preemption is enabled for the cluster, you need to disable it for the queue allocated to the Blaze engine.

Set to FALSE for the queue allocated to the Blaze engine.

For example:

```
<queue name="Blaze">
  <weight>1.0</weight>
  <allowPreemptionFrom>>false</allowPreemptionFrom>
  <schedulingPolicy>fsp</schedulingPolicy>
  <aclSubmitApps>*</aclSubmitApps>
  <aclAdministerApps>*</aclAdministerApps>
</queue>
```

## hbase-site.xml

Configure the following properties in the hbase-site.xml file:

### **zookeeper.znode.parent**

Identifies HBase master and region servers.

Set to the relative path to the znode directory of HBase.

## hive-site.xml

Configure the following properties in the hive-site.xml file:

### **hive.cluster.delegation.token.store.class**

The token store implementation. Required for HiveServer2 high availability and load balancing.

Set to: org.apache.hadoop.hive.thrift.DBTokenStore

### **hive.compactor.initiator.on**

Runs the initiator and cleaner threads on metastore instance. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: TRUE

### **hive.compactor.worker.threads**

The number of worker threads to run in a metastore instance. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: 1

### **hive.enforce.bucketing**

Enables dynamic bucketing while loading to Hive. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: TRUE

### **hive.exec.dynamic.partition**

Enables dynamic partitioned tables for Hive tables. Applicable for Hive versions 0.9 and earlier.

Set to: TRUE

### **hive.exec.dynamic.partition.mode**

Allows all partitions to be dynamic. Required for the Update Strategy transformation in a mapping that writes to a Hive target. Also required if you use Sqoop and define a DDL query to create or replace a partitioned Hive target at run time.

Set to: nonstrict

### **hive.support.concurrency**

Enables table locking in Hive. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: TRUE

### **hive.server2.support.dynamic.service.discovery**

Enables HiveServer2 dynamic service discovery. Required for HiveServer2 high availability.

Set to: TRUE

### **hive.server2.zookeeper.namespace**

The value of the ZooKeeper namespace in the JDBC connection string. Required for HiveServer2 high availability.

Set to: jdbc:hive2://<zookeeper\_ensemble>/  
default;serviceDiscoveryMode=zooKeeper;zooKeeperNamespace=hiveserver2

#### **hive.txn.manager**

Turns on transaction support. Required for an Update Strategy transformation in a mapping that writes to a Hive target.

Set to: org.apache.hadoop.hive.ql.lockmgr.DbTxnManager

#### **hive.zookeeper.quorum**

Comma-separated list of ZooKeeper server host:ports in a cluster. The value of the ZooKeeper ensemble in the JDBC connection string. Required for HiveServer2 high availability.

Set to: jdbc:hive2://<zookeeper\_ensemble>/default;serviceDiscoveryMode=zooKeeper;

### [mapred-site.xml](#)

Configure the following properties in the mapred-site.xml file:

#### **mapreduce.framework.name**

The run-time framework to run MapReduce jobs. Values can be local, classic, or yarn. Required for Sqoop.

Set to: yarn

#### **mapreduce.jobhistory.address**

Location of the MapReduce JobHistory Server. The default port is 10020. Required for Sqoop.

Set to: <MapReduce JobHistory Server>:<port>

#### **yarn.app.mapreduce.am.staging-dir**

The HDFS staging directory used while submitting jobs.

Set to the staging directory path.

### [yarn-site.xml](#)

Configure the following properties in the yarn-site.xml file:

#### **yarn.application.classpath**

Required for dynamic resource allocation.

Add spark\_shuffle.jar to the class path. The .jar file must contain the class "org.apache.spark.network.yarn.YarnShuffleService."

#### **yarn.nodemanager.resource.memory-mb**

The maximum RAM available for each container. Set the maximum memory on the cluster to increase resource memory available to the Blaze engine.

Set to 16 GB if value is less than 16 GB.

#### **yarn.nodemanager.resource.cpu-vcores**

The number of virtual cores for each container. Required for Blaze engine resource allocation.

Set to 10 if the value is less than 10.

#### **yarn.scheduler.minimum-allocation-mb**

The minimum RAM available for each container. Required for Blaze engine resource allocation.

Set to 6 GB if the value is less than 6 GB.

**yarn.nodemanager.vmem-check-enabled**

Disables virtual memory limits for containers. Required for the Blaze and Spark engines.

Set to: false

**yarn.nodemanager.aux-services**

Required for dynamic resource allocation for the Spark engine.

Add an entry for "spark\_shuffle."

**yarn.nodemanager.aux-services.spark\_shuffle.class**

Required for dynamic resource allocation for the Spark engine.

Set to: org.apache.spark.network.yarn.YarnShuffleService

**yarn.resourcemanager.scheduler.class**

Defines the YARN scheduler that the Data Integration Service uses to assign resources.

Set to: org.apache.hadoop.yarn.server.resourcemanager.scheduler

**yarn.node-labels.enabled**

Enables node labeling.

Set to: TRUE

**yarn.node-labels.fs-store.root-dir**

The HDFS location to update node label dynamically.

Set to: <hdfs://[Node name]:[Port]/[Path to store]/[Node labels]/>

## Prepare the Archive File for Import from MapR

After you verify property values in the \*-site.xml files, create a .zip or a .tar file that the Informatica administrator can use to import the cluster configuration into the domain.

Create an archive file that contains the following files from the cluster:

- core-site.xml
- hbase-site.xml. Required only if you access HBase sources and targets.
- hive-site.xml
- mapred-site.xml
- yarn-site.xml

**Note:** To import from MapR, the Informatica administrator must use an archive file.



# Create a Cluster Configuration

After the Hadoop administrator prepares the cluster for import, the Informatica administrator must create a cluster configuration.

## Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from version 10.1.1 or earlier.

A cluster configuration is an object in the domain that contains configuration information about the Hadoop cluster. The cluster configuration enables the Data Integration Service to push mapping logic to the Hadoop environment. Import configuration properties from the Hadoop cluster to create a cluster configuration.

The import process imports values from \*-site.xml files into configuration sets based on the individual \*-site.xml files. When you perform the import, the cluster configuration wizard can create Hadoop, HBase, HDFS, and Hive connection to access the Hadoop environment. If you choose to create the connections, the wizard also associates the cluster configuration with the connections.

**Note:** If you are integrating for the first time and you imported the cluster configuration when you ran the installer, you *must* re-create or refresh the cluster configuration.

## Importing a Hadoop Cluster Configuration from a File

You can import properties from an archive file to create a cluster configuration.

Before you import from the cluster, you must get the archive file from the Hadoop administrator.

1. From the **Connections** tab, click the **ClusterConfigurations** node in the Domain Navigator.
2. From the Actions menu, select **New > Cluster Configuration**.

The **Cluster Configuration** wizard opens.

3. Configure the following properties:

Property	Description
Cluster configuration name	Name of the cluster configuration.
Description	Optional description of the cluster configuration.
Distribution type	The cluster Hadoop distribution type.

Property	Description
Method to import the cluster configuration	Choose <b>Import from file</b> to import properties from an archive file.
Create connections	<p>Choose to create Hadoop, HDFS, Hive, and HBase connections.</p> <p>If you choose to create connections, the Cluster Configuration wizard associates the cluster configuration with each connection that it creates.</p> <p>The Hadoop connection contains default values for properties such as cluster environment variables, cluster path variables, and advanced properties. Based on the cluster environment and the functionality that you use, you can add to the default values or change the default values of these properties. For a list of Hadoop connection properties to configure, see <a href="#">“Configuring Hadoop Connection Properties” on page 269</a>.</p> <p>If you do not choose to create connections, you must manually create them and associate the cluster configuration with them.</p> <p><b>Important:</b> When the wizard creates the Hive connection, it populates the Metadata Connection String and the Data Access Connection String properties with the value from the hive.metastore.uris property. If the Hive metastore and HiveServer2 are running on different nodes, you must update the Metadata Connection String to point to the HiveServer2 host.</p>

- Click **Browse** to select a file. Select the file and click **Open**.
- Click **Next** and verify the cluster configuration information on the summary page.

## Verify or Refresh the Cluster Configuration

You might need to refresh the cluster configuration or update the distribution version in the cluster configuration when you upgrade.

### Perform this task in the following situation:

- You upgraded from version 10.2 or later.

### Verify the Cluster Configuration

The cluster configuration contains a property for the distribution version. The verification task depends on the version you upgraded:

#### Upgrade from 10.2

If you upgraded from 10.2 and you changed the distribution version, you need to verify the distribution version in the General properties of the cluster configuration.

#### Upgrade from 10.2.1

Effective in version 10.2.1, Informatica assigns a default version to each Hadoop distribution type. If you configure the cluster configuration to use the default version, the upgrade process upgrades to the assigned default version if the version changes. If you have not upgraded your Hadoop distribution to Informatica's default version, you need to update the distribution version property.

For example, suppose the assigned default Hadoop distribution version for 10.2.1 is  $n$ , and for 10.2.2 is  $n+1$ . If the cluster configuration uses the default supported Hadoop version of  $n$ , the upgraded cluster

configuration uses the default version of  $n+1$ . If you have not upgraded the distribution in the Hadoop environment you need to change the cluster configuration to use version  $n$ .

If you configure the cluster configuration to use a distribution version that is not the default version, you need to update the distribution version property in the following circumstances:

- Informatica dropped support for the distribution version.
- You changed the distribution version.

### Refresh the Cluster Configuration

If you updated any of the \*-site.xml files noted in the topic to prepare for cluster import, you need to refresh the cluster configuration in the Administrator tool.

## Verify JDBC Drivers for Sqoop Connectivity

Verify that you have the JDBC drivers to access JDBC-compliant databases in the Hadoop environment. You might need separate drivers for metadata import and for run-time processing.

Perform this task in the following situations:
<ul style="list-style-type: none"><li>- You are integrating for the first time.</li><li>- You upgraded from version 10.2.1 or earlier.</li></ul>



You download drivers based on design-time and run-time requirements:

- **Design-time.** To import metadata, you can use the DataDirect drivers packaged with the Informatica installer if they are available. If they are not available, use any Type 4 JDBC driver that the database vendor recommends.
- **Run-time.** To run mappings, use any Type 4 JDBC driver that the database vendor recommends. Some distributions support other drivers to use Sqoop connectors. You cannot use the DataDirect drivers for run-time processing.

## Verify Design-time Drivers

Use the DataDirect JDBC drivers packaged with the Informatica installer to import metadata from JDBC-compliant databases. If the DataDirect JDBC drivers are not available for a specific JDBC-compliant database, download the Type 4 JDBC driver associated with that database.

Copy the JDBC driver .jar files to the following location on the Developer tool machine:

```
<Informatica installation directory>\clients\externaljdbcjars
```

## Verify Run-time Drivers

Verify run-time drivers for mappings that access JDBC-compliant databases in the Hadoop environment. Use any Type 4 JDBC driver that the database vendor recommends.

1. Download Type 4 JDBC drivers associated with the JCBC-compliant databases that you want to access.

2. To use Sqoop TDCH MapR Connector for Teradata, download the following files:
  - sqoop-connector-tdch-1.1-mapr-1707.jar from <https://repository.mapr.com/nexus/content/groups/mapr-public/org/apache/sqoop/connector/sqoop-connector-tdch/>
  - terajdbc4.jar and tdgssconfig.jar from <http://downloads.teradata.com/download/connectivity/jdbc-driver>
  - The MapR Connector for Teradata .jar file from the Teradata website.
3. To optimize the Sqoop mapping performance on the Spark engine while writing data to an HDFS complex file target of the Parquet format, download the following .jar files:
  - parquet-hadoop-bundle-1.6.0.jar from <https://repo1.maven.org/maven2/com/twitter/parquet-hadoop-bundle/1.6.0/>
  - parquet-avro-1.6.0.jar from <https://repo1.maven.org/maven2/com/twitter/parquet-avro/1.6.0/>
  - parquet-column-1.5.0.jar from <https://repo1.maven.org/maven2/com/twitter/parquet-column/1.5.0/>
4. Copy all of the .jar files to the following directory on the machine where the Data Integration Service runs:  
`<Informatica installation directory>\externaljdbcjars`  
Changes take effect after you recycle the Data Integration Service. At run time, the Data Integration Service copies the .jar files to the Hadoop distribution cache so that the .jar files are accessible to all nodes in the cluster.

## Generate MapR Tickets

To run mappings on a MapR cluster that uses Kerberos or MapR Ticket authentication with information in Hive tables, generate a MapR ticket for the Data Integration Service user.

The Data Integration Service user requires an account on the MapR cluster and a MapR ticket on the application service machines that require access to MapR. When the MapR cluster uses both Kerberos and Ticket authentication, you generate a single ticket for the Data Integration Service user for both authentication systems.

After you generate and save MapR tickets, you perform additional steps to configure the Data Integration Service, the Metadata Access Service, and the Analyst Service to communicate with the MapR cluster.

## Generate Tickets

After you create a MapR user account for the Data Integration Service user, name the ticket file using the following naming convention:

```
maprticket_<user name>
```

For example, for a user ID 1234, a MapR ticket file named maprticket\_1234 is generated.

### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from version 10.1.1 or earlier.

Save the ticket on the machines that host the Data Integration Service, the Metadata Access Service, and the Analyst Service. The Data Integration Service and the Analyst Service access the ticket at run time. The Metadata Access Service access the ticket for the Developer tool at design time.

By default, the services access the ticket in the /tmp directory. If you save the ticket to any other location, you must configure the MAPR\_TICKETFILE\_LOCATION environment variable in the service properties.

## Configure the Data Integration Service

When the MapR cluster is secured with Kerberos or MapR Ticket authentication, edit Data Integration Service properties to enable communication between the Informatica domain and the cluster.

### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from version 10.1.1 or earlier.

In the Administrator tool Domain Navigator, select the Data Integration Service to configure, and then select the **Processes** tab.

In the **Environment Variables** area, configure the following property to define the Kerberos authentication protocol:

Property	Value
JAVA_OPTS	<code>-Dhadoop.login=&lt;MAPR_ECOSYSTEM_LOGIN_OPTS&gt; -Dhttps.protocols=TLSv1.2</code>  where <MAPR_ECOSYSTEM_LOGIN_OPTS> is the value of the MAPR_ECOSYSTEM_LOGIN_OPTS property in the file <code>/opt/mapr/conf/env.sh</code> .
MAPR_HOME	MapR client directory on the machine that runs the Data Integration Service. For example, <code>opt/mapr</code> Required if you want to fetch a MapR Streams data object.
MAPR_TICKETFILE_LOCATION	Required when the MapR cluster uses Kerberos or MapR Ticket authentication. Location of the MapR ticket file if you saved it to a directory other than /tmp. For example:  <code>/export/home/username1/Keytabs_and_krb5conf/Tickets/project1/maprticket_30103</code>

Changes take effect when you restart the Data Integration Service.

## Configure the Metadata Access Service

When the MapR cluster is secured with MapR Kerberos or ticketed authentication, edit Metadata Access Service properties to enable communication between the Developer tool and the cluster.

### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from version 10.2 or earlier.

In the Administrator tool Domain Navigator, select the Metadata Access Service to configure, and then select the **Processes** tab.

In the **Environment Variables** area, configure the following property to define the Kerberos authentication protocol:

Property	Value
JAVA_OPTS	<code>-Dhadoop.login=&lt;MAPR_ECOSYSTEM_LOGIN_OPTS&gt; -Dhttps.protocols=TLSv1.2</code>  where <MAPR_ECOSYSTEM_LOGIN_OPTS> is the value of the MAPR_ECOSYSTEM_LOGIN_OPTS property in the file <code>/opt/mapr/conf/env.sh</code> .
MAPR_TICKETFILE_LOCATION	Required when the MapR cluster uses Kerberos or MapR Ticket authentication. Location of the MapR ticket file if you saved it to a directory other than <code>/tmp</code> . For example,  <code>/export/home/username1/Keytabs_and_krb5conf/Tickets/project1/maprticket_30103</code>

Changes take effect when you restart the Metadata Access Service.

## Configure the Analyst Service

If you use the Analyst tool to profile data in Hive data objects, configure properties on the Analyst Service to enable communication between the Analyst tool and the cluster, including testing of the Hive connection.

### Perform this task in the following situations:

- You are integrating for the first time.
- You upgraded from version 10.1.1 or earlier.

In the Administrator tool Domain Navigator, select the Analyst Service to configure, then select the **Processes** tab.

In the **Environment Variables** area, configure the following property to define the Kerberos authentication protocol:

Property	Value
JAVA_OPTS	-Dhadoop.login=hybrid -Dhttps.protocols=TLSv1.2
MAPR_TICKETFILE_LOCATION	Required when the MapR cluster uses Kerberos or MapR Ticket authentication. Location of the MapR ticket file if you saved it to a directory other than /tmp. For example,  /export/home/username1/Keytabs_and_krb5conf/Tickets/project1/maprticket_30103
LD_LIBRARY_PATH	The location of Hadoop libraries. For example,  <Informatica installation directory>/java/jre/lib:<Informatica installation directory>/services/shared/bin:<Informatica installation directory>/server/bin:<Informatica installation directory>/services/shared/hadoop/<MapR location>/lib/native/Linux-amd64-64

Changes take effect when you restart the Analyst Service.

## Complete Upgrade Tasks

If you upgraded the Informatica platform from version 10.2.2, update connections.

### Update Connections

You might need to update connections based on the version you are upgrading from.

If you did not create connections when you created the cluster configuration, you need to update the connections.

### Configure the Hadoop Connection

To use properties that you customized in the `hadoopEnv.properties` file, you must configure the Hadoop connection properties such as cluster environment variables, cluster path variables, and advanced properties.

#### Perform this task in the following situation:

- You upgraded from version 10.1.1 or earlier.

When you run the Informatica upgrade, the installer backs up the existing `hadoopEnv.properties` file. You can find the backup `hadoopEnv.properties` file in the following location:

```
<Previous Informatica installation directory>/services/shared/hadoop/<Hadoop distribution name>_<version>/infaConf
```

Edit the Hadoop connection in the Administrator tool or the Developer tool to include any properties that you manually configured in the `hadoopEnv.properties` file. The Hadoop connection contains default values for properties such as cluster environment and path variables and advanced properties. You can update the default values to match the properties in the `hadoopEnv.properties` file.

## Replace the Connections with New Connections

If you created connections when you imported the cluster configuration, you need to replace connections in mappings with the new connections.

### Perform this task in the following situation:

- You upgraded from version 10.1.1 or earlier.

The method that you use to replace connections in mappings depends on the type of connection.

#### Hadoop connection

Run the following commands to replace the connections:

- `infacmd dis replaceMappingHadoopRuntimeConnections`. Replaces connections associated with mappings that are deployed in applications.
- `infacmd mrs replaceMappingHadoopRuntimeConnections`. Replaces connections associated with mappings that you run from the Developer tool.

For information about the `infacmd` commands, see the *Informatica Command Reference*.

#### Hive, HDFS, and HBase connections

You must replace the connections manually.

## Complete Connection Upgrade

If you *did not* create connections when you imported the cluster configuration, you must update connection properties for Hadoop, Hive, HDFS, and HBase connections.

### Perform this task in the following situations:

- You upgraded from version 10.2.2 or earlier.

Perform the following tasks to update the connections:

#### Update changed properties

Review connections that you created in a previous release to update the values for connection properties. For example, if you added nodes to the cluster or if you updated the distribution version, you might need to verify host names, URIs, or port numbers for some of the properties.

#### Associate the cluster configuration

The Hadoop, Hive, HDFS, and HBase connections must be associated with a cluster configuration. Complete the following tasks:

1. Run `infacmd isp listConnections` to identify the connections that you need to upgrade. Use `-ct` to list connections of a particular type.
2. Run `infacmd isp UpdateConnection` to associate the cluster configuration with the connection. Use `-cn` to name the connection and `-o clusterConfigID` to associate the cluster configuration with the connection.



For more information about `infacmd`, see the *Informatica Command Reference*.

## Replace Hive Run-time Connections with Hadoop Connections

Effective in version 10.2.2, Big Data Management dropped support for the Hive engine and Hive run-time connections. If you used Hive connections to run mappings on the Hadoop cluster, you must generate Hadoop connections from the Hive connections.

### Perform this task in the following situations:

- You upgraded from version 10.1.1 or earlier.
- The Hive connections are configured to run mappings in the Hadoop environment.

The upgrade process generates a connection name for the Hadoop connection and replaces the connection name in the mappings. It does not create the physical connection object. When the upgrade is complete, you must run a command to generate the connection. You generate Hadoop connections from Hive connections that are configured to run mappings in the Hadoop environment.

You must generate Hadoop connections from Hive connections that are configured to run mappings in the Hadoop environment.

1. Run `infacmd isp generateHadoopConnectionFromHiveConnection` to generate a Hadoop connection from a Hive connection that is configured to run in the Hadoop environment.  
The command names the connection as follows: "Autogen\_<Hive connection name>." If the connection name exceeds the 128 character limit, the command fails.
2. If the command fails, perform the following tasks:
  - a. Rename the connection to meet the character limit and run the command again.
  - b. Run `infacmd dis replaceMappingHadoopRuntimeConnections` to replace connections associated with mappings that are deployed in applications.
  - c. Run `infacmd mrs replaceMappingHadoopRuntimeConnections` to replace connections associated with mappings that you run from the Developer tool.
3. If the Hive connection was parameterized, you must update the connection names in the parameter file. Verify that the Hive sources, Hive targets, and the Hive engine parameters are updated with the correct connection name.
4. If any properties changed in the cluster, such as host names, URIs, or port numbers, you must update the properties in the connections.

For information about the `infacmd` commands, see the *Informatica® Command Reference*.

# Part II: Databricks Integration

This part contains the following chapters:

- [Introduction to Databricks Integration, 187](#)
- [Before You Begin Databricks Integration, 191](#)
- [Databricks Integration Tasks, 203](#)

## CHAPTER 10

# Introduction to Databricks Integration

This chapter includes the following topics:

- [Databricks Integration Overview, 187](#)
- [Run-time Process on the Databricks Spark Engine, 188](#)
- [Databricks Integration Task Flow, 189](#)

## Databricks Integration Overview

Data Engineering Integration can connect to Databricks on Azure and on AWS. Databricks is an analytics cloud platform that you can use with Microsoft Azure cloud services or Amazon Web Services. Databricks incorporates the open-source Apache Spark cluster technologies and capabilities.

The Data Integration Service installs the binaries required to integrate the Informatica domain with the Databricks environment. The integration requires Informatica connection objects and cluster configurations. A cluster configuration is a domain object that contains configuration parameters that you import from the Databricks cluster. You then associate the cluster configuration with connections to access the Databricks environment.

**Note:** With the following exceptions, all of the functionality described in this article applies to Informatica 10.4 and later releases:

- Informatica 10.5 adds support for warm pool access and the Sequence Generator transformation.
- Informatica 10.5.2 adds support for Databricks schema evolution and custom parameters.

For more information about these features, see the *Data Engineering Integration User Guide*.

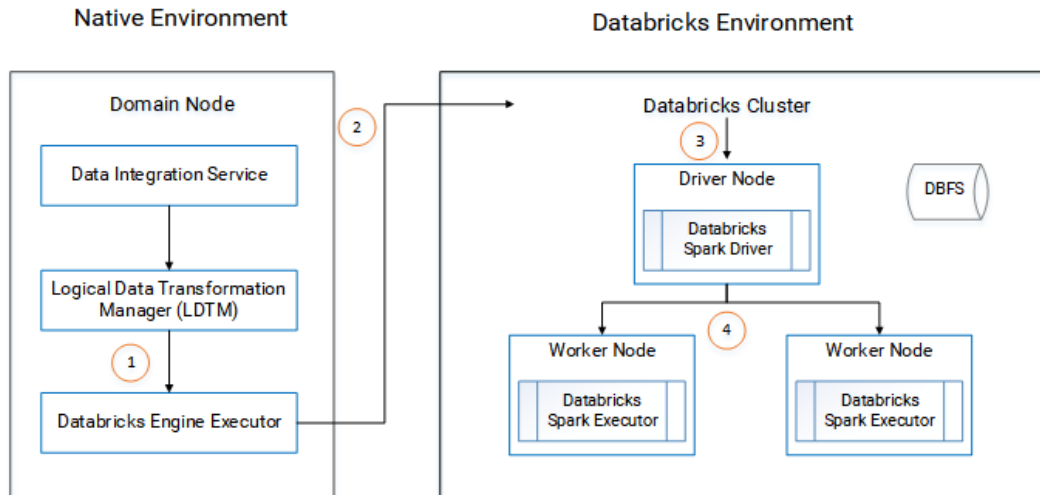
Perform the following tasks to integrate the Informatica domain with the Databricks environment:

1. Install or upgrade to the current Informatica version.
2. Perform pre-import tasks, such as verifying system requirements and permissions.
3. Import the cluster configuration into the domain.
4. Create a Databricks connection to run mappings within the Databricks environment.

# Run-time Process on the Databricks Spark Engine

When you run a job on the Databricks Spark engine, the Data Integration Service pushes the processing to the Databricks cluster, and the Databricks Spark engine runs the job.

The following image shows the components of the Informatica and the Databricks environments:



1. The Logical Data Transformation Manager translates the mapping into a Scala program, packages it as an application, and sends it to the Databricks Engine Executor on the Data Integration Service machine.
2. The Databricks Engine Executor submits the application through REST API to the Databricks cluster, requests to run the application, and stages files for access during run time.
3. The Databricks cluster passes the request to the Databricks Spark driver on the driver node.
4. The Databricks Spark driver distributes the job to one or more Databricks Spark executors that reside on worker nodes.
5. The executors run the job and stage run-time data to the Databricks File System (DBFS) of the workspace.

## Native Environment

The integration with Databricks requires tools, services, and a repository database in the Informatica domain.

### Clients and Tools

When the Informatica domain is integrated with Databricks, you can use the following tools:

#### Informatica Administrator

Use the Administrator tool to manage the Informatica domain and application services. You can also create objects such as connections, cluster configurations, and cloud provisioning configurations to enable data engineering operations.

#### The Developer tool

Use the Developer tool to import sources and targets and create mappings to run in the Databricks environment.

### Application Services

The domain integration with Databricks uses the following services:

**Data Integration Service**

The Data Integration Service can process mappings in the native environment, or it can push the processing to the Databricks environment. The Data Integration Service retrieves metadata from the Model repository when you run a mapping.

**Model Repository Service**

The Model Repository Service manages the Model repository. All requests to save or access Model repository metadata go through the Model repository.

**Model Repository**

The Model repository stores mappings that you create and manage in the Developer tool.

## Databricks Environment

Integration with the Databricks environment includes the following components:

**Databricks Spark engine**

The Databricks run-time engine based on the open-source Apache Spark engine.

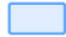

**Databricks File System (DBFS)**

A distributed file system installed on Databricks Runtime clusters. Run-time data is staged in the DBFS and is persisted to a mounted Blob storage container.

## Databricks Integration Task Flow

The following diagram shows the task flow to integrate the Informatica domain with Azure Databricks:

## Before You Begin

-  Databricks Environment
-  Domain Environment

Verify system requirements for the Databricks environment.

Verify system requirements for the domain environment.

Configure preemption for concurrent jobs.

Configure storage access.

Prepare staging directories.

Prepare for token authentication.

Configure the Data Integration Service.

## Integration Tasks

Create a Databricks cluster configuration.

Configure the Databricks connection.

## CHAPTER 11

# Before You Begin Databricks Integration

This chapter includes the following topics:

- [Read the Release Notes, 191](#)
- [Verify System Requirements, 191](#)
- [Configure Preemption for Concurrent Jobs, 192](#)
- [Configure Storage Access, 193](#)
- [Create a Staging Directory for Binary Archive Files, 200](#)
- [Create a Staging Directory for Run-time Processing, 201](#)
- [Prepare for Token Authentication, 201](#)
- [Configure the Data Integration Service, 201](#)
- [Install Python Libraries, 202](#)

## Read the Release Notes

Read the Release Notes for updates to the installation and upgrade process. You can also find information about known and fixed issues for the release.

Find the Release Notes on the Informatica [documentation portal](#).

## Verify System Requirements

Verify that your environment meets the following minimum system requirements for the integration:

### **Informatica services**

Install and configure the Informatica services and the Developer tool. Verify that the domain contains a Model Repository Service and a Data Integration Service.

### Domain access to Databricks

Verify access domain access to Databricks through one of the following methods:

- VPN is enabled between the Informatica domain and the Azure cloud network.
- The Informatica domain and Databricks are installed on the same Azure vnet or AWS VPC.
- The Informatica domain and Databricks are installed on different Azure vnets or AWS VPCs with peering enabled.

For more information about product requirements and supported platforms, see the Product Availability Matrix on Informatica Network:

<https://network.informatica.com/community/informatica-network/product-availability-matrices>

## Configure Preemption for Concurrent Jobs

Configure the Databricks cluster to improve concurrency of jobs.

When you submit a job to Databricks, it allocates resources to run the job. If it does not have enough resources, it puts the job in a queue. Pending jobs fail if resources do not become available before the timeout of 30 minutes.

You can configure preemption on the cluster to control the amount of resources that Databricks allocates to each job, thereby allowing more jobs to run concurrently. You can also configure the timeout for the queue and the interval at which the Databricks Spark engine checks for available resources.

Configure the following Spark properties for the Databricks Spark engine:

### **spark.databricks.preemption.enabled**

Enables the Spark scheduler for preemption. Default is false.

Set to: true

### **spark.databricks.preemption.threshold**

A percentage of resources that are allocated to each submitted job. The job runs with the allocated resources until completion. Default is 0.5, or 50 percent.

Set to a value lower than default, such as 0.1.

### **spark.databricks.preemption.timeout**

The number of seconds that a job remains in the queue before failing. Default is 30.

Set to: 1,800.

**Note:** If you set a value higher than 1,800, Databricks ignores the value and uses the maximum timeout of 1,800.

### **spark.databricks.preemption.interval**

The number of seconds to check for available resources to assign to a job in the queue. Default is 5.

Set to a value lower than the timeout.

Changes take effect after you restart the cluster.

**Important:** Informatica integrates with Databricks, supporting standard concurrency clusters. Standard concurrency clusters have a maximum queue time of 30 minutes, and jobs fail when the timeout is reached. The maximum queue time cannot be extended. Setting the preemption threshold allows more jobs to run concurrently, but with a lower percentage of allocated resources, the jobs can take longer to run. Also, configuring the environment for preemption does not ensure that all jobs will run. In addition to configuring



preemption, you might choose to run cluster workflows to create ephemeral clusters that create the cluster, run the job, and then delete the cluster. For more information about Databricks concurrency, see the [Databricks documentation](#).

## Configure Storage Access

You can integrate Data Engineering Integration with Databricks on the Amazon Web Services (AWS) or Microsoft Azure platforms.

Based on the platform and on the storage type used by the Databricks cluster, you can configure the storage key or get the client credentials of the service principal to access cluster storage. Add the configuration to the Spark configuration on the Databricks cluster.

### Configure S3 and Redshift Authentication and Encryption on AWS

If the Databricks cluster uses the AWS S3 service to access sources or targets stored in S3 buckets or the Redshift data warehouse, configuration depends on your choices for authentication and encryption.

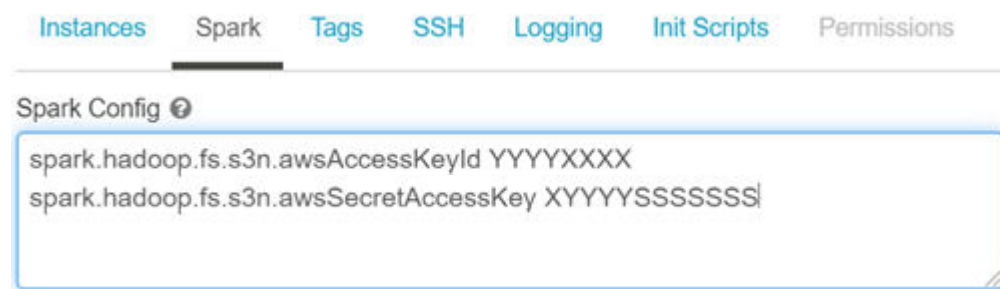
#### Authentication Options

Configure the following authentication options:

##### Key Pair Authentication

If the data resources use AWS key pair authentication, add the access key and the secret key in the Spark tab of the Advanced Configuration section of the Databricks cluster configuration page. Separate each key value with a space.

The image below shows an example of the keys in the text entry pane of the Spark tab:



##### IAM Role Authentication

If the data resources use IAM role authentication, verify that the configuration meets one of the following requirements:

- The S3 bucket belongs to the AWS account in which the Databricks cluster resides.
- The S3 bucket belongs to a different AWS account than the one in which the Databricks cluster resides, and you enabled a cross-account policy to allow the cluster to access the bucket.

For more information about using IAM roles with a Databricks cluster, see the [AWS documentation](#).

## Encryption Options

Choose from the following types of encryption. Each can be combined with either of the two options for authentication.

Configure the following properties in the Spark configuration tab of the cluster:

### Server-Side S3 Encryption (SSE-S3)

Configure the following property to enable SSE-S3 encryption:

```
spark.hadoop.fs.s3a.server-side-encryption-algorithm AES256
```

**Note:** If you use both key pair authentication and SSE-S3 encryption, then add this property in the Spark configuration tab after the first two lines for key pair authentication. For example:

```
spark.hadoop.fs.s3n.awsAccessKeyId YYYYYXXX
spark.hadoop.fs.s3n.awsSecretAccessKey YYYYYSSSSSS
spark.hadoop.fs.s3a.server-side-encryption-algorithm AES256
```

### Server-Side Encryption with KMS (SSE-KMS)

Configure the following properties to enable SSE-KMS encryption:

```
spark.hadoop.fs.s3a.server-side-encryption-kms-master-key-id arn:aws:kms:us-west-
XXXXX:key/XXXXXXXXXXXXXXXXXXXX
spark.hadoop.fs.s3a.server-side-encryption-algorithm aws:kms
spark.hadoop.fs.s3a.impl com.data bricks.s3a.S3AFileSystem
```

**Note:** If you use both key pair authentication and SSE-KMS encryption, then add these properties in the Spark configuration tab after the first two lines for key pair authentication. For example:

```
spark.hadoop.fs.s3n.awsAccessKeyId YYYYYXXX
spark.hadoop.fs.s3n.awsSecretAccessKey YYYYYSSSSSS
spark.hadoop.fs.s3a.server-side-encryption-kms-master-key-id arn:aws:kms:us-west-
XXXXX:key/XXXXXXXXXXXXXXXXXXXX
spark.hadoop.fs.s3a.server-side-encryption-algorithm aws:kms
spark.hadoop.fs.s3a.impl com.data bricks.s3a.S3AFileSystem
```

## Configure AWS Roles and Policies to Access S3 Resources

Configure AWS roles and policies to enable the Data Integration Service to access S3 resources and run mappings on the Databricks cluster.

Perform these steps when the Data Integration Service is deployed on the AWS platform.

The steps in this section are based on the [Databricks documentation](#).

### Step 1. Create an IAM Role and Policy for S3 Access

Using S3 resources with mappings that run on a Databricks cluster requires you to create IAM roles and policies in your AWS account.

1. Log in to the AWS account that has administrator access to the Databricks cluster that you want to integrate with Data Engineering Integration.
2. Optionally create a new IAM role to correspond to an S3 access policy. You can create a policy or use an existing one.

Use the following steps to create the IAM role:

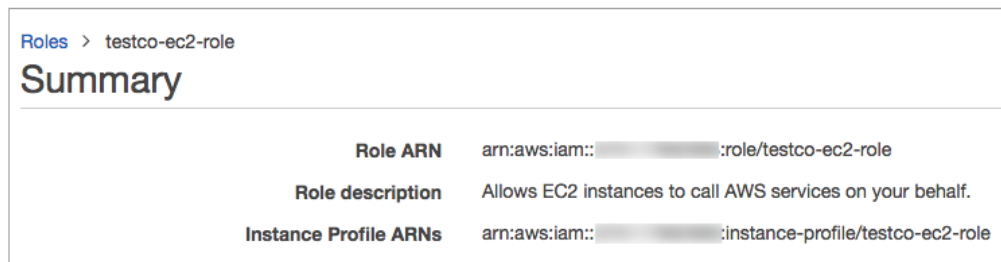
- a. Select the IAM service, then **Roles > Create Role**.  
The Create Roles wizard opens.
- b. Under **Select type of trusted entity**, select **AWS service**.

- c. Under **Choose the service that will use this role**, select **EC2**.
  - d. Click **Next: Permissions > Next: Tags > Next: Review**.
  - e. In the Role name field, type a role name, then click **Create role**.
3. Grant access to the S3 bucket that contains the resources for the mapping to access. To do this, attach an inline policy to the role that you want to use.
  - a. In the **Permissions** tab, click **Inline policy**.
  - b. Select the **JSON** tab.
  - c. Paste the following JSON statement:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "s3:ListBucket"
      ],
      "Resource": [
        "arn:aws:s3:::<s3-bucket-name>"
      ]
    },
    {
      "Effect": "Allow",
      "Action": [
        "s3:PutObject",
        "s3:GetObject",
        "s3:DeleteObject",
        "s3:PutObjectAcl"
      ],
      "Resource": [
        "arn:aws:s3:::<s3-bucket-name>/*"
      ]
    }
  ]
}
```

- d. Edit the pasted JSON statement with the name of the S3 bucket, then click **Review**.
  - e. In the **Name** field, type a name for the policy, then click **Create policy**.

The **Summary** window opens.
4. In the **Summary** window, select and copy the Instance Profile ARN string for use later in this process. The following image shows the **Summary** window with redacted ARN strings:



## Step 2. Configure a Policy for the Target S3 Bucket

Configure the target S3 bucket with a policy that includes account, role, and the bucket name.

You can create a new policy, or add to an existing policy.

1. Copy the following statement:

```
{
  "Sid": "AWS Databricks Policy for s3 bucket - put, get, delete",
  "Effect": "Allow",
  "Principal": {
    "AWS": "arn:aws:iam::<aws-account-id-databricks>:role/<iam-role-for-s3-access>"
  },
  "Action": [
    "s3:PutObject",
    "s3:GetObject",
    "s3:DeleteObject"
  ],
  "Resource": "arn:aws:s3::<s3-bucket-name>/*"
},
{
  "Sid": "AWS Databricks Policy for s3 bucket - list, get bucket location",
  "Effect": "Allow",
  "Principal": {
    "AWS": "arn:aws:iam::<aws-account-id-databricks>:role/<iam-role-for-s3-access>"
  },
  "Action": [
    "s3:ListBucket",
    "s3:GetBucketLocation"
  ],
  "Resource": "arn:aws:s3::<s3-bucket-name>"
}
```

2. Configure the following elements in the statement:

Element	Description
aws-account-id-databricks	Account ID of the AWS account in which you are configuring this integration.
iam-role-for-s3-access	The IAM role that you created in Create an IAM Role and Policy to Access an S3 Bucket.
s3-bucket-name	S3 bucket name.

3. Click **Save**.

## Step 3. Add IAM Roles to the EC2 Policy and Databricks

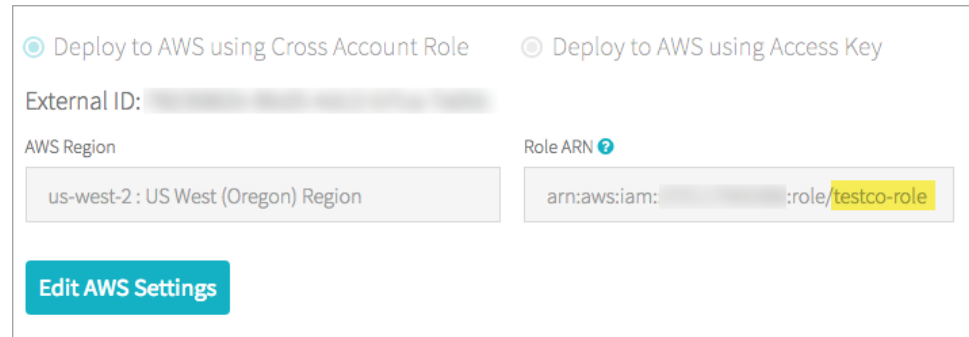
Add IAM roles to EC2 policy and the Databricks account.

**Note:** The IAM role that you add in steps 1 and 2 below is different from the S3 IAM that you created in "Step 1. Create an IAM Role and Policy to Access an S3 Bucket."

1. Add the IAM role that corresponds to the Databricks account to the EC2 instance policy.
  - a. Log into the Databricks account and select the **AWS Account** tab.

- b. Copy the role name at the end of the Role ARN.

The following image shows an example of the Role ARN pane. In this example, the role name to copy is `testco-role`:



2. Add the Role ARN to the EC2 policy.

Modify the EC2 instance policy to allow Databricks to pass the IAM role you copied in step 1 above to the EC2 instances for the Databricks Spark clusters:

- a. In the AWS console, go to the IAM service and select the **Policies** tab.
  - b. Search for `AwsDatabricksUserCreationPolicy`.
  - c. Click **Edit Policy**.

The policy opens in a JSON statement editing pane.

- d. Find the section of the statement that contains `"Action": "iam:PassRole"`.

The following text provides an example of the section to search for:

```
{
  "Effect": "Allow",
  "Action": "iam:PassRole",
  "Resource": "arn:aws:iam::<aws-account-id-databricks>:role/<iam-role-for-s3-access>"
},
```

- e. Paste the IAM role in the place of the string `<iam-role-for-s3-access>` shown in the example.
    - f. Click **Review Policy**, then **Save changes**.

3. Add the Instance Profile ARN to the Databricks account.

- a. In the Databricks Admin Console, click the **IAM Roles** tab.
  - b. Click **Add IAM Role**.
  - c. Paste the Instance Profile ARN string that you created in "Step 1. Create an IAM Role and Policy to Access an S3 Bucket."

The following image shows where to paste the Instance Profile ARN:

Settings / IAM roles / Add IAM Role

< Add IAM Role

Cancel

Add

IAM roles allow you to access your data from Databricks clusters without the need to manage, deploy, or rotate AWS keys.

In order to use an IAM role in Databricks, the access policy used by Databricks to launch clusters must be given the "PassRole" permission for that role.

[Learn more](#)

Instance Profile ARN ⓘ

arn:aws:iam::[REDACTED]:instance-profile/testco-ec2-role

Paste here

IAM Role Name

testco-ec2-role

d. Click **Add**.

## Step 4. Launch a Databricks Cluster with the S3 IAM Role

When you create or launch a cluster, use the S3 IAM role that is contained in the Instance Profile ARN string.

To use the IAM role in the operation of the cluster, select it in the cluster Advanced Options:

1. Select or create a Databricks cluster, and open the **Advanced Options** section.
2. Click the **Instances** tab.
3. Click the dropdown list under **IAM Role**.

The following image shows an example of the **IAM Role** dropdown on the Instances tab:

The screenshot shows the AWS IAM console configuration for an EC2 instance profile. The 'Instances' tab is active. The 'Availability Zone' is set to 'us-west-2c'. The 'Spot Bid Price' is set to '100'. The 'EBS Volume Type' is set to 'None', '# Volumes' is '0', and 'Size in GB' is empty. The 'IAM Role' dropdown is open, showing 'None' and 'testco-ec2-role' (selected).

4. Select the role to use to operate the cluster.

## Download and Install the JDBC Driver to Enable Delta Lake Access

If you want to use Delta Lake as a source or target, perform the following steps to download the Databricks JDBC driver to enable the connection between the domain and the Delta Lake resource.

1. In a browser, go to the [Databricks JDBC/ODBC driver page](#) at online Databricks documentation. Fill out the form and submit it.  
After you submit the form, you receive an email message that includes multiple download options.
2. Save the JDBC driver.jar file to the following directory on each machine that hosts the Data Integration Service, and then restart the Informatica domain:  
<Informatica home>/externaljdbcjars
3. Save the JDBC driver .jar file to the following location on each machine that hosts the Developer tool client, and then restart the Developer tool:  
<Developer tool home>/clients/externaljdbcjars

## Configure ADLS Storage Access

If you use ADLS storage, configure some Hadoop credential configuration options as Databricks Spark options.

Configure the following properties in the advanced Spark properties on the Databricks cluster to Delta Lake tables on ADLS Gen1:

Property	Description
<code>dfs.adls.oauth2.access.token.provider.type</code>	Access token provider type. Use the value <code>ClientCredential</code> .
<code>dfs.adls.oauth2.client.id</code>	The application ID associated with the Service Principal required to authorize the service principal and access the storage.
<code>dfs.adls.oauth2.credential</code>	The password required to authorize the service principal and access the storage.
<code>dfs.adls.oauth2.refresh.url</code>	The OAUTH 2.0 token endpoint required to authorize the service principal and access the storage.

See the [Databricks documentation](#) for examples on how to specify these settings in Scala or Python.

The following excerpt shows how these properties appear in the Spark properties window:

```
spark.hadoop.dfs.adls.oauth2.access.token.provider.type ClientCredential
spark.hadoop.dfs.adls.oauth2.client.id <your-service-client-id>
spark.hadoop.dfs.adls.oauth2.credential <your-service-credentials>
spark.hadoop.dfs.adls.oauth2.refresh.url "https://login.microsoftonline.com/<your-
directory-id>/oauth2/token"
```

## Configure WASB Storage Access

If you use WASB storage, you must set an account access key as the Hadoop configuration key.

To use the Azure storage account access key, add "spark.hadoop" as a prefix to the Hadoop configuration key as shown in the following text:

```
spark.hadoop.fs.azure.account.key.<your-storage-account-name>.blob.core.windows.net
<your-storage-account-access-key>
```

## Create a Staging Directory for Binary Archive Files

Optionally, create a directory on DBFS that the Data Integration Service uses to stage the Informatica binary archive files.

By default, the Data Integration Service writes binary archive files to the DBFS directory `/tmp`.

Optionally, you can create a different staging directory for these files. If you create a staging directory, configure the path in the **Cluster Staging Directory** property of the Data Integration Service.

To create a directory on DBFS, see the [Databricks documentation](#).



# Create a Staging Directory for Run-time Processing

When the Databricks Spark engine runs a job, it stores temporary files in a staging directory.

Optionally, you can create a directory on DBFS to stage temporary files during run time. By default, the Data Integration Service uses the DBFS directory at the following path: `/<cluster staging directory>/DATABRICKS`

To create a directory on DBFS, see the [Databricks documentation](#).

## Prepare for Token Authentication

The Data Integration Service uses token-based authentication to provide access to the Databricks environment.

Create a Databricks user to generate the authentication token. Complete the following tasks to prepare for authentication.

1. In the Access Control tab of the Databricks Admin Console, verify that personal access tokens are enabled.
2. Verify that the Databricks environment contains a user to generate the token.
3. Grant permissions to the token user.
  - If you created staging directories, grant permission to access and write to the directories.
  - If you did not create staging directories, grant permission to create directories.
4. Log in to the Databricks Admin Console as the token user.
5. Go to **Access Control > User Settings** and click the **Generate New Token** button to generate a token for the user.

**Important:** The console displays the token only once, at generation. Copy and save the token value for later use.

## Configure the Data Integration Service

Configure the Data Integration Service to integrate with the Databricks environment.

Perform the following pre-integration tasks:

1. If you created a staging directory, configure the path in the Data Integration Service properties.
2. Prepare an installation of Python on the Data Integration Service if you plan to run the Python transformation on a Hadoop cluster.

**Note:** When you plan to run a mapping with a Python transformation on a Databricks cluster, it is not necessary to install Python on the Data Integration Service machine.

### Configure Data Integration Service Properties

The Data Integration Service contains properties that integrate the domain with the Databricks cluster.

Configure the following property in the Data Integration Service:

### Cluster Staging Directory

The directory on the cluster where the Data Integration Service pushes the binaries to integrate the native and non-native environments and to store temporary files during processing. Default is `/tmp`.

## Install Python Libraries

Databricks comes installed with some Python libraries. If you need to install additional third-party Python libraries, use the pip installer for Databricks. The Databricks Spark engine supports Python version 3.

The Databricks cluster provides a preloaded set of Python libraries. In some cases, your Databricks administrator might determine that the workspace requires additional libraries or modules. When additional libraries or modules are required, they must be installed through an init script during cluster creation. See the [Databricks documentation](#).

Perform the following tasks to install third-party Python libraries:

1. Write an init script that includes the Python libraries to install.
2. Upload the script to the DBFS directory. If you use AWS Databricks, you can upload the script to the S3 directory instead.

**Note:** When you create an ephemeral cluster using a cluster workflow, include the init script file location in the advanced properties for the Create Cluster task.

For more information about the installed Python libraries that come with Databricks, refer to the [Databricks documentation](#).

## CHAPTER 12

# Databricks Integration Tasks

This chapter includes the following topics:

- [Create a Databricks Cluster Configuration, 203](#)
- [Configure the Databricks Connection, 206](#)
- [Complete Upgrade Tasks, 206](#)

## Create a Databricks Cluster Configuration

A Databricks cluster configuration is an object in the domain that contains configuration information about the Databricks cluster. The cluster configuration enables the Data Integration Service to push mapping logic to the Databricks environment.

Use the Administrator tool to import configuration properties from the Databricks cluster to create a cluster configuration. You can import configuration properties from the cluster or from a file that contains cluster properties. You can choose to create a Databricks connection when you perform the import.

**Note:**

- Ensure that you integrate a Databricks cluster with only one Informatica domain.
- Ensure that a Databricks workspace has access to the database VM where you create the reference table.

## Importing a Databricks Cluster Configuration from the Cluster

When you import the cluster configuration directly from the cluster, you provide information to connect to the cluster.

Before you import the cluster configuration, get cluster information from the Databricks administrator.

1. From the **Connections** tab, click the **ClusterConfigurations** node in the Domain Navigator.
2. From the Actions menu, select **New > Cluster Configuration**.  
The **Cluster Configuration** wizard opens.

3. Configure the following properties:

Property	Description
Cluster configuration name	Name of the cluster configuration.
Description	Optional description of the cluster configuration.
Distribution type	The distribution type. Choose <b>Databricks</b> .
Method to import the cluster configuration	Choose <b>Import from cluster</b> .
Databricks domain	Domain name of the Databricks deployment.
Databricks access token	The token ID created within Databricks, required for authentication. . <b>Note:</b> If the token has an expiration date, verify that you get a new token from the Databricks administrator before it expires.
Databricks cluster ID	The cluster ID of the Databricks cluster. To find the cluster ID on the Databricks portal, follow these steps: 1. Select <b>Clusters</b> from the object bar on the left side. 2. Select the cluster you want to integrate with the Informatica domain. 3. Click the <b>Spark ID</b> tab and expand the list of <b>Spark Properties</b> . 4. Select the <b>Tags</b> tab.
Create connection	Choose to create a Databricks connection. If you choose to create a connection, the <b>Cluster Configuration</b> wizard associates the cluster configuration with the Databricks connection. If you do not choose to create a connection, you must manually create one and associate the cluster configuration with it.

4. Click **Next** to verify the information on the summary page.

## Importing a Databricks Cluster Configuration from a File

You can import properties from an archive file to create a cluster configuration.

Complete the following tasks to import a Databricks cluster from a file:

1. Get required cluster properties from the Databricks administrator.
2. Create an .xml file with the cluster properties, and compress it into a .zip or .tar file.
3. Log in to the Administrator tool and import the file.

### Create the Import File

To import the cluster configuration from a file, you must create an archive file.

To create the .xml file for import, you must get required information from the Databricks administrator. You can provide any name for the file and store it locally.

The following table describes the properties required to import the cluster information:

Property Name	Description
cluster_name	Name of the Databricks cluster.
cluster_ID	The cluster ID of the Databricks cluster.
baseURL	URL to access the Databricks cluster. This is the domain URL that appears in your browser menu bar. It commonly incorporates your account region. For example, <a href="https://southcentralus.azuredatabricks.net">https://southcentralus.azuredatabricks.net</a> or <a href="https://westus.azuredatabricks.net">https://westus.azuredatabricks.net</a> .
accesstoken	The token ID created within Databricks required for authentication.

Optionally, you can include other properties specific to the Databricks environment.

When you complete the .xml file, compress it into a .zip or .tar file for import.

### Sample Import File

The following text shows a sample import file with the required properties:

```
<?xml version="1.0" encoding="UTF-8"?><configuration>
  <property>
    <name>cluster_name</name>
    <value>my_cluster</value>
  </property>
  <property>
    <name>cluster_id</name>
    <value>0926-294544-bckt123</value>
  </property>
  <property>
    <name>baseURL</name>
    <value>https://<region>.azuredatabricks.net</value>
  </property>
  <property>
    <name>accesstoken</name>
    <value>dapicf76c2d4567c6sldn654fe875936e778</value>
  </property>
</configuration>
```

### Import the Cluster Configuration

After you create the .xml file with the cluster properties, use the Administrator tool to import into the domain and create the cluster configuration.

1. From the **Connections** tab, click the **ClusterConfigurations** node in the Domain Navigator.
2. From the Actions menu, select **New > Cluster Configuration**.  
The **Cluster Configuration** wizard opens.
3. Configure the following properties:

Property	Description
Cluster configuration name	Name of the cluster configuration.
Description	Optional description of the cluster configuration.

Property	Description
Distribution type	The distribution type. Choose <b>Databricks</b> .
Method to import the cluster configuration	Choose <b>Import from file</b> .
Upload configuration archive file	The full path and file name of the file. Click the Browse button to navigate to the file.
Create connection	Choose to create a Databricks connection. If you choose to create a connection, the <b>Cluster Configuration</b> wizard associates the cluster configuration with the Databricks connection. If you do not choose to create a connection, you must manually create one and associate the cluster configuration with it.

- Click **Next** to verify the information on the summary page.

## Configure the Databricks Connection

Databricks connections contain information to connect to the Databricks cluster. If you did not choose to create a Databricks connection when you imported the cluster configuration, you must manually create one.

For information about the Databricks connection properties, see [“Databricks Connection Properties” on page 227](#).

## Complete Upgrade Tasks

If you upgraded the Informatica platform, you need to perform some additional tasks within the Informatica domain.

Based on the version that you upgraded from, you might need to perform the following updates:

### Update date/time data format for sources and targets

#### Perform this task in the following situation:

- You upgraded from version 10.2.2 or earlier.

Effective in version 10.4.0, the Databricks Spark engine reads and writes date/time values in the format YYYY-MM-DD HH24:MM:SS.US.

To continue to use the Databricks Spark engine to read date/time data in other formats, edit mappings to include an Expression transformation. Read the date/time field from the source as a string. In the expression, convert the string to a date/time format.

To write data to a different format, use an Expression transformation to convert the date/time data to a different format. Then write the data to the target as a string.

# APPENDIX A

## Connections Reference

This appendix includes the following topics:

- [Connections Overview, 208](#)
- [Cloud Provisioning Configuration, 208](#)
- [Amazon Redshift Connection Properties, 218](#)
- [Amazon S3 Connection Properties, 220](#)
- [Blockchain Connection Properties, 222](#)
- [Cassandra Connection Properties, 224](#)
- [Confluent Kafka Connection, 225](#)
- [Databricks Connection Properties, 227](#)
- [Google Analytics Connection Properties, 229](#)
- [Google BigQuery Connection Properties, 230](#)
- [Google Cloud Spanner Connection Properties, 232](#)
- [Google Cloud Storage Connection Properties, 233](#)
- [Google PubSub Connection Properties, 234](#)
- [Hadoop Connection Properties, 234](#)
- [HDFS Connection Properties, 241](#)
- [HBase Connection Properties, 243](#)
- [HBase Connection Properties for MapR-DB, 243](#)
- [Hive Connection Properties, 244](#)
- [JDBC Connection Properties, 247](#)
- [JDBC V2 Connection Properties, 254](#)
- [Kafka Connection Properties, 256](#)
- [Kudu Connection Properties , 259](#)
- [Microsoft Azure Blob Storage Connection Properties, 260](#)
- [Microsoft Azure Cosmos DB SQL API Connection Properties, 261](#)
- [Microsoft Azure Data Lake Storage Gen1 Connection Properties, 262](#)
- [Microsoft Azure Data Lake Storage Gen2 Connection Properties, 263](#)
- [Microsoft Azure SQL Data Warehouse Connection Properties, 264](#)
- [Snowflake Connection Properties, 266](#)
- [Creating a Connection to Access Sources or Targets, 267](#)

- [Creating a Hadoop Connection, 267](#)
- [Configuring Hadoop Connection Properties, 269](#)

## Connections Overview

Create a connection to access non-native environments, Hadoop and Databricks. If you access HBase, HDFS, or Hive sources or targets in the Hadoop environment, you must also create those connections. You can create the connections using the Developer tool, Administrator tool, and infacmd.

You can create the following types of connections:

### **Hadoop connection**

Create a Hadoop connection to run mappings in the Hadoop environment.

### **HBase connection**

Create an HBase connection to access HBase. The HBase connection is a NoSQL connection.

### **HDFS connection**

Create an HDFS connection to read data from or write data to the HDFS file system on a Hadoop cluster.

### **Hive connection**

Create a Hive connection to access Hive as a source or target. You can access Hive as a source if the mapping is enabled for the native or Hadoop environment. You can access Hive as a target if the mapping runs on the Blaze engine.

### **JDBC connection**

Create a JDBC connection and configure Sqoop properties in the connection to import and export relational data through Sqoop.

### **Databricks connection**

Create a Databricks connection to run mappings in the Databricks environment.

**Note:** For information about creating connections to other sources or targets such as social media web sites or Teradata, see the respective PowerExchange adapter user guide for information.

## Cloud Provisioning Configuration

The cloud provisioning configuration establishes a relationship between the Create Cluster task and the cluster connection that the workflows use to run mapping tasks. The Create Cluster task must include a reference to the cloud provisioning configuration. In turn, the cloud provisioning configuration points to the cluster connection that you create for use by the cluster workflow.

The properties to populate depend on the Hadoop distribution you choose to build a cluster on. Choose one of the following connection types:

- AWS Cloud Provisioning. Connects to an Amazon EMR cluster on Amazon Web Services.
- Azure Cloud Provisioning. Connects to an HDInsight cluster on the Azure platform.
- Databricks Cloud Provisioning. Connects to a Databricks cluster on the Azure Databricks platform.



# AWS Cloud Provisioning Configuration Properties

The properties in the AWS cloud provisioning configuration enable the Data Integration Service to contact and create resources on the AWS cloud platform.

## General Properties

The following table describes cloud provisioning configuration general properties:

Property	Description
Name	Name of the cloud provisioning configuration.
ID	ID of the cloud provisioning configuration. Default: Same as the cloud provisioning configuration name.
Description.	Optional. Description of the cloud provisioning configuration.
AWS Access Key ID	Optional. ID of the AWS access key, which AWS uses to control REST or HTTP query protocol requests to AWS service APIs. If you do not specify a value, Informatica attempts to follow the Default Credential Provider Chain.
AWS Secret Access Key	Secret component of the AWS access key. Required if you specify the AWS Access Key ID.
Region	Region in which to create the cluster. This must be the region in which the VPC is running. Use AWS region values. For a list of acceptable values, see AWS documentation. <b>Note:</b> The region where you want to create the cluster can be different from the region in which the Informatica domain is installed.

## Permissions

The following table describes cloud provisioning configuration permissions properties:

Property	Description
EMR Role	Name of the service role for the EMR cluster that you create. The role must have sufficient permissions to create a cluster, access S3 resources, and run jobs on the cluster. When the AWS administrator creates this role, they select the "EMR" role. This contains the default AmazonElasticMapReduceRole policy. You can edit the services in this policy.
EC2 Instance Profile	Name of the EC2 instance profile role that controls permissions on processes that run on the cluster. When the AWS administrator creates this role, they select the "EMR Role for EC2" role. This includes S3 access by default.
Auto Scaling Role	Required if you configure auto-scaling for the EMR cluster. This role is created when the AWS administrator configures auto-scaling on any cluster in the VPC. Default: When you leave this field blank, it is equivalent to setting the Auto Scaling role to "Proceed without role" when the AWS administrator creates a cluster in the AWS console.

## EC2 Configuration

The following table describes cloud provisioning configuration EC2 configuration properties:

Property	Description
EC2 Key Pair	EC2 key pair to enable communication with the EMR cluster master node. Optional. This credential enables you to log into the cluster. Configure this property if you intend the cluster to be non-ephemeral.
EC2 Subnet	ID of the subnet on the VPC in which to create the cluster. Use the subnet ID of the EC2 instance where the cluster runs.
Master Security Group	Optional. ID of the security group for the cluster master node. Acts as a virtual firewall to control inbound and outbound traffic to cluster nodes. Security groups are created when the AWS administrator creates and configures a cluster in a VPC. In the AWS console, the property is equivalent to ElasticMapReduce-master. You can use existing security groups, or the AWS administrator might create dedicated security groups for the ephemeral cluster. If you do not specify a value, the cluster applies the default security group for the VPC.
Additional Master Security Groups	Optional. IDs of additional security groups to attach to the cluster master node. Use a comma-separated list of security group IDs.
Core and Task Security Group	Optional. ID of the security group for the cluster core and task nodes. When the AWS administrator creates and configures a cluster In the AWS console, the property is equivalent to the ElasticMapReduce-slave security group If you do not specify a value, the cluster applies the default security group for the VPC.
Additional Core and Task Security Groups	Optional. IDs of additional security groups to attach to cluster core and task nodes. Use a comma-separated list of security group IDs.
Service Access Security Group	EMR managed security group for service access. Required when you provision an EMR cluster in a private subnet.

## General Properties

The following table describes cloud provisioning configuration general properties:

Property	Description
Name	Name of the cloud provisioning configuration.
ID	ID of the cloud provisioning configuration. Default: Same as the cloud provisioning configuration name.
Description.	Optional. Description of the cloud provisioning configuration.
AWS Access Key ID	Optional. ID of the AWS access key, which AWS uses to control REST or HTTP query protocol requests to AWS service APIs. If you do not specify a value, Informatica attempts to follow the Default Credential Provider Chain.

Property	Description
AWS Secret Access Key	Secret component of the AWS access key. Required if you specify the AWS Access Key ID.
Region	Region in which to create the cluster. This must be the region in which the VPC is running. Use AWS region values. For a list of acceptable values, see AWS documentation. <b>Note:</b> The region where you want to create the cluster can be different from the region in which the Informatica domain is installed.

## Permissions

The following table describes cloud provisioning configuration permissions properties:

Property	Description
EMR Role	Name of the service role for the EMR cluster that you create. The role must have sufficient permissions to create a cluster, access S3 resources, and run jobs on the cluster. When the AWS administrator creates this role, they select the "EMR" role. This contains the default AmazonElasticMapReduceRole policy. You can edit the services in this policy.
EC2 Instance Profile	Name of the EC2 instance profile role that controls permissions on processes that run on the cluster. When the AWS administrator creates this role, they select the "EMR Role for EC2" role. This includes S3 access by default.
Auto Scaling Role	Required if you configure auto-scaling for the EMR cluster. This role is created when the AWS administrator configures auto-scaling on any cluster in the VPC. Default: When you leave this field blank, it is equivalent to setting the Auto Scaling role to "Proceed without role" when the AWS administrator creates a cluster in the AWS console.

## EC2 Configuration

The following table describes cloud provisioning configuration EC2 configuration properties:

Property	Description
EC2 Key Pair	EC2 key pair to enable communication with the EMR cluster master node. Optional. This credential enables you to log into the cluster. Configure this property if you intend the cluster to be non-ephemeral.
EC2 Subnet	ID of the subnet on the VPC in which to create the cluster. Use the subnet ID of the EC2 instance where the cluster runs.
Master Security Group	Optional. ID of the security group for the cluster master node. Acts as a virtual firewall to control inbound and outbound traffic to cluster nodes. Security groups are created when the AWS administrator creates and configures a cluster in a VPC. In the AWS console, the property is equivalent to ElasticMapReduce-master. You can use existing security groups, or the AWS administrator might create dedicated security groups for the ephemeral cluster. If you do not specify a value, the cluster applies the default security group for the VPC.

Property	Description
Additional Master Security Groups	Optional. IDs of additional security groups to attach to the cluster master node. Use a comma-separated list of security group IDs.
Core and Task Security Group	Optional. ID of the security group for the cluster core and task nodes. When the AWS administrator creates and configures a cluster in the AWS console, the property is equivalent to the ElasticMapReduce-slave security group. If you do not specify a value, the cluster applies the default security group for the VPC.
Additional Core and Task Security Groups	Optional. IDs of additional security groups to attach to cluster core and task nodes. Use a comma-separated list of security group IDs.
Service Access Security Group	EMR managed security group for service access. Required when you provision an EMR cluster in a private subnet.

## Azure Cloud Provisioning Configuration Properties

The properties in the Azure cloud provisioning configuration enable the Data Integration Service to contact and create resources on the Azure cloud platform.

### Authentication Details

The following table describes authentication properties to configure:

Property	Description
Name	Name of the cloud provisioning configuration.
ID	ID of the cloud provisioning configuration. Default: Same as the cloud provisioning configuration name.
Description	Optional. Description of the cloud provisioning configuration.
Subscription ID	ID of the Azure account to use in the cluster creation process.
Tenant ID	A GUID string associated with the Azure Active Directory.
Client ID	A GUID string that is the same as the Application ID associated with the Service Principal. The Service Principal must be assigned to a role that has permission to create resources in the subscription that you identified in the Subscription ID property.
Client Secret	An octet string that provides a key associated with the client ID.

### Storage Account Details

Choose to configure access to one of the following storage types:

- Azure Data Lake Storage (ADLS). See [Azure documentation](#).
- An Azure Storage Account, known as general or blob storage. See [Azure documentation](#).

The following table describes the information you need to configure Azure Data Lake Storage (ADLS) with the HDInsight cluster:

Property	Description
Azure Data Lake Store Name	Name of the ADLS storage to access. The ADLS storage and the cluster to create must reside in the same region.
Data Lake Service Principal Client ID	A credential that enables programmatic access to ADLS storage. Enables the Informatica domain to communicate with ADLS and run commands and mappings on the HDInsight cluster.  The service principal is an Azure user that meets the following requirements: <ul style="list-style-type: none"> <li>- Permissions to access required directories in ADLS storage.</li> <li>- Certificate-based authentication for ADLS storage.</li> <li>- Key-based authentication for ADLS storage.</li> </ul>
Data Lake Service Principal Certificate Contents	The Base64 encoded text of the public certificate used with the service principal.  Leave this property blank when you create the cloud provisioning configuration. After you save the cloud provisioning configuration, log in to the VM where the Informatica domain is installed and run <code>infacmd ccps updateADLSCertificate</code> to populate this property.
Data Lake Service Principal Certificate Password	Private key for the service principal. This private key must be associated with the service principal certificate.
Data Lake Service Principal Client Secret	An octet string that provides a key associated with the service principal.
Data Lake Service Principal OAUTH Token Endpoint	Endpoint for OAUTH token based authentication.

The following table describes the information you need to configure Azure General Storage, also known as blob storage, with the HDInsight cluster:

Property	Description
Azure Storage Account Name	Name of the storage account to access. Get the value from the Storage Accounts node in the Azure web console. The storage and the cluster to create must reside in the same region.
Azure Storage Account Key	A key to authenticate access to the storage account. To get the value from the Azure web console, select the storage account, then Access Keys. The console displays the account keys.

## Cluster Deployment Details

The following table describes the cluster deployment properties that you configure:

Property	Description
Resource Group	Resource group in which to create the cluster. A resource group is a logical set of Azure resources.
Virtual Network Resource Group	Optional. Resource group to which the virtual network belongs. If you do not specify a resource group, the Data Integration Service assumes that the virtual network is a member of the same resource group as the cluster.
Virtual Network	Name of the virtual network or vnet where you want to create the cluster. Specify a vnet that resides in the resource group that you specified in the Virtual Network Resource Group property. The vnet must be in the same region as the region in which to create the cluster.
Subnet Name	Subnet in which to create the cluster. The subnet must be a part of the vnet that you designated in the previous property. Each vnet can have one or more subnets. The Azure administrator can choose an existing subnet or create one for the cluster.

## External Hive Metastore Details

You can specify the properties to enable the cluster to connect to a Hive metastore database that is external to the cluster.

If you do not specify an existing external database in this dialog box, the cluster creates its own database on the cluster. This database is terminated when the cluster is terminated.

You can use an external relational database like MySQL or Amazon RDS as the Hive metastore database. The external database must be on the same cloud platform as the cluster to create.

The following table describes the Hive metastore database properties that you configure:

Property	Description
Database Name	Name of the Hive metastore database.
Database Server Name	Server on which the database resides. <b>Note:</b> The database server name on the Azure web console commonly includes the suffix <code>database.windows.net</code> . For example: <code>server123xyz.database.windows.net</code> . You can specify the database server name without the suffix and Informatica will automatically append the suffix. For example, you can specify <code>server123xyz</code> .
Database User Name	User name of the account for the domain to use to access the database.
Database Password	Password for the user account.

## Authentication Details

The following table describes authentication properties to configure:

Property	Description
Name	Name of the cloud provisioning configuration.
ID	ID of the cloud provisioning configuration. Default: Same as the cloud provisioning configuration name.
Description	Optional. Description of the cloud provisioning configuration.
Subscription ID	ID of the Azure account to use in the cluster creation process.
Tenant ID	A GUID string associated with the Azure Active Directory.
Client ID	A GUID string that is the same as the Application ID associated with the Service Principal. The Service Principal must be assigned to a role that has permission to create resources in the subscription that you identified in the Subscription ID property.
Client Secret	An octet string that provides a key associated with the client ID.

## Storage Account Details

Choose to configure access to one of the following storage types:

- Azure Data Lake Storage (ADLS). See [Azure documentation](#).
- An Azure Storage Account, known as general or blob storage. See [Azure documentation](#).

The following table describes the information you need to configure Azure Data Lake Storage (ADLS) with the HDInsight cluster:

Property	Description
Azure Data Lake Store Name	Name of the ADLS storage to access. The ADLS storage and the cluster to create must reside in the same region.
Data Lake Service Principal Client ID	A credential that enables programmatic access to ADLS storage. Enables the Informatica domain to communicate with ADLS and run commands and mappings on the HDInsight cluster.  The service principal is an Azure user that meets the following requirements: <ul style="list-style-type: none"><li>- Permissions to access required directories in ADLS storage.</li><li>- Certificate-based authentication for ADLS storage.</li><li>- Key-based authentication for ADLS storage.</li></ul>
Data Lake Service Principal Certificate Contents	The Base64 encoded text of the public certificate used with the service principal.  Leave this property blank when you create the cloud provisioning configuration. After you save the cloud provisioning configuration, log in to the VM where the Informatica domain is installed and run <code>infacmd ccps updateADLSCertificate</code> to populate this property.
Data Lake Service Principal Certificate Password	Private key for the service principal. This private key must be associated with the service principal certificate.

Property	Description
Data Lake Service Principal Client Secret	An octet string that provides a key associated with the service principal.
Data Lake Service Principal OAUTH Token Endpoint	Endpoint for OAUTH token based authentication.

The following table describes the information you need to configure Azure General Storage, also known as blob storage, with the HDInsight cluster:

Property	Description
Azure Storage Account Name	Name of the storage account to access. Get the value from the Storage Accounts node in the Azure web console. The storage and the cluster to create must reside in the same region.
Azure Storage Account Key	A key to authenticate access to the storage account. To get the value from the Azure web console, select the storage account, then Access Keys. The console displays the account keys.

## Cluster Deployment Details

The following table describes the cluster deployment properties that you configure:

Property	Description
Resource Group	Resource group in which to create the cluster. A resource group is a logical set of Azure resources.
Virtual Network Resource Group	Optional. Resource group to which the virtual network belongs. If you do not specify a resource group, the Data Integration Service assumes that the virtual network is a member of the same resource group as the cluster.
Virtual Network	Name of the virtual network or vnet where you want to create the cluster. Specify a vnet that resides in the resource group that you specified in the Virtual Network Resource Group property. The vnet must be in the same region as the region in which to create the cluster.
Subnet Name	Subnet in which to create the cluster. The subnet must be a part of the vnet that you designated in the previous property. Each vnet can have one or more subnets. The Azure administrator can choose an existing subnet or create one for the cluster.

## External Hive Metastore Details

You can specify the properties to enable the cluster to connect to a Hive metastore database that is external to the cluster.

You can use an external relational database like MySQL or Amazon RDS as the Hive metastore database. The external database must be on the same cloud platform as the cluster to create.

If you do not specify an existing external database in this dialog box, the cluster creates its own database on the cluster. This database is terminated when the cluster is terminated.



The following table describes the Hive metastore database properties that you configure:

Property	Description
Database Name	Name of the Hive metastore database.
Database Server Name	Server on which the database resides. <b>Note:</b> The database server name on the Azure web console commonly includes the suffix <code>database.windows.net</code> . For example: <code>server123xyz.database.windows.net</code> . You can specify the database server name without the suffix and Informatica will automatically append the suffix. For example, you can specify <code>server123xyz</code> .
Database User Name	User name of the account for the domain to use to access the database.
Database Password	Password for the user account.

## Databricks Cloud Provisioning Configuration Properties

The properties in the Databricks cloud provisioning configuration enable the Data Integration Service to contact and create resources on the Databricks cloud platform.

The following table describes the Databricks cloud provisioning configuration properties:

Property	Description
Name	Name of the cloud provisioning configuration. <b>Tip:</b> Because the Administrator tool lists cloud provisioning configuration objects with other connections, use a naming convention such as "CPC" as part of the name of the object to help identify it.
ID	The cluster ID of the Databricks cluster.
Description	Optional description of the cloud provisioning configuration.
Databricks domain	Domain name of the Databricks deployment.
Databricks token ID	The token ID created within Databricks required for authentication. <b>Note:</b> If the token has an expiration date, verify that you get a new token from the Databricks administrator before it expires.
Advanced Properties	Advanced properties that are unique to the Databricks cloud provisioning configuration.

### Advanced Properties

Configure the following properties in the **Advanced Properties** of the Databricks configuration section:

#### **infaspark.pythontx.exec**

Required to run a Python transformation on the Databricks Spark engine. Set to the location of the Python executable binary on the worker nodes in the Databricks cluster.

When you provision the cluster at run time, set this property in the Databricks cloud provisioning configuration. Otherwise, set on the Databricks connection.

For example, set to:

```
infaspark.pythontx.exec=/databricks/python3/bin/python3
```

#### **infaspark.pythontx.executorEnv.PYTHONHOME**

Required to run a Python transformation on the Databricks Spark engine. Set to the location of the Python installation directory on the worker nodes in the Databricks cluster.

When you provision the cluster at run time, set this property in the Databricks cloud provisioning configuration. Otherwise, set on the Databricks connection.

For example, set to:

```
infaspark.pythontx.executorEnv.PYTHONHOME=/databricks/python3
```

## Amazon Redshift Connection Properties

When you set up an Amazon Redshift connection, you must configure the connection properties.

The following table describes the Amazon Redshift connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters:~`!\$%^&*()-+= {[}]  \ : ; " ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection.
Type	The connection type. Select Amazon Redshift in the Database.

The **Details** tab contains the connection attributes of the Amazon Redshift connection. The following table describes the connection attributes:

Property	Description
Username	User name of the Amazon Redshift account.
Password	Password for the Amazon Redshift account.
Access Key ID	Amazon S3 bucket access key ID. <b>Note:</b> Required if you do not use AWS Identity and Access Management (IAM) authentication.
Secret Access Key	Amazon S3 bucket secret access key ID. <b>Note:</b> Required if you do not use AWS Identity and Access Management (IAM) authentication.

Property	Description
Master Symmetric Key	Optional. Provide a 256-bit AES encryption key in the Base64 format when you enable client-side encryption. You can generate a key using a third-party tool. If you specify a value, ensure that you specify the encryption type as client side encryption in the advanced target properties.
JDBC URL	Amazon Redshift connection URL.
Cluster Region	Optional. 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 <b>JDBC URL</b> connection property. If you specify a cluster region in both <b>Cluster Region</b> and <b>JDBC URL</b> connection properties, the Data Integration Service ignores the cluster region that you specify in the <b>JDBC URL</b> connection property. To use the cluster region name that you specify in the <b>JDBC URL</b> connection property, select <b>None</b> as the cluster region in this property. Select one of the following cluster regions: Select one of the following regions: <ul style="list-style-type: none"> <li>- Asia Pacific (Mumbai)</li> <li>- Asia Pacific (Seoul)</li> <li>- Asia Pacific (Singapore)</li> <li>- Asia Pacific (Sydney)</li> <li>- Asia Pacific (Tokyo)</li> <li>- AWS GovCloud (US)</li> <li>- Canada (Central)</li> <li>- China (Beijing)</li> <li>- China (Ningxia)</li> <li>- EU (Ireland)</li> <li>- EU (Frankfurt)</li> <li>- EU (London)</li> <li>- EU (Paris)</li> <li>- South America (Sao Paulo)</li> <li>- US East (Ohio)</li> <li>- US East (N. Virginia)</li> <li>- US West (N. California)</li> <li>- US West (Oregon)</li> </ul> Default is None. You can only read data from or write data to the cluster regions supported by AWS SDK used by PowerExchange for Amazon Redshift.
Customer Master Key ID	Optional. Specify the customer master key ID generated by AWS Key Management Service (AWS KMS) or the Amazon Resource Name (ARN) of your custom key for cross-account access. You must generate the customer master key corresponding to the region where Amazon S3 bucket resides. You can specify any of the following values: <b>Customer generated customer master key</b>  Enables client-side or server-side encryption.  <b>Default customer master key</b>  Enables client-side or server-side encryption. Only the administrator user of the account can use the default customer master key ID to enable client-side encryption.

# Amazon S3 Connection Properties

When you set up an Amazon S3 connection, you must configure the connection properties.

The following table describes the Amazon S3 connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters:~`!\$%^&*()-+={ }\ : ; ' ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	Optional. The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection.
Type	The Amazon S3 connection type.
Access Key	Access key to access the Amazon S3 bucket. Provide the access key value based on the following authentication methods: <ul style="list-style-type: none"><li>- Basic authentication: provide the actual access key value.</li><li>- IAM authentication: do not provide the access key value.</li><li>- Temporary security credentials via assume role: provide access key of an IAM user with no permissions to access Amazon S3 bucket.</li></ul>
Secret Key	Secret access key to access the Amazon S3 bucket. The secret key is associated with the access key and uniquely identifies the account. Provide the access key value based on the following authentication methods: <ul style="list-style-type: none"><li>- Basic authentication: provide the actual access secret value.</li><li>- IAM authentication: do not provide the access secret value.</li><li>- Temporary security credentials via assume role: provide access secret of an IAM user with no permissions to access Amazon S3 bucket.</li></ul>
IAM Role ARN	The ARN of the IAM role assumed by the user to use the dynamically generated temporary security credentials. Enter the value of this property if you want to use the temporary security credentials to access the AWS resources. If you want to use the temporary security credentials with IAM authentication, do not provide the Access Key and Secret Key connection properties. If you want to use the temporary security credentials without IAM authentication, you must enter the value of the Access Key and Secret Key connection properties. For more information about how to obtain the ARN of the IAM role, see the AWS documentation.
Folder Path	The complete path to Amazon S3 objects. The path must include the bucket name and any folder name. Do not use a slash at the end of the folder path. For example, <bucket name>/<my folder name>.
Master Symmetric Key	Optional. Provide a 256-bit AES encryption key in the Base64 format when you enable client-side encryption. You can generate a master symmetric key using a third-party tool.

Property	Description
S3 Account Type	<p>The type of the Amazon S3 account.</p> <p>Select <b>Amazon S3 Storage</b> or <b>S3 Compatible Storage</b>.</p> <p>Select the Amazon S3 storage option to use the Amazon S3 services. Select the S3 compatible storage option to specify the endpoint for a third-party storage provider such as Scalify RING.</p> <p>By default, Amazon S3 storage is selected.</p>
REST Endpoint	<p>The S3 storage endpoint.</p> <p>Specify the S3 storage endpoint in HTTP/HTTPS format when you select the S3 compatible storage option. For example, <code>http://s3.isv.scality.com</code>.</p>
Region Name	<p>Select the AWS region in which the bucket you want to access resides.</p> <p>Select one of the following regions:</p> <ul style="list-style-type: none"> <li>- Asia Pacific (Mumbai)</li> <li>- Asia Pacific (Seoul)</li> <li>- Asia Pacific (Singapore)</li> <li>- Asia Pacific (Sydney)</li> <li>- Asia Pacific (Tokyo)</li> <li>- AWS GovCloud (US)</li> <li>- Canada (Central)</li> <li>- China (Beijing)</li> <li>- China (Hong Kong)</li> <li>- China (Ningxia)</li> <li>- EU (Ireland)</li> <li>- EU (Frankfurt)</li> <li>- EU (London)</li> <li>- EU (Paris)</li> <li>- South America (Sao Paulo)</li> <li>- US East (Ohio)</li> <li>- US East (N. Virginia)</li> <li>- US West (N. California)</li> <li>- US West (Oregon)</li> </ul> <p>Default is US East (N. Virginia).</p> <p>Not applicable for S3 compatible storage.</p>
Customer Master Key ID	<p>Optional. Specify the customer master key ID or alias name generated by AWS Key Management Service (AWS KMS) or the Amazon Resource Name (ARN) of your custom key for cross-account access. You must generate the customer master key for the same region where Amazon S3 bucket reside.</p> <p>You can specify any of the following values:</p> <p><b>Customer generated customer master key</b></p> <p>Enables client-side or server-side encryption.</p> <p><b>Default customer master key</b></p> <p>Enables client-side or server-side encryption. Only the administrator user of the account can use the default customer master key ID to enable client-side encryption.</p>
Federated SSO IdP	<p>SAML 2.0-enabled identity provider for the federated user single sign-on to use with the AWS account.</p> <p>PowerExchange for Amazon S3 supports only the ADFS 3.0 identity provider.</p> <p>Select <b>None</b> if you do not want to use federated user single sign-on.</p>

## Federated user single sign-on connection properties

Configure the following properties when you select **ADFS 3.0** in **Federated SSO IdP**:

Property	Description
Federated User Name	User name of the federated user to access the AWS account through the identity provider.
Federated User Password	Password for the federated user to access the AWS account through the identity provider.
IdP SSO URL	Single sign-on URL of the identity provider for AWS.
SAML Identity Provider ARN	ARN of the SAML identity provider that the AWS administrator created to register the identity provider as a trusted provider.
Role ARN	ARN of the IAM role assumed by the federated user.

## Blockchain Connection Properties

When you set up a blockchain connection, you must configure the connection properties.

The following table describes the general connection properties for a blockchain connection:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ]   \ : ; " ' < , > . ? /
ID	The string that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. Enter a string that you can use to identify the connection. The description cannot exceed 4,000 characters.
Swagger File Path	The absolute path of the swagger file path that contains the REST API to communicate with the blockchain. The swagger file must be a JSON file that is stored on the Data Integration Service machine. If the swagger file is in a different file format, such as YAML, convert the file to JSON format.
Base URL	Required. The base URL that is used to access assets on the blockchain.
Auth Type*	Authentication method that the run-time engine uses to connect to the REST server. You can use none, basic, digest, or OAuth.
Auth User ID*	User name to authenticate to the REST server.
Auth Password*	Password for the user name to authenticate to the REST server.

Property	Description
OAuth Consumer Key*	Required for the OAuth authentication type. Client key that is associated with the REST server.
OAuth Consumer Secret*	Required for the OAuth authentication type. Client password to connect to the REST server.
OAuth Token*	Required for the OAuth authentication type. Access token to connect to the REST server.
OAuth Token Secret*	Required for the OAuth authentication type. Password associated with the OAuth token.
Proxy Type*	Type of proxy. You can use no proxy, platform proxy, or custom.
Proxy Details*	Proxy configuration using the format <host>:<port>.
TrustStore File Path*	The absolute path of the truststore file that contains the SSL certificate.
TrustStore Password*	Password for the truststore file.
KeyStore File Path*	The absolute path of the keystore file that contains the keys and certificates required to establish a two-way secure connection with the REST server.
KeyStore Password*	Password for the keystore file.
Advanced Properties	<p>List of advanced properties to access an asset on the blockchain. Specify the advanced properties using name-value pairs that are separated by a semicolon.</p> <p>You can use the following advanced properties:</p> <ul style="list-style-type: none"> <li>- X-API-KEY. Required if you authenticate to the REST server using an API key.</li> </ul> <p>The advanced properties that you configure in the connection override the values for the corresponding advanced properties in the blockchain data object. For example, if the connection and the data object both specify a base URL, the value in the connection overrides the value in the data object.</p> <p><b>Note:</b> The advanced properties have the precedence level, <b>Operation level &gt; Object level &gt; Connection level</b>. The properties configured at the operation level will override the properties configured at the object or connection level.</p>
Cookies	<p>Required based on how the REST API is implemented. List of cookie properties to specify the cookie information that is passed to the REST server. Specify the properties using name-value pairs that are separated by a semicolon.</p> <p>The cookie properties that you configure in the connection override the values for the corresponding cookie properties in the blockchain data object.</p>
<p>* The property is ignored. To use the functionality, configure the property as an advanced property and provide a name-value pair based on the property name in the swagger file.</p> <p>For example, configure the following name-value pair to use basic authorization:</p> <pre>Authorization=Basic &lt;credentials&gt;</pre> <p><b>Note:</b> You cannot use <b>Test Connection</b> to validate a Blockchain connection.</p>	

# Cassandra Connection Properties

When you set up a Cassandra connection, you must configure the connection properties.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

The following table describes the Cassandra connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ]   \ : ; " ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. The ID must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	Optional. The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection. Not applicable for Data Engineering Streaming.
Type	The connection type. Select <b>Cassandra</b> .
Host Name	Host name or IP address of the Cassandra server.
Port	Cassandra server port number. Default is 9042.
User Name	User name to access the Cassandra server.
Password	Password corresponding to the user name to access the Cassandra server.
Default Keyspace	Name of the Cassandra keyspace to use by default.
SQL Identifier Character	Type of character that the database uses to enclose delimited identifiers in SQL or CQL queries. The available characters depend on the database type. Select <b>None</b> if the database uses regular identifiers. When the Data Integration Service generates SQL or CQL queries, the service does not place delimited characters around any identifiers. Select a character if the database uses delimited identifiers. When the Data Integration Service generates SQL or CQL queries, the service encloses delimited identifiers within this character.
SSL Mode	Select <b>disabled</b> . Not applicable for PowerExchange for Cassandra JDBC. SSL mode indicates the encryption type to use for the connection. You can choose a mode from the following SSL modes: <ul style="list-style-type: none"><li>- Disabled</li><li>- One way</li><li>- Two way</li></ul>



Property	Description
SSL Truststore Path	Not applicable for PowerExchange for Cassandra JDBC or when <b>disabled</b> SSL mode is selected. Absolute path and file name of the SSL truststore file that contains certificates of the trusted SSL server.
SSL Truststore Password	Not applicable for PowerExchange for Cassandra JDBC or when <b>disabled</b> SSL mode is selected. Password for the SSL truststore.
SSL Keystore Path	Not applicable for PowerExchange for Cassandra JDBC or when <b>disabled</b> SSL mode is selected. Absolute path and file name of the SSL keystore file that contains private keys and certificates for the SSL server.
SSL Keystore Password	Not applicable for PowerExchange for Cassandra JDBC or when <b>disabled</b> SSL mode is selected. Password for the SSL keystore.
Additional Connection Properties	<p>Enter one or more JDBC connection parameters in the following format:</p> <pre>&lt;param1&gt;=&lt;value&gt;;&lt;param2&gt;=&lt;value&gt;;&lt;param3&gt;=&lt;value&gt;</pre> <p>PowerExchange for Cassandra JDBC supports the following JDBC connection parameters:</p> <ul style="list-style-type: none"> <li>- BinaryColumnLength</li> <li>- DecimalColumnScale</li> <li>- EnableCaseSensitive</li> <li>- EnableNullInsert</li> <li>- EnablePaging</li> <li>- RowsPerPage</li> <li>- StringColumnLength</li> <li>- VTableSeparator</li> </ul>

## Confluent Kafka Connection

The Confluent Kafka connection is a Messaging connection. Use the Confluent Kafka connection to access a Kafka broker or a Confluent Kafka broker as a source or a target. You can create and manage a Confluent Kafka connection in the Developer tool or through infacmd.

When you configure a Confluent Kafka connection, you configure the following properties:

- List of Kafka brokers or Confluent Kafka brokers that the connection reads from or writes to.
- Number of seconds the Integration Service attempts to reconnect to the database if the connection fails.
- Version of the Confluent Kafka messaging broker.

## General Properties

The following table describes the general connection properties for the Confluent Kafka connection:

Property	Description
Name	Name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ] }   \ : ; " ' < , > . ? /
ID	The string that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	Description of the connection. Enter a string that you can use to identify the connection. The description cannot exceed 4,000 characters.
Location	Domain where you want to create the connection. Select the domain name.
Type	Connection type. Select <code>Messaging/ConfluentKafka</code> .

## Confluent Kafka Broker Properties

The following table describes the Kafka broker properties for the Confluent Kafka connection:

Property	Description
Kafka Broker List	Comma-separated list of Confluent Kafka brokers that maintain the configuration of the Confluent Kafka messaging broker. To specify a Confluent Kafka broker, use the following format: <code>&lt;IP address&gt;:&lt;port&gt;</code>
Retry Timeout	Number of seconds after which the Data Integration Service attempts to reconnect to the Confluent Kafka broker to read or write data. If the source or target is not available for the time you specify, the mapping execution stops to avoid any data loss.
Kafka Broker Version	Version of the Confluent Kafka messaging broker.
Additional Connection Properties	Optional. Comma-separated list of connection properties to connect to the Kafka broker.
Schema Registry URL	Location and port of the schema registry provider on which to connect.

### Additional Connection Properties

You can use the following syntax for specifying the additional connection properties:

```
request.timeout.ms=<value>,session.timeout.ms=<value>,  
fetch.max.wait.ms=<value>,heartbeat.interval.ms=<value>,  
security.protocol=SASL_PLAINTEXT,sasl.kerberos.  
service.name=<kerberos name>,sasl.mechanism=GSSAPI,  
sasl.jaas.config=com.sun.security.auth.module.  
Krb5Login Modulerequired useKeyTab=true  
doNotPrompt=true storeKey=true client=true  
keyTab="<Keytab Location>" principal="<principal>"
```

## SSL Properties

The following table describes the SSL properties for the Confluent Kafka connection:

Property	Description
SSL Mode	Optional. SSL mode indicates the encryption type to use for the connection. You can choose one of the following SSL modes: <ul style="list-style-type: none"><li>- Disabled</li><li>- One way</li><li>- Two way</li></ul> The default value is <code>Disabled</code> .
SSL TrustStore File Path	Required when <code>One way</code> SSL mode is selected. Absolute path and file name of the SSL truststore file that contains certificates of the trusted SSL server.
SSL TrustStore Password	Required when <code>One way</code> SSL mode is selected. Password for the SSL truststore.
SSL KeyStore File Path	Required when <code>Two way</code> SSL mode is selected. Absolute path and file name of the SSL keystore file that contains private keys and certificates for the SSL server.
SSL KeyStore Password	Required when <code>Two way</code> SSL mode is selected. Password for the SSL keystore.
Additional Security Properties	Optional. Comma-separated list of connection properties to connect to the Confluent Kafka broker in a secured way.

## Creating a Confluent Kafka Connection Using `infacmd`

You can use the `infacmd` command line program to create a Confluent Kafka connection.

To create a Confluent Kafka connection on UNIX, run the following command:

```
sh infacmd.sh createConnection -dn <domain name> -un <domain user> -pd <domain password>
-cn <connection name> -cid <connection id> -ct ConfluentKafka -o
"kfkBrkList='<host1:port1>,<host2:port2>,<host3:port3>'" kafkabrokerVersion='<version>'
schemaregistryurl='<schema registry URL>'"
```

For more information about the `CreateConnection` command, see the *Informatica Command Reference*.

## Databricks Connection Properties

Use the Databricks connection to run mappings on a Databricks cluster.

A Databricks connection is a cluster type connection. You can create and manage a Databricks connection in the Administrator tool or the Developer tool. You can use `infacmd` to create a Databricks connection. Configure properties in the Databricks connection to enable communication between the Data Integration Service and the Databricks cluster.

The following table describes the general connection properties for the Databricks connection:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters:~`!\$%^&*()-+={[]}\ :;'"<, > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	Optional. The description of the connection. The description cannot exceed 4,000 characters.
Connection Type	Choose Databricks.
Cluster Configuration	Name of the cluster configuration associated with the Databricks environment. Required if you do not configure the cloud provisioning configuration.
Cloud Provisioning Configuration	Name of the cloud provisioning configuration associated with a Databricks cloud platform. Required if you do not configure the cluster configuration.
Staging Directory	The directory where the Databricks Spark engine stages run-time files. If you specify a directory that does not exist, the Data Integration Service creates it at run time. If you do not provide a directory path, the run-time staging files are written to <i>/&lt;cluster staging directory&gt;/DATABRICKS</i> .
Advanced Properties	List of advanced properties that are unique to the Databricks environment. You can configure run-time properties for the Databricks environment in the Data Integration Service and in the Databricks connection. You can override a property configured at a high level by setting the value at a lower level. For example, if you configure a property in the Data Integration Service custom properties, you can override it in the Databricks connection. The Data Integration Service processes property overrides based on the following priorities: 1. Databricks connection advanced properties 2. Data Integration Service custom properties <b>Note:</b> Informatica does not recommend changing these property values before you consult with third-party documentation, Informatica documentation, or Informatica Global Customer Support. If you change a value without knowledge of the property, you might experience performance degradation or other unexpected results.

## Advanced Properties

Configure the following properties in the **Advanced Properties** of the Databricks configuration section:

### **infaspark.json.parser.mode**

Specifies the parser how to handle corrupt JSON records. You can set the value to one of the following modes:

- **DROPMALFORMED.** The parser ignores all corrupted records. Default mode.
- **PERMISSIVE.** The parser accepts non-standard fields as nulls in corrupted records.
- **FAILFAST.** The parser generates an exception when it encounters a corrupted record and the Spark application goes down.

**infaspark.json.parser.multiLine**

Specifies whether the parser can read a multiline record in a JSON file. You can set the value to true or false. Default is false. Applies only to non-native distributions that use Spark version 2.2.x and above.

**infaspark.flatfile.writer.nullValue**

When the Databricks Spark engine writes to a target, it converts null values to empty strings (" "). For example, 12, AB,"",23p09udj.

The Databricks Spark engine can write the empty strings to string columns, but when it tries to write an empty string to a non-string column, the mapping fails with a type mismatch.

To allow the Databricks Spark engine to convert the empty strings back to null values and write to the target, configure the property in the Databricks Spark connection.

Set to: TRUE

**infaspark.pythontx.exec**

Required to run a Python transformation on the Databricks Spark engine. Set to the location of the Python executable binary on the worker nodes in the Databricks cluster.

When you provision the cluster at run time, set this property in the Databricks cloud provisioning configuration. Otherwise, set on the Databricks connection.

For example, set to:

```
infaspark.pythontx.exec=/databricks/python3/bin/python3
```

**infaspark.pythontx.executorEnv.PYTHONHOME**

Required to run a Python transformation on the Databricks Spark engine. Set to the location of the Python installation directory on the worker nodes in the Databricks cluster.

When you provision the cluster at run time, set this property in the Databricks cloud provisioning configuration. Otherwise, set on the Databricks connection.

For example, set to:

```
infaspark.pythontx.executorEnv.PYTHONHOME=/databricks/python3
```

## Google Analytics Connection Properties

When you set up a Google Analytics connection, you must configure the connection properties.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

The following table describes the Google Analytics connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters:~`!\$%^&*() - + = {[]}\ :;'"<, > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. The ID must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	Optional. The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection.
Type	The connection type. Select <b>Google Analytics</b> .
Service Account ID	Specifies the client_email value present in the JSON file that you download after you create a service account.
Service Account Key	Specifies the private_key value present in the JSON file that you download after you create a service account.
APIVersion	API that PowerExchange for Google Analytics uses to read from Google Analytics reports. Select <b>Core Reporting API v3</b> . <b>Note:</b> PowerExchange for Google Analytics does not support Analytics Reporting API v4.

## Google BigQuery Connection Properties

When you set up a Google BigQuery connection, you must configure the connection properties.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

The following table describes the Google BigQuery connection properties:

Property	Description
Service Account ID	Specifies the client_email value present in the JSON file that you download after you create a service account in Google BigQuery.
Service Account Key	Specifies the private_key value present in the JSON file that you download after you create a service account in Google BigQuery.

Property	Description
Connection mode	<p>The mode that you want to use to read data from or write data to Google BigQuery.</p> <p>Select one of the following connection modes:</p> <ul style="list-style-type: none"> <li>- Simple. Flattens each field within the Record data type field as a separate field in the mapping.</li> <li>- Hybrid. Displays all the top-level fields in the Google BigQuery table including Record data type fields. PowerExchange for Google BigQuery displays the top-level Record data type field as a single field of the String data type in the mapping.</li> <li>- Complex. Displays all the columns in the Google BigQuery table as a single field of the String data type in the mapping.</li> </ul> <p>Default is Simple.</p>
Schema Definition File Path	<p>Specifies a directory on the client machine where the must create a JSON file with the sample schema of the Google BigQuery table. The JSON file name is the same as the Google BigQuery table name.</p> <p>Alternatively, you can specify a storage path in Google Cloud Storage where the must create a JSON file with the sample schema of the Google BigQuery table. You can download the JSON file from the specified storage path in Google Cloud Storage to a local machine.</p>
Project ID	<p>Specifies the project_id value present in the JSON file that you download after you create a service account in Google BigQuery.</p> <p>If you have created multiple projects with the same service account, enter the ID of the project that contains the dataset that you want to connect to.</p>
Storage Path	<p>This property applies when you read or write large volumes of data.</p> <p>Path in Google Cloud Storage where the creates a local stage file to store the data temporarily. You can either enter the bucket name or the bucket name and folder name.</p> <p>For example, enter <code>gs://&lt;bucket_name&gt;</code> or <code>gs://&lt;bucket_name&gt;/&lt;folder_name&gt;</code></p>
Dataset ID	Not applicable for PowerExchange for Google BigQuery.
Use Legacy SQL For Custom Query	Not applicable for PowerExchange for Google BigQuery.
Dataset Name for Custom Query	Not applicable for PowerExchange for Google BigQuery.
Region ID	<p>The region name where the Google BigQuery dataset resides.</p> <p>For example, if you want to connect to a Google BigQuery dataset that resides in Las Vegas region, specify <b>us-west4</b> as the <b>Region ID</b>.</p> <p><b>Note:</b> In the <b>Storage Path</b> connection property, ensure that you specify a bucket name or the bucket name and folder name that resides in the same region as the dataset in Google BigQuery.</p> <p>For more information about the regions supported by Google BigQuery, see the following Google BigQuery documentation:<a href="https://cloud.google.com/bigquery/docs/locations">https://cloud.google.com/bigquery/docs/locations</a></p>

Property	Description
Optional Properties	<p>Specifies whether you can configure certain source and target functionalities through custom properties.</p> <p>You can select one of the following options:</p> <ul style="list-style-type: none"> <li>- None. Select if you do not want to configure any custom properties.</li> <li>- Required. If you want to specify custom properties to configure the source and target functionalities.</li> </ul> <p>Default is None.</p>
Provide Optional Properties	<p>Comma-separated key-value pairs of custom properties to enable certain source and target functionalities.</p> <p>Appears only when you select <b>Required</b> in the Optional Properties.</p> <p>For more information about the list of custom properties that you can specify, see the Informatica Knowledge Base article: <a href="https://kb.informatica.com/faq/7/Pages/26/632722.aspx">https://kb.informatica.com/faq/7/Pages/26/632722.aspx</a></p>

## Google Cloud Spanner Connection Properties

When you set up a Google Cloud Spanner connection, you must configure the connection properties.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

The following table describes the Google Cloud Spanner connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { }   \ : ; " ' < , > . ? /
ID	<p>String that the Data Integration Service uses to identify the connection.</p> <p>The ID is not case sensitive. The ID must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection.</p> <p>Default value is the connection name.</p>
Description	Optional. The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection.
Type	The connection type. Select Google Cloud Spanner.
Project ID	<p>Specifies the project_id value present in the JSON file that you download after you create a service account.</p> <p>If you have created multiple projects with the same service account, enter the ID of the project that contains the bucket that you want to connect to.</p>
Service Account ID	Specifies the client_email value present in the JSON file that you download after you create a service account.



Property	Description
Service Account Key	Specifies the private_key value present in the JSON file that you download after you create a service account.
Instance ID	Name of the instance that you created in Google Cloud Spanner.

## Google Cloud Storage Connection Properties

When you set up a Google Cloud Storage connection, you must configure the connection properties.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

The following table describes the Google Cloud Storage connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { } ] [ \ : ; " ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. The ID must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	Optional. The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection.
Type	The connection type. Select <b>Google Cloud Storage</b> .
Project ID	Specifies the project_id value present in the JSON file that you download after you create a service account. If you have created multiple projects with the same service account, enter the ID of the project that contains the bucket that you want to connect to.
Service Account ID	Specifies the client_email value present in the JSON file that you download after you create a service account.
Service Account Key	Specifies the private_key value present in the JSON file that you download after you create a service account.

# Google PubSub Connection Properties

When you create a Google PubSub connection, you must configure the connection properties.

The following table describes the Google PubSub connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ] }   \ : ; " ' < , > . ? /
ID	The string that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. Enter a string that you can use to identify the connection. The description cannot exceed 4,000 characters.
Connection Type	The connection type. Select <b>Pub Sub</b> connection type.
Client Email	The <code>client_email</code> value available in the JSON file that you download after you create a service account.
Client Id	The <code>client_id</code> value available in the JSON file that you download after you create a service account.
Private Key Id	The <code>private_key_id</code> value available in the JSON file that you download after you create a service account.
Private Key	The <code>private_key</code> value available in the JSON file that you download after you create a service account.
Project Id	The <code>project_id</code> value available in the JSON file that you download after you create a service account.

# Hadoop Connection Properties

Use the Hadoop connection to configure mappings to run on a Hadoop cluster. A Hadoop connection is a cluster type connection. You can create and manage a Hadoop connection in the Administrator tool or the Developer tool. You can use `infacmd` to create a Hadoop connection. Hadoop connection properties are case sensitive unless otherwise noted.

You can configure run-time properties for the Hadoop environment in the Data Integration Service, the Hadoop connection, and in the mapping. You can override a property configured at a high level by setting the value at a lower level. For example, if you configure a property in the Data Integration Service custom properties, you can override it in the Hadoop connection or in the mapping. The Data Integration Service processes property overrides based on the following priorities:

1. Mapping custom properties set using `infacmd ms runMapping` with the `-cp` option
2. Mapping run-time properties for the Hadoop environment

3. Hadoop connection advanced properties for run-time engines
4. Hadoop connection advanced general properties, environment variables, and classpaths
5. Data Integration Service custom properties

**Note:** When a mapping uses Hive Server 2 to run a job or parts of a job, you cannot override properties that are configured on the cluster level in preSQL or post-SQL queries or SQL override statements. Workaround: Instead of attempting to use the cluster configuration on the domain to override cluster properties, pass the override settings to the JDBC URL. For example: `beeline -u "jdbc:hive2://<domain host>:<port_number>/tpch_text_100" --hiveconf hive.execution.engine=tez`

## Hadoop Cluster Properties

Configure properties in the Hadoop connection to enable communication between the Data Integration Service and the Hadoop cluster.

The following table describes the general connection properties for the Hadoop connection:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ]   \ : ; " ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. Enter a string that you can use to identify the connection. The description cannot exceed 4,000 characters.
Cluster Configuration	The name of the cluster configuration associated with the Hadoop environment. Required if you do not configure the Cloud Provisioning Configuration.
Cloud Provisioning Configuration	Name of the cloud provisioning configuration associated with a cloud platform such as Amazon AWS or Microsoft Azure. Required if you do not configure the Cluster Configuration.

Property	Description
Cluster Environment Variables*	<p>Environment variables that the Hadoop cluster uses.</p> <p>If you use a Cloudera CDH 6.x cluster or a Cloudera CDP cluster, configure the locale setting as cluster environment variables. In Cloudera Manager, you must also add the environment variables to the following YARN property:</p> <pre>yarn.nodemanager.env-whitelist</pre> <p>For example, the variable ORACLE_HOME represents the directory where the Oracle database client software is installed.</p> <p>You can configure run-time properties for the Hadoop environment in the Data Integration Service, the Hadoop connection, and in the mapping. You can override a property configured at a high level by setting the value at a lower level. For example, if you configure a property in the Data Integration Service custom properties, you can override it in the Hadoop connection or in the mapping. The Data Integration Service processes property overrides based on the following priorities:</p> <ol style="list-style-type: none"> <li>1. Mapping custom properties set using infacmd ms runMapping with the -cp option</li> <li>2. Mapping run-time properties for the Hadoop environment</li> <li>3. Hadoop connection advanced properties for run-time engines</li> <li>4. Hadoop connection advanced general properties, environment variables, and classpaths</li> <li>5. Data Integration Service custom properties</li> </ol> <p><b>Note:</b> When a mapping uses Hive Server 2 to run a job or parts of a job, you cannot override properties that are configured on the cluster level in preSQL or post-SQL queries or SQL override statements. Workaround: Instead of attempting to use the cluster configuration on the domain to override cluster properties, pass the override settings to the JDBC URL. For example: <code>beeline -u "jdbc:hive2://&lt;domain host&gt;:&lt;port_number&gt;/tpch_text_100" --hiveconf hive.execution.engine=tez</code></p>
Cluster Library Path*	<p>The path for shared libraries on the cluster.</p> <p>The \$DEFAULT_CLUSTER_LIBRARY_PATH variable contains a list of default directories.</p>
Cluster Classpath*	<p>The classpath to access the Hadoop jar files and the required libraries.</p> <p>The \$DEFAULT_CLUSTER_CLASSPATH variable contains a list of paths to the default jar files and libraries.</p> <p>You can configure run-time properties for the Hadoop environment in the Data Integration Service, the Hadoop connection, and in the mapping. You can override a property configured at a high level by setting the value at a lower level. For example, if you configure a property in the Data Integration Service custom properties, you can override it in the Hadoop connection or in the mapping. The Data Integration Service processes property overrides based on the following priorities:</p> <ol style="list-style-type: none"> <li>1. Mapping custom properties set using infacmd ms runMapping with the -cp option</li> <li>2. Mapping run-time properties for the Hadoop environment</li> <li>3. Hadoop connection advanced properties for run-time engines</li> <li>4. Hadoop connection advanced general properties, environment variables, and classpaths</li> <li>5. Data Integration Service custom properties</li> </ol> <p><b>Note:</b> When a mapping uses Hive Server 2 to run a job or parts of a job, you cannot override properties that are configured on the cluster level in preSQL or post-SQL queries or SQL override statements. Workaround: Instead of attempting to use the cluster configuration on the domain to override cluster properties, pass the override settings to the JDBC URL. For example: <code>beeline -u "jdbc:hive2://&lt;domain host&gt;:&lt;port_number&gt;/tpch_text_100" --hiveconf hive.execution.engine=tez</code></p>

Property	Description
Cluster Executable Path*	The path for executable files on the cluster. The \$DEFAULT_CLUSTER_EXEC_PATH variable contains a list of paths to the default executable files.
* Informatica does not recommend changing these property values before you consult with third-party documentation, Informatica documentation, or Informatica Global Customer Support. If you change a value without knowledge of the property, you might experience performance degradation or other unexpected results.	

## Common Properties

The following table describes the common connection properties that you configure for the Hadoop connection:

Property	Description
Impersonation User Name	Required if the Hadoop cluster uses Kerberos authentication. Hadoop impersonation user. The user name that the Data Integration Service impersonates to run mappings in the Hadoop environment.  Data Engineering Integration supports operating system profiles on all Hadoop distributions. In the Hadoop run-time environment, the Data Integration Service pushes the processing to the Hadoop cluster and the run-time engines run mappings with the operating system profile specified Hadoop impersonation properties.
Temporary Table Compression Codec	Hadoop compression library for a compression codec class name. <b>Note:</b> The Spark engine does not support compression settings for temporary tables. When you run mappings on the Spark engine, the Spark engine stores temporary tables in an uncompressed file format.
Codec Class Name	Codec class name that enables data compression and improves performance on temporary staging tables.
Hive Staging Database Name	Namespace for Hive staging tables. Use the name <code>default</code> for tables that do not have a specified database name.  If you do not configure a namespace, the Data Integration Service uses the Hive database name in the Hive target connection to create staging tables.  When you run a mapping in the native environment to write data to Hive, you must configure the Hive staging database name in the Hive connection. The Data Integration Service ignores the value you configure in the Hadoop connection.
Environment SQL	SQL commands to set the Hadoop environment. The Data Integration Service executes the environment SQL at the beginning of each Hive script generated by a HiveServer2 job.  The following rules and guidelines apply to the usage of environment SQL: <ul style="list-style-type: none"> <li>- You can use environment SQL to define Hadoop or Hive parameters that you want to use in the PreSQL commands or in custom queries.</li> <li>- If you use multiple values for the Environment SQL property, ensure that there is no space between the values.</li> </ul>

Property	Description
Engine Type	<p>The Data Integration Service uses HiveServer2 to process portions of some jobs by running HiveServer2 tasks on the Spark engine. When you import the cluster configuration through the Administrator tool, you can choose to create connections. The engine type property is populated by default based on the distribution.</p> <p>When you manually create a connection, you must configure the engine type.</p> <p>You can specify the engine type based on the following Hadoop distributions:</p> <ul style="list-style-type: none"> <li>- Amazon EMR. Tez</li> <li>- Azure HDI. Tez</li> <li>- Cloudera CDH. MRv2</li> <li>- Cloudera CDP. Tez</li> <li>- Dataproc. MRv2</li> <li>- Hortonworks HDP. Tez</li> <li>- MapR. MRv2</li> </ul>
Advanced Properties	<p>List of advanced properties that are unique to the Hadoop environment. The properties are common to the Blaze and Spark engines. The advanced properties include a list of default properties.</p> <p>You can configure run-time properties for the Hadoop environment in the Data Integration Service, the Hadoop connection, and in the mapping. You can override a property configured at a high level by setting the value at a lower level. For example, if you configure a property in the Data Integration Service custom properties, you can override it in the Hadoop connection or in the mapping. The Data Integration Service processes property overrides based on the following priorities:</p> <ol style="list-style-type: none"> <li>1. Mapping custom properties set using <code>infacmd ms runMapping</code> with the <code>-cp</code> option</li> <li>2. Mapping run-time properties for the Hadoop environment</li> <li>3. Hadoop connection advanced properties for run-time engines</li> <li>4. Hadoop connection advanced general properties, environment variables, and classpaths</li> <li>5. Data Integration Service custom properties</li> </ol> <p><b>Note:</b> When a mapping uses Hive Server 2 to run a job or parts of a job, you cannot override properties that are configured on the cluster level in preSQL or post-SQL queries or SQL override statements. Workaround: Instead of attempting to use the cluster configuration on the domain to override cluster properties, pass the override settings to the JDBC URL. For example: <code>beeline -u "jdbc:hive2://&lt;domain host&gt;:&lt;port_number&gt;/tpch_text_100" --hiveconf hive.execution.engine=tez</code></p> <p><b>Note:</b> Informatica does not recommend changing these property values before you consult with third-party documentation, Informatica documentation, or Informatica Global Customer Support. If you change a value without knowledge of the property, you might experience performance degradation or other unexpected results.</p>

## Reject Directory Properties

The following table describes the connection properties that you configure to the Hadoop Reject Directory.

Property	Description
Write Reject Files to Hadoop	<p>If you use the Blaze engine to run mappings, select the check box to specify a location to move reject files. If checked, the Data Integration Service moves the reject files to the HDFS location listed in the property, Reject File Directory.</p> <p>By default, the Data Integration Service stores the reject files based on the RejectDir system parameter.</p>
Reject File Directory	The directory for Hadoop mapping files on HDFS when you run mappings.

## Blaze Configuration

The following table describes the connection properties that you configure for the Blaze engine:

Property	Description
Blaze Staging Directory	<p>The HDFS file path of the directory that the Blaze engine uses to store temporary files. Verify that the directory exists. The YARN user, Blaze engine user, and mapping impersonation user must have write permission on this directory.</p> <p>Default is <code>/blaze/workdir</code>. If you clear this property, the staging files are written to the Hadoop staging directory <code>/tmp/blaze_&lt;user name&gt;</code>.</p>
Blaze User Name	<p>The owner of the Blaze service and Blaze service logs.</p> <p>When the Hadoop cluster uses Kerberos authentication, the default user is the Data Integration Service SPN user. When the Hadoop cluster does not use Kerberos authentication and the Blaze user is not configured, the default user is the Data Integration Service user.</p>
Minimum Port	The minimum value for the port number range for the Blaze engine. Default is 12300.
Maximum Port	The maximum value for the port number range for the Blaze engine. Default is 12600.
YARN Queue Name	<p>The YARN scheduler queue name used by the Blaze engine that specifies available resources on a cluster.</p> <p><b>Note:</b> If YARN preemption is enabled on the cluster, verify with the Hadoop administrator that preemption is disabled on the queue associated with the Blaze engine.</p>
Blaze Job Monitor Address	<p>The host name and port number for the Blaze Job Monitor.</p> <p>Use the following format:</p> <p><code>&lt;hostname&gt;:&lt;port&gt;</code></p> <p>Where</p> <ul style="list-style-type: none"><li>- <code>&lt;hostname&gt;</code> is the host name or IP address of the Blaze Job Monitor server.</li><li>- <code>&lt;port&gt;</code> is the port on which the Blaze Job Monitor listens for remote procedure calls (RPC).</li></ul> <p>For example, enter: <code>myhostname:9080</code></p>

Property	Description
Blaze YARN Node Label	<p>Node label that determines the node on the Hadoop cluster where the Blaze engine runs. If you do not specify a node label, the Blaze engine runs on the nodes in the default partition.</p> <p>If the Hadoop cluster supports logical operators for node labels, you can specify a list of node labels. To list the node labels, use the operators <code>&amp;&amp;</code> (AND), <code>  </code> (OR), and <code>!</code> (NOT).</p> <p><b>Note:</b> You cannot use node labels on a Cloudera CDH cluster.</p>
Advanced Properties	<p>List of advanced properties that are unique to the Blaze engine. The advanced properties include a list of default properties.</p> <p>You can configure run-time properties for the Hadoop environment in the Data Integration Service, the Hadoop connection, and in the mapping. You can override a property configured at a high level by setting the value at a lower level. For example, if you configure a property in the Data Integration Service custom properties, you can override it in the Hadoop connection or in the mapping. The Data Integration Service processes property overrides based on the following priorities:</p> <ol style="list-style-type: none"> <li>1. Mapping custom properties set using <code>infacmd ms runMapping</code> with the <code>-cp</code> option</li> <li>2. Mapping run-time properties for the Hadoop environment</li> <li>3. Hadoop connection advanced properties for run-time engines</li> <li>4. Hadoop connection advanced general properties, environment variables, and classpaths</li> <li>5. Data Integration Service custom properties</li> </ol> <p><b>Note:</b> When a mapping uses Hive Server 2 to run a job or parts of a job, you cannot override properties that are configured on the cluster level in preSQL or post-SQL queries or SQL override statements. Workaround: Instead of attempting to use the cluster configuration on the domain to override cluster properties, pass the override settings to the JDBC URL. For example: <code>beeline -u "jdbc:hive2://&lt;domain host&gt;:&lt;port_number&gt;/tpch_text_100" --hiveconf hive.execution.engine=tez</code></p> <p><b>Note:</b> Informatica does not recommend changing these property values before you consult with third-party documentation, Informatica documentation, or Informatica Global Customer Support. If you change a value without knowledge of the property, you might experience performance degradation or other unexpected results.</p>

## Spark Configuration

The following table describes the connection properties that you configure for the Spark engine:

Property	Description
Spark Staging Directory	<p>The HDFS file path of the directory that the Spark engine uses to store temporary files for running jobs. The YARN user, Data Integration Service user, and mapping impersonation user must have write permission on this directory.</p> <p>If you do not specify a file path, by default, the temporary files are written to the Hadoop staging directory <code>/tmp/SPARK_&lt;user name&gt;</code>.</p> <p>When you run Sqoop jobs on the Spark engine, the Data Integration Service creates a Sqoop staging directory within the Spark staging directory to store temporary files: <code>&lt;Spark staging directory&gt;/sqoop_staging</code></p>
Spark Event Log Directory	Optional. The HDFS file path of the directory that the Spark engine uses to log events.



Property	Description
YARN Queue Name	The YARN scheduler queue name used by the Spark engine that specifies available resources on a cluster. The name is case sensitive.
Advanced Properties	<p>List of advanced properties that are unique to the Spark engine. The advanced properties include a list of default properties.</p> <p>You can configure run-time properties for the Hadoop environment in the Data Integration Service, the Hadoop connection, and in the mapping. You can override a property configured at a high level by setting the value at a lower level. For example, if you configure a property in the Data Integration Service custom properties, you can override it in the Hadoop connection or in the mapping. The Data Integration Service processes property overrides based on the following priorities:</p> <ol style="list-style-type: none"> <li>1. Mapping custom properties set using infacmd ms runMapping with the -cp option</li> <li>2. Mapping run-time properties for the Hadoop environment</li> <li>3. Hadoop connection advanced properties for run-time engines</li> <li>4. Hadoop connection advanced general properties, environment variables, and classpaths</li> <li>5. Data Integration Service custom properties</li> </ol> <p><b>Note:</b> When a mapping uses Hive Server 2 to run a job or parts of a job, you cannot override properties that are configured on the cluster level in preSQL or post-SQL queries or SQL override statements. Workaround: Instead of attempting to use the cluster configuration on the domain to override cluster properties, pass the override settings to the JDBC URL. For example: <code>beeline -u "jdbc:hive2://&lt;domain host&gt;:&lt;port_number&gt;/tpch_text_100" --hiveconf hive.execution.engine=tez</code></p> <p><b>Note:</b> Informatica does not recommend changing these property values before you consult with third-party documentation, Informatica documentation, or Informatica Global Customer Support. If you change a value without knowledge of the property, you might experience performance degradation or other unexpected results.</p>

## HDFS Connection Properties

Use a Hadoop File System (HDFS) connection to access data in the Hadoop cluster. The HDFS connection is a file system type connection. You can create and manage an HDFS connection in the Administrator tool, Analyst tool, or the Developer tool. HDFS connection properties are case sensitive unless otherwise noted.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

The following table describes HDFS connection properties:

Property	Description
Name	Name of the connection. The name is not case sensitive and must be unique within the domain. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ]   \ : ; " ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. The description cannot exceed 765 characters.
Location	The domain where you want to create the connection. Not valid for the Analyst tool.

Property	Description
Type	The connection type. Default is Hadoop File System.
User Name	User name to access HDFS.
NameNode URI	<p>The URI to access the storage system. You can find the value for <code>fs.defaultFS</code> in the <code>core-site.xml</code> configuration set of the cluster configuration.</p> <p>If you create connections when you import the cluster configuration, the NameNode URI property is populated by default, and it is updated each time you refresh the cluster configuration.</p> <p>If you use a Cloudera CDP Public Cloud compute cluster and the HDFS is on a Cloudera Data Lake cluster, set the property <code>spark.yarn.access.hadoopFileSystems</code> in the Spark properties of the Hadoop Connection to the same value as set here.</p>

## Accessing Multiple Storage Types

Use the NameNode URI property in the connection parameters to connect to various storage types. The following table lists the storage type and the NameNode URI format for the storage type:

Storage	NameNode URI Format
HDFS	<p><code>hdfs://&lt;namenode&gt;:&lt;port&gt;</code></p> <p><b>where:</b></p> <ul style="list-style-type: none"> <li>- <code>&lt;namenode&gt;</code> is the host name or IP address of the NameNode.</li> <li>- <code>&lt;port&gt;</code> is the port that the NameNode listens for remote procedure calls (RPC).</li> </ul> <p><code>hdfs://&lt;nameservice&gt;</code> in case of NameNode high availability.</p>
MapR-FS	<code>maprfs:///</code>
WASB in HDInsight	<p><code>wasb://&lt;container_name&gt;@&lt;account_name&gt;.blob.core.windows.net/&lt;path&gt;</code></p> <p><b>where:</b></p> <ul style="list-style-type: none"> <li>- <code>&lt;container_name&gt;</code> identifies a specific Azure Storage Blob container.</li> </ul> <p><b>Note:</b> <code>&lt;container_name&gt;</code> is optional.</p> <ul style="list-style-type: none"> <li>- <code>&lt;account_name&gt;</code> identifies the Azure Storage Blob object.</li> </ul> <p><b>Example:</b></p> <p><code>wasb://infabdmoffering1storage.blob.core.windows.net/infabdmoffering1cluster/mr-history</code></p>
ADLS in HDInsight	<code>adl://home</code>

When you create a cluster configuration from an Azure HDInsight cluster, the cluster configuration uses either ADLS or WASB as the primary storage. You cannot create a cluster configuration with ADLS or WASB as the secondary storage. You can edit the NameNode URI property in the HDFS connection to connect to a local HDFS location.

# HBase Connection Properties

Use an HBase connection to access HBase. The HBase connection is a NoSQL connection. You can create and manage an HBase connection in the Administrator tool or the Developer tool. HBase connection properties are case sensitive unless otherwise noted.

The following table describes HBase connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ] }   \ : ; " ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection.
Type	The connection type. Select HBase.
Database Type	Type of database that you want to connect to. Select <b>HBase</b> to create a connection for an HBase table.

# HBase Connection Properties for MapR-DB

Use an HBase connection to connect to a MapR-DB table. The HBase connection is a NoSQL connection. You can create and manage an HBase connection in the Administrator tool or the Developer tool. HBase connection properties are case sensitive unless otherwise noted.

The following table describes the HBase connection properties for MapR-DB:

Property	Description
Name	Name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ] }   \ : ; " ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.

Property	Description
Description	Description of the connection. The description cannot exceed 4,000 characters.
Location	Domain where you want to create the connection.
Type	Connection type. Select <b>HBase</b> .
Database Type	Type of database that you want to connect to. Select <b>MapR-DB</b> to create a connection for a MapR-DB table.
Cluster Configuration	The name of the cluster configuration associated with the Hadoop environment.
MapR-DB Database Path	Database path that contains the MapR-DB table that you want to connect to. Enter a valid MapR cluster path.  When you create an HBase data object for MapR-DB, you can browse only tables that exist in the MapR-DB path that you specify in the <b>Database Path</b> field. You cannot access tables that are available in sub-directories in the specified path.  For example, if you specify the path as <code>/user/customers/</code> , you can access the tables in the <code>customers</code> directory. However, if the <code>customers</code> directory contains a sub-directory named <code>regions</code> , you cannot access the tables in the following directory: <code>/user/customers/regions</code>

## Hive Connection Properties

Use the Hive connection to access Hive data. A Hive connection is a database type connection. You can create and manage a Hive connection in the Administrator tool, Analyst tool, or the Developer tool. Hive connection properties are case sensitive unless otherwise noted.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

The following table describes Hive connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: <code>~ ` ! \$ % ^ &amp; * ( ) - + = { [ ] }   \ : ; " ' &lt; , &gt; . ? /</code>
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. The description cannot exceed 4000 characters.

Property	Description
Location	The domain where you want to create the connection. Not valid for the Analyst tool.
Type	The connection type. Select Hive.
LDAP username	<p>LDAP user name of the user that the Data Integration Service impersonates to run mappings on a Hadoop cluster. The user name depends on the JDBC connection string that you specify in the Metadata Connection String or Data Access Connection String for the native environment.</p> <p>If the Hadoop cluster uses Kerberos authentication, the principal name for the JDBC connection string and the user name must be the same. Otherwise, the user name depends on the behavior of the JDBC driver. With Hive JDBC driver, you can specify a user name in many ways and the user name can become a part of the JDBC URL.</p> <p>If the Hadoop cluster does not use Kerberos authentication, the user name depends on the behavior of the JDBC driver.</p> <p>If you do not specify a user name, the Hadoop cluster authenticates jobs based on the following criteria:</p> <ul style="list-style-type: none"> <li>- The Hadoop cluster does not use Kerberos authentication. It authenticates jobs based on the operating system profile user name of the machine that runs the Data Integration Service.</li> <li>- The Hadoop cluster uses Kerberos authentication. It authenticates jobs based on the SPN of the Data Integration Service. LDAP username will be ignored.</li> </ul>
Password	Password for the LDAP username.
Environment SQL	<p>SQL commands to set the Hadoop environment. In native environment type, the Data Integration Service executes the environment SQL each time it creates a connection to a Hive metastore. If you use the Hive connection to run profiles on a Hadoop cluster, the Data Integration Service executes the environment SQL at the beginning of each Hive session.</p> <p>The following rules and guidelines apply to the usage of environment SQL in both connection modes:</p> <ul style="list-style-type: none"> <li>- Use the environment SQL to specify Hive queries.</li> <li>- Use the environment SQL to set the classpath for Hive user-defined functions and then use environment SQL or PreSQL to specify the Hive user-defined functions. You cannot use PreSQL in the data object properties to specify the classpath. If you use Hive user-defined functions, you must copy the .jar files to the following directory:</li> </ul> <pre>&lt;Informatica installation directory&gt;/services/shared/hadoop/ &lt;Hadoop distribution name&gt;/extras/hive-auxjars</pre> <ul style="list-style-type: none"> <li>- You can use environment SQL to define Hadoop or Hive parameters that you want to use in the PreSQL commands or in custom queries.</li> <li>- If you use multiple values for the Environment SQL property, ensure that there is no space between the values.</li> </ul>
SQL Identifier Character	The type of character used to identify special characters and reserved SQL keywords, such as WHERE. The Data Integration Service places the selected character around special characters and reserved SQL keywords. The Data Integration Service also uses this character for the <b>Support mixed-case identifiers</b> property.

## Properties to Access Hive as Source or Target

The following table describes the connection properties that you configure to access Hive as a source or target:

Property	Description
JDBC Driver Class Name	Name of the Hive JDBC driver class. If you leave this option blank, the Developer tool uses the default Apache Hive JDBC driver shipped with the distribution. If the default Apache Hive JDBC driver does not fit your requirements, you can override the Apache Hive JDBC driver with a third-party Hive JDBC driver by specifying the driver class name.
Metadata Connection String	<p>The JDBC connection URI used to access the metadata from the Hadoop server.</p> <p>You can use PowerExchange for Hive to communicate with a HiveServer service or HiveServer2 service. To connect to HiveServer, specify the connection string in the following format:</p> <pre>jdbc:hive2://&lt;hostname&gt;:&lt;port&gt;/&lt;db&gt;</pre> <p>Where</p> <ul style="list-style-type: none"><li>- &lt;hostname&gt; is name or IP address of the machine on which HiveServer2 runs.</li><li>- &lt;port&gt; is the port number on which HiveServer2 listens.</li><li>- &lt;db&gt; is the database name to which you want to connect. If you do not provide the database name, the Data Integration Service uses the default database details.</li></ul> <p>To connect to HiveServer2, use the connection string format that Apache Hive implements for that specific Hadoop Distribution. For more information about Apache Hive connection string formats, see the Apache Hive documentation.</p> <p>For user impersonation, you must add <code>hive.server2.proxy.user=&lt;xyz&gt;</code> to the JDBC connection URI. If you do not configure user impersonation, the current user's credentials are used connect to the HiveServer2.</p> <p>If the Hadoop cluster uses SSL or TLS authentication, you must add <code>ssl=true</code> to the JDBC connection URI. For example: <code>jdbc:hive2://&lt;hostname&gt;:&lt;port&gt;/&lt;db&gt;;ssl=true</code></p> <p>If you use self-signed certificate for SSL or TLS authentication, ensure that the certificate file is available on the client machine and the Data Integration Service machine. For more information, see the <i>Data Engineering Integration Guide</i>.</p>
Bypass Hive JDBC Server	<p>JDBC driver mode. Select the check box to use the embedded JDBC driver mode.</p> <p>To use the JDBC embedded mode, perform the following tasks:</p> <ul style="list-style-type: none"><li>- Verify that Hive client and Informatica services are installed on the same machine.</li><li>- Configure the Hive connection properties to run mappings on a Hadoop cluster.</li></ul> <p>If you choose the non-embedded mode, you must configure the Data Access Connection String. Informatica recommends that you use the JDBC embedded mode.</p>
Fine Grained Authorization	<p>When you select the option to observe fine grained authorization in a Hive source, the mapping observes the following:</p> <ul style="list-style-type: none"><li>- Row and column level restrictions. Applies to Hadoop clusters where Sentry or Ranger security modes are enabled.</li><li>- Data masking rules. Applies to masking rules set on columns containing sensitive data by Dynamic Data Masking.</li></ul> <p>If you do not select the option, the Blaze and Spark engines ignore the restrictions and masking rules, and results include restricted or sensitive data.</p>

Property	Description
Data Access Connection String	<p>The connection string to access data from the Hadoop data store. To connect to HiveServer, specify the non-embedded JDBC mode connection string in the following format:</p> <pre>jdbc:hive2://&lt;hostname&gt;:&lt;port&gt;/&lt;db&gt;</pre> <p>Where</p> <ul style="list-style-type: none"> <li>- &lt;hostname&gt; is name or IP address of the machine on which HiveServer2 runs.</li> <li>- &lt;port&gt; is the port number on which HiveServer2 listens.</li> <li>- &lt;db&gt; is the database to which you want to connect. If you do not provide the database name, the Data Integration Service uses the default database details.</li> </ul> <p>To connect to HiveServer2, use the connection string format that Apache Hive implements for the specific Hadoop Distribution. For more information about Apache Hive connection string formats, see the Apache Hive documentation.</p> <p>For user impersonation, you must add <code>hive.server2.proxy.user=&lt;xyz&gt;</code> to the JDBC connection URI. If you do not configure user impersonation, the current user's credentials are used connect to the HiveServer2.</p> <p>If the Hadoop cluster uses SSL or TLS authentication, you must add <code>ssl=true</code> to the JDBC connection URI. For example: <code>jdbc:hive2://&lt;hostname&gt;:&lt;port&gt;/&lt;db&gt;;ssl=true</code></p> <p>If you use self-signed certificate for SSL or TLS authentication, ensure that the certificate file is available on the client machine and the Data Integration Service machine. For more information, see the <i>Data Engineering Integration Guide</i>.</p>
Hive Staging Directory on HDFS	<p>HDFS directory for Hive staging tables. You must grant execute permission to the Hadoop impersonation user and the mapping impersonation users.</p> <p>This option is applicable and required when you write data to a Hive target in the native environment.</p>
Hive Staging Database Name	<p>Namespace for Hive staging tables.</p> <p>The Hive Staging Database Name is automatically updated from the Data Access Connection String. If you want to override the default name, you need to configure the Hive Staging Database Name in the Hive connection.</p> <p>This option is applicable when you run a mapping in the native environment to write data to a Hive target.</p> <p>If you run the mapping on the Blaze or Spark engine, you do not need to configure the Hive staging database name in the Hive connection. The Data Integration Service uses the value that you configure in the Hadoop connection.</p>

## JDBC Connection Properties

You can use a JDBC connection to access tables in a database. You can create and manage a JDBC connection in the Administrator tool, the Developer tool, or the Analyst tool.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

The following table describes JDBC connection properties:

Property	Description
Database Type	The database type.
Name	Name of the connection. The name is not case sensitive and must be unique within the domain. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ] }   \ : ; " ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. The description cannot exceed 765 characters.
User Name	The database user name.
Password	The password for the database user name.
JDBC Driver Class Name	<p>Name of the JDBC driver class.</p> <p>The following list provides the driver class name that you can enter for the applicable database type:</p> <ul style="list-style-type: none"> <li>- DataDirect JDBC driver class name for Oracle: <code>com.informatica.jdbc.oracle.OracleDriver</code></li> <li>- DataDirect JDBC driver class name for IBM DB2: <code>com.informatica.jdbc.db2.DB2Driver</code></li> <li>- DataDirect JDBC driver class name for Microsoft SQL Server: <code>com.informatica.jdbc.sqlserver.SQLServerDriver</code></li> <li>- DataDirect JDBC driver class name for Sybase ASE: <code>com.informatica.jdbc.sybase.SybaseDriver</code></li> <li>- DataDirect JDBC driver class name for Informix: <code>com.informatica.jdbc.informix.InformixDriver</code></li> <li>- DataDirect JDBC driver class name for MySQL: <code>com.informatica.jdbc.mysql.MySQLDriver</code></li> <li>- JDBC driver for Databricks Delta Lake: the name of the driver that you downloaded from Databricks. For information about the driver, see the topic on configuring storage access in the "Before You Begin Databricks Integration" chapter of the <i>Data Engineering Integration Guide</i>.</li> </ul> <p>For more information about which driver class to use with specific databases, see the vendor documentation.</p>
Connection String	<p>Connection string to connect to the database. Use the following connection string:</p> <pre>jdbc:&lt;subprotocol&gt;:&lt;subname&gt;</pre> <p>For more information about the connection string to use with specific drivers, see the vendor documentation.</p>
Environment SQL	<p>Optional. Enter SQL commands to set the database environment when you connect to the database. The Data Integration Service executes the connection environment SQL each time it connects to the database.</p> <p><b>Note:</b> If you enable Sqoop, Sqoop ignores this property.</p>
Transaction SQL	<p>Optional. Enter SQL commands to set the database environment when you connect to the database. The Data Integration Service executes the transaction environment SQL at the beginning of each transaction.</p> <p><b>Note:</b> If you enable Sqoop, Sqoop ignores this property.</p>



Property	Description
SQL Identifier Character	<p>Type of character that the database uses to enclose delimited identifiers in SQL queries. The available characters depend on the database type.</p> <p>Select (None) if the database uses regular identifiers. When the Data Integration Service generates SQL queries, the service does not place delimited characters around any identifiers.</p> <p>Select a character if the database uses delimited identifiers. When the Data Integration Service generates SQL queries, the service encloses delimited identifiers within this character.</p> <p><b>Note:</b> If you enable Sqoop, Sqoop ignores this property.</p>
Support Mixed-case Identifiers	<p>Enable if the database uses case-sensitive identifiers. When enabled, the Data Integration Service encloses all identifiers within the character selected for the <b>SQL Identifier Character</b> property.</p> <p>When the <b>SQL Identifier Character</b> property is set to none, the <b>Support Mixed-case Identifiers</b> property is disabled.</p> <p><b>Note:</b> If you enable Sqoop, Sqoop honors this property when you generate and execute a DDL script to create or replace a target at run time. In all other scenarios, Sqoop ignores this property.</p>
Use Sqoop Connector	<p>Enables Sqoop connectivity for the data object that uses the JDBC connection. The Data Integration Service runs the mapping in the Hadoop run-time environment through Sqoop.</p> <p>You can configure Sqoop connectivity for relational data objects, customized data objects, and logical data objects that are based on a JDBC-compliant database.</p> <p>Select <b>Sqoop v1.x</b> to enable Sqoop connectivity.</p> <p>Default is <b>None</b>.</p>
Sqoop Arguments	<p>Enter the arguments that Sqoop must use to connect to the database. Separate multiple arguments with a space.</p> <p>To run the mapping on the Blaze engine with the Teradata Connector for Hadoop (TDCH) specialized connectors for Sqoop, you must define the TDCH connection factory class in the Sqoop arguments. The connection factory class varies based on the TDCH Sqoop Connector that you want to use.</p> <ul style="list-style-type: none"> <li>- To use Cloudera Connector Powered by Teradata, configure the following Sqoop argument: <ul style="list-style-type: none"> <li>- <code>Dsqoop.connection.factories=com.cloudera.connector.teradata.TeradataManagerFactory</code></li> </ul> </li> <li>- To use Hortonworks Connector for Teradata (powered by the Teradata Connector for Hadoop), configure the following Sqoop argument: <ul style="list-style-type: none"> <li>- <code>Dsqoop.connection.factories=org.apache.sqoop.teradata.TeradataManagerFactory</code></li> </ul> </li> </ul> <p>To run the mapping on the Spark engine, you do not need to define the TDCH connection factory class in the Sqoop arguments. The Data Integration Service invokes the Cloudera Connector Powered by Teradata and Hortonworks Connector for Teradata (powered by the Teradata Connector for Hadoop) by default.</p> <p><b>Note:</b> To run the mapping with a generic JDBC connector instead of the specialized Cloudera or Hortonworks connector, you must define the <code>--driver</code> and <code>--connection-manager</code> Sqoop arguments in the JDBC connection. If you define the <code>--driver</code> and <code>--connection-manager</code> arguments in the Read or Write transformation of the mapping, Sqoop ignores the arguments.</p> <p>If you do not enter Sqoop arguments, the Data Integration Service constructs the Sqoop command based on the JDBC connection properties.</p>

## JDBC Connection String

Specify the connection string in the JDBC connection to connect to the database.

Specify the connection string in the following format:

```
jdbc:<subprotocol>:<subname>
```

Use the following sample connection strings that you can enter for the applicable database type:

- **DataDirect Oracle JDBC driver:**  
`jdbc:informatica:oracle://<host>:<port>;SID=<value>`
- **Oracle JDBC driver:**  
`jdbc:oracle:thin:@//<host>:<port>:<SID>`
- **DataDirect IBM DB2 JDBC driver:**  
`jdbc:informatica:db2://<host>:<port>;DatabaseName=<value>`
- **IBM DB2 JDBC driver:**  
`jdbc:db2://<host>:<port>/<database_name>`
- **DataDirect Microsoft SQL Server JDBC driver:**  
`jdbc:informatica:sqlserver://<host>;DatabaseName=<value>`
- **Microsoft SQL Server JDBC driver:**  
`jdbc:sqlserver://<host>;DatabaseName=<value>`
- **Netezza JDBC driver:**  
`jdbc:netezza://<host>:<port>/<database_name>`
- **Pivotal Greenplum driver:**  
`jdbc:pivotal:greenplum://<host>:<port>;/database_name=<value>`
- **Postgres Greenplum driver:**  
`jdbc:postgresql://<host>:<port>/<database_name>`
- **Teradata JDBC driver:**  
`jdbc:teradata://<host>/database_name=<value>,tmode=<value>,charset=<value>`
- **JDBC driver for Delta Lake:**  
`jdbc:spark://<host name>:443/default;transportMode=http;ssl=1;httpPath=sql/protocolv1/o/<cluster ID>;AuthMech=3;UID=token;PWD=<token string>`

Use the following sample connection strings that you can enter for an SSL-enabled applicable database type:

- **DataDirect Oracle JDBC driver:**  
`jdbc:informatica:oracle://  
<host_name>:<port>;CatalogOptions=6;ServiceName=<service_name>;trustStorePassword=<tru  
ststore_password>;  
keyStorePassword=<keystore_password>;CryptoProtocolVersion=TLSv1.2;keyStore=<keystore_  
location_of_ewallet.p12_file>;  
trustStore=<truststore_location_of_truststore.p12_file>;HostNameInCertificate=<databas  
e_host_name>;encryptionMethod=SSL;  
ValidateServerCertificate=True;`
- **Oracle JDBC driver:**  
`jdbc:oracle:thin:@(DESCRIPTION=(ADDRESS=(PROTOCOL=TCPS)(HOST=<host>)  
(PORT=<port_number>))  
(CONNECT_DATA=(SERVICE_NAME=<service_name>))) "`
- **DataDirect Microsoft SQL Server JDBC driver:**  
`jdbc:informatica:sqlserver://  
<host_name>:<port>;databaseName=<database_name>;EncryptionMethod=SSL;  
CryptoProtocolVersion=<TLSv1.2_or_TLSv1.1_or_TLSv1>;ValidateServerCertificate=false;Tr  
ustStore=<truststore_location>;  
TrustStorePassword=<truststore_password>`
- **Microsoft SQL Server JDBC driver:**  
`jdbc:sqlserver://  
<host_name>:<port>;databaseName=<database_name>;integratedSecurity=false;encrypt=true;  
trustServerCertificate=true;  
TrustStore=/  
<truststore_location>;TrustStorePassword=<truststore_password>;user=<user_name>;passwo  
rd=<password>`

For more information about the connection string to use with specific drivers, see the vendor documentation.

## Sqoop Connection-Level Arguments

In the JDBC connection, you can define the arguments that Sqoop must use to connect to the database. The Data Integration Service merges the arguments that you specify with the default command that it constructs based on the JDBC connection properties. The arguments that you specify take precedence over the JDBC connection properties.

If you want to use the same driver to import metadata and run the mapping, and do not want to specify any additional Sqoop arguments, select **Sqoop v1.x** from the **Use Sqoop Version** list and leave the **Sqoop Arguments** field empty in the JDBC connection. The Data Integration Service constructs the Sqoop command based on the JDBC connection properties that you specify.

However, if you want to use a different driver for run-time tasks or specify additional run-time Sqoop arguments, select **Sqoop v1.x** from the **Use Sqoop Version** list and specify the arguments in the **Sqoop Arguments** field.

A mapping that contains an Update Strategy transformation cannot use a Sqoop-enabled JDBC connection to write to a target. To run the mapping, disable the Sqoop connector in the Write transformation.

You can configure the following Sqoop arguments in the JDBC connection:

### driver

Defines the JDBC driver class that Sqoop must use to connect to the database.

Use the following syntax:

```
--driver <JDBC driver class>
```

For example, use the following syntax depending on the database type that you want to connect to:

- **Aurora:** `--driver com.mysql.jdbc.Driver`
- **Greenplum:** `--driver org.postgresql.Driver`
- **IBM DB2:** `--driver com.ibm.db2.jcc.DB2Driver`
- **IBM DB2 z/OS:** `--driver com.ibm.db2.jcc.DB2Driver`
- **Microsoft SQL Server:** `--driver com.microsoft.sqlserver.jdbc.SQLServerDriver`
- **Netezza:** `--driver org.netezza.Driver`
- **Oracle:** `--driver oracle.jdbc.driver.OracleDriver`
- **Teradata:** `--driver com.teradata.jdbc.TeraDriver`

### connect

Defines the JDBC connection string that Sqoop must use to connect to the database. The JDBC connection string must be based on the driver that you define in the driver argument.

Use the following syntax:

```
--connect <JDBC connection string>
```

For example, use the following syntax depending on the database type that you want to connect to:

- **Aurora:** `--connect "jdbc:mysql://<host_name>:<port>/<schema_name>"`
- **Greenplum:** `--connect jdbc:postgresql://<host_name>:<port>/<database_name>`
- **IBM DB2:** `--connect jdbc:db2://<host_name>:<port>/<database_name>`
- **IBM DB2 z/OS:** `--connect jdbc:db2://<host_name>:<port>/<database_name>`

- **Microsoft SQL Server:** `--connect jdbc:sqlserver://<host_name>:<port or named_instance>;databaseName=<database_name>`
- **Netezza:** `--connect "jdbc:netezza://<database_server_name>:<port>/<database_name>;schema=<schema_name>"`
- **Oracle:** `--connect jdbc:oracle:thin:@<database_host_name>:<database_port>:<database_SID>`
- **Teradata:** `--connect jdbc:teradata://<host_name>/database=<database_name>`

Use the following syntax to connect to an SSL-enabled database:

```
--connect <JDBC connection string>
```

For example, use the following syntax depending on the database type that you want to connect to:

- **Microsoft SQL Server:** `--connect jdbc:sqlserver://<host_name>:<port>;databaseName=<database_name>;integratedSecurity=false;encrypt=true;trustServerCertificate=true;TrustStore=/<truststore_location>;TrustStorePassword=<truststore_password>;user=<user_name>;password=<password>`
- **Oracle:** `--connect jdbc:oracle:thin:@(DESCRIPTION=(ADDRESS=(PROTOCOL=TCPS) (HOST=<host>) (PORT=<port_number>)) (CONNECT_DATA=(SERVICE_NAME=<service_name>))) "`

#### connection-param-file

Defines the extra JDBC parameters through a property file that Sqoop must use to connect to the database. The contents of this file are parsed as standard Java properties and passed into the driver when you create a connection.

Use the following syntax:

```
--connection-param-file <parameter file name>
```

For example, use the following syntax to use the parameter file when you connect to the Oracle database.

```
--connection-param-file param_file
```

#### connection-manager

Defines the connection manager class name that Sqoop must use to connect to the database.

Use the following syntax:

```
--connection-manager <connection manager class name>
```

For example, use the following syntax to use the generic JDBC manager class name:

```
--connection-manager org.apache.sqoop.manager.GenericJdbcManager
```

#### direct

When you read data from or write data to Oracle, you can configure the `direct` argument to enable Sqoop to use OraOop. OraOop is a specialized Sqoop plug-in for Oracle that uses native protocols to connect to the Oracle database. When you configure OraOop, the performance improves.

You can configure OraOop when you run Sqoop mappings on the Spark engine.

Use the following syntax:

```
--direct
```

When you use OraOop, you must use the following syntax to specify multiple arguments:

```
-D<argument=value> -D<argument=value>
```

**Note:** If you specify multiple arguments and include a space character between -D and the argument name-value pair, Sqoop considers only the first argument and ignores the remaining arguments.

If you do not direct the job to a specific queue, the Spark engine uses the default queue.

#### **-Dsqoop.connection.factories**

To run the mapping on the Blaze engine with the Teradata Connector for Hadoop (TDCH) specialized connectors for Sqoop, you must configure the -Dsqoop.connection.factories argument. Use the argument to define the TDCH connection factory class that Sqoop must use. The connection factory class varies based on the TDCH Sqoop Connector that you want to use.

- To use Cloudera Connector Powered by Teradata, configure the -Dsqoop.connection.factories argument as follows:  
`-Dsqoop.connection.factories=com.cloudera.connector.teradata.TeradataManagerFactory`
- To use Hortonworks Connector for Teradata (powered by the Teradata Connector for Hadoop), configure the -Dsqoop.connection.factories argument as follows:  
`-Dsqoop.connection.factories=org.apache.sqoop.teradata.TeradataManagerFactory`

**Note:** To run the mapping on the Spark engine, you do not need to configure the -Dsqoop.connection.factories argument. The Data Integration Service invokes Cloudera Connector Powered by Teradata and Hortonworks Connector for Teradata (powered by the Teradata Connector for Hadoop) by default.

#### **--infaoptimize**

Use this argument to disable the performance optimization of Sqoop pass-through mappings on the Spark engine.

When you run a Sqoop pass-through mapping on the Spark engine, the Data Integration Service optimizes mapping performance in the following scenarios:

- You read data from a Sqoop source and write data to a Hive target that uses the Text format.
- You read data from a Sqoop source and write data to an HDFS target that uses the Flat, Avro, or Parquet format.

If you want to disable the performance optimization, set the --infaoptimize argument to false. For example, if you see data type issues after you run an optimized Sqoop mapping, you can disable the performance optimization.

Use the following syntax:

```
--infaoptimize false
```

For a complete list of the Sqoop arguments that you can configure, see the Sqoop documentation.

## Delta Lake JDBC Connection Properties

To enable the domain to access Delta Lake resources in the cloud platform environment, you must manually create and configure a JDBC connection.

Configure the following connection properties:

Property	Description
Name	Type a name for the connection. Example: DatabricksDeltaCxn
ID	Default: Automatically populated with the connection name. Changing this default value is optional.
Username	Type the following value to use the Databricks user token: <code>token</code>
Password	Value of the token that you configured for the Databricks user.
JDBC Driver Class Name	Type the following value: <code>com.simba.spark.jdbc4.Driver</code>
Connection String	<p>Connection to the Delta Lake resource. This connection string contains all the information that the domain needs to connect to the resource.</p> <p>The connection string contains the following elements:</p> <ul style="list-style-type: none"><li>- <code>jdbc:spark://&lt;server host name&gt;</code>.</li><li>- Port number.</li><li>- Transport mode.</li><li>- <code>ssl</code>. Use <code>1</code> to enable SSL.</li><li>- <code>httpPath</code>.</li><li>- <code>UID</code>. User ID that will be used to run jobs on the cluster. Use <code>token</code>.</li><li>- <code>PWD</code>. Value of the token that you configured for the Databricks user.</li></ul> <p>Example:</p> <pre>jdbc:spark://westus.azuredatabricks.net:443/default;transportMode=http;ssl=1;httpPath=sql/protocolv1/o/1654523072521724/0123-456789-films259;AuthMech=3;UID=token;PWD=&lt;token string&gt;</pre> <p>To get the value of these parameters from the Advanced Options area of the cluster configuration settings:</p> <ol style="list-style-type: none"><li>1. In the Databricks environment, select Clusters.</li><li>2. Select the cluster to connect to.</li><li>3. Expand the Advanced Options and click the <b>JDBC/ODBC</b> tab.</li></ol> <p>For more information about the JDBC connection string for Databricks, see the <a href="#">Databricks documentation</a>.</p>

## JDBC V2 Connection Properties

When you set up a JDBC V2 connection, you must configure the connection properties.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

The following table describes the JDBC V2 connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters:~`!\$%^&*()-+={} \:;"'<,>.?/
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection.
Type	The connection type. Select JDBC V2.

The **Details** tab contains the connection attributes of the JDBC V2 connection. The following table describes the connection attributes:

Property	Description
Username	The database user name. User name with permissions to access the database that supports the Type 4 JDBC driver.
Password	The password for the database user name.
Schema Name	Optional. The schema name to connect in the database. If you do not specify the schema name, all the schemas available in the database are listed.
JDBC Driver Class Name	Name of the JDBC driver class. The following list provides the driver class name that you can enter for the applicable database type: <ul style="list-style-type: none"> <li>- JDBC driver class name for Azure SQL Database: com.microsoft.sqlserver.jdbc.SQLServerDriver</li> <li>- JDBC driver class name for Aurora PostgreSQL: org.postgresql.Driver</li> <li>- JDBC driver class name for SAP HANA Database: com.sap.db.jdbc.Driver</li> </ul> For more information about which driver class to use with specific databases, see the third-party vendor documentation.
Connection String	Connection string to connect to the database. Use the following connection string: jdbc:<subprotocol>:<subname>  The following list provides sample connection strings that you can enter for the applicable database type: <ul style="list-style-type: none"> <li>- Connection string for Azure SQL Database JDBC driver: jdbc:sqlserver://&lt;host&gt;:&lt;port&gt;;database=&lt;database_name&gt;</li> <li>- Connection string for Aurora PostgreSQL JDBC driver: jdbc:postgresql://&lt;host&gt;:&lt;port&gt;[/&lt;database_name&gt;]</li> <li>- Connection string for SAP HANA Database driver: jdbc:sap://&lt;host&gt;:&lt;port&gt;/?databaseName=&lt;Database_Name&gt;</li> </ul> For more information about the connection string to use with specific drivers, see the third-party vendor documentation.

Property	Description
Sub Type	<p>The database type to which you want to connect.</p> <p>You can select from the following database types to connect:</p> <ul style="list-style-type: none"> <li>- <b>Azure SQL Database</b>. Connects to Azure SQL database.</li> <li>- <b>PostgreSQL</b>. Connects to Aurora PostgreSQL database.</li> <li>- <b>SAP HANA Database</b>. Connects to SAP HANA database.</li> <li>- <b>Others</b> . Connects to any database that supports the Type 4 JDBC driver.</li> </ul>
Support Mixed-case Identifiers	<p>Enable if the database uses case-sensitive identifiers. When enabled, the Data Integration Service encloses all identifiers within the character selected for the <b>SQL Identifier Character</b> property.</p> <p>For example, Aurora PostgreSQL database supports mixed-cased characters. You must enable this property to connect to the Aurora PostgreSQL database.</p> <p>When the <b>SQL Identifier Character</b> property is set to none, the <b>Support Mixed-case Identifiers</b> property is disabled.</p>
SQL Identifier Character	<p>Type of character that the database uses to enclose delimited identifiers in SQL queries. The available characters depend on the database type.</p> <p>Select (None) if the database uses regular identifiers. When the Data Integration Service generates SQL queries, the service does not place delimited characters around any identifiers.</p> <p>Select a character if the database uses delimited identifiers. When the Data Integration Service generates SQL queries, the service encloses delimited identifiers within this character.</p> <p><b>Note:</b> Select <b>SQL Identifier Character</b> as <b>None</b> when you specify the SAP HANA Database subtype.</p>

## Kafka Connection Properties

The Kafka connection is a messaging connection. Use the Kafka connection to access Kafka as a message target. You can create and manage a Kafka connection in the Developer tool or through infacmd.

The following table describes the general connection properties for the Kafka connection:

Property	Description
Name	<p>The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters:</p> <p>~ ` ! \$ % ^ &amp; * ( ) - + = { [ ] }   \ : ; " ' &lt; , &gt; . ? /</p>
ID	<p>The string that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.</p>
Description	<p>The description of the connection. Enter a string that you can use to identify the connection. The description cannot exceed 4,000 characters.</p>
Location	<p>The domain where you want to create the connection.</p>
Type	<p>The connection type.</p>



The following table describes the Kafka broker properties for the Kafka connection:

Property	Description
Kafka Broker List	Comma-separated list of Kafka brokers which maintains the configuration of the Kafka messaging broker. To specify a Kafka broker, use the following format: <IP Address>:<port>
ZooKeeper Host Port List	Optional. Comma-separated list of Apache ZooKeeper which maintains the configuration of the Kafka messaging broker. To specify the ZooKeeper, use the following format: <IP Address>:<port>
Retry Timeout	Number of seconds the Integration Service attempts to reconnect to the Kafka broker to write data. If the source or target is not available for the time you specify, the mapping execution stops to avoid any data loss.
Kafka Broker Version	Configure the Kafka messaging broker version to 0.10.1.x-2.0.0.

## General Properties

The following table describes the general connection properties for the Kafka connection:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ] }   \ : ; " ' < , > . ? /
ID	The string that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. Enter a string that you can use to identify the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection. Select the domain name.
Type	The connection type. Select Messaging/Kafka.

## Kafka Broker Properties

The following table describes the Kafka broker properties for the Kafka connection:

Property	Description
Kafka Broker List	<p>Comma-separated list of Kafka brokers which maintain the configuration of the Kafka messaging broker.</p> <p>To specify a Kafka broker, use the following format:</p> <pre>&lt;IP Address&gt;:&lt;port&gt;</pre>
Retry Timeout	<p>Number of seconds after which the Integration Service attempts to reconnect to the Kafka broker to read or write data. If the source or target is not available for the time you specify, the mapping execution stops to avoid any data loss.</p>
Kafka Broker Version	<p>Configure the Kafka messaging broker version to Apache 0.10.1.1 &amp; above.</p>
Additional Connection Properties	<p>Optional. Comma-separated list of connection properties to connect to the Kafka broker.</p> <p>For example, you can use the following syntax:</p> <pre>request.timeout.ms=&lt;value&gt;,session.timeout.ms=&lt;value&gt;, fetch.max.wait.ms=&lt;value&gt;,heartbeat.interval.ms=&lt;value&gt;, security.protocol=SASL_PLAINTEXT,sasl.kerberos. service.name=&lt;kerberos_name&gt;,sasl.mechanism=GSSAPI, sasl.jaas.config=com.sun.security.auth.module. Krb5Login Modulerequired useKeyTab=true doNotPrompt=true storeKey=true client=true keyTab="&lt;Keytab Location&gt;" principal="&lt;principal&gt;";</pre> <p>To reduce the time taken to connect to the Kafka broker, ensure that you set the following properties:</p> <ul style="list-style-type: none"><li>- request.timeout.ms</li><li>- session.timeout.ms</li><li>- fetch.max.wait.ms</li><li>- heartbeat.interval.ms</li></ul> <p>To connect to the Kafka broker in a secured way, ensure that you set one of the following values for the security.protocol property:</p> <ul style="list-style-type: none"><li>- SASL_SSL</li><li>- SSL</li></ul> <p>The default value of security.protocol property is SASL_PLAINTEXT.</p> <p><b>Technical Preview:</b> The Additional Connection Properties is available for technical preview. Technical preview functionality is supported but is unwarranted and is not production-ready. Informatica recommends that you use in non-production environments only.</p> <p>For more information about the connection properties, see <a href="https://kafka.apache.org/documentation/">https://kafka.apache.org/documentation/</a>.</p>

## SSL Properties

The following table describes the SSL properties for the Kafka connection:

Property	Description
SSL Mode	Required. SSL mode indicates the encryption type to use for the connection. You can choose a mode from the following SSL modes: <ul style="list-style-type: none"><li>- Disabled</li><li>- One way</li><li>- Two way</li></ul>
SSL TrustStore File Path	Required when <b>One way</b> SSL mode is selected. Absolute path and file name of the SSL truststore file that contains certificates of the trusted SSL server.
SSL TrustStore Password	Required when <b>One way</b> SSL mode is selected. Password for the SSL truststore.
SSL KeyStore File Path	Required when <b>Two way</b> SSL mode is selected. Absolute path and file name of the SSL keystore file that contains private keys and certificates for the SSL server.
SSL KeyStore Password	Required when <b>Two way</b> SSL mode is selected. Password for the SSL keystore.

## Creating a Kafka Connection Using infacmd

You can use the infacmd command line program to create a Kafka connection.

To create a Kafka connection on UNIX, run the following command:

```
sh infacmd.sh createConnection -dn <domain name> -un <domain user> -pd <domain password> -cn  
<connection name> -cid <connection id> -ct Kafka -o  
kfkBrkList=<host1:port1>,<host2:port2>,<host3:port3> kafkabrokerVersion=<version>  
additionalConnectionProperties=<additional properties>
```

For more information about the CreateConnection command, see the *Informatica Command Reference*.

## Kudu Connection Properties

Use a Kudu connection to access Kudu.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

You can create and manage a Kudu connection in the Administrator tool or the Developer tool. The following table describes the Kudu connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { } ]   \ : ; " ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection.
Type	The connection type. Select Kudu.

The following table describes the properties for metadata access:

Property	Description
Kudu Master URLs	The URLs of the Kudu master tables.
Kudu Library Version	The version number of the Kudu library.
Cluster Configuration	The Hadoop cluster that you use for the connection.

## Microsoft Azure Blob Storage Connection Properties

Use a Microsoft Azure SQL Blob Storage connection to access a Microsoft Azure Blob Storage.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

You can create and manage a Microsoft Azure Blob Storage connection in the Administrator tool or the Developer tool. The following table describes the Microsoft Azure Blob Storage connection properties:

Property	Description
Name	Name of the Microsoft Azure Blob Storage connection.
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	Description of the connection.

Property	Description
Location	The domain where you want to create the connection.
Type	Type of connection. Select Azure Blob Storage.

The **Connection Details** tab contains the connection attributes of the Microsoft Azure Blob Storage connection. The following table describes the connection attributes:

Property	Description
Account Name	Name of the Microsoft Azure Storage account.
Authorization Type	Authorization type. You can select any of the following authorization mechanisms: <ul style="list-style-type: none"> <li>- Shared Key Authorization</li> <li>- Shared Access Signatures</li> </ul>
Account Key	Microsoft Azure Storage access key. Applies when you select shared key authorization.
SAS Token	SAS URI with SAS token that you generate on Microsoft Azure portal for your account. Applies when you select shared access signature authorization type. <b>Note:</b> You must provide a valid SAS URI with a valid SAS token.
Container Name	The root container or sub-folders with the absolute path. <b>Note:</b> To import complex files, specify only the root container.
Endpoint Suffix	Type of Microsoft Azure end-points. You can select any of the following end-points: <ul style="list-style-type: none"> <li>- <code>core.windows.net</code>: Default</li> <li>- <code>core.usgovcloudapi.net</code>: To select the US government Microsoft Azure end-points</li> <li>- <code>core.chinacloudapi.cn</code>: Not applicable</li> </ul>

## Microsoft Azure Cosmos DB SQL API Connection Properties

Use a Microsoft Azure Cosmos DB connection to connect to the Cosmos DB database. When you create a Microsoft Azure Cosmos DB connection, you enter information for metadata and data access.

The following table describes the Microsoft Azure Cosmos DB connection properties:

Property	Description
Name	Name of the Cosmos DB connection.
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.

Property	Description
Description	Description of the connection. The description cannot exceed 765 characters.
Location	The project or folder in the Model repository where you want to store the Cosmos DB connection.
Type	Select Microsoft Azure Cosmos DB SQL API.
Cosmos DB URI	The URI of Microsoft Azure Cosmos DB account.
Key	The primary and secondary key to which provides you complete administrative access to the resources within Microsoft Azure Cosmos DB account.
Database	Name of the database that contains the collections from which you want to read or write JSON documents.

**Note:** You can find the Cosmos DB URI and Key values in the **Keys** settings on Azure portal. Contact your Azure administrator for more details.

## Microsoft Azure Data Lake Storage Gen1 Connection Properties

Use a Microsoft Azure Data Lake Storage Gen1 connection to access a Microsoft Azure Data Lake Storage Gen1.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

You can create and manage a Microsoft Azure Data Lake Storage Gen1 connection in the Administrator tool or the Developer tool. The following table describes the Microsoft Azure Data Lake Storage Gen1 connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ]   \ : ; " ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection.
Type	The connection type. Select Microsoft Azure Data Lake Storage Gen1.

The following table describes the properties for metadata access:

Property	Description
ADLS Account Name	The name of the Microsoft Azure Data Lake Storage Gen1.
ClientID	The ID of your application to complete the OAuth Authentication in the Active Directory.
Client Secret	The client secret key to complete the OAuth Authentication in the Active Directory.
Directory	The Microsoft Azure Data Lake Storage Gen1 directory that you use to read data or write data. The default is root directory.
AuthEndpoint	The OAuth 2.0 token endpoint from where access code is generated based on based on the Client ID and Client secret is completed.

For more information about creating a client ID, client secret, and auth end point, contact the Azure administrator or see Microsoft Azure Data Lake Storage Gen1 documentation.

## Microsoft Azure Data Lake Storage Gen2 Connection Properties

Use a Microsoft Azure Data Lake Storage Gen2 connection to access a Microsoft Azure Data Lake Storage Gen2.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

You can create and manage a Microsoft Azure Data Lake Storage Gen2 connection in the Administrator tool or the Developer tool. The following table describes the Microsoft Azure Data Lake Storage Gen2 connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ]   \ ; ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection.
Type	The connection type. Select Microsoft Azure Data Lake Storage Gen2.

The following table describes the properties for metadata access:

Property	Description
Account Name	The Microsoft Azure Data Lake Storage Gen2 account name or the service name.
Authentication Type	Authentication type to access the Microsoft Azure Data Lake Storage Gen2 account. Select one of the following options: <ul style="list-style-type: none"><li>- <b>Service Principal Authentication.</b> Select to use the client ID, client secret, and tenant ID to connect to Microsoft Azure Data Lake Storage Gen2.</li><li>- <b>Shared Key Authentication.</b> Select to use the account key to connect to Microsoft Azure Data Lake Storage Gen2.</li></ul> <b>Note:</b> You cannot use shared key authentication in a streaming mapping.
Client ID	Applicable for Service Principal Authentication. The ID of your application to complete the OAuth Authentication in the Azure Active Directory (AD).
Client Secret	Applicable for Service Principal Authentication. The client secret key to complete the OAuth Authentication in the Azure Active Directory.
Tenant ID	Applicable for Service Principal Authentication. The Directory ID of the Azure Active Directory.
Account Key	Applicable for Shared Key Authentication. The account key for the Microsoft Azure Data Lake Storage Gen2 account.
File System Name	The name of an existing file system in the Microsoft Azure Data Lake Storage Gen2.
Directory Path	The path of an existing directory without the file system name. There is no default directory. You can select one of the following syntax: <ul style="list-style-type: none"><li>- / for root directory.</li><li>- /dir1</li><li>- dir1/dir2</li></ul>
Adls Gen2 End-point	Type of Microsoft Azure endpoints. You can select any of the following endpoints: <ul style="list-style-type: none"><li>- <code>core.windows.net</code>: Default</li><li>- <code>core.usgovcloudapi.net</code>: To select the Azure Government endpoints</li></ul>

For more information about creating a client ID, client secret, tenant ID, and file system name, contact the Azure administrator or see Microsoft Azure Data Lake Storage Gen2 documentation.

## Microsoft Azure SQL Data Warehouse Connection Properties

Use a Microsoft Azure SQL Data Warehouse connection to access a Microsoft Azure SQL Data Warehouse.

**Note:** The order of the connection properties might vary depending on the tool where you view them.



You can create and manage a Microsoft Azure SQL Data Warehouse connection in the Administrator tool or the Developer tool. The following table describes the Microsoft Azure SQL Data Warehouse connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ]   \ : ; " ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection.
Type	The connection type. Select Azure SQL Data Warehouse.

The following table describes the properties for metadata access:

Property	Description
Azure DW JDBC URL	Microsoft Azure Data Warehouse JDBC connection string. For example, you can enter the following connection string: <code>jdbc:sqlserver:// &lt;Server&gt;.database.windows.net:1433;database=&lt;Database&gt;</code> The Administrator can download the URL from Microsoft Azure portal.
Azure DW JDBC Username	User name to connect to the Microsoft Azure SQL Data Warehouse account. You must have permission to read, write, and truncate data in Microsoft Azure SQL Data Warehouse.
Azure DW JDBC Password	Password to connect to the Microsoft Azure SQL Data Warehouse account.
Azure DW Schema Name	Name of the schema in Microsoft Azure SQL Data Warehouse.
Azure Storage Type	Type of Azure storage to stage the files. You can select any of the following storage type: - Azure Blob. Default. To use Microsoft Azure Blob Storage to stage the files. - ADLS Gen2. To use Microsoft Azure Data Lake Storage Gen2 as storage to stage the files.
Azure Blob Account Name	Name of the Microsoft Azure Storage account to stage the files.
Azure Blob Account Key	The key that authenticates the access to the Blob storage account.
ADLS Gen2 Storage Account Name	Name of the Microsoft Azure Data Lake Storage Gen2 account to stage the files.
ADLS Gen2 Account Key	Microsoft Azure Data Lake Storage Gen2 access key to stage the files.

Property	Description
Blob End-point	<p>Type of Microsoft Azure endpoints.</p> <p>Select one of the following options:</p> <ul style="list-style-type: none"> <li>- <code>core.windows.net</code>: Default</li> <li>- <code>core.usgovcloudapi.net</code>: Select to access the US government Microsoft Azure Data Warehouse endpoints.</li> <li>- <code>core.chinacloudapi.cn</code>: Select to access a Microsoft Azure Data Warehouse endpoint in the China region.</li> </ul> <p>You can configure the US government Microsoft Azure end-points when a mapping runs in the native environment and on the Spark engine.</p>
VNet Rule	Enable to connect to a Microsoft Azure SQL Data Warehouse endpoint residing in a virtual network (VNet).

## Snowflake Connection Properties

When you set up a Snowflake connection, you must configure the connection properties.

**Note:** The order of the connection properties might vary depending on the tool where you view them.

The following table describes the Snowflake connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * ( ) - + = { [ ]   \ ; , " ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. The ID must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	Optional. The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection.
Type	The connection type. Select Snowflake.
Username	The user name to connect to the Snowflake account.
Password	The password to connect to the Snowflake account.
Account	The name of the Snowflake account.
Warehouse	The Snowflake warehouse name.

Property	Description
Role	The Snowflake role assigned to the user.
Additional JDBC URL Parameters	<p>Enter one or more JDBC connection parameters in the following format:</p> <pre>&lt;param1&gt;=&lt;value&gt;&amp;&lt;param2&gt;=&lt;value&gt;&amp;&lt;param3&gt;=&lt;value&gt;...</pre> <p>For example:</p> <pre>user=jon&amp;warehouse=mywh&amp;db=mydb&amp;schema=public</pre> <p>To access Snowflake through Okta SSO authentication, enter the web-based IdP implementing SAML 2.0 protocol in the following format:</p> <pre>authenticator=https://&lt;Your_Okta_Account_Name&gt;.okta.com</pre> <p><b>Note:</b> Microsoft ADFS is not supported.</p> <p>For more information about configuring Okta authentication, see the following website:  <a href="https://docs.snowflake.net/manuals/user-guide/admin-security-fed-auth-configure-snowflake.html#configuring-snowflake-to-use-federated-authentication">https://docs.snowflake.net/manuals/user-guide/admin-security-fed-auth-configure-snowflake.html#configuring-snowflake-to-use-federated-authentication</a></p>

## Creating a Connection to Access Sources or Targets

Create connections before you import data objects, preview data, and profile data.

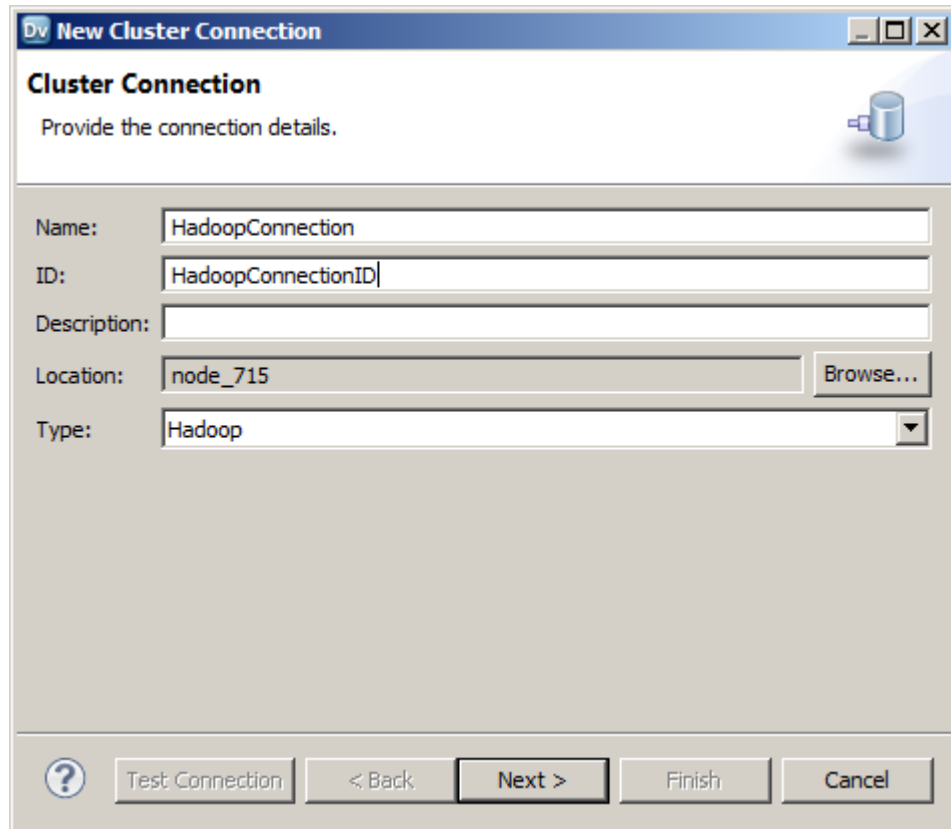
1. Within the Administrator tool click **Manage > Connections**.
2. Select **Actions > New > Connection**.
3. Select the type of connection that you want to create:
  - To select an HBase connection, select **NoSQL > HBase**.
  - To select an HDFS connection, select **File Systems > Hadoop File System**.
  - To select a Hive connection, select **Database > Hive**.
  - To select a JDBC connection, select **Database > JDBC**.
4. Click **OK**.
5. Enter a connection name, ID, and optional description.
6. Configure the connection properties. For a Hive connection, you must choose the **Access Hive as a source or target** option to use Hive as a source or a target.
7. Click **Test Connection** to verify the connection.
8. Click **Finish**.

## Creating a Hadoop Connection

Create a Hadoop connection before you run a mapping in the Hadoop environment.

1. Click **Window > Preferences**.
2. Select **Informatica > Connections**.

3. Expand the domain in the **Available Connections** list.
4. Select the **Cluster** connection type in the **Available Connections** list and click **Add**.  
The **New Cluster Connection** dialog box appears.
5. Enter the general properties for the connection.



The image shows a Windows-style dialog box titled "New Cluster Connection". The title bar includes a small "Dv" icon and standard window controls. The main area has a header "Cluster Connection" with a subtitle "Provide the connection details." and a database icon. Below this are several input fields: "Name:" with the value "HadoopConnection", "ID:" with the value "HadoopConnectionID", "Description:" (empty), "Location:" with the value "node\_715" and a "Browse..." button, and "Type:" with a dropdown menu showing "Hadoop". At the bottom, there is a row of buttons: a help button (question mark), "Test Connection", "< Back", "Next >", "Finish", and "Cancel".

6. Click **Next**.
7. Enter the Hadoop cluster properties, common properties, and the reject directory properties.
8. Click **Next**.
9. Click **Next**.  
Effective in version 10.2.2, Informatica dropped support for the Hive engine. Do not enter Hive configuration properties.
10. Enter configuration properties for the Blaze engine and click **Next**.
11. Enter configuration properties for the Spark engine and click **Finish**.

# Configuring Hadoop Connection Properties

When you create a Hadoop connection, default values are assigned to cluster environment variables, cluster path properties, and advanced properties. You can add or edit values for these properties. You can also reset to default values.

You can configure the following Hadoop connection properties based on the cluster environment and functionality that you use:

- Cluster Environment Variables
- Cluster Library Path
- Common Advanced Properties
- Blaze Engine Advanced Properties
- Spark Engine Advanced Properties

**Note:** Informatica does not recommend changing these property values before you consult with third-party documentation, Informatica documentation, or Informatica Global Customer Support. If you change a value without knowledge of the property, you might experience performance degradation or other unexpected results.

To reset to default values, delete the property values. For example, if you delete the values of an edited Cluster Library Path property, the value resets to the default \$DEFAULT\_CLUSTER\_LIBRARY\_PATH.

## Cluster Environment Variables

Cluster Environment Variables property lists the environment variables that the cluster uses. Each environment variable contains a name and a value. You can add environment variables or edit environment variables.

To edit the property in the text box, use the following format with &: to separate each name-value pair:

```
<name1>=<value1>[&:<name2>=<value2>...&:<nameN>=<valueN>]
```

Configure the following environment variables in the **Cluster Environment Variables** property:

### **HADOOP\_NODE\_JDK\_HOME**

Represents the directory from which you run the cluster services and the JDK version that the cluster nodes use. Required to run the Java transformation in the Hadoop environment and Sqoop mappings on the Blaze engine. Default is /usr/java/default. The JDK version that the Data Integration Service uses must be compatible with the JDK version on the cluster.

Set to <cluster JDK home>/jdk<version>.

For example, HADOOP\_NODE\_JDK\_HOME=<cluster JDK home>/jdk<version>.

## Cluster Library Path

Cluster Library Path property is a list of path variables for shared libraries on the cluster. You can add or edit library path variables.

To edit the property in the text box, use the following format with : to separate each path variable:

```
<variable1>[:<variable2>...:<variableN>]
```

Configure the library path variables in the **Cluster Library Path** property.

## Common Advanced Properties

Common advanced properties are a list of advanced or custom properties that are unique to the Hadoop environment. The properties are common to the Blaze and Spark engines. Each property contains a name and a value. You can add or edit advanced properties.

To edit the property in the text box, use the following format with &: to separate each name-value pair:

```
<name1>=<value1>[&:<name2>=<value2>...&:<nameN>=<valueN>]
```

Configure the following property in the **Advanced Properties** of the common properties section:

### **infapdo.java.opts**

List of Java options to customize the Java run-time environment. The property contains default values.

If mappings in a MapR environment contain a Consolidation transformation or a Match transformation, change the following value:

- -Xmx512M. Specifies the maximum size for the Java virtual memory. Default is 512 MB. Increase the value to at least 700 MB.

For example, `infapdo.java.opts=-Xmx700M`

## Blaze Engine Advanced Properties

Blaze advanced properties are a list of advanced or custom properties that are unique to the Blaze engine. Each property contains a name and a value. You can add or edit advanced properties.

To edit the property in the text box, use the following format with &: to separate each name-value pair:

```
<name1>=<value1>[&:<name2>=<value2>...&:<nameN>=<valueN>]
```

Configure the following properties in the **Advanced Properties** of the Blaze configuration section:

### **infagrid.cadi.namespace**

Namespace for the Data Integration Service to use. Required to set up multiple Blaze instances.

Set to <unique namespace>.

For example, `infagrid.cadi.namespace=TestUser1_namespace`

### **infagrid.blaze.console.jsfport**

JSF port for the Blaze engine console. Use a port number that no other cluster processes use. Required to set up multiple Blaze instances.

Set to <unique JSF port value>.

For example, `infagrid.blaze.console.jsfport=9090`

### **infagrid.blaze.console.httpport**

HTTP port for the Blaze engine console. Use a port number that no other cluster processes use. Required to set up multiple Blaze instances.

Set to <unique HTTP port value>.

For example, `infagrid.blaze.console.httpport=9091`

### **infagrid.node.local.root.log.dir**

Path for the Blaze service logs. Default is /tmp/infra/logs/blaze. Required to set up multiple Blaze instances.

Verify that all blaze users have write permission on /tmp.

Set to <local Blaze services log directory>.

For example, `infagrid.node.local.root.log.dir=<directory path>`

#### **infacal.hadoop.logs.directory**

Path in HDFS for the persistent Blaze logs. Default is `/var/log/hadoop-yarn/apps/informatica`. Required to set up multiple Blaze instances.

Set to <persistent log directory path>.

For example, `infacal.hadoop.logs.directory=<directory path>`

#### **infagrid.node.hadoop.local.root.log.dir**

Path in the Hadoop connection for the service log directory.

Set to <service log directory path>.

For example, `infagrid.node.local.root.log.dir=$HADOOP_NODE_INFA_HOME/blazeLogs`

## Spark Advanced Properties

Spark advanced properties are a list of advanced or custom properties that are unique to the Spark engine. Each property contains a name and a value. You can add or edit advanced properties. Each property contains a name and a value. You can add or edit advanced properties.

Configure the following properties in the **Advanced Properties** of the Spark configuration section:

To edit the property in the text box, use the following format with `&:` to separate each name-value pair:

```
<name1>=<value1>[&:<name2>=<value2>...&:<nameN>=<valueN>]
```

#### **infasjs.env.spark.context-settings.passthrough.spark.dynamicAllocation.executorIdleTimeout**

Maximum time that an Spark Jobserver executor node can be idle before it is removed. Increase the value to assist in debugging data preview jobs that use the Spark engine.

You can specify the time in seconds, minutes, or hours using the suffix *s*, *m*, or *h*, respectively. If you do not specify a time unit, the property uses milliseconds.

If you disable dynamic resource allocation, this property is not used.

Default is 120s.

#### **infasjs.env.spark.jobserver.max-jobs-per-context**

Maximum number of Spark jobs that can run simultaneously on a Spark context. If you increase the value of this property, you might need to allocate more resources by increasing `spark.executor.cores` and `spark.executor.memory`.

Default is 10.

#### **infasjs.env.spark.jobserver.sparkJobTimeoutInMinutes**

Maximum time in minutes that a Spark job can run on a Spark context before the Spark Jobserver cancels the job. Increase the value to assist in debugging data preview jobs that use the Spark engine.

Default is 15.

#### **infaspark.class.log.level.map**

Logging level for specific classes in the Spark driver or executor. When you configure this property, it overrides the tracing level you set for the mapping.

Set the value of this property to a JSON string in the following format: `{"<fully qualified class name>": "<log level>"}`

Join multiple class logging level statements with a comma. You can use the following logging levels: FATAL, WARN, INFO, DEBUG, ALL.

For example, set to:

```
infaspark.class.log.level.map={"org.apache.spark.deploy.yarn.ApplicationMaster":"TRACE", "org.apache.spark.deploy.security.HadoopFSDelegationTokenProvider":"DEBUG"}
```

#### **infaspark.driver.cluster.mode.extraJavaOptions**

List of extra Java options for the Spark driver that runs inside the cluster. Required for streaming mappings to read from or write to a Kafka cluster that uses Kerberos authentication.

For example, set to:

```
infaspark.driver.cluster.mode.extraJavaOptions=  
-Djava.security.egd=file:/dev/./urandom  
-XX:MaxMetaspaceSize=256M -Djavax.security.auth.useSubjectCredsOnly=true  
-Djava.security.krb5.conf=/<path to keytab file>/krb5.conf  
-Djava.security.auth.login.config=<path to jaas config>/kafka_client_jaas.config
```

To configure the property for a specific user, you can include the following lines of code:

```
infaspark.driver.cluster.mode.extraJavaOptions =  
-Djava.security.egd=file:/dev/./urandom  
-XX:MaxMetaspaceSize=256M -XX:+UseG1GC -XX:MaxGCPauseMillis=500  
-Djava.security.krb5.conf=/etc/krb5.conf
```

#### **infaspark.driver.log.level**

Logging level for the Spark driver logs. When you configure this property, it overrides the tracing level you set for the mapping.

Set the value to one of the following levels: FATAL, WARN, INFO, DEBUG, ALL.

#### **infaspark.executor.extraJavaOptions**

List of extra Java options for the Spark executor. Required for streaming mappings to read from or write to a Kafka cluster that uses Kerberos authentication.

For example, set to:

```
infaspark.executor.extraJavaOptions=  
-Djava.security.egd=file:/dev/./urandom  
-XX:MaxMetaspaceSize=256M -Djavax.security.auth.useSubjectCredsOnly=true  
-Djava.security.krb5.conf=/<path to krb5.conf file>/krb5.conf  
-Djava.security.auth.login.config=/<path to jAAS config>/kafka_client_jaas.config
```

To configure the property for a specific user, you can include the following lines of code:

```
infaspark.executor.extraJavaOptions =  
-Djava.security.egd=file:/dev/./urandom  
-XX:MaxMetaspaceSize=256M -XX:+UseG1GC -XX:MaxGCPauseMillis=500  
-Djava.security.krb5.conf=/etc/krb5.conf
```

#### **infaspark.executor.log.level**

Logging level for the Spark executor logs. When you configure this property, it overrides the tracing level you set for the mapping.

Set the value to one of the following levels: FATAL, WARN, INFO, DEBUG, ALL.

#### **infaspark.flatfile.writer.nullValue**

When the Databricks Spark engine writes to a target, it converts null values to empty strings (" "). For example, 12, AB,"",23p09udj.

The Databricks Spark engine can write the empty strings to string columns, but when it tries to write an empty string to a non-string column, the mapping fails with a type mismatch.



To allow the Databricks Spark engine to convert the empty strings back to null values and write to the target, configure the property in the Databricks Spark connection.

Set to: TRUE

#### **infaspark.json.parser.mode**

Specifies the parser how to handle corrupt JSON records. You can set the value to one of the following modes:

- DROPMALFORMED. The parser ignores all corrupted records. Default mode.
- PERMISSIVE. The parser accepts non-standard fields as nulls in corrupted records.
- FAILFAST. The parser generates an exception when it encounters a corrupted record and the Spark application goes down.

#### **infaspark.json.parser.multiLine**

Specifies whether the parser can read a multiline record in a JSON file. You can set the value to true or false. Default is false. Applies only to non-native distributions that use Spark version 2.2.x and above.

#### **infaspark.pythontx.exec**

Required to run a Python transformation on the Spark engine for Data Engineering Integration. The location of the Python executable binary on the worker nodes in the Hadoop cluster.

For example, set to:

```
infaspark.pythontx.exec=/usr/bin/python3.4
```

If you use the installation of Python on the Data Integration Service machine, set the value to the Python executable binary in the Informatica installation directory on the Data Integration Service machine.

For example, set to:

```
infaspark.pythontx.exec=INFA_HOME/services/shared/spark/python/lib/python3.4
```

#### **infaspark.pythontx.executorEnv.LD\_PRELOAD**

Required to run a Python transformation on the Spark engine for Data Engineering Streaming. The location of the Python shared library in the Python installation folder on the Data Integration Service machine.

For example, set to:

```
infaspark.pythontx.executorEnv.LD_PRELOAD=
INFA_HOME/services/shared/spark/python/lib/libpython3.6m.so
```

#### **infaspark.pythontx.executorEnv.PYTHONHOME**

Required to run a Python transformation on the Spark engine for Data Engineering Integration and Data Engineering Streaming. The location of the Python installation directory on the worker nodes in the Hadoop cluster.

For example, set to:

```
infaspark.pythontx.executorEnv.PYTHONHOME=/usr
```

If you use the installation of Python on the Data Integration Service machine, use the location of the Python installation directory on the Data Integration Service machine.

For example, set to:

```
infaspark.pythontx.executorEnv.PYTHONHOME=
INFA_HOME/services/shared/spark/python/
```

**infaspark.pythontx.submit.lib.JEP\_HOME**

Required to run a Python transformation on the Spark engine for Data Engineering Streaming. The location of the Jep package in the Python installation folder on the Data Integration Service machine.

For example, set to:

```
infaspark.pythontx.submit.lib.JEP_HOME=
INFA_HOME/services/shared/spark/python/lib/python3.6/site-packages/jep/
```

**infaspark.useHiveWarehouseAPI**

Enables the Hive Warehouse Connector. Set to TRUE.

For example, `infaspark.useHiveWarehouseAPI=true`.

**spark.authenticate**

Enables authentication for the Spark service on Hadoop. Required for Spark encryption.

Set to TRUE.

For example, `spark.authenticate=TRUE`

**spark.authenticate.enableSaslEncryption**

Enables encrypted communication when SASL authentication is enabled. Required if Spark encryption uses SASL authentication.

Set to TRUE.

For example, `spark.authenticate.enableSaslEncryption=TRUE`

**spark.datasource.hive.warehouse.load.staging.dir**

Directory for the temporary HDFS files used for batch writes to Hive. Required when you enable the Hive Warehouse Connector.

For example, set to `/tmp`

**spark.datasource.hive.warehouse.metastoreUri**

URI for the Hive metastore. Required when you enable the Hive Warehouse Connector. Use the value for *hive.metastore.uris* from the *hive\_site.xml* cluster configuration properties.

For example, set the value to `thrift://mycluster-1.com:9083` .

**spark.driver.cores**

Indicates the number of cores that each driver uses to run jobs on the Spark engine.

Set to: `spark.driver.cores=1`

**spark.driver.extraJavaOptions**

List of extra Java options for the Spark driver.

When you write date/time data within a complex data type to a Hive target using a Hortonworks HDP 3.1 cluster, append the following value to the property: `-Duser.timezone=UTC`

**spark.driver.memory**

Indicates the amount of driver process memory that the Spark engine uses to run jobs.

Recommended value: Allocate at least 256 MB for every data source.

Set to: `spark.driver.memory=3G`

**spark.executor.cores**

Indicates the number of cores that each executor process uses to run tasklets on the Spark engine.

Set to: `spark.executor.cores=1`

#### **spark.executor.extraJavaOptions**

List of extra Java options for the Spark executor.

When you write date/time data within a complex data type to a Hive target using a Hortonworks HDP 3.1 cluster, append the following value to the property: `-Duser.timezone=UTC`

#### **spark.executor.instances**

Indicates the number of instances that each executor process uses to run tasklets on the Spark engine.

Set to: `spark.executor.instances=1`

#### **spark.executor.memory**

Indicates the amount of memory that each executor process uses to run tasklets on the Spark engine.

Set to: `spark.executor.memory=3G`

#### **spark.hadoop.hive.llap.daemon.service.hosts**

Application name for the LLAP service. Required when you enable the Hive Warehouse Connector. Use the value for `hive.llap.daemon.service.hosts` from the `hive_site.xml` cluster configuration properties.

#### **spark.hadoop.hive.zookeeper.quorum**

Zookeeper hosts used by Hive LLAP. Required when you enable the Hive Warehouse Connector. Use the value for `hive.zookeeper.quorum` from the `hive_site.xml` cluster configuration properties.

#### **spark.hadoop.validateOutputSpecs**

Validates if the HBase table exists. Required for streaming mappings to write to a HBase target in an Amazon EMR cluster. Set the value to false.

#### **spark.scheduler.maxRegisteredResourcesWaitingTime**

The number of milliseconds to wait for resources to register before scheduling a task. Default is 30000. Decrease the value to reduce delays before starting the Spark job execution. Required to improve performance for mappings on the Spark engine.

Set to 15000.

For example, `spark.scheduler.maxRegisteredResourcesWaitingTime=15000`

#### **spark.scheduler.minRegisteredResourcesRatio**

The minimum ratio of registered resources to acquire before task scheduling begins. Default is 0.8. Decrease the value to reduce any delay before starting the Spark job execution. Required to improve performance for mappings on the Spark engine.

Set to: 0.5

For example, `spark.scheduler.minRegisteredResourcesRatio=0.5`

#### **spark.shuffle.encryption.enabled**

Enables encrypted communication when authentication is enabled. Required for Spark encryption.

Set to TRUE.

For example, `spark.shuffle.encryption.enabled=TRUE`

#### **spark.sql.hive.hiveserver2.jdbc.url**

URL for HiveServer2 Interactive. Required to use the Hive Warehouse Connector. Use the value in Ambari for HiveServer2 JDBC URL.

**spark.yarn.access.hadoopFileSystems**

Comma-separated list of external file systems that the Spark service can access. By default, the Spark service has access to the file systems listed in `fs.defaultFS` in the `core-site.xml` configuration set of the cluster configuration. Set this property to give the Spark service access to additional file systems.

If you run a mapping on a Cloudera CDP Public Cloud cluster and you use an HDFS on a Cloudera Data Lake cluster, you must allow access to that file system. Append the value for the property `fs.defaultFS` found in `core-site.xml` on the Data Lake cluster. For example:

```
spark.yarn.access.hadoopFileSystems=hdfs://infarndcdppamd1-master1.infarndc.src9-  
1tfl.cloudera.site:8020
```

# INDEX

## A

### ADLS

- access on Databricks [193](#)
- Gen1 Access from HDI [25](#), [26](#)
- Gen2 Access from HDI [26](#)

### Amazon AWS [209](#)

### Amazon EMR

- Hadoop administrator tasks [49](#)
- Hadoop administrator tasks [49](#)
- S3 access policies [60](#)

### Amazon Glue [55](#)

### Amazon Redshift connection

- properties [218](#)

### Amazon S3 connection

- properties [220](#)

### Analyst Service

- configuration for MapR [182](#)

### architecture

- Data Engineering Integration [15](#)
- Data Engineering with Databricks [188](#)
- Hadoop environment [15](#)

### authentication

- Databricks [201](#)

### authorization

- Apache ranger
- HDInsight [36](#)

### Azure

- configuration [212](#)

### Azure HDInsight

- Hadoop administrator tasks [74](#), [80](#)
- hosts file requirement [38](#)
- port requirements [19](#)
- storage access [25](#), [26](#)

## B

### Blaze

- Blaze user [23](#)

### Blaze engine

- create a user account [28](#)
- port requirements [19](#)
- connection properties [234](#)
- directories to create [28](#)

### blockchain

- connection properties [222](#)

## C

### Cassandra connections

- properties [224](#)

### cloud provisioning configuration

- Databricks properties [217](#)
- Amazon AWS properties [209](#)

### cloud provisioning configuration (*continued*)

- Microsoft Azure properties [212](#)

### Cloudera CDH

- Hadoop administrator tasks [98](#)

- locale [107](#)

### Cloudera CDP

- Hadoop administrator tasks [116](#)

- locale [126](#)

### cluster configuration

- create [56](#), [82](#), [103](#), [121](#), [155](#), [177](#)
- import from a cluster [82](#), [103](#), [122](#), [155](#)
- import from a file [56](#), [83](#), [104](#), [123](#), [138](#), [156](#), [177](#)

### cluster integration [14](#)

### cluster workflow

- cloud provisioning connection [208](#)

### component architecture

- clients and tools [15](#)

### components

- Informatica with Databricks environment [188](#)

### configuring

- lzo compression format [39](#)

### Confluent Kafka connection

- Confluent Kafka broker properties [226](#)

- create using infacmd [227](#)

- general properties [226](#)

### connecting to a cluster [82](#), [103](#), [122](#), [155](#)

### Connection

- details [260](#)

- properties [260](#)

### connection properties

- Databricks [227](#)

- blockchain [222](#)

### connections

- properties [208](#), [234](#)

- Google PubSub [234](#)

- HBase [208](#)

- HDFS [208](#)

- Hive [208](#)

- JDBC [208](#)

- to an SSL-enabled cluster [31](#)

### Cosmos DB connection

- creating [261](#)

### creating

- Cosmos DB connection [261](#)

### Custom Hadoop OS Path

- configuring [40](#)

## D

### data engineering

- application services [15](#)

### data engineering integration

- repositories [16](#)

### Data Engineering Integration

- integration with Informatica products [16](#)

- Data Integration Service
  - prerequisites [40](#)
  - configuration for MapR [181](#)
- Databricks
  - authentication [201](#)
  - cloud provisioning configuration [217](#)
  - components [189](#)
  - import file [204](#)
  - import from file [205](#)
  - run-time staging directory [201](#)
  - staging directory [200](#)
  - storage access [193](#)
  - import from cluster [203](#)
- Databricks connection
  - configure [206](#)
- Databricks connection properties [227](#)
- Databricks integration
  - overview [187](#)
  - system requirements [191](#)
- Delta Lake
  - configure access to [199](#)
  - connection [254](#)
- disk space
  - requirements [19](#)

## E

- Enterprise Security Package (ESP)
  - enabling [80](#)
- ephemeral clusters
  - cloud provisioning connection [208](#)

## G

- Google Analytics connections
  - properties [229](#)
- Google BigQuery connection
  - properties [230](#)
- Google Cloud Spanner connection
  - properties [232](#)
- Google Cloud Storage connections
  - properties [233](#)
- Google PubSub
  - connection properties [234](#)

## H

- Hadoop [208](#)
- Hadoop administrator
  - prerequisite tasks for Amazon EMR [49](#)
  - prerequisite tasks for Azure HDInsight [74](#)
  - prerequisite tasks for Cloudera CDH [98](#)
  - prerequisite tasks for Cloudera CDP [116](#)
  - prerequisite tasks for Hortonworks HDP [149](#)
  - prerequisite tasks for MapR [171](#)
- Hadoop administrator tasks
  - Amazon EMR [49](#)
  - Azure HDInsight [74](#)
  - Cloudera CDH [98](#)
  - Cloudera CDP [116](#)
  - configure \*-site.files [49](#), [74](#), [98](#), [116](#), [132](#), [149](#), [172](#)
  - Google Dataproc [132](#)
  - Hortonworks HDP [149](#)
  - MapR [172](#)

- Hadoop connections
  - creating [267](#)
- Hadoop operating system
  - on Data Integration Service [37](#)
- Hadoop staging user [23](#)
- HBase connections
  - MapR-DB properties [243](#)
  - properties [243](#)
- HDFS connections
  - creating [267](#)
  - properties [241](#)
- HDFS staging directory [27](#)
- HDInsight
  - security [36](#)
- Hive
  - Amazon Glue as metastore [55](#)
- Hive connections
  - creating [267](#)
  - properties [244](#)
- Hive metastore
  - Amazon Glue [55](#)
- Hive pushdown
  - connection properties [234](#)
- Hortonworks HDP
  - Hadoop administrator tasks [149](#)
- hosts file
  - Azure HDInsight [38](#)

## I

- impersonation user [23](#)
- import file
  - Databricks [204](#)
- import from cluster
  - Databricks [203](#)
- import from file
  - Databricks [205](#)
- Informatica user [23](#)
- install
  - Jep [41](#)
  - Python [41](#)
  - Python transformation [41](#)
- installation
  - MapR client [171](#)

## J

- JDBC
  - Sqoop connectivity [160](#)
  - driver for Delta Lake [199](#)
- JDBC connection
  - to Delta Lake resource [254](#)
- JDBC connections
  - connection string [249](#)
  - properties [247](#)
- JDBC V2 connection
  - properties [254](#)

## K

- Kafka connection
  - create using infacmd [259](#)
  - general properties [257](#)
  - Kafka broker properties [258](#)

Kudu connection  
properties [259](#)

## M

MapR  
Hadoop administrator tasks [171](#), [172](#)  
Analyst Service configuration [182](#)  
Data Integration Service configuration [181](#)  
Metadata Access Service configuration [182](#)  
tickets [180](#)  
MapR client  
installing [171](#)  
Messaging connection  
Confluent Kafka connection [225](#)  
Metadata Access Service  
configuration for MapR [182](#)  
Microsoft Azure [212](#)  
Microsoft Azure Data Lake Storage Gen1 connection  
properties [262](#)  
Microsoft Azure Data Lake Storage Gen2 connection  
properties [263](#)  
Microsoft Azure SQL Data Warehouse connection  
properties [264](#)

## O

operating system (OS) profile  
user [23](#)  
operating system profile  
configuration, Data Integration Service [40](#)  
pmsuid, Data Integration Service [40](#)  
overview [14](#)

## P

permissions  
Blaze engine user [28](#)  
pmsuid  
description [40](#)  
ports  
Amazon EMR requirements [19](#)  
Azure HDInsight requirements [19](#)  
Blaze engine requirements [19](#)  
Prerequisite  
download Hadoop operating system [37](#)  
prerequisites  
create directories for the Blaze engine [28](#)  
disk space [19](#)  
Hadoop administrator tasks. [49](#), [74](#), [98](#), [116](#), [149](#), [171](#)  
verify system requirements [18](#)  
Data Integration Service properties [40](#)  
uninstall [21](#)  
verify product installations [18](#)  
process [14](#)  
product installations  
prerequisites [18](#)

## R

reject file directory  
HDFS [29](#)

## S

S3 access policies [60](#)  
S3 storage  
access on Databricks [193](#)  
security  
Active Directory (AD) [30](#)  
Apache Ranger [36](#)  
Hadoop staging user [23](#)  
impersonation user [23](#)  
Informatica user [23](#)  
Kerberos [30](#)  
keytab file [30](#), [35](#)  
KMS [30](#)  
ktutil [35](#)  
OAUTH token [30](#), [34](#)  
staging user [23](#)  
security certificates [31](#), [32](#)  
Service Principal Name (SPN) [23](#)  
Snowflake connection  
properties [266](#)  
Spark deploy mode  
Hadoop connection properties [234](#)  
Spark engine  
connection properties [234](#)  
Spark Event Log directory  
Hadoop connection properties [234](#)  
Spark execution parameters  
Hadoop connection properties [234](#)  
Spark HDFS staging directory  
Hadoop connection properties [234](#)  
Sqoop  
JDBC drivers [160](#)  
Sqoop connection arguments  
-Dsquop.connection.factories [251](#)  
connect [251](#)  
direct [251](#)  
driver [251](#)  
SSL [31](#), [32](#)  
staging directory  
Databricks [200](#)  
HDFS [27](#)  
staging user [23](#)  
system requirements  
Databricks integration [191](#)  
prerequisites [18](#)

## T

TDCH connection factory  
-Dsquop.connection.factories [251](#)

## U

uninstall  
prerequisite [21](#)  
user accounts  
MapR [180](#)  
users [23](#)

## W

WASB  
access on Databricks [193](#)  
storage account access [74](#), [80](#), [191](#)

WASBS

storage account access [74](#), [80](#)