



Informatica® Big Data Management
10.2.1

Big Data Management User Guide

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

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

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.

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.

Revision: 1

Publication Date: 2019-10-23

Table of Contents

Preface	12
Informatica Resources.	12
Informatica Network.	12
Informatica Knowledge Base.	12
Informatica Documentation.	12
Informatica Product Availability Matrixes.	13
Informatica Velocity.	13
Informatica Marketplace.	13
Informatica Global Customer Support.	13
 Chapter 1: Introduction to Informatica Big Data Management.....	 14
Informatica Big Data Management Overview.	14
Example.	15
Big Data Management Tasks	15
Read from and Write to Big Data Sources and Targets.	15
Perform Data Discovery.	16
Perform Data Lineage on Big Data Sources.	17
Stream Machine Data.	17
Process Streamed Data in Real Time.	17
Manage Big Data Relationships.	18
Use a Cluster Workflow to Create Clusters on a Cloud Platform.	18
Big Data Management Component Architecture.	19
Clients and Tools.	19
Application Services.	20
Repositories.	20
Hadoop Environment.	21
Hadoop Utilities.	21
Big Data Management Engines.	22
Blaze Engine Architecture.	23
Spark Engine Architecture.	24
Hive Engine Architecture.	25
Big Data Process.	26
Step 1. Collect the Data.	26
Step 2. Cleanse the Data.	26
Step 3. Transform the Data.	27
Step 4. Process the Data.	27
Step 5. Monitor Jobs.	27
 Chapter 2: Mappings in the Hadoop Environment.....	 28
Mappings in the Hadoop Environment Overview.	28

Mapping Run-time Properties.	29
Validation Environments.	29
Execution Environment.	31
Updating Run-time Properties for Multiple Mappings.	34
Data Warehouse Optimization Mapping Example	34
Sqoop Mappings in a Hadoop Environment.	36
Sqoop Mapping-Level Arguments.	36
Configuring Sqoop Properties in the Mapping.	38
Rules and Guidelines for Mappings in a Hadoop Environment.	39
Workflows that Run Mappings in a Hadoop Environment.	40
Configuring a Mapping to Run in a Hadoop Environment.	40
Mapping Execution Plans.	41
Blaze Engine Execution Plan Details.	41
Spark Engine Execution Plan Details.	43
Hive Engine Execution Plan Details.	43
Viewing the Execution Plan for a Mapping in the Developer Tool.	44
Optimization for the Hadoop Environment.	44
Blaze Engine High Availability.	45
Enabling Data Compression on Temporary Staging Tables.	45
Parallel Sorting.	46
Truncating Partitions in a Hive Target.	47
Scheduling, Queuing, and Node Labeling.	47
Spark Engine Optimization for Sqoop Pass-Through Mappings.	50
Troubleshooting a Mapping in a Hadoop Environment.	50
High Precision Decimal Data Type on the Hive Engine.	51
Precision Loss Due to Data Overflow.	51
Precision Loss Due to Intermediate Calculations.	52
Precision and Scale in Multiplication User-Defined Functions.	52
Chapter 3: Mapping Sources in the Hadoop Environment.	54
Sources in a Hadoop Environment.	54
Complex File Sources.	55
Flat File Sources.	55
Generate the Source File Name.	56
Hive Sources.	56
PreSQL and PostgreSQL Commands.	57
Rules and Guidelines for Hive Sources on the Blaze Engine.	58
Intelligent Structure Model Sources.	59
Relational Sources.	60
Sqoop Sources.	60
Rules and Guidelines for Sqoop Sources.	60
Rules and Guidelines for Sqoop Queries.	61

Chapter 4: Mapping Targets in the Hadoop Environment.....	62
Targets in a Hadoop Environment.	62
Complex File Targets.	63
Flat File Targets.	63
HDFS Flat File Targets.	64
Hive Targets.	64
PreSQL and PostSQL Commands.	65
Truncating Hive Targets.	65
Updating Hive Targets with an Update Strategy Transformation.	66
Rules and Guidelines for Hive Targets on the Blaze Engine.	66
Relational Targets.	68
Sqoop Targets.	68
Rules and Guidelines for Sqoop Targets.	68
 Chapter 5: Mapping Transformations in the Hadoop Environment.....	 70
Overview of Mapping Transformations in the Hadoop Environment.	70
Address Validator Transformation in the Hadoop Environment.	73
Aggregator Transformation in the Hadoop Environment.	73
Aggregator Transformation Support on the Blaze Engine.	73
Aggregator Transformation Support on the Spark Engine.	74
Aggregator Transformation Support on the Hive Engine.	74
Case Converter Transformation in the Hadoop Environment.	75
Classifier Transformation in the Hadoop Environment.	75
Comparison Transformation in the Hadoop Environment.	75
Consolidation Transformation in the Hadoop Environment.	75
Data Masking Transformation in the Hadoop Environment.	75
Data Processor Transformation.	76
Data Processor Transformation Support on the Blaze Engine.	76
Data Processor Transformation Support on the Hive Engine.	76
Decision Transformation in the Hadoop Environment.	76
Expression Transformation in the Hadoop Environment.	77
Expression Transformation Support on the Blaze Engine.	77
Expression Transformation Support on the Spark Engine.	77
Expression Transformation Support on the Hive Engine.	77
Filter Transformation in the Hadoop Environment.	77
Filter Transformation Support on the Blaze Engine.	77
Java Transformation in the Hadoop Environment.	78
Java Transformation Support on the Blaze Engine.	78
Java Transformation Support on the Spark Engine.	78
Java Transformation Support on the Hive Engine.	79
Joiner Transformation in the Hadoop Environment	80
Joiner Transformation Support on the Blaze Engine.	80

Joiner Transformation Support on the Spark Engine.	80
Joiner Transformation Support on the Hive Engine.	80
Key Generator Transformation in the Hadoop Environment.	81
Labeler Transformation in the Hadoop Environment.	81
Labeler Transformation Support on the Hive Engine.	81
Lookup Transformation in the Hadoop Environment.	81
Lookup Transformation Support on the Blaze Engine.	81
Lookup Transformation Support on the Spark Engine.	81
Lookup Transformation Support on the Hive Engine.	82
Match Transformation in the Hadoop Environment.	82
Match Transformation Support on the Blaze Engine.	82
Match Transformation Support on the Spark Engine.	83
Match Transformation Support on the Hive Engine.	83
Merge Transformation in the Hadoop Environment.	83
Normalizer Transformation in the Hadoop Environment.	83
Parser Transformation in the Hadoop Environment.	83
Parser Transformation Support on the Hive Engine.	83
Python Transformation in the Hadoop Environment.	84
Python Transformation Support on the Spark Engine.	84
Rank Transformation in the Hadoop Environment.	84
Rank Transformation Support on the Blaze Engine.	84
Rank Transformation Support on the Spark Engine.	84
Rank Transformation Support on the Hive Engine.	85
Router Transformation in the Hadoop Environment.	85
Sequence Generator Transformation in the Hadoop Environment.	85
Sequence Generator Transformation Support on the Blaze Engine.	85
Sequence Generator Transformation Support on the Spark Engine.	85
Sorter Transformation in the Hadoop Environment.	86
Sorter Transformation Support on the Blaze Engine.	86
Sorter Transformation Support on the Spark Engine.	86
Sorter Transformation Support on the Hive Engine.	86
Standardizer Transformation in the Hadoop Environment.	87
Union Transformation in the Hadoop Environment.	87
Update Strategy Transformation in the Hadoop Environment.	87
Update Strategy Transformation Support on the Blaze Engine.	87
Update Strategy Transformation Support on the Spark Engine.	88
Update Strategy Transformation Support on the Hive Engine.	89
Weighted Average Transformation in the Hadoop Environment.	90
 Chapter 6: Processing Hierarchical Data on the Spark Engine.	 91
Processing Hierarchical Data on the Spark Engine Overview.	91
How to Develop a Mapping to Process Hierarchical Data.	92
Complex Data Types.	94

Array Data Type.	95
Map Data Type.	96
Struct Data Type.	97
Rules and Guidelines for Complex Data Types.	98
Complex Ports.	99
Complex Ports in Transformations.	100
Rules and Guidelines for Complex Ports.	100
Creating a Complex Port.	101
Complex Data Type Definitions.	101
Nested Data Type Definitions.	103
Rules and Guidelines for Complex Data Type Definitions.	103
Creating a Complex Data Type Definition.	103
Importing a Complex Data Type Definition.	104
Type Configuration.	106
Changing the Type Configuration for an Array Port.	106
Changing the Type Configuration for a Map Port.	108
Specifying the Type Configuration for a Struct Port.	109
Complex Operators.	110
Extracting an Array Element Using a Subscript Operator.	111
Extracting a Struct Element Using the Dot Operator.	111
Complex Functions.	112
Chapter 7: Configuring Transformations to Process Hierarchical Data.	114
Hierarchical Data Conversion.	114
Convert Relational or Hierarchical Data to Struct Data.	115
Creating a Struct Port.	115
Convert Relational or Hierarchical Data to Nested Struct Data.	117
Creating A Nested Complex Port.	118
Extract Elements from Hierarchical Data.	125
Extracting Elements from a Complex Port.	125
Flatten Hierarchical Data.	127
Flattening a Complex Port.	128
Chapter 8: Processing Unstructured and Semi-structured Data with an Intelligent Structure Model.	130
Processing Unstructured and Semi-structured Data with Intelligent Structure Model Overview.	130
Intelligent Structure Discovery Process.	131
Use Case.	131
Using an Intelligent Structure Model in a Mapping.	132
Rules and Guidelines for Intelligent Structure Models.	133
How to Develop a Mapping to Process Data with an Intelligent Structure Model	133
Mapping Example.	134
Before You Start.	135

Creating an Informatica Intelligent Cloud Services Account.	136
Creating an Intelligent Structure Model.	136
Exporting an Intelligent Structure Model.	137
Checking for Data Loss.	137

Chapter 9: Stateful Computing on the Spark Engine..... 138

Stateful Computing on the Spark Engine Overview.	138
Windowing Configuration.	139
Frame.	139
Partition and Order Keys.	140
Rules and Guidelines for Windowing Configuration.	142
Window Functions.	142
LEAD.	143
LAG.	143
Aggregate Functions as Window Functions.	143
Rules and Guidelines for Window Functions.	147
Windowing Examples.	147
Financial Plans Example.	147
GPS Pings Example.	149
Aggregate Function as Window Function Example.	151

Chapter 10: Monitoring Mappings in the Hadoop Environment..... 154

Monitoring Mappings in the Hadoop Environment Overview.	154
Hadoop Environment Logs.	155
YARN Web User Interface.	155
Accessing the Monitoring URL.	156
Viewing Hadoop Environment Logs in the Administrator Tool.	157
Monitoring a Mapping.	158
Blaze Engine Monitoring.	160
Blaze Job Monitoring Application.	161
Blaze Summary Report.	162
Blaze Engine Logs.	166
Viewing Blaze Logs.	167
Orchestrator Sunset Time.	168
Troubleshooting Blaze Monitoring.	168
Spark Engine Monitoring.	169
Viewing Hive Tasks.	172
Spark Engine Logs.	172
Viewing Spark Logs	173
Troubleshooting Spark Engine Monitoring.	173
Hive Engine Monitoring.	173
Summary Statistics.	174
Execution Statistics.	174

Monitoring with MapReduce Hive Engine.	176
Monitoring with Tez Hive Engine.	177
Hive Engine Logs.	178
Chapter 11: Mappings in the Native Environment.	181
Mappings in the Native Environment Overview.	181
Data Processor Mappings.	181
HDFS Mappings.	182
HDFS Data Extraction Mapping Example.	182
Hive Mappings.	183
Hive Mapping Example.	184
Social Media Mappings.	184
Twitter Mapping Example.	185
Chapter 12: Profiles.	186
Profiles Overview.	186
Native Environment.	186
Hadoop Environment.	187
Column Profiles for Sqoop Data Sources.	187
Creating a Single Data Object Profile in Informatica Developer.	188
Creating an Enterprise Discovery Profile in Informatica Developer.	189
Creating a Column Profile in Informatica Analyst.	190
Creating an Enterprise Discovery Profile in Informatica Analyst.	191
Creating a Scorecard in Informatica Analyst.	192
Monitoring a Profile.	193
Troubleshooting.	194
Chapter 13: Native Environment Optimization.	195
Native Environment Optimization Overview.	195
Processing Big Data on a Grid.	195
Data Integration Service Grid.	196
Grid Optimization.	196
Processing Big Data on Partitions.	196
Partitioned Model Repository Mappings.	196
Partition Optimization.	197
High Availability.	197
Chapter 14: Cluster Workflows.	199
Cluster Workflows Overview.	199
Cluster Workflow Components.	200
Create Cluster Task	200
Cloud Provisioning Configuration.	200
Hadoop Connection	200

Mapping and Other Workflow Tasks.	201
Delete Cluster Task.	201
Cluster Workflows Process	201
Administrator Tasks.	202
Create the Cluster Workflow	202
Workflow Task Run-Time Behavior.	202
Configure the Cluster Workflow	203
Amazon EMR Advanced Properties	204
Azure HDInsight Advanced Properties for the Create Cluster Task.	207
Configure the Create Cluster Task to Run Mappings on the Blaze Engine.	208
Configure the Cluster to Use an External RDS as the Hive Metastore Database.	209
Create Other Workflow Tasks.	210
Add a Delete Cluster Task.	211
Deploy and Run the Workflow.	211
Monitoring Azure HDInsight Cluster Workflow Jobs	211
Appendix A: Connections.	212
Connections.	212
Cloud Provisioning Configuration.	213
AWS Cloud Provisioning Configuration Properties.	213
Azure Cloud Provisioning Configuration Properties.	215
Hadoop Connection Properties.	218
Hadoop Cluster Properties.	218
Common Properties.	219
Reject Directory Properties.	220
Hive Pushdown Configuration.	221
Blaze Configuration.	222
Spark Configuration.	223
HDFS Connection Properties.	223
HBase Connection Properties.	225
HBase Connection Properties for MapR-DB.	226
Hive Connection Properties.	226
JDBC Connection Properties.	230
Sqoop Connection-Level Arguments.	233
Creating a Connection to Access Sources or Targets.	235
Creating a Hadoop Connection.	236
Configuring Hadoop Connection Properties.	237
Cluster Environment Variables.	237
Cluster Library Path.	239
Cluster ClassPath.	240
Cluster Executable Path.	240
Common Advanced Properties.	241
Hive Engine Advanced Properties.	241

Blaze Engine Advanced Properties.	241
Spark Engine Advanced Properties.	242
Appendix B: Data Type Reference.	245
Data Type Reference Overview.	245
Transformation Data Type Support in a Hadoop Environment.	246
Complex File and Transformation Data Types.	246
Avro and Transformation Data Types.	247
JSON and Transformation Data Types.	248
Parquet and Transformation Data Types.	248
Hive Data Types and Transformation Data Types.	250
Hive Complex Data Types.	251
Sqoop Data Types.	252
Aurora Data Types.	252
IBM DB2 and DB2 for z/OS Data Types.	252
Greenplum Data Types.	253
Microsoft SQL Server Data Types.	253
Netezza Data Types.	254
Oracle Data Types.	254
Teradata Data Types.	255
Teradata Data Types with TDCH Specialized Connectors for Sqoop.	255
Appendix C: Function Reference.	257
Function Support in the Hadoop Environment.	257
Function and Data Type Processing.	259
Rules and Guidelines for Spark Engine Processing.	259
Rules and Guidelines for Hive Engine Processing.	261
Appendix D: Parameter Reference.	263
Parameters Overview.	263
Parameter Usage.	264
Index.	266

Preface

The *Informatica Big Data Management® User Guide* provides information about configuring and running mappings in the native and Hadoop run-time environments.

Informatica Resources

Informatica Network

Informatica Network hosts Informatica Global Customer Support, the Informatica Knowledge Base, and other product resources. To access Informatica Network, visit <https://network.informatica.com>.

As a member, you can:

- Access all of your Informatica resources in one place.
- Search the Knowledge Base for product resources, including documentation, FAQs, and best practices.
- View product availability information.
- 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 search Informatica Network for product resources such as documentation, how-to articles, best practices, and PAMs.

To access the Knowledge Base, visit <https://kb.informatica.com>. If you have questions, comments, or ideas about the Knowledge Base, contact the Informatica Knowledge Base team at KB_Feedback@informatica.com.

Informatica Documentation

To get the latest documentation for your product, browse the Informatica Knowledge Base at https://kb.informatica.com/_layouts/ProductDocumentation/Page/ProductDocumentSearch.aspx.

If you have questions, comments, or ideas about this documentation, contact the Informatica Documentation team through email at infa_documentation@informatica.com.

Informatica Product Availability Matrixes

Product Availability Matrixes (PAMs) indicate the versions of operating systems, databases, and other types of data sources and targets that a product release supports. If you are an Informatica Network member, you can access 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. Developed from the real-world experience of hundreds of data management projects, Informatica Velocity represents the collective knowledge of our consultants who have worked with organizations from around the world to plan, develop, deploy, and maintain successful data management solutions.

If you are an Informatica Network member, you can access 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.

Informatica Marketplace

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

Informatica Global Customer Support

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

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

<http://www.informatica.com/us/services-and-training/support-services/global-support-centers>.

If you are an Informatica Network member, you can use Online Support at <http://network.informatica.com>.

CHAPTER 1

Introduction to Informatica Big Data Management

This chapter includes the following topics:

- [Informatica Big Data Management Overview, 14](#)
- [Big Data Management Tasks , 15](#)
- [Big Data Management Component Architecture, 19](#)
- [Big Data Management Engines, 22](#)
- [Big Data Process, 26](#)

Informatica Big Data Management Overview

Informatica Big Data Management enables your organization to process large, diverse, and fast changing data sets so you can get insights into your data. Use Big Data Management to perform big data integration and transformation without writing or maintaining Apache Hadoop code.

Use Big Data Management to collect diverse data faster, build business logic in a visual environment, and eliminate hand-coding to get insights on your data. Consider implementing a big data project in the following situations:

- The volume of the data that you want to process is greater than 10 terabytes.
- You need to analyze or capture data changes in microseconds.
- The data sources are varied and range from unstructured text to social media data.

You can identify big data sources and perform profiling to determine the quality of the data. You can build the business logic for the data and push this logic to the Hadoop cluster for faster and more efficient processing. You can view the status of the big data processing jobs and view how the big data queries are performing.

You can use multiple product tools and clients such as Informatica Developer (the Developer tool) and Informatica Administrator (the Administrator tool) to access big data functionality. Big Data Management connects to third-party applications such as the Hadoop Distributed File System (HDFS) and NoSQL databases such as HBase on a Hadoop cluster on different Hadoop distributions.

The Developer tool includes the native and Hadoop run-time environments for optimal processing. Use the native run-time environment to process data that is less than 10 terabytes. In the native environment, the Data Integration Service processes the data. The Hadoop run-time environment can optimize mapping performance and process data that is greater than 10 terabytes. In the Hadoop environment, the Data Integration Service pushes the processing to nodes in a Hadoop cluster.

When you run a mapping in the Hadoop environment, you can select to use the Spark engine, the Blaze engine, or the Hive engine to run the mapping.

Example

You are an investment banker who needs to calculate the popularity and risk of stocks and then match stocks to each customer based on the preferences of the customer. Your CIO wants to automate the process of calculating the popularity and risk of each stock, match stocks to each customer, and then send an email with a list of stock recommendations for all customers.

You consider the following requirements for your project:

- The volume of data generated by each stock is greater than 10 terabytes.
- You need to analyze the changes to the stock in microseconds.
- The stock is included in Twitter feeds and company stock trade websites, so you need to analyze these social media sources.

Based on your requirements, you work with the IT department to create mappings to determine the popularity of a stock. One mapping tracks the number of times the stock is included in Twitter feeds, and another mapping tracks the number of times customers inquire about the stock on the company stock trade website.

Big Data Management Tasks

Use Big Data Management when you want to access, analyze, prepare, transform, and stream data faster than traditional data processing environments.

You can use Big Data Management for the following tasks:

- Read from and write to diverse big data sources and targets.
- Perform data replication on a Hadoop cluster.
- Perform data discovery.
- Perform data lineage on big data sources.
- Stream machine data.
- Manage big data relationships.
- Create ephemeral clusters.

Note: The *Informatica Big Data Management User Guide* describes how to run big data mappings in the native environment or the Hadoop environment. For information on specific license and configuration requirements for a task, refer to the related product guides.

Read from and Write to Big Data Sources and Targets

In addition to relational and flat file data, you can access unstructured and semi-structured data, social media data, and data in a Hive or Hadoop Distributed File System (HDFS) environment.

You can access the following types of data:

Transaction data

You can access different types of transaction data, including data from relational database management systems, online transaction processing systems, online analytical processing systems, enterprise resource planning systems, customer relationship management systems, mainframe, and cloud.

Unstructured and semi-structured data

You can use data objects with an intelligent structure model, or Data Processor transformations, to read and transform unstructured and semi-structured data.

You can use data objects with an intelligent structure model to read and transform unstructured and semi-structured data on a Spark engine. For example, you can use a complex file data object with an intelligent structure model in a mapping to parse a Microsoft Excel file to load accounting data into S3 storage buckets. For more information, see [“Processing Unstructured and Semi-structured Data with Intelligent Structure Model Overview” on page 130](#). The intelligent structure model is quickly auto-generated from a representative file and can be easily updated or customized.

Alternatively, you can use the Data Processor transformation in a workflow to parse unstructured and semi-structured data. For example, you can parse a Microsoft Excel file to load customer and order data into relational database tables. Data Processor transformations have broad functionality and format support, but require manual setup. For more information, see the *Data Transformation User Guide*.

You can use HParser with a Data Transformation service to transform complex data into flattened, usable formats for Hive, PIG, and MapReduce processing. HParser processes complex files, such as messaging formats, HTML pages and PDF documents. HParser also transforms formats such as ACORD, HIPAA, HL7, EDI-X12, EDIFACT, AFP, and SWIFT. For more information, see the *Data Transformation HParser Operator Guide*.

Social media data

You can use PowerExchange® adapters for social media to read data from social media web sites like Facebook, Twitter, and LinkedIn. You can also use the PowerExchange for DataSift to extract real-time data from different social media web sites and capture data from DataSift regarding sentiment and language analysis. You can use PowerExchange for Web Content-Kapow to extract data from any web site.

Data in Hadoop

You can use PowerExchange adapters to read data from or write data to Hadoop. For example, you can use PowerExchange for Hive to read data from or write data to Hive. You can use PowerExchange for HDFS to extract data from and load data to HDFS. Also, you can use PowerExchange for HBase to extract data from and load data to HBase.

Data in Amazon Web Services

You can use PowerExchange adapters to read data from or write data to Amazon Web services. For example, you can use PowerExchange for Amazon Redshift to read data from or write data to Amazon Redshift. Also, you can use PowerExchange for Amazon S3 to extract data from and load data to Amazon S3.

For more information about PowerExchange adapters, see the related PowerExchange adapter guides.

Perform Data Discovery

Data discovery is the process of discovering the metadata of source systems that include content, structure, patterns, and data domains. Content refers to data values, frequencies, and data types. Structure includes candidate keys, primary keys, foreign keys, and functional dependencies. The data discovery process offers advanced profiling capabilities.

In the native environment, you can define a profile to analyze data in a single data object or across multiple data objects. In the Hadoop environment, you can push column profiles and the data domain discovery process to the Hadoop cluster.

Run a profile to evaluate the data structure and to verify that data columns contain the types of information you expect. You can drill down on data rows in profiled data. If the profile results reveal problems in the data,

you can apply rules to fix the result set. You can create scorecards to track and measure data quality before and after you apply the rules. If the external source metadata of a profile or scorecard changes, you can synchronize the changes with its data object. You can add comments to profiles so that you can track the profiling process effectively.

For more information, see the *Informatica Data Discovery Guide*.

Perform Data Lineage on Big Data Sources

Perform data lineage analysis in Enterprise Information Catalog for big data sources and targets.

Use Enterprise Information Catalog to create a Cloudera Navigator resource to extract metadata for big data sources and targets and perform data lineage analysis on the metadata. Cloudera Navigator is a data management tool for the Hadoop platform that enables users to track data access for entities and manage metadata about the entities in a Hadoop cluster.

You can create one Cloudera Navigator resource for each Hadoop cluster that is managed by Cloudera Manager. Enterprise Information Catalog extracts metadata about entities from the cluster based on the entity type.

Enterprise Information Catalog extracts metadata for the following entity types:

- HDFS files and directories
- Hive tables, query templates, and executions
- Oozie job templates and executions
- Pig tables, scripts, and script executions
- YARN job templates and executions

Note: Enterprise Information Catalog does not extract metadata for MapReduce job templates or executions.

For more information, see the *Informatica Catalog Administrator Guide*.

Stream Machine Data

You can stream machine data in real time. To stream machine data, use Informatica Edge Data Streaming.

Edge Data Streaming is a highly available, distributed, real-time application that collects and aggregates machine data. You can collect machine data from different types of sources and write to different types of targets. Edge Data Streaming consists of source services that collect data from sources and target services that aggregate and write data to a target.

For more information, see the *Informatica Vibe Data Stream for Machine Data User Guide*.

Process Streamed Data in Real Time

You can process streamed data in real time. To process streams of data in real time and uncover insights in time to meet your business needs, use Informatica Big Data Streaming.

Create Streaming mappings to collect the streamed data, build the business logic for the data, and push the logic to a Spark engine for processing. The Spark engine uses Spark Streaming to process data. The Spark engine reads the data, divides the data into micro batches and publishes it.

For more information, see the *Informatica Big Data Streaming User Guide*.

Manage Big Data Relationships

You can manage big data relationships by integrating data from different sources and indexing and linking the data in a Hadoop environment. Use Big Data Management to integrate data from different sources. Then use the MDM Big Data Relationship Manager to index and link the data in a Hadoop environment.

MDM Big Data Relationship Manager indexes and links the data based on the indexing and matching rules. You can configure rules based on which to link the input records. MDM Big Data Relationship Manager uses the rules to match the input records and then group all the matched records. MDM Big Data Relationship Manager links all the matched records and creates a cluster for each group of the matched records. You can load the indexed and matched record into a repository.

For more information, see the *MDM Big Data Relationship Management User Guide*.

Use a Cluster Workflow to Create Clusters on a Cloud Platform

You can create a workflow in the Developer tool that creates a cluster on a cloud platform and runs Mapping tasks and other tasks.

A cluster workflow contains a Create Cluster task that has configuration properties for a cluster and a reference to cloud provisioning and Hadoop connections.

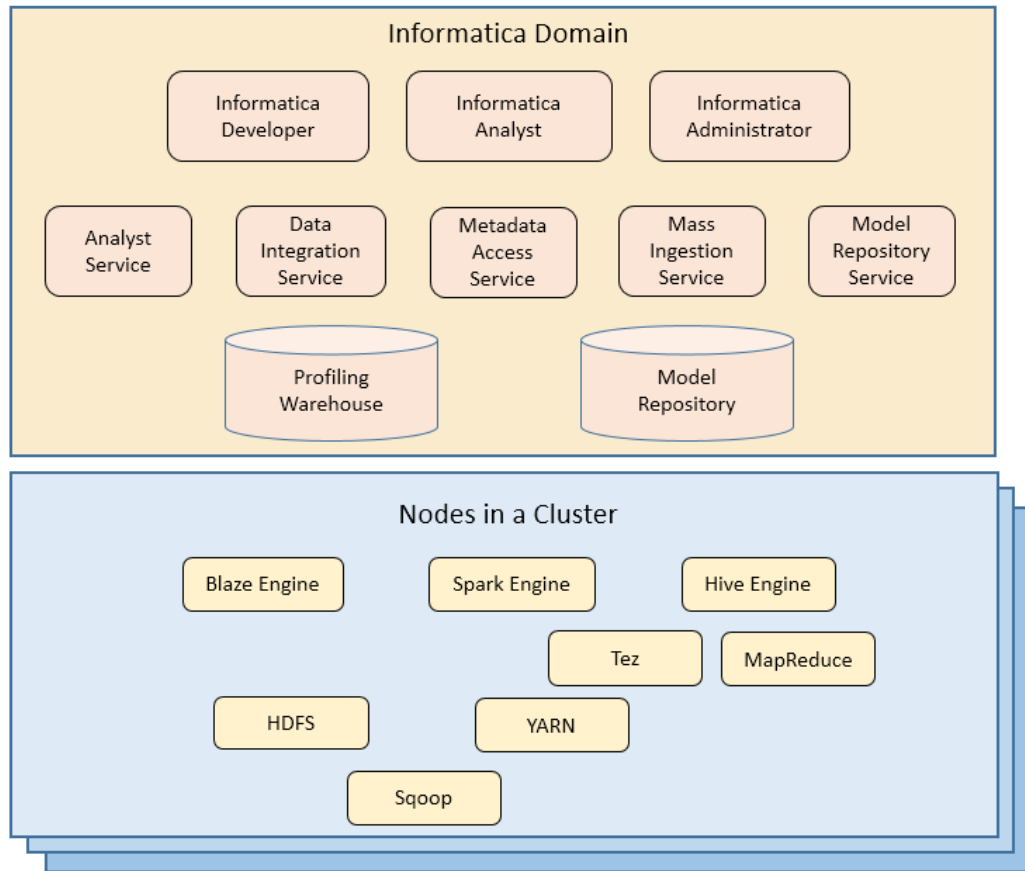
When you deploy and run a cluster workflow, it creates a cluster on a cloud platform and runs Mapping tasks and other tasks on the cluster.

You can optionally include a Delete Cluster task that terminates the cluster when workflow tasks are complete. A cluster that is created and then terminated is called an ephemeral cluster.

Big Data Management Component Architecture

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

The following image shows the components of Big Data Management:



Clients and Tools

Based on your product license, you can use multiple Informatica tools and clients to manage big data projects.

Use the following tools to manage big data projects:

Informatica Administrator

Monitor the status of profile, mapping, and MDM Big Data Relationship Management jobs on the Monitoring tab of the Administrator tool. The Monitoring tab of the Administrator tool is called the Monitoring tool. You can also design a Vibe Data Stream workflow in the Administrator tool.

Informatica Analyst

Create and run profiles on big data sources, and create mapping specifications to collaborate on projects and define business logic that populates a big data target with data.

Informatica Developer

Create and run profiles against big data sources, and run mappings and workflows on the Hadoop cluster from the Developer tool.

Application Services

Big Data Management uses application services in the Informatica domain to process data.

Big Data Management 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 the Hadoop cluster in the Hadoop 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 is a user-managed service that allows the Developer tool to access Hadoop connection information to import and preview metadata. 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. Informatica recommends to create a separate Metadata Access Service instance for each Hadoop distribution. If you use a common Metadata Access Service instance for different Hadoop distributions, you might face exceptions.

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.

Repositories

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

Big Data Management 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.

Hadoop Environment

Big Data Management 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.

Big Data Management 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 pushes mapping and profiling jobs to the Blaze, Spark, or Hive engine in the Hadoop environment.

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

Hadoop Utilities

Big Data Management uses third-party Hadoop utilities such as Sqoop to process data efficiently.

Sqoop is a Hadoop command line program to process data between relational databases and HDFS through MapReduce programs. You can use Sqoop to import and export data. When you use Sqoop, you do not need to install the relational database client and software on any node in the Hadoop cluster.

To use Sqoop, you must configure Sqoop properties in a JDBC connection and run the mapping in the Hadoop environment. You can configure Sqoop connectivity for relational data objects, customized data objects, and logical data objects that are based on a JDBC-compliant database. For example, you can configure Sqoop connectivity for the following databases:

- Aurora
- Greenplum
- IBM DB2
- IBM DB2 for z/OS
- Microsoft SQL Server
- Netezza
- Oracle
- Teradata

The Model Repository Service uses JDBC to import metadata. The Data Integration Service runs the mapping in the Hadoop run-time environment and pushes the job processing to Sqoop. Sqoop then creates map-reduce jobs in the Hadoop cluster, which perform the import and export job in parallel.

Specialized Sqoop Connectors

When you run mappings through Sqoop, you can use the following specialized connectors:

OraOop

You can use OraOop with Sqoop to optimize performance when you read data from or write data to Oracle. OraOop is a specialized Sqoop plug-in for Oracle that uses native protocols to connect to the Oracle database.

You can configure OraOop when you run Sqoop mappings on the Spark and Hive engines.

Teradata Connector for Hadoop (TDCH) Specialized Connectors for Sqoop

You can use the following TDCH specialized connectors for Sqoop to read data from or write data to Teradata:

- Cloudera Connector Powered by Teradata
- Hortonworks Connector for Teradata (powered by the Teradata Connector for Hadoop)
- MapR Connector for Teradata

These connectors are specialized Sqoop plug-ins that Cloudera, Hortonworks, and MapR provide for Teradata. They use native protocols to connect to the Teradata database.

Informatica supports Cloudera Connector Powered by Teradata and Hortonworks Connector for Teradata on the Blaze and Spark engines. When you run Sqoop mappings on the Blaze engine, you must configure these connectors. When you run Sqoop mappings on the Spark engine, the Data Integration Service invokes these connectors by default.

Informatica supports MapR Connector for Teradata on the Spark engine. When you run Sqoop mappings on the Spark engine, the Data Integration Service invokes the connector by default.

Note: For information about running native Teradata mappings with Sqoop, see the *Informatica PowerExchange for Teradata Parallel Transporter API User Guide*.

Big Data Management Engines

When you run a big data mapping, you can choose to run the mapping in the native environment or a Hadoop environment. If you run the mapping in a Hadoop environment, the mapping will run on one of the following job execution engines:

- Blaze engine
- Spark engine
- Hive engine

For more information about how Big Data Management uses each engine to run mappings, workflows, and other tasks, see the chapter about Big Data Management Engines.

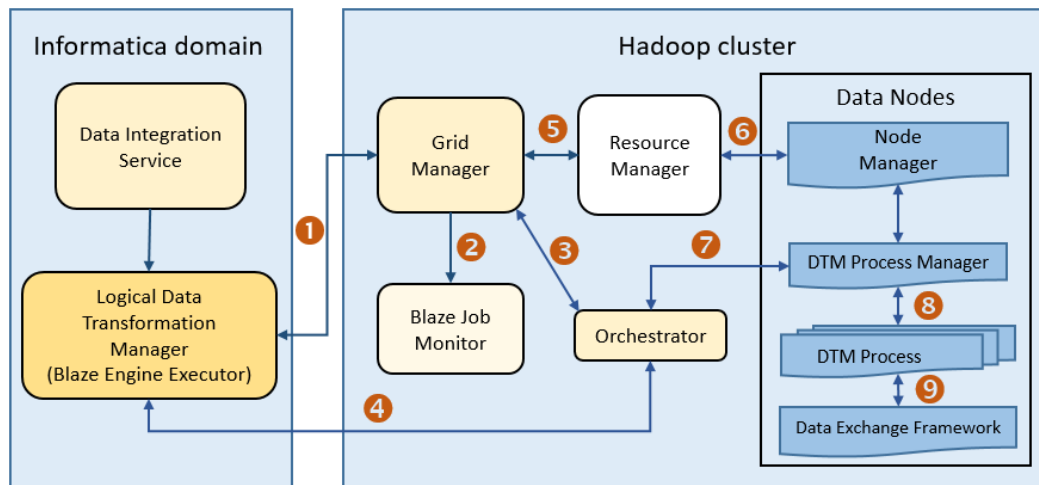
Blaze Engine Architecture

To run a mapping on the Informatica Blaze engine, the Data Integration Service submits jobs to the Blaze engine executor. The Blaze engine executor is a software component that enables communication between the Data Integration Service and the Blaze engine components on the Hadoop cluster.

The following Blaze engine components appear on the Hadoop cluster:

- Grid Manager. Manages tasks for batch processing.
- Orchestrator. Schedules and processes parallel data processing tasks on a cluster.
- Blaze Job Monitor. Monitors Blaze engine jobs on a cluster.
- DTM Process Manager. Manages the DTM Processes.
- DTM Processes. An operating system process started to run DTM instances.
- Data Exchange Framework. Shuffles data between different processes that process the data on cluster nodes.

The following image shows how a Hadoop cluster processes jobs sent from the Blaze engine executor:



The following events occur when the Data Integration Service submits jobs to the Blaze engine executor:

1. The Blaze Engine Executor communicates with the Grid Manager to initialize Blaze engine components on the Hadoop cluster, and it queries the Grid Manager for an available Orchestrator.
2. The Grid Manager starts the Blaze Job Monitor.
3. The Grid Manager starts the Orchestrator and sends Orchestrator information back to the LDTM.
4. The LDTM communicates with the Orchestrator.
5. The Grid Manager communicates with the Resource Manager for available resources for the Orchestrator.
6. The Resource Manager handles resource allocation on the data nodes through the Node Manager.
7. The Orchestrator sends the tasks to the DTM Processes through the DTM Process Manager.
8. The DTM Process Manager continually communicates with the DTM Processes.
9. The DTM Processes continually communicate with the Data Exchange Framework to send and receive data across processing units that run on the cluster nodes.

Application Timeline Server

The Hadoop Application Timeline Server collects basic information about completed application processes. The Timeline Server also provides information about completed and running YARN applications.

The Grid Manager starts the Application Timeline Server in the Yarn configuration by default.

The Blaze engine uses the Application Timeline Server to store the Blaze Job Monitor status. On Hadoop distributions where the Timeline Server is not enabled by default, the Grid Manager attempts to start the Application Timeline Server process on the current node.

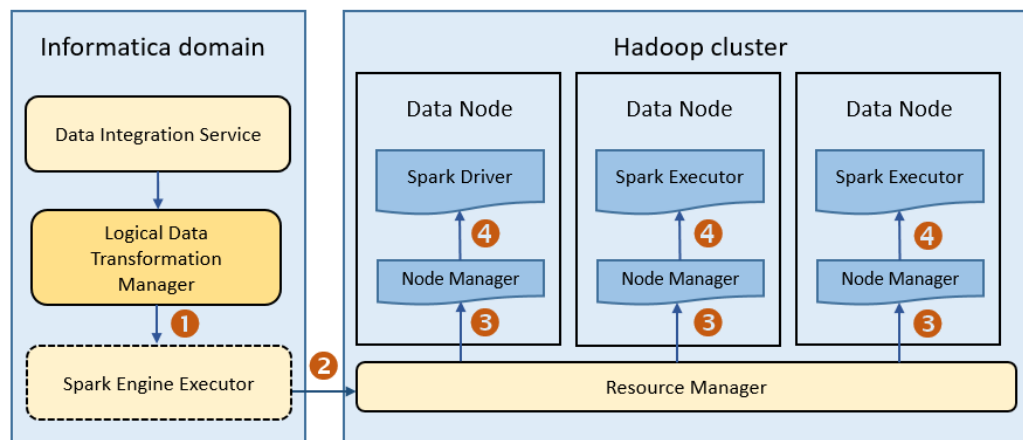
If you do not enable the Application Timeline Server on secured Kerberos clusters, the Grid Manager attempts to start the Application Timeline Server process in HTTP mode.

Spark Engine Architecture

The Data Integration Service can use the Spark engine on a Hadoop cluster to run Model repository mappings.

To run a mapping on the Spark engine, the Data Integration Service sends a mapping application to the Spark executor. The Spark executor submits the job to the Hadoop cluster to run.

The following image shows how a Hadoop cluster processes jobs sent from the Spark executor:



The following events occur when Data Integration Service runs a mapping on the Spark engine:

1. The Logical Data Transformation Manager translates the mapping into a Scala program, packages it as an application, and sends it to the Spark executor.
2. The Spark executor submits the application to the Resource Manager in the Hadoop cluster and requests resources to run the application.

Note: When you run mappings on the HDInsight cluster, the Spark executor launches a spark-submit script. The script requests resources to run the application.

3. The Resource Manager identifies the Node Managers that can provide resources, and it assigns jobs to the data nodes.
4. Driver and Executor processes are launched in data nodes where the Spark application runs.

Hive Engine Architecture

The Data Integration Service can use the Hive engine to run Model repository mappings or profiles on a Hadoop cluster.

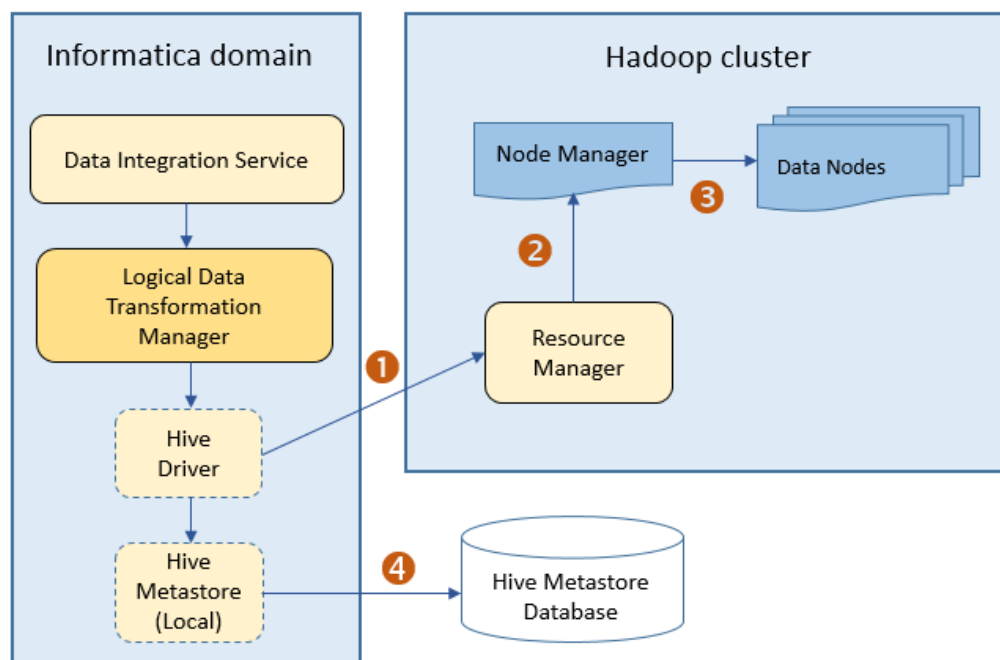
To run a mapping or profile on the Hive engine, the Data Integration Service creates HiveQL queries based on the transformation or profiling logic. The Data Integration Service submits the HiveQL queries to the Hive driver. The Hive driver converts the HiveQL queries to MapReduce or Tez jobs, and then sends the jobs to the Hadoop cluster.

Note: Effective in version 10.2.1, the MapReduce mode of the Hive run-time engine is deprecated, and Informatica will drop support for it in a future release. The Tez mode remains supported.

The Tez engine can process jobs on Hortonworks HDP, Azure HDInsight, and Amazon Elastic MapReduce. To use Cloudera CDH including Apache Hadoop, the jobs can process only on the MapReduce engine.

When you run a mapping on the Spark engine that launches Hive tasks, the mapping runs either on the MapReduce or on the Tez engines. For example, Hortonworks HDP cluster launches Hive tasks on MapReduce or Tez engines. A Cloudera CDH cluster launches Hive tasks on MapReduce engine.

The following image shows the architecture of how a Hadoop cluster processes MapReduce or Tez jobs sent from the Hive driver:



The following events occur when the Hive driver sends jobs to the Hadoop cluster:

1. The Hive driver sends the MapReduce or Tez jobs to the Resource Manager in the Hadoop cluster.
2. The Resource Manager sends the jobs request to the Node Manager that retrieves a list of data nodes that can process the MapReduce or Tez jobs.
3. The Node Manager assigns MapReduce or Tez jobs to the data nodes.
4. The Hive driver also connects to the Hive metadata database through the Hive metastore to determine where to create temporary tables. The Hive driver uses temporary tables to process the data. The Hive driver removes temporary tables after completing the task.

Big Data Process

As part of a big data project, you collect the data from diverse data sources. You can perform profiling, cleansing, and matching for the data. You build the business logic for the data and push the transformed data to the data warehouse. Then you can perform business intelligence on a view of the data.

Based on your big data project requirements, you can perform the following high-level tasks:

1. Collect the data.
2. Cleanse the data
3. Transform the data.
4. Process the data.
5. Monitor jobs.

Step 1. Collect the Data

Identify the data sources from which you need to collect the data.

Big Data Management provides several ways to access your data in and out of Hadoop based on the data types, data volumes, and data latencies in the data.

You can use PowerExchange adapters to connect to multiple big data sources. You can schedule batch loads to move data from multiple source systems to HDFS without the need to stage the data. You can move changed data from relational and mainframe systems into HDFS or the Hive warehouse. For real-time data feeds, you can move data off message queues and into HDFS.

You can collect the following types of data:

- Transactional
- Interactive
- Log file
- Sensor device
- Document and file
- Industry format

Step 2. Cleanse the Data

Cleanse the data by profiling, cleaning, and matching your data. You can view data lineage for the data.

You can perform data profiling to view missing values and descriptive statistics to identify outliers and anomalies in your data. You can view value and pattern frequencies to isolate inconsistencies or unexpected patterns in your data. You can drill down on the inconsistent data to view results across the entire data set.

You can automate the discovery of data domains and relationships between them. You can discover sensitive data such as social security numbers and credit card numbers so that you can mask the data for compliance.

After you are satisfied with the quality of your data, you can also create a business glossary from your data. You can use the Analyst tool or Developer tool to perform data profiling tasks. Use the Analyst tool to perform data discovery tasks. Use Metadata Manager to perform data lineage tasks.

Step 3. Transform the Data

You can build the business logic to parse data in the Developer tool. Eliminate the need for hand-coding the transformation logic by using pre-built Informatica transformations to transform data.

Step 4. Process the Data

Based on your business logic, you can determine the optimal run-time environment to process your data. If your data is less than 10 terabytes, consider processing your data in the native environment. If your data is greater than 10 terabytes, consider processing your data in the Hadoop environment.

Step 5. Monitor Jobs

Monitor the status of your processing jobs. You can view monitoring statistics for your processing jobs in the Monitoring tool. After your processing jobs complete you can get business intelligence and analytics from your data.

CHAPTER 2

Mappings in the Hadoop Environment

This chapter includes the following topics:

- [Mappings in the Hadoop Environment Overview, 28](#)
- [Mapping Run-time Properties, 29](#)
- [Data Warehouse Optimization Mapping Example , 34](#)
- [Sqoop Mappings in a Hadoop Environment, 36](#)
- [Rules and Guidelines for Mappings in a Hadoop Environment, 39](#)
- [Workflows that Run Mappings in a Hadoop Environment, 40](#)
- [Configuring a Mapping to Run in a Hadoop Environment, 40](#)
- [Mapping Execution Plans, 41](#)
- [Optimization for the Hadoop Environment, 44](#)
- [Troubleshooting a Mapping in a Hadoop Environment, 50](#)
- [High Precision Decimal Data Type on the Hive Engine, 51](#)

Mappings in the Hadoop Environment Overview

Configure the Hadoop run-time environment in the Developer tool to optimize mapping performance and process data that is greater than 10 terabytes. In the Hadoop environment, the Data Integration Service pushes the processing to nodes on a Hadoop cluster. When you select the Hadoop environment, you can also select the engine to push the mapping logic to the Hadoop cluster.

You can run standalone mappings, mappings that are a part of a workflow in the Hadoop environment.

Based on the mapping logic, the Hadoop environment can use the following engines to push processing to nodes on a Hadoop cluster:

- Informatica Blaze engine. An Informatica proprietary engine for distributed processing on Hadoop.
- Spark engine. A high performance engine for batch processing that can run on a Hadoop cluster or on a Spark standalone mode cluster.
- Hive engine. A batch processing engine that uses Hadoop technology such as MapReduce or Tez.

Note: Effective in version 10.2.1, the MapReduce mode of the Hive run-time engine is deprecated, and Informatica will drop support for it in a future release. The Tez mode remains supported.

When you configure the mapping, Informatica recommends that you select all engines. The Data Integration Service determines the best engine to run the mapping during validation. You can also choose to select which engine the Data Integration Service uses. You might select an engine based on whether an engine supports a particular transformation or based on the format in which the engine returns data.

When you run a mapping in the Hadoop environment, you must configure a Hadoop connection for the mapping. When you edit the Hadoop connection, you can set the run-time properties for the Hadoop environment and the properties for the engine that runs the mapping.

You can view the execution plan for a mapping to run in the Hadoop environment. View the execution plan for the engine that the Data Integration Service selects to run the mapping.

You can monitor Hive queries and the Hadoop jobs in the Monitoring tool. Monitor the jobs on a Hadoop cluster with the YARN Web User Interface or the Blaze Job Monitor web application.

The Data Integration Service logs messages from the DTM, the Blaze engine, the Spark engine, and the Hive engine in the run-time log files.

Mapping Run-time Properties

The mapping run-time properties depend on the environment that you select for the mapping.

The mapping properties contains the **Validation Environments** area and an **Execution Environment** area. The properties in the **Validation Environment** indicate whether the Developer tool validates the mapping definition for the native execution environment, the Hadoop execution environment, or both. When you run a mapping in the native environment, the Data Integration Service processes the mapping.

When you run a mapping in the Hadoop environment, the Data Integration Service pushes the mapping execution to the Hadoop cluster through a Hadoop connection. The Hadoop cluster processes the mapping.

The following image shows the mapping **Run-time** properties in a Hadoop environment:

The screenshot displays the 'Properties' window in Informatica Developer, specifically the 'Run-time' tab. The 'Validation Environments' section on the left lists five environments: Native, Hadoop, Hive, Blaze, and Spark. The 'Hadoop' and 'Spark' checkboxes are selected. The 'Execution Environment' section on the right is set to 'Hadoop'. It contains several sub-sections: 'Hadoop' with 'Connection' set to 'Auto Deploy' and 'Runtime Properties' set to 'On the Data Integration Service Machine'; 'Pushdown Configuration' with 'Pushdown Type' set to 'None' and 'Pushdown Compatibility' set to 'Rows with the same key cannot be reordered'; and 'Source Configuration' with 'Maximum Rows Read' set to 'Read All Rows', 'Maximum Runtime Interval' set to 'Run Indefinitely', and 'State Store' set to 'StateStore (Parameter)'.

Name	Value
Native	<input type="checkbox"/>
Hadoop	<input checked="" type="checkbox"/>
Hive	<input type="checkbox"/>
Blaze	<input type="checkbox"/>
Spark	<input checked="" type="checkbox"/>

Name	Value
Hadoop	
Connection	Auto Deploy
Runtime Properties	On the Data Integration Service Machine
Reject File Directory	On the Data Integration Service Machine
Pushdown Configuration	
Pushdown Type	None
Pushdown Compatibility	Rows with the same key cannot be reordered
Source Configuration	
Maximum Rows Read	Read All Rows
Maximum Runtime Interval	Run Indefinitely
State Store	StateStore (Parameter)

Validation Environments

The properties in the **Validation Environments** indicate whether the Developer tool validates the mapping definition for the native execution environment or the Hadoop execution environment.

You can configure the following properties for the **Validation Environments**:

Native

Default environment. The Data Integration Service runs the mapping in a native environment.

Hadoop

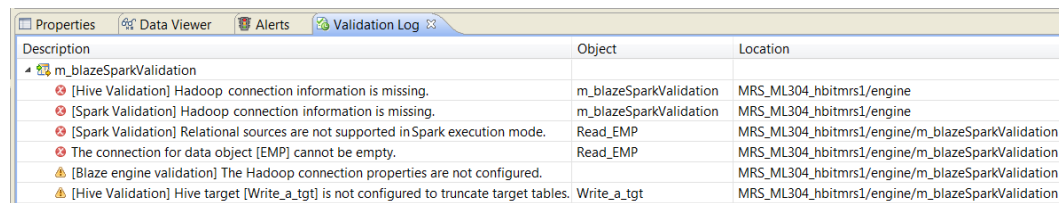
Run the mapping in the Hadoop environment. The Data Integration Service pushes the transformation logic to the Hadoop cluster through a Hadoop connection. The Hadoop cluster processes the data.

Select the engine to process the mapping. You can select the Blaze, Spark, or Hive engines.

You can use a mapping parameter to indicate the execution environment. When you select the execution environment, click **Assign Parameter**. Configure a string parameter. Set the default value to Native or Hadoop.

When you validate the mapping, validation occurs for each engine that you choose in the **Validation Environments**. The validation log might contain validation errors specific to each engine. If the mapping is valid for at least one mapping, the mapping is valid. The errors for the other engines appear in the validation log as warnings. If the mapping is valid for multiple Hadoop engines, you can view the execution plan to determine which engine will run the job. You can view the execution plan in the **Data Viewer** view.

The following image shows validation errors for the Blaze, Spark, and Hive engines:



The screenshot shows a 'Validation Log' window with a table of validation results. The table has three columns: 'Description', 'Object', and 'Location'. The 'Description' column contains various error messages, some with icons (red X for error, yellow triangle for warning). The 'Object' column lists the objects involved in the errors, and the 'Location' column shows the path to the mapping or data object.

Description	Object	Location
✖ [Hive Validation] Hadoop connection information is missing.	m_blazeSparkValidation	MRS_ML304_hbitmrs1/engine
✖ [Spark Validation] Hadoop connection information is missing.	m_blazeSparkValidation	MRS_ML304_hbitmrs1/engine
✖ [Spark Validation] Relational sources are not supported in Spark execution mode.	Read_EMP	MRS_ML304_hbitmrs1/engine/m_blazeSparkValidation
✖ The connection for data object [EMP] cannot be empty.	Read_EMP	MRS_ML304_hbitmrs1/engine/m_blazeSparkValidation
⚠ [Blaze engine validation] The Hadoop connection properties are not configured.		MRS_ML304_hbitmrs1/engine/m_blazeSparkValidation
⚠ [Hive Validation] Hive target [Write_a_tgt] is not configured to truncate target tables. Write_a_tgt	Write_a_tgt	MRS_ML304_hbitmrs1/engine/m_blazeSparkValidation

Execution Environment

Configure Hadoop properties, Pushdown Configuration properties, and Source Configuration properties in the **Execution Environment** area.

Configure the following properties in a Hadoop Execution Environment:

Name	Value
Connection	Defines the connection information that the Data Integration Service requires to push the mapping execution to the Hadoop cluster. Select the Hadoop connection to run the mapping in the Hadoop cluster. You can assign a user-defined parameter for the Hadoop connection.
Runtime Properties	<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">1. Mapping custom properties set using <code>infacmd ms runMapping</code> with the <code>-cp</code> option2. Mapping run-time properties for the Hadoop environment3. Hadoop connection advanced properties for run-time engines4. Hadoop connection advanced general properties, environment variables, and classpaths5. Data Integration Service custom properties
Reject File Directory	<p>The directory for Hadoop mapping files on HDFS when you run mappings in the Hadoop environment. The Blaze engine can write reject files to the Hadoop environment for flat file, HDFS, and Hive targets. The Spark and Hive engines can write reject files to the Hadoop environment for flat file and HDFS targets.</p> <p>Choose one of the following options:</p> <ul style="list-style-type: none">- On the Data Integration Service machine. The Data Integration Service stores the reject files based on the <code>RejectDir</code> system parameter.- On the Hadoop Cluster. The reject files are moved to the reject directory configured in the Hadoop connection. If the directory is not configured, the mapping will fail.- Defer to the Hadoop Connection. The reject files are moved based on whether the reject directory is enabled in the Hadoop connection properties. If the reject directory is enabled, the reject files are moved to the reject directory configured in the Hadoop connection. Otherwise, the Data Integration Service stores the reject files based on the <code>RejectDir</code> system parameter.

You can configure the following pushdown configuration properties:

Name	Value
Pushdown type	Choose one of the following options: <ul style="list-style-type: none">- None. Select no pushdown type for the mapping.- Source. The Data Integration Service tries to push down transformation logic to the source database.- Full. The Data Integration Service pushes the full transformation logic to the source database.
Pushdown Compatibility	Optionally, if you choose full pushdown optimization and the mapping contains an Update Strategy transformation, you can choose a pushdown compatibility option or assign a pushdown compatibility parameter. Choose one of the following options: <ul style="list-style-type: none">- Multiple rows do not have the same key. The transformation connected to the Update Strategy transformation receives multiple rows without the same key. The Data Integration Service can push the transformation logic to the target.- Multiple rows with the same key can be reordered. The target transformation connected to the Update Strategy transformation receives multiple rows with the same key that can be reordered. The Data Integration Service can push the Update Strategy transformation to the Hadoop environment.- Multiple rows with the same key cannot be reordered. The target transformation connected to the Update Strategy transformation receives multiple rows with the same key that cannot be reordered. The Data Integration Service cannot push the Update Strategy transformation to the Hadoop environment.

You can configure the following source properties:

Name	Value
Maximum Rows Read	Reserved for future use.
Maximum Runtime Interval	Reserved for future use.
State Store	Reserved for future use.

Parsing JSON Records on the Spark Engine

In the mapping run-time properties, you can configure how the Spark engine parses corrupt records and multiline records when it reads from JSON sources in a mapping.

Configure the following Hadoop run-time properties:

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 throws 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 Hadoop distributions that use Spark version 2.2.x.

Reject File Directory

You can write reject files to the Data Integration Service machine or to the Hadoop cluster. Or, you can defer to the Hadoop connection configuration. The Blaze engine can write reject files to the Hadoop environment for flat file, HDFS, and Hive targets. The Spark and Hive engines can write reject files to the Hadoop environment for flat file and HDFS targets.

If you configure the mapping run-time properties to defer to the Hadoop connection, the reject files for all mappings with this configuration are moved based on whether you choose to write reject files to Hadoop for the active Hadoop connection. You do not need to change the mapping run-time properties manually to change the reject file directory.

For example, if the reject files are currently moved to the Data Integration Service machine and you want to move them to the directory configured in the Hadoop connection, edit the Hadoop connection properties to write reject files to Hadoop. The reject files of all mappings that are configured to defer to the Hadoop connection are moved to the configured directory.

You might also want to choose to defer to the Hadoop connection when the connection is parameterized to alternate between multiple Hadoop connections. For example, the parameter might alternate between one Hadoop connection that is configured to move reject files to the Data Integration Service machine and another Hadoop connection that is configured to move reject files to the directory configured in the Hadoop connection. If you choose to defer to the Hadoop connection, the reject files are moved depending on the active Hadoop connection in the connection parameter.

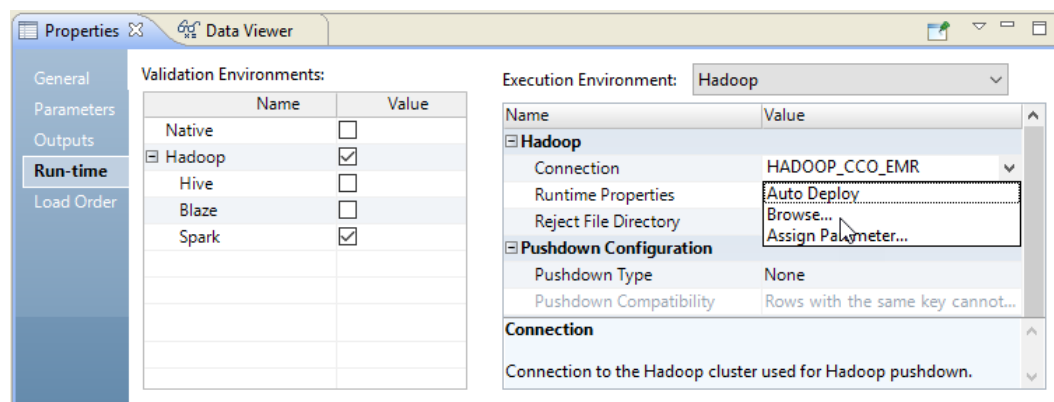
Changing the Hadoop Cluster for a Mapping Run

You can change the Hadoop cluster that a mapping runs on.

You might want to change the Hadoop cluster to run the mapping on a different Hadoop cluster type or different version.

To change the Hadoop cluster that a mapping runs on, click the connection where it appears in the mapping Run-time properties. Then click **Browse**, and select a Hadoop connection that uses a cluster configuration for the cluster you want the mapping to run on.

The following image shows the **Browse** choice for a Hadoop connection:



After you change the connection that the mapping uses, you must restart the Data Integration Service for the change to take effect.

Updating Run-time Properties for Multiple Mappings

You can enable or disable the validation environment or set the execution environment for multiple mappings. You can update multiple mappings that you run from the Developer tool or mappings that are deployed to a Data Integration Service. Use the Command Line Interface to perform these updates.

The following table describes the commands to update mapping run-time properties:

Command	Description
dis disableMappingValidationEnvironment	Disables the mapping validation environment for mappings that are deployed to the Data Integration Service.
mrs disableMappingValidationEnvironment	Disables the mapping validation environment for mappings that you run from the Developer tool.
dis enableMappingValidationEnvironment	Enables a mapping validation environment for mappings that are deployed to the Data Integration Service.
mrs enableMappingValidationEnvironment	Enables a mapping validation environment for mappings that you run from the Developer tool.
dis setMappingExecutionEnvironment	Specifies the mapping execution environment for mappings that are deployed to the Data Integration Service.
mrs setMappingExecutionEnvironment	Specifies the mapping execution environment for mappings that you run from the Developer tool.

Data Warehouse Optimization Mapping Example

You can optimize an enterprise data warehouse with the Hadoop system to store more terabytes of data cheaply in the warehouse.

For example, you need to analyze customer portfolios by processing the records that have changed in a 24-hour time period. You can offload the data on Hadoop, find the customer records that have been inserted, deleted, and updated in the last 24 hours, and then update those records in your data warehouse. You can capture these changes even if the number of columns change or if the keys change in the source files.

To capture the changes, you can create the following mappings in the Developer tool:

Mapping_Day1

Create a mapping to read customer data from flat files in a local file system and write to an HDFS target for the first 24-hour period.

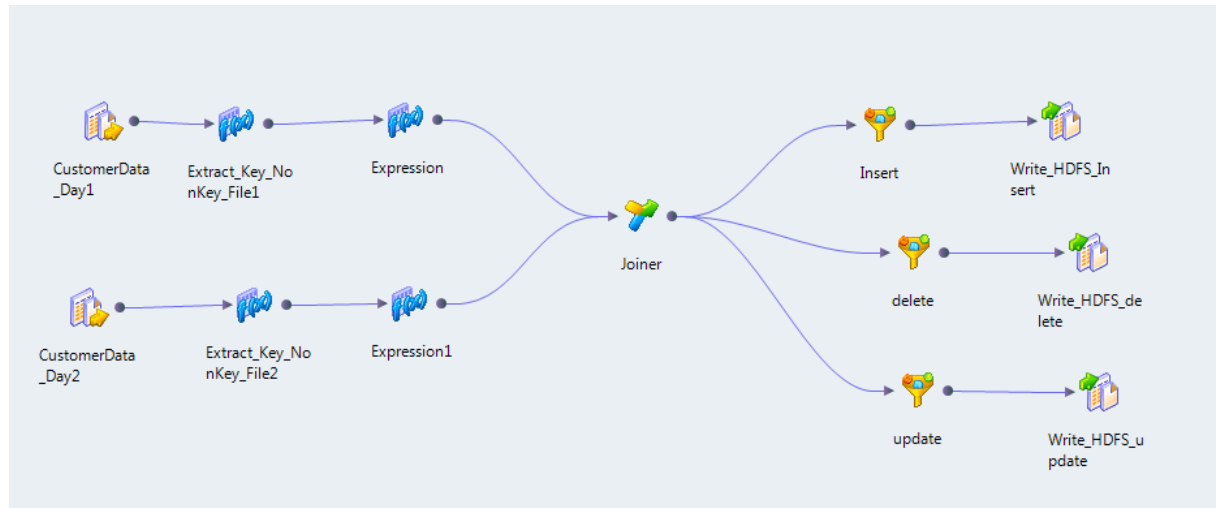
Mapping_Day2

Create a mapping to read customer data from flat files in a local file system and write to an HDFS target for the next 24-hour period.

m_CDC_DWHOptimization

Create a mapping to capture the changed data. The mapping reads data from HDFS and identifies the data that has changed. To increase performance, you configure the mapping to run on Hadoop cluster nodes in a Hadoop environment.

The following image shows the mapping m_CDC_DWHOptimization:



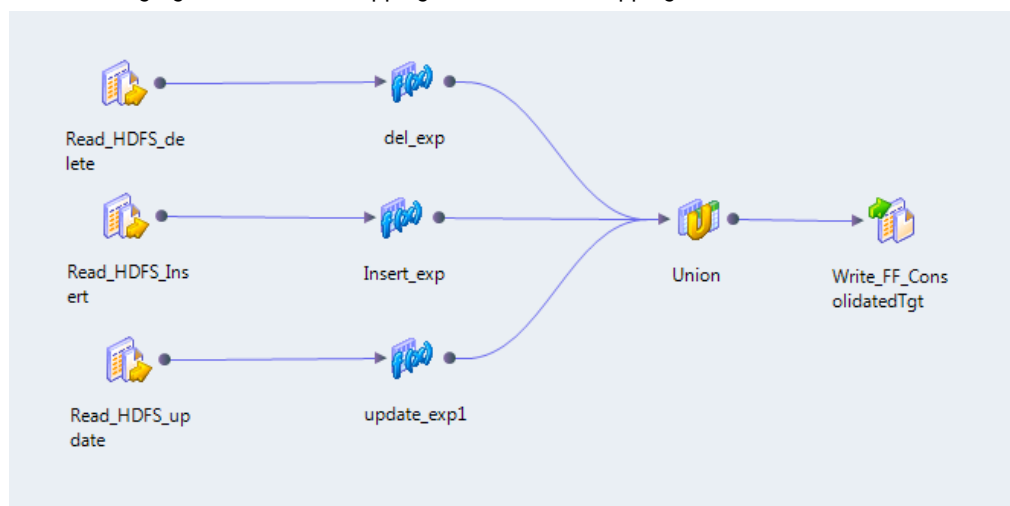
The mapping contains the following objects:

- Read transformations. Transformations that read data from HDFS files that were the targets of Mapping_Day1 and Mapping_Day2. The Data Integration Service reads all of the data as a single column.
- Expression transformations. Extract a key from the non-key values in the data. The expressions use the INSTR function and SUBSTR function to perform the extraction of key values.
- Joiner transformation. Performs a full outer join on the two sources based on the keys generated by the Expression transformations.
- Filter transformations. Use the output of the Joiner transformation to filter rows based on whether or not the rows should be updated, deleted, or inserted.
- Write transformations. Transformations that write the data to three HDFS files based on whether the data is inserted, deleted, or updated.

Consolidated_Mapping

Create a mapping to consolidate the data in the HDFS files and load the data to the data warehouse.

The following figure shows the mapping Consolidated_Mapping:

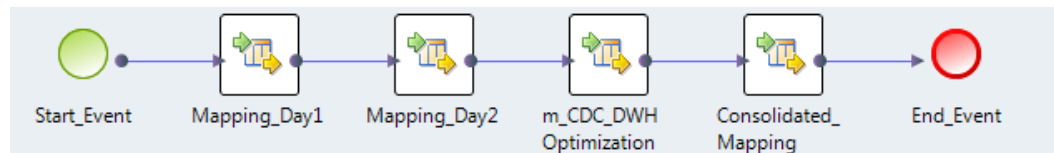


The mapping contains the following objects:

- Read transformations. Transformations that read data from HDFS files that were the target of the previous mapping are the sources of this mapping.
- Expression transformations. Add the deleted, updated, or inserted tags to the data rows.
- Union transformation. Combines the records.
- Write transformation. Transformation that writes data to the flat file that acts as a staging location on the local file system.

You can open each mapping and right-click to run the mapping. To run all mappings in sequence, use a workflow.

The following image shows the example Data Warehouse Optimization workflow:



To run the workflow, use the `infacmd wfs startWorkflow` command.

Sqoop Mappings in a Hadoop Environment

After you copy the Type 4 JDBC driver .jar files required for Sqoop connectivity to the `externaljdbcjars` directory, enable Sqoop in a JDBC connection, and import a Sqoop source or Sqoop target, you can create a mapping. You can then run the Sqoop mapping in the Hadoop run-time environment with a Hadoop connection. You can run Sqoop mappings on the Blaze, Spark, and Hive engines.

If you use Cloudera Connector Powered by Teradata or Hortonworks Connector for Teradata, you can run mappings on the Blaze or Spark engines. If you use MapR Connector for Teradata, you can run mappings on the Spark engine.

In the mapping, you can specify additional Sqoop arguments and disable the Sqoop connector.

Note: If you add or delete a Type 4 JDBC driver .jar file required for Sqoop connectivity from the `externaljdbcjars` directory, changes take effect after you restart the Data Integration Service. If you run the mapping on the Blaze engine, changes take effect after you restart the Data Integration Service and Blaze Grid Manager. When you run the mapping for the first time, you do not need to restart the Data Integration Service and Blaze Grid Manager. You need to restart the Data Integration Service and Blaze Grid Manager only for the subsequent mapping runs.

Sqoop Mapping-Level Arguments

If a data object uses Sqoop, you can click the corresponding **Read** transformation or **Write** transformation in the Sqoop mapping to define the arguments that Sqoop must use to process the data. The Data Integration Service merges the additional Sqoop arguments that you specify in the mapping with the arguments that you specified in the JDBC connection and constructs the Sqoop command.

The Sqoop arguments that you specify in the mapping take precedence over the arguments that you specified in the JDBC connection. However, if you do not enable the Sqoop connector in the JDBC connection but enable the Sqoop connector in the mapping, the Data Integration Service does not run the mapping through Sqoop. The Data Integration Service runs the mapping through JDBC.

You can configure the following Sqoop arguments in a Sqoop mapping:

- m or num-mappers
- split-by
- batch
- infaoptimize

For a complete list of the Sqoop arguments that you can configure, see the Sqoop documentation.

m or num-mappers

The m or num-mappers argument defines the number of map tasks that Sqoop must use to import and export data in parallel.

Use the following syntax:

```
-m <number of map tasks>
--num-mappers <number of map tasks>
```

If you configure the m argument or num-mappers argument, you must also configure the split-by argument to specify the column based on which Sqoop must split the work units.

Use the m argument or num-mappers argument to increase the degree of parallelism. You might have to test different values for optimal performance.

When you configure the m argument or num-mappers argument and run Sqoop mappings on the Spark or Blaze engines, Sqoop dynamically creates partitions based on the file size.

Note: If you configure the num-mappers argument to export data on the Blaze or Spark engine, Sqoop ignores the argument. Sqoop creates map tasks based on the number of intermediate files that the Blaze or Spark engine creates.

split-by

The split-by argument defines the column based on which Sqoop splits work units.

Use the following syntax:

```
--split-by <column_name>
```

You can configure the split-by argument to improve the performance. If the primary key does not have an even distribution of values between the minimum and maximum range, you can configure the split-by argument to specify another column that has a balanced distribution of data to split the work units.

If you do not define the split-by column, Sqoop splits work units based on the following criteria:

- If the data object contains a single primary key, Sqoop uses the primary key as the split-by column.
- If the data object contains a composite primary key, Sqoop defaults to the behavior of handling composite primary keys without the split-by argument. See the Sqoop documentation for more information.
- If the data object does not contain a primary key, the value of the m argument and num-mappers argument default to 1.

Rules and Guidelines for the split-by Argument

Consider the following restrictions when you configure the split-by argument:

- If you configure the split-by argument and the split-by column contains NULL values, Sqoop does not import the rows that contain NULL values. However, the mapping runs successfully and no error is written in the YARN log.

- If you configure the split-by argument and the split-by column contains special characters, the Sqoop import process fails.
- The split-by argument is required in the following scenarios:
 - You use Cloudera Connector Powered by Teradata or Hortonworks Connector for Teradata, and the Teradata table does not contain a primary key.
 - You create a custom query to override the default query when you import data from a Sqoop source.

batch

The batch argument indicates that Sqoop must export data in batches.

Use the following syntax:

```
--batch
```

You can configure the batch argument to improve the performance.

infaoptimize

The infaoptimize argument defines whether you want 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
```

Configuring Sqoop Properties in the Mapping

You can specify additional Sqoop arguments and disable the Sqoop connector at the mapping level. The Sqoop arguments that you specify at the mapping level take precedence over the arguments that you specified in the JDBC connection.

1. Open the mapping that contains the data object for which you want to configure Sqoop properties.
2. Select the Read or Write transformation that is associated with the data object.
3. Click the **Advanced** tab.
4. To disable the Sqoop connector for the data object, select the **Disable Sqoop Connector** check box.
5. Perform one of the following steps:
 - To specify additional Sqoop import arguments for the data object, enter the import arguments in the **Additional Sqoop Import Arguments** text box.
 - To specify additional Sqoop export arguments for the data object, enter the export arguments in the **Additional Sqoop Export Arguments** text box.

The Data Integration Service merges the additional Sqoop arguments that you specified in the mapping with the arguments that you specified in the JDBC connection and constructs the Sqoop command. The Data Integration Service then invokes Sqoop on a Hadoop node.

Rules and Guidelines for Mappings in a Hadoop Environment

You can run mappings in a Hadoop environment. When you run mappings in a Hadoop environment, some differences in processing and configuration apply.

The following processing differences apply to mappings in a Hadoop environment:

- A mapping is run in high precision mode in a Hadoop environment for Hive 0.11 and above.
- In a Hadoop environment, sources that have data errors in a column result in a null value for the column. In the native environment, the Data Integration Service does not process the rows that have data errors in a column.
- When you cancel a mapping that reads from a flat file source, the file copy process that copies flat file data to HDFS may continue to run. The Data Integration Service logs the command to kill this process in the Hive session log, and cleans up any data copied to HDFS. Optionally, you can run the command to kill the file copy process.
- When you set a limit on the number of rows read from the source for a Blaze mapping, the Data Integration Service runs the mapping with the Hive engine instead of the Blaze engine.

The following configuration differences apply to mappings in a Hadoop environment:

- Set the optimizer level to none or minimal if a mapping validates but fails to run. If you set the optimizer level to use cost-based or semi-join optimization methods, the Data Integration Service ignores this at run-time and uses the default.
- The Spark engine does not honor the early projection optimization method in all cases. If the Data Integration Service removes the links between unused ports, the Spark engine might reconnect the ports.

When the Spark engine runs a mapping, it processes jobs on the cluster using HiveServer2 in the following cases:

- The mapping writes to a target that is a Hive table bucketed on fields of type char or varchar.
- The mapping reads from or writes to Hive transaction-enabled tables.
- The mapping reads from or writes to Hive tables where column-level security is enabled.
- The mapping writes to a Hive target and is configured to create or replace the table at run time.

Workflows that Run Mappings in a Hadoop Environment

You can add a mapping that you configured to run in a Hadoop environment to a Mapping task in a workflow. When you deploy and run the workflow, the Mapping task runs the mapping.

You might decide to run a mapping from a workflow so that you can make decisions during the workflow run. You can configure a workflow to run multiple mappings in sequence or in parallel. You can configure a workflow to send emails that notify users about the status of the Mapping tasks.

When a Mapping task runs a mapping configured to run in a Hadoop environment, do not assign the Mapping task outputs to workflow variables. Mappings that run in a Hadoop environment do not provide the total number of target, source, and error rows. When a Mapping task includes a mapping that runs in a Hadoop environment, the task outputs contain a value of zero (0).

Configuring a Mapping to Run in a Hadoop Environment

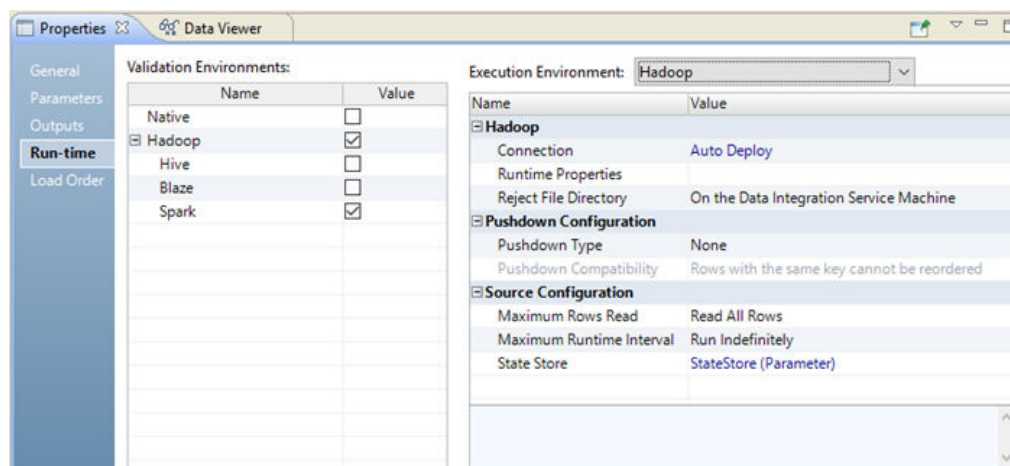
You can configure a mapping to run in a Hadoop environment. To configure a mapping, you must select a validation environment and an execution environment.

Note: Hive engine is deprecated and Informatica plans to drop support in a future release.

1. Select a mapping from a project or folder from the **Object Explorer** view to open in the editor.
2. In the **Properties** view, select the **Run-time** tab.
3. Select **Hadoop** as the value for the validation environment.

By default, the Blaze and Spark engines are selected. Disable the engines that you do not want to use.

4. Select **Hadoop** for the execution environment.



5. In the Hadoop environment, select **Connection** and use the drop down in the value field to browse for a connection or create a connection parameter:
 - To select a connection, click **Browse** and select a connection.

- To create a connection parameter, click **Assign Parameter**.
6. Configure the rest of the properties for the Hadoop execution environment.
 7. Right-click an empty area in the editor and click **Validate**.
The Developer tool validates the mapping.
 8. View validation errors on the **Validation Log** tab.
 9. Click the **Data Viewer** view.
 10. Click **Show Execution Plan** to view the execution plan for the mapping.

Mapping Execution Plans

The Data Integration Service generates an execution plan to run mappings on a Blaze, Spark, or Hive engine. The Data Integration Service translates the mapping logic into code that the run-time engine can execute. You can view the plan in the Developer tool before you run the mapping and in the Administrator tool after you run the mapping.

The Data Integration Service generates mapping execution plans to run on the following engines:

Informatica Blaze engine

The Blaze engine execution plan simplifies the mapping into segments. It contains tasks to start the mapping, run the mapping, and clean up the temporary tables and files. It contains multiple tasklets and the task recovery strategy. It also contains pre- and post-grid task preparation commands for each mapping before running the main mapping on a Hadoop cluster. A pre-grid task can include a task such as copying data to HDFS. A post-grid task can include tasks such as cleaning up temporary files or copying data from HDFS.

Spark engine

The Spark execution plan shows the run-time Scala code that runs the mapping logic. A translation engine translates the mapping into an internal representation of the logic. The internal representation is rendered into Scala code that accesses the Spark API. You can view the Scala code in the execution plan to debug the logic.

Hive engine

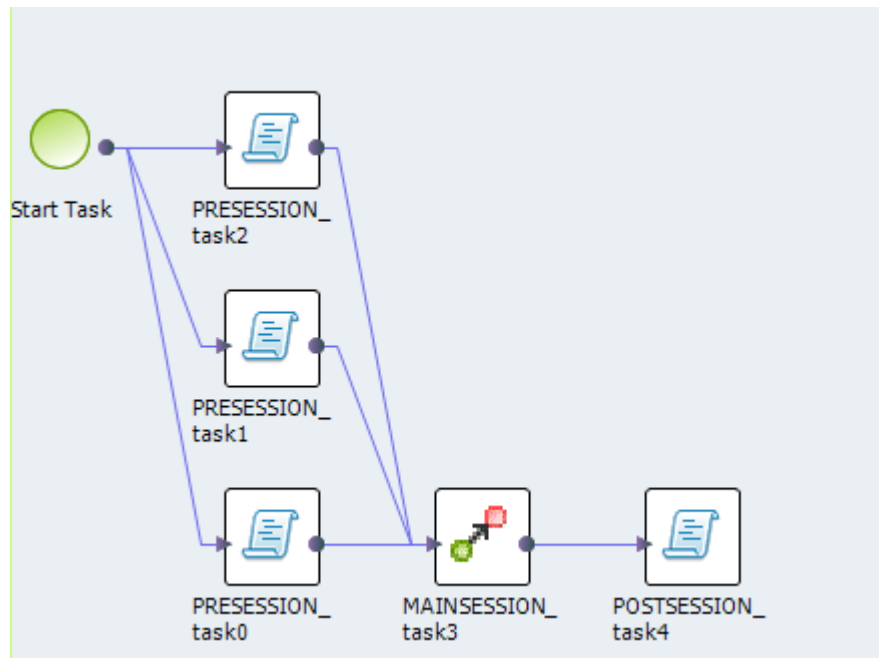
The Hive execution plan is a series of Hive queries. The plan contains tasks to start the mapping, run the mapping, and clean up the temporary tables and files. You can view the Hive execution plan that the Data Integration Service generates before you run the mapping. When the Data Integration Service pushes the mapping to the Hive engine, it has a Hive executor that can process the mapping. The Hive executor simplifies the mapping to an equivalent mapping with a reduced set of instructions and generates a Hive execution plan.

Blaze Engine Execution Plan Details

You can view details of the Blaze engine execution plan in the Administrator tool and Developer tool.

In the Developer tool, the Blaze engine execution plan appears as a workflow. You can click on each component in the workflow to get the details.

The following image shows the Blaze execution plan in the Developer tool:



The Blaze engine execution plan workflow contains the following components:

- Start task. The workflow start task.
- Command task. The pre-processing or post-processing task for local data.
- Grid mapping. An Informatica mapping that the Blaze engine compiles and distributes across a cluster of nodes.
- Grid task. A parallel processing job request sent by the Blaze engine executor to the Grid Manager.
- Grid segment. Segment of a grid mapping that is contained in a grid task.
- Tasklet. A partition of a grid segment that runs on a separate DTM.

In the Administrator tool, the Blaze engine execution plan appears as a script.

The following image shows the Blaze execution script:

Test - XraKb1DXEeWg		Properties	Blaze Execution Plan	Summary
Script Id	Script			
MAINSESSION_task3	<p>Execution scriptStep [MAINSESSION_task3], type [GridTaskStepImpl]. With "from" step(s): PRESESSION_task0, PRESESSION_task1, PRESESSION_task2. With "to" step(s): POSTSESSION_task4.</p> <p>Grid mapping task has totally [3] substeps:</p> <p>Execution step [submapping-2], type [SegmentStepImpl]. With no "from" step. With "to" step(s): submapping-3.</p> <p>Included instances: Read_IN_OUT[SourceTx], DETarget_Joiner_G1[TargetTx],</p> <p>Execution step [submapping-1], type [SegmentStepImpl]. With no "from" step. With "to" step(s): submapping-3.</p> <p>Included instances: DETarget_Joiner_G0[TargetTx], Read_IN_OUT1[SourceTx],</p> <p>Execution step [submapping-3], type [SegmentStepImpl]. With "from" step(s): submapping-1, submapping-2. With no "to" step.</p> <p>Included instances: Write_IN_OUT[TargetTx], DESource_Joiner_G1[SourceTx], Joiner[JoinerTx],</p> <p>DESource_Joiner_G0[SourceTx],</p>			

In the Administrator tool, the Blaze engine execution plan has the following details:

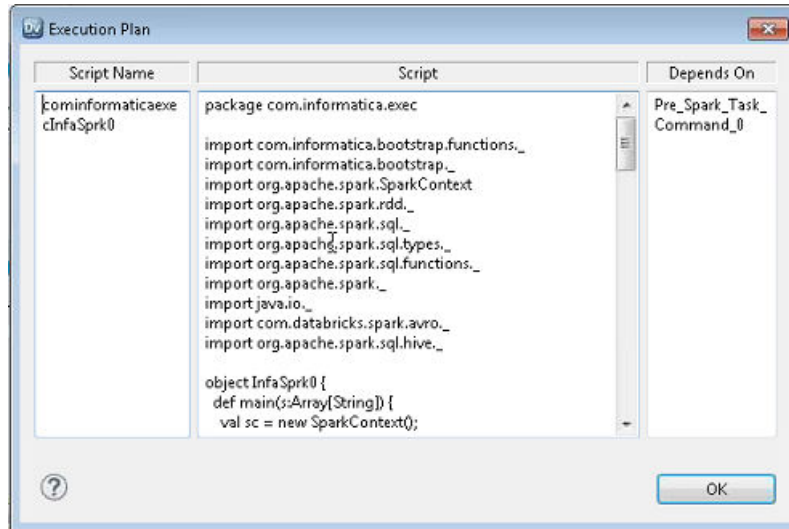
- Script ID. Unique identifier for the Blaze engine script.
- Script. Blaze engine script that the Data Integration Service generates based on the mapping logic.
- Depends on. Tasks that the script depends on. Tasks include other scripts and Data Integration Service tasks, like the Start task.

Spark Engine Execution Plan Details

You can view the details of a Spark engine execution plan from the Administrator tool or Developer tool.

The Spark engine execution plan shows the Scala code to run in the Hadoop cluster.

The following image shows the execution plan for a mapping to run on the Spark engine:



The Spark engine execution plan has the following details:

- Script ID. Unique identifier for the Spark engine script.
- Script. Scala code that the Data Integration Service generates based on the mapping logic.
- Depends on. Tasks that the script depends on. Tasks include other scripts and Data Integration Service tasks.

Hive Engine Execution Plan Details

You can view the details of a Hive engine execution plan for a mapping from the Administrator tool or Developer tool.

The following table describes the properties of a Hive engine execution plan:

Property	Description
Script Name	Name of the Hive script.
Script	Hive script that the Data Integration Service generates based on the mapping logic.
Depends On	Tasks that the script depends on. Tasks include other scripts and Data Integration Service tasks, like the Start task.

Note: Effective in version 10.2.1, the MapReduce mode of the Hive run-time engine is deprecated, and Informatica will drop support for it in a future release. The Tez mode remains supported.

When you choose to run a mapping in the Hadoop environment, the Blaze and Spark run-time engines are selected by default.

Previously, the Hive run-time engine was also selected.

If you select Hive to run a mapping, the Data Integration Service will use Tez. You can use the Tez engine only on the following Hadoop distributions:

- Amazon EMR
- Azure HDInsight
- Hortonworks HDP

In a future release, when Informatica drops support for MapReduce, the Data Integration Service will ignore the Hive engine selection and run the mapping on Blaze or Spark.

Viewing the Execution Plan for a Mapping in the Developer Tool

You can view the Hive or Blaze engine execution plan for a mapping that runs in a Hadoop environment. You do not have to run the mapping to view the execution plan in the Developer tool.

Note: You can also view the execution plan in the Administrator tool.

1. In the Developer tool, open the mapping.
2. Select the **Data Viewer** view.
3. Select **Show Execution Plan**.

The **Data Viewer** view shows the details for the execution plan.

Optimization for the Hadoop Environment

Optimize the Hadoop environment to increase performance.

You can optimize the Hadoop environment in the following ways:

Configure a highly available Hadoop cluster.

You can configure the Data Integration Service and the Developer tool to read from and write to a highly available Hadoop cluster. The steps to configure a highly available Hadoop cluster depend on the type of Hadoop distribution. For more information about configuration steps for a Hadoop distribution, see the *Informatica Big Data Management Hadoop Integration Guide*.

Compress data on temporary staging tables.

You can enable data compression on temporary staging tables to increase mapping performance.

Run mappings on the Blaze engine.

Run mappings on the highly available Blaze engine. The Blaze engine enables restart and recovery of grid tasks and tasklets by default.

Perform parallel sorts.

When you use a Sorter transformation in a mapping, the Data Integration Service enables parallel sorting by default when it pushes the mapping logic to the Hadoop cluster. Parallel sorting improves mapping performance.

Partition Joiner transformations.

When you use a Joiner transformation in a Blaze engine mapping, the Data Integration Service can apply map-side join optimization to improve mapping performance. The Data Integration Service applies map-side join optimization if the master table is smaller than the detail table. When the Data Integration Service applies map-side join optimization, it moves the data to the Joiner transformation without the cost of shuffling the data.

Truncate partitions in a Hive target.

You can truncate partitions in a Hive target to increase performance. To truncate partitions in a Hive target, you must choose to both truncate the partition in the Hive target and truncate the target table.

Assign resources on Hadoop clusters.

You can use schedulers to assign resources on a Hadoop cluster. You can use a capacity scheduler or a fair scheduler depending on the needs of your organization.

Configure YARN queues to share resources on Hadoop clusters.

You can configure YARN queues to redirect jobs on the Hadoop cluster to specific queues. The queue where a job is assigned defines the resources that are allocated to perform the job.

Label nodes in a Hadoop cluster.

You can label nodes in a Hadoop cluster to divide the cluster into partitions that have specific characteristics.

Optimize Sqoop mappings on the Spark engine.

The Data Integration Service can optimize the performance of Sqoop pass-through mappings that run on the Spark engine.

Blaze Engine High Availability

The Blaze engine is a highly available engine that determines the best possible recovery strategy for grid tasks and tasklets.

Based on the size of the grid task, the Blaze engine attempts to apply the following recovery strategy:

- No high availability. The Blaze engine does not apply a recovery strategy.
- Full restart. Restarts the grid task.

Enabling Data Compression on Temporary Staging Tables

To optimize performance when you run a mapping in the Hadoop environment, you can enable data compression on temporary staging tables. When you enable data compression on temporary staging tables, mapping performance might increase.

To enable data compression on temporary staging tables, complete the following steps:

1. Configure the Hive connection to use the codec class name that the Hadoop cluster uses to enable compression on temporary staging tables.
2. Configure the Hadoop cluster to enable compression on temporary staging tables.

Hadoop provides following compression libraries for the following compression codec class names:

Compression Library	Codec Class Name	Performance Recommendation
Zlib	org.apache.hadoop.io.compress.DefaultCodec	n/a
Gzip	org.apache.hadoop.io.compress.GzipCodec	n/a
Snappy	org.apache.hadoop.io.compress.SnappyCodec	Recommended for best performance.

Compression Library	Codec Class Name	Performance Recommendation
Bz2	org.apache.hadoop.io.compress.BZip2Codec	Not recommended. Degrades performance.
LZO	com.hadoop.compression.lzo.LzoCodec	n/a

Step 1. Configure the Hive Connection to Enable Data Compression on Temporary Staging Tables

Use the Administrator tool or the Developer tool to configure the Hive connection. You can edit the Hive connection properties to configure the codec class name that enables data compression on temporary staging tables.

1. In the Hive connection properties, edit the properties to run mappings in a Hadoop cluster.
2. Select **Temporary Table Compression Codec**.
3. Choose to select a predefined codec class name or enter a custom codec class name.
 - To select a predefined codec class name, select a compression library from the list.
 - To enter a custom codec class name, select custom from the list and enter the codec class name that matches the codec class name in the Hadoop cluster.

Step 2. Configure the Hadoop Cluster to Enable Compression on Temporary Staging Tables

To enable compression on temporary staging tables, you must install a compression codec on the Hadoop cluster.

For more information about how to install a compression codec, refer to the Apache Hadoop or Hive documentation.

1. Verify that the native libraries for the compression codec class name are installed on every node on the cluster.
2. To include the compression codec class name that you want to use, update the property `io.compression.codecs` in `core-site.xml`. The value for this property is a comma separated list of all the codec class names supported on the cluster.
3. Verify that the Hadoop-native libraries for the compression codec class name that you want to use are installed on every node on the cluster.
4. Verify that the `LD_LIBRARY_PATH` variable on the Hadoop cluster includes the locations of both the native and Hadoop-native libraries where you installed the compression codec.

Parallel Sorting

To improve mapping performance, the Data Integration Service enables parallel sorting by default in a mapping that has a Sorter transformation and a flat file target.

The Data Integration Service enables parallel sorting for mappings in a Hadoop environment based on the following rules and guidelines:

- The mapping does not include another transformation between the Sorter transformation and the target.
- The data type of the sort keys does not change between the Sorter transformation and the target.
- Each sort key in the Sorter transformation must be linked to a column in the target.

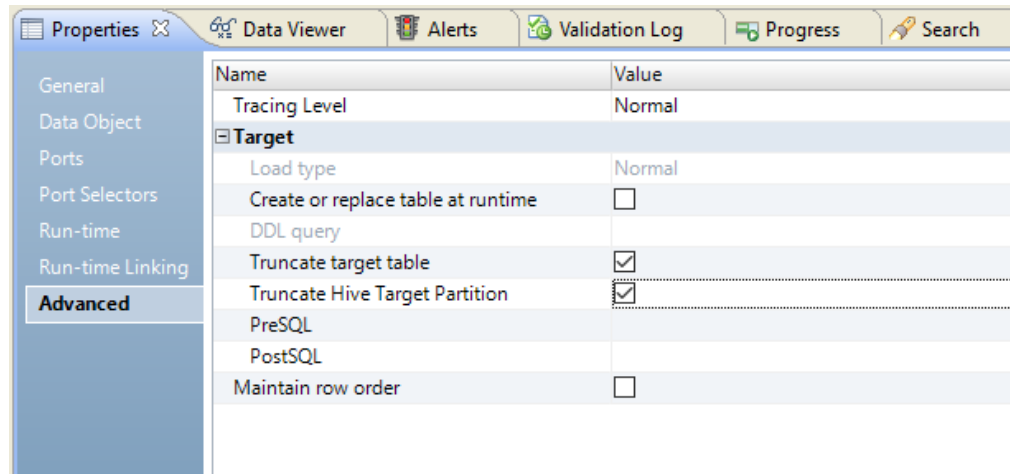
Truncating Partitions in a Hive Target

To truncate partitions in a Hive target, you must edit the write properties for the customized data object that you created for the Hive target in the Developer tool.

You can truncate partitions in a Hive target when you use the Blaze or Spark run-time engines to run the mapping.

1. Open the customized data object in the editor.
2. To edit write properties, select the **Input** transformation in the **Write** view, and then select the **Advanced** properties.

The following image shows the **Advanced** properties tab:



3. Select **Truncate Hive Target Partition**.
4. Select **Truncate target table**.

Scheduling, Queuing, and Node Labeling

You can use YARN schedulers, YARN queues, and node labels to optimize performance when you run a mapping in the Hadoop environment.

A YARN scheduler assigns resources to YARN applications on the Hadoop cluster while honoring organizational policies on sharing resources. You can configure YARN to use a fair scheduler or a capacity scheduler. A fair scheduler shares resources evenly among all jobs running on the cluster over time. A capacity scheduler allows multiple organizations to share a large cluster and distributes resources based on capacity allocations. The capacity scheduler guarantees each organization a certain capacity and distributes any excess capacity that is underutilized.

YARN queues are organizing structures for YARN schedulers and allow multiple tenants to share a cluster. The capacity of each queue specifies the percentage of cluster resources that are available for applications submitted to the queue. You can redirect Blaze, Spark, Hive, and Sqoop jobs to specific YARN queues.

Node labels allow YARN queues to run on specific nodes in a cluster. You can use node labels to partition a cluster into sub-clusters such that jobs run on nodes with specific characteristics. For example, you might label nodes that process data faster compared to other nodes. Nodes that are not labeled belong to the default partition. You can associate the node labels with capacity scheduler queues.

You can also use the node labels to configure the Blaze engine. When you use node labels to configure the Blaze engine, you can specify the nodes on the Hadoop cluster where you want the Blaze engine to run.

Note: You must install and configure Big Data Management for every node on the cluster, even if the cluster is not part of the queue you are using.

Enable Scheduling and Node Labeling

To enable scheduling and node labeling in the Hadoop environment, update the `yarn-site.xml` file in the domain environment.

Configure the following properties:

yarn.resourcemanager.scheduler.class

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

```
<property>
  <name>yarn.resourcemanager.scheduler.class</name>
  <value>org.apache.hadoop.yarn.server.resourcemanager.scheduler.[Scheduler Type].
[Scheduler Type]Scheduler</value>
</property>
```

For example:

```
<property>
  <name>yarn.resourcemanager.scheduler.class</name>

  <value>org.apache.hadoop.yarn.server.resourcemanager.scheduler.capacity.CapacityScheduler</value>
</property>
```

yarn.node-labels.enabled

Enables node labeling.

```
<property>
  <name>yarn.node-labels.enabled</name>
  <value>TRUE</value>
</property>
```

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

The HDFS location to update the node label dynamically.

```
<property>
  <name>yarn.node-labels.fs-store.root-dir</name>
  <value>hdfs://[Node name]:[Port]/[Path to store]/[Node labels]</value>
</property>
```

Define YARN Queues

You can define multiple YARN queues to redirect jobs to specific queues. You can redirect Blaze, Spark, Hive, and Sqoop jobs.

To redirect jobs on the Blaze engine to a specific queue, configure the following Blaze configuration property in the Hadoop connection:

Property	Description
YARN Queue Name	The YARN scheduler queue name used by the Blaze engine that specifies available resources on a cluster.

To redirect jobs on the Spark engine to a specific queue, configure the following Spark configuration property in the Hadoop connection:

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.

To redirect jobs on the Hive engine to a specific queue, configure the following Hive connection property:

Property	Description
Data Access Connection String	<p>The Hive connection string to specify the queue name for Hive SQL override mappings on the Blaze engine.</p> <p>Use the following format:</p> <ul style="list-style-type: none">- <code>MapReduce.mapred.job.queue.name=<YARN queue name></code>- <code>Tez.tez.queue.name=<YARN queue name></code> <p>For example, <code>jdbc:hive2://business.com:10000/default;principal=hive/_HOST@INFACRB?mapred.job.queue.name=root.test</code></p>

To redirect Sqoop jobs to a specific queue, configure the following JDBC connection property:

Property	Description
Sqoop Arguments	<p>The Sqoop connection-level argument to direct a MapReduce job to a specific YARN queue.</p> <p>Use the following format:</p> <p><code>-Dmapred.job.queue.name=<YARN queue name></code></p> <p>If you do not direct the job to a specific queue, the Spark engine uses the default queue.</p>

Configure the Blaze Engine to Use Node Labels

If you enable node labeling in the Hadoop environment, you can use node labels to run Blaze components as YARN applications on specific cluster nodes.

To run Blaze components on the labeled cluster nodes, specify the node labels when you configure the Blaze engine.

To specify node labels on the Blaze engine, list the node labels in the following Hadoop connection property:

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>&&</code> (AND), <code> </code> (OR), and <code>!</code> (NOT).</p>

Note: When the Blaze engine uses node labels, Blaze components might be redundant on the labeled nodes. If a node contains multiple labels and you specify the labels in different Hadoop connections, multiple Grid Manager, Orchestrator, or Job Monitor instances might run on the same node.

Spark Engine Optimization for Sqoop Pass-Through Mappings

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 in the JDBC connection or Sqoop mapping. 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
```

Rules and Guidelines for Sqoop Spark Engine Optimization

Consider the following rules and guidelines when you run Sqoop mappings on the Spark engine:

- The Data Integration Service does not optimize mapping performance in the following scenarios:
 - There are unconnected ports between the source and target in the mapping.
 - The data types of the source and target in the mapping do not match.
 - You write data to a partitioned Hive target table.
 - You run a mapping on an Azure HDInsight cluster that uses WASB to write data to an HDFS complex file target of the Parquet format.
- If you configure Hive-specific Sqoop arguments to write data to a Hive target, Sqoop ignores the arguments.
- If you configure a delimiter for a Hive target table that is different from the default delimiter, Sqoop ignores the delimiter.

Troubleshooting a Mapping in a Hadoop Environment

When I run a mapping with a Hive source or a Hive target on a different cluster, the Data Integration Service fails to push the mapping to Hadoop with the following error: Failed to execute query [exec0_query_6] with error code [10], error message [FAILED: Error in semantic analysis: Line 1:181 Table not found customer_eur], and SQL state [42000]].

When you run a mapping in a Hadoop environment, the Hive connection selected for the Hive source or Hive target, and the mapping must be on the same Hive metastore.

When I run a mapping with a Hadoop distribution on MapReduce 2, the Administrator tool shows the percentage of completed reduce tasks as 0% instead of 100%.

Verify that the Hadoop jobs have reduce tasks.

When the Hadoop distribution is on MapReduce 2 and the Hadoop jobs do not contain reducer tasks, the Administrator tool shows the percentage of completed reduce tasks as 0%.

When the Hadoop distribution is on MapReduce 2 and the Hadoop jobs contain reducer tasks, the Administrator tool shows the percentage of completed reduce tasks as 100%.

When I run mappings with SQL overrides concurrently, the mappings hang.

There are not enough available resources because the cluster is being shared across different engines.

Configure YARN to use the capacity scheduler and use different YARN scheduler queues for Blaze, Spark, and Hive.

When I configure a mapping to create a partitioned Hive table, the mapping fails with the error "Need to specify partition columns because the destination table is partitioned."

This issue happens because of internal Informatica requirements for a query that is designed to create a Hive partitioned table. For details and a workaround, see [Knowledge Base article 516266](#).

High Precision Decimal Data Type on the Hive Engine

In high-precision mode, the Hive engine supports decimal data types with precision up to 38 digits and a maximum scale of 38. The scale must be less than the precision.

When you run mappings in high-precision mode, the Hive engine maintains the precision at a maximum of 38 digits. As a result, precision loss can occur due to data overflow or intermediate calculations.

To minimize precision loss, configure the precision and scale to accurately reflect the data.

Precision Loss Due to Data Overflow

Precision loss can occur for decimal data types that run in high-precision mode on the Hive engine.

In high precision mode, transformations on the Hive engine support precision up to 38 digits. If the result of an expression has precision that is higher than 38, data overflow occurs. If data overflow occurs, excess digits are truncated in the result of the expression. To prevent data overflow, the Hive engine preserves the integral part of the result and adjusts the scale to maintain the precision at a maximum value of 38.

To maintain the precision at a maximum value of 38, the Hive engine subtracts digits of precision until the precision is 38. The following rules describe how the Hive engine adjusts the scale:

- The Hive engine subtracts the same number of digits from the scale if the number of digits subtracted from the precision is less than or equal to the scale.
- The Hive engine does not change the scale if the number of digits subtracted from the precision is greater than the scale.

When the Hive engine reduces the scale, the result is less accurate.

For example, the following expression multiplies two decimal data types `dec(38,2)` and `dec(2,2)`:

```
dec(38,2) * dec(2,2)
```

The precision and scale of the result is calculated according to the following expression:

```
dec(38+2,2+2) = dec(40,4)
```

Since the precision of the result is greater than 38, the Hive engine adjusts the scale to return the precision at a maximum value of 38. The following expression shows how the Hive engine adjusts the precision and scale:

```
dec(40-2,4-2)  
= dec(38,2)
```

Precision Loss Due to Intermediate Calculations

When mappings run in high-precision mode on the Hive engine, additional precision loss can occur in expressions if precision loss occurs during intermediate calculations. The precision loss that occurs during intermediate calculations can affect the return value of the expression.

Precision loss might occur during intermediate calculations if intermediate results are rounded.

For example, you might have the following expression:

```
206.85 * 891.94 * 0.01
```

The values 206.85 and 891.94 are decimal data types dec(38,2). The value 0.01 is decimal data type dec(2,2).

The expression can be computed with or without intermediate calculations.

Computing without Intermediate Calculations

If you compute the expression without intermediate calculations, the entire expression is computed at once:

```
206.85 * 891.94 * 0.01 = 1844.97789 = 1844.98
```

The result is a decimal data type dec(38,2), so the result 1844.97789 is rounded using a scale of 2. The final result is 1844.98.

Computing with Intermediate Calculations

If you compute the expression with intermediate calculations, the result of the expression might differ.

An expression might require intermediate calculations if the expression contains a user-defined function. For example, a user-defined function might require that the expression is calculated using the following steps:

```
Step 1. 206.85 * 0.01 = 2.0685 = 2.07
Step 2. 2.07 * 891.94 = 1846.3158 = 1846.32
```

In Step 1, the Hive engine processes the intermediate calculation. Since the intermediate result is a decimal data type dec(38,2), the intermediate result 2.0685 is rounded using a scale of 2. The resulting value 2.07 is used in the following calculation in Step 2. The result in Step 2 is another decimal data type dec(38,2), so the result in Step 2 is also rounded using a scale of 2. The final result is 1846.32.

Note that the value 1846.32 computed with intermediate calculations differs from the value 1844.98 that is computed without intermediate calculations. The precision loss that occurs during the intermediate calculations skews the return value.

Precision and Scale in Multiplication User-Defined Functions

When the Hive engine processes multiplication user-defined functions in high precision mode, the Hive engine calculates the precision and scale of the result based on the precision and scale of the input ports.

The Hive engine uses the following rules to calculate the precision and scale of the result of the user-defined function:

- If the difference between the precision and scale is greater than or equal to 32, the maximum scale of the result can be 6.
- If the difference between the precision and scale is less than 32, the maximum scale can be greater than 6.
- If the scale is greater than 6, the maximum difference between the precision and scale is 32.
- If the scale is less than 6, the difference between the precision and scale can be greater than 32 but less than 38.

If the Hive engine cannot represent the result, data overflow occurs. When data overflow occurs, the Hive engine writes NULL to the target.

For example, you might use a user-defined function to multiply the following decimal inputs, `dec(38,10)` and `dec(38,6)`:

```
123456789123456789.1234567891 * 123456789123456789.123456
```

The precision and scale of the result is (38,6), but the multiplication result is a decimal with more than 38 digits of precision. Since the Hive engine cannot represent the multiplication result as a decimal data type `dec(38,6)`, the Hive engine writes NULL to the target.

CHAPTER 3

Mapping Sources in the Hadoop Environment

This chapter includes the following topics:

- [Sources in a Hadoop Environment, 54](#)
- [Complex File Sources, 55](#)
- [Flat File Sources, 55](#)
- [Hive Sources, 56](#)
- [Intelligent Structure Model Sources, 59](#)
- [Relational Sources, 60](#)
- [Sqoop Sources, 60](#)

Sources in a Hadoop Environment

You can push a mapping to the Hadoop environment that includes a source from the native environment or from the Hadoop environment. Some sources have limitations when you reference them in the Hadoop environment.

You can run mappings with the following sources in a Hadoop environment:

- Complex file sources
- Flat file sources
- HBase sources
- Hive sources
- Intelligent structure model sources
- ODBC sources
- Relational sources
- Sqoop sources

When a mapping runs in the Hadoop environment, an HDFS source or a Hive source cannot reside on a remote cluster. A remote cluster is a cluster that is remote from the machine that the Hadoop connection references in the mapping.

Complex File Sources

A mapping that runs in the Hadoop environment can process complex files.

You can read files from the local file system or from HDFS. To read large volumes of data, you can connect a complex file source to read data from a directory of files that have the same format and properties. You can read compressed binary files.

A mapping that runs on the Blaze engine or the Hive engine can contain a Data Processor transformation. You can include a complex file data object without a Data Processor transformation to read complex files that are flat files. If the complex file is a hierarchical file, you must connect the complex file data object to a Data Processor transformation.

A mapping that runs on the Spark engine can process hierarchical data through complex data types. Use a complex file data object that represents the complex files in the Hadoop Distributed File System. If the complex file contains hierarchical data, you must enable the read operation to project columns as complex data types.

The following table shows the complex files that a mapping can process in the Hadoop environment:

File Type	Format	Blaze Engine	Spark Engine	Hive Engine
Avro	Flat	Supported	Supported	Supported
Avro	Hierarchical	Supported*	Supported**	Supported*
JSON	Flat	Supported*	Supported	Supported*
JSON	Hierarchical	Supported*	Supported**	Supported*
ORC	Flat	Not supported	Supported	Not supported
ORC	Hierarchical	Not supported	Not supported	Not supported
Parquet	Flat	Supported	Supported	Supported
Parquet	Hierarchical	Supported*	Supported**	Supported*
XML	Flat	Supported*	Not supported	Supported*
XML	Hierarchical	Supported*	Not supported	Supported*
* The complex file data object must be connected to a Data Processor transformation.				
** The complex file read operation must be enabled to project columns as complex data type.				

Flat File Sources

A mapping that is running in a Hadoop environment can read a flat file source from a native environment.

Consider the following limitations when you configure the mapping to read a flat file source:

- You cannot use an indirect source type.
- The row size in a flat file source cannot exceed 190 MB.

- You cannot use a command to generate or to transform flat file data and send the output to the flat file reader at run time.

Generate the Source File Name

You can generate the source file name for the flat file data object. The content of the file name column remains consistent across different modes of execution.

When you push processing to the specific engine for the required file types, the file name column returns the path based on the following formats:

Run-time Engine	Type of Files Processes	Returned Path
Hive	HDFS source files	<staged path><HDFS file path> For example, hdfs://host name:port/hive/warehouse/ff.txt
Hive	Flat files in the local system	<local file path> For example, /home/devbld/Desktop/ff.txt
Blaze	Flat files in the local system	<staged path><local file path> For example, hdfs://host name:port/hive/warehouse/home/devbld/Desktop/ff.txt
Spark	HDFS source files	hdfs://<host name>:<port>/<file name path> For example, hdfs://host name:port/hive/warehouse/ff.txt
Spark	Flat files in the local system	<local file path> For example, /home/devbld/Desktop/ff.txt

The file name column returns the content in the following format for a high availability cluster: hdfs://<host name>/<file name path>

For example, hdfs://irl dv:5008/hive/warehouse/ff.txt

Hive Sources

You can include Hive sources in an Informatica mapping that runs in the Hadoop environment.

Consider the following limitations when you configure a Hive source in a mapping that runs in the Hadoop environment:

- A mapping fails to run when you have Unicode characters in a Hive source definition.
- The third-party Hive JDBC driver does not return the correct precision and scale values for the Decimal data type. As a result, when you import Hive tables with a Decimal data type into the Developer tool, the Decimal data type precision is set to 38 and the scale is set to 0. Consider the following configuration rules and guidelines based on the version of Hive:
 - Hive 0.11. Accept the default precision and scale for the Decimal data type in the Developer tool.

- Hive 0.12. Accept the default precision and scale for the Decimal data type in the Developer tool.
- Hive 0.12 with Cloudera CDH 5.0. You can configure the precision and scale fields for source columns with the Decimal data type in the Developer tool.
- Hive 0.13 and above. You can configure the precision and scale fields for source columns with the Decimal data type in the Developer tool.
- Hive 0.14 or above. The precision and scale used for the Decimal data type in the Hive database also appears in the Developer tool.

A mapping that runs on the Spark engine can have partitioned Hive source tables and bucketed sources.

PreSQL and PostSQL Commands

You can create SQL commands for Hive sources. You can execute the SQL commands to execute SQL statements such as insert, update, and delete on the Hive source.

PreSQL is an SQL command that runs against the Hive source before the mapping reads from the source. PostSQL is an SQL command that runs against the Hive source after the mapping writes to the target.

You can use PreSQL and PostSQL on the Spark engine. The Data Integration Service does not validate PreSQL or PostSQL commands for a Hive source.

Note: You can manually validate the SQL by running the following query in a Hive command line utility:

```
CREATE VIEW <table name> (<port list>) AS <SQL>
```

where:

- <table name> is a name of your choice
- <port list> is the comma-delimited list of ports in the source
- <SQL> is the query to validate

Pre-Mapping SQL Commands

PreSQL is an SQL command that runs against a Hive source before the mapping reads from the source.

For example, you might use a Hive source in a mapping. The data stored in the Hive source changes regularly and you must update the data in the Hive source before the mapping reads from the source to make sure that the mapping reads the latest records. To update the Hive source, you can configure a PreSQL command.

Post-Mapping SQL Commands

PostSQL is an SQL command that runs against a Hive source after the mapping writes to the target.

For example, you might use a Hive source in a mapping. After the mapping writes to a target, you might want to delete the stage records stored in the Hive source. You want to run the command only after the mapping writes the data to the target to make sure that the data is not removed prematurely from the Hive source. To delete the records in the Hive source table after the mapping writes to the target, you can configure a PostSQL command.

Rules and Guidelines for Pre- and Post-Mapping SQL Commands

Consider the following restrictions when you run PreSQL and PostSQL commands against Hive sources:

- When you create an SQL override on a Hive source, you must enclose keywords or special characters in backtick (') characters.
- When you run a mapping with a Hive source in the Hadoop environment, references to a local path in pre-mapping SQL commands are relative to the Data Integration Service node. When you run a mapping with a Hive source in the native environment, references to local path in pre-mapping SQL commands are relative to the Hive server node.

Rules and Guidelines for Hive Sources on the Blaze Engine

You can include Hive sources in an Informatica mapping that runs on the Blaze engine.

Consider the following rules and guidelines when you configure a Hive source in a mapping that runs on the Blaze engine:

- Hive sources for a Blaze mapping include the TEXT, Sequence, Avro, RCfile, ORC, and Parquet storage formats.
- A mapping that runs on the Blaze engine can have bucketed Hive sources and Hive ACID tables.
- Hive ACID tables must be bucketed.
- The Blaze engine supports Hive tables that are enabled for locking.
- Hive sources can contain quoted identifiers in Hive table names, column names, and schema names.
- The TEXT storage format in a Hive source for a Blaze mapping can support ASCII characters as column delimiters and the newline characters as a row separator. You cannot use hex values of ASCII characters. For example, use a semicolon (;) instead of 3B.
- You can define an SQL override in the Hive source for a Blaze mapping.
- The Blaze engine can read from an RCFile as a Hive source. To read from an RCFile table, you must create the table with the `SerDe` clause.
- The Blaze engine can read from Hive tables that are compressed. To read from a compressed Hive table, you must set the `TBLPROPERTIES` clause.

RCFile as Hive Tables

The Blaze engine can read and write to RCFile as Hive tables. However, the Blaze engine supports only the ColumnarSerDe `SerDe`. In Hortonworks, the default `SerDe` for an RCFile is `LazyBinaryColumnarSerDe`. To read and write to an RCFile table, you must create the table by specifying the `SerDe` as

```
org.apache.hadoop.hive.serde2.columnar.ColumnarSerDe.
```

For example:

```
CREATE TABLE TEST_RCFile
(id int, name string)
ROW FORMAT SERDE
'org.apache.hadoop.hive.serde2.columnar.ColumnarSerDe' STORED AS RCFILE;
```

You can also set the default RCFile `SerDe` from the Ambari or Cloudera manager. Set the property `hive.default.rcfile.serde` to `org.apache.hadoop.hive.serde2.columnar.ColumnarSerDe`.

Compressed Hive Tables

The Blaze engine can read and write to Hive tables that are compressed. However, to read from a compressed Hive table or write to a Hive table in compressed format, you must set the `TBLPROPERTIES` clause as follows:

- When you create the table, set the table properties:

```
TBLPROPERTIES ('property_name'='property_value')
```
- If the table already exists, alter the table to set the table properties:

```
ALTER TABLE table_name SET TBLPROPERTIES ('property_name' = 'property_value');
```

The property name and value are not case sensitive. Depending on the file format, the table property can take different values.

The following table lists the property names and values for different file formats:

File Format	Table Property Name	Table Property Values
Avro	avro.compression	BZIP2, deflate, Snappy
ORC	orc.compress	Snappy, ZLIB
Parquet	parquet.compression	GZIP, Snappy
RCFile	rcfile.compression	Snappy, ZLIB
Sequence	sequencefile.compression	BZIP2, GZIP, LZ4, Snappy
Text	text.compression	BZIP2, GZIP, LZ4, Snappy

Note: Unlike the Hive engine, the Blaze engine does not write data in the default ZLIB compressed format when it writes to a Hive target stored as ORC format. To write in a compressed format, alter the table to set the TBLPROPERTIES clause to use ZLIB or Snappy compression for the ORC file format.

The following text shows sample commands to create table and alter table:

- Create table:

```
create table CBO_3T_JOINS_CUSTOMER_HIVE_SEQ_GZIP
(C_CUSTKEY DECIMAL(38,0), C_NAME STRING, C_ADDRESS STRING,
C_PHONE STRING, C_ACCTBAL DECIMAL(10,2),
C_MKTSEGMENT VARCHAR(10), C_COMMENT VARCHAR(117))
partitioned by (C_NATIONKEY DECIMAL(38,0))
TBLPROPERTIES ('sequencefile.compression'='gzip')
stored as SEQUENCEFILE;
```

- Alter table:

```
ALTER TABLE table_name
SET TBLPROPERTIES (avro.compression='BZIP2');
```

Intelligent Structure Model Sources

Intelligent structure model sources are valid in mappings in a Hadoop environment.

You can include the following sources in an Informatica mapping that runs in a Hadoop environment on a Spark engine:

- complex file
- Amazon S3
- Microsoft Azure Blob

Note: Intelligent structure model in data objects is available for technical preview. Technical preview functionality is supported but is unwarranted and is not production-ready. Informatica recommends that you use intelligent structure model in the Microsoft Azure Blob Storage in non-production environments only.

Relational Sources

Relational sources are valid in mappings that run in a Hadoop environment if you use the Hive engine or the Blaze engine. The Spark engine cannot run mappings with relational resources.

The Data Integration Service does not run pre-mapping SQL commands or post-mapping SQL commands against relational sources. You cannot validate and run a mapping with PreSQL or PostSQL properties for a relational source in a Hadoop environment.

The Data Integration Service can use multiple partitions to read from the following relational sources:

- IBM DB2
- Oracle

Note: You do not have to set maximum parallelism for the Data Integration Service to use multiple partitions in the Hadoop environment.

Sqoop Sources

Sqoop sources are valid in mappings in a Hadoop environment.

You can include a JDBC-compliant database as a Sqoop source in an Informatica mapping that runs in a Hadoop environment.

For example, you can include the following sources in a Sqoop mapping:

- Aurora
- Greenplum
- IBM DB2
- IBM DB2 for z/OS
- Microsoft SQL Server
- Netezza
- Oracle
- Teradata

Rules and Guidelines for Sqoop Sources

Consider the following rules and guidelines when you configure a Sqoop source in a mapping:

- If you specify a sort condition in a mapping, the Data Integration Service ignores the Order By condition.
- You cannot sort columns in a Sqoop source.
- You cannot read distinct rows from a Sqoop source.
- When you read data from an Oracle source through Sqoop and run the mapping on the Blaze or Spark engine, Sqoop treats the owner name as case sensitive.
- Sqoop uses the values that you configure in the **User Name** and **Password** fields of the JDBC connection. If you configure the --username or --password argument in a JDBC connection or mapping, Sqoop ignores the arguments. If you create a password file to access a database, Sqoop ignores the password file.

Rules and Guidelines for Sqoop Queries

Consider the following rules and guidelines when you configure a Sqoop query in a mapping:

- To override the default query in a mapping with an advanced query, you must define a mapping parameter and set its value to \$CONDITIONS. You must then include \$CONDITIONS in the WHERE clause of the custom query.
- If you define a custom query, you must verify that the metadata of the custom query matches the metadata of the source object. Otherwise, Sqoop might write blank values to the target.
- When you enable OraOop and configure an advanced query to read data from an Oracle source through Sqoop, the mapping fails on the Spark engine.

CHAPTER 4

Mapping Targets in the Hadoop Environment

This chapter includes the following topics:

- [Targets in a Hadoop Environment, 62](#)
- [Complex File Targets, 63](#)
- [Flat File Targets, 63](#)
- [HDFS Flat File Targets, 64](#)
- [Hive Targets, 64](#)
- [Relational Targets, 68](#)
- [Sqoop Targets, 68](#)

Targets in a Hadoop Environment

You can push a mapping to the Hadoop environment that includes a target from the native environment or from the Hadoop environment. Some sources have limitations when you reference them in the Hadoop environment.

You can run mappings with the following targets in a Hadoop environment:

- Complex files
- Flat file (native)
- Greenplum
- HBase
- HDFS flat file
- Hive
- IBM DB2
- Netezza
- ODBC
- Oracle
- Sqoop targets
- Teradata

A mapping that runs with the Spark engine can have partitioned Hive target tables and bucketed targets.

When a mapping runs in the Hadoop environment, an HDFS target or a Hive target cannot reside on a remote cluster. A remote cluster is a cluster that is remote from the machine that the Hadoop connection references in the mapping.

Complex File Targets

A mapping that runs in the Hadoop environment can process complex files.

The following table shows the complex files that a mapping can process in the Hadoop environment:

File Type	Format	Blaze Engine	Spark Engine	Hive Engine
Avro	Flat	Supported	Supported	Supported
Avro	Hierarchical	Supported*	Supported**	Supported*
JSON	Flat	Supported*	Supported	Supported*
JSON	Hierarchical	Supported*	Supported**	Supported*
ORC	Flat	Not supported	Supported	Not supported
ORC	Hierarchical	Not supported	Not supported	Not supported
Parquet	Flat	Supported	Supported	Supported
Parquet	Hierarchical	Supported*	Supported**	Supported*
XML	Flat	Supported*	Not supported	Supported*
XML	Hierarchical	Supported*	Not supported	Supported*
* The complex file writer object must be connected to a Data Processor transformation.				
** The complex file writer object must be enabled to project columns as complex data type.				

Flat File Targets

A mapping that is running in a Hadoop environment can write to a flat file target that is in a native environment.

Consider the following limitations when you configure a flat file target in a mapping that runs in a Hadoop environment:

- The Data Integration Service truncates the target files and reject files before writing the data. When you use a flat file target, you cannot append output data to target files and reject files.
- The Data Integration Service can write to a file output for a flat file target. When you have a flat file target in a mapping, you cannot write data to a command.

HDFS Flat File Targets

HDFS flat file targets are valid in mappings that run in a Hadoop environment.

When you use a HDFS flat file target in a mapping, you must specify the full path that includes the output file directory and file name. The Data Integration Service might generate multiple output files in the output directory when you run the mapping in a Hadoop environment.

Hive Targets

A mapping that is running in the Hadoop environment can write to a Hive target.

A Hive target can be an internal table or an external table. Internal Hive tables are managed by Hive and are also known as managed tables. External Hive tables are managed by an external source such as HDFS, Amazon S3, Azure Blob, WASB, or ADLS.

Consider the following restrictions when you configure a Hive target in a mapping that runs in the Hadoop environment:

- A mapping fails to run when you use Unicode characters in a Hive target definition.
- When you set up a dynamic target for a partitioned Hive table, the value used for the partition is the final column in the table. If the table has a dynamic partition column, the final column of the table is the dynamic partition column. To use a different column for the partition, move it to the last column of the table. If the table has multiple partition columns, the dynamic partition values are selected from the last columns of the upstream transformation. You can use an Expression transformation to reorder the columns if necessary.

When a mapping creates or replaces a Hive table, the type of table that the mapping creates depends on the run-time engine that you use to run the mapping.

The following table shows the table type for each run-time engine:

Run-Time Engine	Resulting Table Type
Blaze	MANAGED_TABLE
Spark	EXTERNAL_TABLE
Hive	MANAGED_TABLE

You can design a mapping to truncate an internal or external Hive table that is bucketed and partitioned.

In a mapping that runs on the Spark engine or the Blaze engine, you can create a custom DDL query that creates or replaces a Hive table at run time. However, with the Blaze engine, you cannot use a backtick (`) character in the DDL query. The backtick character is required in HiveQL when you include special characters or keywords in a query.

The Spark engine can write to bucketed Hive targets. Bucketing and partitioning of Hive tables can improve performance by reducing data shuffling and sorting.

PreSQL and PostSQL Commands

You can create SQL commands for Hive targets. You can execute the SQL commands to execute SQL statements such as insert, update, and delete on the Hive target.

PreSQL is an SQL command that runs against the Hive target before the mapping reads from the source.

PostSQL is an SQL command that runs against the Hive target after the mapping writes to the target.

You can use PreSQL and PostSQL on the Spark engine. The Data Integration Service does not validate PreSQL or PostSQL commands for a Hive target.

Pre-Mapping SQL Commands

PreSQL is an SQL command that runs against a Hive target before the mapping reads from a source.

For example, you might use a Hive target in a mapping. The data stored in the Hive target contains old records from the previous day and you must delete the old records in the Hive target before you run the mapping that writes new records to the target. To delete the old records in the Hive target, you can configure a PreSQL command.

Post-Mapping SQL Commands

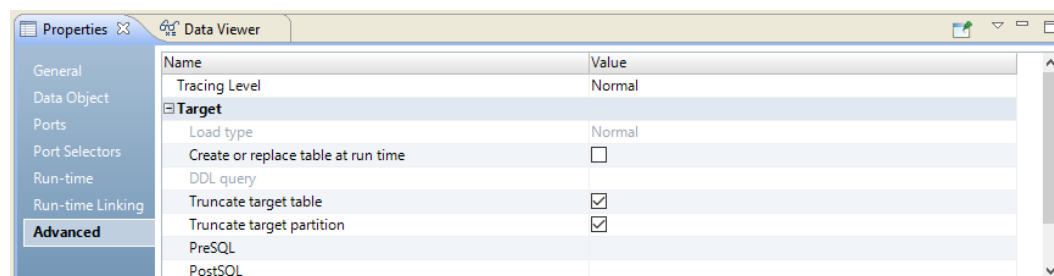
PostSQL is an SQL command that runs against a Hive source after the mapping writes to the target.

For example, you might use a Hive target in a mapping that runs in a test environment. After the mapping writes to the target, you might want to delete the records stored in the Hive target so that you can test the same Hive target again in the next mapping run. To delete the records in the Hive target, you can configure a PostSQL command.

Truncating Hive Targets

Truncate Hive target tables to delete the table contents. You can truncate internal and external Hive tables in the Hadoop environment.

Truncate a Hive table in the Hive table properties. The following image shows the Hive table properties:



To truncate the entire Hive table, choose the option to truncate the target table. To truncate a partition in the Hive table, you must choose to truncate the target table and to truncate the partition in the Hive target table.

Consider the following restrictions when you truncate a Hive table in a mapping that runs in the Hadoop environment:

- The Data Integration Service can truncate the partition in the Hive target in which the data is being inserted. You must choose to both truncate the partition in the Hive target and truncate the target table.
- You must truncate the target table to overwrite data to a Hive table with Hive version 0.7. The Data Integration Service ignores write, update override, delete, insert, and update strategy properties when it writes data to a Hive target.

Updating Hive Targets with an Update Strategy Transformation

For mappings that run on the Spark engine, you can use Hive MERGE statements to perform Update Strategy tasks. When a query uses a MERGE statement instead of INSERT, UPDATE or DELETE statements, processing is more efficient.

To use Hive MERGE, select **true** for the option in the Advanced Properties of the Update Strategy transformation.

The mapping ignores the Hive MERGE option and the Data Integration Service uses INSERT, UPDATE and DELETE to perform the operation under the following scenarios:

- The mapping runs on Blaze or Hive.
- In scenarios where MERGE is restricted by Hive implementation on particular Hadoop distributions.

The mapping log contains results of the operation, including whether restrictions affected results.

When the update affects partitioning or bucketing columns, updates to the columns are omitted.

Note: The Developer tool and the Data Integration Service do not validate against this restriction. If the Update Strategy expression violates these restrictions, the mapping might produce unexpected results.

Rules and Guidelines for Hive Targets on the Blaze Engine

You can include Hive targets in an Informatica mapping that runs on the Blaze engine.

Consider the following rules and guidelines when you configure a Hive target in a mapping that runs on the Blaze engine:

- A mapping that runs on the Blaze engine can have partitioned and bucketed Hive tables as targets. However, if you append data to a bucketed table, the Blaze engine overwrites the data in the bucketed target.
- Mappings that run on the Blaze engine can read and write to sorted targets.
- The Blaze engine supports Hive tables that are enabled for locking.
- The Blaze engine can create or replace Hive target tables.
- A mapping that runs on the Blaze engine can write to Hive ACID tables. To write to a Hive ACID table, the mapping must contain an Update Strategy transformation connected to the Hive target. The update strategy expression must flag each row for insert.
- The Blaze engine can write to Hive tables that are compressed. To write to a Hive table in compressed format, you must set the `TBLPROPERTIES` clause.
- For a mapping that writes to a Hive target on the Blaze engine, you cannot use update or delete update strategy properties.

RCFile as Hive Tables

The Blaze engine can read and write to RCFile as Hive tables. However, the Blaze engine supports only the ColumnarSerDe SerDe. In Hortonworks, the default SerDe for an RCFile is LazyBinaryColumnarSerDe. To read and write to an RCFile table, you must create the table by specifying the SerDe as

`org.apache.hadoop.hive.serde2.columnar.ColumnarSerDe`.

For example:

```
CREATE TABLE TEST_RCFile
(id int, name string)
ROW FORMAT SERDE
'org.apache.hadoop.hive.serde2.columnar.ColumnarSerDe' STORED AS RCFILE;
```

You can also set the default RCFile SerDe from the Ambari or Cloudera manager. Set the property `hive.default.rcfile.serde` to `org.apache.hadoop.hive.serde2.columnar.ColumnarSerDe`.

Compressed Hive Tables

The Blaze engine can read and write to Hive tables that are compressed. However, to read from a compressed Hive table or write to a Hive table in compressed format, you must set the `TBLPROPERTIES` clause as follows:

- When you create the table, set the table properties:
`TBLPROPERTIES ('property_name'='property_value')`
- If the table already exists, alter the table to set the table properties:
`ALTER TABLE table_name SET TBLPROPERTIES ('property_name' = 'property_value');`

The property name and value are not case sensitive. Depending on the file format, the table property can take different values.

The following table lists the property names and values for different file formats:

File Format	Table Property Name	Table Property Values
Avro	avro.compression	BZIP2, deflate, Snappy
ORC	orc.compress	Snappy, ZLIB
Parquet	parquet.compression	GZIP, Snappy
RCFile	rcfile.compression	Snappy, ZLIB
Sequence	sequencefile.compression	BZIP2, GZIP, LZ4, Snappy
Text	text.compression	BZIP2, GZIP, LZ4, Snappy

Note: Unlike the Hive engine, the Blaze engine does not write data in the default ZLIB compressed format when it writes to a Hive target stored as ORC format. To write in a compressed format, alter the table to set the `TBLPROPERTIES` clause to use ZLIB or Snappy compression for the ORC file format.

The following text shows sample commands to create table and alter table:

- Create table:

```
create table CBO_3T_JOINS_CUSTOMER_HIVE_SEQ_GZIP
(C_CUSTKEY DECIMAL(38,0), C_NAME STRING, C_ADDRESS STRING,
C_PHONE STRING, C_ACCTBAL DECIMAL(10,2),
C_MKTSEGMENT VARCHAR(10), C_COMMENT VARCHAR(117))
partitioned by (C_NATIONKEY DECIMAL(38,0))
TBLPROPERTIES ('sequencefile.compression'='gzip')
stored as SEQUENCEFILE;
```
- Alter table:

```
ALTER TABLE table_name
SET TBLPROPERTIES ('avro.compression'='BZIP2');
```

Relational Targets

Relational targets are valid in mappings in a Hadoop environment if you use the Hive or Blaze engine. The Spark engine cannot run mappings with relational targets.

The Data Integration Service does not run pre-mapping SQL commands or post-mapping SQL commands against relational targets in a Hadoop environment. You cannot validate and run a mapping with PreSQL or PostSQL properties for a relational target in a Hadoop environment.

The Data Integration Service can use multiple partitions to write to the following relational targets:

- IBM DB2
- Oracle

Note: You do not have to set maximum parallelism for the Data Integration Service to use multiple partitions in the Hadoop environment.

Sqoop Targets

Sqoop targets are valid in mappings in a Hadoop environment.

You can include a JDBC-compliant database as a Sqoop target in an Informatica mapping that runs in a Hadoop environment:

For example, you can include the following targets in a Sqoop mapping:

- Aurora
- Greenplum
- IBM DB2
- IBM DB2 for z/OS
- Microsoft SQL Server
- Netezza
- Oracle
- Teradata

You can insert data. You cannot update or delete data in a target. If you configure update arguments, Sqoop ignores them.

Rules and Guidelines for Sqoop Targets

Consider the following rules and guidelines when you configure a Sqoop target in a mapping:

- If a column name or table name contains a special character, the Sqoop export process fails.
- If you configure the **Maintain Row Order** property for a Sqoop target, the Data Integration Service ignores the property.
- If a mapping contains a Sqoop source, an Aggregator transformation, and a flat file target, you must disable the **Maintain Row Order** property for the target. Otherwise, the mapping fails.
- When you run a Sqoop mapping on the Blaze engine, verify that you have not deleted any target port from the mapping. Otherwise, the mapping fails.

- When you export null data to a Microsoft SQL Server column that is defined as not null, the Data Integration Service fails the Sqoop mapping on the Blaze engine instead of rejecting and writing the null data to the bad file.
- When you write data to an Oracle target through Sqoop and run the mapping on the Blaze or Spark engine, Sqoop treats the owner name as case sensitive.
- Sqoop uses the values that you configure in the **User Name** and **Password** fields of the JDBC connection. If you configure the --username or --password argument in a JDBC connection or mapping, Sqoop ignores the arguments. If you create a password file to access a database, Sqoop ignores the password file.

CHAPTER 5

Mapping Transformations in the Hadoop Environment

This chapter includes information about transformation support in the Hadoop environment.

Overview of Mapping Transformations in the Hadoop Environment

Due to the differences between native environment and Hadoop environment, only certain transformations are valid or are valid with restrictions in the Hadoop environment. Some functions, expressions, data types, and variable fields are not valid in the Hadoop environment.

Consider the following processing differences that can affect whether transformations and transformation behavior are valid or are valid with restrictions in the Hadoop environment:

- Hadoop uses distributed processing and processes data on different nodes. Each node does not have access to the data that is being processed on other nodes. As a result, the Hadoop execution engine might not be able to determine the order in which the data originated.
- Each of the run-time engines in the Hadoop environment can process mapping logic differently.

The following table lists transformations and support for different engines in a Hadoop environment:

Transformation	Supported Engines
<i>Transformations not listed in this table are not supported in the Hadoop environment.</i>	
Address Validator	<ul style="list-style-type: none">- Blaze- Spark*- Hive
Aggregator	<ul style="list-style-type: none">- Blaze- Spark- Hive
Case Converter	<ul style="list-style-type: none">- Blaze- Spark*- Hive

Transformation	Supported Engines
Classifier	<ul style="list-style-type: none"> - Blaze - Spark* - Hive
Comparison	<ul style="list-style-type: none"> - Blaze - Spark* - Hive
Consolidation	<ul style="list-style-type: none"> - Blaze - Spark* - Hive
Data Masking	<ul style="list-style-type: none"> - Blaze - Spark - Hive
Data Processor	<ul style="list-style-type: none"> - Blaze - Hive
Decision	<ul style="list-style-type: none"> - Blaze - Spark* - Hive
Expression	<ul style="list-style-type: none"> - Blaze - Spark - Hive
Filter	<ul style="list-style-type: none"> - Blaze - Spark - Hive
Java	<ul style="list-style-type: none"> - Blaze - Spark - Hive
Joiner	<ul style="list-style-type: none"> - Blaze - Spark - Hive
Key Generator	<ul style="list-style-type: none"> - Blaze - Spark* - Hive
Labeler	<ul style="list-style-type: none"> - Blaze - Spark* - Hive
Lookup	<ul style="list-style-type: none"> - Blaze - Spark - Hive
Match	<ul style="list-style-type: none"> - Blaze - Spark* - Hive

Transformation	Supported Engines
Merge	<ul style="list-style-type: none"> - Blaze - Spark* - Hive
Normalizer	<ul style="list-style-type: none"> - Blaze - Spark - Hive
Parser	<ul style="list-style-type: none"> - Blaze - Spark* - Hive
Python	<ul style="list-style-type: none"> - Spark
Rank	<ul style="list-style-type: none"> - Blaze - Spark - Hive
Router	<ul style="list-style-type: none"> - Blaze - Spark - Hive
Sequence Generator	<ul style="list-style-type: none"> - Blaze - Spark*
Sorter	<ul style="list-style-type: none"> - Blaze - Spark - Hive
Standardizer	<ul style="list-style-type: none"> - Blaze - Spark* - Hive
Union	<ul style="list-style-type: none"> - Blaze - Spark - Hive
Update Strategy	<ul style="list-style-type: none"> - Blaze - Spark* - Hive
Weighted Average	<ul style="list-style-type: none"> - Blaze - Spark* - Hive
<p>*Not supported for Big Data Streaming on the Spark engine. For more information about Big Data Streaming transformations, see the <i>Informatica Big Data Streaming User Guide</i>.</p>	

Address Validator Transformation in the Hadoop Environment

The Address Validator transformation processing behavior is the same on the Blaze, Spark, and Hive engines.

The Address Validator transformation is supported with the following restrictions:

- The Address Validator transformation cannot generate a certification report.
- The Address Validator transformation fails validation when it is configured to run in Interactive mode or Suggestion List mode.

Aggregrator Transformation in the Hadoop Environment

The Aggregator transformation processing in the Hadoop environment depends on the engine that runs the transformation.

Aggregator Transformation Support on the Blaze Engine

Some processing rules for the Blaze engine differ from the processing rules for the Data Integration Service.

Mapping Validation

Mapping validation fails in the following situations:

- The transformation contains stateful variable ports.
- The transformation contains unsupported functions in an expression.

Aggregate Functions

If you use a port in an expression in the Aggregator transformation but you do not use the port within an aggregate function, the Blaze engine might use any row in the port to process the expression.

The row that the Blaze engine uses might not be the last row in the port. Hadoop execution is distributed, and thus the Blaze engine might not be able to determine the last row in the port.

Data Cache Optimization

The data cache for the Aggregator transformation is optimized to use variable length to store binary and string data types that pass through the Aggregator transformation. The optimization is enabled for record sizes up to 8 MB. If the record size is greater than 8 MB, variable length optimization is disabled.

When variable length is used to store data that passes through the Aggregator transformation in the data cache, the Aggregator transformation is optimized to use sorted input and a pass-through Sorter transformation is inserted before the Aggregator transformation in the run-time mapping.

To view the Sorter transformation, view the optimized mapping or view the execution plan in the Blaze validation environment.

During data cache optimization, the data cache and the index cache for the Aggregator transformation are set to Auto. The sorter cache for the Sorter transformation is set to the same size as the data cache for the Aggregator transformation. To configure the sorter cache, you must configure the size of the data cache for the Aggregator transformation.

Aggregator Transformation Support on the Spark Engine

Some processing rules for the Spark engine differ from the processing rules for the Data Integration Service.

Mapping Validation

Mapping validation fails in the following situations:

- The transformation contains stateful variable ports.
- The transformation contains unsupported functions in an expression.

Aggregate Functions

If you use a port in an expression in the Aggregator transformation but you do not use the port within an aggregate function, the Spark engine might use any row in the port to process the expression.

The row that the Spark engine uses might not be the last row in the port. Hadoop execution is distributed, and thus the Spark engine might not be able to determine the last row in the port.

Data Cache Optimization

You cannot optimize the data cache for the transformation to store data using variable length.

Aggregator Transformation Support on the Hive Engine

Some processing rules for the Hive engine differ from the processing rules for the Data Integration Service.

Mapping Validation

Mapping validation fails in the following situations:

- The transformation contains stateful variable ports.
- The transformation contains unsupported functions in an expression.

Aggregate Functions

If you use a port in an expression in the Aggregator transformation but you do not use the port within an aggregate function, the Hive engine might use any row in the port to process the expression.

The row that the Hive engine uses might not be the last row in the port. Hadoop execution is distributed, and thus the Hive engine might not be able to determine the last row in the port.

Data Cache Optimization

You cannot optimize the data cache for the transformation to store data using variable length.

Case Converter Transformation in the Hadoop Environment

The Case Converter transformation is supported without restrictions on the Blaze, Spark, and Hive engines.

Classifier Transformation in the Hadoop Environment

The Classifier transformation is supported without restrictions on the Blaze, Spark, and Hive engines.

Comparison Transformation in the Hadoop Environment

The Comparison transformation is supported without restrictions on the Blaze, Spark, and Hive engines.

Consolidation Transformation in the Hadoop Environment

The Consolidation transformation processing behavior is the same on the Blaze, Spark, and Hive engines.

The Consolidation transformation might process data differently in the native environment and in a Hadoop environment.

The transformation might demonstrate the following differences in behavior:

- The transformation might process records in a different order in each environment.
- The transformation might identify a different record as the survivor in each environment.

Data Masking Transformation in the Hadoop Environment

The Data Masking transformation behavior is the same on the Blaze, Spark, and Hive engines.

The Data Masking transformation might process data differently in the native environment and in a Hadoop environment.

Mapping validation fails in the following situations in a Hadoop environment:

- The transformation is configured for repeatable expression masking.
- The transformation is configured for unique repeatable substitution masking.

Data Processor Transformation

The Data Processor transformation processing in the Hadoop environment depends on the engine that runs the transformation.

The Data Processor transformation is not supported on the Spark engine.

Data Processor Transformation Support on the Blaze Engine

Mapping validation fails when the transformation **Data processor mode** is set to **Input Mapping** or **Service and Input Mapping**.

Data Processor Transformation Support on the Hive Engine

Mapping validation fails in the following situations:

- The transformation contains more than one input port.
- The transformation contains pass-through ports.

Decision Transformation in the Hadoop Environment

The Decision transformation processing in the Hadoop environment depends on the engine that runs the transformation.

You can run Decision transformation on the following engines in a Hadoop environment:

- Blaze engine. Supported without restrictions.
- Spark engine. You must configure the Decision transformation properties to be partitionable.
- Hive engine. Supported without restrictions.

Expression Transformation in the Hadoop Environment

The Expression transformation processing in the Hadoop environment depends on the engine that runs the transformation.

Expression Transformation Support on the Blaze Engine

Mapping validation fails in the following situations:

- The transformation contains stateful variable ports.
- The transformation contains unsupported functions in an expression.

An Expression transformation with a user-defined function returns a null value for rows that have an exception error in the function.

Expression Transformation Support on the Spark Engine

Mapping validation fails in the following situations:

- The transformation contains stateful variable ports.
- The transformation contains unsupported functions in an expression.

If an expression results in numerical errors, such as division by zero or SQRT of a negative number, it returns an infinite or an NaN value. In the native environment, the expression returns null values and the rows do not appear in the output.

Expression Transformation Support on the Hive Engine

Mapping validation fails in the following situations:

- The transformation contains stateful variable ports.
- The transformation contains unsupported functions in an expression.

An Expression transformation with a user-defined function returns a null value for rows that have an exception error in the function.

Filter Transformation in the Hadoop Environment

The Filter transformation processing in the Hadoop environment depends on the engine that runs the transformation.

The Filter transformation is supported without restrictions on the Spark and Hive engines.

Filter Transformation Support on the Blaze Engine

When a mapping contains a Filter transformation on a partitioned column of a Hive source, the Blaze engine can read only the partitions that contain data that satisfies the filter condition. To push the filter to the Hive source, configure the Filter transformation to be the next transformation in the mapping after the source.

Java Transformation in the Hadoop Environment

The Java transformation processing in the Hadoop environment depends on the engine that runs the transformation.

Java Transformation Support on the Blaze Engine

To use external .jar files in a Java transformation, perform the following steps:

1. Copy external .jar files to the Informatica installation directory in the Data Integration Service machine at the following location: `<Informatica installation directory>/services/shared/jars`. Then recycle the Data Integration Service.
2. On the machine that hosts the Developer tool where you develop and run the mapping that contains the Java transformation:
 - a. Copy external .jar files to a directory on the local machine.
 - b. Edit the Java transformation to include an import statement pointing to the local .jar files.
 - c. Update the classpath in the Java transformation.
 - d. Compile the transformation.

Java Transformation Support on the Spark Engine

You can use complex data types to process hierarchical data.

Some processing rules for the Spark engine differ from the processing rules for the Data Integration Service.

General Restrictions

The Java transformation is supported with the following restrictions on the Spark engine:

- The Java code in the transformation cannot write output to standard output when you push transformation logic to Hadoop. The Java code can write output to standard error which appears in the log files.
- For date/time values, the Spark engine supports the precision of up to microseconds. If a date/time value contains nanoseconds, the trailing digits are truncated.

Partitioning

The Java transformation has the following restrictions when used with partitioning:

- The Partitionable property must be enabled in the Java transformation. The transformation cannot run in one partition.
- The following restrictions apply to the Transformation Scope property:
 - The value Transaction for transformation scope is not valid.
 - If you enable an input port for partition key, the transformation scope must be set to All Input.
 - Stateless must be enabled if the transformation scope is row.

Mapping Validation

Mapping validation fails in the following situations:

- You reference an unconnected Lookup transformation from an expression within a Java transformation.
- You select a port of a complex data type as the partition or sort key.

- You enable nanosecond processing in date/time and the Java transformation contains a port of complex data type with an element of a date/time type. For example, a port of type `array<data/time>` is not valid if you enable nanosecond processing in date/time.
- When you enable high precision, a validation error occurs in the following situations:
 - The Java transformation contains a port of a decimal data type.
 - The Java transformation contains a complex data type with an element of a decimal data type.

Using External .jar Files

To use external .jar files in a Java transformation, perform the following steps:

1. Copy external .jar files to the Informatica installation directory in the Data Integration Service machine at the following location:
`<Informatica installation directory>/services/shared/jars`
2. Recycle the Data Integration Service.
3. On the machine that hosts the Developer tool where you develop and run the mapping that contains the Java transformation:
 - a. Copy external .jar files to a directory on the local machine.
 - b. Edit the Java transformation to include an import statement pointing to the local .jar files.
 - c. Update the classpath in the Java transformation.
 - d. Compile the transformation.

Setting the JDK Path

To use complex ports in the Java transformation and to run Java user code directly on the Spark engine, you must set the JDK path.

In the Administrator tool, configure the following execution option for the Data Integration Service:

Property	Description
JDK Home Directory	<p>The JDK installation directory on the machine that runs the Data Integration Service. Changes take effect after you recycle the Data Integration Service.</p> <p>The JDK version that the Data Integration Service uses must be compatible with the JRE version on the cluster.</p> <p>For example, enter a value such as <code>/usr/java/default</code>.</p> <p>Default is blank.</p>

Java Transformation Support on the Hive Engine

You can enable the Stateless advanced property when you run mappings in a Hadoop environment.

The Java code in the transformation cannot write output to standard output when you push transformation logic to Hadoop. The Java code can write output to standard error which appears in the log files.

Some processing rules for the Hive engine differ from the processing rules for the Data Integration Service.

Partitioning

You can optimize the transformation for faster processing when you enable an input port as a partition key and sort key. The data is partitioned across the reducer tasks and the output is partially sorted.

The following restrictions apply to the Transformation Scope property:

- The value Transaction for transformation scope is not valid.
- If transformation scope is set to Row, a Java transformation is run by mapper script.
- If you enable an input port for partition Key, the transformation scope is set to All Input. When the transformation scope is set to All Input, a Java transformation is run by the reducer script and you must set at least one input field as a group-by field for the reducer key.

Using External .jar Files

To use external .jar files in a Java transformation, perform the following steps:

1. Copy external .jar files to the Informatica installation directory in the Data Integration Service machine at the following location: `<Informatica installation directory>/services/shared/jars`. Then recycle the Data Integration Service.
2. On the machine that hosts the Developer tool where you develop and run the mapping that contains the Java transformation:
 - a. Copy external .jar files to a directory on the local machine.
 - b. Edit the Java transformation to include an import statement pointing to the local .jar files.
 - c. Update the classpath in the Java transformation.
 - d. Compile the transformation.

Joiner Transformation in the Hadoop Environment

The Joiner transformation processing in the Hadoop environment depends on the engine that runs the transformation.

Joiner Transformation Support on the Blaze Engine

Mapping validation fails in the following situations:

- The transformation contains an inequality join and map-side join is disabled.
- The Joiner transformation expression references an unconnected Lookup transformation.

Map-side join is disabled when the Joiner transformation is configured for detail outer join or full outer join.

Joiner Transformation Support on the Spark Engine

Mapping validation fails in the following situations:

- Case sensitivity is disabled.
- The join condition is of binary data type or contains binary expressions.

Joiner Transformation Support on the Hive Engine

Mapping validation fails in the following situations:

- The transformation contains an inequality join.

Key Generator Transformation in the Hadoop Environment

The Key Generator transformation is supported without restrictions on the Blaze, Spark, and Hive engines.

Labeler Transformation in the Hadoop Environment

The Labeler transformation processing in the Hadoop environment depends on the engine that runs the transformation.

The Labeler transformation is supported without restrictions on the Blaze and Spark engines.

Labeler Transformation Support on the Hive Engine

A mapping with a Labeler transformation fails to run on a cluster node when the following conditions are true:

- The transformation reads reference data from a probabilistic model file.
- The version of the Java Development Kit on the cluster node differs from the version of the Java Development Kit that the Informatica installation specifies for probabilistic analysis.

Lookup Transformation in the Hadoop Environment

The Lookup transformation processing in the Hadoop environment depends on the engine that runs the transformation.

Lookup Transformation Support on the Blaze Engine

Mapping validation fails in the following situations:

- The cache is configured to be shared, named, persistent, dynamic, or uncached. The cache must be a static cache.

If you add a data object that uses Sqoop as a Lookup transformation in a mapping, the Data Integration Service does not run the mapping through Sqoop. It runs the mapping through JDBC.

Lookup Transformation Support on the Spark Engine

Some processing rules for the Spark engine differ from the processing rules for the Data Integration Service.

Mapping Validation

Mapping validation fails in the following situations:

- Case sensitivity is disabled.

- The lookup condition in the Lookup transformation contains binary data type.
- The cache is configured to be shared, named, persistent, dynamic, or uncached. The cache must be a static cache.

The mapping fails in the following situation:

- The transformation is unconnected and used with a Joiner or Java transformation.

Multiple Matches

When you choose to return the first, last, or any value on multiple matches, the Lookup transformation returns any value.

If you configure the transformation to report an error on multiple matches, the Spark engine drops the duplicate rows and does not include the rows in the logs.

Lookup Transformation Support on the Hive Engine

If you add a data object that uses Sqoop as a Lookup transformation in a mapping, the Data Integration Service does not run the mapping through Sqoop. It runs the mapping through JDBC.

When you run a mapping that contains a Lookup transformation, the Data Integration Service creates lookup cache .jar files. Hive copies the lookup cache .jar files to the following temporary directory: `/tmp/<user_name>/hive_resources`. The Hive parameter `hive.downloaded.resources.dir` determines the location of the temporary directory. You can delete the lookup cache .jar files specified in the LDTM log after the mapping completes to retrieve disk space.

Mapping Validation

Mapping validation fails in the following situations:

- The cache is configured to be shared, named, persistent, dynamic, or uncached. The cache must be a static cache.
- The lookup is a relational Hive data source.

Mappings fail in the following situations:

- The lookup is unconnected.

Match Transformation in the Hadoop Environment

The Match transformation processing in the Hadoop environment depends on the engine that runs the transformation.

Match Transformation Support on the Blaze Engine

Mapping validation fails when the Match transformation is configured to write identity index data to database tables.

A Match transformation generates cluster ID values differently in native and Hadoop environments. In a Hadoop environment, the transformation appends a group ID value to the cluster ID.

Match Transformation Support on the Spark Engine

Mapping validation fails when the Match transformation is configured to write identity index data to database tables.

A Match transformation generates cluster ID values differently in native and Hadoop environments. In a Hadoop environment, the transformation appends a group ID value to the cluster ID.

Match Transformation Support on the Hive Engine

Mapping validation fails if a Match transformation specifies an identity match type.

A Match transformation generates cluster ID values differently in native and Hadoop environments. In a Hadoop environment, the transformation appends a group ID value to the cluster ID.

Merge Transformation in the Hadoop Environment

The Merge transformation is supported without restrictions on the Blaze, Spark, and Hive engines.

Normalizer Transformation in the Hadoop Environment

The Normalizer transformation is supported without restrictions on the Blaze, Spark, and Hive engines.

Parser Transformation in the Hadoop Environment

The Parser transformation processing in the Hadoop environment depends on the engine that runs the transformation.

The Parser transformation is supported without restrictions on the Blaze and Spark engines.

Parser Transformation Support on the Hive Engine

A mapping with a Parser transformation fails to run on a cluster node when the following conditions are true:

- The transformation reads reference data from a probabilistic model file.
- The version of the Java Development Kit on the cluster node differs from the version of the Java Development Kit that the Informatica installation specifies for probabilistic analysis.

Python Transformation in the Hadoop Environment

The Python transformation is supported with restrictions on the Spark engine. The Python transformation is not supported on the Blaze and Hive engines, and in the native environment.

Python Transformation Support on the Spark Engine

Mapping validation fails if a user-defined default value is assigned to an output port.

Mapping execution fails in the following situations:

- An output port is not assigned a value in the Python code.
- The data types in corresponding input and output ports are not the same, and the Python code does not convert the data type in the input port to the data type in the output port.

Note: The Data Integration Service does not validate Python code.

Rank Transformation in the Hadoop Environment

The Rank transformation processing in the Hadoop environment depends on the engine that runs the transformation.

Rank Transformation Support on the Blaze Engine

The data cache for the Rank transformation is optimized to use variable length to store binary and string data types that pass through the Rank transformation. The optimization is enabled for record sizes up to 8 MB. If the record size is greater than 8 MB, variable length optimization is disabled.

When variable length is used to store data that passes through the Rank transformation in the data cache, the Rank transformation is optimized to use sorted input and a pass-through Sorter transformation is inserted before the Rank transformation in the run-time mapping.

To view the Sorter transformation, view the optimized mapping or view the execution plan in the Blaze validation environment.

During data cache optimization, the data cache and the index cache for the Rank transformation are set to Auto. The sorter cache for the Sorter transformation is set to the same size as the data cache for the Rank transformation. To configure the sorter cache, you must configure the size of the data cache for the Rank transformation.

Rank Transformation Support on the Spark Engine

Some processing rules for the Spark engine differ from the processing rules for the Data Integration Service.

Mapping Validation

Mapping validation fails in the following situations:

- Case sensitivity is disabled.
- The rank port is of binary data type.

Data Cache Optimization

You cannot optimize the data cache for the transformation to store data using variable length.

Rank Transformation Support on the Hive Engine

Some processing rules for the Hive engine differ from the processing rules for the Data Integration Service.

Mapping Validation

Mapping validation fails in the following situations:

- Case sensitivity is disabled.

Data Cache Optimization

You cannot optimize the data cache for the transformation to store data using variable length.

Router Transformation in the Hadoop Environment

The Router transformation is supported without restrictions on the Blaze, Spark, and Hive engines.

Sequence Generator Transformation in the Hadoop Environment

The Sequence Generator transformation processing in the Hadoop environment depends on the engine that runs the transformation.

The Sequence Generator transformation is not supported on the Hive engine.

Sequence Generator Transformation Support on the Blaze Engine

A mapping with a Sequence Generator transformation consumes significant resources when the following conditions are true:

- You set the **Maintain Row Order** property on the transformation to *true*.
- The mapping runs in a single partition.

Sequence Generator Transformation Support on the Spark Engine

The Sequence Generator transformation does not maintain row order in output data. If you enable the **Maintain Row Order** property on the transformation, the Data Integration Service ignores the property.

Sorter Transformation in the Hadoop Environment

The Sorter transformation processing in the Hadoop environment depends on the engine that runs the transformation.

Sorter Transformation Support on the Blaze Engine

Some processing rules for the Blaze engine differ from the processing rules for the Data Integration Service.

Mapping Validation

Mapping validation fails in the following situations:

- The target is configured to maintain row order and the Sorter transformation is not connected directly to a flat file target.

Global Sort

The Blaze engine can perform global sorts when the following conditions are true:

- The Sorter transformation is connected directly to flat file targets.
- The target is configured to maintain row order.
- The sort key is not a binary data type.

If any of the conditions are not true, the Blaze engine performs a local sort.

Data Cache Optimization

If a Sorter transformation is inserted before an Aggregator or Rank transformation to optimize the Aggregator or Rank data cache, the size of the sorter cache is the same size as the data cache for the Aggregator or Rank transformation. To configure the sorter cache, you must configure the size of the data cache for the Aggregator or Rank transformation.

Sorter Transformation Support on the Spark Engine

Mapping validation fails when case sensitivity is disabled.

The Data Integration Service logs a warning and ignores the Sorter transformation in the following situations:

- There is a type mismatch in between the target and the Sorter transformation sort keys.
- The transformation contains sort keys that are not connected to the target.
- The Write transformation is not configured to maintain row order.
- The transformation is not directly upstream from the Write transformation.

The Data Integration Service treats null values as low even if you configure the transformation to treat null values as high.

Data Cache Optimization

You cannot optimize the sorter cache to store data using variable length.

Sorter Transformation Support on the Hive Engine

Mapping validation fails when case sensitivity is disabled.

The Data Integration Service logs a warning and ignores the Sorter transformation in the following situations:

- There is a type mismatch between the Sorter transformation and the target.
- The transformation contains sort keys that are not connected to the target.
- The Write transformation is not configured to maintain row order.
- The transformation is not directly upstream from the Write transformation.

The Data Integration Service treats null values as low even if you configure the transformation to treat null values as high.

Data Cache Optimization

You cannot optimize the sorter cache to store data using variable length.

Standardizer Transformation in the Hadoop Environment

The Standardizer transformation is supported without restrictions on the Blaze, Spark, and Hive engines.

Union Transformation in the Hadoop Environment

The Union transformation is supported without restrictions on the Blaze, Spark, and Hive engines.

Update Strategy Transformation in the Hadoop Environment

The Update Strategy transformation processing in the Hadoop environment depends on the engine that runs the transformation.

Note: The Update Strategy transformation is supported only on Hadoop distributions that support Hive ACID.

Update Strategy Transformation Support on the Blaze Engine

You can use the Update Strategy transformation on the Hadoop distributions that support Hive ACID.

Some processing rules for the Blaze engine differ from the processing rules for the Data Integration Service.

General Restrictions

If the Update Strategy transformation receives multiple update rows for the same primary key value, the transformation selects one random row to update the target.

If multiple Update Strategy transformations write to different instances of the same target, the target data might be unpredictable.

The Blaze engine executes operations in the following order: deletes, updates, inserts. It does not process rows in the same order as the Update Strategy transformation receives them.

Hive targets always perform Update as Update operations. Hive targets do not support Update Else Insert or Update as Insert.

Mapping Validation and Compile Validation

Mapping validation fails in the following situations:

- The Update Strategy transformation is connected to more than one target.
- The Update Strategy transformation is not located immediately before the target.
- The Update Strategy target is not a Hive target.
- The Update Strategy transformation target is an external ACID table.
- The target does not contain a primary key.
- The Hive target property to truncate the target table at run time is enabled.
- The Hive target property to create or replace the target table at run time is enabled.

The mapping fails in the following situation:

- The target is not ORC bucketed.
- The Hive target is modified to have fewer rows than the actual table.

Compile validation errors occur and the mapping execution stops in the following situations:

- The Hive version is earlier than 0.14.
- The target table is not enabled for transactions.

Using Hive Target Tables

To use a Hive target table with an Update Strategy transformation, you must create the Hive target table with the following clause in the Hive Data Definition Language: `TBLPROPERTIES ("transactional"="true")`.

To use an Update Strategy transformation with a Hive target, verify that the following properties are configured in the `hive-site.xml` configuration set associated with the Hadoop connection:

```
hive.support.concurrency      true
hive.enforce.bucketing       true
hive.exec.dynamic.partition.mode nonstrict
hive.txn.manager              org.apache.hadoop.hive.q1.lockmgr.DbTxnManager
hive.compactor.initiator.on   true
hive.compactor.worker.threads 1
```

Update Strategy Transformation Support on the Spark Engine

You can use the Update Strategy transformation on the Hadoop distributions that support Hive ACID.

Some processing rules for the Spark engine differ from the processing rules for the Data Integration Service.

General Restrictions

The Update Strategy transformation does not forward rejected rows to the next transformation.

If the Update Strategy transformation receives multiple update rows for the same primary key value, the transformation selects one random row to update the target.

If multiple Update Strategy transformations write to different instances of the same target, the target data might be unpredictable.

If the mapping runs on the Spark engine, you can choose the Use Hive Merge option. The option has the following restrictions:

- If you run the mapping with Blaze or Hive, the Hive MERGE option is ignored. The Data Integration Service runs the transformation with the previous implementation.
- A single row for delete or update cannot match multiple rows in the target. When the mapping violates this restriction, the mapping fails with a runtime error.
- If you configure the Update Strategy expression to update partitioning or bucketing columns, the mapping ignores the Hive MERGE option. In addition, the mapping does not update the columns.

Note: The Developer tool and the Data Integration Service do not validate against these restrictions. If the expression or the mapping violates these restrictions, the mapping might run, but the results will not be as expected.

Hive targets always perform Update as Update operations. Hive targets do not support Update Else Insert or Update as Insert.

Mapping Validation

Mapping validation fails in the following situations:

- The Update Strategy transformation is connected to more than one target.
- The Update Strategy transformation is not located immediately before the target.
- The Update Strategy target is not a Hive target.
- The Update Strategy transformation target is an external ACID table.
- The target does not contain a connected primary key.
- The Hive target property to truncate the target table at run time is enabled.
- The Hive target property to create or replace the target table at run time is enabled.

The mapping fails in the following situations:

- The target table is not enabled for transactions.
- The target is not ORC bucketed.

Using Hive Target Tables

To use a Hive target table with an Update Strategy transformation, you must create the Hive target table with the following clause in the Hive Data Definition Language: `TBLPROPERTIES ("transactional"="true")`.

To use an Update Strategy transformation with a Hive target, verify that the following properties are configured in the `hive-site.xml` configuration set associated with the Hadoop connection:

```
hive.support.concurrency      true
hive.enforce.bucketing       true
hive.exec.dynamic.partition.mode nonstrict
hive.txn.manager              org.apache.hadoop.hive.ql.lockmgr.DbTxnManager
hive.compactor.initiator.on   true
hive.compactor.worker.threads 1
```

Update Strategy Transformation Support on the Hive Engine

The Update Strategy transformation is supported only on Hadoop distributions that support Hive ACID.

Some processing rules for the Hive engine differ from the processing rules for the Data Integration Service.

General Restrictions

Hive targets always perform Update as Update operations. Hive targets do not support Update Else Insert or Update as Insert.

When the Update Strategy transformation receives multiple update rows for the same key, the results might differ.

Mapping Validation and Compile Validation

Mapping validation fails in the following situations:

- The transformation is connected to more than one target.
- The transformation is not connected directly to the target.

The mapping fails in the following situations:

- The target is not ORC bucketed.

Compile validation errors occur and the mapping execution stops in the following situations:

- The target is not a Hive target on the same cluster.
- The Hive version is earlier than 0.14.
- A primary key is not configured.

Weighted Average Transformation in the Hadoop Environment

The Weighted Average transformation is supported without restrictions on the Blaze, Spark, and Hive engines.

CHAPTER 6

Processing Hierarchical Data on the Spark Engine

This chapter includes the following topics:

- [Processing Hierarchical Data on the Spark Engine Overview, 91](#)
- [How to Develop a Mapping to Process Hierarchical Data, 92](#)
- [Complex Data Types, 94](#)
- [Complex Ports, 99](#)
- [Complex Data Type Definitions, 101](#)
- [Type Configuration, 106](#)
- [Complex Operators, 110](#)
- [Complex Functions, 112](#)

Processing Hierarchical Data on the Spark Engine Overview

You can use complex data types, such as array, struct, and map, in mappings that run on the Spark engine. With complex data types, the Spark engine directly reads, processes, and writes hierarchical data in complex files.

The Spark engine can process hierarchical data in Avro, JSON, and Parquet complex files. The Spark engine uses complex data types to represent the native data types for hierarchical data in complex files. For example, a hierarchical data of type record in an Avro file is represented as a struct data type on the Spark engine.

You can develop mappings for the following hierarchical data processing scenarios:

- To generate and modify hierarchical data.
- To transform relational data to hierarchical data.
- To transform hierarchical data to relational data.
- To convert data from one complex file format to another. For example, read hierarchical data from an Avro source and write to a JSON target.

To read from and write to complex files, you create complex file data objects. Configure the read and write operations for the complex file data object to project columns as complex data types. Read and Write transformations based on these complex file data objects can read and write hierarchical data.

Configure the following objects and transformation properties in a mapping to process hierarchical data:

- Complex ports. To pass hierarchical data in a mapping, create complex ports. You create complex ports by assigning complex data types to ports.
- Complex data type definitions. To process hierarchical data of type struct, create or import complex data type definitions that represent the schema of struct data.
- Type configuration. To define the properties of a complex port, specify or change the type configuration.
- Complex operators and functions. To generate or modify hierarchical data, create expressions using complex operators and functions.

You can also use hierarchical conversion wizards to simplify some of the mapping development tasks.

How to Develop a Mapping to Process Hierarchical Data

Develop a mapping with complex ports, operators, and functions to process hierarchical data on the Spark engine.

Note: The tasks and the order in which you perform the tasks to develop the mapping depend on the mapping scenario.

The following list outlines the high-level tasks to develop and run a mapping to read, write, and process hierarchical data in complex files.

Create an HDFS connection.

Create a Hadoop Distributed File System (HDFS) connection to access data in complex files that are stored in the HDFS. You can create and manage an HDFS connection in the Administrator tool or the Developer tool.

Create a complex file data object.

1. Create a complex file data object to represent the complex files in the HDFS as sources or targets. The Developer tool creates the read and write operations when you create the complex file data object.
2. Configure the complex file data object properties.
3. In the read and write operations, enable the column file properties to project columns in the complex files as complex data types.

Create a mapping and add mapping objects.

1. Create a mapping, and add Read and Write transformations.
To read from and write to a complex file, add Read and Write transformations based on the complex file data object.
To write to an Avro or Parquet file, you can also create a complex file target from an existing transformation in the mapping.
2. Based on the mapping logic, add other transformations that are supported on the Spark engine.

Generate struct data.

Based on the mapping scenario, use one of the hierarchical conversion wizards to generate struct data. You can also perform the following steps manually:

Create or import complex data type definitions for struct ports.

1. Create or import complex data type definitions that represent the schema of the struct data. The complex data type definitions are stored in the type definition library, which is a Model repository object. The default name of the type definition library is `Type_Definition_Library`.
2. If a mapping uses one or more mapplets, rename the type definition libraries in the mapping and the mapplets to ensure that the names are unique.

Create and configure struct ports in transformations.

1. Create ports in transformations and assign struct complex data type.
2. Specify the type configuration for the struct ports.
You must reference a complex data type definition for the struct port.
3. Create expressions with complex functions to generate struct data.

Modify struct data.

You can convert struct data to relational or hierarchical data. If the struct data contains elements of primitive data types, you can extract the elements as relational data. If the struct data contains elements of complex data types, you can extract the elements as hierarchical data. Based on the mapping scenario, use one of the hierarchical conversion wizards to modify struct data. You can also perform the following steps manually.

1. Create output ports with port properties that match the element of the struct data that you want to extract.
2. Create expressions with complex operators or complex functions to modify the struct data.

Generate array data.

1. Create ports in transformations and assign array complex data type.
2. Specify the type configuration for the array ports.
3. Create expressions with complex functions to generate array data.

Modify array data.

You can convert array data to relational or hierarchical data. If the array data contains elements of primitive data types, you can extract the elements as relational data. If the array data contains elements of complex data types, you can extract the elements as hierarchical data. Based on the mapping scenario, use one of the hierarchical conversion wizards to modify array data. You can also perform the following steps manually:

1. Create output ports with port properties that match the element of the array data that you want to extract.
2. Create expressions with complex operators or complex functions to modify the array data.

Configure the transformations.

Link the ports and configure the transformation properties based on the mapping logic.

Configure the mapping to run on the Spark engine.

Configure the following mapping run-time properties:

1. Select Hadoop as the validation environment and Spark as the engine.
2. Select Hadoop as the execution environment and select a Hadoop connection.

Validate and run the mapping on the Spark engine.

1. Validate the mapping to fix any errors.
2. Optionally, view the Spark engine execution plan to debug the logic.

3. Run the mapping.

Complex Data Types

A complex data type is a transformation data type that represents multiple data values in a single column position. The data values are called elements. Elements in a complex data type can be of primitive or complex data types. Use complex data types to process hierarchical data in mappings that run on the Spark engine.

Transformation data types include the following data types:

Primitive data type

A transformation data type that represents a single data value in a single column position. Data types such as decimal, integer, and string are primitive data types. You can assign primitive data types to ports in any transformation.

Complex data type

A transformation data type that represents multiple data values in a single column position. Data types such as array, map, and struct are complex data types. You can assign complex data types to ports in some transformations for the Spark engine.

Nested data type

A complex data type that contains elements of complex data types. Complex data types such as an array of structs or a struct with an array of other structs are nested data types.

The following table lists the complex data types:

Complex Data Type	Description
array	An array is an ordered collection of elements. The elements can be of a primitive or complex data type. All elements in the array must be of the same data type.
map	A map is an unordered collection of key-value pair elements. The key must be of a primitive data type. The value can be of a primitive or complex data type.
struct	A struct is a collection of elements of different data types. The elements can be of primitive or complex data types. Struct has a schema that defines the structure of the data.

The following image shows primitive, complex, and nested data types assigned to ports in a transformation:

	Name	Type	Type Configuration	
1	emp_name	string	N/A	→ 1
2	emp_sal	decimal	N/A	
3	emp_phone	array	string []	→ 2
4	emp_id_dept	map	< integer, string >	
5	emp_address	struct	(typedef_adrs)	→ 3
6	emp_bonus	array	(typedef_perf) []	

1. Primitive data types
2. Complex data types
3. Nested data type

The ports emp_name and emp_sal are of primitive data types. The ports emp_phone, emp_id_dept, and emp_address are of complex data types. The port emp_bonus is of a nested data type. The array contains elements of type struct.

Array Data Type

An array data type represents an ordered collection of elements. To pass, generate, or process array data, assign array data type to ports.

An array is a zero-based indexed list. An array index indicates the position of the array element. For example, the array index 0 indicates the first element in an array. The transformation language includes operators to access array elements and functions to generate and process array data.

An array can be one-dimensional or multidimensional. A one-dimensional array is a linear array. A multidimensional array is an array of arrays. Array transformation data types can have up to five dimensions.

Format

```
array <data_type> []
```

The following table describes the arguments for this data type:

Argument	Description
array	Name of the array column or port.
data_type	Data type of the elements in an array. The elements can be primitive data types or complex data types. All elements in the array must be of the same data type.
[]	Dimension of the array represented as subscript. A single subscript [] represents a one-dimensional array. Two subscripts [] [] represent a two-dimensional array. Elements in each dimension are of the same data type.

The elements of an array do not have names. The number of elements in an array can be different for each row.

Array Examples

One-dimensional array

The following array column represents a one-dimensional array of string elements that contains customer phone numbers:

```
custphone string[]
```

The following example shows data values for the custphone column:

custphone

```
[205-128-6478,722-515-2889]
```

```
[107-081-0961,718-051-8116]
```

```
[107-031-0961,NULL]
```

Two-dimensional array

The following array column represents a two-dimensional array of string elements that contains customer work and personal email addresses.

```
email_work_pers string[][]
```

The following example shows data values for the email_work_pers column:

email_work_pers

```
[john_baer@xyz.com,jbaer@xyz.com][john.baer@fgh.com,jbaer@ijk.com]
```

```
[bobbi_apperley@xyz.com,bapperl@xyz.com][apperlbob@fgh.com,bobbi@ijk.com]
```

```
[linda_bender@xyz.com,lbender@xyz.com][l.bender@fgh.com,NULL]
```

Map Data Type

A map data type represents an unordered collection of key-value pair elements. A map element is a key and value pair that maps one thing to another. To pass, generate, or process map data, assign map data type to ports.

The key must be of a primitive data type. The value can be of a primitive or complex data type. A map data type with values of a complex data type is a nested map. A nested map can contain up to three levels of nesting of map data type.

The transformation language includes subscript operator to access map elements. It also includes functions to generate and process map data.

Format

```
map <primitive_type -> data_type>
```


The following table describes the arguments for this data type:

Argument	Description
map	Name of the map column or port.
primitive_type	Data type of the key in a map element. The key must be of a primitive data type.
data_type	Data type of the value in a map element. The value can be of a primitive or complex data type.

Map Example

The following map column represents map data with an integer key and a string value to map customer ids with customer names:

```
custid_name <integer -> string>
```

The following example shows data values for the custid_name column:

custid_name

```
<26745 -> 'John Baer'>
```

```
<56743 -> 'Bobbi Apperley'>
```

```
<32879 -> 'Linda Bender'>
```

Struct Data Type

A struct data type represents a collection of elements of different data types. A struct data type has an associated schema that defines the structure of the data. To pass, generate, or process struct data, assign struct data type to ports.

The schema for the struct data type determines the element names and the number of elements in the struct data. The schema also determines the order of elements in the struct data. Informatica uses complex data type definitions to represent the schema of struct data.

The transformation language includes operators to access struct elements. It also includes functions to generate and process struct data and to modify the schema of the data.

Format

```
struct {element_name1:value1 [, element_name2:value2, ...]}
```

Schema for the struct is of the following format:

```
schema {element_name1:data_type1 [, element_name2:data_type2, ...]}
```

The following table describes the arguments for this data type:

Argument	Description
struct	Name of the struct column or port.
schema	A definition of the structure of data. Schema is a name-type pair that determines the name and data type of the struct elements.
element_name	Name of the struct element.
value	Value of the struct element.
data_type	Data type of the element value. The element values can be of a primitive or complex data type. Each element in the struct can have a different data type.

Struct Example

The following schema is for struct data to store customer addresses:

```
address
{st_number:integer,st_name:string,city:string,state:string,zip:string}
```

The following example shows struct data values for the `cust_address` column:

cust_address

```
{st_number:154,st_name:Addison Ave,city:Redwood City,state:CA,zip:94065}

{st_number:204,st_name:Ellis St,city:Mountain View,state:CA,zip:94043}

{st_number:357,st_name:First St,city:Sunnyvale,state:CA,zip:94085}
```

Rules and Guidelines for Complex Data Types

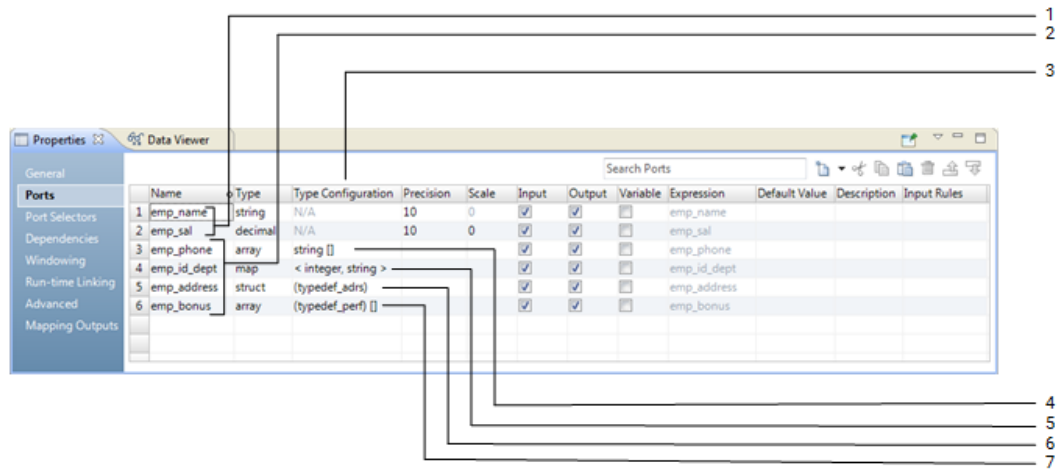
Consider the following rules and guidelines when you work with complex data types:

- A nested data type can contain up to 10 levels of nesting.
- A nested map can contain up to three levels of nesting of map data types.
- An array data type cannot directly contain an element of type array. Use multidimensional arrays to create a nested array. For example, an array with two dimensions is an array of arrays.
- A multidimensional array can contain up to five levels of nesting. The array dimension determines the levels of nesting.
- Each array in a multidimensional array must have elements of the same data type.

Complex Ports

A complex port is a port that is assigned a complex data type. Based on the complex data type, you must specify the complex port properties. Use complex ports in transformations to pass or process hierarchical data in a mapping.

The following image shows the complex ports and complex port properties on the Ports tab for a transformation:



1. Port
2. Complex port
3. Type configuration
4. Type configuration for an array port
5. Type configuration for a map port
6. Type configuration for a struct port
7. Type configuration for a port of nested data type

Based on the data type, a transformation can include the following ports and port properties:

Port

A port of a primitive data type that you can create in any transformation.

Complex port

A port of a complex or nested data type that you can create in some transformations. Array, map, and struct ports are complex ports. Based on the complex data type, you specify the complex port properties in the type configuration column.

Type configuration

A set of properties that you specify for the complex port. The type configuration determines the data type of the complex data type elements or the schema of the data. You specify the data type of the elements for array and map ports. You specify a complex data type definition for the struct port.

Type configuration for an array port

Properties that determine the data type of the array elements. In the image, the array port emp_phone is a one-dimensional array with an ordered collection of string elements. An array with string elements is also called an array of strings.

Type configuration for a map port

Properties that determine the data type of the key-value pair of the map elements. In the image, the map port `emp_id_dept` is an unordered collection of key-value pairs of type integer and string.

Type configuration for a struct port

Properties that determine the schema of the data. To represent the schema, you create or import a complex data type definition. In the image, the struct port `emp_address` references a complex data type definition `typedef_adrs`.

Type configuration for a port of nested data type

Properties that determine the nested data type. In the image, the array port `emp_bonus` is a one-dimensional array with an ordered collection of struct elements. The struct elements reference a complex data type definition `typedef_bonus`. An array with struct elements is also called an array of structs.

Complex Ports in Transformations

You can create complex ports in some transformations that are supported on the Spark engine. Read and Write transformations can represent ports that pass hierarchical data as complex data types.

You can create complex ports in the following transformations:

- Aggregator
- Expression
- Filter
- Java
- Joiner
- Lookup
- Normalizer
- Router
- Sorter
- Union

The Read and Write transformations can read and write hierarchical data in complex files. To read and write hierarchical data, the Read and Write transformations must meet the following requirements:

- The transformation must be based on a complex file data object.
- The data object read and write operations must project columns as complex data types.

Rules and Guidelines for Complex Ports

Consider the following rules and guidelines when you work with complex ports:

- Aggregator transformation. You cannot define a group by value as a complex port.
- Filter transformation. You cannot use the operators `>`, `<`, `>=`, and `<=` in a filter condition to compare data in complex ports.
- Joiner transformation. You cannot use the operators `>`, `<`, `>=`, and `<=` in a join condition to compare data in complex ports.
- Lookup transformation. You cannot use a complex port in a lookup condition.
- Rank transformation. You cannot define a group by or rank value as a complex port.

- Router transformation. You cannot use the operators `>`, `<`, `>=`, and `<=` in a group filter condition to compare data in complex ports.
- Sorter transformation. You cannot define a sort key value as a complex port.
- You can use complex operators to specify an element of a complex port that is of a primitive data type. For example, an array port "emp_names" contains string elements. You can define a group by value as `emp_names[0]`, which is of type string.

Creating a Complex Port

Create complex ports in transformations to pass or process hierarchical data in mappings that run on the Spark engine.

1. Select the transformation in the mapping.
2. Create a port.
3. In the **Type** column for the port, select a complex data type.

The complex data type for the port appears in the Type column.

After you create a complex port, specify the type configuration for the complex port.

Complex Data Type Definitions

A complex data type definition represents the schema of the data. Use complex data type definitions to define the schema of the data for a struct port. If the complex port is of type struct or contains elements of type struct, you must specify the type configuration to reference a complex data type definition.

You can create or import complex data type definitions. You import complex data type definitions from complex files. Complex data type definitions are stored in the type definition library, which is a Model repository object.

When you create a mapping or a mapplet, the Developer tool creates a type definition library with the name `Type_Definition_Library`. You can view and rename the type definition library and the complex data type definitions in the Outline view of a mapping. The complex data type definitions appear on the type definition library tab of the mapping editor. The tab name is based on the name of the type definition library. You create or import complex data type definitions on the type definition library tab. When a mapping uses one or more mapplets, rename the type definition libraries in the mapping and the mapplets to ensure that the names are unique.

Type definition library and complex data type definitions are specific to a mapping or mapplet. A mapping or a mapplet can have one or more complex data type definitions. You can reference the same complex data type definition in one or more complex ports in the mapping or mapplet.

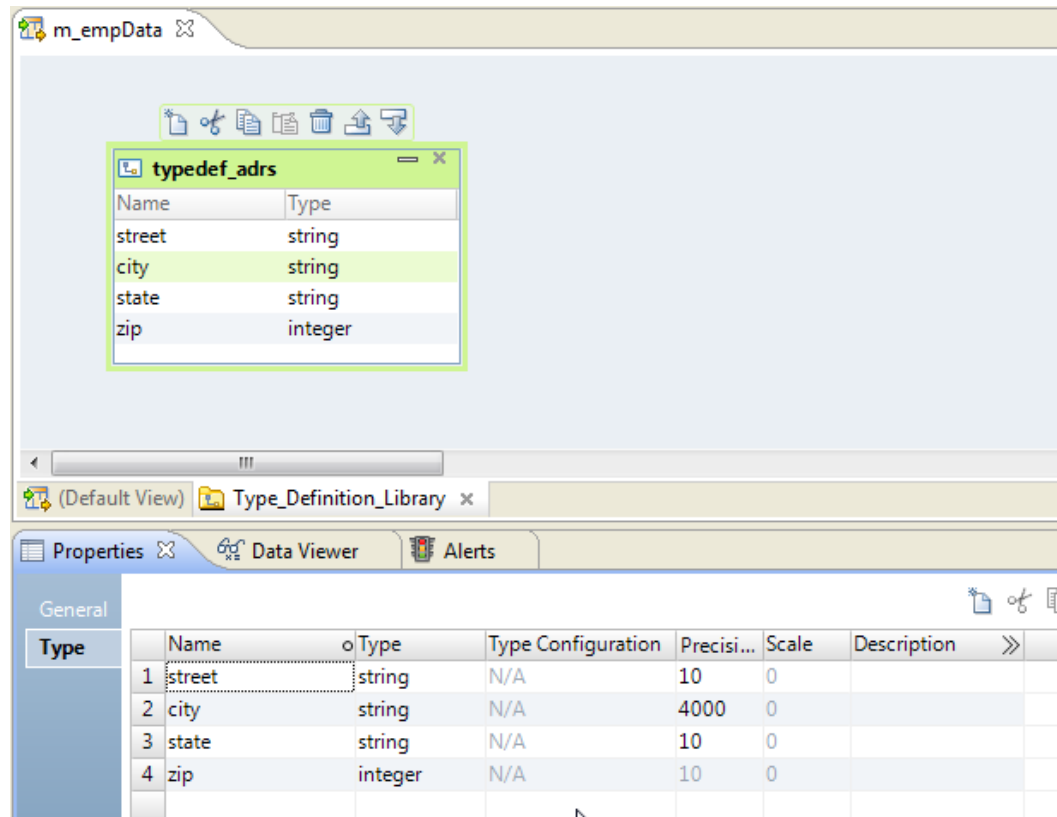
Complex Data Type Definition Example

The complex ports `work_address` and `home_address` of type struct have the following schema:

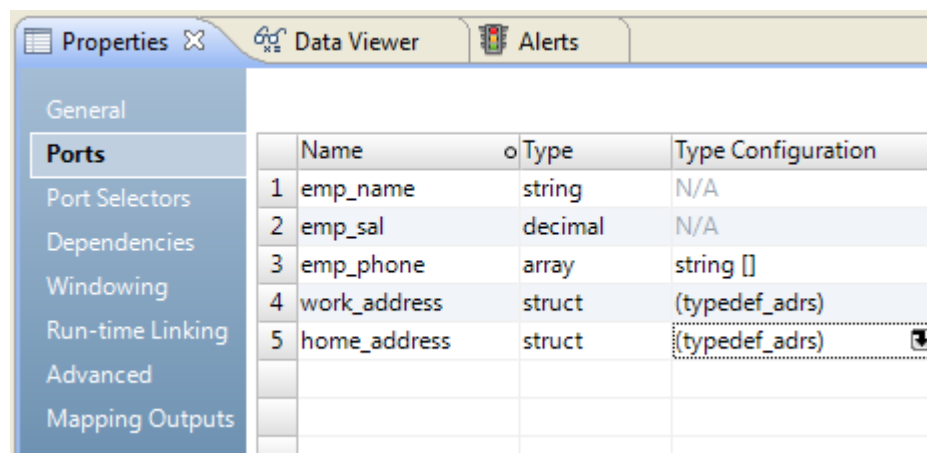
```
{street:string, city:string, state:string, zip:integer}
```

In the mapping, you create a complex data type definition `typedef_adrs` that defines the schema of the data. Then, you specify the type configuration of the struct ports `work_address` and `home_address` to use the `typedef_adrs` complex data type definition.

The following image shows a complex data type definition typedef_adrs:



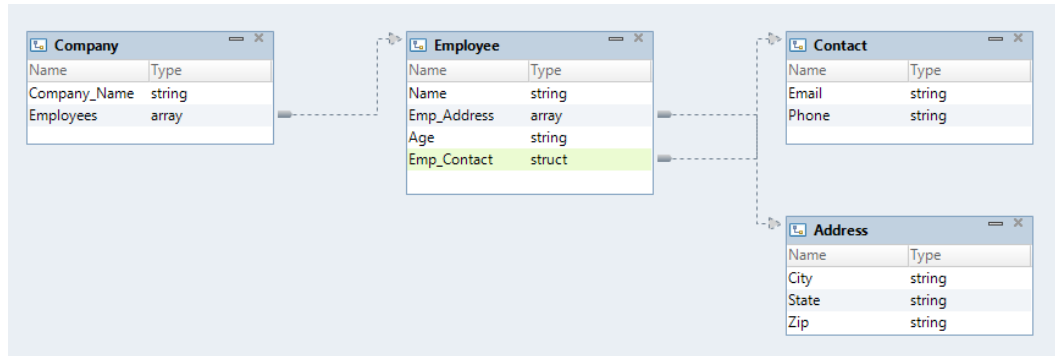
The following image shows two struct ports work_address and home_address that reference the complex data type definition typedef_adrs:



Nested Data Type Definitions

Elements of a complex data type definition can reference one or more complex data type definitions in the type definition library. Such complex data type definitions are called nested data type definitions.

The following image shows a nested data type definition Company on the type definition library tab:



The nested data type definition Company references the following complex data type definitions:

- In the complex data type definition Company, the array element Employees references the complex data type definition Employee.
- In the complex data type definition Employee, the Emp_Address element references the complex data type definition Address.
- In the complex data type definition Employee, the Emp_Contact element references the complex data type definition Contact.

Note: In a recursive data type definition, one of the complex data type definitions at any level is the same as any of its parents. You cannot reference a recursive data type definition to a struct port or a struct element in a complex port.

Rules and Guidelines for Complex Data Type Definitions

Consider the following rules and guidelines when you work with complex data type definitions:

- A struct port or a struct element in a complex port must reference a complex data type definition.
- You cannot reference a complex data type definition in one mapping to a complex port in another mapping.
- You cannot reference a recursive data type definition to a struct port or a struct element in a complex port.
- Changes to the elements in a complex data type definition are automatically propagated to complex ports that reference the complex data type definition.
- If you change the name of the type definition library or the complex data type definition, you must update the name in expressions that use complex functions such as STRUCT_AS, RESPEC, and CAST.

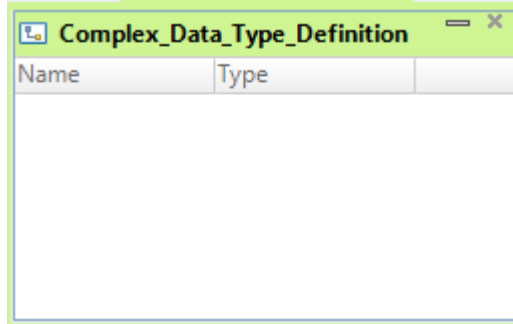
Creating a Complex Data Type Definition

A complex data type definition represents the schema of the data of type struct. You can create a complex data type definition for a mapping or a mapplet on the Type Definition Library tab of the mapping editor.

1. In the Developer tool, open the mapping or mapplet where you want to import a complex data type definition.

2. In the **Outline** view, select the **Type Definition Library** to view the **Type_Definition_Library** tab in the mapping editor.
3. Right-click an empty area of the editor and click **New Complex Data Type Definition**.

An empty complex data type definition appears on the Type Definition Library tab of the mapping editor.



4. Optionally, change the name of the complex data type definition.
5. Click the **New** button to add elements to the data type definition with a name and a type.
The data type of the element can be of primitive or complex data type.

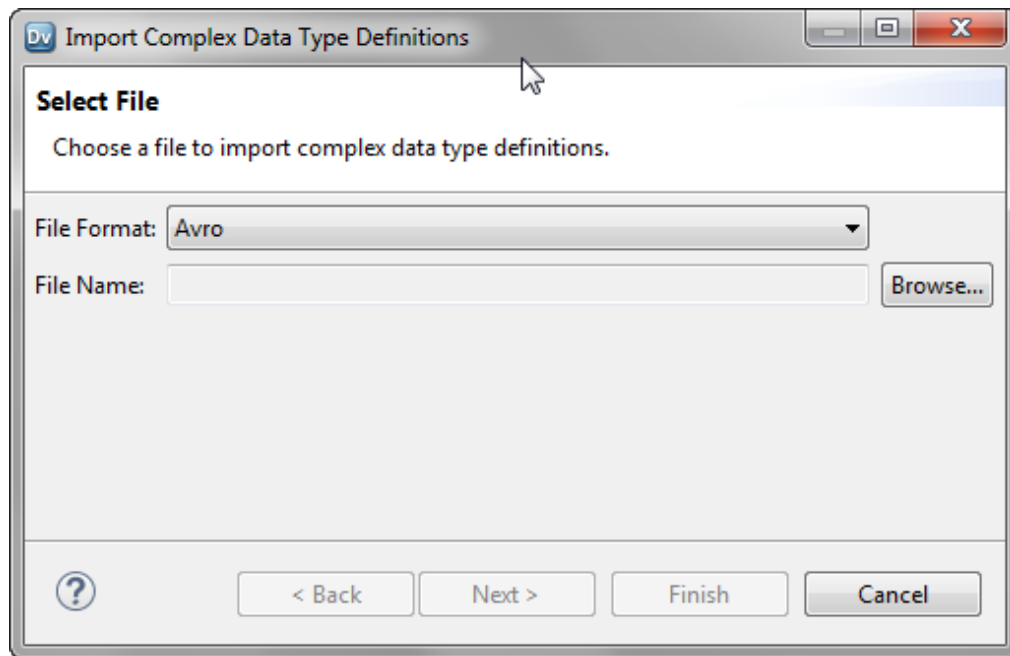
Importing a Complex Data Type Definition

A complex data type definition represents the schema of the data of type struct. You can import the schema of the hierarchical data in the complex file to the type definition library as a complex data type definition.

Import complex data type definitions from an Avro, JSON, or Parquet schema file. For example, you can import an Avro schema from an .avsc file as a complex data type definition.

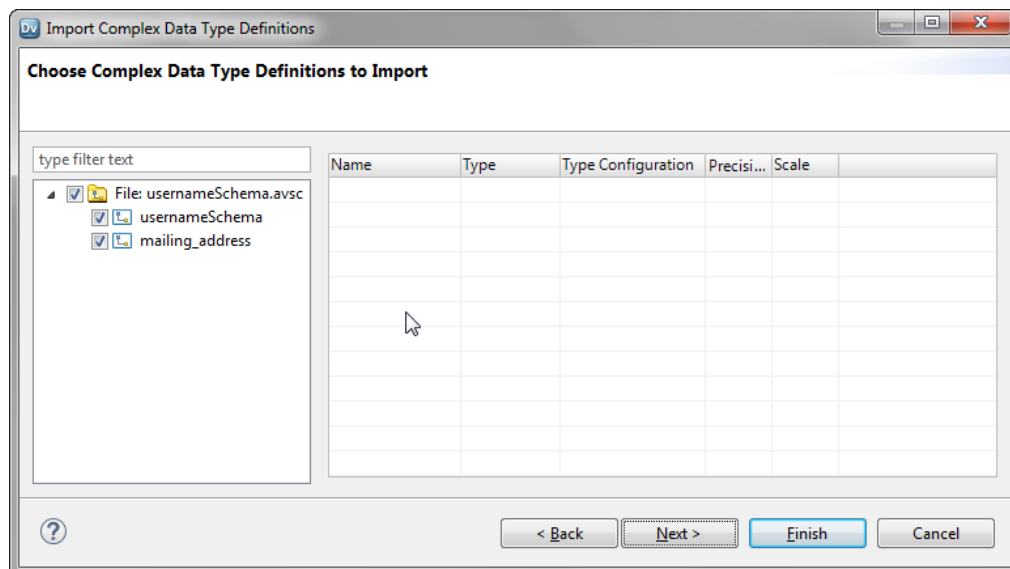
1. In the Developer tool, open the mapping or mapplet where you want to import a complex data type definition.
2. In the **Outline** view, select the **Type Definition Library** to view the **Type_Definition_Library** tab in the mapping editor.
3. Right-click an empty area of the editor and click **Import Complex Data Type Definitions**.

The **Import Complex Data Type Definitions** dialog box appears.



4. Select a complex file format from the **File Format** list.
5. Click **Browse** to select the complex file schema from which you want to import the complex data type definition.
6. Navigate to the complex file schema and click **Open**.
7. Click **Next**.

The **Choose Complex Data Type Definitions to Import** page appears.



8. Select one or more schemas from the list to import.
9. Click **Next**.

The **Complex Data Type Definitions to Import** page appears.

10. Review the complex data type definitions that you want to import and click **Finish**.

The complex data type definition appears in the **Type Definition Library** tab of the mapping or the mapplet.

Type Configuration

A type configuration is a set of complex port properties that specify the data type of the complex data type elements or the schema of the data. After you create a complex port, you specify or change the type configuration for the complex port on the Ports tab.

The type configuration for a complex port depends on the complex data type assigned to the port. You must specify a complex data type definition for a struct port or a struct element in a complex port. Array and map ports have a default type configuration that you can change.

The following table lists the type configuration properties for complex ports:

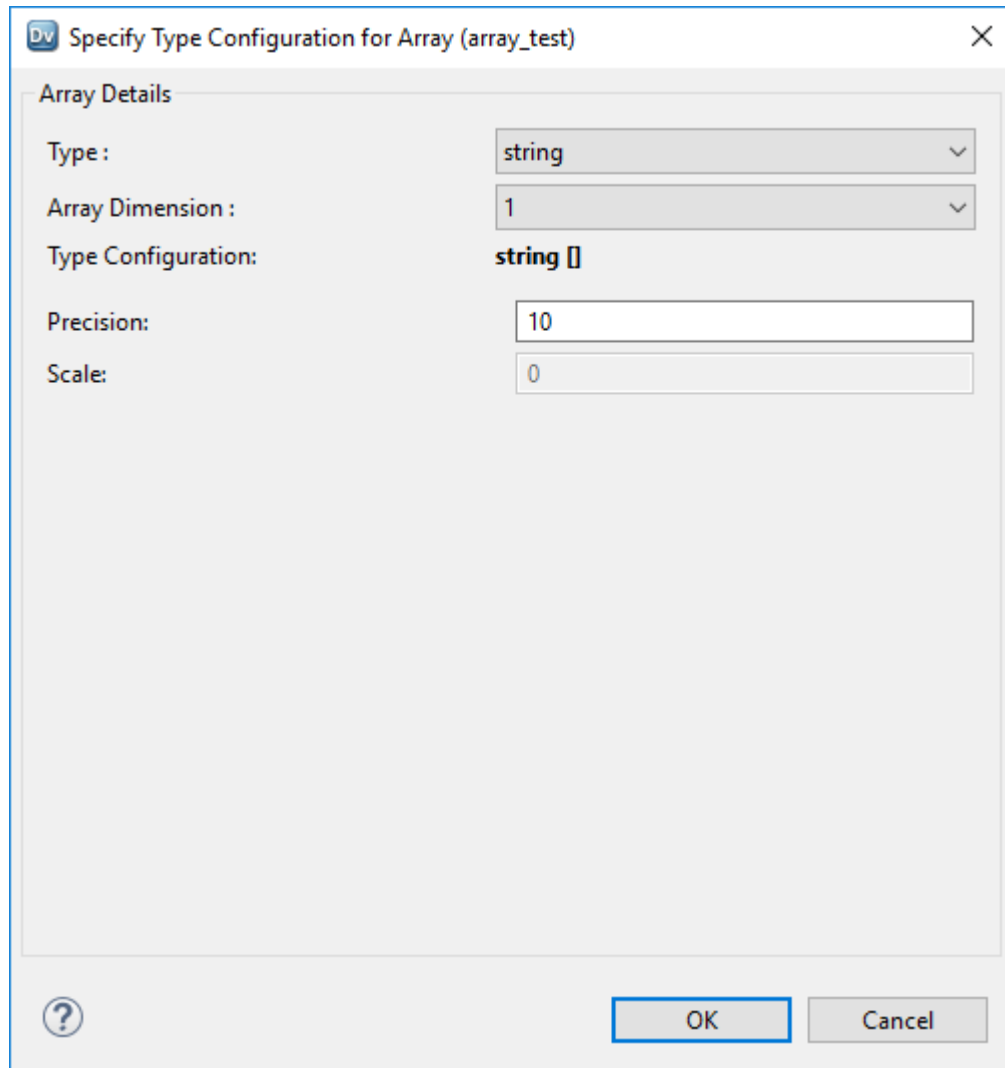
Complex Port	Type Configuration Properties
array	<ul style="list-style-type: none">- Type. Data type of the elements of an array. Default is string.- Array Dimension. Number of levels for an array. Default is 1.- Precision. Maximum number of significant digits for numeric data types, or maximum number of characters for string data types.- Scale. Maximum number of digits after the decimal point of numeric values.
map	<ul style="list-style-type: none">- Type. Data type of the key-value pair. Default is string for key and value.- Precision. Maximum number of significant digits for numeric data types, or maximum number of characters for string data types.- Scale. Maximum number of digits after the decimal point of numeric values.
struct	<ul style="list-style-type: none">- Complex data type definition. The schema of the data that you created or imported as complex data type definition in the type definition library.

Changing the Type Configuration for an Array Port

The default type configuration for an array is `string[]`, which is a one-dimensional array of strings. You can change the type configuration for the array port. Specify the data type of the array elements and the array dimension.

1. In the transformation, select the array port for which you want to specify the type configuration.
2. In the **Type Configuration** column, click the **Open** button.

The **Specify Type Configuration for Array** dialog box appears.

The image shows a dialog box titled "Specify Type Configuration for Array (array_test)". It has a close button (X) in the top right corner. The dialog is divided into two main sections. The top section, labeled "Array Details", contains five configuration options: "Type:" with a dropdown menu showing "string"; "Array Dimension:" with a dropdown menu showing "1"; "Type Configuration:" displaying "string []"; "Precision:" with a text input field containing "10"; and "Scale:" with a text input field containing "0". The bottom section of the dialog contains a help icon (question mark in a circle) on the left and two buttons, "OK" and "Cancel", on the right.

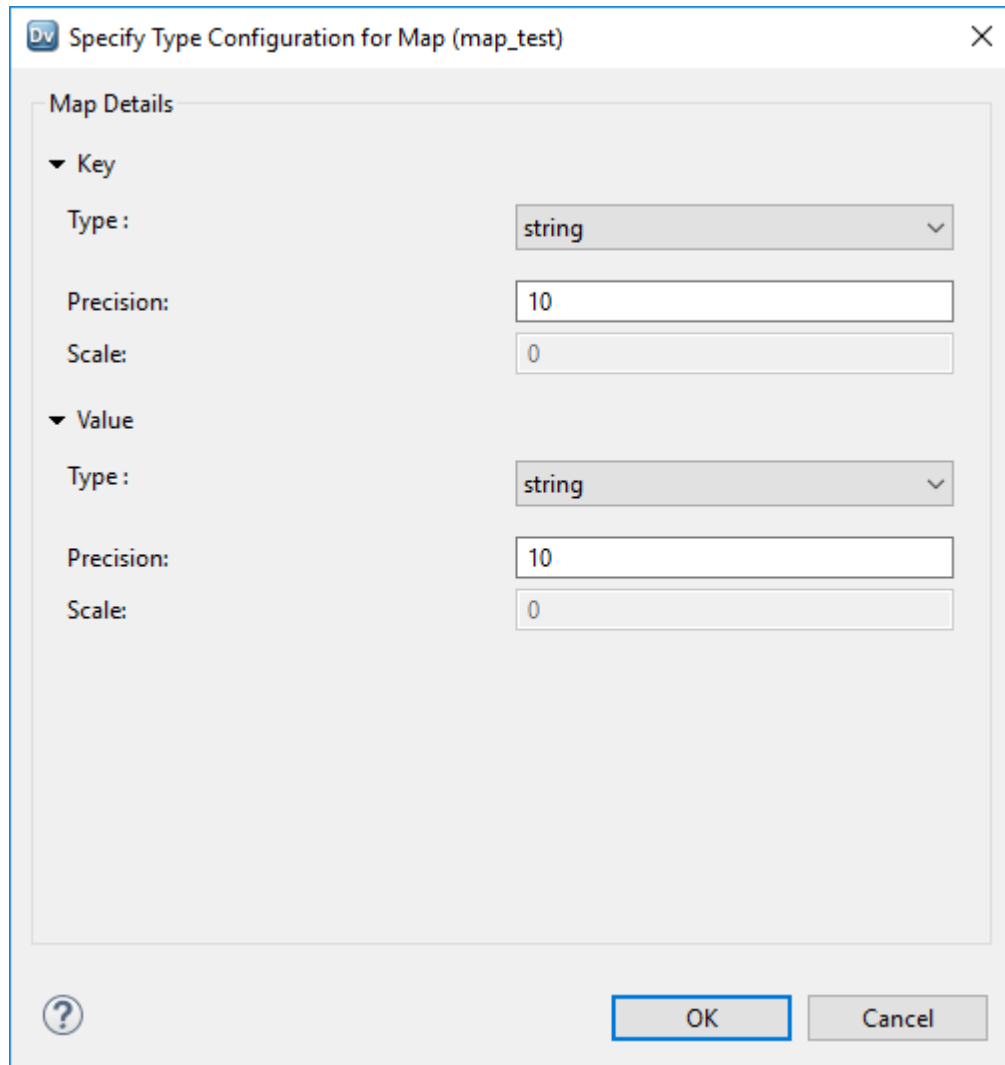
3. In the **Type** drop-down menu, change the data type of the elements in the array.
For example, if the elements in the array are of type integer, select **integer**.
The type configuration value appears as **integer []**.
4. Optionally, in the **Array Dimension** drop-down menu, select a number to create a nested array.
For example, array dimension 3 creates a nested array of 3 levels. The type configuration value appears as **integer [] [] []**.
5. Optionally, configure the following properties:
 - **Precision.** Maximum number of significant digits for numeric data types, or maximum number of characters for string data types.
 - **Scale.** Maximum number of digits after the decimal point of numeric values.
6. Click **OK**.
On the **Ports** tab, the specified type configuration appears in the **Type Configuration** column of the array port.

Changing the Type Configuration for a Map Port

The default type configuration for a map is `<string, string>`. You can change the type configuration for a map port. Specify the data type of the key-value pair.

1. In the transformation, select the map port for which you want to set the type configuration.
2. In the **Type Configuration** column, click the **Open** button.

The **Type Configuration** dialog box appears:



The dialog box titled "Specify Type Configuration for Map (map_test)" contains a "Map Details" section. Under "Key", there are three fields: "Type:" with a dropdown menu set to "string", "Precision:" with a text box containing "10", and "Scale:" with a text box containing "0". Under "Value", there are three fields: "Type:" with a dropdown menu set to "string", "Precision:" with a text box containing "10", and "Scale:" with a text box containing "0". At the bottom right are "OK" and "Cancel" buttons, and at the bottom left is a help icon (question mark in a circle).

3. Specify the properties for the key and value of the map data type.
 - a. In the **Type** drop-down menu, select the data type.
 - b. Optionally, in the **Precision** box, specify the maximum number of significant digits for numeric data types, or maximum number of characters for string data types.
 - c. Optionally, in the **Scale** box, specify the maximum number of digits after the decimal point of numeric values.
4. Click **OK**.

On the **Ports** tab, the specified type configuration appears in the **Type Configuration** column of the map port.

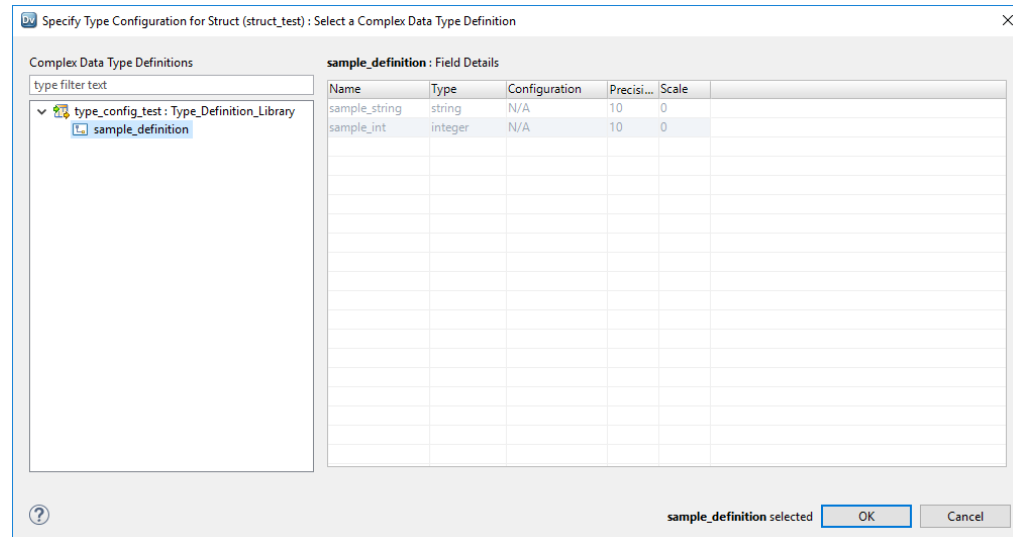
Specifying the Type Configuration for a Struct Port

The type configuration for a struct port requires a complex data type definition in the type definition library. Specify the complex data type definition that represents the schema of the struct data.

Before you specify the type configuration for the struct port, create or import a complex data type definition. You can also reference an existing complex data type definition in the type definition library.

1. In the transformation, select the struct port for which you want to specify the type configuration.
2. In the **Type Configuration** column, click the **Open** button.

The **Type Configuration** dialog box appears:



3. Select a complex data type definition from the list of definitions in the type definition library for the mapping or the mapplet.
4. Click **OK**.

On the **Ports** tab, the specified type configuration appears in the **Type Configuration** column of the struct port.

Complex Operators

A complex operator is a type of operator to access elements in a complex data type. The transformation language includes complex operators for complex data types. Use complex operators to access or extract elements in hierarchical data.

The following table lists the complex operators:

Complex Operators	Description
[]	<p>Subscript operator. Use to access array and map elements.</p> <p>Syntax: <code>array[index]</code></p> <ul style="list-style-type: none">- <code>array</code> is the array from which you want to access an element.- <code>index</code> is the position of an element within an array. <p>For example, use <code>[0]</code> to access the first element in an one-dimensional array.</p> <p>Syntax: <code>map[key]</code></p> <ul style="list-style-type: none">- <code>map</code> is the map from which you want to access the value corresponding to a key.- <code>key</code> is the map key for which you want to retrieve the map value.
.	<p>Dot operator. Use to access struct elements.</p> <p>Syntax: <code>struct.element_name</code></p> <ul style="list-style-type: none">- <code>struct</code> is the struct from which you want to access an element.- <code>element_name</code> is the name of the struct element.

Use complex operators in expressions to convert hierarchical data to relational data.

For more information about the operator syntax, return values, and examples, see the *Informatica Developer Transformation Language Reference*.

Complex Operator Examples

- Convert array data in a JSON file to relational data.

A JSON file has a `quarterly_sales` column that is of the array type with integer elements. You develop a mapping to transform the hierarchical data to relational data. The array port `quarterly_sales` contains four elements. Use subscript operators in expressions to extract the four elements into the four integer ports, `q1_sales`, `q2_sales`, `q3_sales`, and `q4_sales`.

The expressions for the integer output port are as follows:

```
quarterly_sales[0]
quarterly_sales[1]
quarterly_sales[2]
quarterly_sales[3]
```

- Convert struct data in a Parquet file to relational data.

A Parquet file has an `address` column that contains customer address as a struct type. You develop a mapping to transform hierarchical data to relational data. The struct port `address` contains three elements `city`, `state`, and `zip`, of type string. Use dot operators in expressions to extract the three struct elements into the three string ports, `city`, `state`, and `zip`.

The expressions for the string output port are as follows:

```
address.city
address.state
address.zip
```

Extracting an Array Element Using a Subscript Operator

Use the subscript operator in expressions to extract one or more elements of an array.

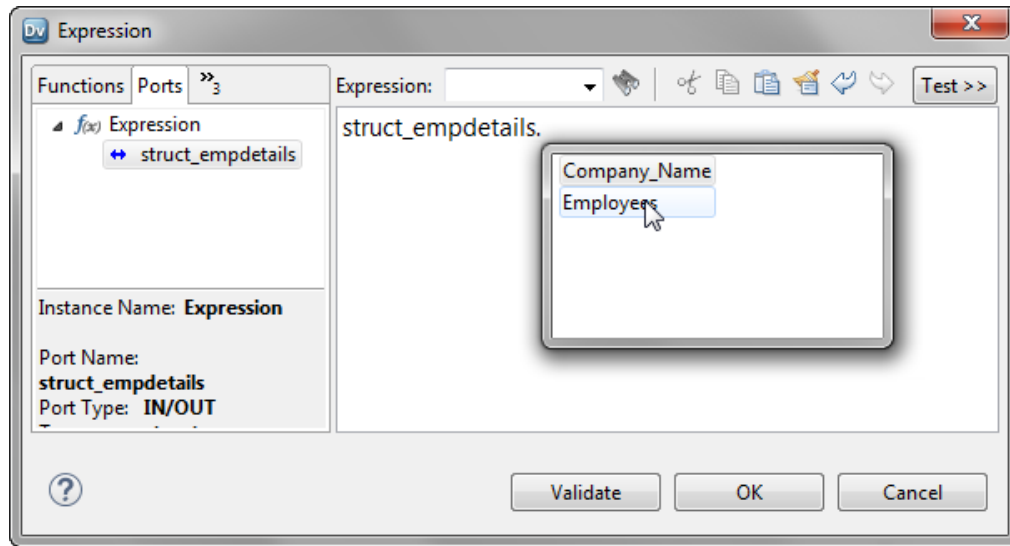
1. In the transformation, create an output port of the same data type as the return value of the array element that you want to extract.
To extract more than one element in an array, create an array output port. The data type of the elements in the type configuration for the array port must match the data type of the return value of the array elements.
2. In the **Expression** column for the output port, click the **Open** button.
The **Expression Editor** opens.
3. Delete the existing expression in the editor.
4. Click the **Ports** tab and select an array port.
5. Enter one or more index values within the subscript operator.
To extract array elements in a multidimensional array, enter index values within each subscript operator.
6. Click **Validate** to validate the expression.
7. Click **OK** to exit the **Validate Expression** dialog box.
8. Click **OK** to exit the **Expression Editor**.

Extracting a Struct Element Using the Dot Operator

Use the dot operator in expressions to extract elements of a struct.

1. In the transformation, create an output port of the same data type as the return value of the struct element that you want to extract.
2. In the **Expression** column for the output port, click the **Open** button.
The Expression Editor opens.
3. Delete the existing expression in the editor.
4. Click the **Ports** tab and select the struct port.
5. Enter a dot after the port name.

The Developer tool displays a list of elements in the struct.



6. Select the element that you want to extract.
7. Click **Validate** to validate the expression.
8. Click **OK** to exit the Validate Expression dialog box.
9. Click **OK** to exit the Expression Editor.

Complex Functions

A complex function is a type of pre-defined function in which the value of the input or the return type is of a complex data type. The transformation language includes complex functions for complex data types. Use complex functions to generate and process hierarchical data.

The following table describes the complex functions for array data type:

Complex Function	Description
ARRAY (<i>element1</i> [<i>element2</i>] ...)	Generates an array with the specified elements. The data type of the argument determines the data type of the array.
COLLECT_LIST (<i>value</i>)	Returns an array with elements in the specified port. The data type of the argument determines the data type of the array. COLLECT_LIST is an aggregate function.
CONCAT_ARRAY ('', <i>array_of_strings</i>)	Concatenates string elements in an array based on a separator that you specify and returns a string.
SIZE (<i>array</i>)	Returns the size of the array.

The following table describes the complex functions for map data type:

Complex Function	Description
MAP (<i>key1</i> , <i>value1</i> [, <i>key2</i> , <i>value2</i>] ...)	Generates a map with elements based on the specified key-value pair. The data type of arguments determines the data type of the map elements.
MAP_FROM_ARRAYS (<i>key_array</i> , <i>value_array</i>)	Generates a map from the specified key and value arrays. The data type of arguments determines the data type of the map elements.
MAP_KEYS (<i>map</i>)	Returns an array of keys for the specified map.
MAP_VALUES (<i>map</i>)	Returns an array of values for the specified map.
COLLECT_MAP (<i>value</i>)	Returns a map with elements based on the specified arguments. The data type of the argument determines the data type of the map. COLLECT_LIST is an aggregate function.
SIZE (<i>map</i>)	Returns the size of the map.

The following table describes the complex functions for struct data type:

Complex Function	Description
STRUCT_AS (<i>type_definition</i> , <i>struct</i>)	Generates a struct with a schema based on the specified complex data type definition and the values you pass as argument.
STRUCT (<i>element_name1</i> , <i>value1</i>)	Generates a struct with the specified names and values for the elements. The data type of the value is based on the specified value argument.
RESPEC (<i>type_definition</i> , <i>struct</i>)	Renames each element of the given struct value based on the names of the elements in the specified complex data type definition.
CAST (<i>type_definition</i> , <i>struct</i>)	Changes the data type and the name of each element for the given struct value based on the corresponding data type and name of the element in the specified complex data type definition.

For more information about the function syntax, return value, and examples, see the *Informatica Developer Transformation Language Reference*.

CHAPTER 7

Configuring Transformations to Process Hierarchical Data

This chapter includes the following topics:

- [Hierarchical Data Conversion, 114](#)
- [Convert Relational or Hierarchical Data to Struct Data, 115](#)
- [Convert Relational or Hierarchical Data to Nested Struct Data, 117](#)
- [Extract Elements from Hierarchical Data, 125](#)
- [Flatten Hierarchical Data, 127](#)

Hierarchical Data Conversion

In the Developer tool, use **Hierarchical Conversion** wizards that simplify the tasks for developing a mapping to process hierarchical data on the Spark engine. These wizards add one or more transformations to the mapping and create output ports. The expressions in the output ports use complex functions and operators to convert relational data to hierarchical data or hierarchical data to relational data.

The following table lists the hierarchical conversion wizards and when to use these wizards:

Hierarchical Conversion Wizard	When to Use
Create Struct Port	To convert relational and hierarchical data to struct data. For example, you want to convert data in one or more columns of any data type to hierarchical data column of type struct.
Create Nested Complex Port	To convert relational and hierarchical data in two tables to a nested struct data. You can select one or more columns of any data type. For example, you want to convert data in columns from two tables to a hierarchical data column of type struct with a nested schema. The wizard creates a struct with an array of structs.

Hierarchical Conversion Wizard	When to Use
Extract from Complex Port	To convert hierarchical data to relational data or modify the hierarchical data. For example, you want to perform the following conversions: <ul style="list-style-type: none"> - Convert one or more elements of primitive data types in a hierarchical data to relational data columns. - Convert one or more elements of complex data types in a hierarchical data to hierarchical data columns.
Flatten Complex Port	To convert hierarchical data to relational data. For example, you want to convert one or more elements of a hierarchical data to relational data columns.

Convert Relational or Hierarchical Data to Struct Data

You can convert relational or hierarchical data in one or more columns to hierarchical data of type struct. Use the **Create Struct Port** wizard in the Developer tool to perform the conversion.

For example, a relational table contains three columns city, state, and zip. You create a mapping to convert the data in the three columns to one hierarchical column. Select the three ports in the transformation and use the Create Struct Port wizard to generate struct data with the selected ports as its elements.

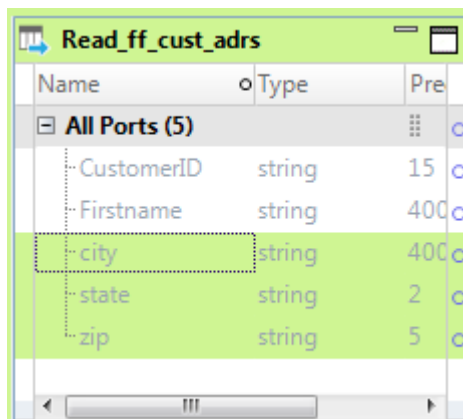
The wizard performs the following tasks:

- Creates a complex data type definition based on the ports that you select.
- Adds an Expression transformation to the mapping.
- Creates a struct output port to represent the struct data.
- Creates an expression that uses the STRUCT_AS function to generate struct data.

Creating a Struct Port

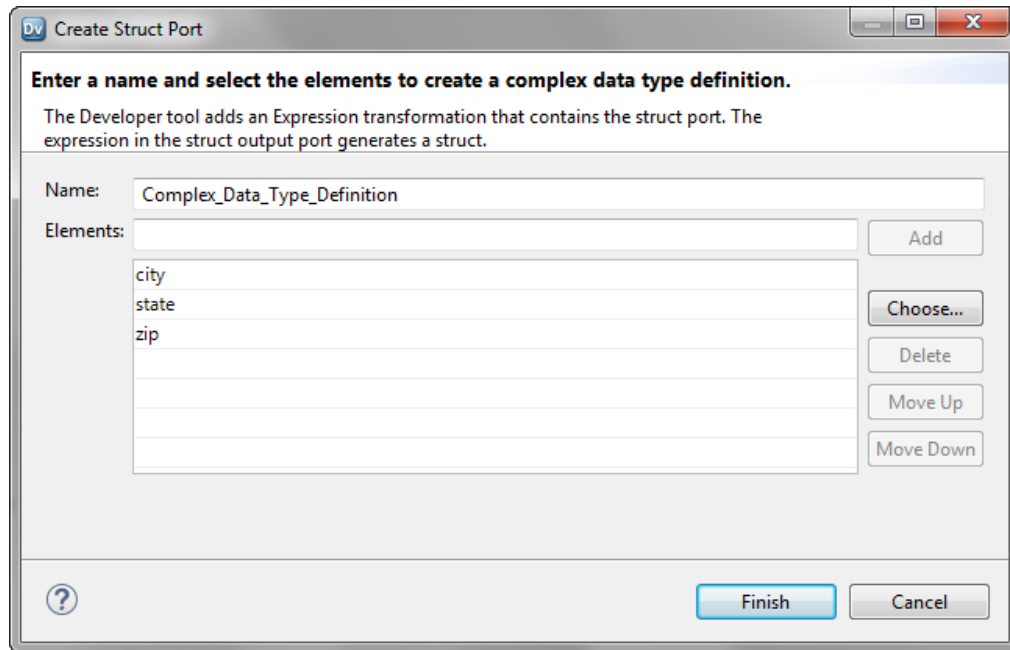
Use the **Create Struct Port** wizard to convert data that passes through one or more ports to struct data.

1. In the transformation, select one or more ports that you want to convert as elements of the struct data.
The ports you select also determine the elements of the complex data type definition.



2. Right-click the selected ports, and select **Hierarchical Conversions > Create Struct Port**.

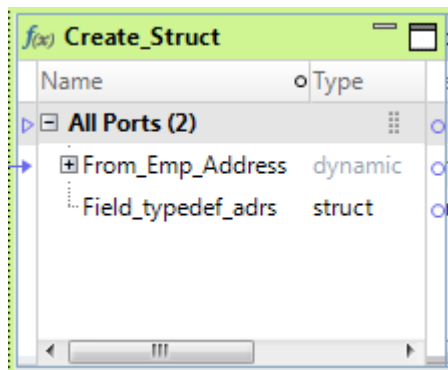
The **Create Struct Port** wizard appears with the list of ports that you selected.



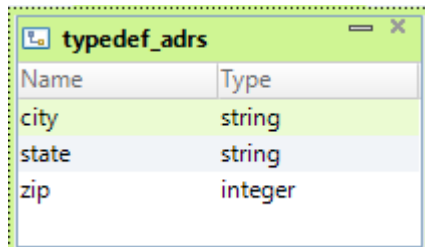
3. Optionally, in the **Name** box, change the name of the complex data type definition.
For example, typedef_address.
4. Optionally, click **Choose** to select other ports in the transformation.
5. Click **Finish**.

You can see the following changes in the mapping:

- The mapping contains a new Expression transformation `Create_Struct` with a struct output port and a dynamic port with ports from the upstream transformation.



- The type definition library contains the new complex data type definition.



- The struct output port references the complex data type definition.
- The struct output port contains an expression with the STRUCT_AS function. For example,

```
STRUCT_AS(:Type.Type_Definition_Library.typedef_address,city,state,zip)
```

Convert Relational or Hierarchical Data to Nested Struct Data

You can convert relational or hierarchical data in one or more columns in two transformations to nested struct data. Use the **Create Nested Complex Port** wizard in the Developer tool to perform the conversion.

The wizard converts relational data from two tables to struct data with a nested data type definition.

For example, a relational table contains employee bank account details. Another table contains the bank details. You create a mapping to convert data in the two tables into a hierarchical format that the payment system can process. Select the ports in the two transformations and use the Create Nested Complex Port wizard. The wizard generates struct data that contains an array element for each bank. The array contains a struct element for the employee bank account details.

Relational input

The Emp_Details table contains the following columns: employee ID, name, address, bank id, and bank account number.

The Bank_Details table contains the following columns: bank id, bank name, SWIFT code.

Struct output

```
banks{
  bank_swift:integer
  [bank_name{account_number:integer,employee_id:integer,name:string,address:string}]
}
```

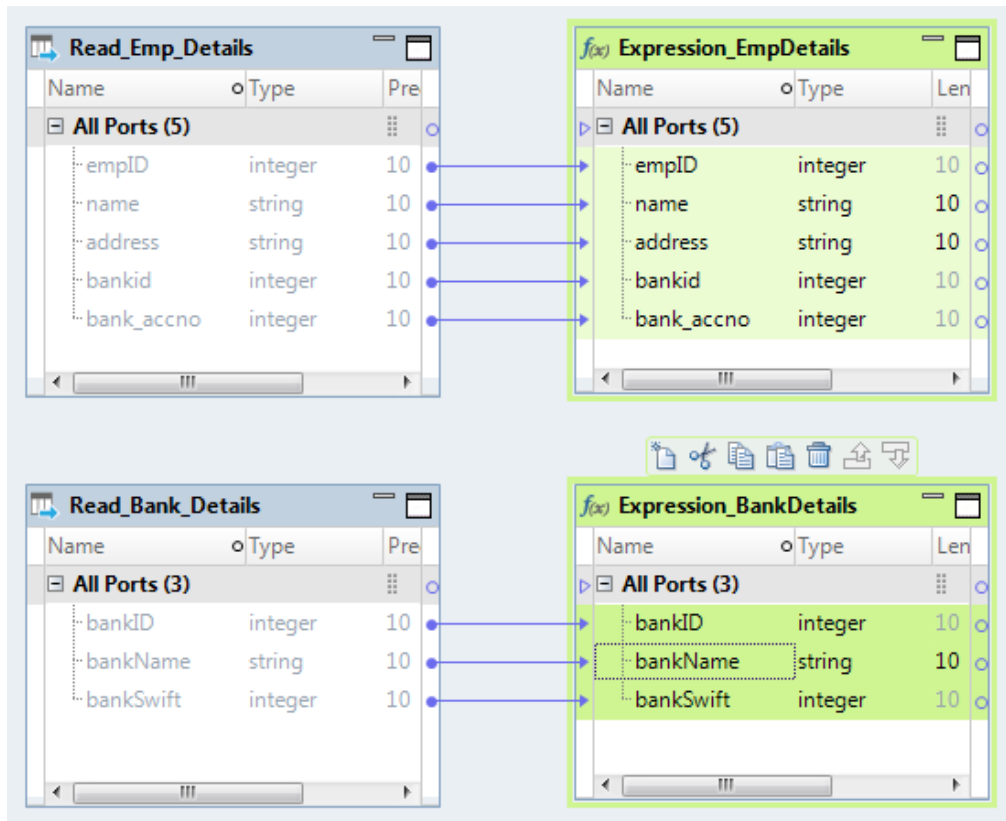
The wizard performs the following tasks:

- Adds a Joiner transformation to join source data from the parent and child tables based on the join keys. The transformation that contains ports for the parent complex data type definition is the parent transformation, and the other transformation is the child transformation.
- Creates a child complex data type definition and adds an Expression transformation that contains the child struct port. The expression in the struct port generates a struct based on the child complex data type definition.
- Adds an Aggregator transformation that contains an array port. The expression in the array port generates an array with the child struct as its element, and groups the struct values based on the group by port.
- Creates a parent complex data type definition and adds an Expression transformation that contains the nested complex port. The expression in the nested complex port generates a struct with an array of structs as one of its elements.

Creating A Nested Complex Port

Use the **Create Nested Complex Port** wizard to convert data that passes through one or more ports in two transformations to a struct data that references a nested data type definition. You specify the transformations from which the wizard creates the parent and child complex data type definitions.

1. Open the mapping or maplet in which you want to create a nested complex port.
2. Press Ctrl and select the ports in the two transformations that you want to convert as elements of the nested complex port.



3. Right-click the selected ports, and select **Hierarchical Conversions > Create Nested Complex Port**.

The **Create Nested Complex Port** wizard appears.

4. Choose the transformation that contains ports for the parent complex data type definition and select the join keys to join the tables.

5. Click **Next**.

The wizard page to create child complex data type definition appears.

Create Nested Complex Port

Enter a name and select the elements to create the child complex data type definition.

The Developer tool adds an Expression transformation that contains the child struct port. The expression in the output port generates a struct.

Name:

Elements:

empID	
name	
address	
bankid	
bank_accno	

Buttons: Add, Choose..., Delete, Move Up, Move Down

Navigation: ?

6. Optionally, change the name of the child complex data type definition and make any changes to the elements.

Create Nested Complex Port

Enter a name and select the elements to create the child complex data type definition.

The Developer tool adds an Expression transformation that contains the child struct port. The expression in the output port generates a struct.

Name:

Elements:

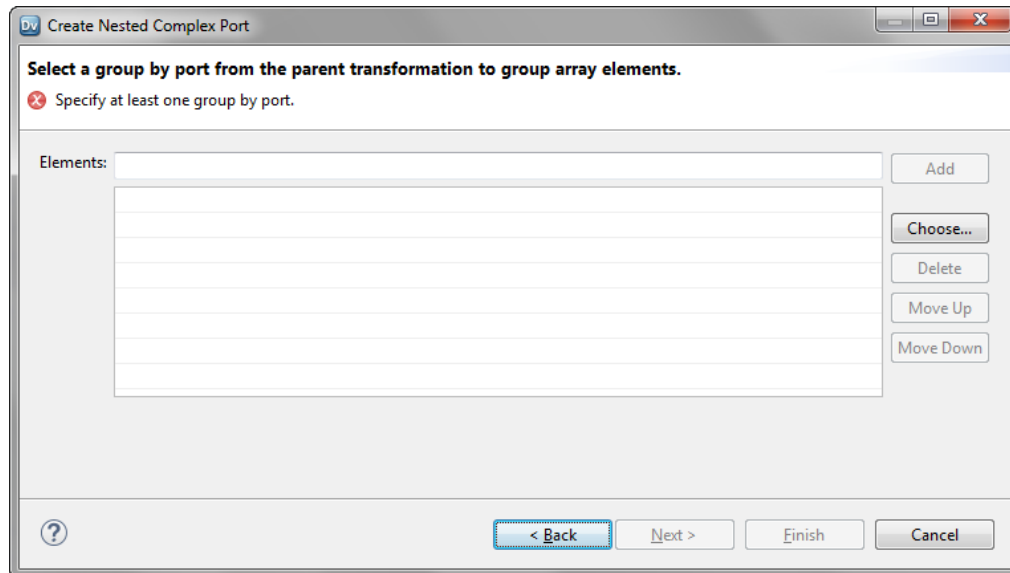
empID	
name	
address	
bank_accno	

Buttons: Add, Choose..., Delete, Move Up, Move Down

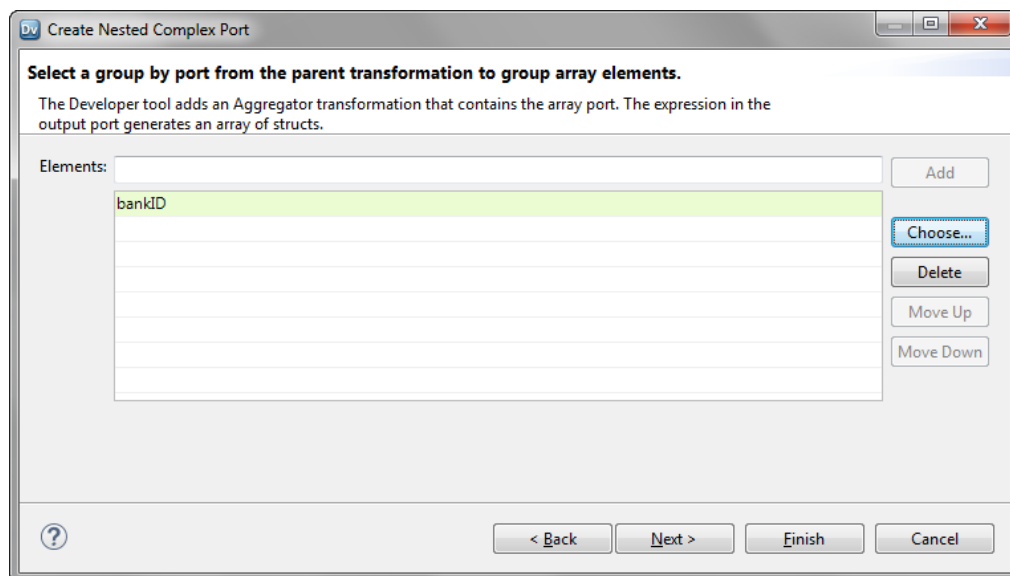
Navigation: ?

7. Click **Next**.

The wizard page to select group by port from the parent transformation appears.



8. Click **Choose**.
The **Ports** dialog box appears.
9. Select a group by port from the parent transformation and click **OK**.
The selected port appears on the wizard page.



10. Click **Next**.

The final wizard page to create parent complex data type definition and nested complex port appears.

Enter a name and select the elements to create the parent complex data type definition.

The Developer tool adds an Expression transformation that contains the nested complex port, which is a struct port with an array of structs. The expression in the output port generates a struct.

Name:

Elements:

bankID	
bankName	
bankSwift	

Buttons: Add, Choose..., Delete, Move Up, Move Down

Bottom buttons: < Back, Next >, Finish, Cancel

11. Optionally, change the name of the parent complex data type definition and make any changes to the elements.

You cannot delete the port that you selected as the group by port.

Enter a name and select the elements to create the parent complex data type definition.

The Developer tool adds an Expression transformation that contains the nested complex port, which is a struct port with an array of structs. The expression in the output port generates a struct.

Name:

Elements:

bankID	
bankName	
bankSwift	

Buttons: Add, Choose..., Delete, Move Up, Move Down

Bottom buttons: < Back, Next >, Finish, Cancel

Cannot delete the element because it is used as the group by port.

12. Click **Finish**.

You can see the following changes in the mapping:

- The mapping contains a new Joiner transformation `Join_<transformation1>_<transformation1>` that joins the two tables.

Name	Type
Output (2)	
From_Expression_BankDetails	dynamic
From_Expression_EmpDetails	dynamic
Master (1)	
From_Expression_BankDetails	dynamic
Detail (1)	
From_Expression_EmpDetails	dynamic

- The mapping contains a new Expression transformation `Create_ChildStruct` with a struct output port and a dynamic port with ports from the upstream transformation.

Name	Type
All Ports (2)	
From_Join_Expression_B...	dynamic
Field_typedef_empinfo	struct

- The type definition library contains the new child complex data type definition. The struct output port references the child complex data type definition.

Name	Type
empID	integer
name	string
address	string
bank_accno	integer

The output port contains an expression with the `STRUCT_AS` function:

```
STRUCT_AS(:Type.Type_Definition_Library.<child_typedef>,<comma_delimited_child_elements>)
```

- The mapping contains a new Aggregator transformation `Agg_ArrayOfStruct` with the group by port to group the child structs into an array.

Name	Type
All Ports (2)	
From_Create_ChildStruct	dynamic
Field_typedef_empinfo_List	array

- The mapping contains a new Expression transformation `Final_NestedComplexPort` with a struct output port and a dynamic port with ports from the upstream transformation.

Name	Type
All Ports (2)	
From_Agg_ArrayOfStruct	dynamic
Field_typedef_bank	struct

- The type definition library contains the new parent complex data type definition, which is a nested data type definition. The struct output port references the parent complex data type definition.

Name	Type
bankID1	integer
bankName	string
bankSwift	integer
Field_typedef_empinfo_List	array

Name	Type
empID	integer
name	string
address	string
bank_accno	integer

The output port contains an expression with the `STRUCT_AS` function:

```
STRUCT_AS(:Type.Type_Definition_Library.<parent_typedef>,<comma_delimited_struct_elements>)
```

Extract Elements from Hierarchical Data

You can extract elements of a primitive or complex data type from hierarchical data. Use the **Extract from Complex Port** wizard in the Developer tool to perform the conversion.

Based on the data type of the elements in the complex port that you select, the wizard performs the following conversions:

- If the elements are of primitive data types, the wizard converts the element to a relational data.
- If the elements are of complex data types, the wizard converts the element to a hierarchical data.

The wizard adds an Expression transformation that contains one or more extracted ports. The number of elements that you select to extract determine the number of output ports in the transformation. The data type of the element determines the data type of the output port. The expression in the output ports use complex operators to extract elements from the complex port.

The following table describes the wizard conversions based on the data type of the complex port that you select:

Complex Port Data Type	Wizard Conversion
array	Relational data. If you specify an array index, the wizard extracts an element in the array. Default is 0. The wizard extracts the first element of the array. Hierarchical data. If you do not specify an array index, the wizard extracts the entire array.
struct	Relational data. The wizard extracts the element in the struct that you selected.
array of structs	Relational data. If you select an element in the struct, the wizard extracts the element. The wizard requires an array index value. Default is 0. Hierarchical data. If you select the array and specify an array index, the wizard extracts the struct element in the array. If you do not specify an array index, the wizard extracts the entire array.
array of maps	Relational data. If you select the key or value element in the array and specify an array index, the wizard extracts the element. Hierarchical data. If you select the array and specify an array index, the wizard extracts the map element in the array. If you do not specify an array index, the wizard extracts the entire array.
struct of structs	Relational data. If you select an element in the parent or the child struct, the wizard extracts the element. Hierarchical data. If you select the parent struct or child struct, the wizard extracts that struct.
struct of maps	Hierarchical data. If you select the map element, the wizard extracts the map data. If you select the struct, the wizard extracts the struct data in another port.

Extracting Elements from a Complex Port

Use the **Extract from Complex Port** wizard to extract elements from hierarchical data.

1. In the transformation, select a complex port from which you want to extract elements.
2. Right-click the selected ports, and select **Hierarchical Conversions > Extract from Complex Port**.

The **Extract from Complex Port** wizard appears with the list of elements in the complex port.

Name	Type	Type Configuration	Select Elements	Array Index
company	struct	(Company)	<input type="checkbox"/>	
Company_...	string	N/A	<input type="checkbox"/>	
Employees	array	(Employee) []	<input type="checkbox"/>	
Name	string	N/A	<input type="checkbox"/>	
Emp_Addr...	array	(Address) []	<input type="checkbox"/>	
City	string	N/A	<input type="checkbox"/>	
State	string	N/A	<input type="checkbox"/>	
Zip	string	N/A	<input type="checkbox"/>	
Age	string	N/A	<input type="checkbox"/>	
Emp_Cont...	struct	(Contact)	<input type="checkbox"/>	
Email	string	N/A	<input type="checkbox"/>	
Phone	string	N/A	<input type="checkbox"/>	

3. In the **Select Elements** column, select the check box for each element that you want to extract.
4. To extract an array element, specify an array index in the **Array Index** column.

If you delete the array index, the wizard extracts the entire array. To extract an element in the struct from a complex port for an array of structs, you must specify an array index value.

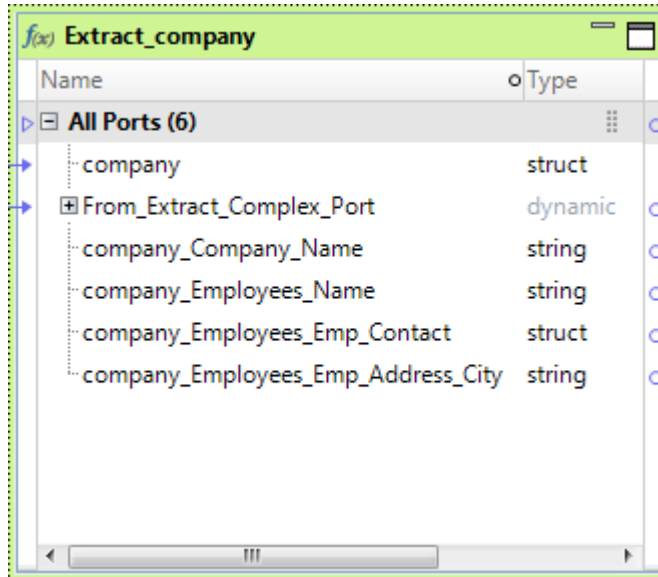
Name	Type	Type Configuration	Select Elements	Array Index
company	struct	(Company)	<input type="checkbox"/>	
Company_...	string	N/A	<input checked="" type="checkbox"/>	
Employees	array	(Employee) []	<input checked="" type="checkbox"/>	[0]
Name	string	N/A	<input checked="" type="checkbox"/>	
Emp_Addr...	array	(Address) []	<input checked="" type="checkbox"/>	[0]
City	string	N/A	<input checked="" type="checkbox"/>	
State	string	N/A	<input type="checkbox"/>	
Zip	string	N/A	<input type="checkbox"/>	
Age	string	N/A	<input checked="" type="checkbox"/>	
Emp_Cont...	struct	(Contact)	<input checked="" type="checkbox"/>	
Email	string	N/A	<input type="checkbox"/>	
Phone	string	N/A	<input type="checkbox"/>	

5. Click **Finish**.

You can see the following changes in the mapping:

- The mapping contains a new Expression transformation `Extract_<complex_port_name>` with the following ports:
 - The complex port from which you want to extract elements as the input port.
 - One or more output ports for the extracted elements. The number of elements that you selected to extract determines the number of output ports in the transformation.

- A dynamic port with ports from the upstream transformation.



- The output ports contain expressions that use complex operators to extract elements. The following image shows the expressions for the output ports in the Expression column on the Ports tab:

Search Ports									
	Name	Type	Type Configuration	Precisi...	Scale	Input	Output	Variable	Expression
1	company	struct	(Company)			<input checked="" type="checkbox"/>	<input type="checkbox"/>		company
2	From_Extract_Complex_Port	dynamic	N/A			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>		From_Extract_Complex_Port
3	company_Company_Name	string	N/A	10	0	<input type="checkbox"/>	<input checked="" type="checkbox"/>		company.Company_Name
4	company_Employees_Name	string	N/A	10	0	<input type="checkbox"/>	<input checked="" type="checkbox"/>		company.Employees[0].Name
5	company_Employees_Emp_Contact	struct	(Contact)			<input type="checkbox"/>	<input checked="" type="checkbox"/>		company.Employees[0].Emp_Contact
6	company_Employees_Emp_Address_City	string	N/A	10	0	<input type="checkbox"/>	<input checked="" type="checkbox"/>		company.Employees[0].Emp_Address[0].City

Flatten Hierarchical Data

You can flatten elements of hierarchical data into relational data. Use the **Flatten Complex Port** wizard in the Developer tool to perform the conversion.

The wizard converts hierarchical data to relational data. When you have hierarchical data with nested data type, you can select specific elements or all elements of complex data type to flatten.

Based on the data type of the complex port, the wizard performs the following tasks:

struct

- Adds an Expression transformation with flattened output ports. The expression for the output ports uses the dot operator to extract elements in the struct.
- Adds a final Expression transformation that contains a dynamic port with all ports from the upstream transformation including the flattened struct ports.

array

Adds a Normalizer transformation with flattened output ports. The wizard flattens the array field in the Normalizer view.

Adds a final Expression transformation that contains a dynamic port with all ports from the upstream transformation including the flattened array ports.

nested data type

Adds one or more Expression and Normalizer transformations with flattened output ports. If you select a child element of a nested complex port, the wizard flattens both the parent and child elements.

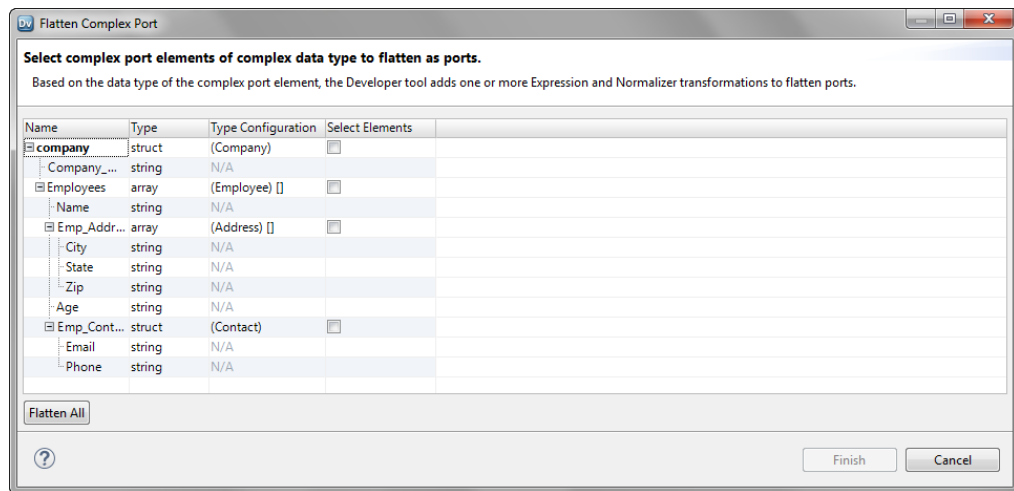
Adds a final Expression transformation that contains a dynamic port with all ports from the upstream transformation including the flattened ports.

Flattening a Complex Port

Use the **Flatten Complex Port** wizard to convert hierarchical data to relational data.

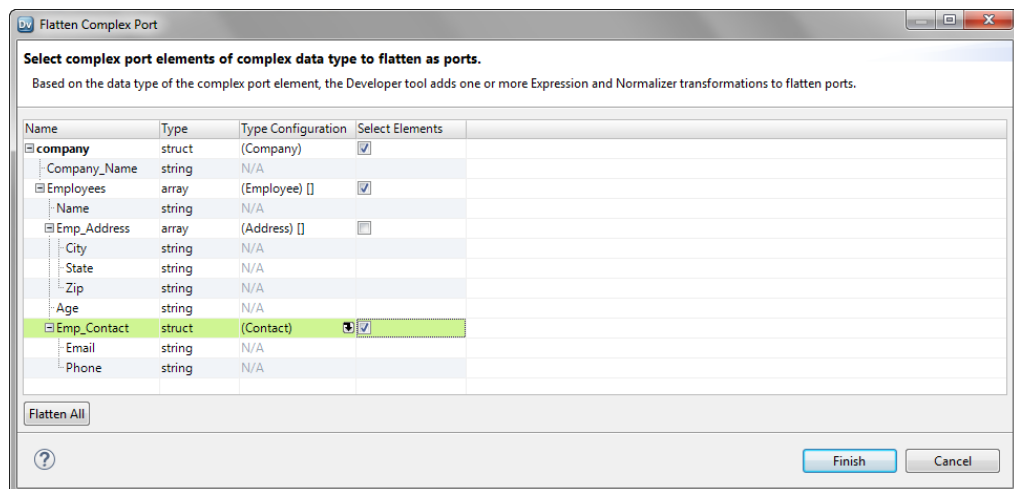
1. In the transformation, select a complex port that you want to flatten.
2. Right-click the selected ports, and select **Hierarchical Conversions > Flatten Complex Port**.

The **Flatten Complex Port** wizard appears with the list of elements in the complex port.



3. In the **Select Elements** column, select the check box for each element of a struct or an array data type that you want to extract.

If you select a child element of a nested complex port, the wizard automatically selects the parent element to flatten.



4. Optionally, click **Flatten All** to flatten all elements of struct or data type in the complex port.

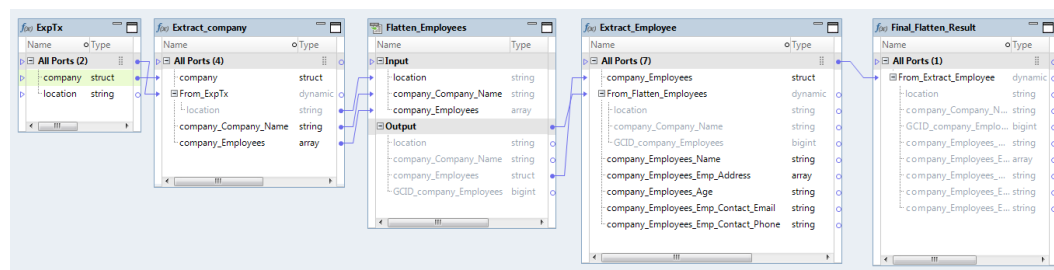
The wizard selects all elements of struct or data type in the complex port to flatten.

5. Click **Finish**.

You can see the following changes in the mapping:

- For each struct element that you selected to flatten, the mapping contains a new Expression transformation `Extract_<element_of_struct_type>` with the following ports:
 - The complex port from which you want to extract elements.
 - One or more output ports for the extracted elements.
 - The output ports contain expressions that use complex operators to extract elements.
 - A dynamic port with ports from the upstream transformation
- For each array element that you selected to flatten, the mapping contains a new Normalizer transformation `Flatten_<element_of_array_type>` with an input group and an output group. The output group contains the flattened normalized fields for the array element.
- A final Expression transformation `Final_Flatten_Result` with a dynamic port that contains the flattened ports for the complex port and other ports in the upstream transformation.

The following image shows an example of mapping changes:



CHAPTER 8

Processing Unstructured and Semi-structured Data with an Intelligent Structure Model

This chapter includes the following topics:

- [Processing Unstructured and Semi-structured Data with Intelligent Structure Model Overview, 130](#)
- [Intelligent Structure Discovery Process, 131](#)
- [Use Case, 131](#)
- [Using an Intelligent Structure Model in a Mapping, 132](#)
- [Rules and Guidelines for Intelligent Structure Models, 133](#)
- [How to Develop a Mapping to Process Data with an Intelligent Structure Model , 133](#)
- [Before You Start, 135](#)
- [Creating an Informatica Intelligent Cloud Services Account, 136](#)
- [Creating an Intelligent Structure Model, 136](#)
- [Exporting an Intelligent Structure Model, 137](#)
- [Checking for Data Loss, 137](#)

Processing Unstructured and Semi-structured Data with Intelligent Structure Model Overview

You can use CLAIRE™ Intelligent Structure Discovery to parse semi-structured or structured data in mappings that run on the Spark engine.

Long, complex files with little or no structure can be difficult to understand much less parse. CLAIRE™ Intelligent Structure Discovery can automatically discover the structure in unstructured data.

CLAIRE™ uses machine learning algorithms to decipher data in semi-structured or unstructured data files and create a model of the underlying structure of the data. You can generate an Intelligent structure model, a model of the pattern, repetitions, relationships, and types of fields of data discovered in a file, in Informatica Intelligent Cloud Services.

To use the model, you export it from Data Integration, and then can associate it with a data object in a Big Data Management mapping. You can run the mapping on the Spark engine to process the data. The mapping

uses the Intelligent structure model to extract and parse data from input files based on the structure expressed in the model.

Intelligent Structure Discovery Process

You can create a CLAIRE™ Intelligent structure model in Intelligent Structure Discovery. Intelligent Structure Discovery is a service in Data Integration.

When you provide a sample file, Intelligent Structure Discovery determines the underlying structure of the information and creates a model of the structure. After you create an Intelligent structure model you can view, edit, and refine it. For example, you can select to exclude or combine structure elements. You can normalize repeating groups.

When you finish refining the model, you can export it and then associate it with a data object in a Big Data Management mapping.

The following image shows the process by which Intelligent Structure Discovery deciphers the underlying patterns of data and creates a model of the data patterns.



You can create models for semi-structured data from Microsoft Excel, Microsoft Word tables, PDF forms, and CSV files, or unstructured text files. You can also create models for structured data such as XML and JSON files.

You can quickly model data for files whose structure is very hard, time consuming, and costly to find, such as log files, clickstreams, customer web access, error text files, or other internet, sensor, or device data that does not follow industry standards.

Use Case

You work in an operations group for an insurance company. Your team wants to process web logs to identify operations and security issues.

Your back-end system collects data on system access and security activities in various branches of the company. When the data is collected, the data is stored in the corporate data center and hosted in Amazon S3 storage.

Your team wants to understand the types of operations issues that have caused the most errors and system downtime in the past few weeks. You want to store data afterwards for auditing purposes.

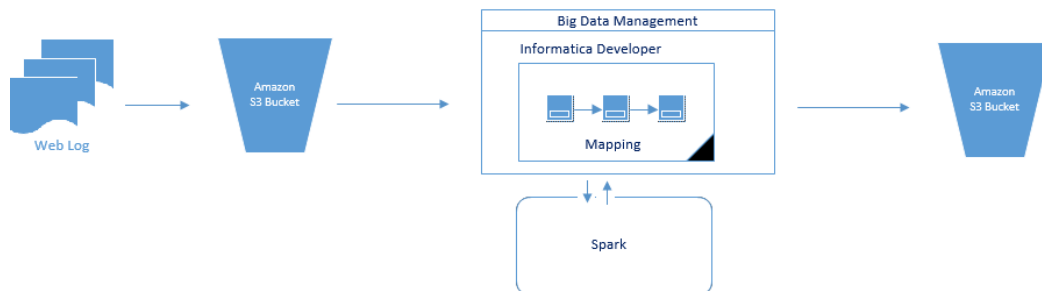
Before your data analysts can begin working with the data, you need to parse the data in Amazon S3 input buckets and produce actionable data. But you cannot spend the time and resources required to sift through the data to create models of analysis. You might have to develop numerous mappings and parameter sets to parse the data to make sure that actionable data is created from the weblogs.

Instead of manually creating individual transformations, your team can use automatically generated Intelligent structure models to determine the relevant data sets. You create an Intelligent structure model in Intelligent Structure Discovery, an application in Data Integration that uses machine learning algorithms to decipher data in structured or unstructured data files and discover the underlying structure of the data.

Intelligent Structure Discovery creates an Intelligent structure model that represents the input file data structure. You create a mapping with a data object that uses the intelligent structure model to output actionable data sets.

After the mapping fetches data from Amazon S3 input buckets, the mapping processes the data with an Intelligent structure model to prepare the data, and can then write the data to Amazon S3 output buckets.

The following image shows the process to fetch the data from Amazon S3 input buckets, parse and prepare the data, and then write the data to Amazon S3 output buckets. Analysts can then use the data to handle security issues and improve operations management.



Using an Intelligent Structure Model in a Mapping

To use an Intelligent structure model in a mapping, you associate it with a data object. The data object can then be used to process files that are similar to the one used to create the model.

You can add a Read transformation based on the data object to a mapping. If you want to process the data any further, such as checking data quality, or structuring the data into relational format, you add relevant downstream transformations to the mapping.

When the mapping runs, the Read transformation reads one or more input files and parses the data into fields, arrays, and structs.

Depending on the model and input, the data object output might contain primitive data types, complex data types, or nested data types. For more information about working with these data types in transformations, see [“Complex Data Types” on page 94](#).

Rules and Guidelines for Intelligent Structure Models

Consider the following rules and guidelines when you work with Intelligent structure model:

- If the Intelligent structure model does not match the input file, or only partially matches the input file, there might be a large amount of unidentified data and data loss. It is important therefore to choose a representative file to create the model.
- When you create an Intelligent structure model in Intelligent Structure Discovery, ensure that the names for output groups follow the Developer tool naming conventions.
 - To ensure compatibility, an element name must contain only English letters (A- Z, a-z), numerals (0-9), and underscores.
 - Do not use reserved logical terms such as And, Or, or other reserved names.
 - Do not start element names with a number.
 - Do not use duplicate names for different elements.
- Ensure that the Intelligent structure model is valid before you add it to the Column Projection properties of the data object properties.
- You can only use a Read transformation with an Intelligent structure model in a mapping. Do not create or use a Write transformation with an Intelligent structure model, as the mapping will fail.
- If the Intelligent structure model version is incompatible with the Big Data Management version, import the model in a compatible model version from Data Integration.
- A data object with an Intelligent structure model parse PDF forms, Microsoft Word tables, and XML files whose size is less than the supported Hadoop split size of 256 M.
- An Intelligent structure model can parse the data within PDF form fields but not data outside the fields.
- An Intelligent structure model can parse data within Microsoft Word tables. Other data is unidentified.

How to Develop a Mapping to Process Data with an Intelligent Structure Model

You can create a mapping with a data object that incorporates an Intelligent structure model to parse data. You run the mapping on the Spark engine to process the data.

Note: The tasks and the order in which you perform the tasks to develop the mapping depend on the mapping scenario.

The following list outlines the high-level tasks to develop and run a mapping to read and process data in files of any type that an Intelligent structure model can process, and then write the data to a target.

In Data Integration, create an Intelligent structure model.

Create an Intelligent structure model using a representative file. Export the .amodel file. Save the file locally or copy the file to the relevant file storage system.

For more information, see *Informatica Intelligent Cloud Services Mappings* at the following link:

<https://network.informatica.com/onlinehelp/IICS/prod/CDI/en/index.htm#page/hh-cloud-mappings/Mappings.html>.

In Big Data Management, create a connection.

Create a connection to access data in files that are stored in the relevant system. You can create the following types of connections that will work with the data objects that can incorporate an intelligent structure:

- Hadoop Distributed File System
- Amazon S3
- Microsoft Azure Blob

Create a data object with an Intelligent structure model to read and parse source data.

1. Create a data object with an Intelligent structure model to represent the files stored as sources. You can create the following types of data objects with an Intelligent structure model:
 - complex file
 - Amazon S3
 - Microsoft Azure Blob
2. Configure the data object properties.
3. In the read data object operation, enable the column file properties to project columns in the files as complex data types.

Create a data object to write data to a target.

1. Create a data object to write the data to target storage.
2. Configure the data object properties.

Create a mapping and add mapping objects.

1. Create a mapping.
2. Add a Read transformation based on the data object with the Intelligent structure model.
3. Based on the mapping logic, add other transformations that are supported on the Spark engine. Link the ports and configure the transformation properties based on the mapping logic.
4. Add a Write transformation based on the data object that passes the data to the target storage or output. Link the ports and configure the transformation properties based on the mapping logic.

Configure the mapping to run on the Spark engine.

Configure the following mapping run-time properties:

1. Select Hadoop as the validation environment and Spark as the engine.
2. Select Hadoop as the execution environment and select a Hadoop connection.

Validate and run the mapping on the Spark engine.

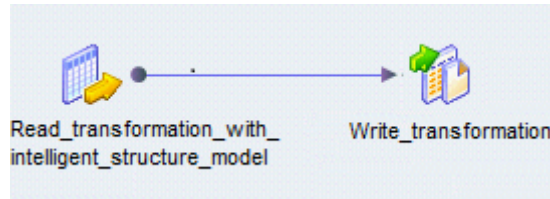
1. Validate the mapping and fix any errors.
2. Optionally, view the Spark engine execution plan to debug the logic.
3. Run the mapping.

Mapping Example

Your organization needs to analyze purchase order details such as customer ID, item codes, and item quantity. The purchase order details are stored in Microsoft Excel spreadsheets in HDFS. The data must be changed into text files for storage. In addition, private customer data must be removed. Create a mapping that reads all the purchase records from the file in HDFS using a data object with an intelligent structure. The mapping must parse the data and write it to a storage target.

You can use the extracted data for auditing.

The following figure shows the example mapping:



You can use the following objects in the HDFS mapping:

HDFS Input

The input object, `Read_transformation_with_intelligent_structure_model`, is a Read transformation that processed a Microsoft Excel file stored in HDFS and creates field output.

Amazon S3 Output

The output object, `Write_transformation`, is a Write transformation that represents an Amazon S3 bucket.

When you run the mapping, the Data Integration Service reads the file in a binary stream and passes it to the Read transformation. The Read transformation extracts the relevant data in accordance to the intelligent structure and passes data to the Write transformation. The Write transformation writes the data to the storage target.

You can configure the mapping to run in a Hadoop run-time environment.

Complete the following tasks to configure the mapping:

1. Create an HDFS connection to read files from the Hadoop cluster.
2. Create a complex file data object read operation. Specify the following parameters:
 - The intelligent structure as the resource in the data object. The intelligent structure was configured so that it does not pass sensitive data.
 - The HDFS file location.
 - The input file folder location in the read data object operation.
3. Drag and drop the complex file data object read operation into a mapping.
4. Create an Amazon S3 connection.
5. Create a Write transformation for the Amazon S3 data object and add it to the mapping.

Before You Start

Before you use Intelligent Structure Discovery to make a model, ensure that you have an active Informatica Intelligent Cloud Services account, have defined user roles through the Administrator, and have a relevant license.

A trial subscription includes access to Intelligent Structure Discovery without a license during the trial period.

For more information about registration and roles, refer to the *Informatica Intelligent Cloud Services Administrator Guide* at the following link:

<https://network.informatica.com/onlinehelp/IICS/prod/admin/en/index.htm>.

Creating an Informatica Intelligent Cloud Services Account

Create an Informatica Intelligent Cloud Services account if you do not already have one.

1. To create an account, access the Informatica Intelligent Cloud Services login page at the following link:
<https://dm-us.informaticacloud.com/identity-service/home>.
2. On the Informatica Intelligent Cloud Services login page, click **Don't have an account?**.
3. In the setup page, enter your contact details and email address. You can select to use your email address as a username.

Note: You will receive a registration email with access details and a link to confirm your account. You must confirm the account within 48 hours or register again.

4. In the registration email, click the account confirmation access link.
The Informatica Intelligent Cloud Services password setup window appears.
5. On the password setup window, define your Informatica Intelligent Cloud Services password, and then click **Log In**.

The Informatica Intelligent Cloud Services service picker appears. You can now access Data Integration.

Creating an Intelligent Structure Model

Create an Intelligent structure model in Data Integration. It is recommended to use a simplified sample file to generate the model. The sample file should have a representative data structure that is similar to the files that you want to parse.

For more information about creating and refining an Intelligent structure model, refer to the *Informatica Intelligent Cloud Services Mappings*.

1. In the Informatica Intelligent Cloud Services service picker, select **Data Integration**.
2. Click **New > Component > Intelligent Structure**.
3. In the **Intelligent Structures** page, click **New**.
4. In the **New Intelligent Structure** page, enter a name and description. You must provide a name for the Intelligent structure model.
5. To create an Intelligent structure model, browse for a sample file and click **Discover Structure**.
Intelligent Structure Discovery creates and displays an Intelligent structure model of the data.
6. To refine the Intelligent structure model and customize the structure of the output data, you can select a data node and select to combine, exclude, flatten, or collapse the node.
7. To save the Intelligent structure model, click **OK**.

Exporting an Intelligent Structure Model

After you create or edit a model, you can export the model to your local drive.

1. In Data Integration, select the **Explore** page.
2. On the **Explore** page, navigate to the project and folder with the Intelligent structure model.
The **Explore** page displays all the assets in the folder.
3. Click to select the row that contains the relevant Intelligent structure model. In the **Actions** menu, select **Edit**. In the **Intelligent Structure Details** panel, click **Edit**.
The Intelligent structure model appears in a separate page.
4. In the **Intelligent Structure** page, find the icon menu in the upper right corner of the Visual Model tab, and click the **Export model** icon.
Intelligent Structure Discovery downloads the `.amodel` file to your local drive.

Checking for Data Loss

To verify if a mapping with an Intelligent structure model lost data when it processed an input file, check the Spark history log.

- To get the log, run the following command on the Spark engine host machine:

```
yarn logs -applicationId<application ID>
```

CHAPTER 9

Stateful Computing on the Spark Engine

This chapter includes the following topics:

- [Stateful Computing on the Spark Engine Overview, 138](#)
- [Windowing Configuration, 139](#)
- [Window Functions, 142](#)
- [Windowing Examples, 147](#)

Stateful Computing on the Spark Engine Overview

You can use window functions in Expression transformations to perform stateful computations on the Spark engine.

A stateful computation is a function that takes some state and returns a value along with a new state.

You can use a window function to perform stateful computations. A window function takes a small subset of a larger data set for processing and analysis.

Window functions operate on a group of rows and calculate a return value for every input row. This characteristic of window functions makes it easy to express data processing tasks that are difficult to express concisely without window functions.

Use window functions to perform the following tasks:

- Retrieve data from upstream or downstream rows.
- Calculate a cumulative sum based on a group of rows.
- Calculate a cumulative average based on a group of rows.

Before you define a window function in an Expression transformation, you must describe the window by configuring the windowing properties. Windowing properties include a frame specification, partition keys, and order keys. The frame specification states which rows are included in the overall calculation for the current row. The partition keys determine which rows are in the same partition. The order keys determine how rows in a partition are ordered.

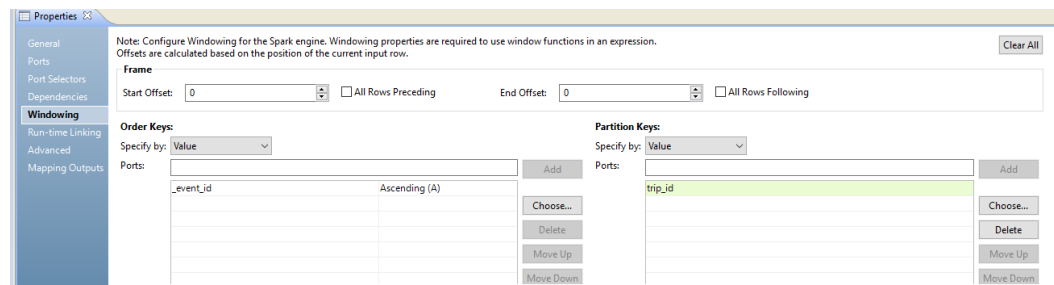
After you configure windowing properties, you define a window function in the Expression transformation. Spark supports the window functions LEAD and LAG. You can also use aggregate functions as window functions in an Expression transformation.

Windowing Configuration

When you include a window function in an Expression transformation, you configure the windowing properties associated with the function. Windowing properties define the partitioning, ordering, and frame boundaries associated with a particular input row.

Configure a transformation for windowing on the Windowing tab.

The following image shows the Windowing tab:



You configure the following groups of properties on the Windowing tab:

Frame

Defines the rows that are included in the frame for the current input row, based on physical offsets from the position of the current input row.

You configure a frame if you use an aggregate function as a window function. The window functions LEAD and LAG reference individual rows and ignore the frame specification.

Partition Keys

Separate the input rows into different partitions. If you do not define partition keys, all rows belong to a single partition.

Order Keys

Define how rows in a partition are ordered. The ports you choose determine the position of a row within a partition. The order key can be ascending or descending. If you do not define order keys, the rows have no particular order.

Frame

The frame determines which rows are included in the calculation for the current input row, based on their relative position to the current row.

If you use an aggregate function instead of LEAD or LAG, you must specify a window frame. LEAD and LAG reference individual row sand ignore the frame specification.

The start offset and end offset describe the number of rows that appear before and after the current input row. An offset of "0" represents the current input row. For example, a start offset of -3 and an end offset of 0 describes a frame including the current input row and the three rows before the current row.

The following image shows a frame with a start offset of -1 and an end offset of 1:

	Type	Category	Revenue	
	Action	Video game	1000	
Current input row →	Arcade	Video game	1000	← 1 PRECEDING
	Sports	Video game	2000	
	Adventure	Video game	3000	← 1 FOLLOWING
	Strategy	Video game	4000	

For every input row, the function performs an aggregate operation on the rows inside the frame. If you configure an aggregate expression like SUM with the preceding frame, the expression calculates the sum of the values within the frame and returns a value of 6000 for the input row.

You can also specify a frame that does not include the current input row. For example, a start offset of 10 and an end offset of 15 describes a frame that includes six total rows, from the tenth to the fifteenth row after the current row.

Note: The start offset must be less than or equal to the end offset.

Offsets of **All Rows Preceding** and **All Rows Following** represent the first row of the partition and the last row of the partition. For example, if the start offset is All Rows Preceding and the end offset is -1, the frame includes one row before the current row and all rows before that.

The following figure illustrates a frame with a start offset of 0 and an end offset of All Rows Following:

	Genre	Recordings	Revenue	
	Jazz	233	5000	
	Gospel	214	1000	
Current input row →	Country	145	2000	
	Ethnic	154	9000	
	Pop	317	4000	
	Rock	237	2100	
	Classical	221	3200	
	EDM	153	950	
	Hip Hop	839	2300	
	Punk	415	7650	

All Rows Following
↓

Partition and Order Keys

Configure partition and order keys to form groups of rows and define the order or sequence of rows within each partition.

Use the following keys to specify how to group and order the rows in a window:

Partition keys

Configure partition keys to define partition boundaries, rather than performing the calculation across all inputs. The window function operates across the rows that fall into the same partition as the current row.

You can specify the partition keys by value or parameter. Select **Value** to use port names. Choose **Parameter** to use a sort key list parameter. A sort key list parameter contains a list of ports to sort by. If you do not specify partition keys, all the data is included in the same partition.

Order keys

Use order keys to determine how rows in a partition are ordered. Order keys define the position of a particular row in a partition.

You can specify the order keys by value or parameter. Select **Value** to use port names. Choose **Parameter** to use a sort key list parameter. A sort key list parameter contains a list of ports to sort by. You must also choose to arrange the data in ascending or descending order. If you do not specify order keys, the rows in a partition are not arranged in any particular order.

Example

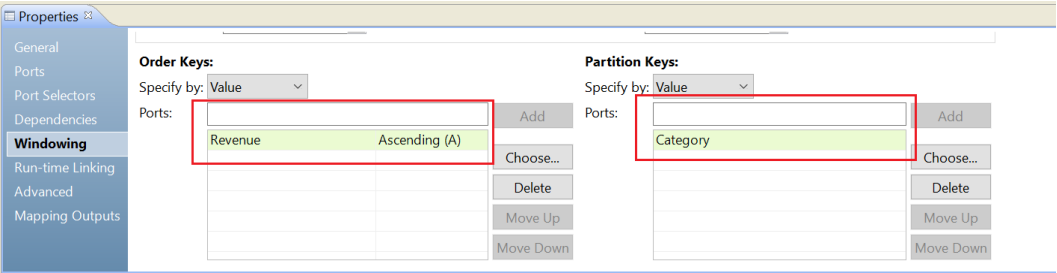
You are the owner of a coffee and tea shop. You want to calculate the best-selling and second best-selling coffee and tea products.

The following table lists the products, the corresponding product categories, and the revenue from each product:

Product	Category	Revenue
Espresso	Coffee	600
Black	Tea	550
Cappuccino	Coffee	500
Americano	Coffee	600
Oolong	Tea	250
Macchiato	Coffee	300
Green	Tea	450
White	Tea	650

You partition the data by category and order the data by descending revenue.

The following image shows the properties you configure on the Windowing tab:



The following table shows the data grouped into two partitions according to category. Within each partition, the revenue is organized in descending order:

Product	Category	Revenue
Espresso	Coffee	600
Americano	Coffee	600
Cappuccino	Coffee	500
Macchiato	Coffee	300
White	Tea	650
Black	Tea	550
Green	Tea	450
Oolong	Tea	250

Based on the partitioning and ordering specifications, you determine that the two best-selling coffees are espresso and Americano, and the two best-selling teas are white and black.

Rules and Guidelines for Windowing Configuration

Certain guidelines apply when you configure a transformation for windowing.

Consider the following rules and guidelines when you define windowing properties for a window function:

- When you configure a frame, the start offset must be less than or equal to the end offset. Otherwise, the frame is not valid.
- Configure a frame specification if you use an aggregate function as a window function. LEAD and LAG operate based on the offset value and ignore the frame specification.
- You cannot use complex ports as partition or order keys.
- You cannot preview the data in a transformation configured for windowing.
- Assign unique port names to partition and order keys to avoid run-time errors.
- The partition and order keys cannot use both a dynamic port and one or more generated ports of the same dynamic port. You must select either the dynamic port or the generated ports.

Window Functions

Window functions calculate a return value for every input row of a table, based on a group of rows.

A window function performs a calculation across a set of table rows that are related to the current row. You can also perform this type of calculation with an aggregate function. But unlike regular aggregate functions, a window function does not group rows into a single output row. The rows retain unique identities.

You can define the LEAD and LAG analytic window functions in an Expression transformation. LEAD and LAG give access to multiple rows within a table, without the need for a self-join.

LEAD

The LEAD function returns data from future rows.

LEAD uses the following syntax:

```
LEAD ( Column name, Offset, Default )
```

LEAD returns the value at an offset number of rows after the current row. Use the LEAD function to compare values in the current row with values in a following row. Use the following arguments with the LEAD function:

- Column name. The column name whose value from the subsequent row is to be returned.
- Offset. The number of rows following the current row from which the data is to be retrieved. For example, an offset of "1" accesses the next immediate row, and an offset of "3" accesses the third row after the current row.
- Default. The default value to be returned if the offset is outside the scope of the partition. If you do not specify a default, the default is NULL.

For more information about the LEAD function, see the *Informatica Transformation Language Reference*.

LAG

The LAG function returns data from preceding rows.

LAG uses the following syntax:

```
LAG ( Column name, Offset, Default )
```

LAG returns the value at an offset number of rows before the current row. Use the LAG function to compare values in the current row with values in a previous row. Use the following arguments with the LAG function:

- Column name. The column name whose value from the prior row is to be returned.
- Offset. The number of rows preceding the current row from which the data is to be retrieved. For example, an offset of "1" accesses the previous row, and an offset of "3" accesses the row that is three rows before the current row.
- Default. The default value to be returned if the offset is outside the scope of the partition. If you do not specify a default, the default is NULL.

For more information about the LAG function, see the *Informatica Transformation Language Reference*.

Aggregate Functions as Window Functions

In addition to LEAD and LAG, you can also use aggregate functions as window functions. When you use aggregate functions like SUM and AVG as window functions, you can perform running calculations that are similar to the stateful functions MOVINGSUM, MOVINGAVG, and CUME. Window functions are more flexible than stateful functions because you can set a specific end offset.

To use an aggregate function as a window function, you must define a frame in the windowing properties. You define a frame to limit the scope of the calculation. The aggregate function performs a calculation across the frame and produces a single value for each row.

Example

You are a lumber salesperson who sold different quantities of wood over the past two years. You want to calculate a running total of sales quantities.

The following table lists each sale ID, the date, and the quantity sold:

Sale_ID	Date	Quantity
30001	2016-08-02	10
10001	2016-12-24	10
10005	2016-12-24	30
40001	2017-01-09	40
10006	2017-01-18	10
20001	2017-02-12	20

A SUM function adds all the values and returns one output value. To get a running total for each row, you can define a frame for the function boundaries.

The following image shows the frame you specify on the Windowing tab:

Note: Configure Windowing for the Spark engine. Windowing properties are required to use window functions in an expression. Offsets are calculated based on the position of the current input row.

Frame

Start Offset: 0 ☒ All Rows Preceding End Offset: 0 ☐ All Rows Following

Order Keys:

Specify by: Value

Ports:

Order Key	Order
Date	Ascending (A)

Partition Keys:

Specify by: Value

Ports:

Partition Key

You configure the following windowing properties:

- Start offset: All Rows Preceding
- End offset: 0
- Order Key: Date Ascending
- Partition Key: Not specified

You define the following aggregate function:

```
SUM (Quantity)
```

SUM adds the quantity in the current row to the quantities in all the rows preceding the current row. The function returns a running total for each row.

The following table lists a running sum for each date:

Sale_ID	Date	Quantity	Total
30001	2016-08-02	10	10
10001	2016-12-24	10	20
10005	2016-12-24	30	50
40001	2017-01-09	40	90

Sale_ID	Date	Quantity	Total
10006	2017-01-18	10	100
20001	2017-02-12	20	120

Aggregate Offsets

An aggregate function performs a calculation on a set of values inside a partition. If the frame offsets are outside the partition, the aggregate function ignores the frame.

If the offsets of a frame are not within the partition or table, the aggregate function calculates within the partition. The function does not return NULL or a default value.

For example, you partition a table by seller ID and you order by quantity. You set the start offset to -3 and the end offset to 4.

The following image shows the partition and frame for the current input row:



The frame includes eight total rows, but the calculation remains within the partition. If you define an AVG function with this frame, the function calculates the average of the quantities inside the partition and returns 18.75.

Nested Aggregate Functions

A nested aggregate function in a window function performs a calculation separately for each partition. A nested aggregate function behaves differently in a window function and an Aggregator transformation.

A nested aggregate function in an Aggregator transformation performs a calculation globally across all rows. A nested aggregate function in an Expression transformation performs a separate calculation for each partition.

For example, you configure the following nested aggregate function in an Aggregator transformation and an Expression transformation:

```
MEDIAN ( COUNT ( P1 ) )
```

You define P2 as the group by port in the Aggregator transformation.

The function takes the median of the count of the following set of data:

P1	P2
10	1
7	1
12	1
11	2
13	2
8	2
10	2

RETURN VALUE: 3.5

When you include nested aggregate functions in an Expression transformation and configure the transformation for windowing, the function performs the calculation separately for each partition.

You partition the data by P2 and specify a frame of all rows preceding and all rows following. The window function performs the following calculations:

1. COUNT (P1) produces one value for every row. COUNT returns the number of rows in the partition that have non-null values.
2. MEDIAN of that value produces the median of a window of values generated by COUNT.

The window function produces the following outputs:

P1	P2	Output
10	1	3
7	1	3
12	1	3
11	2	4
13	2	4
8	2	4
10	2	4

You can nest aggregate functions with multiple window functions. For example:

```
LAG ( LEAD( MAX( FIRST ( p1 ) ) )
```

Note: You can nest any number of the window functions LEAD and LAG, but you cannot nest more than one aggregate function within another aggregate function.

Rules and Guidelines for Window Functions

Certain guidelines apply when you use window functions on the Spark engine.

Consider the following rules and guidelines when you define window functions in a transformation:

- Specify a constant integer as the offset argument in a window function.
- Specify a default argument that is the same data type as the input value.
- You cannot specify a default argument that contains complex data type or a SYSTIMESTAMP argument.
- To use the LEAD and LAG window functions, you must configure partition and order keys in the windowing properties.
- To use an aggregate function as a window function in an Expression transformation, you must configure a frame specification in the windowing properties.

Windowing Examples

The examples in this section demonstrate how to use LEAD, LAG, and other aggregate functions as window functions in an Expression transformation.

Financial Plans Example

You are a banker with information about the financial plans of two of your customers. Each plan has an associated start date.

For each customer, you want to know the expiration date for the current plan based on the activation date of the next plan. The previous plan ends when a new plan starts, so the end date for the previous plan is the start date of the next plan minus one day.

The following table lists the customer codes, the associated plan codes, and the start date of each plan:

CustomerCode	PlanCode	StartDate
C1	00001	2014-10-01
C2	00002	2014-10-01
C2	00002	2014-11-01
C1	00004	2014-10-25
C1	00001	2014-09-01
C1	00003	2014-10-10

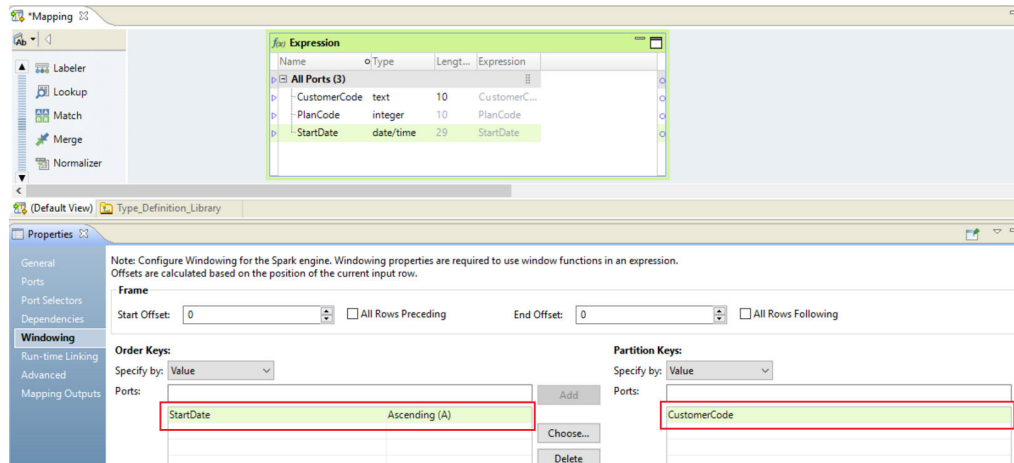
Define partition and order keys

You partition the data by customer code and order the data by ascending start date.

You configure the following windowing properties:

Property	Description
Order key	StartDate Ascending. Arranges the data chronologically by ascending start date.
Partition key	CustomerCode. Groups the rows according to customer code so that calculations are based on individual customers.
Frame	Not specified. Window functions access rows based on the offset argument and ignore the frame specification.

The following image shows the windowing properties you configure on the Windowing tab:



The following table lists the data grouped by customer code and ordered by start date:

CustomerCode	PlanCode	StartDate
C1	00001	2014-09-01
C1	00002	2014-10-01
C1	00003	2014-10-10
C1	00004	2014-10-25
C2	00001	2014-10-01
C2	00002	2014-11-01

The start dates for each customer are arranged in ascending order so that the dates are chronological.

Define a window function

You define a LEAD function to access the subsequent row for every input.

You define the following function on the Ports tab of the Expression transformation:

```
LEAD ( StartDate, 1, '01-Jan-2100' )
```

Where:

- `StartDate` indicates the target column that the function operates on.
- `1` is the offset. This value accesses the next immediate row.
- `01-Jan-2100` is the default value. The expression returns "01-Jan-2100" if the returned value is outside the bounds of the partition.

Define an `ADD_TO_DATE` function

You use an `ADD_TO_DATE` function to subtract one day from the date you accessed.

You define the following expression on the Ports tab of the Expression transformation:

```
ADD_TO_DATE ( LEAD ( StartDate, 1, '01-Jan-2100' ), 'DD', -1, )
```

By subtracting one day from the start date of the next plan, you find the end date of the current plan.

The following table lists the end dates of each plan:

CustomerCode	PlanCode	StartDate	EndDate
C1	00001	2014-09-01	2014-09-30
C1	00002	2014-10-01	2014-10-09
C1	00003	2014-10-10	2014-10-24
C1	00004	2014-10-25	2099-12-31*
C2	00001	2014-10-01	2014-10-31
C2	00002	2014-11-01	2099-12-31*

*The `LEAD` function returned the default value because these plans have not yet ended. The rows were outside the partition, so the `ADD_TO_DATE` function subtracted one day from 01-Jan-2100, returning 2099-12-31.

GPS Pings Example

Your organization receives GPS pings from vehicles that include trip and event IDs and a time stamp. You want to calculate the time difference between each ping and flag the row as skipped if the time difference with the previous row is less than 60 seconds.

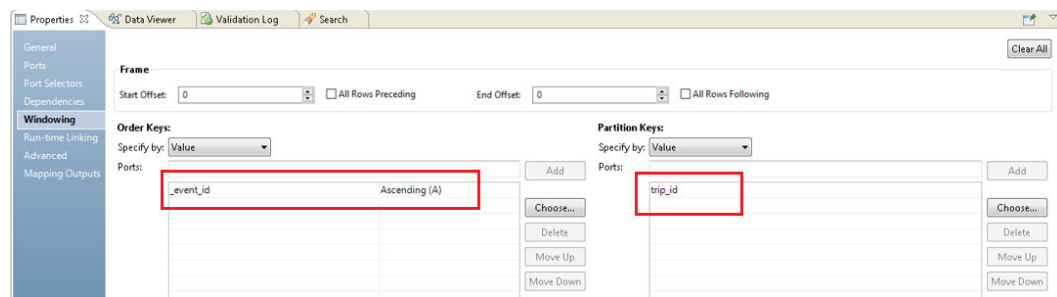
You order the events chronologically and partition the events by trip. You define a window function that accesses the event time from the previous row, and you use an `ADD_TO_DATE` function to calculate the time difference between the two events.

Windowing Properties

You define the following windowing properties on the Windowing tab:

Property	Description
Order key	<code>_event_id</code> Ascending. Arranges the data chronologically by ascending event ID.
Partition key	<code>trip_id</code> . Groups the rows according to trip ID so calculations are based on events from the same trip.
Frame	Not specified. Window functions access rows based on the offset argument and ignore the frame specification.

The following image shows the windowing properties you configure in the Expression transformation:



Window Function

You define the following LAG function to get the event time from the previous row:

```
LAG ( _event_time, 1, NULL )
```

Where:

- `_event_time` is the column name whose value from the previous row is to be returned.
- `1` is the offset. This value represents the row immediately before the current row.
- `NULL` is the default value. The function returns NULL if the return value is outside the bounds of the partition.

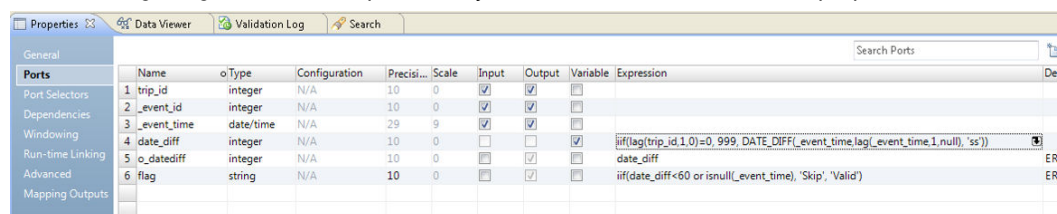
You define the following DATE_DIFF function to calculate the length of time between the two dates:

```
DATE_DIFF ( _event_time, LAG ( _event_time, 1, NULL ), 'ss' )
```

You flag the row as skipped if the DATE_DIFF is less than 60 seconds, or if the `_event_time` is NULL:

```
IF ( DATE_DIFF < 60 or ISNULL ( _event_time ), 'Skip', 'Valid' )
```

The following image shows the expressions you define in the transformation properties:



Output

The transformation produces the following outputs:

Trip ID	Event ID	Event Time	Time Difference	Flag
101	1	2017-05-03 12:00:00	NULL*	Skip
101	2	2017-05-03 12:00:34	34	Skip
101	3	2017-05-03 12:02:00	86	Valid
101	4	2017-05-03 12:02:23	23	Skip
102	1	2017-05-03 12:00:00	NULL*	Skip
102	2	2017-05-03 12:01:56	116	Valid
102	3	2017-05-03 12:02:00	4	Skip
102	4	2017-05-03 13:00:00	3480	Valid
103	1	2017-05-03 12:00:00	NULL*	Skip
103	2	2017-05-03 12:00:12	12	Skip
103	3	2017-05-03 12:01:12	60	Valid

*The rows preceding these rows are outside the bounds of the partition, so the transformation produced NULL values.

Aggregate Function as Window Function Example

You work for a human resources group and you want to compare each of your employees' salaries with the average salary in his or her department:

The following table lists the department names, the employee identification number, and the employee's salary:

Department	Employee	Salary
Development	11	5200
Development	7	4200
Development	9	4500
Development	8	6000
Development	10	5200
Personnel	5	3500
Personnel	2	3900
Sales	3	4800

Department	Employee	Salary
Sales	1	5000
Sales	4	4800

You set an unbounded frame to include all employees in the calculation, and you define an aggregate function to calculate the difference between each employee's salary and the average salary in the department.

Windowing Properties

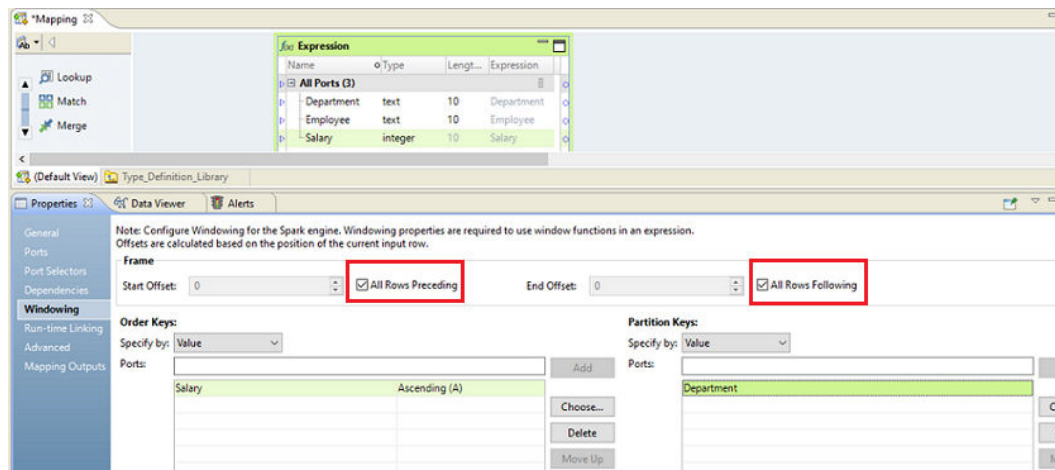
You define the following windowing properties on the Windowing tab:

Property	Description
Order key	Salary Ascending. Arranges the data by increasing salary.
Partition key	Department. Groups the rows according to department.
Start offset	All Rows Preceding
End offset	All Rows Following

With an unbounded frame, the aggregate function includes all partition rows in the calculation.

For example, suppose the current row is the third row. The third row is in the "Development" partition, so the frame includes the third row in addition to all rows before and after the third row in the "Development" partition.

The following image shows the windowing properties you configure in the Expression transformation:



Window Function

An aggregate function acts as a window function when you configure the transformation for windowing.

You define the following aggregate function to calculate the difference between each employee's salary and the average salary in his or her department:

```
Salary - AVG ( Salary ) = Salary_Diff
```


Output

The transformation produces the following salary differences:

Department	Employee	Salary	Salary_Diff
Development	11	5200	-820
Development	7	4200	-520
Development	9	4500	180
Development	8	6000	180
Development	10	5200	980
Personnel	5	3500	200
Personnel	2	3900	200
Sales	3	4800	-66
Sales	1	5000	-66
Sales	4	4800	134

You can identify which employees are making less or more than the average salary for his or her department. Based on this information, you can add other transformations to learn more about your data. For example, you might add a Rank transformation to produce a numerical rank for each employee within his or her department.

CHAPTER 10

Monitoring Mappings in the Hadoop Environment

This chapter includes the following topics:

- [Monitoring Mappings in the Hadoop Environment Overview, 154](#)
- [Hadoop Environment Logs, 155](#)
- [Blaze Engine Monitoring, 160](#)
- [Spark Engine Monitoring, 169](#)
- [Hive Engine Monitoring, 173](#)

Monitoring Mappings in the Hadoop Environment Overview

On the Monitor tab of the Administrator tool, you can view statistics and log events for mappings run in the Hadoop environment. The Monitor tab displays current and historical information about mappings run on Blaze, Spark, and Hive engines.

Use the Summary Statistics view to view graphical summaries of object state and distribution across the Data Integration Services. You can also view graphs of the memory and CPU that the Data Integration Services used to run the objects.

When you run a mapping in the Hadoop environment, the Data Integration Service generates log events. You can view log events relating to different types of errors such as Hadoop connection failures, Hive query failures, Hive command failures, or other Hadoop job failures.

When you run Hive mappings or Spark jobs that launches Hive tasks, you can select Tez or MapReduce engines and monitor statistics.

Hadoop Environment Logs

The Data Integration Service generates log events when you run a mapping in the Hadoop environment.

You can view logs for the Blaze engine, the Spark engine, or the MapReduce or Tez Hive engines. You can view log events relating to different types of errors such as Hadoop connection failures, Hive query failures, Hive command failures, or other Hadoop job failures.

When you run a mapping on the Spark engine, you can view the Scala code in logs that the Logical Data Translation Generator generates from the Informatica mapping.

You can view reject files in the reject file directory specified for the Data Integration Service.

YARN Web User Interface

You can view the applications that ran on a cluster in the YARN web user interface. Click the Monitoring URL for Blaze, Hive, or Spark jobs to access the YARN web user interface.

Blaze, Spark, and Hive engines run on the Hadoop cluster that you configure in the Hadoop connection. The YARN web user interface shows each job that the engine runs as a YARN application.

The following image shows the Application Monitoring page of the YARN web user interface:

The screenshot shows the YARN Web User Interface for the cluster 'psrlInfra.informatica.com:8088/cluster'. The page title is 'All Applications'. On the left, there is a navigation menu with options: Cluster, About, Nodes, Applications, NEW, SAVING, SUBMITTED, ACCEPTED, RUNNING, FINISHED, FAILED, KILLED, Scheduler, and Tools. The main content area displays 'Cluster Metrics' and 'User Metrics for dr.who'. Below these, there is a table of applications. The table has columns: ID, User, Name, Application Type, Queue, Start Time, Finish Time, State, Final Status, Running Containers, and Allocated CPU V-Cores. The table shows several applications, including 'application_1463379223882_0567' and 'application_1463379223882_0566', all in a 'FINISHED' state.

ID	User	Name	Application Type	Queue	Start Time	Finish Time	State	Final Status	Running Containers	Allocated CPU V-Cores
application_1463379223882_0567	Idmui	InfSprk0	SPARK	root.Idmui	Mon May 16 13:11:59 -0700 2016	Mon May 16 13:13:03 -0700 2016	FINISHED	SUCCEEDED	N/A	N/A
application_1463379223882_0566	Idmui	InfSprk0	SPARK	root.Idmui	Mon May 16 13:11:58 -0700 2016	Mon May 16 13:13:01 -0700 2016	FINISHED	SUCCEEDED	N/A	N/A
application_1463379223882_0565	Idmui	InfSprk0	SPARK	root.Idmui	Mon May 16 13:11:56 -0700 2016	Mon May 16 13:13:00 -0700 2016	FINISHED	SUCCEEDED	N/A	N/A
application_1463379223882_0564	Idmui	InfSprk0	SPARK	root.Idmui	Mon May 16 13:11:55 -0700 2016	Mon May 16 13:12:59 -0700 2016	FINISHED	SUCCEEDED	N/A	N/A

The **Application Type** indicates which engine submitted the YARN application.

The application ID is the unique identifier for the application. The application ID is a link to the application summary. The URL is the same as the Monitoring URL in the Administrator tool.

Click the **Logs** link in the application summary to view the application logs on the Hadoop cluster.

The amount of information in the application logs depends on the tracing level that you configure for a mapping in the Developer tool. The following table describes the amount of information that appears in the application logs for each tracing level:

Tracing Level	Messages
None	The log displays FATAL messages. FATAL messages include non-recoverable system failures that cause the service to shut down or become unavailable.
Terse	The log displays FATAL and ERROR code messages. ERROR messages include connection failures, failures to save or retrieve metadata, service errors.
Normal	The log displays FATAL, ERROR, and WARNING messages. WARNING errors include recoverable system failures or warnings.
Verbose initialization	The log displays FATAL, ERROR, WARNING, and INFO messages. INFO messages include system and service change messages.
Verbose data	The log displays FATAL, ERROR, WARNING, INFO, and DEBUG messages. DEBUG messages are user request logs.

Accessing the Monitoring URL

The Monitoring URL opens the Blaze Job Monitor web application or the YARN web user interface. Access the Monitoring URL from the **Execution Statistics** view in the Administrator tool.

1. In the **Monitor** tab of the Administrator tool, click the **Execution Statistics** view.
2. Select **Ad Hoc Jobs** or select a deployed mapping job or workflow from an application in the Navigator.
The list of jobs appears in the contents panel.

3. Select a mapping job and expand the mapping to select a grid task for the mapping.

The Monitoring URL appears in the **Properties** view.

The screenshot shows the 'Ad Hoc Jobs' window in the Administrator tool. It displays a table of jobs with columns: Name, Type, State, Job ID, Started By, Start Time, Elapsed Time, and End Time. The job 'MAINSESSION_task2' is selected. Below the table, the 'Properties' view for 'MAINSESSION_task2' is shown, including a status message 'This grid task is completed.' and a 'General Properties' section with details like Name, Type, Started By, User Security Domain, Start Time, Elapsed Time, End Time, % Task Completed, Monitoring URL, Incoming Task Dependencies, and Outgoing Task Dependencies.

Name	Type	State	Job ID	Started By	Start Time	Elapsed Time	End Time
PassThrough	Mapping	Completed	T2GJgmceE...	Administrator	09/29/2015 19:52:38	00:01:32	09/29/2015 19:54:10
POSTSES...	Command ...	Completed	T2GJgmceE...	Administrator	09/29/2015 19:53:41	00:00:17	09/29/2015 19:53:58
MAINSESSION_task2	Grid Task	Completed	T2GJgmceE...	Administrator	09/29/2015 19:52:57	00:00:43	09/29/2015 19:53:41
PRESESSI...	Command ...	Completed	T2GJgmceE...	Administrator	09/29/2015 19:52:38	00:00:02	09/29/2015 19:52:41

Showing 33 results. ☒ Receive New Job Notifications

MAINSESSION_task2 - T2GJgmceEeWPuPKvpC0s5Q_MAINSESSION_task2

This grid task is completed.

General Properties

Name	MAINSESSION_task2
Type	Grid Task
Started By	Administrator
User Security Domain	Native
Start Time	09/29/2015 19:52:57
Elapsed Time	00:00:43
End Time	09/29/2015 19:53:41
% Task Completed	100
Monitoring URL	http://psrhagadn21.informatica.com:9080/Blaze?tasktype=gridtask&id=qtid-24-1-79555597-4&isParent=false
Incoming Task Dependencies	, PRESESSION_task0PRESESSION_task1
Outgoing Task Dependencies	, POSTSESSION_task3

Viewing Hadoop Environment Logs in the Administrator Tool

You can view log events for a Blaze or Hive mapping from the Monitor tab of the Administrator tool.

1. In the Administrator tool, click the **Monitor** tab.
2. Select the **Execution Statistics** view.
3. In the Navigator, choose to open an ad hoc job, a deployed mapping job, or a workflow.
 - To choose an ad hoc job, expand a Data Integration Service and click **Ad Hoc Jobs**.
 - To choose a deployed mapping job, expand an application and click **Deployed Mapping Jobs**.
 - To choose a workflow, expand an application and click **Workflows**.

The list of jobs appears in the contents panel.

4. Click **Actions > View Logs for Selected Object** to view the run-time logs for the mapping.

The log file shows the results of the Hive queries and Blaze engine queries run by the Data Integration Service. This includes the location of Hive session logs and Hive session history file.

Monitoring a Mapping

You can monitor a mapping that runs in the Hadoop environment.

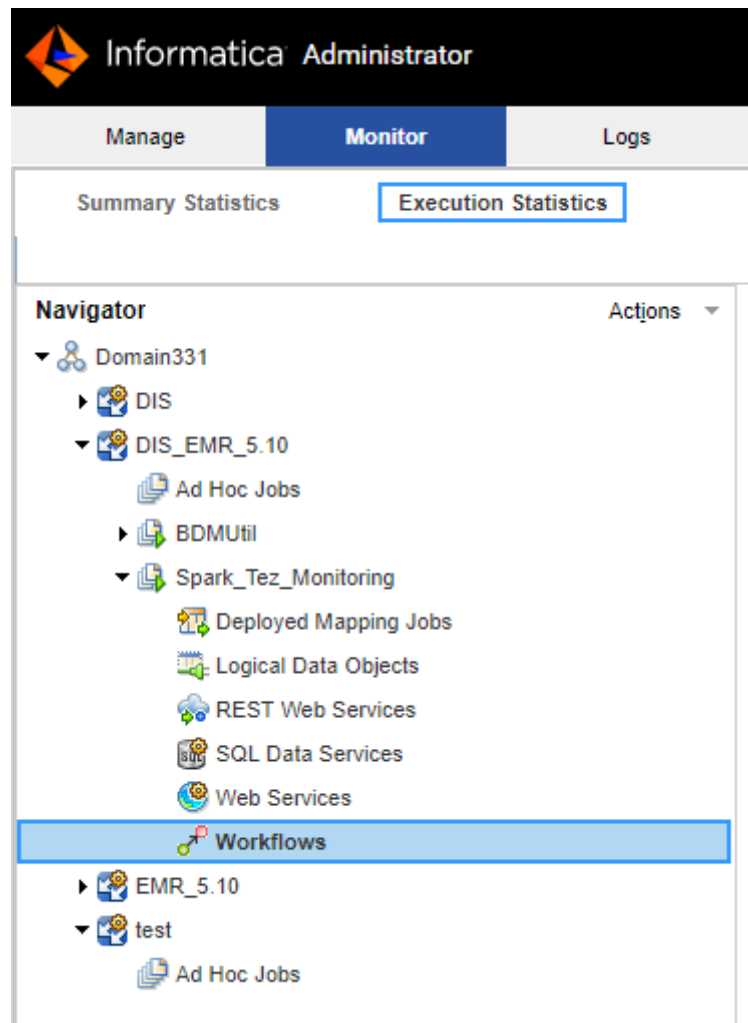
1. In the Administrator tool, click the **Monitor** tab.
2. Select the **Execution Statistics** view.
3. In the Navigator, choose to open an ad hoc job, a deployed mapping job, or a workflow.
 - To choose an ad hoc job, expand a Data Integration Service and click **Ad Hoc Jobs**.

The screenshot shows the Informatica Administrator interface. The top navigation bar includes 'Manage', 'Monitor' (selected), 'Logs', 'Reports', 'Security', and 'Cloud'. Below this, there are tabs for 'Summary Statistics' and 'Execution Statistics' (selected). The left sidebar shows a tree view under 'Domain331' with 'DIS' expanded, showing 'Ad Hoc Jobs' selected. The main pane displays a table of 'Ad Hoc Jobs' with columns: Name, Type, State, Job ID, Started By, Start Time, Elapsed Time, and End Time. The table lists several 'm_FF_FF' mapping jobs, most in a 'Completed' state. Below the table, there are tabs for 'Properties', 'Spark Execution Plan', 'Summary Statistics', 'Detailed Statistics', and 'Historical Statistics'. The 'Properties' tab is active, showing a message 'The Mapping job has completed.' and a 'General Properties' section with details like Name (m_FF_FF), Type (Mapping), Started By (Administrator), User Security Domain (Native), Start Time (04/06/2018 16:23:44), Elapsed Time (00:01:58), End Time (04/06/2018 16:25:43), and Operating System Profile (OSP_usr04).

- To choose a deployed mapping job, expand an application and click **Deployed Mapping Jobs**.

The screenshot shows the Informatica Administrator interface. The top navigation bar includes 'Manage', 'Monitor' (selected), 'Logs', 'Reports', 'Security', and 'Cloud'. Below this, there are tabs for 'Summary Statistics' and 'Execution Statistics' (selected). The left sidebar shows a tree view under 'Domain331' with 'DIS' expanded, showing 'Deployed Mapping Jobs' selected. The main pane displays a table of 'Deployed Mapping Jobs' with columns: Name, Type, State, Job ID, Started By, Start Time, Elapsed Time, and End Time. The table lists various jobs, including 'm_Ustx_Readback' and 'exec0', mostly in a 'Completed' state. Below the table, there are tabs for 'Properties', 'Hive Execution Plan', 'Summary Statistics', 'Detailed Statistics', and 'Historical Statistics'. The 'Properties' tab is active, showing a message 'The Deployed Mapping job has completed.' and a 'General Properties' section with details like Name (m_Ustx_Readback), Type (Deployed Mapping), Started By (Administrator), User Security Domain (Native), Start Time (04/06/2018 21:33:39), Elapsed Time (00:08:17), and End Time (04/06/2018 21:41:57).

- To choose a workflow, expand an application and click **Workflows**.



The list of jobs appears in the contents panel.

4. Click a job to view its properties.

The contents panel shows the default **Properties** view for the job. For a Blaze engine mapping, the Blaze engine monitoring URL appears in the general properties in the details panel. The monitoring URL is a link to the YARN web user interface for Spark jobs.

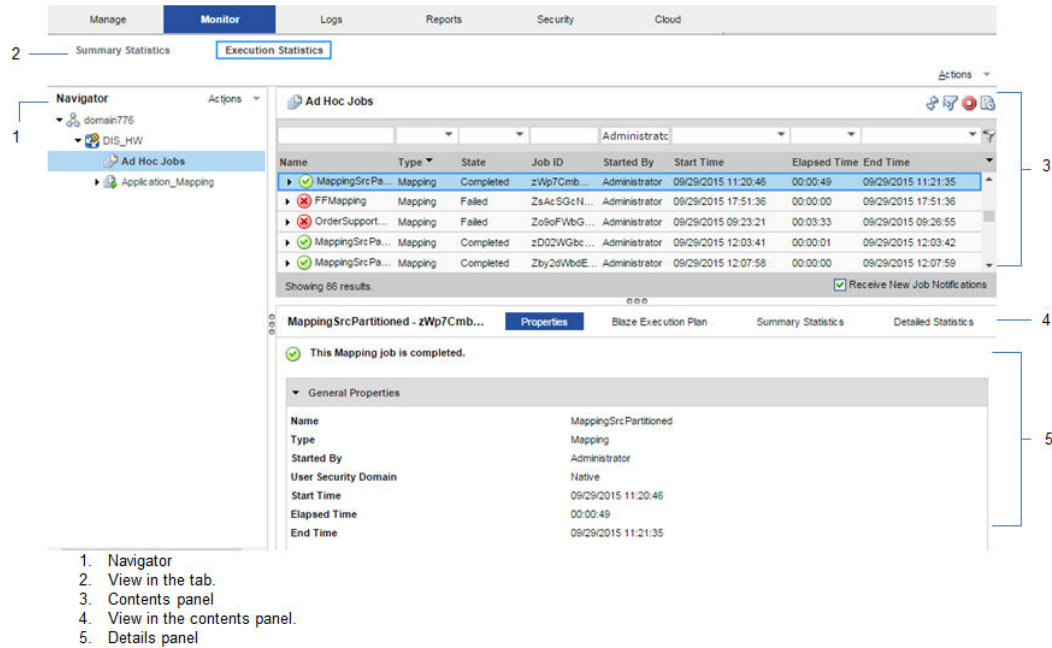
5. Choose a view in the contents panel to view more information about the job:
 - To view the execution plan for the mapping, select the **Execution Plan** view.
 - To view the summary statistics for a job, click the **Summary Statistics** view.
 - To view the detailed statistics for a job, click the **Detailed Statistics** view.

Note: You can view the number of rows processed in the Summary Statistics for a Hive source or target. The remaining values do not appear for Hive sources and targets.

Blaze Engine Monitoring

You can monitor statistics and view log events for a Blaze engine mapping job in the Monitor tab of the Administrator tool. You can also monitor mapping jobs for the Blaze engine in the Blaze Job Monitor web application.

The following image shows the Monitor tab in the Administrator tool:



The Monitor tab has the following views:

Summary Statistics

Use the **Summary Statistics** view to view graphical summaries of object states and distribution across the Data Integration Services. You can also view graphs of the memory and CPU that the Data Integration Services used to run the objects.

Execution Statistics

Use the **Execution Statistics** view to monitor properties, run-time statistics, and run-time reports. In the Navigator, you can expand a Data Integration Service to monitor **Ad Hoc Jobs** or expand an application to monitor deployed mapping jobs or workflows

When you select **Ad Hoc Jobs**, deployed mapping jobs, or workflows from an application in the Navigator of the **Execution Statistics** view, a list of jobs appears in the contents panel. The contents panel displays jobs that are in the queued, running, completed, failed, aborted, and cancelled state. The Data Integration Service submits jobs in the queued state to the cluster when enough resources are available.

The contents panel groups related jobs based on the job type. You can expand a job type to view the related jobs under it.

Access the following views in the **Execution Statistics** view:

Properties

The **Properties** view shows the general properties about the selected job such as name, job type, user who started the job, and start time of the job. You can also monitor jobs on the Hadoop cluster from the Monitoring URL that appears for the mapping in the general properties. The Monitoring URL opens the

Blaze Job Monitor in a web page. The Blaze Job Monitor displays detailed monitoring statistics for a mapping such as the number of grid tasks, grid segments, or tasklets, and recovery attempts for each tasklet.

Blaze Execution Plan

The Blaze execution plan displays the Blaze engine script that the Data Integration Service generates based on the mapping logic. The execution plan includes the tasks that the script depends on. Each script has a unique identifier.

Summary Statistics

The **Summary Statistics** view appears in the details panel when you select a mapping job in the contents panel. The **Summary Statistics** view displays throughput and resource usage statistics for the job.

You can view the following throughput statistics for the job:

- Source. The name of the mapping source file.
- Target name. The name of the target file.
- Rows. The number of rows read for source and target. If the target is Hive, this is the only summary statistic available.
- Average Rows/Sec. Average number of rows read per second for source and target.
- Bytes. Number of bytes read for source and target.
- Average Bytes/Sec. Average number of bytes read per second for source and target.
- First Row Accessed. The date and time when the Data Integration Service started reading the first row in the source file.
- Dropped rows. Number of source rows that the Data Integration Service did not read.

Detailed Statistics

The **Detailed Statistics** view appears in the details panel when you select a mapping job in the contents panel. The **Detailed Statistics** view displays graphs of the throughput and resource usage statistics for the job run.

Blaze Job Monitoring Application

Use the Blaze Job Monitor application to monitor Blaze engine jobs on the Hadoop cluster.

You configure the host that starts the Blaze Job Monitor in the Hadoop connection properties. You might want to configure the Blaze Job Monitor address to avoid conflicts with other users on the same cluster, or if you have access to a limited number of nodes. If you do not configure the Blaze Job Monitor address, the Grid Manager starts the host on the first alphabetical cluster node with a default port of 9080.

The Blaze engine monitoring URL appears in the Monitor tab of the Administrator tool when you view a Blaze engine mapping job. When you click the URL, the Blaze engine monitoring application opens in a web page.

Note: You can also access the Blaze Job Monitor through the LDTM log. After the session load summary, the log displays a list of segments within the grid task. Each segment contains a link to the Blaze Job Monitor. Click on a link to see the execution details of that segment.

You configure the host that starts the Blaze Job Monitor in the Hadoop connection properties. The default address is <hostname>:9080.

The following image shows the Blaze Job Monitor:

Name	Start Time	End Time	Elapsed Time	State	Host Name	Log
g98-499-1-42938595-32_s6_1-0_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	psishaqan21.informatica.com	Log
g98-499-1-42938595-32_s6_1-1_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	psishaqan28.informatica.com	Log
g98-499-1-42938595-32_s6_1-1_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:59 PM	0:1:52	Succeeded	psishaqan23.informatica.com	Log
g98-499-1-42938595-32_s6_1-2_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	psishaqan28.informatica.com	Log
g98-499-1-42938595-32_s6_1-0_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	psishaqan25.informatica.com	Log
g98-499-1-42938595-32_s4_1-0_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	psishaqan21.informatica.com	Log
g98-499-1-42938595-32_s6_1-2_1	Mon Oct 31 2016 2:45:04 PM	Mon Oct 31 2016 2:45:47 PM	0:0:43	Succeeded	psishaqan21.informatica.com	Log
g98-499-1-42938595-32_s1_1-1_1	Mon Oct 31 2016 2:45:04 PM	Mon Oct 31 2016 2:45:07 PM	0:0:3	Succeeded	psishaqan28.informatica.com	Log

Use the **Task History** panel on the left to filter Blaze mapping jobs by the following criteria:

- Grid task. A parallel processing job request sent by the Blaze engine executor to the Grid Manager. You can further filter by all tasks, succeeded tasks, running tasks, or failed tasks.
- Grid segment. Part of a grid mapping that is contained in a grid task.
- Tasklet. A partition of a grid segment that runs on a separate DTM.
- Tasklet Attempts. The number of recovery attempts to restart a tasklet. Click **Log** to view the mapping grid task log.

The Blaze Job Monitor displays the task history for mapping jobs with the same namespace. You can monitor properties for a task such as start time, end time, elapsed time, or state of the task. You can also view log events. If you filter mapping jobs by grid segment, you can mouse over a grid segment to view the logical name of the segment.

By default, the Blaze Job Monitor automatically refreshes the list of tasks every five seconds and reverts to the first page that displays tasks. Disable auto refresh if you want to browse through multiple pages. To turn off automatic refresh, click **Action > Disable Auto Refresh**.

The Blaze Job Monitor displays the first 100,000 grid tasks run in the past seven days. The Blaze Job Monitor displays the grid segments, tasklets, and tasklet attempts for grid tasks that are running and grid tasks that were accessed in the last 30 minutes.

Blaze Summary Report

The Blaze Summary Report displays more detailed statistics about a mapping job. In the Blaze Job Monitor, a green summary report button appears beside the names of successful grid tasks. Click the button to open the Blaze Summary Report.

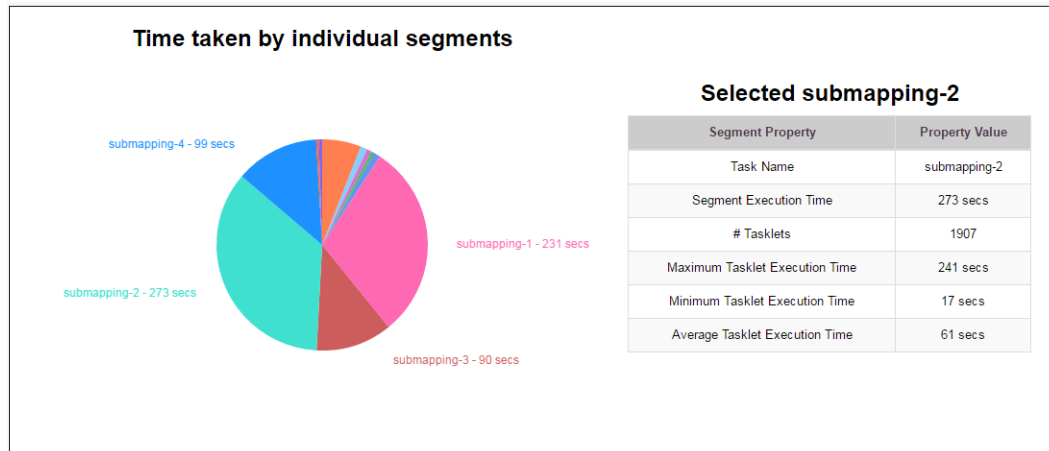
Note: The Blaze Summary Report is available for technical preview. Technical preview functionality is supported but is unwarranted and is not production-ready. Informatica recommends that you use in non-production environments only.

All Grid Tasks

Show 25 entries						
Name	Start Time	End Time	Elapsed Time	State		
gtid-476-1-36233666-2	Mon Oct 17 2016 1:32:42 PM	Mon Oct 17 2016 1:53:18 PM	0:20:36	Succeeded		
gtid-476-1-36233666-1	2016 1:30:51 PM	Mon Oct 17 2016 1:31:20 PM	0:0:28	Succeeded		
gtid-441-1-26795155-4	Mon Oct 17 2016 11:55:06 AM	Mon Oct 17 2016 11:59:44 AM	0:4:37	Succeeded		
gtid-441-1-26795155-3	Mon Oct 17 2016 11:47:10 AM	Mon Oct 17 2016 11:51:51 AM	0:4:40	Succeeded		
gtid-441-1-26795155-2	Mon Oct 17 2016 11:02:37 AM	Mon Oct 17 2016 11:06:18 AM	0:3:40	Succeeded		
gtid-441-1-26795155-1	Mon Oct 17 2016 10:53:35 AM	Mon Oct 17 2016 10:54:27 AM	0:0:51	Failed		
gtid-437-1-25270758-1	Mon Oct 17 2016 10:28:08 AM	Mon Oct 17 2016 10:28:56 AM	0:0:47	Failed		

Time Taken by Individual Segments

A pie chart visually represents the time taken by individual segments contained within a grid task.



When you click on a particular segment in the pie chart, the **Selected Submapping** table displays detailed information about that segment. The table lists the following segment statistics:

- Task Name. The logical name of the selected segment.
- Segment Execution Time. The time taken by the selected segment.
- # Tasklets. The number of tasklets in the selected segment.
- Minimum Tasklet Execution Time. The execution time of the tasklet within the selected segment that took the shortest time to run.
- Maximum Tasklet Execution Time. The execution time of the tasklet within the selected segment that took the longest time to run.
- Average Tasklet Execution Time. The average execution time of all tasklets within the selected segment.

Mapping Properties

The Mapping Properties table lists basic information about the mapping job.

Mapping Properties	
Mapping Property	Property Value
DIS Name	dis_cdh
Informatica Version	10.1.1
Mapping Name	q97_hive_mapping
Total Segments	13
Maximum Segment Execution Time	273 secs
Minimum Segment Execution Time	0 secs
Average Segment Execution Time	59 secs
Mapping Execution Time	0:10:14

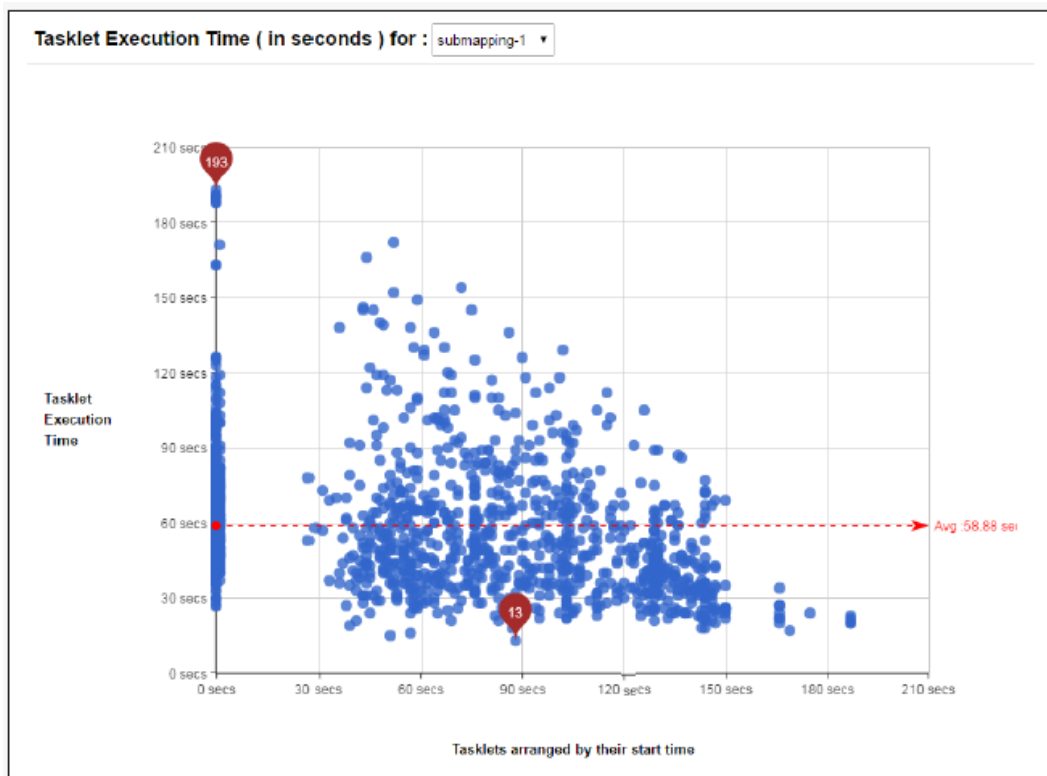
The Mapping Properties table displays the following information:

- The DIS name under which the mapping was run.
- The Informatica version.
- The name of the mapping.
- The total number of segments for the mapping.
- The execution time of the segment that took the longest time to run.
- The execution time of the segment that took the shortest time to run.
- The average execution time of all segments.
- The total mapping execution time.

Tasklet Execution Time

A time series graph displays the execution time of all tasklets within the selected segment.

The x-axis represents the tasklet start time and the y-axis represents the actual tasklet execution time. The red dashed line represents the average execution time for all tasklets, and the two red markers show the minimum and maximum execution times within the segment.



Selected Tasklet Information

When you select a tasklet from the **Tasklet Execution Time** graph, you can see more data about that individual tasklet. This data includes source and target row counts as well as cache information for any

cache-based transformation processed by the tasklet. Click the **Get Detailed Log** button to see a full log of the selected tasklet.

The screenshot displays the 'Selected tasklet' log for tasklet `gtid-299-1-82064486-16_s8_t-394_1`. It features three main sections, each with a 'Show' dropdown set to 10 entries and a search bar.

Source Name Table:

Source Name	# Rows Processed
empty_source1	0
Read_catalog_sales	4,937,484

Showing 1 to 2 of 2 entries. Red arrow points to the '# Rows Processed' column with the text: 'Row counts for all sources processed by tasklet'.

Target Name Table:

Target Name	# Rows Processed
DETarget_Aggregator1_G0	0

Showing 1 to 1 of 1 entries. Red arrow points to the '# Rows Processed' column with the text: 'Row counts for all targets written by tasklet'.

Transformation Name Table:

Transformation Name	Index Cache (bytes)		Data Cache (bytes)	
	Configured	Used	Configured	Used
Joiner1	178,956,970	14,400	357,913,940	6,968

Showing 1 to 1 of 1 entries. Red arrow points to the 'Used' column for Data Cache with the text: 'Cache information for any cache-based transformation processed by tasklet'.

Blaze Engine Logs

The mapping run log appears in the LDTM log on the domain and in the tasklet logs on the Hadoop cluster.

You can find information about the mapping run on the Blaze engine in the following log files:

LDTM log

The LDTM logs the results of the mapping run on the Blaze engine. You can view the LDTM log from the Developer tool or the Monitoring tool for a mapping job.

You can configure the Data Integration Service to log details about the mapping execution to the session log. To enable logging of LDTM mapping execution details, set the log tracing level to verbose initialization or verbose data.

Note: Informatica recommends setting the tracing level to verbose data only for debugging. Do not use verbose data to run jobs concurrently for production.

Mapping execution details include the following information:

- Start time, end time, and state of each task
- Blaze Job Monitor URL
- Number of total, succeeded, and failed/cancelled tasklets
- Number of processed and rejected rows for sources and targets

- Data errors, if any, for transformations in each executed segment

Blaze component and tasklet logs

The Blaze engine stores tasklet and Blaze component log events in temporary and permanent directories on the Hadoop cluster.

The following list describes Blaze log properties and where to configure them:

infagrid.node.local.root.log.dir

Temporary directory for tasklet and component logs. An administrator must create a directory with read, write, and execute permissions on all nodes on the Hadoop cluster.

Configure this property in the Advanced properties of the Blaze configuration in the Hadoop connection. Default is `/tmp/infa/logs/blaze`.

Blaze Staging Directory

Permanent directory on HDFS that contains tasklet log event.

Configure this property in the Blaze Configuration properties of the Hadoop connection.

infagrid.delete.local.log

Boolean property to determine whether to delete tasklet logs from the temporary directory after copying the logs to the permanent directory.

Configure this property in the Advanced properties of the Blaze configuration in the Hadoop connection. To retain the logs in the temporary directory, set to false. Default is true.

yarn.nodemanager.local-dirs

Directory for the DTM process. The Data Integration Service stores logs in this location if you do not configure the `infagrid.node.local.root.log.dir` or Blaze Staging Directory properties.

This property is configured by default when the cluster configuration is imported.

Viewing Blaze Logs

You can view logs for a Blaze mapping from the Blaze Job Monitor.

1. In the Blaze Job Monitor, select a job from the list of jobs.
2. In the row for the selected job, click the **Logs** link.

The screenshot shows the Informatica Blaze Job Monitor interface. On the left is a navigation menu with options like Task History, Grid Tasks, All Succeeded, All Failed, Grid Segments, Tasklets, and Attempts. The main area is titled 'Blaze Job Monitor' and contains a table of 'All Tasklet Attempts'. The table has columns for Name, Start Time, End Time, Elapsed Time, State, Host Name, and Log. There are 10 rows of data, all showing a 'Succeeded' state. A search bar and a 'Show 25 entries' link are at the top of the table.

Name	Start Time	End Time	Elapsed Time	State	Host Name	Log
gtd-499-1-42958595-32_s6_1-0_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	prfhagadn21.informatica.com	Log
gtd-499-1-42958595-32_s6_1-1_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	prfhagadn28.informatica.com	Log
gtd-499-1-42958595-32_s6_1-1_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:59 PM	0:1:52	Succeeded	prfhagadn23.informatica.com	Log
gtd-499-1-42958595-32_s6_1-2_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	prfhagadn26.informatica.com	Log
gtd-499-1-42958595-32_s6_1-0_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	prfhagadn26.informatica.com	Log
gtd-499-1-42958595-32_s4_1-0_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	prfhagadn21.informatica.com	Log
gtd-499-1-42958595-32_s6_1-2_1	Mon Oct 31 2016 2:45:04 PM	Mon Oct 31 2016 2:45:47 PM	0:0:43	Succeeded	prfhagadn21.informatica.com	Log
gtd-499-1-42958595-32_s1_1-1_1	Mon Oct 31 2016 2:45:04 PM	Mon Oct 31 2016 2:45:07 PM	0:0:3	Succeeded	prfhagadn28.informatica.com	Log

The log events appear in another browser window.

Orchestrator Sunset Time

Orchestrator sunset time is the maximum lifetime for an Orchestrator service. Sunset time determines the maximum amount of time that the Blaze engine can run a mapping job. The default sunset time is 24 hours. After 24 hours, the Orchestrator shuts down, which causes the Blaze Grid Manager to shut down.

You can configure the Orchestrator sunset time to be greater than or less than 24 hours. Configure the following property in the Hadoop connection:

Property	Description
<code>infagrid.orchestrator.svc.sunset.time</code>	Maximum lifetime for an Orchestrator service, in hours. Default is 24 hours.

You can also disable sunset by setting the property to 0 or a negative value. If you disable sunset, the Orchestrator never shuts down during a mapping run.

Troubleshooting Blaze Monitoring

When I run a mapping on the Blaze engine and try to view the grid task log, the Blaze Job Monitor does not fetch the full log.

The grid task log might be too large. The Blaze Job Monitor can only fetch up to 2 MB of an aggregated log. The first line of the log reports this information and provides the location of the full log on HDFS. Follow the link to HDFS and search for "aggregated logs for grid mapping." The link to the full log contains the grid task number.

The Blaze Job Monitor will not start.

Check the Hadoop environment logs to locate the issue. If you do not find an issue, stop the Grid Manager with the `infacmd stopBlazeService` command and run the mapping again.

The Monitoring URL does not appear in the Properties view of the Administrator tool.

Locate the URL in the YARN log.

When Blaze processes stop unexpectedly, Blaze does not save logs in the expected location.

When Blaze stops unexpectedly, you can access Blaze service logs through the YARN monitor. Use one of these methods:

- The Grid Manager log contains all Blaze job container IDs and identifies the host on which Blaze ran. Alter the Grid Manager log URL with the container ID and host name of the Blaze host.
- Run the command `yarn logs -applicationID <Blaze Grid Manager Application ID>`.

A Blaze Job Monitor that has been running for several days loses its connection to the Application Timeline Server on the Hortonworks cluster.

The Blaze engine requires a running Application Timeline Server on the cluster. When the Blaze engine starts a mapping run, the Blaze Job Monitor checks the state of the Application Timeline Server. The Grid Manager will start it if it is not running. When the connection to the Application Timeline Server is lost, the Blaze engine attempts to reconnect to it. If the Application Timeline Server stops during a Blaze mapping run, you can restart it by restarting the Grid Manager.

Note: When the Application Timeline Server is configured to run on the cluster by default, the cluster administrator must manually restart it on the cluster.

When a mapping takes more than 24 hours to execute, the mapping fails.

When mappings run on the Blaze engine for more than 24 hours, some mappings might fail because the Orchestrator service has a default sunset time of 24 hours. After 24 hours, the Orchestrator shuts down, which causes the Blaze Grid Manager to shut down.

To increase the sunset time to be more than 24 hours, configure the following property in the Hadoop connection advanced properties:

```
infagrid.orchestrator.svc.sunset.time=[HOURS]
```

You can also disable sunset by setting the property to 0 or a negative value. If you disable sunset, the Blaze Grid Manager never shuts down.

Spark Engine Monitoring

You can monitor statistics and view log events for a Spark engine mapping job in the Monitor tab of the Administrator tool. You can also monitor mapping jobs for the Spark engine in the YARN web user interface.

The following image shows the Monitor tab in the Administrator tool:

The screenshot displays the Administrator tool's Monitor tab. The top navigation bar includes 'Manage', 'Monitor' (selected), 'Logs', 'Reports', 'Security', and 'Cloud'. Below this, the 'Execution Statistics' sub-tab is active. The left sidebar (Navigator) shows a tree structure with 'domain776' expanded, containing 'DIS_HW' and 'Ad Hoc Jobs' (selected). The main content area shows a table of 'Ad Hoc Jobs' with columns: Name, Type, State, Job ID, Started By, Start Time, Elapsed Time, and End Time. The table lists several jobs, with the first one 'MappingSrcPa...' being 'Completed'. Below the table, a 'MappingSrcPartitioned - zWp7Cmb...' job is selected, and its 'Properties' panel is open, showing details like Name, Type, Started By, User Security Domain, Start Time, Elapsed Time, and End Time. A legend at the bottom left identifies the numbered callouts: 1. Navigator, 2. View in the tab, 3. Contents panel, 4. View in the contents panel, 5. Details panel.

Name	Type	State	Job ID	Started By	Start Time	Elapsed Time	End Time
MappingSrcPa...	Mapping	Completed	zWp7Cmb...	Administrator	09/29/2015 11:20:46	00:00:49	09/29/2015 11:21:35
FFIMapping	Mapping	Failed	ZsAcSGcH...	Administrator	09/29/2015 17:51:36	00:00:00	09/29/2015 17:51:36
OrderSupport...	Mapping	Failed	ZoIoFWbG...	Administrator	09/29/2015 09:23:21	00:03:33	09/29/2015 09:26:55
MappingSrcPa...	Mapping	Completed	zD02WGb...	Administrator	09/29/2015 12:03:41	00:00:01	09/29/2015 12:03:42
MappingSrcPa...	Mapping	Completed	Zby2dVbdE...	Administrator	09/29/2015 12:07:58	00:00:00	09/29/2015 12:07:59

General Properties	
Name	MappingSrcPartitioned
Type	Mapping
Started By	Administrator
User Security Domain	Native
Start Time	09/29/2015 11:20:46
Elapsed Time	00:00:49
End Time	09/29/2015 11:21:35

1. Navigator
2. View in the tab.
3. Contents panel
4. View in the contents panel.
5. Details panel

The Monitor tab has the following views:

Summary Statistics

Use the **Summary Statistics** view to view graphical summaries of object states and distribution across the Data Integration Services. You can also view graphs of the memory and CPU that the Data Integration Services used to run the objects.

Execution Statistics

Use the **Execution Statistics** view to monitor properties, run-time statistics, and run-time reports. In the Navigator, you can expand a Data Integration Service to monitor **Ad Hoc Jobs** or expand an application to monitor deployed mapping jobs or workflows

When you select **Ad Hoc Jobs**, deployed mapping jobs, or workflows from an application in the Navigator of the **Execution Statistics** view, a list of jobs appears in the contents panel. The contents panel displays jobs that are in the queued, running, completed, failed, aborted, and cancelled state. The Data Integration Service submits jobs in the queued state to the cluster when enough resources are available.

The contents panel groups related jobs based on the job type. You can expand a job type to view the related jobs under it.

Access the following views in the **Execution Statistics** view:

Properties

The **Properties** view shows the general properties about the selected job such as name, job type, user who started the job, and start time of the job.

Spark Execution Plan

When you view the Spark execution plan for a mapping, the Data Integration Service translates the mapping to a Scala program and an optional set of commands. The execution plan shows the commands and the Scala program code.

Summary Statistics

The **Summary Statistics** view appears in the details panel when you select a mapping job in the contents panel. The **Summary Statistics** view displays the following throughput statistics for the job:

- Source. The name of the mapping source file.
- Target name. The name of the target file.
- Rows. The number of rows read for source and target.

The following image shows the **Summary Statistics** view in the details panel for a mapping run on the Spark engine:

The screenshot shows the Informatica Administrator interface. The top navigation bar includes 'Manage', 'Monitor', 'Logs', 'Reports', 'Security', and 'Cloud'. The 'Monitor' tab is active, and the 'Execution Statistics' sub-tab is selected. The left sidebar shows the 'Navigator' with a tree view containing 'QA_MERCURY_DOMAIN_1011', 'DIS_HOI', 'DIS_HOP25', 'Ad Hoc Jobs', and 'DIS_TE'. The 'Ad Hoc Jobs' folder is expanded. The main panel displays a table of jobs. The job 'Map_Spark' is selected, and its details are shown in the 'Summary Statistics' view. The 'Summary Statistics' view is highlighted with a red circle. Below it, the 'Spark Run Stages' view is shown, which is highlighted with a yellow box. The 'Spark Run Stages' view displays a table with columns: Source, Rows, Average Rows/Sec, Bytes, Average Bytes/Sec, First Row Accessed, and Dropped Rows. The table has two rows: 'Read_Ratp' and 'Write_Ratp'.

Source	Rows	Average Rows/Sec	Bytes	Average Bytes/Sec	First Row Accessed	Dropped Rows
Read_Ratp	28967753	N/A	N/A	N/A	N/A	-1
Target						
Write_Ratp	28967753	N/A	N/A	N/A	N/A	N/A

You can also view the Spark run stages information in the details pane of the Summary Statistics view on the Execution Statistics Monitor tab. It appears as a list after the sources and targets.

The **Spark Run Stages** displays the absolute counts and throughput of rows and bytes related to the Spark application stage statistics. Rows refer to the number of rows that the stage writes, and bytes refer to the bytes broadcasted in the stage.

The following image displays the Spark Run Stages:

The screenshot shows the 'Ad Hoc Jobs' interface. The top table lists several mapping jobs, all in a 'Completed' state. The selected job is 'q3_hive_mapping_default_with_monitor'. Below this, the 'Detailed Statistics' tab is active, showing a hierarchical view of the job's execution. It includes a 'Throughput' section with source and target statistics, and a 'Spark Run Stages' section with detailed performance metrics for each stage.

Name	Type	State	Job ID	Started By	Start Time	Elapsed Time	End Time
q3_hive_mapping_default_with_monitor	Mapping	Completed	TdyKyvXEvSBHOKyMm...	Administrator	03/20/2018 13:48:43	00:03:49	03/20/2018 13:52:32
q3_hive_mapping_default_with_monitor	Mapping	Completed	ESHy4uXEvSBHOKyMmDQ	Administrator	03/20/2018 13:47:06	00:03:56	03/20/2018 13:51:03
q3_hive_mapping_default_with_monitor	Mapping	Completed	DVYLSwTEwQzHyJMPQI...	Administrator	03/20/2018 13:18:20	00:03:52	03/20/2018 13:22:13
q3_hive_mapping_default_with_monitor	Mapping	Completed	LJ3bCwREwQzHyJMPQI...	Administrator	03/20/2018 13:04:53	00:03:19	03/20/2018 13:08:13
hdfs_filter	Mapping	Completed	AMSRHwREwQzHyJMPQI...	Administrator	03/20/2018 11:52:02	00:02:09	03/20/2018 11:54:12

Source		Rows	Average Rows/Sec
Read_dsm		52000	
Read_state_sim		73049	1974
Read_state_sales		115203420	294000

Target		Rows	Average Rows/Sec
Write_tgt_c3		144	144

Spark Run Stages	Rows	Average Rows/Sec	Bytes	Average Bytes/Sec
Stage_0_IntoSpark0	73049	1974	518000	14000
Stage_1_IntoSpark0	73049	N/A	518000	N/A
Stage_2_IntoSpark0	52000	52000	500240	500240

For example, the Spark Run Stages column contains the Spark application staged information starting with stage_<ID>. In the example, Stage_0 shows the statistics related to the Spark run stage with ID=0 in the Spark application.

Consider when the Spark engine reads source data that includes a self-join with verbose data enabled. In this scenario, the optimized mapping from the Spark application does not contain any information on the second instance of the same source in the Spark engine logs.

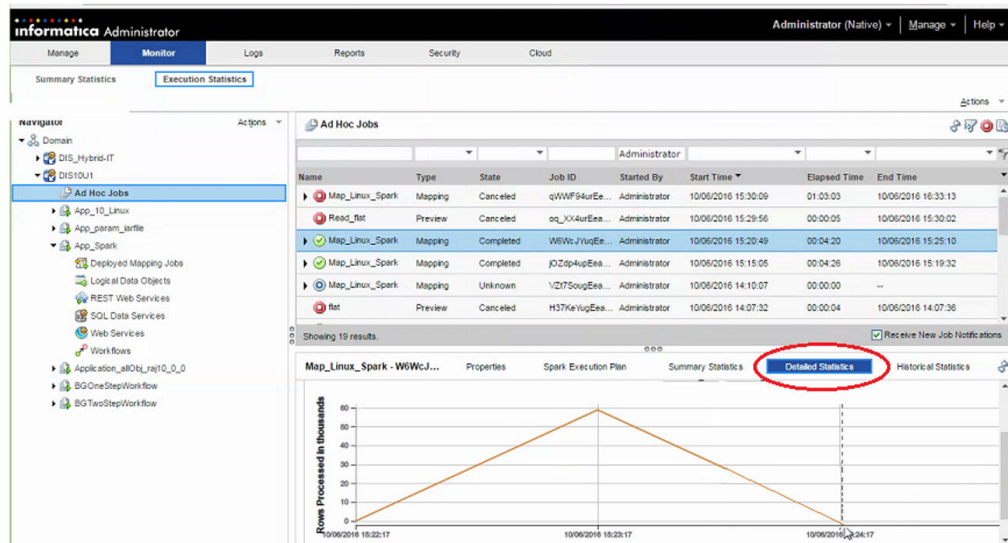
Consider when you read data from the temporary table and the Hive query of the customized data object leads to the shuffling of the data. In this scenario, the filtered source statistics appear instead of reading from the temporary source table in the Spark engine log.

When you run a mapping with Spark monitoring enabled, performance varies based on the mapping complexity. It can take up to three times longer than usual processing time with monitoring enabled. By default, monitoring is disabled.

Detailed Statistics

The **Detailed Statistics** view appears in the details panel when you select a mapping job in the contents panel. The **Detailed Statistics** view displays a graph of the row count for the job run.

The following image shows the **Detailed Statistics** view in the details panel for a mapping run on the Spark engine:



Viewing Hive Tasks

When you have a Hive source with a transactional table, you can view the Hive task associated with the Spark job.

When you run a mapping on Spark that launches Hive tasks, you can view the Hive query statistics in the session log and in the Administrator tool for monitoring along with the Spark application. For example, you can monitor information related to the Update Strategy transformation and SQL authorization associated to the mapping on Spark.

You can view the Summary Statistics for a Hive task in the Administrator tool. The Spark statistics continue to appear. When the Spark engine launches a Hive task, you can see Source Load Summary and Target Load summary including Spark data frame with Hive task statistics. Otherwise, when you have only a Spark task, the Source Load Summary and Target Load Summary do not appear in the session log.

Under Target Load Summary, all Hive instances will be prefixed with 'Hive_Target_'. You can see same instance name in the Administrator tool.

Spark Engine Logs

The Spark engine logs appear in the LDTM log. The LDTM logs the results of the Spark engine execution plan run for the mapping. You can view the LDTM log from the Developer tool or the Monitoring tool for a mapping job.

The log for the Spark engine shows the step to translate the mapping to an internal format, steps to optimize the mapping, steps to render the mapping to Spark code, and the steps to submit the code to the Spark executor. The logs also show the Scala code that the Logical Data Translation Generator creates from the mapping logic.

When you run Sqoop mappings on the Spark engine, the Data Integration Service prints the Sqoop log events in the mapping log.

Viewing Spark Logs

You can view logs for a Spark mapping from the YARN web user interface.

1. In the YARN web user interface, click an application ID to view.
 2. Click the application **Details**.
 3. Click the **Logs** URL in the application details to view the logs for the application instance.
- The log events appear in another browser window.

Troubleshooting Spark Engine Monitoring

Do I need to configure a port for Spark Engine Monitoring?

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. 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.

Hive Engine Monitoring

You can monitor statistics and view log events for a Hive engine mapping job in the Monitor tab of the Administrator tool.

The following image shows the Monitor tab in the Administrator tool:

The screenshot displays the Administrator tool's Monitor tab. The interface includes a top navigation bar with tabs: Manage, Monitor (selected), Logs, Reports, Security, and Cloud. Below this, there are sub-tabs: Summary Statistics and Execution Statistics (selected). On the left, a Navigator pane shows a tree structure with 'domain776' expanded, containing 'DIS_HW' and 'Ad Hoc Jobs'. The 'Ad Hoc Jobs' folder is selected. The main area shows a table of 'Ad Hoc Jobs' with columns: Name, Type, State, Job ID, Started By, Start Time, Elapsed Time, and End Time. The table lists several mapping jobs, with the first one 'MappingSrcPa...' being highlighted. Below the table, there are tabs for 'Properties', 'Hive Execution Plan', 'Summary Statistics', and 'Detailed Statistics'. The 'Properties' tab is active, showing a message 'This Mapping job is completed.' and a 'General Properties' section with details like Name, Type, Started By, User Security Domain, Start Time, Elapsed Time, and End Time. Numbered callouts 1 through 5 point to specific UI elements: 1. Navigator, 2. Summary Statistics View and Execution Statistics View, 3. Contents Panel, 4. View on the contents panel, 5. Details panel.

Name	Type	State	Job ID	Started By	Start Time	Elapsed Time	End Time
MappingSrcPa...	Mapping	Completed	zWp7Cmb...	Administrator	09/29/2015 11:20:46	00:00:49	09/29/2015 11:21:35
FFMapping	Mapping	Failed	ZsAcSGcN...	Administrator	09/29/2015 17:51:36	00:00:00	09/29/2015 17:51:36
OrderSupport...	Mapping	Failed	Zo8oFWbG...	Administrator	09/29/2015 09:23:21	00:03:33	09/29/2015 09:26:55
MappingSrcPa...	Mapping	Completed	zD02WGb...	Administrator	09/29/2015 12:03:41	00:00:01	09/29/2015 12:03:42
MappingSrcPa...	Mapping	Completed	Zby2dWbE...	Administrator	09/29/2015 12:07:58	00:00:00	09/29/2015 12:07:59

General Properties	
Name	MappingSrcPartitioned - zWp7Cmb...
Type	Mapping
Started By	Administrator
User Security Domain	Native
Start Time	09/29/2015 11:20:46
Elapsed Time	00:00:49
End Time	09/29/2015 11:21:35

1. Navigator
2. Summary Statistics View and Execution Statistics View
3. Contents Panel
4. View on the contents panel
5. Details panel

The Monitor tab has the Summary Statistics and Execution Statistics views.

Summary Statistics

Use the **Summary Statistics** view to view graphical summaries of object states and distribution across the Data Integration Services. You can also view graphs of the memory and CPU that the Data Integration Services used to run the objects.

Execution Statistics

Use the **Execution Statistics** view to monitor properties, run-time statistics, and run-time reports. In the Navigator, you can expand a Data Integration Service to monitor **Ad Hoc Jobs** or expand an application to monitor deployed mapping jobs or workflows

When you select **Ad Hoc Jobs**, deployed mapping jobs, or workflows from an application in the Navigator of the **Execution Statistics** view, a list of jobs appears in the contents panel. The contents panel displays jobs that are in the queued, running, completed, failed, aborted, and cancelled state. The Data Integration Service submits jobs in the queued state to the cluster when resources are available.

The contents panel groups related jobs based on the job type. You can expand a job type to view the related jobs under it.

Access the following views on the content panel under the **Execution Statistics** view:

- Properties
- Hive Execution Plan
- Summary Statistics
- Detailed Statistics

Properties

The **Properties** view on the content panel shows the general properties about the selected job such as name, job type, user who started the job, and start time of the job.

Hive Execution Plan

The Hive execution plan displays the Hive script that the Data Integration Service generates based on the mapping logic. The execution plan includes the Hive queries and Hive commands. Each script has a unique identifier.

Summary Statistics

The **Summary Statistics** view appears in the details panel when you select a mapping job in the contents panel. The **Summary Statistics** view displays throughput and resource usage statistics for the job.

You can view the following throughput statistics for the job:

- Source. The name of the mapping source instance.
- Target name. The name of the target instance.
- Rows. The number of rows read for source and target. If the target is Hive, this is the only summary statistic available.
- Average Rows/Sec. Average number of rows read per second for source and target.
- Bytes. Number of bytes read for source and target.
- Average Bytes/Sec. Average number of bytes read per second for source and target.

- First Row Accessed. The date and time when the Data Integration Service started reading the first row in the source file.
- Dropped rows. Number of source rows that the Data Integration Service did not read.

You can view the throughput statistics for the job in the details pane in the following image:

The screenshot shows the Informatica Administrator interface. The 'Ad Hoc Jobs' list is visible, showing various jobs with their status (Completed, Failed). The 'Summary Statistics' for the job 'm_EachObjectOwnedByDiffUser' are displayed, showing throughput for 'AllHiveSourceTables' and 'Write_test_target_impuser1'.

Source	Rows	Average Rows/Sec	Bytes	Average Bytes/Sec	First Row Accessed	Dropped Rows
AllHiveSourceTables	11	0	0	0	N/A	0

Target	Rows	Average Rows/Sec	Bytes	Average Bytes/Sec	Rejected Rows
Write_test_target_impuser1	5	0	0	0	0

The Hive summary statistics include a row called "AllHiveSourceTables." This row includes records read from the following sources for the MapReduce engine:

- Original Hive sources in the mapping.
- Staging Hive tables defined by the Hive engine.
- Staging data between two linked MapReduce jobs in each query.

If the LDTM session includes one Tez job, the "AllHiveSourceTables" statistics only includes original Hive sources in the mapping.

Note: The AllHiveSourceTables statistics only includes the original Hive sources in a mapping for the Tez job.

When a mapping contains customized data objects or logical data objects, the summary statistics display the original source data instead of the customized data objects or logical data objects in the Administrator tool and in the session log. The Hive driver reads data from the original source data.

You can view the Tez job statistics in the Administrator tool when reading and writing to Hive tables that the Spark engine launches in any of the following scenarios:

- You have resources present in the Amazon buckets.
- You have transactional Hive tables.
- You have table columns secured with fine-grained SQL authorization.

Incorrect statistics appears for all the Hive sources and targets indicating zero rows for average rows for each second, bytes, average bytes for each second, and rejected rows. You can see that only processed rows contain correct values, and the remaining columns will contain either 0 or N/A.

When an Update Strategy transformation runs on the Hive engine, the Summary Statistics for the target table instance combines the number of inserted rows processed, deleted rows processed, and twice the number of updated rows processed. The update operations are handled as separate delete and insert operations.

Detailed Statistics

The **Detailed Statistics** view appears in the details panel when you select a mapping job in the contents panel. The **Detailed Statistics** view displays graphs of the throughput and resource usage statistics for the job run.

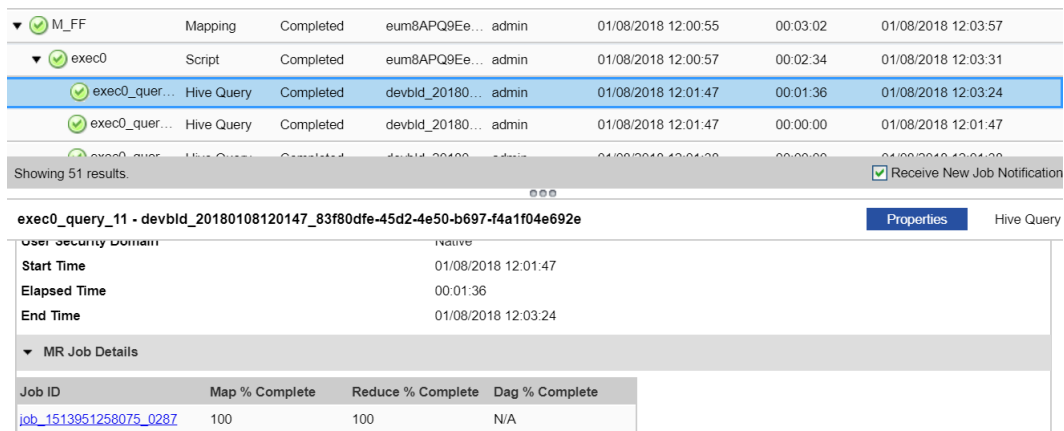
Monitoring with MapReduce Hive Engine

You can monitor the MapReduce Hive engine.

You can also monitor and view Hive tasks that use MapReduce to run Spark jobs. Or, you can monitor MapReduce engines for Hive mappings.

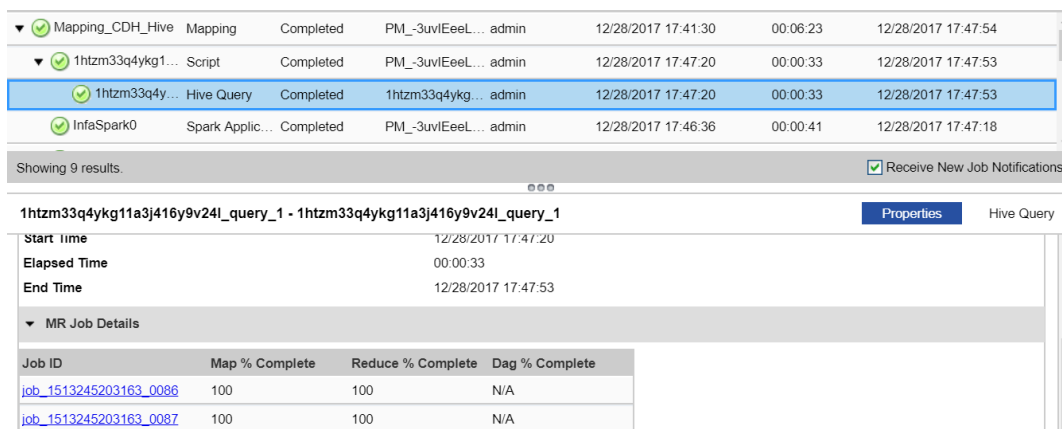
Note: Effective in version 10.2.1, the MapReduce mode of the Hive run-time engine is deprecated, and Informatica will drop support for it in a future release. The Tez mode remains supported.

The following image shows the MapReduce Hive Query properties on the Monitor tab in the Administrator tool:



Job ID	Map % Complete	Reduce % Complete	Dag % Complete
job_1513951258075_0287	100	100	N/A

The following image shows a Hive task that uses MapReduce to run Spark jobs:



Job ID	Map % Complete	Reduce % Complete	Dag % Complete
job_1513245203163_0086	100	100	N/A
job_1513245203163_0087	100	100	N/A

You can view the following information under the MR Job details for MapReduce:

Property	Applicable Values	Description
Job ID	Job_<name>	You can select the link under Job ID to view the application cluster For example, if the Job ID property contains a value starting with the prefix job_ in the MR Job Details pane, the naming convention indicates that the MapReduce engine is in use.
Map % Complete	0 - 100	You can specify a value from 0 through 100 for MapReduce.
Reduce % Complete	0 - 100	You can specify a value from 0 through 100 for MapReduce.
DAG % Complete	N/A	DAG % is not applicable for MapReduce.

Monitoring with Tez Hive Engine

You can monitor Tez Hive engine.

Tez uses YARN timeline as its application history store. Tez stores most of its lifecycle information into the history store, such as all the DAG information. You can monitor the Tez engine information, such as DAG % complete.

Tez relies on the application time line server as a backing store for the application data generated during the lifetime of a YARN application. Tez interfaces with the application timeline server and displays both a live and historical view of the Tez application inside a Tez web application.

The following image shows the Tez Hive Query properties on the Monitor tab in the Administrator tool:

The screenshot shows the Tez Hive Query properties in the Administrator tool. The top table lists several Hive queries, with 'exec0_query_7' selected. Below this, the 'Properties' tab for 'exec0_query_7' is shown, displaying general properties and MR Job Details.

Job ID	Map % Complete	Reduce % Complete	Dag % Complete
application_1513139742016_0083	N/A	N/A	100

You can monitor and view Hive tasks that use Tez to run Spark jobs. Or, you can monitor Tez engines for Hive mappings.

The following image shows a Hive task that uses Tez to run Spark jobs:

✔	TST_Monitoring_S...	Mapping	Completed	IsSjl_FLEeeGy...	Administrator	01/04/2018 18:04:21	00:01:55	01/04/2018 18:06:16	
▼	✔ 7axej03w5j911...	Script	Completed	IsSjl_FLEeeGy...	Administrator	01/04/2018 18:05:57	00:00:17	01/04/2018 18:06:15	
	✔ 7axej03w5j911...	Hive Query	Completed	7axej03w5j911...	Administrator	01/04/2018 18:05:57	00:00:17	01/04/2018 18:06:15	
	✔ InfaSpark0	Spark Applic...	Completed	IsSjl_FLEeeGy...	Administrator	01/04/2018 18:04:48	00:01:06	01/04/2018 18:05:55	
▶	✔ 1mnl17lqrwc1b1...	Script	Completed	IsSjl_FLEeeGy...	Administrator	01/04/2018 18:04:28	00:00:19	01/04/2018 18:04:47	
Showing 116 results.									
<div>◻◻◻</div>									
7axej03w5j911pos4ojk7qo9l_query_1 - 7axej03w5j911pos4ojk7qo9l_query_1								<div>Properties</div>	Hive Query
Start Time		01/04/2018 18:05:57							
Elapsed Time		00:00:17							
End Time		01/04/2018 18:06:15							
▼ MR Job Details									
Job ID		Map % Complete	Reduce % Complete	Dag % Com					
application_1513950138003_4036		N/A	N/A	100					

You can view the following information under the MR Job details for Tez:

Property	Applicable Values	Description
Job ID	Application_<name>	<p>You can select the link under Job ID to view the application cluster.</p> <p>For example, if the Job ID property contains a value starting with the prefix <code>application_</code> in the MR Job Details pane, the naming convention indicates that the Tez engine is in use.</p> <p>You can click the link under Job ID to view the application cluster. If you click the Tracking URL for the Tez job, you get redirected to the Hadoop Resource Manager. If you then click History, you can view the Tez view, which is provided by the Hadoop distribution in Ambari.</p> <p>For each application ID, there are multiple DAGs information.</p>
Map % Complete	N/A	Map % is not applicable for Tez.
Reduce % Complete	N/A	Reduce % is not applicable for Tez.
DAG % Complete	0 - 100	You can specify a value from 0 through 100 for Tez.

When you specify a query in Hive, the script launches a Hadoop job, such as INSERT or DELETE query. Or, the script launches a Hive query. If the script launches no Hadoop jobs, it appears blank for the following fields, such as Job ID, reduce % complete, and DAG % complete.

Note: If the active Resource Manager goes down during a mapping run on the Tez engine, the Tez monitoring statistics might become unavailable for Hive jobs or Spark jobs that use HiveServer 2 tasks.

Hive Engine Logs

The Hive engine logs appear in the LDTM log and the Hive session log.

You can find the information about Hive engine log events in the following log files:

LDTM log

The LDTM logs the results of the Hive queries run for the mapping. You can view the LDTM log from the Developer tool or the Administrator tool for a mapping job.

Hive session log

When you have a Hive script in the Hive execution plan of a mapping, the Data Integration Service opens a Hive session to run the Hive queries.

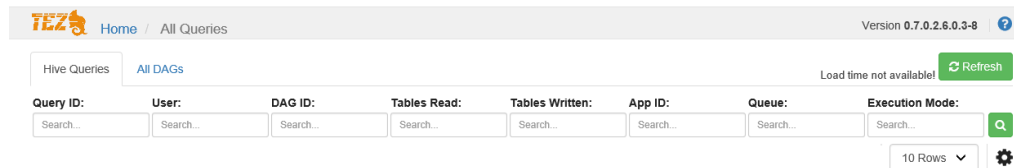
A Hive session updates a log file in the following directory on the Data Integration Service node:

```
<Informatica installation directory>/tomcat/bin/disTemp/.
```

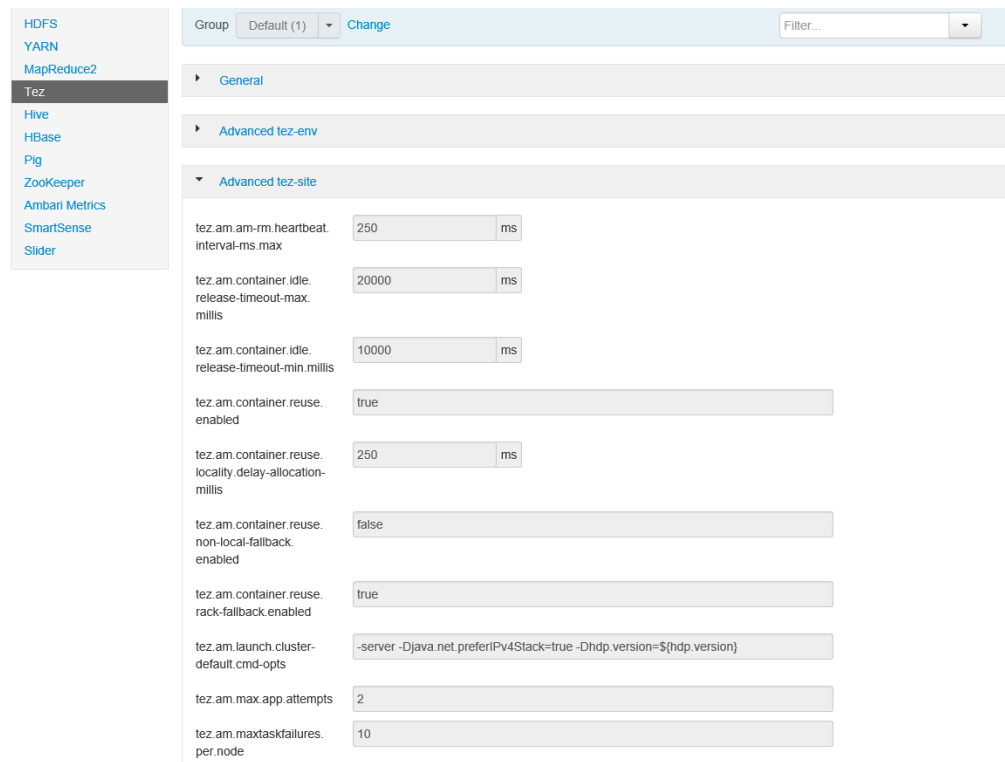
The full path to the Hive session log appears in the LDTM log.

You can view information about DAG vertices in the Tez job link and in the session log. The Tez layout and views might differ based on the selected configurations for the Tez specific properties.

The following image shows the Tez Hive query properties in Tez:



The following image shows the advanced Tez properties in Tez:



The following image shows the advanced Tez properties related to DAG, vertex, and task counts:

tez.session.am.dag.submit.timeout.secs	600	seconds
tez.session.client.timeout.secs	-1	seconds
tez.shuffle-vertex-manager.max-src-fraction	0.4	
tez.shuffle-vertex-manager.min-src-fraction	0.2	
tez.staging-dir	/tmp/\${user.name}/staging	
tez.task.am.heartbeat.counter.interval.ms.max	4000	ms
tez.task.generate.counters.per.io	<input checked="" type="checkbox"/>	
tez.task.get-task.sleep.interval.ms.max	200	ms
tez.task.launch.cluster.default.cmd-opts	-server -Djava.net.preferIPv4Stack=true -Dhdp.version=\${hdp.version}	
tez.task.max-events-per-heartbeat	500	
tez.tez-ui.history-url.base	http://dtmhd26001.informatica.com:8080/#/main/view/TEZ/tez_cluster_instance	
tez.use.cluster.hadoop.libs	<input type="checkbox"/>	

► [Custom tez-site](#)

The monitoring properties appear in the Hive session log under Mapping Status Report when enabled for verbose data or verbose initialization for Tez.

To get DAG tracking URL in the workflow log, you have to update the `tez.tez-ui.history-url.base` with the following value in the HDInsights cluster:

```
<host address>:<port>/#/main/view/TEZ/tez_cluster_instance.
```

For example, a complete DAG URL is as follows:

```
https://ivlhdp584.informatica.com:8443/#/main/view/TEZ/tez_cluster_instance?viewPath=%2F%23%2Fdag%2Fdag_1520917602092_9282_1
```

CHAPTER 11

Mappings in the Native Environment

This chapter includes the following topics:

- [Mappings in the Native Environment Overview, 181](#)
- [Data Processor Mappings, 181](#)
- [HDFS Mappings, 182](#)
- [Hive Mappings, 183](#)
- [Social Media Mappings, 184](#)

Mappings in the Native Environment Overview

You can run a mapping in the native environment. In the native environment, the Data Integration Service runs the mapping from the Developer tool. You can run standalone mappings or mappings that are a part of a workflow.

In the native environment, you can read and process data from large unstructured and semi-structured files, Hive, or social media web sites. You can include the following objects in the mappings:

- Hive sources
- Flat file sources or targets in the local system or in HDFS
- Complex file sources in the local system or in HDFS
- Data Processor transformations to process unstructured and semi-structured file formats
- Social media sources

Data Processor Mappings

The Data Processor transformation processes unstructured and semi-structured file formats in a mapping. It converts source data to flat CSV records that MapReduce applications can process.

You can configure the Data Processor transformation to process messaging formats, relational data, HTML pages, XML, JSON, AVRO, Parquet, Cobol, Microsoft Excel, Microsoft Word, and PDF documents. You can

also configure it to transform structured formats such as ACORD, HIPAA, HL7, EDI-X12, EDIFACT, AFP, and SWIFT.

For example, an application produces hundreds of data files per second and writes the files to a directory. You can create a mapping that extracts the files from the directory, passes them to a Data Processor transformation, and writes the data to a target.

HDFS Mappings

Create an HDFS mapping to read or write to HDFS.

You can read and write fixed-width and delimited file formats. You can read or write compressed files. You can read text files and binary file formats such as sequence file from HDFS. You can specify the compression format of the files. You can use the binary stream output of the complex file data object as input to a Data Processor transformation to parse the file.

You can define the following objects in an HDFS mapping:

- Flat file data object or complex file data object operation as the source to read data from HDFS.
- Transformations.
- Flat file data object as the target to write data to HDFS or any target.

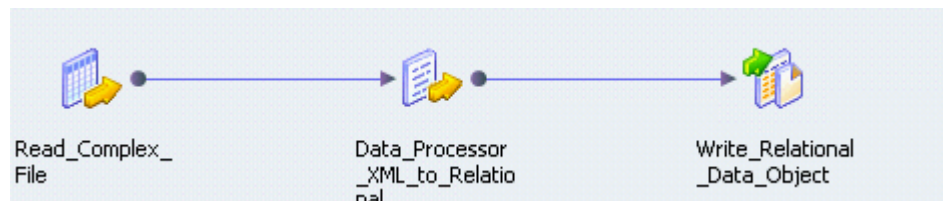
Validate and run the mapping. You can deploy the mapping and run it or add the mapping to a Mapping task in a workflow.

HDFS Data Extraction Mapping Example

Your organization needs to analyze purchase order details such as customer ID, item codes, and item quantity. The purchase order details are stored in a semi-structured compressed XML file in HDFS. The hierarchical data includes a purchase order parent hierarchy level and a customer contact details child hierarchy level. Create a mapping that reads all the purchase records from the file in HDFS. The mapping must convert the hierarchical data to relational data and write it to a relational target.

You can use the extracted data for business analytics.

The following figure shows the example mapping:



You can use the following objects in the HDFS mapping:

HDFS Input

The input object, Read_Complex_File, is a Read transformation that represents a compressed XML file stored in HDFS.

Data Processor Transformation

The Data Processor transformation, Data_Processor_XML_to_Relational, parses the XML file and provides a relational output.

Relational Output

The output object, `Write_Relational_Data_Object`, is a Write transformation that represents a table in an Oracle database.

When you run the mapping, the Data Integration Service reads the file in a binary stream and passes it to the Data Processor transformation. The Data Processor transformation parses the specified file and provides a relational output. The output is written to the relational target.

You can configure the mapping to run in a native or Hadoop run-time environment.

Complete the following tasks to configure the mapping:

1. Create an HDFS connection to read files from the Hadoop cluster.
2. Create a complex file data object read operation. Specify the following parameters:
 - The file as the resource in the data object.
 - The file compression format.
 - The HDFS file location.
3. Optionally, you can specify the input format that the Mapper uses to read the file.
4. Drag and drop the complex file data object read operation into a mapping.
5. Create a Data Processor transformation. Configure the following properties in the Data Processor transformation:
 - An input port set to buffer input and binary data type.
 - Relational output ports depending on the number of columns you want in the relational output. Specify the port size for the ports. Use an XML schema reference that describes the XML hierarchy. Specify the normalized output that you want. For example, you can specify `PurchaseOrderNumber_Key` as a generated key that relates the Purchase Orders output group to a Customer Details group.
 - Create a Streamer object and specify Streamer as a startup component.
6. Create a relational connection to an Oracle database.
7. Import a relational data object.
8. Create a write transformation for the relational data object and add it to the mapping.

Hive Mappings

Based on the mapping environment, you can read data from or write data to Hive.

In a native environment, you can read data from Hive. To read data from Hive, complete the following steps:

1. Create a Hive connection.
2. Configure the Hive connection mode to access Hive as a source or target.
3. Use the Hive connection to create a data object to read from Hive.
4. Add the data object to a mapping and configure the mapping to run in the native environment.

You can write to Hive in a Hadoop environment. To write data to Hive, complete the following steps:

1. Create a Hive connection.
2. Configure the Hive connection mode to access Hive as a source or target.

3. Use the Hive connection to create a data object to write to Hive.
4. Add the data object to a mapping and configure the mapping to run in the Hadoop environment.

You can define the following types of objects in a Hive mapping:

- A Read Transformation to read data from Hive
- Transformations
- A target or an SQL data service. You can write to Hive if you run the mapping in a Hadoop cluster.

Validate and run the mapping. You can deploy the mapping and run it or add the mapping to a Mapping task in a workflow.

Hive Mapping Example

Your organization, HypoMarket Corporation, needs to analyze customer data. Create a mapping that reads all the customer records. Create an SQL data service to make a virtual database available for end users to query.

You can use the following objects in a Hive mapping:

Hive input

The input file is a Hive table that contains the customer names and contact details.

Create a relational data object. Configure the Hive connection and specify the table that contains the customer data as a resource for the data object. Drag the data object into a mapping as a read data object.

SQL Data Service output

Create an SQL data service in the Developer tool. To make it available to end users, include it in an application, and deploy the application to a Data Integration Service. When the application is running, connect to the SQL data service from a third-party client tool by supplying a connect string.

You can run SQL queries through the client tool to access the customer data.

Social Media Mappings

Create mappings to read social media data from sources such as Facebook and LinkedIn.

You can extract social media data and load them to a target in the native environment only. You can choose to parse this data or use the data for data mining and analysis.

To process or analyze the data in Hadoop, you must first move the data to a relational or flat file target and then run the mapping in the Hadoop cluster.

You can use the following Informatica adapters in the Developer tool:

- PowerExchange for DataSift
- PowerExchange for Facebook
- PowerExchange for LinkedIn
- PowerExchange for Twitter
- PowerExchange for Web Content-Kapow Katalyst

Review the respective PowerExchange adapter documentation for more information.

Twitter Mapping Example

Your organization, Hypomarket Corporation, needs to review all the tweets that mention your product HypoBasket with a positive attitude since the time you released the product in February 2012.

Create a mapping that identifies tweets that contain the word HypoBasket and writes those records to a table.

You can use the following objects in a Twitter mapping:

Twitter input

The mapping source is a Twitter data object that contains the resource Search.

Create a physical data object and add the data object to the mapping. Add the Search resource to the physical data object. Modify the query parameter with the following query:

```
QUERY=HypoBasket:)&since:2012-02-01
```

Sorter transformation

Optionally, sort the data based on the timestamp.

Add a Sorter transformation to the mapping. Specify the timestamp as the sort key with direction as ascending.

Mapping output

Add a relational data object to the mapping as a target.

After you run the mapping, Data Integration Service writes the extracted tweets to the target table. You can use text analytics and sentiment analysis tools to analyze the tweets.

CHAPTER 12

Profiles

This chapter includes the following topics:

- [Profiles Overview, 186](#)
- [Native Environment, 186](#)
- [Hadoop Environment, 187](#)
- [Creating a Single Data Object Profile in Informatica Developer, 188](#)
- [Creating an Enterprise Discovery Profile in Informatica Developer, 189](#)
- [Creating a Column Profile in Informatica Analyst, 190](#)
- [Creating an Enterprise Discovery Profile in Informatica Analyst, 191](#)
- [Creating a Scorecard in Informatica Analyst, 192](#)
- [Monitoring a Profile, 193](#)
- [Troubleshooting, 194](#)

Profiles Overview

You can create and run column profiles, enterprise discovery profiles, and scorecards in the native run-time environment or Hadoop run-time environment.

When you create or edit a profile or scorecard, you can choose the run-time environment. After you choose the run-time environment, the Developer tool or the Analyst tool sets the run-time environment in the profile definition. To process the profiles quickly, you can choose the Hadoop run-time environment.

Native Environment

In Informatica Developer, you can run single object profiles, multiple object profiles, and enterprise discovery profiles in the native environment. In Informatica Analyst, you can run column profiles, enterprise discovery profiles, and scorecards in the native environment.

When you run a profile in the native run-time environment, the Analyst tool or Developer tool submits the profile jobs to the Profiling Service Module. The Profiling Service Module then breaks down the profile jobs into a set of mappings. The Data Integration Service runs these mappings on the same machine where the Data Integration Service runs and writes the profile results to the profiling warehouse. By default, all profiles run in the native run-time environment.

You can use native sources to create and run profiles in the native environment. A native data source is a non-Hadoop source, such as a flat file, relational source, or mainframe source. You can also run a profile on a mapping specification or a logical data source with a Hive or HDFS data source in the native environment.

Hadoop Environment

You can run profiles and scorecards in the Hadoop environment on the Hive engine or Blaze engine. You can choose Hive or Hadoop option to run the profiles in the Hadoop run-time environment. After you choose the Hive option and select a Hadoop connection, the Data Integration Service pushes the profile logic to the Hive engine on the Hadoop cluster to run profiles. After you choose the Hadoop option and select a Hadoop connection, the Data Integration Service pushes the profile logic to the Blaze engine on the Hadoop cluster to run profiles.

When you run a profile in the Hadoop environment, the Analyst tool submits the profile jobs to the Profiling Service Module. The Profiling Service Module then breaks down the profile jobs into a set of mappings. The Data Integration Service pushes the mappings to the Hadoop environment through the Hadoop connection. The Hive engine or Blaze engine processes the mappings and the Data Integration Service writes the profile results to the profiling warehouse.

In the Developer tool, you can run single object profiles and multiple object profiles, and enterprise discovery profiles on the Blaze engine. In the Analyst tool, you can run column profiles, enterprise discovery profiles, and scorecards on the Blaze engine.

Column Profiles for Sqoop Data Sources

You can run a column profile on data objects that use Sqoop. You can select the Hive or Hadoop run-time environment to run the column profiles.

On the Hive engine, to run a column profile on a relational data object that uses Sqoop, you must set the Sqoop argument `m` to 1 in the JDBC connection. Use the following syntax:

```
-m 1
```

When you run a column profile on a logical data object or customized data object, you can configure the `num-mappers` argument to achieve parallelism and optimize performance. You must also configure the `split-by` argument to specify the column based on which Sqoop must split the work units.

Use the following syntax:

```
--split-by <column_name>
```

If the primary key does not have an even distribution of values between the minimum and maximum range, you can configure the `split-by` argument to specify another column that has a balanced distribution of data to split the work units.

If you do not define the `split-by` column, Sqoop splits work units based on the following criteria:

- If the data object contains a single primary key, Sqoop uses the primary key as the `split-by` column.
- If the data object contains a composite primary key, Sqoop defaults to the behavior of handling composite primary keys without the `split-by` argument. See the Sqoop documentation for more information.
- If a data object contains two tables with an identical column, you must define the `split-by` column with a table-qualified name. For example, if the table name is `CUSTOMER` and the column name is `FULL_NAME`, define the `split-by` column as follows:

```
--split-by CUSTOMER.FULL_NAME
```

- If the data object does not contain a primary key, the value of the `m` argument and `num-mappers` argument default to 1.

When you use Cloudera Connector Powered by Teradata or Hortonworks Connector for Teradata and the Teradata table does not contain a primary key, the `split-by` argument is required.

Creating a Single Data Object Profile in Informatica Developer

You can create a single data object profile for one or more columns in a data object and store the profile object in the Model repository.

1. In the **Object Explorer** view, select the data object you want to profile.
2. Click **File > New > Profile** to open the profile wizard.
3. Select **Profile** and click **Next**.
4. Enter a name for the profile and verify the project location. If required, browse to a new location.
5. Optionally, enter a text description of the profile.
6. Verify that the name of the data object you selected appears in the **Data Objects** section.
7. Click **Next**.
8. Configure the profile operations that you want to perform. You can configure the following operations:
 - Column profiling
 - Primary key discovery
 - Functional dependency discovery
 - Data domain discovery

Note: To enable a profile operation, select **Enabled as part of the "Run Profile" action** for that operation. Column profiling is enabled by default.

9. Review the options for your profile.
You can edit the column selection for all profile types. Review the filter and sampling options for column profiles. You can review the inference options for primary key, functional dependency, and data domain discovery. You can also review data domain selection for data domain discovery.
10. Review the drill-down options, and edit them if necessary. By default, the **Enable Row Drilldown** option is selected. You can edit drill-down options for column profiles. The options also determine whether drill-down operations read from the data source or from staged data, and whether the profile stores result data from previous profile runs.
11. In the **Run Settings** section, choose a run-time environment. Choose **Native**, **Hive (deprecated)**, or **Hadoop** as the run-time environment. When you choose the Hive or Hadoop option, select a Hadoop connection.
12. Click **Finish**.

Creating an Enterprise Discovery Profile in Informatica Developer

You can create a profile on multiple data sources under multiple connections. The Developer tool creates individual profile tasks for each source.

1. In the **Object Explorer** view, select multiple data objects you want to run a profile on.
2. Click **File > New > Profile** to open the profile wizard.
3. Select **Enterprise Discovery Profile** and click **Next**.
4. Enter a name for the profile and verify the project location. If required, browse to a new location.
5. Verify that the name of the data objects you selected appears within the **Data Objects** section. Click **Choose** to select more data objects, if required.
6. Click **Next**.

The **Add Resources to Profile Definition** pane appears. You can select multiple, external relational connections and data sources from this pane.

7. Click **Choose** to open the **Select Resources** dialog box.

The **Resources** pane lists all the internal and external connections and data objects under the Informatica domain.

8. Click **OK** to close the dialog box.
9. Click **Next**.

10. Configure the profile types that you want to run. You can configure the following profile types:

- Data domain discovery
- Column profile
- Primary key profile
- Foreign key profile

Note: Select **Enabled as part of "Run Enterprise Discovery Profile" action** for the profile types that you want to run as part of the enterprise discovery profile. Column profiling is enabled by default.

11. Review the options for the profile.

You can edit the sampling options for column profiles. You can also edit the inference options for data domain, primary key, and foreign key profiles.

12. Select **Create profiles**.

The Developer tool creates profiles for each individual data source.

13. Select **Run enterprise discovery profile on finish** to run the profile when you complete the profile configuration. If you enabled all the profiling operations, the Developer tool runs column, data domain, and primary key profiles on all selected data sources. Then, the Developer tool runs a foreign key profile across all the data sources.

14. Click **Finish**.

After you run an enterprise discovery profile, you need to refresh the Model Repository Service before viewing the results. This step is required as the import of metadata for external connections happens in the Model repository. You need to refresh the Model Repository Service so that the Developer tool reflects the changes to the Model repository.

Creating a Column Profile in Informatica Analyst

You can create a custom profile or default profile. When you create a custom profile, you can configure the columns, sample rows, and drill-down options. When you create a default profile, the column profile and data domain discovery runs on the entire data set with all the data domains.

1. In the **Discovery** workspace, click **Profile**, or select **New > Profile** from the header area.

Note: You can right-click on the data object in the **Library** workspace and create a profile. In this profile, the profile name, location name, and data object are extracted from the data object properties. You can create a default profile or customize the settings to create a custom profile.

The **New Profile** wizard appears.

2. The **Single source** option is selected by default. Click **Next**.
3. In the **Specify General Properties** screen, enter a name and an optional description for the profile. In the Location field, select the project or folder where you want to create the profile. Click **Next**.
4. In the **Select Source** screen, click **Choose** to select a data object, or click **New** to import a data object. Click **Next**.

- In the **Choose Data Object** dialog box, select a data object. Click **OK**.
The Properties pane displays the properties of the selected data object. The Data Preview pane displays the columns in the data object.
- In the **New Data Object** dialog box, you can choose a connection, schema, table, or view to create a profile on, select a location, and create a folder to import the data object. Click **OK**.

5. In the **Select Source** screen, select the columns that you want to run a profile on. Optionally, select **Name** to select all the columns. Click **Next**.

All the columns are selected by default. The Analyst tool lists column properties, such as the name, data type, precision, scale, nullable, and participates in the primary key for each column.

6. In the **Specify Settings** screen, choose to run a column profile, data domain discovery, or a column profile and data domain discovery. By default, column profile option is selected.
 - Choose **Run column profile** to run a column profile.
 - Choose **Run data domain discovery** to perform data domain discovery. In the **Data domain** pane, select the data domains that you want to discover, select a conformance criteria, and select the columns for data domain discovery in the **Edit columns selection for data domain discovery** dialog box.
 - Choose **Run column profile** and **Run data domain discovery** to run the column profile and data domain discovery. Select the data domain options in the **Data domain** pane.

Note: By default, the columns that you select is for column profile and data domain discovery. Click **Edit** to select or deselect columns for data domain discovery.
 - Choose **Data**, **Columns**, or **Data and Columns** to run data domain discovery on.
 - Choose a sampling option. You can choose **All rows (complete analysis)**, **Sample first**, **Random sample**, or **Random sample (auto)** as a sampling option in the **Run profile on** pane. This option applies to column profile and data domain discovery.
 - Choose a drilldown option. You can choose **Live** or **Staged** drilldown option, or you can choose **Off** to disable drilldown in the **Drilldown** pane. Optionally, click **Select Columns** to select columns to drill down on. You can choose to omit data type and data domain inference for columns with an approved data type or data domain.

- Choose **Native**, **Hive (deprecated)**, or **Hadoop** option as the run-time environment. If you choose the Hive or Hadoop option, click **Choose** to select a Hadoop connection in the **Select a Hadoop Connection** dialog box.
7. Click **Next**.
The **Specify Rules and Filters** screen opens.
 8. In the **Specify Rules and Filters** screen, you can perform the following tasks:
 - Create, edit, or delete a rule. You can apply existing rules to the profile.
 - Create, edit, or delete a filter.

Note: When you create a scorecard on this profile, you can reuse the filters that you create for the profile.
 9. Click **Save and Finish** to create the profile, or click **Save and Run** to create and run the profile.

Creating an Enterprise Discovery Profile in Informatica Analyst

You can run column profile and data domain discovery as part of enterprise discovery in Informatica Analyst.

1. In the **Discovery** workspace, select **New > Profile**.
The **New Profile** wizard appears.
2. Select **Enterprise Discovery**. Click **Next**.
The **Specify General Properties** tab appears.
3. In the **Specify General Properties** tab, enter a name for the enterprise discovery profile and an optional description. In the Location field, select the project or folder where you want to create the profile. Click **Next**.
The **Select Data Objects** tab appears.
4. In the **Select Data Objects** tab, click **Choose**.
The **Choose Data objects** dialog box appears.
5. In the **Choose Data objects** dialog box, choose one or more data objects to add to the profile. Click **Save**.
The data objects appear in the **Data Objects** pane.
6. Click **Next**.
The **Select Resources** tab appears.
7. In the **Select Resources** tab, click **Choose** to open the **Select Resources** tab.
You can import data from multiple relational data sources.
8. In the **Select Resources** tab, select the connections, schemas, tables, and views that you want to include in the profile. Click **Save**.
The left pane in the dialog box lists all the internal and external connections, schemas, tables, and views under the Informatica domain.
The resources appear in the **Resource** pane.
9. Click **Next**.
The **Specify Settings** tab appears.

10. In the **Specify Settings** tab, you can configure the column profile options and data domain discovery options. Click **Save and Finish** to save the enterprise discovery profile, or click **Save and Run** to run the profile.

You can perform the following tasks in the **Specify Settings** tab.

- Enable data domain discovery. Click **Choose** to select data domains that you want to discover from the **Choose Data Domains** dialog box. The selected data domains appear in the **Data Domains for Data Domain Discovery** pane.
- Run data domain on data, column name, or on both data and column name.
- Select all the rows in the data source, or choose a maximum number of rows to run domain discovery on.
- Choose a minimum conformance percentage or specify the minimum number of conforming rows for data domain discovery.
- Enable column profile settings and select all rows or first few rows in the data source for the column profile. You can exclude data type inference for columns with approved data types in the column profile.
- Choose **Native** or **Hadoop** as the run-time environment.

You can view the enterprise discovery results under the **Summary** and **Profiles** tabs.

Creating a Scorecard in Informatica Analyst

Create a scorecard and add columns from a profile to the scorecard. You must run a profile before you add columns to the scorecard.

1. In the **Library** workspace, select the project or folder that contains the profile.
2. Click the profile to open the profile.

The profile results appear in the summary view in the **Discovery** workspace.

3. Click **Actions > Add to scorecard**.

The **Add to Scorecard** wizard appears.

4. In the **Add to Scorecard** screen, you can choose to create a new scorecard, or edit an existing scorecard to add the columns to a predefined scorecard. The **New Scorecard** option is selected by default. Click **Next**.

5. In the **Step 2 of 8** screen, enter a name for the scorecard. Optionally, you can enter a description for the scorecard. Select the project and folder where you want to save the scorecard. Click **Next**.

By default, the scorecard wizard selects the columns and rules defined in the profile. You cannot add columns that are not included in the profile.

6. In the **Step 3 of 8** screen, select the columns and rules that you want to add to the scorecard as metrics. Optionally, click the check box in the left column header to select all columns. Optionally, select **Column Name** to sort column names. Click **Next**.

7. In the **Step 4 of 8** screen, you can add a filter to the metric.

You can apply the filter that you created for the profile to the metrics, or create a new filter. Select a metric in the **Metric Filters** pane, and click the **Manage Filters** icon to open the **Edit Filter: column name** dialog box. In the **Edit Filter: column name** dialog box, you can choose to perform one of the following tasks:

- Choose a filter that you created for the profile. Click **Next**.
- Select an existing filter. Click the edit icon to edit the filter in the **Edit Filter** dialog box. Click **Next**.
- Click the plus (+) icon to create filters in the **New Filter** dialog box. Click **Next**.

Optionally, you can choose to apply the selected filters to all the metrics in the scorecard.

The filter appears in the **Metric Filters** pane.

8. In the **Step 4 of 8** screen, click **Next**.
9. In the **Step 5 of 8** screen, select each metric in the **Metrics** pane to perform the following tasks:
 - Configure valid values. In the **Score using: Values** pane, select one or more values in the **Available Values** pane, and click the right arrow button to move them to the **Valid Values** pane. The total number of valid values for a metric appears at the top of the **Available Values** pane.
 - Configure metric thresholds. In the **Metric Thresholds** pane, set the thresholds for **Good**, **Acceptable**, and **Unacceptable** scores.
 - Configure the cost of invalid data. To assign a constant value to the cost for the metric, select **Fixed Cost**. To attach a numeric column as a variable cost to the metric, select **Variable Cost**, and click **Select Column** to select a numeric column. Optionally, click **Change Cost Unit** to change the unit of cost. If you do not want to configure the cost of invalid data for the metric, choose **None**.
10. Click **Next**.
11. In the **Step 6 of 8** screen, you can select a metric group to which you can add the metrics, or create a new metric group. To create a new metric group, click the group icon. Click **Next**.
12. In the **Step 7 of 8** screen, specify the weights for the metrics in the group and thresholds for the group.
13. In the **Step 8 of 8** screen, select **Native** or **Hadoop** run-time environment option to run the scorecard. If you choose the Hadoop option, click **Browse** to choose a Hadoop connection to run the profile on the Blaze engine.
14. Click **Save** to save the scorecard, or click **Save & Run** to save and run the scorecard.

The scorecard appears in the **Scorecard** workspace.

Monitoring a Profile

You can monitor a profile in the Administrator tool.

1. In the Administrator tool, click the **Monitor** tab.
2. In the Navigator workspace, select **Jobs**.
3. Select a profiling job.
4. In the **Summary Statistics** tab, you can view the general properties of the profile, summary statistics, and detailed statistics of the profile.
5. Click the **Execution Statistics** tab to view execution statistics for the profile.

Troubleshooting

Can I drill down on profile results if I run a profile in the Hadoop environment?

Yes, except for profiles in which you have set the option to drill down on staged data.

I get the following error message when I run a profile in the Hadoop environment: "[LDTM_1055] The Integration Service failed to generate a Hive workflow for mapping [Profile_CUSTOMER_INFO12_14258652520457390]." How do I resolve this?

This error can result from a data source, rule transformation, or run-time environment that is not supported in the Hadoop environment. For more information about objects that are not valid in the Hadoop environment, see the Mappings in a Hadoop Environment chapter.

You can change the data source, rule, or run-time environment and run the profile again. View the profile log file for more information on the error.

I see "N/A" in the profile results for all columns after I run a profile. How do I resolve this?

Verify that the profiling results are in the profiling warehouse. If you do not see the profile results, verify that the database path is accurate in the Cluster Environment Variables property of the Hadoop connection. You can also verify the database path from the Hadoop job tracker on the Monitoring tab of the Administrator tool.

After I run a profile on a Hive source, I do not see the results. When I verify the Hadoop job tracker, I see the following error when I open the profile job: "XML Parsing Error: no element found." What does this mean?

The Hive data source does not have any record and is empty. The data source must have a minimum of one row of data for successful profile run.

After I run a profile on a Hive source, I cannot view some of the column patterns. Why?

When you import a Hive source, the Developer tool sets the precision for string columns to 4000. The Developer tool cannot derive the pattern for a string column with a precision greater than 255. To resolve this issue, set the precision of these string columns in the data source to 255 and run the profile again.

When I run a profile on large Hadoop sources, the profile job fails and I get an "execution failed" error. What can be the possible cause?

One of the causes can be a connection issue. Perform the following steps to identify and resolve the connection issue:

1. Open the Hadoop job tracker.
2. Identify the profile job and open it to view the MapReduce jobs.
3. Click the hyperlink for the failed job to view the error message. If the error message contains the text "java.net.ConnectException: Connection refused", the problem occurred because of an issue with the Hadoop cluster. Contact your network administrator to resolve the issue.

CHAPTER 13

Native Environment Optimization

This chapter includes the following topics:

- [Native Environment Optimization Overview, 195](#)
- [Processing Big Data on a Grid, 195](#)
- [Processing Big Data on Partitions, 196](#)
- [High Availability, 197](#)

Native Environment Optimization Overview

You can optimize the native environment to increase performance. To increase performance, you can configure the Data Integration Service to run on a grid and to use multiple partitions to process data. You can also enable high availability to ensure that the domain can continue running despite temporary network, hardware, or service failures.

You can run profiles, sessions, and workflows on a grid to increase the processing bandwidth. A grid is an alias assigned to a group of nodes that run profiles, sessions, and workflows. When you enable grid, the Data Integration Service runs a service process on each available node of the grid to increase performance and scalability.

You can also run mapping with partitioning to increase performance. When you run a partitioned session or a partitioned mapping, the Data Integration Service performs the extract, transformation, and load for each partition in parallel.

You can configure high availability for the domain. High availability eliminates a single point of failure in a domain and provides minimal service interruption in the event of failure.

Processing Big Data on a Grid

You can run an Integration Service on a grid to increase the processing bandwidth. When you enable grid, the Integration Service runs a service process on each available node of the grid to increase performance and scalability.

Big data may require additional bandwidth to process large amounts of data. For example, when you run a Model repository profile on an extremely large data set, the Data Integration Service grid splits the profile into multiple mappings and runs the mappings simultaneously on different nodes in the grid.

Data Integration Service Grid

You can run Model repository mappings and profiles on a Data Integration Service grid.

When you run mappings on a grid, the Data Integration Service distributes the mappings to multiple DTM processes on nodes in the grid. When you run a profile on a grid, the Data Integration Service splits the profile into multiple mappings and distributes the mappings to multiple DTM processes on nodes in the grid.

For more information about the Data Integration Service grid, see the *Informatica Administrator Guide*.

Grid Optimization

You can optimize the grid to increase performance and scalability of the Data Integration Service.

To optimize the grid, complete the following task:

Add nodes to the grid.

Add nodes to the grid to increase processing bandwidth of the Data Integration Service.

Processing Big Data on Partitions

You can run a Model repository mapping with partitioning to increase performance. When you run a mapping configured with partitioning, the Data Integration Service performs the extract, transformation, and load for each partition in parallel.

Mappings that process large data sets can take a long time to process and can cause low data throughput. When you configure partitioning, the Data Integration Service uses additional threads to process the session or mapping which can increase performance.

Partitioned Model Repository Mappings

You can enable the Data Integration Service to use multiple partitions to process Model repository mappings.

If the nodes where mappings run have multiple CPUs, you can enable the Data Integration Service to maximize parallelism when it runs mappings. When you maximize parallelism, the Data Integration Service dynamically divides the underlying data into partitions and processes all of the partitions concurrently.

Optionally, developers can set a maximum parallelism value for a mapping in the Developer tool. By default, the maximum parallelism for each mapping is set to Auto. Each mapping uses the maximum parallelism value defined for the Data Integration Service. Developers can change the maximum parallelism value in the mapping run-time properties to define a maximum value for a particular mapping. When maximum parallelism is set to different integer values for the Data Integration Service and the mapping, the Data Integration Service uses the minimum value.

For more information, see the *Informatica Application Services Guide* and the *Informatica Developer Mapping Guide*.

Partition Optimization

You can optimize the partitioning of Model repository mappings to increase performance. You can add more partitions, select the best performing partition types, use more CPUs, and optimize the source or target database for partitioning.

To optimize partitioning, perform the following tasks:

Increase the number of partitions.

When you configure Model repository mappings, you increase the number of partitions when you increase the maximum parallelism value for the Data Integration Service or the mapping.

Increase the number of partitions to enable the Data Integration Service to create multiple connections to sources and process partitions of source data concurrently. Increasing the number of partitions increases the number of threads, which also increases the load on the Data Integration Service nodes. If the Data Integration Service node or nodes contain ample CPU bandwidth, processing rows of data concurrently can increase performance.

Note: If you use a single-node Data Integration Service and the Data Integration Service uses a large number of partitions in a session or mapping that processes large amounts of data, you can overload the system.

Use multiple CPUs.

If you have a symmetric multi-processing (SMP) platform, you can use multiple CPUs to concurrently process partitions of data.

Optimize the source database for partitioning.

You can optimize the source database for partitioning. For example, you can tune the database, enable parallel queries, separate data into different tablespaces, and group sorted data.

Optimize the target database for partitioning.

You can optimize the target database for partitioning. For example, you can enable parallel inserts into the database, separate data into different tablespaces, and increase the maximum number of sessions allowed to the database.

High Availability

High availability eliminates a single point of failure in an Informatica domain and provides minimal service interruption in the event of failure. When you configure high availability for a domain, the domain can continue running despite temporary network, hardware, or service failures. You can configure high availability for the domain, application services, and application clients.

The following high availability components make services highly available in an Informatica domain:

- **Resilience.** An Informatica domain can tolerate temporary connection failures until either the resilience timeout expires or the failure is fixed.
- **Restart and failover.** A process can restart on the same node or on a backup node after the process becomes unavailable.
- **Recovery.** Operations can complete after a service is interrupted. After a service process restarts or fails over, it restores the service state and recovers operations.

When you plan a highly available Informatica environment, consider the differences between internal Informatica components and systems that are external to Informatica. Internal components include the

Service Manager, application services, and command line programs. External systems include the network, hardware, database management systems, FTP servers, message queues, and shared storage.

High availability features for the Informatica environment are available based on your license.

CHAPTER 14

Cluster Workflows

This chapter includes the following topics:

- [Cluster Workflows Overview, 199](#)
- [Cluster Workflow Components, 200](#)
- [Cluster Workflows Process , 201](#)
- [Administrator Tasks, 202](#)
- [Create the Cluster Workflow , 202](#)
- [Deploy and Run the Workflow, 211](#)

Cluster Workflows Overview

You can use a workflow to create a cluster that runs Mapping and other tasks on a cloud platform cluster.

A cluster workflow contains a Create Cluster task that you configure with information about the cluster to create. The cluster workflow uses other elements that enable communication between the Data Integration Service and the cloud platform, such as a cloud provisioning configuration and a Hadoop connection.

If you want to create an ephemeral cluster, you can include a Delete Cluster task. An ephemeral cluster is a cloud platform cluster that you create and use for running mappings and other tasks, then terminate when tasks are complete to save cloud platform resources.

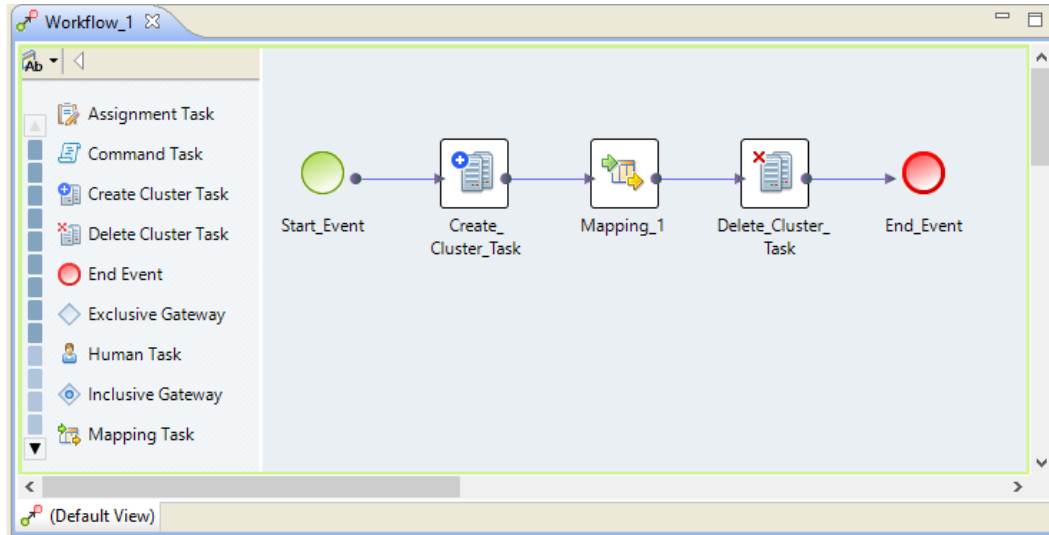
You can use a cluster workflow with the Amazon EMR or Microsoft Azure HDInsight cloud platforms.

Note: To create a cluster on Cloudera Altus, you create a workflow with Command tasks that perform the tasks that a cluster workflow automates. For more information about creating a cluster on Cloudera Altus, see the article "Implementing Informatica Big Data Management with Ephemeral Clusters on a Cloudera Altus Cluster" on the Informatica Network.

Cluster Workflow Components

The cluster workflow creates an ephemeral cluster that includes a Create Cluster task, at least one Mapping task, and a Delete Cluster task.

The following image shows a sample cluster workflow:



Create Cluster Task

The Create Cluster task contains all the settings that Amazon EMR or Azure HDInsight requires to create a cluster with a master node and worker nodes. It also contains a reference to a Hadoop connection.

Create one Create Cluster task for a cluster workflow.

When you create a cluster workflow, you drag a Create Cluster task into the workflow editor, and then configure task properties.

Cloud Provisioning Configuration

The cloud provisioning configuration associates the Create Cluster task with the Hadoop connection associated with the workflow.

Configure the cloud provisioning configuration with information about the cloud platform account. To configure the cloud provisioning configuration, use the Administrator tool.

The Create Cluster task must include a reference to the cloud provisioning configuration.

For more information, see the *Big Data Management Administrator Guide*.

Hadoop Connection

Create a dedicated Hadoop connection to use with the cluster.

The Hadoop connection saves property values for the Data Integration Service to use for cluster workflow operations. When you run a cluster workflow, the Data Integration Service creates temporary Hadoop connections based on these values.

Mapping and Other Workflow Tasks

After the Create Cluster task, add other workflow tasks to run on the cluster.

Add one or more Mapping tasks to the workflow. You can also include Command and other workflow tasks.

Before you add mappings to Mapping tasks, you prepare the mappings. Configure the Hadoop Connection property in each mapping to designate whether to run the mapping on the cluster that the workflow creates, or on another cluster.

Delete Cluster Task

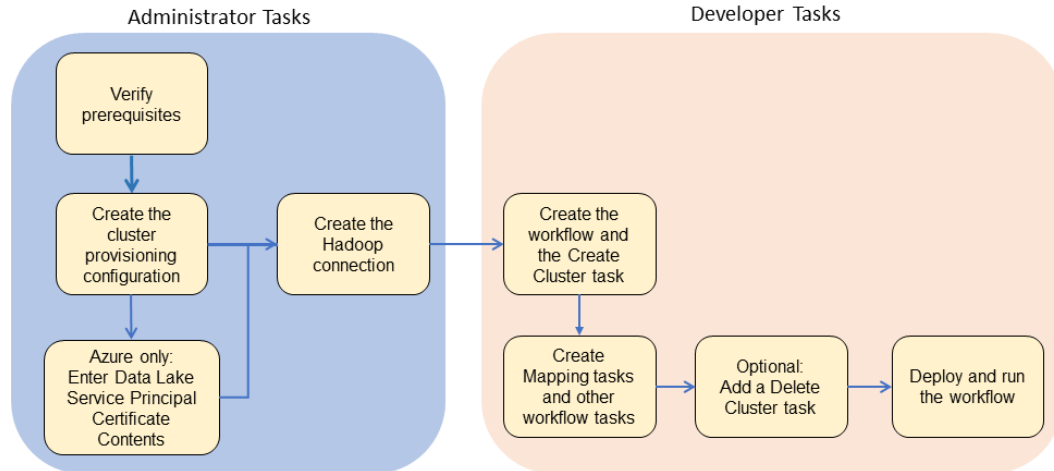
The Delete Cluster task terminates the cluster, and deletes the cluster and other resources that the cluster workflow creates.

To create an ephemeral cluster, add a Delete Cluster task to the workflow. When you do not include it, the cluster continues to run until you terminate it. You can use `infacmd ccps deleteCluster` to terminate the cluster at any time.

Cluster Workflows Process

Creation of a cluster workflow requires administrator and developer tasks.

The following image shows the process to create, configure, and run a cluster workflow:



First, an administrator completes the following steps:

1. Verify domain and cloud platform prerequisites.
2. Create the cluster provisioning configuration on the domain.
3. Create a Hadoop connection for the cluster workflow to use.

Then, a developer completes the following steps:

1. Create the workflow and the Create Cluster task.
2. Create Mapping tasks and associate them with mappings that you prepared. Optionally add Command tasks and other tasks to the workflow.

3. Optionally add a Delete Cluster task to the workflow.
4. Deploy and run the workflow.

Administrator Tasks

Before the developer can create a cluster workflow, an administrator performs tasks on the Informatica domain.

Before you create a cluster workflow, an administrator must complete the following tasks:

Step 1. Verify prerequisites.

On the cloud platform, install an Informatica domain. Ensure the domain has permission to access cloud storage resources and to create a cluster on the cloud platform. Enable DNS resolution from an on-premises domain.

Step 2. Create the cloud provisioning configuration.

Use the Administrator tool to create a cloud provisioning configuration. The cloud provisioning configuration contains all of the information that the Data Integration Service requires to contact and create resources on the cloud platform.

Step 3. Create a Hadoop connection.

Create a dedicated Hadoop connection to use with the cluster workflow. Populate the Cloud Provisioning Configuration property with the name of the cloud provisioning configuration you created for the workflow. The workflow uses settings in the Hadoop connection to run jobs on the cluster.

Create the Cluster Workflow

Create the cluster workflow and add and configure workflow elements.

The cluster workflow must have a Create Cluster task and at least one Mapping task. In addition, you can add Command and other workflow tasks. You can also add a Delete Cluster task, so that the cluster becomes an ephemeral cluster.

Workflow Task Run-Time Behavior

Set mapping and Mapping task properties to specify where the workflow runs Mapping tasks.

You can create a cluster workflow to run some mappings on the cluster that the workflow creates, and other mappings on another cluster.

To specify where each mapping runs, configure options in the mapping and the Mapping task.

Run the mapping on the cluster that the workflow creates.

The following table describes the run-time behavior based on how you configure the mapping and Mapping task:

Mapping Hadoop Connection Property	Mapping Task Cluster Identifier Property	Run Time Behavior
Auto Deploy	Auto Deploy	The Data Integration Service generates temporary Hadoop connections based on the values in the Hadoop connection associated with the workflow, and uses the temporary connections to run mappings on the cluster.
<Hadoop connection name>	Auto Deploy	The Mapping task Cluster Identifier property overrides the mapping Hadoop connection property. You might want to do this if you want to maintain the mapping Hadoop connection property value.

Run the mapping on another cluster.

The following table describes the run-time behavior based on how you configure the mapping and Mapping task:

Mapping Hadoop Connection Property	Mapping Task Cluster Identifier Property	Run Time Behavior
<Hadoop connection name>	Blank	The Mapping task Cluster Identifier property takes input from the Hadoop connection and runs the mapping on the cluster identified in the cloud configuration property of the Hadoop connection.

Configure the Cluster Workflow

A cluster workflow must have one Create Cluster task.

1. In the Developer tool, create a workflow.
2. From the palette of tasks, drag a Create Cluster task to the workflow editor.
3. Complete the Create Cluster task general properties.

Property	Description
Name	Task name.
Description	Optional description.

Property	Description
Connection Name	Name of the cloud provisioning configuration to use with the workflow.
Connection Type	Choose from the following options: <ul style="list-style-type: none"> - Amazon EMR. Create an Amazon EMR cluster. - HDInsight. Create an Azure HDInsight cluster.

- Configure task input and output properties.

Input properties

The Create Cluster task does not require any unique values for task input properties.

Output properties

Set the Cluster Identifier property to the default value, AutoDeployCluster.

Note: The Cluster Identifier property of the Create Cluster task overrides the Cluster Identifier property of the Mapping task.

- Set the advanced properties that correspond to your cloud platform.
 - To create a cluster on Amazon Web Services, see [“Amazon EMR Advanced Properties ” on page 204.](#)
 - To create a cluster on Microsoft Azure, see [“Azure HDInsight Advanced Properties for the Create Cluster Task” on page 207.](#)
- Configure the Software Settings property in the advanced properties if you want to perform the following optional tasks:
 - To run mappings on the Blaze engine, see [“Configure the Create Cluster Task to Run Mappings on the Blaze Engine” on page 208.](#)
 - To configure an external relational database as the Hive metastore database, see [“Configure the Cluster to Use an External RDS as the Hive Metastore Database” on page 209.](#)
- Connect the workflow Start_Event to the Create Cluster task.

Amazon EMR Advanced Properties

Set the advanced properties for an Amazon EMR cluster.

General Options

The following table describes general options that you can set for an EMR cluster:

Property	Description
Cluster Name	Name of the cluster to create.
Release Version	EMR version to run on the cluster. Enter the AWS version tag string to designate the version. For example: <code>emr-5.8.0</code> Default is Latest version supported.

Property	Description
Connection Name	Name of the Hadoop connection that you configured for use with the cluster workflow.
S3 Log URI	Optional. S3 location of logs for cluster creation. Format: s3://<bucket name>/<folder name> If you do not supply a location, no cluster logs will be stored.

Master Instance Group Options

The following table describes master instance group options that you can set for an EMR cluster:

Property	Description
Master Instance Type	Master node EC2 instance type. You can specify any available EC2 instance type. Default is m4.4xlarge.
Master Instance Maximum Spot Price	Maximum spot price for the master node. Setting this property changes the purchasing option of the master instance group to Spot instead of On-demand.

Core Instance Group Options

The following table describes core instance group options that you can set for an EMR cluster:

Property	Description
Core Instance Type	Core node EC2 instance type. You can specify any available EC2 instance type. Default is m4.4xlarge.
Core Instance Count	Number of core EC2 instances to create in the cluster. Default is 2.
Core Instance Maximum Spot Price	Maximum spot price for core nodes. Setting this property changes the purchasing option of the core instance group to Spot instead of On-demand.
Core Auto-Scaling Policy	Optional. Auto-scaling policy for core instances. Type the policy JSON statement here, or provide a path to a file that contains a JSON statement. Format: file:\\<path_to_policy_config_file>

Task Instance Group Options

The following table describes task instance group options that you can set for an EMR cluster:

Property	Description
Task Instance Type	Task node EC2 instance type. You can specify any available EC2 instance type. Default is m4.4xlarge.
Task Instance Count	Number of task EC2 instances to create in the cluster. Default is 2.
Task Instance Maximum Spot Price	Maximum spot price for task nodes. Setting this property changes the purchasing option of the task instance group to Spot instead of On-demand.
Task Auto-Scaling Policy	Optional. Auto-scaling policy for task instances. Type the policy JSON statement here, or provide a path to a file that contains a JSON statement. Format: <code>file:\\<path_to_policy_config_file></code>

Additional Options

The following table describes additional options that you can set for an EMR cluster:

Property	Description
Applications	Optional. Applications to add to the default applications that AWS installs. AWS installs certain applications when it creates an EMR cluster. In addition, you can specify additional applications. Select additional applications from the drop-down list. This field is equivalent to the Software Configuration list in the AWS EMR cluster creation wizard.
Tags	Optional. Tags to propagate to cluster EC2 instances. Tags assist in identifying EC2 instances. Format: <code>TagName1=TagValue1, TagName2=TagValue2</code>
Software Settings	Optional. Custom configurations to apply to the applications installed on the cluster. This field is equivalent to the Edit Software Settings field in the AWS cluster creation wizard. You can use this as a method to modify the software configuration on the cluster. Type the configuration JSON statement here, or provide a path to a file that contains a JSON statement. Format: <code>file:\\<path_to_custom_config_file></code>
Steps	Optional. Commands to run after cluster creation. For example, you can use this to run Linux commands or HDFS or Hive Hadoop commands. This field is equivalent to the Add Steps field in the AWS cluster creation wizard. Type the command statement here, or provide a path to a file that contains a JSON statement. Format: <code>file:\\<path_to_command_file></code>
Bootstrap Actions	Optional. Actions to perform after EC2 instances are running, and before applications are installed. Type the JSON statement here, or provide a path to a file that contains a JSON statement. Format: <code>file:\\<path_to_policy_config_file></code>

Azure HDInsight Advanced Properties for the Create Cluster Task

The following table describes the Advanced properties for a Microsoft Azure HDInsight cluster:

Property	Description
Cluster Name	Name of the cluster to create.
Azure Cluster Type	Type of the cluster to be created. Choose one of the options in the drop-down list. Default is Hadoop.
HDInsight version	HDInsight version to run on the cluster. Enter the HDInsight version tag string to designate the version. Default is the latest version supported.
Azure Cluster Location	Use the drop-down list to choose the location in which to create the cluster.
Head Node VM Size	Size of the head node instance to create. Default is Standard_D12_v2.
Number of Worker Node Instances	Number of worker node instances to create in the cluster. Default is 2.
Worker Node VM Size	Size of the worker node instance to create. Default is Standard_D13_v2.
Default Storage Type	Primary storage type to be used for the cluster. Choose one of the following options: <ul style="list-style-type: none">- Azure Data Lake Store- Azure BLOB storage account Default is BLOB storage
Default Storage Container or Root Mount Path	Default container for data. Type one of the following paths: <ul style="list-style-type: none">- For ADLS storage, type the path to the storage. For example, you can type <code>storage-name</code> or <code>storage-name/folder-name</code>.- For blob storage, type the path to the container. Format: <code>/path/</code>
Log Location	Optional. Path to the directory to store workflow event logs. Default is <code>/app-logs</code> .

Property	Description
Attach External Hive Metastore	If you select this option, the workflow attaches an external Hive metastore to the cluster if you configured an external Hive metastore in the cloud provisioning configuration.
Bootstrap JSON String	<p>JSON statement to run during cluster creation. You can use this statement to configure cluster details. For example, you could configure Hadoop properties on the cluster, add tags to cluster resources, or run script actions.</p> <p>Choose one of the following methods to populate the property:</p> <ul style="list-style-type: none"> - Type the JSON statement. Use the following format: <pre>{ "core-site" : { "<sample_property_key1>": "<sample_property_val1>", "<sample_property_key2>": "<sample_property_val2>" }, "tags": { "<tag_key>": "<tag_val>" }, "scriptActions": [{ "name": "setenvironmentvariable", "uri": "scriptActionUri", "parameters": "headnode" }] }</pre> - Provide a path to a file that contains a JSON statement. Format: <pre>file://<path_to_bootstrap_file></pre>

Configure the Create Cluster Task to Run Mappings on the Blaze Engine

If you want to use the Blaze engine to run mappings on the cloud platform cluster, you must set cluster configuration properties in the Software Setting property of the Create Cluster task.

Configure the Create Cluster task to set configuration properties in *-site.xml files on the cluster. Hadoop clusters run based on these settings.

The following text shows sample configuration of the Software Settings property:

```
[
  {
    "Classification": "yarn-site",
    "Properties": {
      "yarn.scheduler.minimum-allocation-mb": "250",
      "yarn.scheduler.maximum-allocation-mb": "8192",
      "yarn.nodemanager.resource.memory-mb": "16000",
      "yarn.nodemanager.resource.cpu-vcores": "12"
    }
  },
  {
    "Classification": "core-site",
    "Properties": {
      "hadoop.proxyuser.<DIS/OSPUSER>.groups": "<group names>",
      "hadoop.proxyuser.<DIS/OSPUSER>.hosts": "*"
    }
  }
]
```

yarn-site

yarn.scheduler.minimum-allocation-mb

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

yarn.scheduler.maximum-allocation-mb

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

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.

yarn.nodemanager.resource.cpu-vcores

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

core-site**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.

hadoop.proxyuser.<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 the property to " *" to allow impersonation from any host. This is required to run a Blaze mapping on a cloud platform cluster.

Configure the Cluster to Use an External RDS as the Hive Metastore Database

If you want to use a relational database on the cloud platform as the Hive metastore database for the cluster, you must set cluster configuration properties in the Software Setting property of the Create Cluster task.

Configure the Create Cluster task to set configuration properties in the hive-site.xml configuration file on the cluster. Use a text file to specify hive-site settings, and specify the path to the file in the Software Settings property.

The following text shows sample configuration of the Software Settings property:

```
{
  "Classification": "hive-site",
  "Properties": {
    "javax.jdo.option.ConnectionURL": "jdbc:mysql://<RDS_HOST>:<PORT>\\/  
<USER_SCHEMA>?createDatabaseIfNotExist=true",
    "javax.jdo.option.ConnectionDriverName": "<JDBC driver name>",
    "javax.jdo.option.ConnectionUserName": "<USER>",
    "javax.jdo.option.ConnectionPassword": "<USER>"
  }
}
```

Example:

```
Example:
{
  "Classification": "hive-site",
  "Properties": {
    "javax.jdo.option.ConnectionURL": "jdbc:mysql://<host name>:<port number>\\/  
hive?createDatabaseIfNotExist=true",
    "javax.jdo.option.ConnectionDriverName": "org.mariadb.jdbc.Driver",
    "javax.jdo.option.ConnectionUserName": "hive",
    "javax.jdo.option.ConnectionPassword": "hive"
  }
}
```

[hive-site](#)

javax.jdo.option.ConnectionURL

JDBC connection string for the data store.

javax.jdo.option.ConnectionDriverName

JDBC driver class name for the data store. Specify a JDBC driver that is compatible with the cloud platform.

javax.jdo.option.ConnectionUserName

User name to use to connect to the database.

javax.jdo.option.ConnectionPassword

Password for the database user account.

Create Other Workflow Tasks

Populate the cluster workflow with at least one Mapping task. You can add Command or other workflow tasks and events.

1. Drag a Mapping task from the task list to the workflow editor.
The **Mapping Task** dialog box opens.
2. Name the Mapping task.
3. Select a mapping to run with the Mapping task. Click **Browse** next to the Mapping property, select a mapping, and click **Finish**.
4. Optionally select a parameter set to associate with the Mapping task. Click **Browse** next to the Parameter Set property, select a parameter set, and click **Finish**.

For more information on how to use parameter sets with mappings, see the *Informatica Developer Mapping Guide*.

5. Optionally complete Input and Output properties.
The Mapping task does not require any unique values for input or output properties.
6. Configure the Cluster Identifier property in Advanced properties.
The Cluster Identifier property designates the cluster to use to run the Mapping task.

The following table describes values for Cluster Identifier properties:

Value	Description
Blank (no value)	Run the mapping runs on the cluster configured in the Hadoop connection associated with the mapping.
AutoDeploy	Run the mapping on the cluster that the workflow creates. When you choose this option, it also populates the Cluster Identifier property in the Create Cluster task with the value Set to AutoDeployCluster. Default is AutoDeploy.
(Assign to task input)	Select this option to accept input from another source than the Create Cluster task. If you choose this option, enter a parameter value in the Cluster Identifier property of the Mapping task Input properties tab.

7. Click **Finish** to the Mapping task.

8. Optionally add more Mapping and other tasks to the workflow.

You can include any other workflow tasks in a cluster workflow. For example, you might want to add a Command task to perform tasks after a Mapping task runs.

Add a Delete Cluster Task

To create an ephemeral cluster, add a Delete Cluster task.

The Delete Cluster task terminates the cluster and deletes the cluster and other resources that the cluster workflow creates.

If you do not add a Delete Cluster task, the cluster that the workflow creates remains running when the workflow ends. You can delete the cluster at any time.

1. Drag a Delete Cluster task to the workflow editor.
2. In the General properties, optionally rename the Delete Cluster task.
3. Connect the final task in the workflow to the Delete Cluster task, and connect the Delete Cluster task to the workflow End_Event.

You can also use `infacmd ccps deleteCluster` to delete a cloud cluster.

Deploy and Run the Workflow

After you complete the cluster workflow, deploy and run the workflow.

You can monitor cluster provisioning tasks on the AWS or Azure web console. If you configured a log location, view the logs at the location that you configured in the Create Cluster task properties.

You can also monitor Data Integration Service jobs in the Administrator tool.

Note: After the workflow begins executing tasks, the task to provision the cluster may take several minutes.

Monitoring Azure HDInsight Cluster Workflow Jobs

You can access mapping log URLs through the **Monitoring** tab in the Administrator tool to monitor workflow jobs that run on an Azure HDInsight cluster. The log location depends on the run-time engine that each mapping uses.

Blaze and Spark engines

To access the monitoring URL for mappings that run on Blaze or Spark, expand the workflow and the mapping in the **Monitoring** tab. Select the Grid Task and view the value for the Monitoring URL property in the lower pane. Use this path to find the log.

Hive engine

To access the monitoring URL for mappings that run on Hive, expand the workflow and the mapping in the **Monitoring** tab. Select a Hive Query job, and then expand the MR Job Details node in the lower pane. The Job ID is hyperlinked, but clicking on the link does not lead to the log. To find the job monitoring log, copy the URL path and read it to find the log. Repeat the steps for each Hive Query job.

APPENDIX A

Connections

This appendix includes the following topics:

- [Connections, 212](#)
- [Cloud Provisioning Configuration, 213](#)
- [Hadoop Connection Properties, 218](#)
- [HDFS Connection Properties, 223](#)
- [HBase Connection Properties, 225](#)
- [HBase Connection Properties for MapR-DB, 226](#)
- [Hive Connection Properties, 226](#)
- [JDBC Connection Properties, 230](#)
- [Creating a Connection to Access Sources or Targets, 235](#)
- [Creating a Hadoop Connection, 236](#)
- [Configuring Hadoop Connection Properties, 237](#)

Connections

Define a Hadoop connection to run a mapping in the Hadoop environment. Depending on the sources and targets, define connections to access data in HBase, HDFS, Hive, or relational databases. 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. If you select the mapping validation environment or the execution environment as Hadoop, select the Hadoop connection. Before you run mappings in the Hadoop environment, review the information in this guide about rules and guidelines for mappings that you can run 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 or Hive engine.

JDBC connection

Create a JDBC connection and configure Sqoop properties in the connection to import and export relational data through Sqoop.

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 Hadoop connection that the workflow uses 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 Hadoop 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.

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.

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. Note: 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"> - Permissions to access required directories in ADLS storage. - Certificate-based authentication for ADLS storage. - Key-based authentication for ADLS storage.
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.

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. Note: 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.

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.

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.
Cluster Environment Variables*	Environment variables that the Hadoop cluster uses. For example, the variable <code>ORACLE_HOME</code> represents the directory where the Oracle database client software is installed. 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 <code>infacmd ms runMapping</code> with the <code>-cp</code> 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
Cluster Library Path*	The path for shared libraries on the cluster. The <code>\$DEFAULT_CLUSTER_LIBRARY_PATH</code> variable contains a list of default directories.

Property	Description
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"> 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
Cluster Executable Path*	<p>The path for executable files on the cluster.</p> <p>The \$DEFAULT_CLUSTER_EXEC_PATH variable contains a list of paths to the default executable files.</p>
<p>* 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>	

Common Properties

The following table describes the common connection properties that you configure for the Hadoop connection:

Property	Description
Impersonation User Name	<p>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.</p> <p>The Data Integration Service runs mappings based on the user that is configured. Refer the following order to determine which user the Data Integration Services uses to run mappings:</p> <ol style="list-style-type: none"> 1. Operating system profile user. The mapping runs with the operating system profile user if the profile user is configured. If there is no operating system profile user, the mapping runs with the Hadoop impersonation user. 2. Hadoop impersonation user. The mapping runs with the Hadoop impersonation user if the operating system profile user is not configured. If the Hadoop impersonation user is not configured, the Data Integration Service runs mappings with the Data Integration Service user. 3. Informatica services user. The mapping runs with the operating user that starts the Informatica daemon if the operating system profile user and the Hadoop impersonation user are not configured.
Temporary Table Compression Codec	<p>Hadoop compression library for a compression codec class name.</p> <p>Note: 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.</p>
Codec Class Name	Codec class name that enables data compression and improves performance on temporary staging tables.

Property	Description
Hive Staging Database Name	<p>Namespace for Hive staging tables. Use the name <code>default</code> for tables that do not have a specified database name.</p> <p>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.</p>
Advanced Properties	<p>List of advanced properties that are unique to the Hadoop environment. The properties are common to the Blaze, Spark, and Hive 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"> 1. Mapping custom properties set using <code>infacmd ms runMapping</code> with the <code>-cp</code> 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 <p>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.</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 <code>RejectDir</code> system parameter.</p>
Reject File Directory	The directory for Hadoop mapping files on HDFS when you run mappings.

Hive Pushdown Configuration

The following table describes the connection properties that you configure for the Hive engine:

Property	Description
Environment SQL	<p>SQL commands to set the Hadoop environment. The Data Integration Service executes the environment SQL at the beginning of each Hive script generated in a Hive execution plan.</p> <p>The following rules and guidelines apply to the usage of environment SQL:</p> <ul style="list-style-type: none">- Use the environment SQL to specify Hive queries.- 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:<Informatica installation directory>/services/shared/hadoop/<Hadoop distribution name>/extras/hive-auxjars- You can use environment SQL to define Hadoop or Hive parameters that you want to use in the PreSQL commands or in custom queries.- If you use multiple values for the environment SQL, ensure that there is no space between the values.
Hive Warehouse Directory	<p>Optional. The absolute HDFS file path of the default database for the warehouse that is local to the cluster.</p> <p>If you do not configure the Hive warehouse directory, the Hive engine first tries to write to the directory specified in the cluster configuration property <code>hive.metastore.warehouse.dir</code>. If the cluster configuration does not have the property, the Hive engine writes to the default directory <code>/user/hive/warehouse</code>.</p>
Engine Type	<p>The engine that the Hadoop environment uses to run a mapping on the Hadoop cluster. You can choose MRv2 or Tez. You can select Tez if it is configured for Amazon EMR, Azure HDInsight, or Hortonworks HDP.</p>
Advanced Properties	<p>List of advanced properties that are unique to the Hive 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">1. Mapping custom properties set using <code>infacmd ms runMapping</code> with the <code>-cp</code> option2. Mapping run-time properties for the Hadoop environment3. Hadoop connection advanced properties for run-time engines4. Hadoop connection advanced general properties, environment variables, and classpaths5. Data Integration Service custom properties <p>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.</p>

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_<user name></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	The YARN scheduler queue name used by the Blaze engine that specifies available resources on a cluster.
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><hostname>:<port></code></p> <p>Where</p> <ul style="list-style-type: none">- <code><hostname></code> is the host name or IP address of the Blaze Job Monitor server.- <code><port></code> is the port on which the Blaze Job Monitor listens for remote procedure calls (RPC). <p>For example, enter: <code>myhostname:9080</code></p>
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>&&</code> (AND), <code> </code> (OR), and <code>!</code> (NOT).</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">1. Mapping custom properties set using <code>infacmd ms runMapping</code> with the <code>-cp</code> option2. Mapping run-time properties for the Hadoop environment3. Hadoop connection advanced properties for run-time engines4. Hadoop connection advanced general properties, environment variables, and classpaths5. Data Integration Service custom properties <p>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.</p>

Spark Configuration

The following table describes the connection properties that you configure for the Spark engine:

Property	Description
Spark Staging Directory	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. By default, the temporary files are written to the Hadoop staging directory <code>/tmp/spark_<user name></code> .
Spark Event Log Directory	Optional. The HDFS file path of the directory that the Spark engine uses to log events.
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">1. Mapping custom properties set using <code>infacmd ms runMapping</code> with the <code>-cp</code> option2. Mapping run-time properties for the Hadoop environment3. Hadoop connection advanced properties for run-time engines4. Hadoop connection advanced general properties, environment variables, and classpaths5. Data Integration Service custom properties <p>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.</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.
Type	The connection type. Default is Hadoop File System.
User Name	User name to access HDFS.
NameNode URI	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. Note: 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. If you manually set this property or override the value, the refresh operation does not update this property.

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	hdfs://<namenode>:<port> where: - <namenode> is the host name or IP address of the NameNode. - <port> is the port that the NameNode listens for remote procedure calls (RPC). hdfs://<nameservice> in case of NameNode high availability.
MapR-FS	maprfs:///

Storage	NameNode URI Format
WASB in HDInsight	<p>wasb://<container_name>@<account_name>.blob.core.windows.net/<path></p> <p>where:</p> <ul style="list-style-type: none"> - <container_name> identifies a specific Azure Storage Blob container. <p>Note: <container_name> is optional.</p> <ul style="list-style-type: none"> - <account_name> identifies the Azure Storage Blob object. <p>Example:</p> <p>wasb://infabdmoffering1storage.blob.core.windows.net/infabdmoffering1cluster/mr-history</p>
ADLS in HDInsight	adl://home

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	<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>~ ` ! \$ % ^ & * () - + = { [] } \ : ; " ' < , > . ? /</p>
ID	<p>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. 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. Select HBase.</p>
Database Type	<p>Type of database that you want to connect to.</p> <p>Select HBase to create a connection for an HBase table.</p>

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.
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 HBase .
Database Type	Type of database that you want to connect to. Select MapR-DB 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 Database Path 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	<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>~ ` ! \$ % ^ & * () - + = { [] } \ : ; " ' < , > . ? /</p>
ID	<p>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. The description cannot exceed 4000 characters.</p>
Location	<p>The domain where you want to create the connection. Not valid for the Analyst tool.</p>
Type	<p>The connection type. Select Hive.</p>
User Name	<p>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 runs Hortonworks HDP, you must provide a user name. If you use Tez to run mappings, you must provide the user account for the Data Integration Service. If you do not use Tez to run mappings, you can use an impersonation user account.</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"> - 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. - The Hadoop cluster uses Kerberos authentication. It authenticates jobs based on the SPN of the Data Integration Service. User Name will be ignored.
Password	<p>Password for the user name.</p>

Property	Description
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"> - Use the environment SQL to specify Hive queries. - 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: <code><Informatica installation directory>/services/shared/hadoop/ <Hadoop distribution name>/extras/hive-auxjars</code> - You can use environment SQL to define Hadoop or Hive parameters that you want to use in the PreSQL commands or in custom queries. - If you use multiple values for the Environment SQL property, ensure that there is no space between the values.
SQL Identifier Character	<p>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 Support mixed-case identifiers property.</p>

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. You can use PowerExchange for Hive to communicate with a HiveServer service or HiveServer2 service.</p> <p>To connect to HiveServer, specify the connection string in the following format: <code>jdbc:hive2://<hostname>:<port>/<db></code></p> <p>Where</p> <ul style="list-style-type: none">- <hostname> is name or IP address of the machine on which HiveServer2 runs.- <port> is the port number on which HiveServer2 listens.- <db> 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. <p>To connect to HiveServer 2, 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=<xyz></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://<hostname>:<port>/<db>;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>Informatica Big Data Management Hadoop 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">- Verify that Hive client and Informatica services are installed on the same machine.- Configure the Hive connection properties to run mappings on a Hadoop cluster. <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>

Property	Description
Observe Fine Grained SQL Authorization	When you select the option to observe fine-grained SQL authorization in a Hive source, the mapping observes row and column-level restrictions on data access. If you do not select the option, the Blaze and Spark engines ignore the restrictions, and results include restricted data. Applicable to Hadoop clusters where Sentry or Ranger security modes are enabled.
Data Access Connection String	<p>The connection string to access data from the Hadoop data store.</p> <p>To connect to HiveServer, specify the non-embedded JDBC mode connection string in the following format:</p> <pre>jdbc:hive2://<hostname>:<port>/<db></pre> <p>Where</p> <ul style="list-style-type: none"> - <hostname> is name or IP address of the machine on which HiveServer2 runs. - <port> is the port number on which HiveServer2 listens. - <db> 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. <p>To connect to HiveServer 2, 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=<xyz></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://<hostname>:<port>/<db>;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>Informatica Big Data Management Hadoop Integration Guide</i>.</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.

Property	Description
User Name	<p>The database user name.</p> <p>If you configure Sqoop, Sqoop uses the user name that you configure in this field. If you configure the --username argument in a JDBC connection or mapping, Sqoop ignores the argument.</p>
Password	<p>The password for the database user name.</p> <p>If you configure Sqoop, Sqoop uses the password that you configure in this field. If you configure the --password argument in a JDBC connection or mapping, Sqoop ignores the argument.</p>
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"> - DataDirect JDBC driver class name for Oracle: com.informatica.jdbc.oracle.OracleDriver - DataDirect JDBC driver class name for IBM DB2: com.informatica.jdbc.db2.DB2Driver - DataDirect JDBC driver class name for Microsoft SQL Server: com.informatica.jdbc.sqlserver.SQLServerDriver - DataDirect JDBC driver class name for Sybase ASE: com.informatica.jdbc.sybase.SybaseDriver - DataDirect JDBC driver class name for Informix: com.informatica.jdbc.informix.InformixDriver - DataDirect JDBC driver class name for MySQL: com.informatica.jdbc.mysql.MySQLDriver <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:<subprotocol>:<subname></pre> <p>The following list provides sample connection strings that you can enter for the applicable database type:</p> <ul style="list-style-type: none"> - Connection string for DataDirect Oracle JDBC driver: jdbc:informatica:oracle://<host>:<port>;SID=<value> - Connection string for Oracle JDBC driver: jdbc:oracle:thin:@//<host>:<port>:<SID> - Connection string for DataDirect IBM DB2 JDBC driver: jdbc:informatica:db2://<host>:<port>;DatabaseName=<value> - Connection string for IBM DB2 JDBC driver: jdbc:db2://<host>:<port>/<database_name> - Connection string for DataDirect Microsoft SQL Server JDBC driver: jdbc:informatica:sqlserver://<host>;DatabaseName=<value> - Connection string for Microsoft SQL Server JDBC driver: jdbc:sqlserver://<host>;DatabaseName=<value> - Connection string for Netezza JDBC driver: jdbc:netezza://<host>:<port>/<database_name> - Connection string for Pivotal Greenplum driver: jdbc:pivotal:greenplum://<host>:<port>;/database_name=<value> - Connection string for Postgres Greenplum driver: jdbc:postgresql://<host>:<port>/<database_name> - Connection string for Teradata JDBC driver: jdbc:teradata://<host>/database_name=<value>,tmode=<value>,charset=<value> <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>Note: If you enable Sqoop, Sqoop ignores this property.</p>

Property	Description
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>Note: If you enable Sqoop, Sqoop ignores this property.</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>Note: 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 SQL Identifier Character property.</p> <p>When the SQL Identifier Character property is set to none, the Support Mixed-case Identifiers property is disabled.</p> <p>Note: 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 Sqoop v1.x to enable Sqoop connectivity.</p> <p>Default is None.</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"> - To use Cloudera Connector Powered by Teradata, configure the following Sqoop argument: <ul style="list-style-type: none"> - <code>Dsqoop.connection.factories=com.cloudera.connector.teradata.TeradataManagerFactory</code> - To use Hortonworks Connector for Teradata (powered by the Teradata Connector for Hadoop), configure the following Sqoop argument: <ul style="list-style-type: none"> - <code>Dsqoop.connection.factories=org.apache.sqoop.teradata.TeradataManagerFactory</code> <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>Note: 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> <p>On the Hive engine, to run a column profile on a relational data object that uses Sqoop, set the Sqoop argument <code>m</code> to 1. Use the following syntax:</p> <pre>-m 1</pre>

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.

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>`

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 and Hive engines.

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.

To direct a MapReduce job to a specific YARN queue, configure the following argument:

```
-Dmapred.job.queue.name=<YARN queue name>
```

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.

Creating a Connection to Access Sources or Targets

Create an HBase, HDFS, Hive, or JDBC connection before you import data objects, preview data, and profile data.

1. Click **Window > Preferences**.
2. Select **Informatica > Connections**.
3. Expand the domain in the **Available Connections** list.
4. 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**.
5. Click **Add**.
6. Enter a connection name and optional description.
7. Click **Next**.
8. 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. The **Access Hive to run mappings in Hadoop cluster** options is no more applicable. To use the Hive driver to run mappings in the Hadoop cluster, use a Hadoop connection.
9. Click **Test Connection** to verify the connection.
10. 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.

New Cluster Connection

Cluster Connection

Provide the connection details.

Name:

ID:

Description:

Location:

Type:

6. Click **Next**.
7. Enter the Hadoop cluster properties, common properties, and the reject directory properties.
8. Click **Next**.
9. Enter configuration properties for the Hive engine and click **Next**.
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
- Cluster ClassPath
- Cluster Executable Path
- Common Advanced Properties
- Hive Engine 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. You must use JDK version 1.7 or later. Default is /usr/java/default. The JDK version that the Data Integration Service uses must be compatible with the JRE version on the cluster.

Set to <cluster JDK home>/jdk<version>.

For example, HADOOP_NODE_JDK_HOME=<cluster JDK home>/jdk<version>.

DB2_HOME

Specifies the DB2 home directory. Required to run mappings with DB2 sources and targets on the Hive engine.

Set to /databases/db2<version>.

For example, DB2_HOME=/databases/db2V10.5_64BIT.

DB2INSTANCE

Specifies the DB2 database instance name. Required to run mappings with DB2 sources and targets on the Hive engine.

Set to <DB2 instance name>.

For example, `DB2INSTANCE=db10inst`.

DB2CODEPAGE

Specifies the code page configured in the DB2 instance. Required to run mappings with DB2 sources and targets on the Hive engine.

Set to <DB2 instance code page>.

For example, `DB2CODEPAGE="1208"`.

GPHOME_LOADERS

Represents the directory to the Greenplum libraries. Required to run Greenplum mappings on the Hive engine.

Set to <Greenplum libraries directory>.

For example, `GPHOME_LOADERS=opt/thirdparty/`.

PYTHONPATH

Represents the directory to the Python path libraries. Required to run Greenplum mappings on the Hive engine.

Set to <Python path libraries directory>.

For example, `PYTHONPATH=$GPHOME_LOADERS/bin/ext`.

NZ_HOME

Represents the directory that contains the Netezza client libraries. Required to run Netezza mappings on the Hive or Blaze engine.

Set to <Netezza client library directory>.

For example, `NZ_HOME=/opt/thirdparty/netezza`.

NZ_ODBC_INI_PATH

Represents the directory that contains the odbc.ini file. Required to run Netezza mappings on the Hive or Blaze engine.

Set to <odbc.ini file path>.

For example, `NZ_ODBC_INI_PATH=/opt/ODBCINI`.

ODBCINI

Represents the path and file name of the odbc.ini file.

- Required to run Netezza mappings on the Hive or Blaze engine.
Set to <odbc.ini file path>/<file name>.

For example, `ODBCINI=/opt/ODBCINI/odbc.ini`.

- Required to run mappings with ODBC sources and targets on the Hive engine.
Set to <odbc.ini file path>/<file name>.

For example, `ODBCINI=$HADOOP_NODE_INFA_HOME/ODBC7.1/odbc.ini`.

ODBC_HOME

Specifies the ODBC home directory. Required to run mappings with ODBC sources and targets on the Hive engine.

Set to <odbc home directory>.

For example, `ODBC_HOME=$HADOOP_NODE_INFA_HOME/ODBC7.1.`

ORACLE_HOME

Specifies the Oracle home directory. Required to run mappings with Oracle sources and targets on the Hive engine.

Set to <Oracle home directory>.

For example, `ORACLE_HOME=/databases/oracle12.1.0_64BIT.`

TNS_ADMIN

Specifies the directory to the Oracle client `tnsnames.ora` configuration files. Required to run mappings with Oracle sources and targets on the Hive engine.

Set to <tnsnames.ora config files directory>.

For example, `TNS_ADMIN=/opt/ora_tns.`

HADOOP_CLASSPATH

Represents the directory to the TDCH libraries. Required to run Teradata mappings through TDCH on the Hive engine.

Set to <TDCH libraries directory>.

For example,

```
/opt/cloudera/parcels/CDH-5.13.0-1.cdh5.13.0.p0.29/lib/hive/conf
/opt/cloudera/parcels/CDH-5.13.0-1.cdh5.13.0.p0.29/lib/hive/lib/*
/usr/lib/tdch/1.5/lib/*
```

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 following library path variables in the **Cluster Library Path** property:

\$DB2_HOME/lib64

Represents the directory to the DB2 libraries. Required to run mappings with DB2 sources and targets on the Hive engine.

\$GPHOME_LOADERS/lib

Represents the path to the Greenplum libraries. Required to run Greenplum mappings on the Hive engine.

\$GPHOME_LOADERS/ext/python/lib

Represents the path to the Python libraries. Required to run Greenplum mappings on the Hive engine.

\$NZ_HOME/lib64

Represents the path to the Netezza libraries. Required to run Netezza mappings on the Hive or Blaze engine.

\$ORACLE_HOME/lib

Represents the directory to the Oracle libraries. Required to run mappings with Oracle sources and targets on the Hive engine.

/usr/lib/tdch/1.5/lib/*

The path to the TDCH libraries directory. Required to run Teradata mappings through TDCH on the Hive engine.

Cluster ClassPath

Cluster ClassPath property is a list of classpath variables to access the Hadoop jar files and the required libraries on the cluster. You can add or edit classpath variables.

To edit the property in the text box, use the following format with : to separate each path variable:

```
<variable1>[:<variable2>...:<variableN>]
```

Configure the following classpath variable in the **Cluster ClassPath** property:

/usr/lib/tdch/1.5/lib/*

Path to the TDCH libraries directory. Required to run Teradata mappings through TDCH on the Hive engine.

Cluster Executable Path

Cluster Executable Path property is a list of path variables to access executable files on the cluster. You can add or edit executable 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 following library path variables in the **Cluster Executable Path** property:

\$DB2_HOME/bin

Represents the directory to the DB2 binaries. Required to run mappings with DB2 sources and targets on the Hive engine.

\$GPHOME_LOADERS/bin

Represents the path to the Greenplum binaries. Required to run Greenplum mappings on the Hive engine.

\$GPHOME_LOADERS/ext/python/bin

Represents the path to the Python binaries. Required to run Greenplum mappings on the Hive engine.

\$ORACLE_HOME/bin

Represents the path to the Oracle binaries. Required to run mappings with Oracle sources and targets on the Hive engine.

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, Spark, and Hive 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`

Hive Engine Advanced Properties

Hive advanced properties are a list of advanced or custom properties that are unique to the Hive 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>]
```

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/infa/logs/blaze`. Required to set up multiple Blaze instances.

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>`

Spark Engine 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.

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 Spark configuration section:

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.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.authenticate.sasl.encryption.aes.enabled

Enables AES support when SASL authentication is enabled. Required if Spark encryption uses SASL authentication.

Set to TRUE.

For example, `spark.authenticate.sasl.encryption.aes.enabled=TRUE`

infaspark.pythontx.executorEnv.LD_PRELOAD

The location of the Python shared library in the Python installation folder on the Data Integration Service machine. Required to run a Python transformation on the Spark engine.

For example, set to:

```
infaspark.pythontx.executorEnv.LD_PRELOAD=  
<Informatica installation directory>/services/shared/spark/python/lib/  
libpython3.6m.so
```

infaspark.pythontx.submit.lib.JEP_HOME

The location of the Jep package in the Python installation folder on the Data Integration Service machine. Required to run a Python transformation on the Spark engine.

For example, set to:

```
infaspark.pythontx.submit.lib.JEP_HOME=  
<Informatica installation directory>/services/shared/spark/python/lib/python3.6/site-  
packages/jep/
```

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.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
```

APPENDIX B

Data Type Reference

This appendix includes the following topics:

- [Data Type Reference Overview, 245](#)
- [Transformation Data Type Support in a Hadoop Environment, 246](#)
- [Complex File and Transformation Data Types, 246](#)
- [Hive Data Types and Transformation Data Types, 250](#)
- [Sqoop Data Types, 252](#)

Data Type Reference Overview

Informatica Developer uses the following data types in mappings that run on the Hadoop cluster:

Native data types

Native data types are specific to the Hadoop sources and targets used as a physical data object. Native data types appear in the physical data object column properties.

Transformation data types

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

Transformation data types include the following data types:

- Primitive data type. Represents a single data value in a single column position.
- Complex data type. Represents multiple data values in a single column position. Use complex data types in mappings that run on the Spark engine to process hierarchical data in complex files.

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

Transformation Data Type Support in a Hadoop Environment

The following table shows the Informatica transformation data type support in a Hadoop environment:

Transformation Data Type	Support
Array	Supported*
Bigint	Supported
Binary	Supported
Date/Time	Supported
Decimal	Supported
Double	Supported
Integer	Supported
Map	Supported*
String	Supported
Struct	Supported*
Text	Supported
timestampWithTZ	Not supported
* Supported only on the Spark engine.	

Complex File and Transformation Data Types

You can use complex data types in mappings to process hierarchical data in complex files.

You can use complex data types in the following complex files in mappings that run on the Spark engine:

- Avro
- JavaScript Object Notation (JSON)
- Parquet

Avro and Transformation Data Types

Apache Avro data types map to transformation data types that the Data Integration Service uses to move data across platforms.

The following table compares Avro data types and transformation data types:

Avro	Transformation	Description
Array	Array	Unlimited number of characters.
Boolean	Integer	-2,147,483,648 to 2,147,483,647 Precision of 10, scale of 0
Bytes	Binary	1 to 104,857,600 bytes
Double	Double	Precision of 15 digits
Fixed	Binary	1 to 104,857,600 bytes
Float	Double	Precision of 15 digits
Int	Integer	-2,147,483,648 to 2,147,483,647 Precision of 10, scale of 0
Long	Bigint	-9,223,372,036,854,775,808 to 9,223,372,036,854,775,807 Precision of 19, scale of 0
Map	Map	Unlimited number of characters.
Record	Struct	Unlimited number of characters.
String	String	1 to 104,857,600 characters
Union	Corresponding data type in a union of ["primitive_type complex_type", "null"] or ["null", "primitive_type complex_type"].	Dependent on primitive or complex data type.

Avro Union Data Type

A union indicates that a field might have more than one data type. For example, a union might indicate that a field can be a string or a null. A union is represented as a JSON array containing the data types.

The Developer tool only interprets a union of ["primitive_type|complex_type", "null"] or ["null", "primitive_type|complex_type"]. The Avro data type converts to the corresponding transformation data type. The Developer tool ignores the null.

Unsupported Avro Data Types

The Developer tool does not support the following Avro data types:

- enum
- null

JSON and Transformation Data Types

JavaScript Object Notation data types map to transformation data types that the Data Integration Service uses to move data across platforms.

The following table compares JSON data types and transformation data types:

JSON	Transformation	Description
Array	Array	Unlimited number of characters.
Double	Double	Precision of 15 digits
Integer	Integer	-2,147,483,648 to 2,147,483,647 Precision of 10, scale of 0
Object	Struct	Unlimited number of characters.
String	String	1 to 104,857,600 characters

Unsupported JSON Data Types

The Developer tool does not support the following JSON data types:

- date/timestamp
- enum
- union

Parquet and Transformation Data Types

Apache Parquet data types map to transformation data types that the Data Integration Service uses to move data across platforms.

The following table compares Parquet data types and transformation data types:

Parquet	Transformation	Description
Binary	Binary	1 to 104,857,600 bytes
Binary (UTF8)	String	1 to 104,857,600 characters
Boolean	Integer	-2,147,483,648 to 2,147,483,647 Precision of 10, scale of 0
Double	Double	Precision of 15 digits

Parquet	Transformation	Description
Fixed Length Byte Array	Decimal	<p>Decimal value with declared precision and scale. Scale must be less than or equal to precision.</p> <p>For transformations that support precision up to 38 digits, the precision is 1 to 38 digits, and the scale is 0 to 38.</p> <p>For transformations that support precision up to 28 digits, the precision is 1 to 28 digits, and the scale is 0 to 28.</p> <p>If you specify the precision greater than the maximum number of digits, the Data Integration Service converts decimal values to double in high precision mode.</p>
Float	Double	Precision of 15 digits
group (LIST)	Array	Unlimited number of characters.
Int32	Integer	-2,147,483,648 to 2,147,483,647 Precision of 10, scale of 0
Int64	Bigint	-9,223,372,036,854,775,808 to 9,223,372,036,854,775,807 Precision of 19, scale of 0
Int64 (TIMESTAMP_MILLIS)	Date/Time	Jan 1, 0001 A.D. to Dec 31, 9999 A.D. Precision of 29, scale of 9 (precision to the nanosecond) Combined date/time value.
Int96	Date/Time	Jan 1, 0001 A.D. to Dec 31, 9999 A.D. Precision of 29, scale of 9 (precision to the nanosecond) Combined date/time value.
Map	Map	Unlimited number of characters.
Struct	Struct	Unlimited number of characters.
Union	Corresponding primitive data type in a union of ["primitive_type", "null"] or ["null", "primitive_type"].	Dependent on primitive data type.

Parquet Union Data Type

A union indicates that a field might have more than one data type. For example, a union might indicate that a field can be a string or a null. A union is represented as a JSON array containing the data types.

The Developer tool only interprets a union of ["primitive_type", "null"] or ["null", "primitive_type"]. The Parquet data type converts to the corresponding transformation data type. The Developer tool ignores the null.

Unsupported Parquet Data Types

The Developer tool does not support the following Parquet data types:

- int96 (TIMESTAMP_MILLIS)

Hive Data Types and Transformation Data Types

The following table lists the Hive data types that Data Integration Service supports and the corresponding transformation data types:

Hive Data Type	Transformation Data Type	Range and Description
Binary	Binary	1 to 104,857,600 bytes. You can read and write data of Binary data type in a Hadoop environment. You can use the user-defined functions to transform the binary data type.
Tiny Int	Integer	-32,768 to 32,767
Integer	Integer	-2,147,483,648 to 2,147,483,647 Precision 10, scale 0
Bigint	Bigint	-9,223,372,036,854,775,808 to 9,223,372,036,854,775,807 Precision 19, scale 0
Decimal	Decimal	<p>Precision 1 to 28, scale 0 to 28</p> <p>For transformations that support precision up to 38 digits, the precision is 1 to 38 digits, and the scale is 0 to 38.</p> <p>For transformations that support precision up to 28 digits, the precision is 1 to 28 digits, and the scale is 0 to 28.</p> <p>For transformations that support precision up to 38 digits, the precision is 1 to 38 digits, and the scale is 0 to 38.</p> <p>For transformations that support precision up to 28 digits, the precision is 1 to 28 digits, and the scale is 0 to 28.</p> <p>If a mapping is not enabled for high precision, the Data Integration Service converts all decimal values to double values.</p> <p>If a mapping is enabled for high precision, the Data Integration Service converts decimal values with precision greater than 38 digits to double values.</p>
Double	Double	Precision 15
Float	Double	Precision 15
String	String	1 to 104,857,600 characters

Hive Data Type	Transformation Data Type	Range and Description
Boolean	Integer	1 or 0 The default transformation type for boolean is integer. You can also set this to string data type with values of True and False.
Array	String	1 to 104,857,600 characters
Struct	String	1 to 104,857,600 characters
Map	String	1 to 104,857,600 characters
Timestamp	datetime	The time stamp format is YYYY-MM-DD HH:MM:SS.ffffffff. Precision 29, scale 9.
Date	datetime	0000-0101 to 999912-31. Hive date format is YYYY-MM-DD. Precision 10, scale 0.
Char	String	1 to 255 characters
Varchar	String	1 to 65355 characters

Hive Complex Data Types

Hive complex data types such as arrays, maps, and structs are a composite of primitive or complex data types. Informatica Developer represents complex data types with the string data type and uses delimiters to separate the elements of the complex data type.

Note: Hive complex data types in a Hive source or Hive target are not supported when you run mappings on a Hadoop cluster.

The following table describes the transformation types and delimiters that are used to represent the complex data types:

Complex Data Type	Description
Array	The elements in the array are of string data type. The elements in the array are delimited by commas. For example, an array of <code>fruits</code> is represented as <code>[apple,banana,orange]</code> .
Map	Maps contain key-value pairs and are represented as pairs of strings and integers delimited by the <code>=</code> character. String and integer pairs are delimited by commas. For example, a map of <code>fruits</code> is represented as <code>[1=apple,2=banana,3=orange]</code> .
Struct	Structs are represented as pairs of strings and integers delimited by the <code>:</code> character. String and integer pairs are delimited by commas. For example, a struct of <code>fruits</code> is represented as <code>[1,apple]</code> .

Sqoop Data Types

When you use Sqoop, some variations apply in the processing. Sqoop supports a subset of data types that database vendors support.

Aurora Data Types

Informatica supports the following Aurora data types when you use Sqoop:

- Binary
- Bit
- Blob (supported only for import)
- Char
- Date
- Datetime
- Decimal
- Double
- Enum
- Float
- Integer
- Numeric
- Real
- Set
- Text
- Time
- Timestamp
- Varbinary
- Varchar

IBM DB2 and DB2 for z/OS Data Types

Informatica supports the following IBM DB2 and DB2 for z/OS data types when you use Sqoop:

- Bigint
- Blob (supported only for import)
- Char
- Clob
- Date
- DBClob
- DecFloat (supported only for import)
- Decimal
- Double (supported only for DB2 for z/OS)
- Float (supported only for DB2)

- Graphic
- Integer
- LongVargraphic (supported only for DB2)
- Numeric
- Real
- Smallint
- Time
- Timestamp
- Varchar
- Vargraphic
- XML (supported only for import)

Greenplum Data Types

Informatica supports the following Greenplum data types when you use Sqoop:

- Bigint
- Bigserial
- Bytea
- Date
- Decimal
- Double
- Integer
- Nchar
- Numeric
- Nvarchar
- Real
- Serial
- Smallint
- Text
- Time
- Timestamp

Microsoft SQL Server Data Types

Informatica supports the following Microsoft SQL Server data types when you use Sqoop:

- Bigint
- Bit
- Char
- Datetime
- Decimal
- Float

- INT
- Money
- Numeric
- Real
- Smalldatetime
- Smallint
- Smallmoney
- Text
- Time
- Tinyint
- Varchar

Rules and Guidelines for Sqoop Microsoft SQL Server Data Types

Consider the following rules and guidelines when you configure Microsoft SQL Server data types in a Sqoop mapping:

- If you create or replace the target table at run time and run the mapping on the Blaze or Spark engine to export Bigint data, the mapping fails.
- If you run a Sqoop mapping to export time data, Sqoop does not export milliseconds.

Netezza Data Types

Informatica supports the following Netezza data types when you use Sqoop:

- Bigint
- Blob (supported only for import)
- Byteint
- Char
- Date
- Double
- Float4
- Float8
- Number
- Timestamp
- Varchar

Oracle Data Types

Informatica supports the following Oracle data types when you use Sqoop:

- Blob (supported only for import)
- Char
- Date
- Float
- Long

- Nchar (supported if you configure OraOop)
- Nvarchar (supported if you configure OraOop)
- Number(P,S)
- Timestamp
- Varchar
- Varchar2

Rules and Guidelines for Sqoop Oracle Data Types

Consider the following rules and guidelines when you configure Oracle data types in a Sqoop mapping:

- If you run a Sqoop mapping on the Blaze engine to export Oracle float data, Sqoop truncates the data.
- If you run a Sqoop mapping on the Blaze engine to export Oracle timestamp data with nanoseconds, Sqoop writes only three digits to the target.
- If you configure OraOop and run a Sqoop mapping on the Spark engine to export Oracle timestamp data, Sqoop writes only three digits to the target.

Teradata Data Types

Informatica supports the following Teradata data types when you use Sqoop:

- Bigint (supported only for import)
- Blob (supported only for import)
- Byteint
- Char
- Clob
- Date
- Decimal
- Double
- Float
- Integer
- Number
- Numeric
- Real
- Smallint
- Time
- **Note:** If you run a Sqoop mapping to export time data, Sqoop does not export milliseconds.
- Timestamp
- Varchar

Teradata Data Types with TDCH Specialized Connectors for Sqoop

Informatica supports the following Teradata data types when you use the Cloudera Connector Powered by Teradata, Hortonworks Connector for Teradata, and MapR Connector for Teradata with Sqoop:

- Bigint

- Byte (supported only by Hortonworks Connector for Teradata and MapR Connector for Teradata)
- Byteint
- Character
- Date
- Decimal
- Double Precision/Float/Real
- Integer
- Number(P,S)
- Numeric
- Smallint
- Time (supported only by Cloudera Connector Powered by Teradata and Hortonworks Connector for Teradata)
- Timestamp
- Varchar
- Varbyte (supported only by Hortonworks Connector for Teradata and MapR Connector for Teradata)

APPENDIX C

Function Reference

This appendix includes the following topics:

- [Function Support in the Hadoop Environment, 257](#)
- [Function and Data Type Processing, 259](#)

Function Support in the Hadoop Environment

Some Informatica transformation language functions that are valid in the native environment might be restricted or unsupported in the Hadoop environment. Functions not listed in this table are supported on all engines without restrictions.

Important: When you push a mapping to the Hadoop environment, the engine that processes the mapping uses a set of rules different from the Data Integration Service. As a result, the mapping results can vary based on the rules that the engine uses.

The following table lists functions and levels of support for functions on different engines in the Hadoop environment.

Function	Blaze Engine	Spark Engine	Hive Engine
ABORT	Not supported	Not supported	Not supported
ANY	Supported	Supported	Supported
ARRAY	Not supported	Supported	Not supported
AVG	Supported	Supported	Supported
CAST	Not supported	Supported	Not supported
COLLECT_LIST	Not supported	Supported	Not supported
COLLECT_MAP	Not supported	Supported	Not supported
CONCAT_ARRAY	Not supported	Supported	Not supported
COUNT	Supported	Supported	Supported
CREATE_TIMESTAMP_TZ	Supported	Not supported	Not supported

Function	Blaze Engine	Spark Engine	Hive Engine
CUME	Not supported	Not supported	Not supported
DECOMPRESS	Supported	Supported	Supported
ERROR	Not supported	Not supported	Not supported
FIRST	Not supported	Not supported	Not supported
GET_TIMEZONE	Supported	Not supported	Not supported
GET_TIMESTAMP	Supported	Not supported	Not supported
LAST	Not supported	Not supported	Not supported
LAG	Not supported	Supported	Not supported
LEAD	Not supported	Supported	Not supported
MAP	Not supported	Supported	Not supported
MAP_FROM_ARRAYS	Not supported	Supported	Not supported
MAP_KEYS	Not supported	Supported	Not supported
MAP_VALUES	Not supported	Supported	Not supported
MAX (Dates)	Supported	Supported	Not supported
MAX (Numbers)	Supported	Supported	Supported
MAX (String)	Supported	Supported	Supported
METAPHONE	Supported	Supported	Supported
MIN (Dates)	Supported	Supported	Not supported
MIN (Numbers)	Supported	Supported	Supported
MIN (String)	Supported	Supported	Supported
MOVINGAVG	Not supported	Not supported	Not supported
MOVINGSUM	Not supported	Not supported	Not supported
PERCENTILE	Supported	Supported	Supported
RESPEC	Not supported	Supported	Not supported
SIZE	Not supported	Supported	Not supported
STDDEV	Supported	Supported	Supported
STRUCT	Not supported	Supported	Not supported

Function	Blaze Engine	Spark Engine	Hive Engine
STRUCT_AS	Not supported	Supported	Not supported
SUM	Supported	Supported	Supported
SYSTIMESTAMP	Supported	Supported	Supported
TO_DATE	Supported	Supported	Supported
TO_DECIMAL	Supported	Supported	Supported
TO_DECIMAL38	Supported	Supported	Supported
TO_TIMESTAMP_TZ	Supported	Not supported	Supported
UUID4	Supported	Supported	Supported
UUID_UNPARSE	Supported	Supported	Supported
VARIANCE	Supported	Supported	Supported

For more information, see the *Informatica Transformation Language Reference*.

Function and Data Type Processing

When you run a mapping in a Hadoop environment, the Hadoop engine might process Informatica functions and data types differently from the Data Integration Service. Also, each of the run-time engines in the Hadoop environment can process Informatica functions and data types differently. Therefore, some variations apply in the processing and validity of functions and data types, and mapping results can vary.

Rules and Guidelines for Spark Engine Processing

Some restrictions and guidelines apply to processing Informatica functions on the Spark engine.

Important: When you push a mapping to the Hadoop environment, the engine that processes the mapping uses a set of rules different from the Data Integration Service. As a result, the mapping results can vary based on the rules that the engine uses. This topic contains some processing differences that Informatica discovered through internal testing and usage. Informatica does not test all the rules of the third-party engines and cannot provide an extensive list of the differences.

Consider the following rules and guidelines for function and data type processing on the Spark engine:

- The Spark engine and the Data Integration Service process overflow values differently. The Spark engine processing rules might differ from the rules that the Data Integration Service uses. As a result, mapping results can vary between the native and Hadoop environment when the Spark engine processes an overflow. Consider the following processing variation for Spark:
 - If an expression results in numerical errors, such as division by zero or SQRT of a negative number, it returns an infinite or an NaN value. In the native environment, the expression returns null values and the rows do not appear in the output.

- The Spark engine and the Data Integration Service process data type conversions differently. As a result, mapping results can vary between the native and Hadoop environment when the Spark engine performs a data type conversion. Consider the following processing variations for Spark:
 - The results of arithmetic operations on floating point types, such as Decimal, can vary up to 0.1 percent between the native environment and a Hadoop environment.
 - The Spark engine ignores the scale argument of the TO_DECIMAL function. The function returns a value with the same scale as the input value.
 - When the scale of a double or decimal value is smaller than the configured scale, the Spark engine trims the trailing zeros.
 - The Spark engine cannot process dates to the nanosecond. It can return a precision for date/time data up to the microsecond.
- The Spark engine does not support high precision. If you enable high precision, the Spark engine processes data in low-precision mode.
- If you use Hive 2.3, the Spark engine guarantees scale values.
 For example, when the Spark engine processes the decimal 1.1234567 with scale 9 using Hive 2.3, the output is 1.123456700. If you do not use Hive 2.3, the output is 1.1234567.
- The Hadoop environment treats "/n" values as null values. If an aggregate function contains empty or NULL values, the Hadoop environment includes these values while performing an aggregate calculation.
- Mapping validation fails if you configure SYSTIMESTAMP with a variable value, such as a port name. The function can either include no argument or the precision to which you want to retrieve the timestamp value.
- Mapping validation fails if an output port contains a Timestamp with Time Zone data type.
- Avoid including single and nested functions in an Aggregator transformation. The Data Integration Service fails the mapping in the native environment. It can push the processing to the Hadoop environment, but you might get unexpected results. Informatica recommends creating multiple transformations to perform the aggregation.
- You cannot preview data for a transformation that is configured for windowing.
- The Spark METAPHONE function uses phonetic encoders from the `org.apache.commons.codec.language` library. When the Spark engine runs a mapping, the METAPHONE function can produce an output that is different from the output in the native environment. The following table shows some examples:

String	Data Integration Service	Spark Engine
Might	MFT	MT
High	HF	H

- If you use the TO_DATE function on the Spark engine to process a string written in ISO standard format, you must add *T* to the date string and **T** to the format string. The following expression shows an example that uses the TO_DATE function to convert a string written in the ISO standard format YYYY-MM-DDTHH24:MI:SS:

```
TO_DATE('2017-11-03*T*12:45:00','YYYY-MM-DD**T**HH24:MI:SS')
```

The following table shows how the function converts the string:

ISO Standard Format	RETURN VALUE
2017-11-03T12:45:00	Nov 03 2017 12:45:00

- The UUID4 function is supported only when used as an argument in UUID_UNPARSE or ENC_BASE64.
- The UUID_UNPARSE function is supported only when the argument is UUID4().

Rules and Guidelines for Hive Engine Processing

Some restrictions and guidelines apply to processing Informatica functions on the Hive engine.

Important: When you push a mapping to the Hadoop environment, the engine that processes the mapping uses a set of rules different from the Data Integration Service. As a result, the mapping results can vary based on the rules that the engine uses. This topic contains some processing differences that Informatica discovered through internal testing and usage. Informatica does not test all the rules of the third-party engines and cannot provide an extensive list of the differences.

Consider the following rules and guidelines for function and data type processing on the Hive engine:

- The Hive engine and the Data Integration Service process overflow values differently. The Hive engine processing rules might differ from the rules that the Data Integration Service uses. As a result, mapping results can vary between the native and Hadoop environment when the Hive engine processes an overflow. Consider the following processing variations for Hive:
 - Hive uses a maximum or minimum value for integer and bigint data when there is data overflow during data type conversion.
 - If an expression results in numerical errors, such as division by zero or SQRT of a negative number, it returns an infinite or an NaN value. In the native environment, the expression returns null values and the rows do not appear in the output.
- The Hive engine and the Data Integration Service process data type conversions differently. As a result, mapping results can vary between the native and Hadoop environment when the Hive engine performs a data type conversion. Consider the following processing variations for Hive:
 - The results of arithmetic operations on floating point types, such as Decimal, can vary up to 0.1 percent between the native environment and a Hadoop environment.
 - You can use high precision Decimal data type with Hive 0.11 and above. When you run mappings on the Hive engine, the Data Integration Service converts decimal values with a precision greater than 38 digits to double values. When you run mappings that do not have high precision enabled, the Data Integration Service converts decimal values to double values.
 - When the Data Integration Service converts a decimal with a precision of 10 and a scale of 3 to a string data type and writes to a flat file target, the results can differ between the native environment and a Hadoop environment. For example, on the Hive engine, HDFS writes the output string for the decimal 19711025 with a precision of 10 and a scale of 3 as 1971. The Data Integration Service sends the output string for the decimal 19711025 with a precision of 10 and a scale of 3 as 1971.000.
 - The results of arithmetic operations on floating point types, such as a Double, can vary up to 0.1 percent between the Data Integration Service and the Hive engine.
 - When you run a mapping with a Hive target that uses the Double data type, the Data Integration Service processes the double data up to 17 digits after the decimal point.
- The Hadoop environment treats "/n" values as null values. If an aggregate function contains empty or NULL values, the Hadoop environment includes these values while performing an aggregate calculation.

- Avoid including single and nested functions in an Aggregator transformation. The Data Integration Service fails the mapping in the native environment. It can push the processing to the Hadoop environment, but you might get unexpected results. Informatica recommends creating multiple transformations to perform the aggregation.
- The UUID4 function is supported only when used as an argument in UUID_UNPARSE or ENC_BASE64.
- The UUID_UNPARSE function is supported only when the argument is UUID4().

APPENDIX D

Parameter Reference

This appendix includes the following topics:

- [Parameters Overview, 263](#)
- [Parameter Usage, 264](#)

Parameters Overview

A mapping parameter represents a constant value that you can change between mapping runs. Use parameters to change the values of connections, file directories, expression components, port lists, port links, and task properties. You can use system parameters or user-defined parameters.

System parameters are built-in parameters for a Data Integration Service. System parameters define the directories where the Data Integration Service stores log files, cache files, reject files, source files, target files, and temporary files. An administrator defines the system parameter default values for a Data Integration Service in the Administrator tool.

User-defined parameters are parameters that you define in transformations, mappings, or workflows. Create user-defined parameters to rerun a mapping with different connection, flat file, cache file, temporary file, expression, ports, or reference table values.

You can override parameter values using a parameter set or a parameter file. A parameter set is a repository object that contains mapping parameter values. A parameter file is an XML file that contains parameter values. When you run the mapping with a parameter set or a parameter file, the Data Integration Service uses the parameter values defined in the parameter set or parameter file instead of the default parameter values you configured in the transformation, mapping, or workflow.

You can use the following parameters to represent additional properties in the Hadoop environment:

Parameters for sources and targets

You can use parameters to represent additional properties for the following big data sources and targets:

- Complex file
- Flat file
- HBase
- HDFS
- Hive

Parameters for the Hadoop connection and run-time environment

You can set the Hive version, run-time environment, and Hadoop connection with a parameter.

For more information about mapping parameters, see the *Informatica Developer Mapping Guide*.

Parameter Usage

Use parameters for big data sources or target properties, connection properties, and run-time environment properties.

Big Data Sources and Targets

Hive sources

You can configure the following parameters for Hive Read transformation properties:

- Connection. Configure this parameter on the **Run-time** tab.
- Owner. Configure this parameter on the **Run-time** tab.
- Resource. Configure this parameter on the **Run-time** tab.
- Joiner queries. Configure this parameter on the **Query** tab.
- Filter queries. Configure this parameter on the **Query** tab.
- PreSQL commands. Configure this parameter on the **Advanced** tab.
- PostSQL commands. Configure this parameter on the **Advanced** tab.
- Constraints. Configure this parameter on the **Advanced** tab.

HBase sources and targets

You can configure the following parameters for HBase Read and Write transformation properties:

- Connection. Configure this parameter on the **Overview** tab.
- Date Time Format for the Read or Write data object. Configure this parameter on the **Advanced** tab.

Complex file sources and targets

You can configure the following parameters for complex file Read and Write transformation properties:

- Connection. Configure this parameter on the **Overview** tab.
- Data object read operation. Configure the following parameters on the **Advanced** tab:
 - File Path
 - File Format.
 - Input Format
 - Compression Format
 - Custom Compression Codec properties
- Data object write operation. Configure the following parameters on the **Advanced** tab:
 - File Name
 - File Format
 - Output Format
 - Output Key Class

- Output Value Class
- Compression Format
- Custom Compression Codec
- Sequence File Compression Type

Flat file on HDFS sources and targets

You can configure the following parameters for a flat file on HDFS Read and Write transformation properties:

- Data object read operation. Configure the following parameters on the **Run-time** tab:
 - Source File Name
 - Source File Directory
- Data object write operation. Configure the following parameters on the **Run-time** tab:
 - Output File Directory
 - Output File Name

Hadoop connection and run-time environment

You can configure the following mapping parameters on the **Run-time** tab for a mapping in the Hadoop environment:

- Hive version.
- Run-time environment.
- Hadoop connection.

INDEX

A

- accessing elements
 - array [110](#)
 - map [110](#)
 - struct [110](#)
- Address Validator transformation
 - Hadoop environment [73](#)
- aggregate window function
 - example [143](#)
 - nesting [145](#)
 - offsets [145](#)
- Aggregator transformation
 - Blaze engine [73](#)
 - Hadoop environment [73](#)
 - Hive engine [74](#)
 - Spark engine [74](#)
- Amazon AWS [213](#)
- architecture
 - Big Data Management [19](#)
 - Hadoop environment [21](#)
- array
 - complex data type [95](#)
 - accessing elements [110](#)
 - dimension [95](#)
 - example [95](#)
 - extracting element [111](#)
 - format [95](#)
 - generating [112](#)
 - index [95](#)
 - multi-dimensional [95](#)
 - one-dimensional [95](#)
- array functions
 - ARRAY [112](#)
 - COLLECT_LIST [112](#)
 - CONCAT_ARRAY [112](#)
 - SIZE [112](#)
- array port
 - type configuration [106](#)
- AutoDeploy [202](#), [203](#), [210](#)
- Azure
 - configuration [215](#)

B

- big data
 - access [15](#)
 - application services [20](#)
 - big data process [26](#)
 - data lineage [17](#)
 - repositories [20](#)
- big data process
 - collect your data [26](#)
- Big Data Streaming
 - description [17](#)

- Blaze engine
 - Blaze engine architecture [23](#)
 - connection properties [218](#)
 - mapping properties [164](#)
 - monitoring [160–162](#)
 - segment time [163](#)
 - summary report [162–164](#), [166](#)
 - tasklet execution time [164](#)
 - tasklet information [166](#)
- Blaze execution plan
 - monitoring [158](#)
- Blaze Job Monitor
 - logging [167](#)

C

- Case Converter transformation
 - Hadoop environment [75](#)
- Classifier transformation
 - Hadoop environment [75](#)
- cloud platform clusters [199](#)
- cloud provisioning configuration
 - Amazon AWS properties [213](#)
 - Microsoft Azure properties [215](#)
- cluster workflow
 - cloud provisioning connection [213](#)
- cluster workflows
 - administrator tasks [202](#)
 - Azure job monitoring [211](#)
 - cloud provisioning configuration [200](#)
 - components [200](#), [202](#), [203](#), [210](#), [211](#)
 - Create Cluster task [203](#)
 - Delete Cluster task [201](#), [211](#)
 - external RDS as the Hive metastore db [209](#)
 - Hadoop connection [200](#)
 - Mapping and other tasks [201](#)
 - mappings in cluster workflows [202](#), [210](#)
 - monitoring [211](#)
 - Overview [199](#)
 - running [211](#)
 - Running mappings on Blaze engine [208](#)
 - sample [200](#)
 - workflow tasks [210](#), [211](#)
- ClusterIdentifier property [203](#)
- Comparison transformation
 - Hadoop environment [75](#)
- complex data type
 - array [94](#)
 - map [94](#)
 - struct [94](#)
- complex data type definition
 - creating [103](#)
 - example [101](#)
 - importing [104](#)
 - nested data type definition [103](#)

- complex data type definition (*continued*)
 - recursive data type definition [103](#)
 - struct [101](#)
 - type definition library [101](#)
- complex data types
 - array [95](#)
 - map [96](#)
 - struct [97](#)
- complex file formats
 - Avro [247](#)
 - JSON [248](#)
 - overview [246](#)
 - Parquet [248](#)
- complex file sources
 - in Hadoop environment [55](#)
- complex functions
 - ARRAY [112](#)
 - CAST [112](#)
 - COLLECT_LIST [112](#)
 - COLLECT_MAP [112](#)
 - CONCAT_ARRAY [112](#)
 - MAP [112](#)
 - MAP_FROM_ARRAYS [112](#)
 - MAP_KEYS [112](#)
 - MAP_VALUES [112](#)
 - RESPEC [112](#)
 - SIZE [112](#)
 - STRUCT [112](#)
 - STRUCT_AS [112](#)
- complex operator
 - dot operator [110](#)
 - example [110](#)
 - subscript operator [110](#)
- complex port
 - array [99](#)
 - creating [101](#)
 - map [99](#)
 - nested [99](#)
 - Read transformation [100](#)
 - struct [99](#)
 - transformations [100](#)
 - type configuration [99](#), [106](#)
 - Write transformation [100](#)
- component architecture
 - clients and tools [19](#)
- configure a mapping
 - Hadoop [40](#)
- configuring
 - array port properties [106](#)
 - map port properties [108](#)
 - struct port
 - properties [109](#)
- connections
 - properties [213](#), [218](#)
 - HBase [212](#)
 - HDFS [212](#)
 - Hive [212](#)
 - JDBC [212](#)
- Consolidation transformation
 - Hadoop environment [75](#)
- conversion
 - hierarchical to relational [114](#)
- Create Cluster task [200](#)
- creating a column profile
 - profiles [190](#)

D

- Data Discovery
 - description [16](#)
- Data Integration Service grid [196](#)
- Data Masking transformation
 - Hadoop environment [75](#)
- data object
 - processing data with an intelligent structure [132](#)
- data object profiles
 - creating a single profile [188](#)
 - enterprise discovery [189](#)
- Data Processor transformation
 - Blaze engine [76](#)
- data type
 - complex data type [94](#)
 - nested data type [94](#)
 - primitive data type [94](#)
- data types
 - complex files [246](#)
 - Hive [250](#)
 - Hive complex data types [251](#)
 - processing in a Hadoop environment [259](#)
 - support [246](#)
- decimal
 - in high precision [51](#)
 - precision loss [51](#), [52](#)
- Decision transformation
 - Hadoop environment [76](#)
- Delete Cluster task [211](#)
- dot operator
 - extracting struct element [111](#)

E

- Edge Data Streaming
 - description [17](#)
- enterprise discovery
 - running in Informatica Analyst [191](#)
- ephemeral clusters
 - cloud provisioning connection [213](#)
 - Overview [199](#)
- example
 - aggregate window function [143](#)
 - array [95](#)
 - complex data type definition [101](#)
 - complex operator [110](#)
 - dot operator [110](#)
 - map [96](#)
 - nested data type definition [103](#)
 - partition and order keys [140](#)
 - struct [97](#)
 - subscript operator [110](#)
 - windowing [147](#), [149](#)
- execution plan
 - Spark engine [43](#)
- Expression transformation
 - Hadoop environment [77](#)
 - Hive engine [77](#)
- external RDS as the Hive metastore [209](#)

F

- Filter transformation
 - Hadoop environment [77](#)
- flat file source [56](#)

- flat file sources
 - in Hadoop environment [55](#)
- functions
 - processing in a Hadoop environment [259](#)

G

- grid
 - Data Integration Service [196](#)
 - description [195](#)
 - optimization [196](#)

H

- Hadoop [212](#)
- Hadoop connection for cluster workflows [200](#)
- Hadoop connections
 - creating [236](#)
- Hadoop environment
 - transformation support [70](#)
 - complex file sources [55](#)
 - flat file limitations [55](#)
 - flat file targets [63](#)
 - Hive targets [64](#)
 - Intelligent structure model sources restrictions [59](#)
 - logs [155](#)
 - optimization [44](#)
 - parameter usage [264](#)
 - parameters [263](#)
 - relational sources [60](#)
 - Sqoop sources restrictions [60](#)
 - Sqoop targets restrictions [68](#)
 - valid sources [54](#)
- Hadoop execution plan
 - description, for mapping [28](#)
 - overview [41](#)
- Hadoop mapping
 - run-time properties [29](#)
- hadoop utilities
 - Sqoop [21](#)
- HBase connections
 - MapR-DB properties [226](#)
 - properties [225](#)
- HDFS connections
 - creating [235](#)
 - properties [223](#)
- HDFS mappings
 - data extraction example [182](#)
 - description [182](#)
- hierarchical conversion
 - generate nested struct [117](#)
 - generate struct [115](#)
 - hierarchical to hierarchical [125](#)
 - hierarchical to relational [125](#), [127](#)
 - relational to hierarchical [115](#), [117](#)
- hierarchical data
 - complex files [91](#)
 - converting [114](#)
 - extracting elements [125](#)
 - flattening elements [127](#)
 - how to process [92](#)
 - modifying [125](#), [127](#)
 - processing [91](#), [130](#)
 - Spark engine [91](#)
- high availability
 - description [197](#)
- Hive
 - target limitations [64](#)
- Hive connections
 - creating [235](#)
 - properties [226](#)
- hive engine
 - high precision [52](#)
 - udf [52](#)
 - user-defined function [52](#)
- Hive engine
 - data type processing [261](#)
 - function processing [261](#)
 - rules and guidelines [261](#)
 - high precision [51](#)
 - Hive engine architecture [25](#)
 - Hive engine execution plan [43](#)
 - monitoring [173](#)
 - precision loss [51](#)
- Hive execution plan
 - monitoring [158](#)
- Hive mappings
 - description [183](#)
 - workflows [40](#)
- Hive metastore [209](#)
- Hive pushdown
 - connection properties [218](#)
- Hive query
 - description, for mapping [28](#)
- Hive query plan
 - viewing, for mapping [44](#)
 - viewing, for profile [193](#)
- Hive script
 - description, for mapping [28](#)
- Hive sources
 - postSQL [57](#)
 - preSQL [57](#)
 - with Blaze engine [58](#)
 - with Informatica mappings [56](#)
- Hive target
 - truncate [65](#)
- Hive targets
 - postSQL [65](#)
 - preSQL [65](#)
 - with Blaze engine [66](#)
- how to
 - process data with an intelligent structure [133](#)
 - process hierarchical data [92](#)

I

- Informatica Big Data Management
 - overview [14](#)
- Informatica engine
 - Informatica engine execution plan [41](#)
- Informatica Intelligent Cloud Services account
 - creating [136](#)
- intelligent structure
 - development [131](#)
 - how to process data [133](#)
 - scenario [131](#)
- Intelligent structure mappings
 - data conversion example [134](#)
- intelligent structure model
 - creating [136](#)
 - exporting [137](#)
- Intelligent structure model sources
 - in Hadoop environment [59](#)

J

- Java transformation
 - Hadoop environment [78](#)
 - Hive engine [79](#)
 - Spark engine [78](#)
- JDBC connections
 - properties [230](#)
 - Sqoop configuration [230](#)
- Joiner transformation
 - Hadoop environment [80](#)

K

- Key Generator transformation
 - Hadoop environment [81](#)

L

- Labeler transformation
 - Hadoop environment [81](#)
- logging
 - mapping run on Hadoop [167](#)
 - Spark engine [173](#)
- logs
 - Blaze engine [166](#)
 - Hadoop environment [155](#)
 - Hive engine [178](#)
 - Spark engine [172](#)
- logs URL
 - YARN web user interface [173](#)
- Lookup transformation
 - Hadoop environment [81](#)
 - Hive engine [82](#)
 - Spark engine [81](#)

M

- map
 - complex data type [96](#)
 - accessing elements [110](#)
 - example [96](#)
 - format [96](#)
 - generating [112](#)
 - key-value pair [96](#)
- map functions
 - COLLECT_MAP [112](#)
 - MAP [112](#)
 - MAP_FROM_ARRAYS [112](#)
 - MAP_KEYS [112](#)
 - MAP_VALUES [112](#)
 - SIZE [112](#)
- map port
 - type configuration [108](#)
- mapping
 - intelligent structure [132](#)
- mapping example
 - Hive [184](#)
 - Twitter [185](#)
- mapping execution plans
 - overview [41](#)
- mapping run on Hadoop
 - logging [157](#)
 - monitoring [158](#)
 - overview [28](#)

- Mapping tasks for cluster workflows [201](#)
- Match transformation
 - Hadoop environment [82](#)
- MDM Big Data Relationship Management
 - description [18](#)
- Merge transformation
 - Hadoop environment [83](#)
- Microsoft Azure [215](#)
- monitoring cluster workflow jobs [211](#)
- Monitoring URL
 - Blaze and Spark jobs [156](#)

N

- native environment
 - high availability [197](#)
 - mappings [181](#)
 - optimization [195](#)
 - partitioning [196](#)
- nested data type definition
 - example [103](#)
- nested struct
 - generating [117](#)
- node labeling
 - Blaze [49](#)
- Normalizer transformation
 - Hadoop environment [83](#)

O

- optimization
 - compress temporary staging tables [45](#)
 - node labeling [48](#)
 - queuing [48](#)
 - scheduling [48](#)
 - truncate partitions [47](#)
- overview
 - processing data with an intelligent structure model [130](#)
 - processing hierarchical data [91](#)
 - processing structured and unstructured data [130](#)

P

- parameters
 - Hadoop environment [263](#)
- Parser transformation
 - Hadoop environment [83](#)
- partitioning
 - description [196](#)
 - optimization [197](#)
- postSQL
 - for Hive sources [57](#)
 - for Hive targets [65](#)
- precision loss
 - Hive engine [52](#)
 - in high precision [51](#)
 - intermediate calculations [52](#)
- preSQL
 - for Hive sources [57](#)
 - for Hive targets [65](#)
- processing data
 - Spark engine [130](#)
- processing hierarchical data
 - overview [91](#)

- processing with an intelligent structure
 - in a mapping [132](#)
 - running on Spark [132](#)
- processing with an intelligent structure model
 - overview [130](#)
- profile run on Blaze engine
 - Overview [186](#)
- profile run on Hadoop
 - monitoring [193](#)
- profile run on Hive
 - Overview [186](#)
- profiles
 - creating a column profile [190](#)
- Python transformation
 - Hadoop environment [84](#)

R

- Rank transformation
 - Hadoop environment [84](#)
 - Hive engine [85](#)
 - Spark engine [84](#)
- registration
 - Informatica Intelligent Cloud Services account [136](#)
- relational to hierarchical [114](#)
- Router transformation
 - Hadoop environment [85](#)
- rules and guidelines
 - Spark engine [259](#)
 - Hive engine [261](#)
 - window functions [147](#)
 - windowing properties [142](#)
- run-time properties
 - Hadoop mapping [29](#)

S

- Sequence Generator transformation
 - Hadoop environment [85](#)
- social media mappings
 - description [184](#)
- Sorter transformation
 - Blaze engine [86](#)
 - Hadoop environment [86](#)
 - Hive engine [86](#)
 - Spark engine [86](#)
- source file name
 - for flat file sources [56](#)
- sources
 - in Hadoop environment [54](#)
- Spark
 - execution plans [43](#)
- Spark deploy mode
 - Hadoop connection properties [218](#)
- Spark engine
 - data type processing [259](#)
 - function processing [259](#)
 - rules and guidelines [259](#)
 - connection properties [218](#)
 - hierarchical data [91](#)
 - monitoring [169](#)
 - processing data [130](#)
- Spark Event Log directory
 - Hadoop connection properties [218](#)
- Spark execution parameters
 - Hadoop connection properties [218](#)

- Spark HDFS staging directory
 - Hadoop connection properties [218](#)
- Sqoop configuration
 - mapping properties [38](#)
 - profiling [187](#)
- Sqoop connection arguments
 - Dsqoop.connection.factories [233](#)
 - connect [233](#)
 - direct [233](#)
 - driver [233](#)
- Sqoop connectivity
 - supported data types [252](#)
- Sqoop data types
 - Aurora [252](#)
 - Greenplum [253](#)
 - IBM DB2 [252](#)
 - IBM DB2 for z/OS [252](#)
 - Microsoft SQL Server [253](#)
 - Netezza [254](#)
 - Oracle [254](#)
 - Teradata [255](#)
 - Teradata Data Types with TDCH Sqoop Specialized Connectors [255](#)
- Sqoop mapping arguments
 - batch [38](#)
 - m [37](#)
 - num-mappers [37](#)
 - split-by [37](#)
- Sqoop mappings
 - overview [36](#)
 - supported engines [36](#)
- Sqoop sources
 - in Hadoop environment [60](#)
- Sqoop targets
 - in Hadoop environment [68](#)
- Standardizer transformation
 - Hadoop environment [87](#)
- stateful computing
 - overview [138](#)
- struct
 - complex data type [97](#)
 - accessing elements [110](#)
 - changing data type [112](#)
 - changing schema [112](#)
 - complex data type definition [101](#)
 - example [97](#)
 - extracting element [111](#)
 - format [97](#)
 - generating [112](#), [115](#)
 - name-type pair [97](#)
 - schema [97](#)
- struct functions
 - CAST [112](#)
 - RESPEC [112](#)
 - STRUCT [112](#)
 - STRUCT_AS [112](#)
- struct port
 - complex data type definition [109](#)
 - type configuration [109](#)
- subscript operator
 - extracting array element [111](#)

T

- targets
 - flat files in Hadoop mapping [63](#)
- TDCH connection factory
 - Dsqoop.connection.factories [233](#)

- third-party tools
 - hadoop cluster [21](#)
- transformations
 - in Hadoop environment [70](#)
- type configuration
 - array [106](#)
 - map [106](#)
 - struct [106](#)
 - struct port [109](#)
- type definition library
 - complex data type definition [101](#)
 - mapping [101](#)
 - maplet [101](#)
 - nested data type definition [103](#)

U

- Union transformation
 - Hadoop environment [87](#)
- unstructured data
 - processing [130](#)
- Update Strategy transformation
 - Blaze engine [87](#)
 - Hadoop environment [87](#)
 - Hive engine [89](#)
 - Spark engine [88](#)
- utilities
 - hadoop cluster [21](#)

V

- validation environments
 - Hadoop mapping [29](#)

W

- Weighted Average transformation
 - Hadoop environment [90](#)
- window functions
 - aggregate functions [143](#)
 - LAG [143](#)
 - LEAD [143](#)
 - overview [142](#)
 - rules and guidelines [147](#)
- windowing
 - example [147](#), [149](#), [151](#)
 - overview [138](#)
 - properties [139](#)
- windowing properties
 - frame [139](#), [145](#)
 - order [139](#), [140](#)
 - partition [139](#), [140](#)
 - rules and guidelines [142](#)
- workflows
 - Hive mappings [40](#)

Y

- YARN
 - queues [47](#)
 - scheduler
 - capacity [47](#)
 - fair [47](#)
- YARN web user interface
 - description [155](#)