



Informatica® Big Data Management
10.1.1 Update 2

User Guide

© Copyright Informatica LLC 2012, 2018

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.

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, including without limitation: Copyright DataDirect Technologies. All rights reserved. Copyright © Sun Microsystems. All rights reserved. Copyright © RSA Security Inc. All Rights Reserved. Copyright © Ordinal Technology Corp. All rights reserved. Copyright © Aandacht c.v. All rights reserved. Copyright Genivia, Inc. All rights reserved. Copyright Isomorphic Software. All rights reserved. Copyright © Meta Integration Technology, Inc. All rights reserved. Copyright © Intalio. All rights reserved. Copyright © Oracle. All rights reserved. Copyright © Adobe Systems Incorporated. All rights reserved. Copyright © DataArt, Inc. All rights reserved. Copyright © ComponentSource. All rights reserved. Copyright © Microsoft Corporation. All rights reserved. Copyright © Rogue Wave Software, Inc. All rights reserved. Copyright © Teradata Corporation. All rights reserved. Copyright © Yahoo! Inc. All rights reserved. Copyright © Glyph & Cog, LLC. All rights reserved. Copyright © Thinkmap, Inc. All rights reserved. Copyright © Clearpace Software Limited. All rights reserved. Copyright © Information Builders, Inc. All rights reserved. Copyright © OSS Nokalva, Inc. All rights reserved. Copyright Edifecs, Inc. All rights reserved. Copyright Cleo Communications, Inc. All rights reserved. Copyright © International Organization for Standardization 1986. All rights reserved. Copyright © ej-technologies GmbH. All rights reserved. Copyright © Jaspersoft Corporation. All rights reserved. Copyright © International Business Machines Corporation. All rights reserved. Copyright © yWorks GmbH. All rights reserved. Copyright © Lucent Technologies. All rights reserved. Copyright © University of Toronto. All rights reserved. Copyright © Daniel Veillard. All rights reserved. Copyright © Unicode, Inc. Copyright IBM Corp. All rights reserved. Copyright © MicroQuill Software Publishing, Inc. All rights reserved. Copyright © PassMark Software Pty Ltd. All rights reserved. Copyright © LogiXML, Inc. All rights reserved. Copyright © 2003-2010 Lorenzi Davide, All rights reserved. Copyright © Red Hat, Inc. All rights reserved. Copyright © The Board of Trustees of the Leland Stanford Junior University. All rights reserved. Copyright © EMC Corporation. All rights reserved. Copyright © Flexera Software. All rights reserved. Copyright © Jinfonet Software. All rights reserved. Copyright © Apple Inc. All rights reserved. Copyright © Telerik Inc. All rights reserved. Copyright © BEA Systems. All rights reserved. Copyright © PDFlib GmbH. All rights reserved. Copyright © Orientation in Objects GmbH. All rights reserved. Copyright © Tanuki Software, Ltd. All rights reserved. Copyright © Ricebridge. All rights reserved. Copyright © Sencha, Inc. All rights reserved. Copyright © Scalable Systems, Inc. All rights reserved. Copyright © jqWidgets. All rights reserved. Copyright © Tableau Software, Inc. All rights reserved. Copyright © MaxMind, Inc. All Rights Reserved. Copyright © TMatte Software s.r.o. All rights reserved. Copyright © MapR Technologies Inc. All rights reserved. Copyright © Amazon Corporate LLC. All rights reserved. Copyright © Highsoft. All rights reserved. Copyright © Python Software Foundation. All rights reserved. Copyright © BeOpen.com. All rights reserved. Copyright © CNRI. All rights reserved.

This product includes software developed by the Apache Software Foundation (<http://www.apache.org/>), and/or other software which is licensed under various versions of the Apache License (the "License"). You may obtain a copy of these Licenses at <http://www.apache.org/licenses/>. Unless required by applicable law or agreed to in writing, software distributed under these Licenses is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the Licenses for the specific language governing permissions and limitations under the Licenses.

This product includes software which was developed by Mozilla (<http://www.mozilla.org/>), software copyright The JBoss Group, LLC, all rights reserved; software copyright © 1999-2006 by Bruno Lowagie and Paulo Soares and other software which is licensed under various versions of the GNU Lesser General Public License Agreement, which may be found at <http://www.gnu.org/licenses/lgpl.html>. The materials are provided free of charge by Informatica, "as-is", without warranty of any kind, either express or implied, including but not limited to the implied warranties of merchantability and fitness for a particular purpose.

The product includes ACE(TM) and TAO(TM) software copyrighted by Douglas C. Schmidt and his research group at Washington University, University of California, Irvine, and Vanderbilt University, Copyright (©) 1993-2006, all rights reserved.

This product includes software developed by the OpenSSL Project for use in the OpenSSL Toolkit (copyright The OpenSSL Project. All Rights Reserved) and redistribution of this software is subject to terms available at <http://www.openssl.org> and <http://www.openssl.org/source/license.html>.

This product includes Curl software which is Copyright 1996-2013, Daniel Stenberg, <daniel@haxx.se>. All Rights Reserved. Permissions and limitations regarding this software are subject to terms available at <http://curl.haxx.se/docs/copyright.html>. Permission to use, copy, modify, and distribute this software for any purpose with or without fee is hereby granted, provided that the above copyright notice and this permission notice appear in all copies.

The product includes software copyright 2001-2005 (©) MetaStuff, Ltd. All Rights Reserved. Permissions and limitations regarding this software are subject to terms available at <http://www.dom4j.org/license.html>.

The product includes software copyright © 2004-2007, The Dojo Foundation. All Rights Reserved. Permissions and limitations regarding this software are subject to terms available at <http://dojotoolkit.org/license>.

This product includes ICU software which is copyright International Business Machines Corporation and others. All rights reserved. Permissions and limitations regarding this software are subject to terms available at <http://source.icu-project.org/repos/icu/icu/trunk/license.html>.

This product includes software copyright © 1996-2006 Per Bothner. All rights reserved. Your right to use such materials is set forth in the license which may be found at <http://www.gnu.org/software/kawa/Software-License.html>.

This product includes OSSP UUID software which is Copyright © 2002 Ralf S. Engelschall, Copyright © 2002 The OSSP Project Copyright © 2002 Cable & Wireless Deutschland. Permissions and limitations regarding this software are subject to terms available at <http://www.opensource.org/licenses/mit-license.php>.

This product includes software developed by Boost (<http://www.boost.org/>) or under the Boost software license. Permissions and limitations regarding this software are subject to terms available at http://www.boost.org/LICENSE_1_0.txt.

This product includes software copyright © 1997-2007 University of Cambridge. Permissions and limitations regarding this software are subject to terms available at <http://www.pcre.org/license.txt>.

This product includes software copyright © 2007 The Eclipse Foundation. All Rights Reserved. Permissions and limitations regarding this software are subject to terms available at <http://www.eclipse.org/org/documents/epl-v10.php> and at <http://www.eclipse.org/org/documents/edl-v10.php>.

This product includes software licensed under the terms at <http://www.tcl.tk/software/tcltk/license.html>, <http://www.bosrup.com/web/overlib/?License>, <http://www.stlport.org/doc/license.html>, <http://asm.ow2.org/license.html>, <http://www.cryptix.org/LICENSE.TXT>, <http://hsqldb.org/web/hsqLicense.html>, <http://httpunit.sourceforge.net/doc/license.html>, <http://jung.sourceforge.net/license.txt>, http://www.gzip.org/zlib/zlib_license.html, <http://www.openldap.org/software/release/license.html>, <http://www.libssh2.org>, <http://slf4j.org/license.html>, <http://www.sente.ch/software/OpenSourceLicense.html>, <http://fusesource.com/downloads/license-agreements/fuse-message-broker-v-5-3-license-agreement>, <http://antlr.org/license.html>, <http://aopalliance.sourceforge.net/>, <http://www.bouncycastle.org/license.html>, <http://www.jgraph.com/jgraphdownload.html>, <http://www.jcraft.com/jsch/LICENSE.txt>, http://jotm.objectweb.org/bsd_license.html, <http://www.w3.org/Consortium/Legal/2002/copyright-software-20021231>, <http://www.slf4j.org/license.html>, <http://nanoxml.sourceforge.net/orig/copyright.html>, <http://www.json.org/license.html>, <http://forge.ow2.org/projects/javaservice/>, <http://www.postgresql.org/about/license.html>, <http://www.sqlite.org/copyright.html>, <http://www.tcl.tk/software/tcltk/license.html>, <http://www.jaxen.org/faq.html>, <http://www.jdom.org/docs/faq.html>, <http://www.slf4j.org/license.html>, <http://www.iodbc.org/dataspace/iodbc/wiki/IODBC/License>, <http://www.keplerproject.org/md5/license.html>, <http://www.toedter.com/en/jcalendar/license.html>, <http://www.edankert.com/bounce/index.html>, <http://www.net-snmp.org/about/license.html>, <http://www.openmdx.org/#FAQ>, http://www.php.net/license/3_01.txt, <http://srp.stanford.edu/license.txt>;

<http://www.schneider.com/blowfish.html>; <http://www.jmock.org/license.html>; <http://xsom.java.net>; <http://benalman.com/about/license/>; <https://github.com/CreateJS/EaselJS/blob/master/src/easeljs/display/Bitmap.js>; <http://www.h2database.com/html/license.html#summary>; <http://jsoncpp.sourceforge.net/LICENSE>; <http://jdbc.postgresql.org/license.html>; <http://protobuf.googlecode.com/svn/trunk/src/google/protobuf/descriptor.proto>; <https://github.com/rantav/hector/blob/master/LICENSE>; <http://web.mit.edu/Kerberos/krb5-current/doc/mitK5license.html>; <http://jibx.sourceforge.net/jibx-license.html>; <https://github.com/lyokato/libgeohash/blob/master/LICENSE>; <https://github.com/hjiang/jsonxx/blob/master/LICENSE>; <https://code.google.com/p/lz4/>; <https://github.com/jedisct1/libsodium/blob/master/LICENSE>; <http://one-jar.sourceforge.net/index.php?page=documents&file=license>; <https://github.com/EsotericSoftware/kryo/blob/master/license.txt>; <http://www.scala-lang.org/license.html>; <https://github.com/tinkerpop/blueprints/blob/master/LICENSE.txt>; <http://gee.cs.oswego.edu/dl/classes/EDU/oswego/cs/dl/util/concurrent/intro.html>; <https://aws.amazon.com/asl/>; <https://github.com/twbs/bootstrap/blob/master/LICENSE>; <https://sourceforge.net/p/xmlunit/code/HEAD/tree/trunk/LICENSE.txt>; <https://github.com/documentcloud/underscore-contrib/blob/master/LICENSE>, and <https://github.com/apache/hbase/blob/master/LICENSE.txt>.

This product includes software licensed under the Academic Free License (<http://www.opensource.org/licenses/afl-3.0.php>), the Common Development and Distribution License (<http://www.opensource.org/licenses/cddl1.php>), the Common Public License (<http://www.opensource.org/licenses/cpl1.0.php>), the Sun Binary Code License Agreement Supplemental License Terms, the BSD License (<http://www.opensource.org/licenses/bsd-license.php>), the new BSD License (<http://opensource.org/licenses/BSD-3-Clause>), the MIT License (<http://www.opensource.org/licenses/mit-license.php>), the Artistic License (<http://www.opensource.org/licenses/artistic-license-1.0>) and the Initial Developer's Public License Version 1.0 (<http://www.firebirdsql.org/en/initial-developer-s-public-license-version-1-0/>).

This product includes software copyright © 2003-2006 Joe Walnes, 2006-2007 XStream Committers. All rights reserved. Permissions and limitations regarding this software are subject to terms available at <http://xstream.codehaus.org/license.html>. This product includes software developed by the Indiana University Extreme! Lab. For further information please visit <http://www.extreme.indiana.edu/>.

This product includes software Copyright (c) 2013 Frank Balluffi and Markus Moeller. All rights reserved. Permissions and limitations regarding this software are subject to terms of the MIT license.

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

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

NOTICES

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

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

The information in this documentation is subject to change without notice. If you find any problems in this documentation, please report them to us in writing at Informatica LLC 2100 Seaport Blvd. Redwood City, CA 94063.

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: 2018-07-03

Table of Contents

Preface	9
Informatica Resources.	9
Informatica Network.	9
Informatica Knowledge Base.	9
Informatica Documentation.	9
Informatica Product Availability Matrixes.	10
Informatica Velocity.	10
Informatica Marketplace.	10
Informatica Global Customer Support.	10
 Chapter 1: Introduction to Informatica Big Data Management.....	 11
Informatica Big Data Management Overview.	11
Example.	12
Big Data Management Tasks	12
Read from and Write to Big Data Sources and Targets.	12
Perform Data Discovery.	13
Perform Data Lineage on Big Data Sources.	14
Stream Machine Data.	14
Manage Big Data Relationships.	14
Big Data Management Component Architecture.	15
Clients and Tools.	15
Application Services.	16
Repositories.	16
Hadoop Environment.	16
Hadoop Utilities.	17
Big Data Management Engines.	18
Blaze Engine Architecture.	18
Spark Engine Architecture.	19
Hive Engine Architecture.	20
Big Data Process.	21
Step 1. Collect the Data.	21
Step 2. Cleanse the Data.	21
Step 3. Transform the Data.	22
Step 4. Process the Data.	22
Step 5. Monitor Jobs.	22
 Chapter 2: Connections.....	 23
Connections.	23
Hadoop Connection Properties.	24
HDFS Connection Properties.	33

HBase Connection Properties.	34
Hive Connection Properties.	35
JDBC Connection Properties.	41
Sqoop Connection-Level Arguments.	43
Creating a Connection to Access Sources or Targets.	45
Creating a Hadoop Connection.	46
Chapter 3: Mappings in the Hadoop Environment.	48
Mappings in the Hadoop Environment Overview.	48
Mapping Run-time Properties.	49
Validation Environments.	49
Execution Environment.	50
Data Warehouse Optimization Mapping Example	51
Sqoop Mappings in a Hadoop Environment.	53
Sqoop Mapping-Level Arguments.	54
Configuring Sqoop Properties in the Mapping.	55
Rules and Guidelines for Mappings in a Hadoop Environment.	56
Workflows that Run Mappings in a Hadoop Environment.	56
Configuring a Mapping to Run in a Hadoop Environment.	57
Mapping Execution Plans.	57
Blaze Engine Execution Plan Details.	58
Spark Engine Execution Plan Details.	59
Hive Engine Execution Plan Details.	60
Viewing the Execution Plan for a Mapping in the Developer Tool.	60
Optimization for the Hadoop Environment.	60
Blaze Engine High Availability.	61
Truncating Partitions in a Hive Target.	61
Enabling Data Compression on Temporary Staging Tables.	62
Scheduling, Queuing, and Node Labeling.	63
Parallel Sorting.	65
Troubleshooting a Mapping in a Hadoop Environment.	65
Chapter 4: Mapping Objects in the Hadoop Environment.	67
Sources in a Hadoop Environment.	67
Flat File Sources.	68
Hive Sources.	68
Complex File Sources.	71
Relational Sources.	71
Sqoop Sources.	72
Targets in a Hadoop Environment.	73
Flat File Targets.	73
HDFS Flat File Targets.	73
Hive Targets.	74

Complex File Targets.	76
Relational Targets.	76
Sqoop Targets.	77
Transformations in a Hadoop Environment.	77
Transformation Support on the Blaze Engine.	79
Transformation Support on the Spark Engine.	82
Transformation Support on the Hive Engine.	83
Function and Data Type Processing.	86
Rules and Guidelines for Spark Engine Processing.	86
Rules and Guidelines for Hive Engine Processing.	87
Chapter 5: Monitoring Mappings in the Hadoop Environment.	88
Monitoring Mappings in the Hadoop Environment Overview.	88
Hadoop Environment Logs.	88
YARN Web User Interface.	89
Accessing the Monitoring URL.	89
Viewing Hadoop Environment Logs in the Administrator Tool.	90
Monitoring a Mapping.	91
Blaze Engine Monitoring.	92
Blaze Job Monitoring Application.	93
Blaze Summary Report.	94
Blaze Engine Logs.	98
Viewing Blaze Logs.	99
Troubleshooting Blaze Monitoring.	100
Spark Engine Monitoring.	100
Spark Engine Logs.	102
Viewing Spark Logs.	102
Hive Engine Monitoring.	103
Hive Engine Logs.	104
Chapter 6: Mappings in the Native Environment.	105
Mappings in the Native Environment Overview.	105
Data Processor Mappings.	105
HDFS Mappings.	106
HDFS Data Extraction Mapping Example.	106
Hive Mappings.	107
Hive Mapping Example.	108
Social Media Mappings.	108
Twitter Mapping Example.	109
Chapter 7: Profiles.	110
Profiles Overview.	110
Native Environment.	110

Hadoop Environment.	111
Column Profiles for Sqoop Data Sources.	111
Creating a Single Data Object Profile in Informatica Developer.	112
Creating an Enterprise Discovery Profile in Informatica Developer.	113
Creating a Column Profile in Informatica Analyst.	114
Creating an Enterprise Discovery Profile in Informatica Analyst.	115
Creating a Scorecard in Informatica Analyst.	116
Monitoring a Profile.	117
Troubleshooting.	117
Chapter 8: Native Environment Optimization.	119
Native Environment Optimization Overview.	119
Processing Big Data on a Grid.	119
Data Integration Service Grid.	120
Grid Optimization.	120
Processing Big Data on Partitions.	120
Partitioned Model Repository Mappings.	120
Partition Optimization.	121
High Availability.	121
Appendix A: Data Type Reference.	123
Data Type Reference Overview.	123
Transformation Data Type Support in a Hadoop Environment.	123
Hive Data Types and Transformation Data Types.	124
Hive Complex Data Types.	126
Sqoop Data Types.	126
Aurora Data Types.	126
IBM DB2 and DB2 for z/OS Data Types.	127
Greenplum Data Types.	127
Microsoft SQL Server Data Types.	128
Netezza Data Types.	128
Oracle Data Types.	129
Teradata Data Types.	129
Teradata Data Types with TDCH Specialized Connectors for Sqoop.	130
Appendix B: Function Reference.	131
Function Support in a Hadoop Environment.	131
Appendix C: Parameter Reference.	134
Parameters Overview.	134
Parameter Usage.	135

Appendix D: Multiple Blaze Instances on a Cluster.....	137
Overview.	137
Step 1. Prepare the Hadoop Cluster for the Blaze Engine.	138
Create a Blaze User Account.	138
Create Blaze Engine Directories and Grant Permissions.	139
Grant Permissions on the Hive Source Database.	139
Step 2. Configure Data Integration Service Properties.	139
Step 3. Update hadoopEnv.properties.	141
Step 4. Create a Hadoop Connection.	143
Step 5. Set Mapping Preferences.	145
Result.	146
Index.	147

Preface

The *Informatica Big Data Management® User Guide* provides information about how to configure Informatica products for Hadoop.

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, 11](#)
- [Big Data Management Tasks , 12](#)
- [Big Data Management Component Architecture, 15](#)
- [Big Data Management Engines, 18](#)
- [Big Data Process, 21](#)

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.

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 parser transformations to read and transform unstructured and semi-structured data. For example, you can use the Data Processor transformation in a workflow to parse a Microsoft Word file to load customer and order data into relational database tables.

You can use HParser 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 Live Data Map 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. Live Data Map extracts metadata about entities from the cluster based on the entity type.

Live Data Map 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: Live Data Map does not extract metadata for MapReduce job templates or executions.

For more information, see the *Live Data Map Administrator Guide*.

Stream Machine Data

You can stream machine data in real time. To stream machine data, use Informatica Vibe Data Stream for Machine Data (Vibe Data Stream).

Vibe Data Stream 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. Vibe Data Stream 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*.

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.

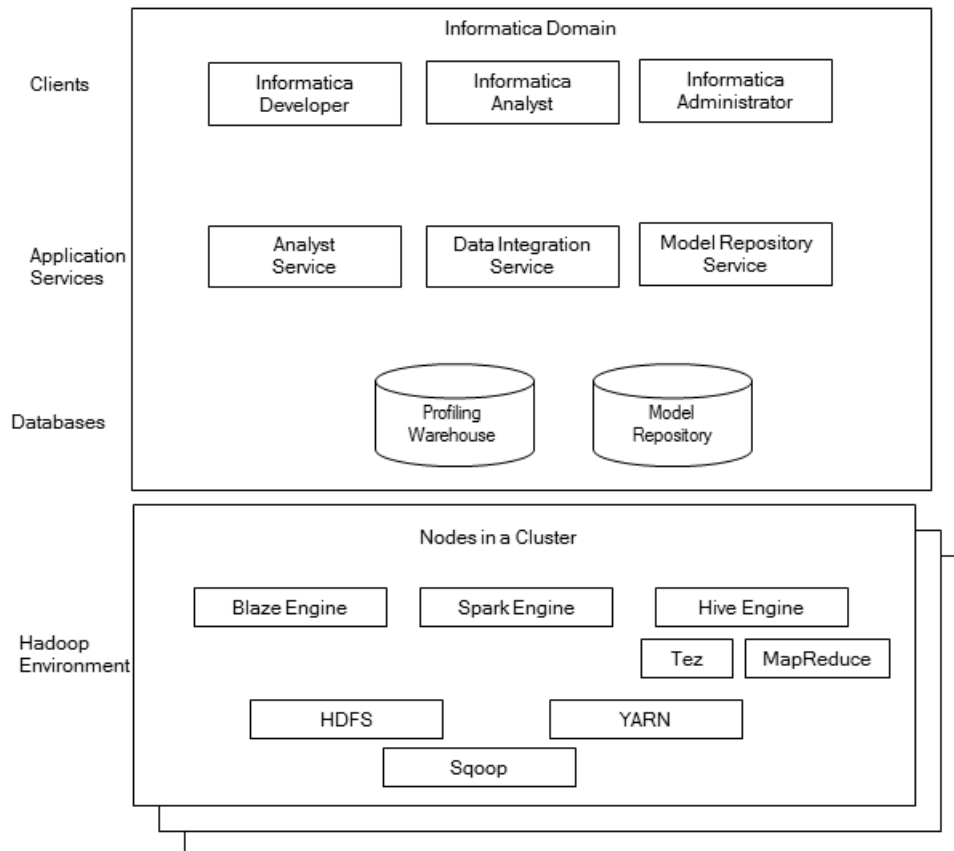
MDM Big Data Relationship Manager indexes and links the data based on the indexing and matching rules. You can configure rules based on which to link the input records. MDM Big Data Relationship Manager uses the rules to match the input records and then group all the matched records. MDM Big Data Relationship 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*.

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.

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 connects to Hadoop clusters that are distributed by third parties. 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 Hadoop as a data source and push job processing to the Hadoop cluster. It can also connect to HDFS, which enables high performance access to files across the cluster. It can connect to Hive, which is a data warehouse that connects to HDFS and uses SQL-like queries to run MapReduce jobs on Hadoop, or YARN, which can manage Hadoop clusters more efficiently. 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.

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 when you read data from or write data to Teradata:

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

Cloudera Connector Powered by Teradata and Hortonworks Connector for Teradata are specialized Sqoop plug-ins that Cloudera and Hortonworks provide for Teradata. These TDCH Sqoop Connectors use native protocols to connect to the Teradata database.

You can configure the Cloudera Connector Powered by Teradata and Hortonworks Connector for Teradata when you run Sqoop mappings on the Blaze engine.

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 the environment to run the mapping in, native environment or Hadoop environment. If you run the mapping in the Hadoop environment, the mapping will run on the Blaze engine, the Spark engine, or the Hive engine.

When you validate a mapping, you can validate it against one or all of the engines. The Developer tool returns validation messages for each engine. You can then choose to run the mapping in the native environment or in the Hadoop environment. When you run the mapping in the native environment, the Data Integration Service processes the mapping logic. When you run the mapping in the Hadoop environment, the Data Integration Service uses a proprietary rule-based methodology to determine the best engine to run the mapping. The rule-based methodology evaluates the mapping sources and the mapping logic to determine the engine. The Data Integration Service translates the mapping logic into code that the engine can process, and it transfers the code to the engine.

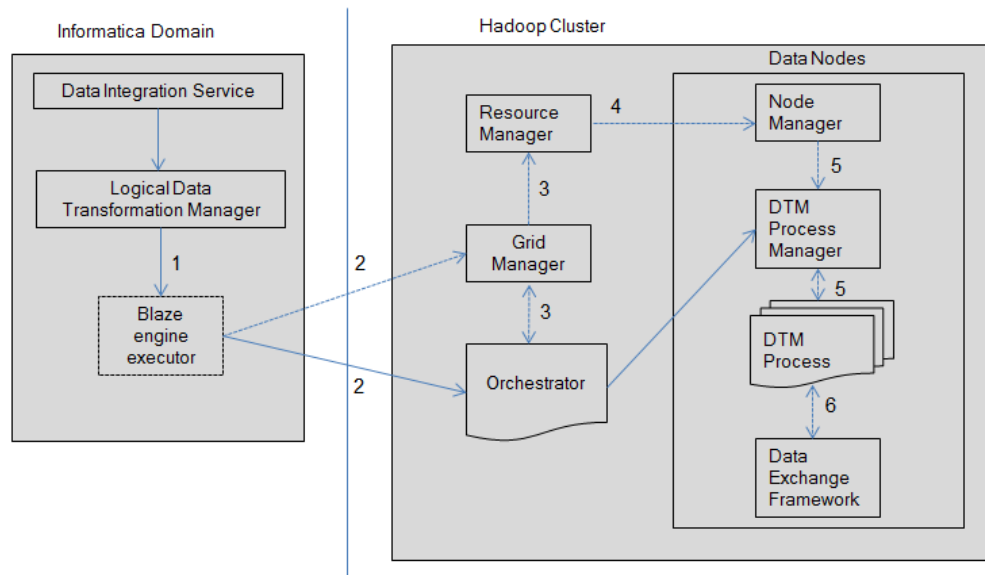
Blaze Engine Architecture

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.
- 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 receives a job request from the LDTM.

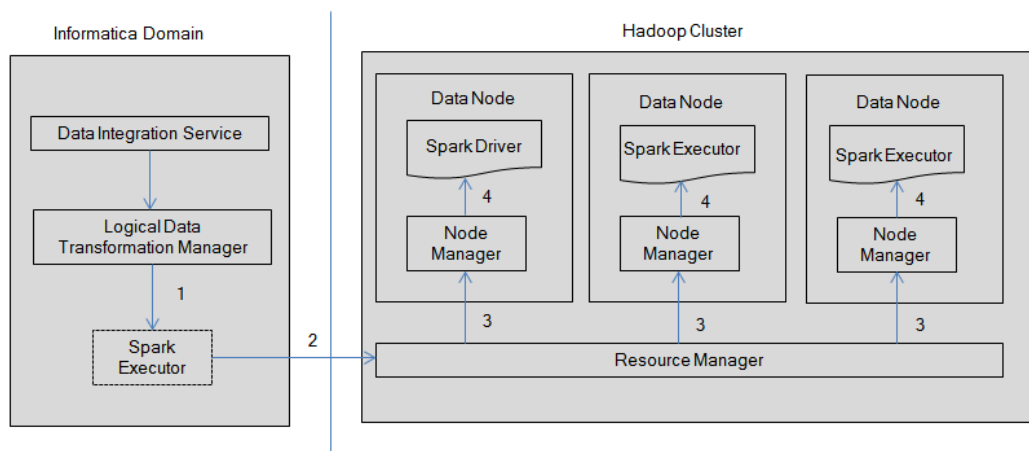
2. It initializes the communication with the Grid Manager to initialize Blaze engine components on the Hadoop cluster, and it queries the Grid Manager for an available Orchestrator.
3. The Orchestrator communicates with the Grid Manager and the Resource Manager for available resources on the Hadoop cluster.
4. The Resource Manager sends a job request to the Node Manager on the Hadoop cluster.
5. The Node Manager sends the tasks to the DTM Processes through the DTM Process Manager.
6. The DTM Processes communicate with the Data Exchange Framework to send and receive data across processing units that run on the cluster nodes.

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:

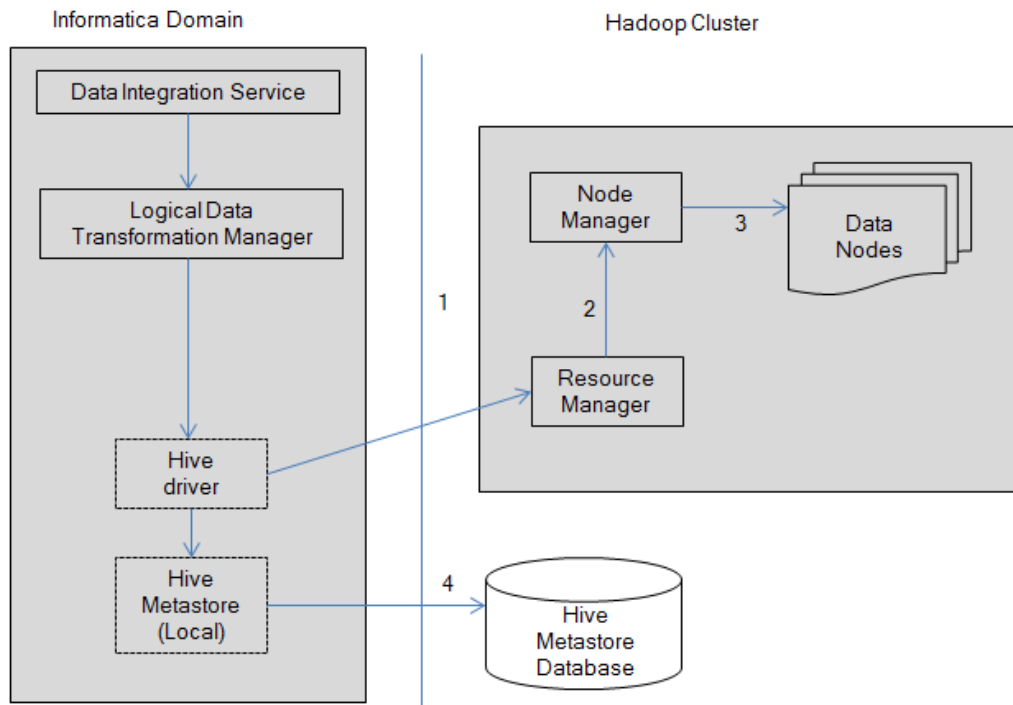
1. The Logical Data Transformation Manager translates the mapping into a Scala program, packages it as an application, and sends it to the Spark executor.
2. The Spark executor submits the application to the Resource Manager in the Hadoop cluster and requests resources to run the application.
Note: When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

The following diagram shows the architecture of how a Hadoop cluster processes MapReduce 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 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 jobs.
3. The Node Manager assigns MapReduce 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

Connections

This chapter includes the following topics:

- [Connections, 23](#)
- [Hadoop Connection Properties, 24](#)
- [HDFS Connection Properties, 33](#)
- [HBase Connection Properties, 34](#)
- [Hive Connection Properties, 35](#)
- [JDBC Connection Properties, 41](#)
- [Creating a Connection to Access Sources or Targets, 45](#)
- [Creating a Hadoop Connection, 46](#)

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

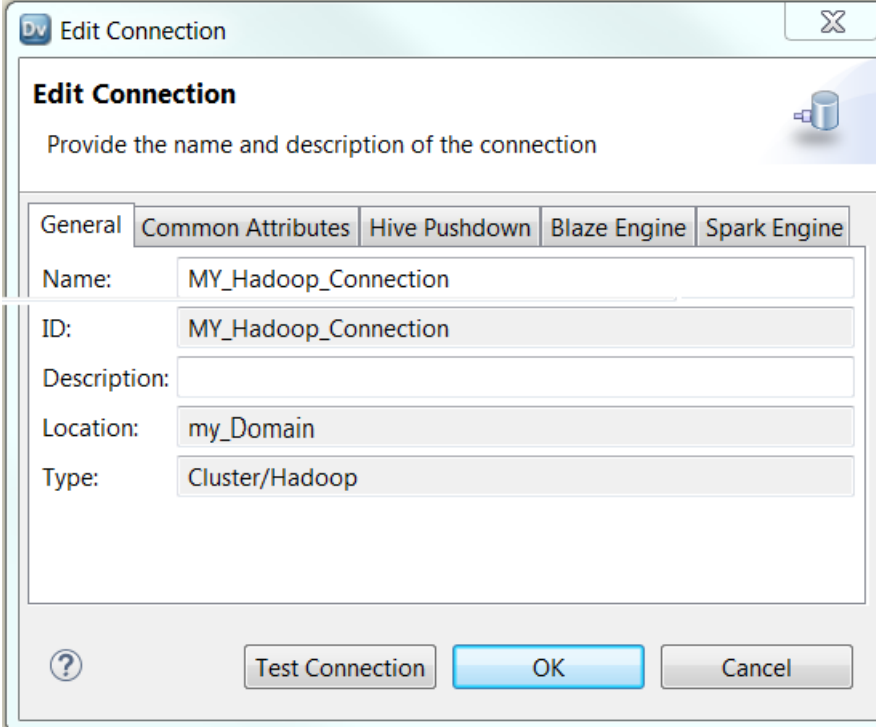
Create a JDBC connection and configure Sqoop properties in the connection to import and export relational data through Sqoop.

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.

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.

The following image shows the Hadoop connection properties:



The screenshot shows a dialog box titled "Edit Connection" with a close button (X) in the top right corner. Below the title bar, there is a section titled "Edit Connection" with a database icon and the instruction "Provide the name and description of the connection". Below this, there are five tabs: "General", "Common Attributes", "Hive Pushdown", "Blaze Engine", and "Spark Engine". The "General" tab is selected. The "General" tab contains the following fields:

- Name: MY_Hadoop_Connection
- ID: MY_Hadoop_Connection
- Description: (empty field)
- Location: my_Domain
- Type: Cluster/Hadoop

At the bottom of the dialog box, there is a question mark icon, a "Test Connection" button, an "OK" button, and a "Cancel" button.

General Properties

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.
Location	The domain where you want to create the connection. Select the domain name.
Type	The connection type. Select Hadoop.

Common Attributes - Hadoop Cluster Properties

The following table describes the connection properties that you configure for the Hadoop cluster:

Property	Description
Resource Manager Address	<p>The service within Hadoop that submits requests for resources or spawns YARN applications.</p> <p>Use the following format:</p> <pre><hostname>:<port></pre> <p>Where</p> <ul style="list-style-type: none">- <hostname> is the host name or IP address of the Yarn resource manager.- <port> is the port on which the Yarn resource manager listens for remote procedure calls (RPC). <p>For example, enter: <code>myhostame:8032</code></p> <p>You can also get the Resource Manager Address property from <code>yarn-site.xml</code> located in the following directory on the Hadoop cluster: <code>/etc/hadoop/conf/</code></p> <p>The Resource Manager Address appears as the following property in <code>yarn-site.xml</code>:</p> <pre><property> <name>yarn.resourcemanager.address</name> <value>hostname:port</value> <description>The address of the applications manager interface in the Resource Manager.</description> </property></pre> <p>Optionally, if the <code>yarn.resourcemanager.address</code> property is not configured in <code>yarn-site.xml</code>, you can find the host name from the <code>yarn.resourcemanager.hostname</code> or <code>yarn.resourcemanager.scheduler.address</code> properties in <code>yarn-site.xml</code>. You can then configure the Resource Manager Address in the Hadoop connection with the following value: <code>hostname:8032</code></p>
Default File System URI	<p>The URI to access the default Hadoop Distributed File System.</p> <p>Use the following connection URI:</p> <pre>hdfs://<node name>:<port></pre> <p>Where</p> <ul style="list-style-type: none">- <node name> is the host name or IP address of the NameNode.- <port> is the port on which the NameNode listens for remote procedure calls (RPC). <p>For example, enter: <code>hdfs://myhostname:8020/</code></p> <p>You can also get the Default File System URI property from <code>core-site.xml</code> located in the following directory on the Hadoop cluster: <code>/etc/hadoop/conf/</code></p> <p>Use the value from the <code>fs.defaultFS</code> property found in <code>core-site.xml</code>.</p> <p>For example, use the following value:</p> <pre><property> <name>fs.defaultFS</name> <value>hdfs://localhost:8020</value> </property></pre> <p>If the Hadoop cluster runs MapR, use the following URI to access the MapR File system: <code>maprfs:///</code>.</p> <p>The Azure HDInsight File System default file system can be Windows Azure Storage Blob (WASB) or Azure Data Lake Store (ADLS).</p> <p>If the cluster uses WASB storage, use the following string to specify the URI:</p> <pre>wasb://<container_name>@<account_name>.blob.core.windows.net/<path></pre> <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.

Property	Description
	<p>Example:</p> <pre>wasb://infabdmoffering1storage.blob.core.windows.net/infabdmoffering1cluster/mr-history</pre> <p>If the cluster uses ADLS storage, use the following format to specify the URI: <code>adl://home</code></p> <p>The following is the <code>fs.defaultFS</code> property as it appears in <code>hdfs-site.xml</code>:</p> <pre><property>fs.defaultFS</property> <value>adl://home</value></pre>

Common Attributes - Common Properties

The following table describes the common connection properties that you configure for the Hadoop connection:

Property	Description
Impersonation User Name	<p>User name of the user that the Data Integration Service impersonates to run mappings on a Hadoop cluster.</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.</p> <p>Note: You must use user impersonation for the Hadoop connection if the Hadoop cluster uses Kerberos authentication.</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 operating system profile user name of the machine that runs the Data Integration Service.</p>
Temporary Table Compression Codec	Hadoop compression library for a compression codec class name.
Codec Class Name	Codec class name that enables data compression and improves performance on temporary staging tables.
Hadoop Connection Custom Properties	<p>Custom properties that are unique to the Hadoop connection. You can specify multiple properties.</p> <p>Use the following format:</p> <pre><property1>=<value></pre> <p>Where</p> <ul style="list-style-type: none"> - <code><property1></code> is a Blaze, Hive, or Hadoop property. - <code><value></code> is the value of the Hive or Hadoop property. <p>To specify multiple properties use <code>&</code> as the property separator.</p> <p>Use custom properties only at the request of Informatica Global Customer Support.</p>

Hive Pushdown - Hive Pushdown Configuration

The following table describes the connection properties that you configure to push mapping logic to the Hadoop cluster:

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. The path must be the fully qualified path to the JAR files used for user-defined functions. Set the parameter <code>hive.aux.jars.path</code> with all the entries in <code>infapdo.aux.jars.path</code> and the path to the JAR files for user-defined functions.- 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. The following sample text shows two values that can be used for the Environment SQL property: <pre>set hive.execution.engine='tez';set hive.exec.dynamic.partition.mode='nonstrict';</pre>
Database Name	<p>Namespace for tables. Use the name <code>default</code> for tables that do not have a specified database name.</p>
Hive Warehouse Directory on HDFS	<p>The absolute HDFS file path of the default database for the warehouse that is local to the cluster. For example, the following file path specifies a local warehouse: <code>/user/hive/warehouse</code></p> <p>For Cloudera CDH, if the Metastore Execution Mode is remote, then the file path must match the file path specified by the Hive Metastore Service on the Hadoop cluster.</p> <p>You can get the value for the Hive Warehouse Directory on HDFS from the <code>hive.metastore.warehouse.dir</code> property in <code>hive-site.xml</code> located in the following directory on the Hadoop cluster: <code>/etc/hadoop/conf/</code></p> <p>For example, use the following value:</p> <pre><property> <name>hive.metastore.warehouse.dir</name> <value>/usr/hive/warehouse </value> <description>location of the warehouse directory</description> </property></pre> <p>For MapR, <code>hive-site.xml</code> is located in the following directory: <code>/opt/mapr/hive/<hive version>/conf</code>.</p>

Hive Pushdown - Hive Configuration

You can use the values for Hive configuration properties from `hive-site.xml` or `mapred-site.xml` located in the following directory on the Hadoop cluster: `/etc/hadoop/conf/`.

The following table describes the connection properties that you configure for the Hive engine:

Property	Description
Metastore Execution Mode	<p>Controls whether to connect to a remote metastore or a local metastore. By default, local is selected. For a local metastore, you must specify the Metastore Database URI, Metastore Database Driver, Username, and Password. For a remote metastore, you must specify only the Remote Metastore URI.</p> <p>You can get the value for the Metastore Execution Mode from hive-site.xml. The Metastore Execution Mode appears as the following property in hive-site.xml:</p> <pre><property> <name>hive.metastore.local</name> <value>true</true> </property></pre> <p>Note: The <code>hive.metastore.local</code> property is deprecated in hive-site.xml for Hive server versions 0.9 and above. If the <code>hive.metastore.local</code> property does not exist but the <code>hive.metastore.uris</code> property exists, and you know that the Hive server has started, you can set the connection to a remote metastore.</p>
Metastore Database URI	<p>The JDBC connection URI used to access the data store in a local metastore setup. Use the following connection URI:</p> <pre>jdbc:<datastore type>://<node name>:<port>/<database name></pre> <p>where</p> <ul style="list-style-type: none"> - <code><node name></code> is the host name or IP address of the data store. - <code><data store type></code> is the type of the data store. - <code><port></code> is the port on which the data store listens for remote procedure calls (RPC). - <code><database name></code> is the name of the database. <p>For example, the following URI specifies a local metastore that uses MySQL as a data store:</p> <pre>jdbc:mysql://hostname23:3306/metastore</pre> <p>You can get the value for the Metastore Database URI from hive-site.xml. The Metastore Database URI appears as the following property in hive-site.xml:</p> <pre><property> <name>javax.jdo.option.ConnectionURL</name> <value>jdbc:mysql://MYHOST/metastore</value> </property></pre>
Metastore Database Driver	<p>Driver class name for the JDBC data store. For example, the following class name specifies a MySQL driver:</p> <pre>com.mysql.jdbc.Driver</pre> <p>You can get the value for the Metastore Database Driver from hive-site.xml. The Metastore Database Driver appears as the following property in hive-site.xml:</p> <pre><property> <name>javax.jdo.option.ConnectionDriverName</name> <value>com.mysql.jdbc.Driver</value> </property></pre>
Metastore Database User Name	<p>The metastore database user name.</p> <p>You can get the value for the Metastore Database User Name from hive-site.xml. The Metastore Database User Name appears as the following property in hive-site.xml:</p> <pre><property> <name>javax.jdo.option.ConnectionUserName</name> <value>hiveuser</value> </property></pre>

Property	Description
Metastore Database Password	<p>The password for the metastore user name.</p> <p>You can get the value for the Metastore Database Password from hive-site.xml. The Metastore Database Password appears as the following property in hive-site.xml:</p> <pre><property> <name>javax.jdo.option.ConnectionPassword</name> <value>password</value> </property></pre>
Remote Metastore URI	<p>The metastore URI used to access metadata in a remote metastore setup. For a remote metastore, you must specify the Thrift server details.</p> <p>Use the following connection URI: thrift://<hostname>:<port></p> <p>Where</p> <ul style="list-style-type: none"> - <hostname> is name or IP address of the Thrift metastore server. - <port> is the port on which the Thrift server is listening. <p>For example, enter: thrift://myhostname:9083/</p> <p>You can get the value for the Remote Metastore URI from hive-site.xml. The Remote Metastore URI appears as the following property in hive-site.xml:</p> <pre><property> <name>hive.metastore.uris</name> <value>thrift://<n.n.n.n>:9083</value> <description> IP address or fully-qualified domain name and port of the metastore host</description> </property></pre>

Property	Description
Engine Type	<p>The engine that the Hadoop environment uses to run a mapping on the Hadoop cluster. Select a value from the drop down list.</p> <p>For example select: MRv2</p> <p>To set the engine type in the Hadoop connection, you must get the value for the <code>mapreduce.framework.name</code> property from <code>mapred-site.xml</code> located in the following directory on the Hadoop cluster: <code>/etc/hadoop/conf/</code></p> <p>If the value for <code>mapreduce.framework.name</code> is <code>classic</code>, select <code>mrsv1</code> as the engine type in the Hadoop connection.</p> <p>If the value for <code>mapreduce.framework.name</code> is <code>yarn</code>, you can select the <code>mrsv2</code> or <code>tez</code> as the engine type in the Hadoop connection. Do not select Tez if Tez is not configured for the Hadoop cluster.</p> <p>You can also set the value for the engine type in <code>hive-site.xml</code>. The engine type appears as the following property in <code>hive-site.xml</code>:</p> <pre><property> <name>hive.execution.engine</name> <value>tez</value> <description>Chooses execution engine. Options are: mr (MapReduce, default) or tez (Hadoop 2 only)</description> </property></pre>
Job Monitoring URL	<p>The URL for the MapReduce JobHistory server. You can use the URL for the JobTracker URI if you use MapReduce version 1.</p> <p>Use the following format:</p> <pre><hostname>:<port></pre> <p>Where</p> <ul style="list-style-type: none"> - <code><hostname></code> is the host name or IP address of the JobHistory server. - <code><port></code> is the port on which the JobHistory server listens for remote procedure calls (RPC). <p>For example, enter: <code>myhostname:8021</code></p> <p>You can get the value for the Job Monitoring URL from <code>mapred-site.xml</code>. The Job Monitoring URL appears as the following property in <code>mapred-site.xml</code>:</p> <pre><property> <name>mapred.job.tracker</name> <value>myhostname:8021 </value> <description>The host and port that the MapReduce job tracker runs at.</description> </property></pre>

Blaze Engine

The following table describes the connection properties that you configure for the Blaze engine:

Property	Description
Temporary Working Directory on HDFS	<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>For example, enter: <code>/blaze/workdir</code></p>
Blaze Service User Name	The operating system profile user name for the Blaze engine.
Minimum Port	<p>The minimum value for the port number range for the Blaze engine.</p> <p>For example, enter: <code>12300</code></p>

Property	Description
Maximum Port	The maximum value for the port number range for the Blaze engine. For example, enter: 12600
Yarn Queue Name	The YARN scheduler queue name used by the Blaze engine that specifies available resources on a cluster. The name is case sensitive.
Blaze Service Custom Properties	Custom properties that are unique to the Blaze engine. You can specify multiple properties. Use the following format: <property1>=<value> Where <ul style="list-style-type: none"> - <property1> is a Blaze engine optimization property. - <value> is the value of the Blaze engine optimization property. To enter multiple properties, separate each name-value pair with the following text: & : . Use custom properties only at the request of Informatica Global Customer Support.

Spark Engine

The following table describes the connection properties that you configure for the Spark engine:

Property	Description
Spark HDFS Staging Directory	The HDFS file path of the directory that the Spark engine uses to store temporary files for running jobs. The YARN user, Spark engine user, and mapping impersonation user must have write permission on this directory.
Spark Event Log Directory	Optional. The HDFS file path of the directory that the Spark engine uses to log events. The Data Integration Service accesses the Spark event log directory to retrieve final source and target statistics when a mapping completes. These statistics appear on the Summary Statistics tab and the Detailed Statistics tab of the Monitoring tool. If you do not configure the Spark event log directory, the statistics might be incomplete in the Monitoring tool.
Spark Execution Parameters	An optional list of configuration parameters to apply to the Spark engine. You can change the default Spark configuration properties values, such as <code>spark.executor.memory</code> or <code>spark.driver.cores</code> . Use the following format: <property1>=<value> <ul style="list-style-type: none"> - <property1> is a Spark configuration property. - <value> is the value of the property. For example, you can configure a YARN scheduler queue name that specifies available resources on a cluster: <code>spark.yarn.queue=TestQ</code> To enter multiple properties, separate each name-value pair with the following text: & :

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	<p>The URI to access HDFS.</p> <p>Use the following format to specify the NameNode URI in Cloudera and Hortonworks distributions:</p> <pre>hdfs://<namenode>:<port></pre> <p>Where</p> <ul style="list-style-type: none">- <namenode> is the host name or IP address of the NameNode.- <port> is the port that the NameNode listens for remote procedure calls (RPC). <p>Use one of the following formats to specify the NameNode URI in MapR distribution:</p> <ul style="list-style-type: none">- maprfs:///- maprfs:///mapr/my.cluster.com/ <p>Where my.cluster.com is the cluster name that you specify in the mapr-clusters.conf file.</p>

HBase Connection Properties

Use an HBase connection to access HBase. The HBase connection is a NoSQL connection. You can create and manage an HBase connection in the Administrator tool or the Developer tool. Hbase connection properties are case sensitive unless otherwise noted.

The following table describes HBase connection properties:

Property	Description
Name	The name of the connection. The name is not case sensitive and must be unique within the domain. You can change this property after you create the connection. The name cannot exceed 128 characters, contain spaces, or contain the following special characters: ~ ` ! \$ % ^ & * () - + = { [] } \ : ; " ' < , > . ? /
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. The description cannot exceed 4,000 characters.
Location	The domain where you want to create the connection.
Type	The connection type. Select HBase.
ZooKeeper Host(s)	Name of the machine that hosts the ZooKeeper server. The name is case sensitive. When the ZooKeeper runs in the replicated mode, specify a comma-separated list of servers in the ZooKeeper quorum servers. If the TCP connection to the server breaks, the client connects to a different server in the quorum.
ZooKeeper Port	Port number of the machine that hosts the ZooKeeper server.
Enable Kerberos Connection	Enables the Informatica domain to communicate with the HBase master server or region server that uses Kerberos authentication.

Property	Description
HBase Master Principal	<p>Service Principal Name (SPN) of the HBase master server. Enables the ZooKeeper server to communicate with an HBase master server that uses Kerberos authentication.</p> <p>Enter a string in the following format:</p> <pre>hbase/<domain.name>@<YOUR-REALM></pre> <p>Where:</p> <ul style="list-style-type: none"> - domain.name is the domain name of the machine that hosts the HBase master server. - YOUR-REALM is the Kerberos realm.
HBase Region Server Principal	<p>Service Principal Name (SPN) of the HBase region server. Enables the ZooKeeper server to communicate with an HBase region server that uses Kerberos authentication.</p> <p>Enter a string in the following format:</p> <pre>hbase_rs/<domain.name>@<YOUR-REALM></pre> <p>Where:</p> <ul style="list-style-type: none"> - domain.name is the domain name of the machine that hosts the HBase master server. - YOUR-REALM is the Kerberos realm.

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> <pre>~ ` ! \$ % ^ & * () - + = { []] \ : ; " ' < , > . ? /</pre>
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>

Property	Description
Type	The connection type. Select Hive.
Connection Modes	<p>Hive connection mode. Select at least one of the following options:</p> <ul style="list-style-type: none"> - Access Hive as a source or target. Select this option if you want to use Hive as a source or a target. - Use Hive to run mappings in Hadoop cluster. Select this option if you want to use the Hive driver to run mappings in the Hadoop cluster.

Property	Description
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.
Common Attributes to Both the Modes: 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 in the 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. The path must be the fully qualified path to the JAR files used for user-defined functions. Set the parameter hive.aux.jars.path with all the entries in infapdo.aux.jars.path and the path to the JAR files for user-defined functions. - 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. The following sample text shows two values that can be used for the Environment SQL: <pre>set hive.execution.engine='tez';set hive.exec.dynamic.partition.mode='nonstrict';</pre> <p>If you use the Hive connection to run profiles in the Hadoop cluster, the Data Integration service executes only the environment SQL of the Hive connection. If the Hive sources and targets are on different clusters, the Data Integration Service does not execute the different environment SQL commands for the connections of the Hive source or target.</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
Metadata Connection String	<p>The JDBC connection URI used to access the metadata from the Hadoop server.</p> <p>You can use PowerExchange for Hive to communicate with a HiveServer service or HiveServer2 service.</p> <p>To connect to HiveServer, specify the 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 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>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>Big Data Management Installation and Configuration 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 in the 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>
Observe Fine Grained SQL Authorization	<p>When you select the option to observe fine-grained SQL authentication in a Hive source, the mapping observes row and column-level restrictions on data access. If you do not select the option, the Blaze run-time engine ignores the restrictions, and results include restricted data.</p>
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>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>Big Data Management Installation and Configuration Guide</i>.</p>

Properties to Run Mappings in Hadoop Cluster

The following table describes the Hive connection properties that you configure when you want to use the Hive connection to run Informatica mappings in the Hadoop cluster:

Property	Description
Database Name	Namespace for tables. Use the name <code>default</code> for tables that do not have a specified database name.
Default FS URI	<p>The URI to access the default Hadoop Distributed File System.</p> <p>Use the following connection URI:</p> <pre>hdfs://<node name>:<port></pre> <p>Where</p> <ul style="list-style-type: none">- <code><node name></code> is the host name or IP address of the NameNode.- <code><port></code> is the port on which the NameNode listens for remote procedure calls (RPC). <p>If the Hadoop cluster runs MapR, use the following URI to access the MapR File system: <code>maprfs:///</code>.</p>
JobTracker/Yarn Resource Manager URI	<p>The service within Hadoop that submits the MapReduce tasks to specific nodes in the cluster.</p> <p>Use the following format:</p> <pre><hostname>:<port></pre> <p>Where</p> <ul style="list-style-type: none">- <code><hostname></code> is the host name or IP address of the JobTracker or Yarn resource manager.- <code><port></code> is the port on which the JobTracker or Yarn resource manager listens for remote procedure calls (RPC). <p>If the cluster uses MapR with YARN, use the value specified in the <code>yarn.resourcemanager.address</code> property in <code>yarn-site.xml</code>. You can find <code>yarn-site.xml</code> in the following directory on the NameNode of the cluster: <code>/opt/mapr/hadoop/hadoop-2.5.1/etc/hadoop</code>.</p> <p>MapR with MapReduce 1 supports a highly available JobTracker. If you are using MapR distribution, define the JobTracker URI in the following format: <code>maprfs:///</code></p>
Hive Warehouse Directory on HDFS	<p>The absolute HDFS file path of the default database for the warehouse that is local to the cluster. For example, the following file path specifies a local warehouse:</p> <pre>/user/hive/warehouse</pre> <p>For Cloudera CDH, if the Metastore Execution Mode is remote, then the file path must match the file path specified by the Hive Metastore Service on the Hadoop cluster.</p> <p>For MapR, use the value specified for the <code>hive.metastore.warehouse.dir</code> property in <code>hive-site.xml</code>. You can find <code>hive-site.xml</code> in the following directory on the node that runs HiveServer2: <code>/opt/mapr/hive/hive-0.13/conf</code>.</p>

Property	Description
Advanced Hive/Hadoop Properties	<p>Configures or overrides Hive or Hadoop cluster properties in hive-site.xml on the machine on which the Data Integration Service runs. You can specify multiple properties.</p> <p>Select Edit to specify the name and value for the property. The property appears in the following format:</p> <pre><property1>=<value></pre> <p>Where</p> <ul style="list-style-type: none"> - <property1> is a Hive or Hadoop property in hive-site.xml. - <value> is the value of the Hive or Hadoop property. <p>When you specify multiple properties, &: appears as the property separator.</p> <p>The maximum length for the format is 1 MB.</p> <p>If you enter a required property for a Hive connection, it overrides the property that you configure in the Advanced Hive/Hadoop Properties.</p> <p>The Data Integration Service adds or sets these properties for each map-reduce job. You can verify these properties in the JobConf of each mapper and reducer job. Access the JobConf of each job from the Jobtracker URL under each map-reduce job.</p> <p>The Data Integration Service writes messages for these properties to the Data Integration Service logs. The Data Integration Service must have the log tracing level set to log each row or have the log tracing level set to verbose initialization tracing.</p> <p>For example, specify the following properties to control and limit the number of reducers to run a mapping job:</p> <pre>mapred.reduce.tasks=2&hive.exec.reducers.max=10</pre>
Temporary Table Compression Codec	Hadoop compression library for a compression codec class name.
Codec Class Name	Codec class name that enables data compression and improves performance on temporary staging tables.
Metastore Execution Mode	Controls whether to connect to a remote metastore or a local metastore. By default, local is selected. For a local metastore, you must specify the Metastore Database URI, Driver, Username, and Password. For a remote metastore, you must specify only the Remote Metastore URI.
Metastore Database URI	<p>The JDBC connection URI used to access the data store in a local metastore setup. Use the following connection URI:</p> <pre>jdbc:<datastore type>://<node name>:<port>/<database name></pre> <p>where</p> <ul style="list-style-type: none"> - <node name> is the host name or IP address of the data store. - <data store type> is the type of the data store. - <port> is the port on which the data store listens for remote procedure calls (RPC). - <database name> is the name of the database. <p>For example, the following URI specifies a local metastore that uses MySQL as a data store:</p> <pre>jdbc:mysql://hostname23:3306/metastore</pre> <p>For MapR, use the value specified for the <code>javax.jdo.option.ConnectionURL</code> property in <code>hive-site.xml</code>. You can find <code>hive-site.xml</code> in the following directory on the node where HiveServer 2 runs: <code>/opt/mapr/hive/hive-0.13/conf</code>.</p>

Property	Description
Metastore Database Driver	<p>Driver class name for the JDBC data store. For example, the following class name specifies a MySQL driver:</p> <pre>com.mysql.jdbc.Driver</pre> <p>For MapR, use the value specified for the <code>javax.jdo.option.ConnectionDriverName</code> property in <code>hive-site.xml</code>. You can find <code>hive-site.xml</code> in the following directory on the node where HiveServer 2 runs: <code>/opt/mapr/hive/hive-0.13/conf</code>.</p>
Metastore Database Username	<p>The metastore database user name.</p> <p>For MapR, use the value specified for the <code>javax.jdo.option.ConnectionUserName</code> property in <code>hive-site.xml</code>. You can find <code>hive-site.xml</code> in the following directory on the node where HiveServer 2 runs: <code>/opt/mapr/hive/hive-0.13/conf</code>.</p>
Metastore Database Password	<p>The password for the metastore user name.</p> <p>For MapR, use the value specified for the <code>javax.jdo.option.ConnectionPassword</code> property in <code>hive-site.xml</code>. You can find <code>hive-site.xml</code> in the following directory on the node where HiveServer 2 runs: <code>/opt/mapr/hive/hive-0.13/conf</code>.</p>
Remote Metastore URI	<p>The metastore URI used to access metadata in a remote metastore setup. For a remote metastore, you must specify the Thrift server details.</p> <p>Use the following connection URI:</p> <pre>thrift://<hostname>:<port></pre> <p>Where</p> <ul style="list-style-type: none"> - <code><hostname></code> is name or IP address of the Thrift metastore server. - <code><port></code> is the port on which the Thrift server is listening. <p>For MapR, use the value specified for the <code>hive.metastore.uris</code> property in <code>hive-site.xml</code>. You can find <code>hive-site.xml</code> in the following directory on the node where HiveServer 2 runs: <code>/opt/mapr/hive/hive-0.13/conf</code>.</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	<p>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:</p> <pre>~ ` ! \$ % ^ & * () - + = { [] } \ : ; " ' < , > . ? /</pre>

Property	Description
ID	String that the Data Integration Service uses to identify the connection. The ID is not case sensitive. It must be 255 characters or less and must be unique in the domain. You cannot change this property after you create the connection. Default value is the connection name.
Description	The description of the connection. The description cannot exceed 765 characters.
User Name	The database user name.
Password	The password for the database user name.
JDBC Driver Class Name	<p>Name of the JDBC driver class.</p> <p>The following list provides the driver class name that you can enter for the applicable database type:</p> <ul style="list-style-type: none"> - DataDirect JDBC driver class name for Oracle: <code>com.informatica.jdbc.oracle.OracleDriver</code> - DataDirect JDBC driver class name for IBM DB2: <code>com.informatica.jdbc.db2.DB2Driver</code> - DataDirect JDBC driver class name for Microsoft SQL Server: <code>com.informatica.jdbc.sqlserver.SQLServerDriver</code> - DataDirect JDBC driver class name for Sybase ASE: <code>com.informatica.jdbc.sybase.SybaseDriver</code> - DataDirect JDBC driver class name for Informix: <code>com.informatica.jdbc.informix.InformixDriver</code> - DataDirect JDBC driver class name for MySQL: <code>com.informatica.jdbc.mysql.MySQLDriver</code> <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> <p><code>jdbc:<subprotocol>:<subname></code></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>
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>

Property	Description
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 read data from or write data to Teradata through Teradata Connector for Hadoop (TDCH) specialized connectors for Sqoop, 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 the 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 the 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>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>`

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 read data from or write data to Teradata through Teradata Connector for Hadoop (TDCH) specialized connectors for Sqoop, you can 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 the Cloudera Connector Powered by Teradata, configure the `-Dsqoop.connection.factories` argument as follows:
`-Dsqoop.connection.factories=com.cloudera.connector.teradata.TeradataManagerFactory`
- To use the 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`

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 Hive connection mode and specify the commands for environment SQL. The SQL commands apply to both the connection modes. Select at least one of the following connection modes:

Option	Description
Access Hive as a source or target	Select this option if you want to use Hive as a source or a target.
Access Hive to run mappings in Hadoop cluster	Select this option if you want to use the Hive driver to run mappings in the Hadoop cluster.

9. Click **Test Connection** to verify the connection.

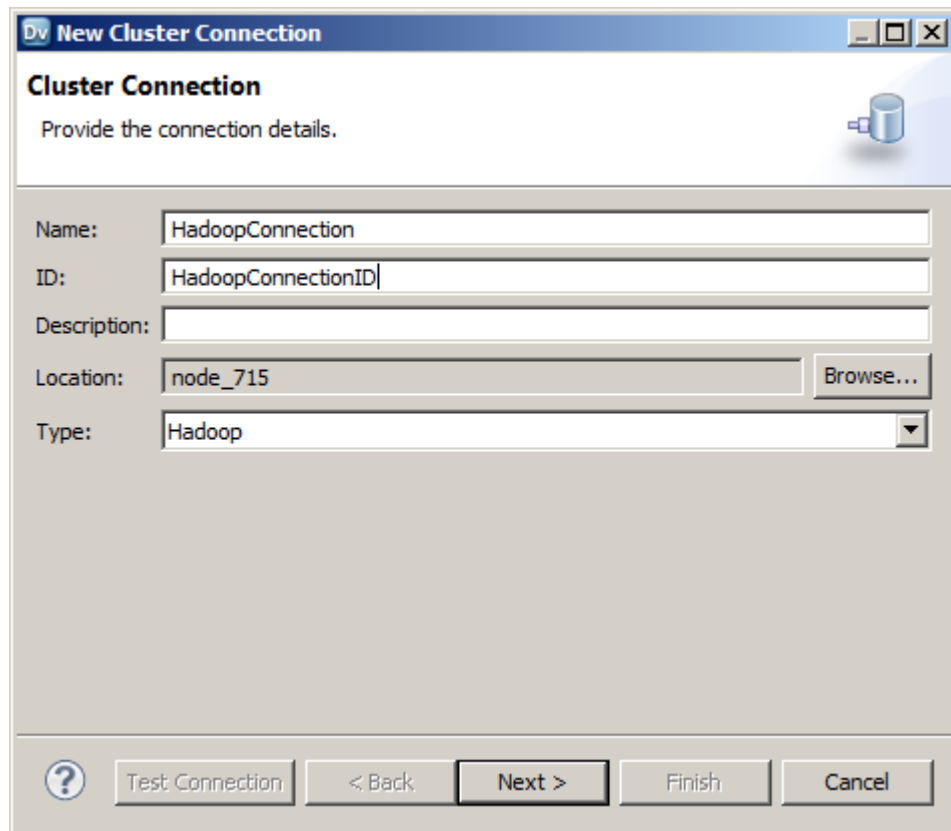
You can test a Hive connection that is configured to access Hive data. You cannot test a Hive connection that is configured to run Informatica mappings in the Hadoop cluster.

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.

The image shows a screenshot of the 'New Cluster Connection' dialog box in Informatica. The dialog has a title bar with 'Dv New Cluster Connection' and standard window controls. Below the title bar, the text 'Cluster Connection' is displayed, followed by the instruction 'Provide the connection details.' and a small database icon. The main area contains several input fields: 'Name:' with the value 'HadoopConnection', 'ID:' with the value 'HadoopConnectionID', 'Description:' (empty), 'Location:' with the value 'node_715' and a 'Browse...' button to its right, and 'Type:' with a dropdown menu showing 'Hadoop'. At the bottom, there is a row of buttons: a help button (question mark icon), 'Test Connection', '< Back', 'Next >' (which is highlighted with a black border), 'Finish', and 'Cancel'.

6. Click **Next**.
7. Enter the Hadoop cluster properties and the common properties for the Hadoop connection.
8. Click **Next**.
9. Enter the Hive pushdown configuration properties and the Hive configuration.
10. Click **Next**.
11. If you are using the Blaze engine, enter the properties for the Blaze engine.
12. If you are using the Spark engine, enter the properties for the Spark engine.

13. Click **Finish**.

CHAPTER 3

Mappings in the Hadoop Environment

This chapter includes the following topics:

- [Mappings in the Hadoop Environment Overview, 48](#)
- [Mapping Run-time Properties, 49](#)
- [Data Warehouse Optimization Mapping Example , 51](#)
- [Sqoop Mappings in a Hadoop Environment, 53](#)
- [Rules and Guidelines for Mappings in a Hadoop Environment, 56](#)
- [Workflows that Run Mappings in a Hadoop Environment, 56](#)
- [Configuring a Mapping to Run in a Hadoop Environment, 57](#)
- [Mapping Execution Plans, 57](#)
- [Optimization for the Hadoop Environment, 60](#)
- [Troubleshooting a Mapping in a Hadoop Environment, 65](#)

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.

You can select which engine the Data Integration Service uses. Informatica recommends that you select all engines. When you select more than one engine, the Data Integration Service determines the best engine to run the mapping during validation.

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:

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

Name	Value
Connection	spark1
Execution Parameters	
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 Inter...	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 Hive connection. The Hadoop cluster processes the data. Select the Hive on MapReduce engine, the Blaze engine, or the Spark engine to process the mapping. The Hadoop connection must contain the configuration properties for each engine that you choose. If you choose Hive on MapReduce engine, you can also select the Hive version. Select a version number from

the list or assign a parameter to the Hive version. The parameter must be a string that contains a version from the Hive version list. If you use the Blaze engine, you cannot clear the **Hive on MapReduce** engine.

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 Hive.

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 engine, the Spark engine, and the Hive on MapReduce engine:

Description	Object	Location
<ul style="list-style-type: none"> [Hive Validation] Hadoop connection information is missing. [Spark Validation] Hadoop connection information is missing. [Spark Validation] Relational sources are not supported in Spark execution mode. The connection for data object [EMP] cannot be empty. [Blaze engine validation] The Hadoop connection properties are not configured. [Hive Validation] Hive target [Write_a_tgt] is not configured to truncate target tables. 	<ul style="list-style-type: none"> m_blazeSparkValidation m_blazeSparkValidation Read_EMP Read_EMP Write_a_tgt 	<ul style="list-style-type: none"> MRS_ML304_hbitmrs1/engine MRS_ML304_hbitmrs1/engine MRS_ML304_hbitmrs1/engine/m_blazeSparkValidation MRS_ML304_hbitmrs1/engine/m_blazeSparkValidation MRS_ML304_hbitmrs1/engine/m_blazeSparkValidation 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.
Execution Parameters	Overrides the Hadoop custom properties or the Spark default configuration parameters for the mapping. An execution parameter in the mapping properties overrides the same execution parameter that you enter in the Hadoop connection.

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 pushes the full transformation logic to the source database.
Pushdown Capability	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.

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.

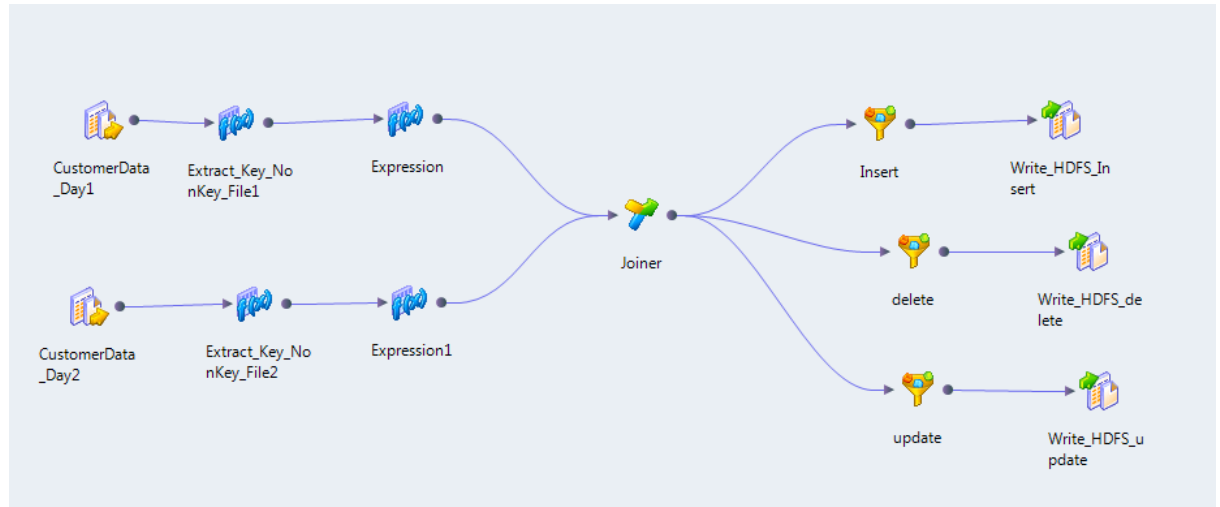
Mapping_Day2

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:



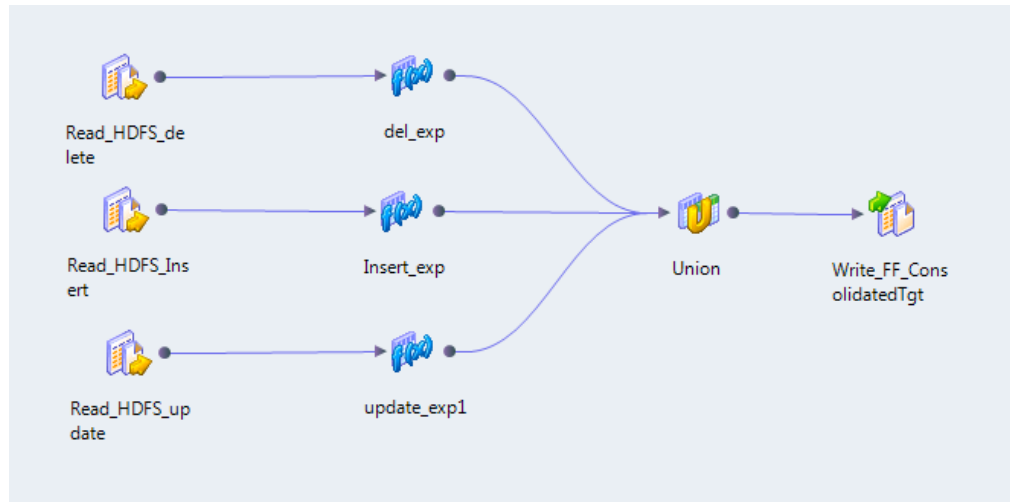
The mapping contains the following objects:

- 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:

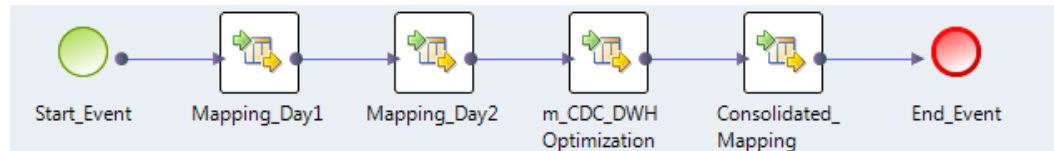


The mapping contains the following objects:

- 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 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.

Note: You can run Sqoop mappings on the Spark engine only when you want to read data from or write data to Oracle databases. If you use the Cloudera Connector Powered by Teradata or Hortonworks Connector for Teradata, you must run the mappings on the Blaze engine.

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

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`

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 engine, Sqoop ignores the argument. Sqoop creates map tasks based on the number of intermediate files that the Blaze 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 the 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.

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.
- Mappings that contain a Hive source or a Hive target must use the same Hive connection to push the mapping to Hadoop.
- The Data Integration Service ignores the data file block size configured for HDFS files in the `hdfs-site.xml` file. The Data Integration Service uses a default data file block size of 64 MB for HDFS files. To change the data file block size, copy `/usr/lib/hadoop/conf/hdfs-site.xml` to the following location in the Hadoop distribution directory for the Data Integration Service node: `/opt/Informatica/services/shared/hadoop/[Hadoop_distribution_name]/conf`. You can also update the data file block size in the following file: `/opt/Informatica/services/shared/hadoop/[Hadoop_distribution_name]/conf/hive-default.xml`.

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).

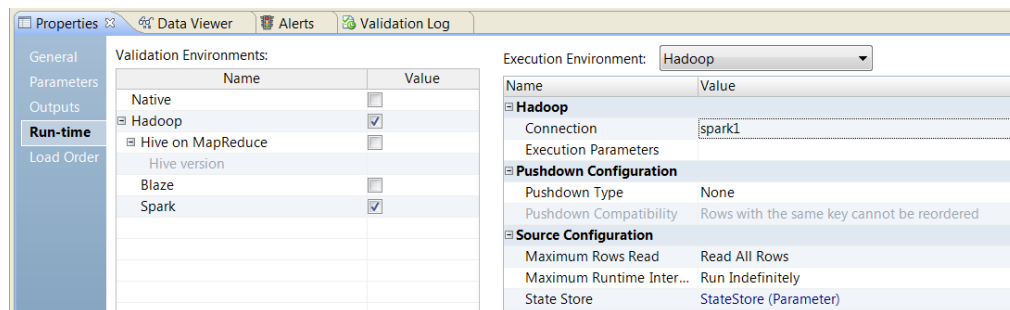
Configuring a Mapping to Run in a Hadoop Environment

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

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.

The Hive on MapReduce, the Blaze, and the Spark engines are selected by default. To only use the Hive on MapReduce engine, clear the other engines. If you use the Blaze engine, you cannot clear the Hive on MapReduce engine.

4. In the execution environment, select **Hadoop**.



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. Optionally, select **Execution Parameters** to override a Hadoop custom property or a Spark default configuration parameter.
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

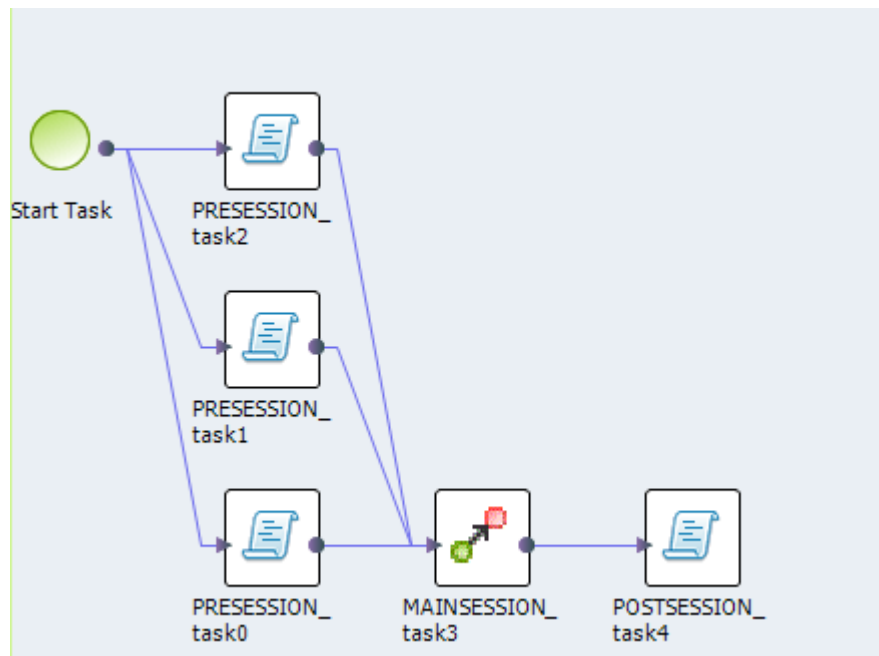
The Hive execution plan is a series of Hive queries. The plan contains tasks to start the mapping, run the mapping, and clean up the temporary tables and files. You can view the Hive execution plan that the Data Integration Service generates before you run the mapping. When the Data Integration Service 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_OUT[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:

Script Name	Script	Depends On
cominformaticaexe cInfaSprk0	<pre>package com.informatica.exec import com.informatica.bootstrap.functions._ import com.informatica.bootstrap._ import org.apache.spark.SparkContext import org.apache.spark.rdd._ import org.apache.spark.sql._ import org.apache.spark.sql.types._ import org.apache.spark.sql.functions._ import org.apache.spark._ import java.io._ import com.databricks.spark.avro._ import org.apache.spark.sql.hive._ object InfaSprk0 { def main(args: Array[String]) { val sc = new SparkContext();</pre>	Pre_Spark_Task_ Command_0

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.

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. To view the execution plan in the Developer tool, select the **Data Viewer** view for the mapping and click **Show Execution Plan**.
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

You can optimize the Hadoop environment and the Hadoop cluster to increase performance.

You can optimize the Hadoop environment and the Hadoop cluster 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 Installation and Configuration 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 with some restrictions.

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

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 not apply a recovery strategy.
- Full restart. Restarts the grid task.

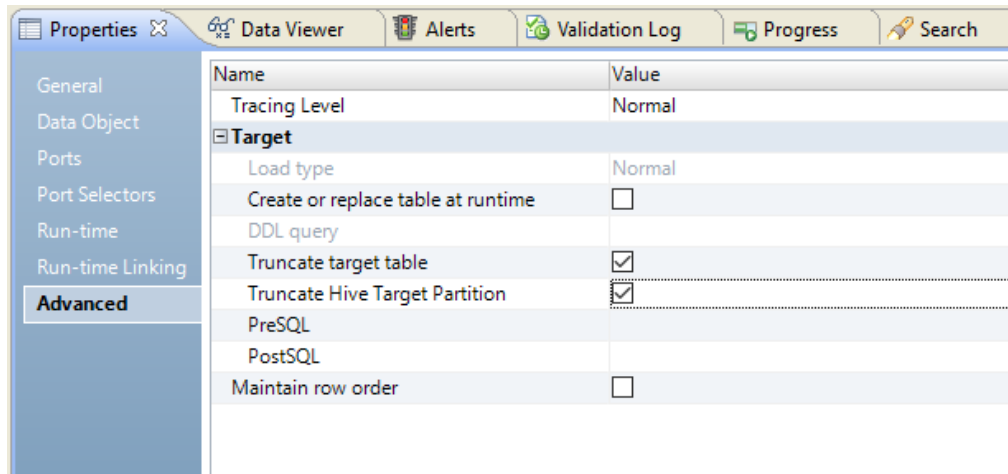
Truncating Partitions in a Hive Target

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**.

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.
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.

Scheduling, Queuing, and Node Labeling

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

A scheduler assigns resources on the cluster to applications that need them, while honoring organizational policies on sharing resources. You can configure YARN to use the capacity scheduler or the fair scheduler. The capacity scheduler allows multiple organizations to share a large cluster and distributes resources based on capacity allocations. The fair scheduler shares resources evenly among all jobs running on the cluster.

Queues are the organizing structure for YARN schedulers, allowing multiple tenants to share the cluster. The capacity of each queue specifies the percentage of cluster resources that are available for applications submitted to the queue. You can direct the Blaze and Spark engines to a YARN scheduler queue.

You can use node labels to run YARN applications on cluster nodes. Node labels partition a cluster into sub-clusters so that jobs can run on nodes with specific characteristics. You can then associate node labels with capacity scheduler queues.

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.

Scheduling and Node Labeling Configuration

Update the `yarn-site.xml` file on the domain environment to enable scheduling and node labeling in the Hadoop 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.CapacitySche
duler></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>
```

Queuing Configuration

The following table describes the Hadoop connection properties you configure to direct a job to a specific YARN scheduler queue:

Property	Description
YARN Queue Name	The Blaze engine connection property to specify a YARN scheduler queue name. The name is case sensitive.
Spark Execution Parameter	The Spark engine connection property to specify a YARN scheduler queue name. Use the following format: <code>spark.yarn.queue=<YARN queue name></code>

The following table describes the Hive connection property you configure to direct a job to a specific YARN scheduler queue:

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@INFAKRB?mapred.job.queue.name=root.test</code></p>

The following table describes the JDBC connection property you configure to direct a job to a specific YARN scheduler queue:

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>

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.

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 MapR 2.1.2 distribution that processes large amounts of data, monitoring the mapping from the Administrator tool stops.

You can check the Hadoop task tracker log to see if there a timeout that results in Hadoop job tracker and Hadoop task tracker losing connection. To continuously monitor the mapping from the Administrator tool, increase the virtual memory to 640 MB in the `hadoopEnv.properties` file. The default is 512 MB. For example, `infapdo.java.opts=-Xmx640M -XX:GCTimeRatio=34 -XX:`

```
+UseConcMarkSweepGC -XX:+UseParNewGC -XX:ParallelGCThreads=2 -XX:NewRatio=2 -  
Djava.library.path=$HADOOP_NODE_INFA_HOME/services/shared/bin:  
$HADOOP_NODE_HADOOP_DIST/lib/native/Linux-amd64-64 -Djava.security.egd=file:/dev/./  
urandom -Dmapr.library.flatclass
```

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%.

CHAPTER 4

Mapping Objects in the Hadoop Environment

This chapter includes the following topics:

- [Sources in a Hadoop Environment, 67](#)
- [Targets in a Hadoop Environment, 73](#)
- [Transformations in a Hadoop Environment, 77](#)
- [Function and Data Type Processing, 86](#)

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:

- Flat file (native)
- HBase
- HDFS complex file
- HDFS flat file
- Hive
- IBM DB2
- Netezza
- ODBC
- Oracle
- Sqoop sources
- Teradata

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.

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:

Pushdown Processing 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 High-Availability cluster: hdfs://<host name>/<file name path>

For example, hdfs://irl1dv: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:

- The Data Integration Service can run pre-mapping SQL commands against the source database before it reads from a Hive source. When you create a 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.
- A mapping fails to validate when you configure post-mapping SQL commands. The Data Integration Service does not run post-mapping SQL commands against a Hive source.
- 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.

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, RC, 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

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');
```

Complex File Sources

A mapping that is running 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 reader 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 reader object to a Data Processor transformation.

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*	Not supported	Supported*
JSON	Flat	Supported*	Not supported	Supported*
JSON	Hierarchical	Supported*	Not 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*	Not supported	Supported*
XML	Flat	Supported*	Not supported	Supported*
XML	Hierarchical	Supported*	Not supported	Supported*
* The complex file reader object must be connected to a Data Processor transformation.				

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 create a password file to access a database, to run the mapping successfully, you must ensure that the password file exists on HDFS.
- 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.
- 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 enable OraOop and configure an advanced query to read data from an Oracle source through Sqoop, the mapping fails on the Spark engine.
- 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.
- If you configure the --username or --password argument in a JDBC connection or mapping, Sqoop ignores the arguments.

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 but it cannot have 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.

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.

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

- The Data Integration Service does not run pre-mapping or post-mapping SQL commands against a Hive target. You cannot validate and run a mapping with PreSQL or PostSQL properties for a Hive target.
- A mapping fails to run if the Hive target definition differs in the number and order of the columns from the relational table in the Hive database.
- A mapping fails to run when you use Unicode characters in a Hive target definition.
- 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.
- 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.

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.

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

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.
- The Blaze engine supports Hive tables that are enabled for locking.
- The Blaze engine can create or replace Hive target tables and truncate the partition in the Hive target table. You must choose to both truncate the partition in the Hive target and truncate the target table.
- 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 an RCFile as a Hive target. To write to an RCFile table, you must create the table with the `serDe` clause.

- 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.

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

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');
```

Complex File Targets

A mapping that is running 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*	Not supported	Supported*
JSON	Flat	Supported*	Not supported	Supported*
JSON	Hierarchical	Supported*	Not 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*	Not 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.				

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 you create a password file to access a Sqoop target, to run the mapping successfully, you must ensure that the password file exists on HDFS.
- 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.
- 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.
- If you configure the --username or --password argument in a JDBC connection or mapping, Sqoop ignores the arguments.

Transformations in a Hadoop Environment

Due to the differences between native environment and Hadoop environment only certain transformations are valid or valid with restrictions in the Hadoop environment. The Data Integration Service does not process

transformations that contain functions, expressions, data types, and variable fields that are not valid in a Hadoop environment.

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

Transformation	Blaze Engine	Spark Engine	Hive Engine
Address Validator	Supported with restrictions	Not supported	Supported with restrictions
Aggregator	Supported with restrictions	Supported with restrictions	Supported with restrictions
Case Converter	Supported	Not supported	Supported
Comparison	Supported	Not supported	Supported
Consolidation	Supported with restrictions	Not supported	Supported with restrictions
Data Masking	Supported with restrictions	Not supported	Supported with restrictions
Data Processor	Supported with restrictions	Not supported	Supported with restrictions
Decision	Supported	Not supported	Supported
Expression	Supported with restrictions	Supported with restrictions	Supported with restrictions
Filter	Supported	Supported	Supported
Java	Supported	Supported with restrictions	Supported with restrictions
Joiner	Supported with restrictions	Supported with restrictions	Supported with restrictions
Key Generator	Supported	Not supported	Not supported
Labeler	Supported	Not supported	Supported
Lookup	Supported with restrictions	Supported with restrictions	Supported with restrictions
Match	Supported with restrictions	Not supported	Supported with restrictions
Merge	Supported	Not supported	Supported
Normalizer	Supported	Not supported	Supported
Parser	Supported	Not supported	Supported
Rank	Supported	Not supported	Supported with restrictions
Router	Supported	Supported	Supported
Sequence Generator	Supported with restrictions	Not supported	Not supported

Transformation	Blaze Engine	Spark Engine	Hive Engine
Sorter	Supported	Supported with restrictions	Supported with restrictions
SQL	Not supported	Not supported	Supported with restrictions
Standardizer	Supported	Not supported	Supported
Union	Supported	Supported	Supported
Update Strategy	Supported with restrictions	Not supported	Supported with restrictions
Weighted Average	Supported	Not supported	Supported

Transformation Support on the Blaze Engine

Some restrictions and guidelines apply to processing transformations on the Blaze engine.

The following table describes rules and guidelines for processing transformations on the Blaze engine:

Transformation	Rules and Guidelines
Address Validator	The Address Validator transformation cannot generate a certification report.
Aggregator	Mapping validation fails in the following situations: <ul style="list-style-type: none"> - The transformation contains stateful variable ports. - The transformation contains unsupported functions in an expression.
Case Converter	Supported without restrictions.
Comparison	Supported without restrictions.
Data Processor	Mapping validation fails in the following situations. <ul style="list-style-type: none"> - The transformation Data processor mode is set to Input Mapping or Service and Input Mapping.
Decision	Supported without restrictions.
Expression	Mapping validation fails in the following situations: <ul style="list-style-type: none"> - 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	Supported without restrictions. 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	Supported without restrictions.
Joiner	Mapping validation fails in the following situations: <ul style="list-style-type: none"> - The transformation contains an inequality join.

Transformation	Rules and Guidelines
Key Generator	Supported without restrictions.
Labeler	Supported without restrictions.
Lookup	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - The cache is configured to be shared, named, persistent, dynamic, or uncached. The cache must be a static cache. <p>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.</p>
Match	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - The mapping specifies an identity match type. <p>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.</p>
Merge	Supported without restrictions.
Normalizer	Supported without restrictions.
Parser	Supported without restrictions.
Rank	Supported without restrictions.
Router	Supported without restrictions.
Sequence Generator	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - The reset property for the transformation is enabled. <p>If you want the Data Integration Service to reset the sequence data object to the start value for each mapping run, then you must run the mapping in the native environment.</p>
Sorter	<p>The Blaze engine can perform global sorts when the following conditions are true:</p> <ul style="list-style-type: none"> - 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. <p>If any of the conditions are not true, the Blaze engine performs a local sort.</p>
Standardizer	Supported without restrictions.
Union	Supported without restrictions.

Transformation	Rules and Guidelines
Update Strategy	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - The Update Strategy transformation is connected to more than one target. - The Update Strategy transformation is not connected directly to a target. - The Update Strategy transformation target is an external Hive table. - The target is not a Hive target on the same cluster as the Update Strategy transformation. - The target does not contain a primary key. - The Hive target property to truncate the target table at runtime is enabled. - The Hive target property to create or replace the target table at runtime is enabled. <p>The mapping fails in the following situations:</p> <ul style="list-style-type: none"> - The target is not ORC bucketed. - The target is ORC bucketed on all columns. <p>Compile validation errors occur and the mapping execution stops in the following situations:</p> <ul style="list-style-type: none"> - The Hive version is earlier than 0.14. - The target table is not enabled for transactions. <p>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: <code>TBLPROPERTIES ("transactional"="true")</code>.</p> <p>To use an Update Strategy transformation with a Hive target, you must configure the following properties in the <code>hive-site.xml</code> file:</p> <pre>hive.support.concurrency true hive.enforce.bucketing true hive.exec.dynamic.partition.mode true hive.txn.manager org.apache.hadoop.hive.ql.lockmgr.DbTxnManager hive.compactor.initiator.on true hive.compactor.worker.threads 1</pre> <p>Use the following format to configure the properties in the <code>hive-site.xml</code> file:</p> <pre><property> <name>property-name</name> <value>property-value</value> </property></pre> <p>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.</p> <p>The Blaze engine executes transactions in the following order: deletes, updates, inserts. It does not process rows in the same order as the Update Strategy transformation receives them.</p> <p>The Blaze engine performs Update as Update even if the transformation is configured to Update as Insert or Update else Insert.</p>
Weighted Average	Supported without restrictions.
<i>Transformations not listed in this table are not supported.</i>	

Transformation Support on the Spark Engine

Some restrictions and guidelines apply to processing transformations on the Spark engine.

The following table describes rules and guidelines for the transformations that are supported on the Spark engine:

Transformation	Rules and Guidelines
Aggregator	Mapping validation fails in the following situations: <ul style="list-style-type: none">- The transformation contains stateful variable ports.- The transformation contains unsupported functions in an expression.
Expression	Mapping validation fails in the following situations: <ul style="list-style-type: none">- 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.
Filter	Supported without restrictions.
Java	<p>You must copy external .jar files that a Java transformation requires to the Informatica installation directory on the Hadoop cluster at the following location: <code>[\$HADOOP_NODE_INFA_HOME]/services/shared/jars</code>.</p> <p>To run user code directly on the Spark engine, the JDK version that the Data Integration Service uses must be compatible with the JRE version on the cluster. For best performance, create the environment variable <code>DIS_JDK_HOME</code> on the Data Integration Service in the Administrator tool. The environment variable contains the path to the JDK installation folder on the machine running the Data Integration Service. For example, you might enter a value such as <code>/usr/java/default</code>.</p> <p>The Partitionable property must be enabled in the Java transformation. The transformation cannot run in one partition.</p> <p>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.</p> <p>When you enable high precision and the Java transformation contains a field that is a decimal data type, a validation error occurs.</p> <p>The following restrictions apply to the Transformation Scope property:</p> <ul style="list-style-type: none">- 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. <p>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.</p>
Joiner	Mapping validation fails in the following situations: <ul style="list-style-type: none">- Case sensitivity is disabled.- The join condition in the Joiner transformation contains binary data type or binary expressions.

Transformation	Rules and Guidelines
Lookup	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - Case sensitivity is disabled. - The lookup condition in the Lookup transformation contains binary data type. - The transformation is not configured to return all rows that match the condition. - The lookup is a data object. - The cache is configured to be shared, named, persistent, dynamic, or uncached. The cache must be a static cache. <p>The mapping fails in the following situations:</p> <ul style="list-style-type: none"> - The transformation is unconnected. <p>When you use Sqoop and look up data in a Hive table based on a column of the float data type, the Lookup transformation might return incorrect results.</p>
Router	Supported without restrictions.
Sorter	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - Case sensitivity is disabled. <p>The Data Integration Service logs a warning and ignores the Sorter transformation in the following situations:</p> <ul style="list-style-type: none"> - 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. <p>The Data Integration Service treats null values as high even if you configure the transformation to treat null values as low.</p>
Union	Supported without restrictions.
<i>Transformations not listed in this table are not supported.</i>	

Transformation Support on the Hive Engine

Some restrictions and guidelines apply to processing transformations on the Hive engine.

The following table describes rules and guidelines for processing transformations on the Hive engine.

Transformation	Rules and Guidelines
Address Validator	The Address Validator transformation cannot generate a certification report.
Aggregator	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - The transformation contains stateful variable ports. - The transformation contains unsupported functions in an expression.
Case Converter	Supported without restrictions.
Comparison	Supported without restrictions.
Consolidation	<p>The Consolidation transformation might process data differently in the native environment and in a Hadoop environment.</p> <p>The transformation might demonstrate the following differences in behavior:</p> <ul style="list-style-type: none"> - 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.

Transformation	Rules and Guidelines
Data Masking	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - The transformation is configured for repeatable expression masking. - The transformation is configured for unique repeatable substitution masking.
Data Processor	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - The transformation contains more than one input port. - The transformation contains pass-through ports.
Decision	Supported without restrictions.
Expression	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - The transformation contains stateful variable ports. - The transformation contains unsupported functions in an expression. <p>An Expression transformation with a user-defined function returns a null value for rows that have an exception error in the function.</p>
Filter	Supported without restrictions.
Java	<p>You must copy external .jar files that a Java transformation requires to the Informatica installation directory in the Hadoop cluster nodes at the following location: <code>[\$HADOOP_NODE_INFA_HOME]/services/shared/jars</code></p> <p>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.</p> <p>The following restrictions apply to the Transformation Scope property:</p> <ul style="list-style-type: none"> - 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. <p>You can enable the Stateless advanced property when you run mappings in a Hadoop environment.</p> <p>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.</p>
Joiner	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - The transformation contains an inequality join.
Labeler	Supported without restrictions.
Lookup	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - 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. <p>Mappings fail in the following situations:</p> <ul style="list-style-type: none"> - The lookup is unconnected. <p>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.</p> <p>When you a run 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: <code>/tmp/<user_name>/hive_resources</code> . The Hive parameter <code>hive.downloaded.resources.dir</code> 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.</p>

Transformation	Rules and Guidelines
Match	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - The transformation specifies an identity match type. <p>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.</p>
Merge	Supported without restrictions.
Parser	Supported without restrictions.
Rank	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - Case sensitivity is disabled.
Router	Supported without restrictions.
Sorter	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - Case sensitivity is disabled. <p>The Data Integration Service logs a warning and ignores the Sorter transformation in the following situations:</p> <ul style="list-style-type: none"> - 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. <p>The Data Integration Service treats null values as high, even if you configure the transformation to treat null values as low.</p>
Standardizer	Supported without restrictions.
SQL	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - The mapping also contains a Hive source or target.
Union	Supported without restrictions.
Update Strategy	<p>Mapping validation fails in the following situations:</p> <ul style="list-style-type: none"> - The transformation is connected to more than one target. - The transformation is not connected directly to the target. <p>The mapping fails in the following situations:</p> <ul style="list-style-type: none"> - The target is not ORC bucketed. - The target is ORC bucketed on all columns. <p>Compile validation errors occur and the mapping execution stops in the following situations:</p> <ul style="list-style-type: none"> - 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. <p>The Hive engine performs Update as Update even if the transformation is configured to Update as Insert or Update else Insert.</p> <p>When the Update Strategy transformation receives multiple update rows for the same key, the results might differ.</p>
Weighted Average	Supported without restrictions.
<i>Transformations not listed in this table are not supported.</i>	

Function and Data Type Processing

When you run a mapping in a Hadoop environment, the engine that runs the mapping might process Informatica functions and data types differently from the Data Integration Service. Some variations apply in the processing and validity of functions and data types because of differences between the environments. As a result, 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 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 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.
- 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().

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().

CHAPTER 5

Monitoring Mappings in the Hadoop Environment

This chapter includes the following topics:

- [Monitoring Mappings in the Hadoop Environment Overview, 88](#)
- [Hadoop Environment Logs, 88](#)
- [Blaze Engine Monitoring, 92](#)
- [Spark Engine Monitoring, 100](#)
- [Hive Engine Monitoring, 103](#)

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.

The Data Integration Service also generates log events when you run a mapping in the Hadoop environment. 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.

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 Hive on MapReduce engine. 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.

You can view the Scala code 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 a Hadoop cluster. The page title is "All Applications". On the left, there is a sidebar with a "Cluster" section containing links for "About", "Nodes", "Applications", "NEW", "NEW SAVING", "SUBMITTED", "ACCEPTED", "RUNNING", "FINISHED", "FAILED", "KILLED", and "Scheduler". Below this is a "Tools" section. The main content area displays "Cluster Metrics" and "User Metrics for dr.who". The "Cluster Metrics" table shows various statistics for the cluster, including Apps Submitted, Apps Pending, Apps Running, Apps Completed, Containers Running, Memory Used, Memory Total, Memory Reserved, VCoers Used, VCoers Total, VCoers Reserved, Active Nodes, and Decommissioned Nodes. The "User Metrics for dr.who" table shows similar statistics for the user. Below these tables is a list of applications, with columns for ID, User, Name, Application Type, Queue, StartTime, FinishTime, State, FinalStatus, Running Containers, and Allocated CPU VCoers. The applications listed are all SPARK applications, all in a FINISHED state, and all succeeded.

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	VCoers Used	VCoers Total	VCoers Reserved	Active Nodes	Decommissioned Nodes
568	0	0	568	0	0 B	32 GB	0 B	0	32	0	1	0

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Containers Pending	Containers Reserved	Memory Used	Memory Pending	Memory Reserved
0	0	0	568	0	0	0	0 B	0 B	0 B

ID	User	Name	Application Type	Queue	StartTime	FinishTime	State	FinalStatus	Running Containers	Allocated CPU VCoers
application_1463379223882_0568	ldmui	InfSprk0	SPARK	root.ldmui	Mon May 16 13:11:59 -0700 2016	Mon May 16 13:13:03 -0700 2016	FINISHED	SUCCEEDED	N/A	N/A
application_1463379223882_0567	ldmui	InfSprk0	SPARK	root.ldmui	Mon May 16 13:11:58 -0700 2016	Mon May 16 13:13:01 -0700 2016	FINISHED	SUCCEEDED	N/A	N/A
application_1463379223882_0566	ldmui	InfSprk0	SPARK	root.ldmui	Mon May 16 13:11:56 -0700 2016	Mon May 16 13:13:00 -0700 2016	FINISHED	SUCCEEDED	N/A	N/A
application_1463379223882_0565	ldmui	InfSprk0	SPARK	root.ldmui	Mon May 16 13:11:55 -0700 2016	Mon May 16 13:12:59 -0700 2016	FINISHED	SUCCEEDED	N/A	N/A
application_1463379223882_0564	ldmui	InfSprk0	SPARK	root.ldmui	Mon May 16 13:11:54	Mon May 16 13:12:59	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.

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 this task 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**.
 - 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 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.

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 Blaze Monitor tab in the Administrator tool:

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	Zo9oFWbG...	Administrator	09/29/2015 09:23:21	00:03:33	09/29/2015 09:26:55
MappingSrcPa...	Mapping	Completed	zD02WGbc...	Administrator	09/29/2015 12:03:41	00:00:01	09/29/2015 12:03:42
MappingSrcPa...	Mapping	Completed	Zby26WbdE...	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 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.

Note: If you select a Hive mapping in the contents panel, a row called "AllHiveSourceTables" appears in the Summary Statistics view. This row displays the number of rows processed across all sources.

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. 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.

The following image shows the Blaze Job Monitor:

The screenshot shows the Blaze Job Monitor application interface. On the left is a navigation menu with options: Task History, Grid Tasks, Succeeded, Running, Failed, Grid Segments, Tasklets, and Attempts. The main area is titled "All Tasklet Attempts" and contains a table with columns: Name, Start Time, End Time, Elapsed Time, State, Host Name, and Log. The table lists several tasklet attempts, all of which are in a "Succeeded" state. The "Name" column contains identifiers like "grid-499-1-42938595-32_s6_j0_1". The "Start Time" and "End Time" columns show timestamps from October 31, 2016. The "Elapsed Time" column shows durations like "0:0:49". The "State" column shows "Succeeded" with a green bar icon. The "Host Name" column shows "psrhaqdn21.informatica.com". The "Log" column contains a "Log" link for each entry.

Name	Start Time	End Time	Elapsed Time	State	Host Name	Log
grid-499-1-42938595-32_s6_j0_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:49	Succeeded	psrhaqdn21.informatica.com	Log
grid-499-1-42938595-32_s6_j1_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:49	Succeeded	psrhaqdn28.informatica.com	Log
grid-499-1-42938595-32_s6_j2_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:59 PM	0:1:52	Succeeded	psrhaqdn23.informatica.com	Log
grid-499-1-42938595-32_s6_j3_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:49	Succeeded	psrhaqdn28.informatica.com	Log
grid-499-1-42938595-32_s6_j4_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:49	Succeeded	psrhaqdn26.informatica.com	Log
grid-499-1-42938595-32_s6_j5_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:49	Succeeded	psrhaqdn21.informatica.com	Log
grid-499-1-42938595-32_s6_j6_1	Mon Oct 31 2016 2:45:04 PM	Mon Oct 31 2016 2:45:47 PM	0:0:43	Succeeded	psrhaqdn21.informatica.com	Log
grid-499-1-42938595-32_s1_j1_1	Mon Oct 31 2016 2:45:04 PM	Mon Oct 31 2016 2:45:07 PM	0:0:3	Succeeded	psrhaqdn28.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. You can monitor properties for a task such as name, start time, end 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.

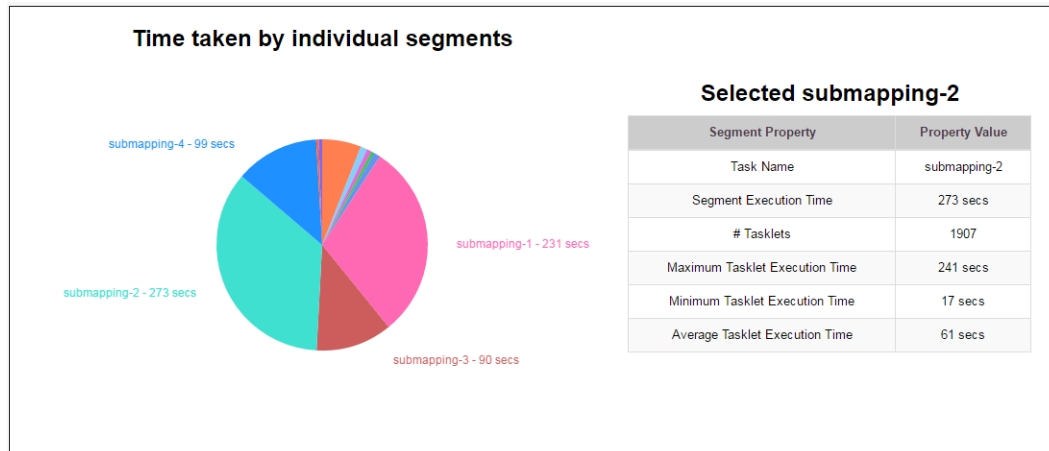
Note: The Blaze Summary Report is in beta. It contains most of the major features, but is not yet complete.

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	Summary Report 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

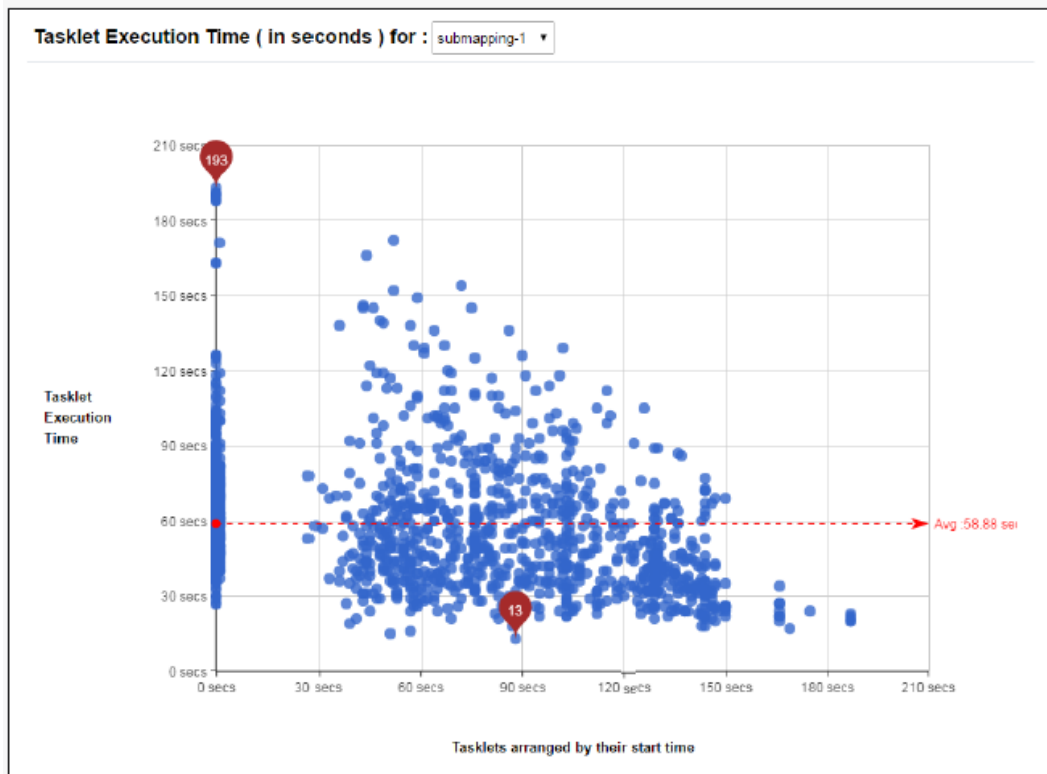
The Mapping Properties table displays the following information:

- 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 tables with red arrows pointing to specific data points:

- Source Table:** Shows sources processed. The `Read_catalog_sales` source has 4,937,484 rows processed. A red arrow points to this value with the text: "Row counts for all sources processed by tasklet".
- Target Table:** Shows targets written. The `DETarget_Aggregator1_G0` target has 0 rows processed. A red arrow points to this value with the text: "Row counts for all targets written by tasklet".
- Transformation Table:** Shows cache usage for transformations. The `Joiner1` transformation has an index cache of 14,400 bytes used and a data cache of 6,968 bytes used. A red arrow points to these values 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. 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 log file directories are specified by properties in the `hadoopEnv.properties` file located in the following location for each Hadoop distribution:

```
<Informatica Installation directory>/services/shared/hadoop/<distribution directory>/  
infaConf
```

The temporary directory is specified by the following property in the `hadoopEnv.properties` file: `infaGrid.node.local.root.log.dir`. An administrator must create a directory with read, write, and execute permissions on all nodes on the Hadoop cluster.

For example, configure the following path for the property:

```
infaGrid.node.local.root.log.dir=$HADOOP_NODE_INFA_HOME/dtmLogs
```

After the mapping completes, the Data Integration Service moves the tasklet log events from the temporary directory to a permanent directory on HDFS. The permanent directory is specified by the following property in the `hadoopEnv.properties` file: `infaCal.hadoop.logs.directory`.

For example, configure the following path for the property:

```
infaCal.hadoop.logs.directory=/var/log/hadoop-yarn/apps/informatica
```

If you want to retain the tasklet logs in the temporary directory, set the value of the following property in the `hadoopEnv.properties` file to `false`: `infaGrid.delete.local.log`

If you do not configure the temporary or permanent directories, the tasklet log events appear in the directory configured for the DTM Process. You can get the directory for the DTM Process from the value for the `yarn.nodemanager.local-dirs` property in `yarn-site.xml` on the cluster node.

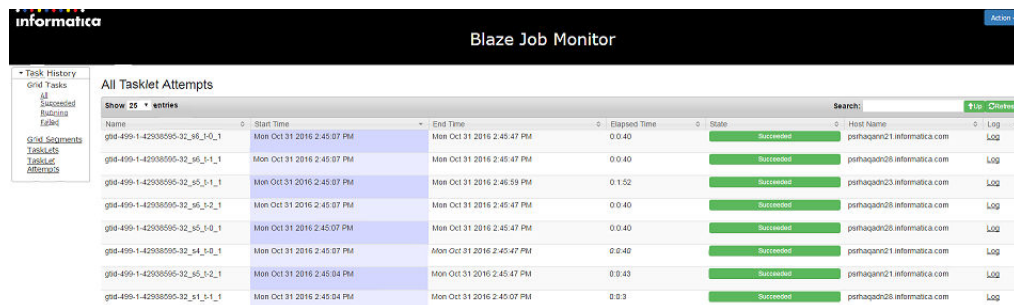
The following sample code describes the `yarn.nodemanager.local-dirs` property:

```
<property>  
<name>yarn.nodemanager.local-dirs</name>  
<value>/var/lib/hadoop-yarn/cache/${user.name}/nm-local-dir</value>  
<description>List of directories to store local files.</description>  
</property>
```

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 pane with links for 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. The 'Log' column contains links to view the logs for each tasklet attempt.

Name	Start Time	End Time	Elapsed Time	State	Host Name	Log
gnd-499-1-42908595-32_36_3-0_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	psrhagand21.informatica.com	Log
gnd-499-1-42908595-32_36_3-1-1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:07 PM	0:0:40	Succeeded	psrhagand26.informatica.com	Log
gnd-499-1-42908595-32_36_3-1-1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:46:59 PM	0:1:52	Succeeded	psrhagand23.informatica.com	Log
gnd-499-1-42908595-32_36_3-2-1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	psrhagand26.informatica.com	Log
gnd-499-1-42908595-32_36_3-0_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	psrhagand26.informatica.com	Log
gnd-499-1-42908595-32_36_3-0_1	Mon Oct 31 2016 2:45:07 PM	Mon Oct 31 2016 2:45:47 PM	0:0:40	Succeeded	psrhagand21.informatica.com	Log
gnd-499-1-42908595-32_36_3-2-1	Mon Oct 31 2016 2:45:04 PM	Mon Oct 31 2016 2:45:47 PM	0:0:43	Succeeded	psrhagand21.informatica.com	Log
gnd-499-1-42908595-32_36_3-1-1	Mon Oct 31 2016 2:45:04 PM	Mon Oct 31 2016 2:45:07 PM	0:0:3	Succeeded	psrhagand26.informatica.com	Log

The log events appear in another browser window.

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 Monitoring URL is not displayed in the Properties view of the Administrator tool.

Locate the URL in the YARN log.

When the connection to the Application Timeline Server is lost, the Blaze engine continues to attempt to reconnect to it.

Restart the Application Timeline Server from the cluster management interface.

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 Spark Monitor tab in the Administrator tool:

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 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.

Note: If you select a Hive mapping in the contents panel, a row called "AllHiveSourceTables" appears in the Summary Statistics view. This row displays the number of rows processed across all sources.

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_HDP25', 'Ad Hoc Jobs', and 'DIS_TE'. The 'Ad Hoc Jobs' folder is expanded, showing a list of jobs. The job 'Map_Spark' is selected. The main panel displays the 'Summary Statistics' view for this job. The view shows a table with columns: Source, Rows, Average Rows/Sec, Bytes, Average Bytes/Sec, First Row Accessed, and Dropped Rows. The data is organized into two sections: 'Source' and 'Target'. The 'Source' section shows 'Read_Ratp' with 28967753 rows. The 'Target' section shows 'Write_Ratp' with 28967753 rows. The 'Summary Statistics' tab is highlighted with a red circle in the original image.

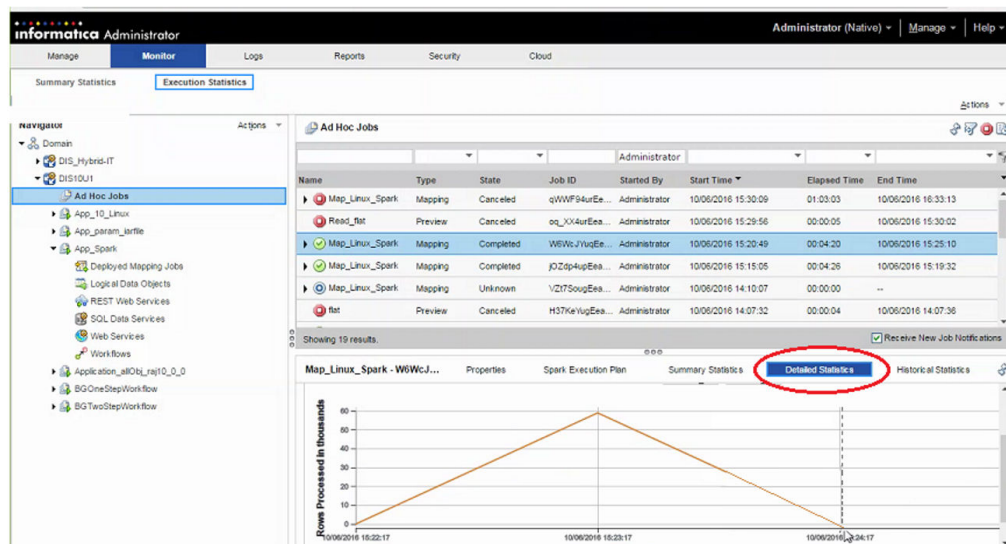
Source	Rows	Average Rows/Sec	Bytes	Average Bytes/Sec	First Row Accessed	Dropped Rows
Read_Ratp	28967753	N/A	N/A	N/A	N/A	-1

Target	Rows	Average Rows/Sec	Bytes	Average Bytes/Sec	Rejected Rows
Write_Ratp	28967753	N/A	N/A	N/A	N/A

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:



Spark Engine Logs

The Spark engine logs appear in the LDTM log.

The LDTM logs the results of the Spark engines 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.

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.

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 Hive Monitor tab in the Administrator tool:

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 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 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.

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 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.

Note: If you select a Hive mapping in the contents panel, a row called "AllHiveSourceTables" appears in the Summary Statistics view. This row displays the number of rows processed across all sources.

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.

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

For every Hive script in the Hive execution plan for 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: `<InformaticaInstallationDir>/tomcat/bin/disTemp/`. The full path to the Hive session log appears in the LDTM log.

CHAPTER 6

Mappings in the Native Environment

This chapter includes the following topics:

- [Mappings in the Native Environment Overview, 105](#)
- [Data Processor Mappings, 105](#)
- [HDFS Mappings, 106](#)
- [Hive Mappings, 107](#)
- [Social Media Mappings, 108](#)

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, HTML pages, XML, 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.

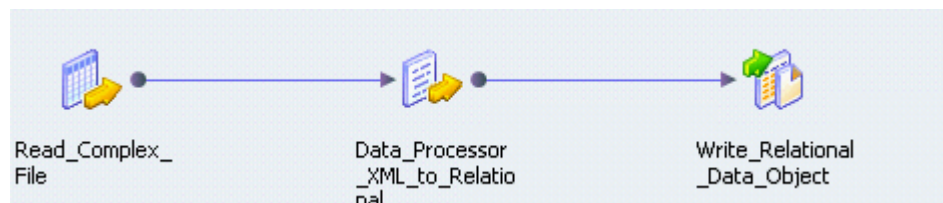
Validate and run the mapping. You can deploy the mapping and run it or add the mapping to a Mapping task in a workflow.

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, Read_Complex_File, is 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, Write_Relational_Data_Object, is 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 7

Profiles

This chapter includes the following topics:

- [Profiles Overview, 110](#)
- [Native Environment, 110](#)
- [Hadoop Environment, 111](#)
- [Creating a Single Data Object Profile in Informatica Developer, 112](#)
- [Creating an Enterprise Discovery Profile in Informatica Developer, 113](#)
- [Creating a Column Profile in Informatica Analyst, 114](#)
- [Creating an Enterprise Discovery Profile in Informatica Analyst, 115](#)
- [Creating a Scorecard in Informatica Analyst, 116](#)
- [Monitoring a Profile, 117](#)
- [Troubleshooting, 117](#)

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 the 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**, 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 with 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 with 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**, 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:

- You can choose a filter that you created for the profile. Click **Next**.
- Select an existing filter. Click 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**.

The filter appears in the **Metric Filters** pane. You can apply the same filter to all the metrics in the scorecard.

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 and configure the valid values from the list of all values in the **Score using: Values** pane. You can perform the following tasks in the **Step 5 of 7** screen:
 - You can select multiple 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.
 - In the **Metric Thresholds** pane, configure metric thresholds.
You can set thresholds for **Good**, **Acceptable**, and **Unacceptable** scores.
 - Select each metric and configure the cost of invalid data. To assign a constant value to the cost for the metric, select **Fixed Cost**. Optionally, click **Change Cost Unit** to change the unit of cost or choose **None**. To attach a numeric column as a variable cost to the metric, select **Variable Cost**, and click **Select Column** to select a numeric column.
10. 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**.
11. In the **Step 7 of 8** screen, specify the weights for the metrics in the group and thresholds for the group.
12. In the **Step 8 of 8** screen, select **Native** or **Hadoop** as the run-time environment to run the scorecard.
13. 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 HadoopEnv.properties file. 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 8

Native Environment Optimization

This chapter includes the following topics:

- [Native Environment Optimization Overview, 119](#)
- [Processing Big Data on a Grid, 119](#)
- [Processing Big Data on Partitions, 120](#)
- [High Availability, 121](#)

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.

APPENDIX A

Data Type Reference

This appendix includes the following topics:

- [Data Type Reference Overview, 123](#)
- [Transformation Data Type Support in a Hadoop Environment, 123](#)
- [Hive Data Types and Transformation Data Types, 124](#)
- [Hive Complex Data Types, 126](#)
- [Sqoop Data Types, 126](#)

Data Type Reference Overview

Informatica Developer uses the following data types in Hive mappings:

- Hive native data types. Hive data types appear in the physical data object column properties.
- Transformation data types. 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.

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
Bigint	Supported
Binary	Supported
Date/Time	Supported

Transformation Data Type	Support
Decimal	Supported
Double	Supported
Integer	Supported
String	Supported
Text	Supported
timestampWithTZ	Not supported

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

Hive Data Type	Transformation Data Type	Range and Description
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
Boolean	Integer	<p>1 or 0</p> <p>The default transformation type for boolean is integer. You can also set this to string data type with values of True and False.</p>
Arrays	String	1 to 104,857,600 characters
Struct	String	1 to 104,857,600 characters
Maps	String	1 to 104,857,600 characters
Timestamp	datetime	The time stamp format is YYYY-MM-DD HH:MM:SS.ffffff. 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 in 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
- 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

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
- Tinyint
- Varchar

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
- Timestamp
- Varchar

Teradata Data Types with TDCH Specialized Connectors for Sqoop

Informatica supports the following Teradata data types when you use the Sqoop Cloudera Connector Powered by Teradata and Hortonworks Connector for Teradata with Sqoop:

- Bigint
- Byte (supported only by Hortonworks Connector for Teradata)
- Byteint
- Character
- Date
- Decimal
- Double Precision/Float/Real
- Integer
- Number(P,S)
- Numeric
- Smallint
- Time
- Timestamp (supported only by Cloudera Connector Powered by Teradata)
- Varchar
- Varbyte (supported only by Hortonworks Connector for Teradata)

APPENDIX B

Function Reference

This appendix includes the following topic:

- [Function Support in a Hadoop Environment, 131](#)

Function Support in a Hadoop Environment

Some Informatica transformation language functions that are valid in the native environment might be supported with restrictions or unsupported in a Hadoop environment.

The following table lists functions and levels of support for functions on different engines in a Hadoop environment:

Function	Blaze Engine	Spark Engine	Hive Engine
ABORT	Not supported	Not supported	Not supported
AES_DECRYPT	Supported	Supported	Supported
AES_ENCRYPT	Supported	Supported	Supported
ANY	Supported	Supported with restrictions	Supported with restrictions
AVG	Supported	Supported with restrictions	Supported with restrictions
COMPRESS	Supported	Supported	Supported
COUNT	Supported	Supported with restrictions	Supported with restrictions
CRC32	Supported	Supported	Supported
CREATE_TIMESTAMP_TZ	Supported	Not supported	Not supported
CUME	Not supported	Not supported	Not supported
DECOMPRESS	Supported	Supported with restrictions	Supported
DEC_BASE64	Supported	Supported	Supported

Function	Blaze Engine	Spark Engine	Hive Engine
ENC_BASE64	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
MAX (Dates)	Supported	Supported	Not supported
MAX (Numbers)	Supported	Supported with restrictions	Supported with restrictions
MAX (String)	Supported	Supported with restrictions	Supported with restrictions
MD5	Supported	Supported	Supported
MIN (Dates)	Supported	Supported	Not supported
MIN (Numbers)	Supported	Supported with restrictions	Supported with restrictions
MIN (String)	Supported	Supported with restrictions	Supported with restrictions
MOVINGAVG	Not supported	Not supported	Not supported
MOVINGSUM	Not supported	Not supported	Not supported
PERCENTILE	Supported	Supported with restrictions	Supported with restrictions
STDDEV	Supported	Supported with restrictions	Supported with restrictions
SUM	Supported	Supported with restrictions	Supported with restrictions
SYSTIMESTAMP	Supported	Supported with restrictions	Supported
TO_DECIMAL	Supported	Supported with restrictions	Supported with restrictions
TO_DECIMAL38	Supported	Supported with restrictions	Supported with restrictions
TO_TIMESTAMP_TZ	Supported	Not supported	Supported
UUID4	Supported	Supported	Supported with restrictions

Function	Blaze Engine	Spark Engine	Hive Engine
UUID_UNPARSE	Supported	Supported	Supported with restrictions
VARIANCE	Supported	Supported with restrictions	Supported with restrictions
<i>Functions not listed in this table are supported on all engines without restrictions.</i>			

APPENDIX C

Parameter Reference

This appendix includes the following topics:

- [Parameters Overview, 134](#)
- [Parameter Usage, 135](#)

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.
- PostgreSQL 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.

APPENDIX D

Multiple Blaze Instances on a Cluster

This appendix includes the following topics:

- [Overview, 137](#)
- [Step 1. Prepare the Hadoop Cluster for the Blaze Engine, 138](#)
- [Step 2. Configure Data Integration Service Properties, 139](#)
- [Step 3. Update `hadoopEnv.properties`, 141](#)
- [Step 4. Create a Hadoop Connection, 143](#)
- [Step 5. Set Mapping Preferences, 145](#)
- [Result, 146](#)

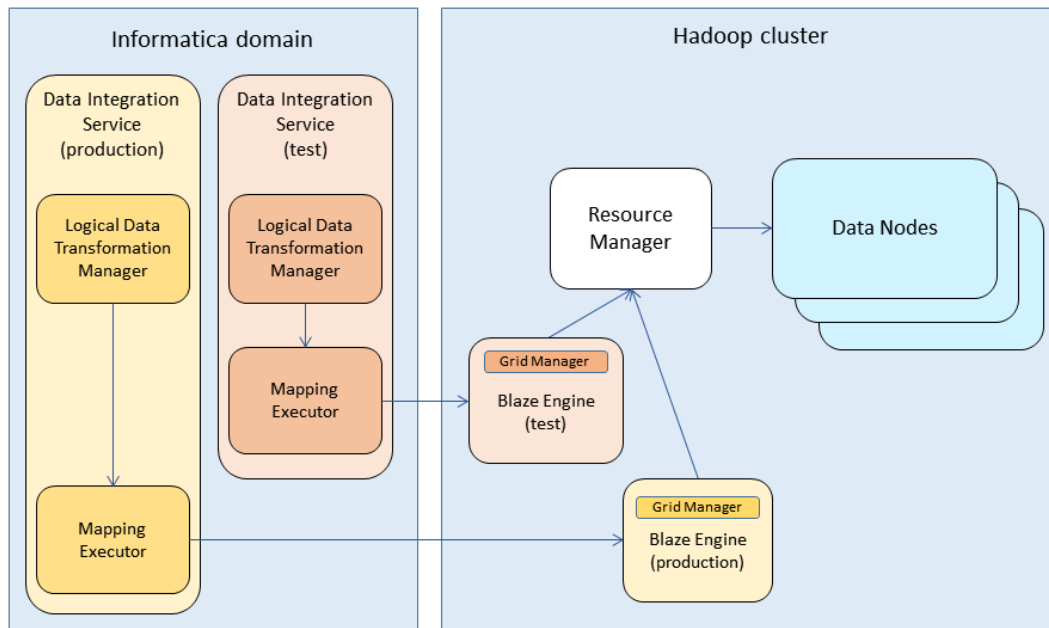
Overview

When you use the Blaze engine to run mappings, Blaze uses a Grid Manager at run time to allot tasks to various nodes in a Hadoop cluster. The Grid Manager aids in resource allocation.

You can use the same Hadoop cluster to stage your test environment and establish a production environment. To control resource use on the cluster, you can establish a separate Blaze instance for testing and another for production.

Each instance requires a separate Grid Manager. You create an additional Grid Manager by performing a series of steps to create separate infrastructure for each Blaze instance, including a unique namespace and a Hadoop connection for each Blaze instance to use.

The following image shows how a separate Data Integration Service on the domain creates a separate Grid Manager on the cluster:



The image shows how separate Data Integration Services use separate Blaze instances. Each instance uses a separate Grid Manager to communicate with the cluster resource manager to balance resources.

Perform the following steps to set up separate Blaze instances:

- Step 1. Prepare the Hadoop cluster for the Blaze engine.
- Step 2. Configure Data Integration Service properties.
- Step 3. Update `hadoopEnv.properties`.
- Step 4. Create a new Hadoop connection.
- Step 5. Set mapping preferences.

Step 1. Prepare the Hadoop Cluster for the Blaze Engine

To run mappings on the Blaze engine, perform the following tasks:

1. Create an account for the Blaze engine user.
2. Create Blaze engine directories and grant permissions.
3. Grant permissions on the Hive source database.

Create a Blaze User Account

On all nodes in the Hadoop cluster, create an operating system user account for the user you want to run the additional Blaze instance. For example, run the following command:

```
useradd testuser1
```

Create Blaze Engine Directories and Grant Permissions

Create the following directories on the Hadoop cluster:

Local services log directory

Create a local services log directory on all nodes in the cluster and grant permissions to the Blaze user account. The `hadoopEnv.properties` file on the domain contains an entry for this directory. The file uses an environment variable, `$HADOOP_NODE_INFA_HOME`, that gets set to the Big Data Management installation directory. The default installation directory is `/opt/Informatica`. For example, run the following commands:

```
hadoop fs mkdir -p /opt/Informatica/blazeLogs
hadoop fs -chmod 777 /opt/Informatica/blazeLogs
```

If you use a different directory name, you must update the following property in the `hadoopEnv.properties` file: `infagrid.node.local.root.log.dir`

HDFS temporary working directory

Create a working directory on HDFS for the Blaze engine and grant permissions to the Blaze user account. For example, run the following commands:

```
hadoop fs mkdir -p /blaze/workdir
hadoop fs -chmod 777 /blaze/workdir
```

When you configure connection properties, you provide the path to this directory. Alternatively, you can create this directory when you create the connection.

Note: This directory is separate from the aggregated persistent log directory.

Verify that a persistent aggregated HDFS log directory exists on the cluster. For example, `/var/log/Hadoop-yarn/apps/Informatica`.

Note: It is not necessary to create a new directory for persistent logs. Both Blaze instances can use the same persistent aggregated HDFS log directory.

Grant Permissions on the Hive Source Database

Grant the Blaze user account `CREATE TABLE` permission on the Hive source database. The `CREATE TABLE` permission is required in the following situations:

- The Hive source table uses SQL standard-based authorization.
- A mapping contains a Lookup transformation with an SQL override.

Step 2. Configure Data Integration Service Properties

Configure Data Integration Service properties to enable two Blaze instances on the Hadoop environment.

You can create a Data Integration Service, or configure one that has not run mappings using the Blaze engine.

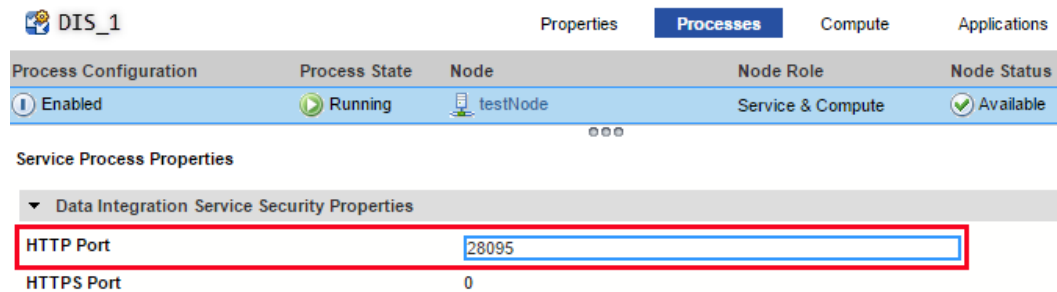
Configure Data Integration Service properties in the Administrator tool.

Data Integration Service Process Properties

Configure the following property on the **Processes** tab:

Property	Description
HTTP Port	The port that the Data Integration Service uses to communicate with the cluster over HTTP. Configure the port with a number that no other process uses.

The following image shows the HTTP Port property:



Data Integration Service Properties

The following table describes the Hadoop properties to configure for the Data Integration Service:

Property	Description
Informatica Home Directory on Hadoop	The Big Data Management home directory on every data node created by the installer. Default is <code>/opt/Informatica</code> .
Hadoop Distribution Directory	<p>The directory containing a collection of Hive and Hadoop .jar files on the cluster from the RPM Install locations. The directory contains the minimum set of .jar files required to process Informatica mappings in a Hadoop environment.</p> <p>You can duplicate an existing Hadoop distribution folder if you want to use the same domain as an existing Grid Manager. For example, if you use the distribution folder <code>cloudera_cdh<version></code> for production, you can duplicate it and name the folder like <code>cloudera_cdh<version>_test</code>.</p> <p>The distribution folders reside in <code><Informatica installation directory>/services/shared/hadoop/</code>. For example, <code><Informatica installation directory>/services/shared/hadoop/cloudera_cdh<version></code>.</p> <p>Note: When you use two Informatica domains to run mappings on the same cluster, do not duplicate the Hadoop distribution folder.</p>
Data Integration Service Hadoop Distribution Directory	<p>The Hadoop distribution directory on the Data Integration Service node.</p> <p>Type <code><Informatica installation directory>/Informatica/services/shared/hadoop/<Hadoop distribution name>_<version number></code>.</p> <p>For example:</p> <pre>../../../../services/shared/hadoop/amazon_emr_5.0.0</pre> <p>Configure the Data Integration Service Hadoop Distribution Directory property with the folder on the domain that contains the Hadoop distribution settings that you want to use for the new Blaze instance.</p> <p>Note: The contents of the Data Integration Service Hadoop distribution directory must be identical to Hadoop distribution directories on the data nodes.</p>

Step 3. Update hadoopEnv.properties

Update the `hadoopEnv.properties` file on each node where the Data Integration Service runs to configure an additional Blaze instance.

Open `hadoopEnv.properties` and back it up before you configure it. You can find the `hadoopEnv.properties` file in the following location:

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

Optionally Create a New Namespace

When the machine where the Data Integration Service runs contains two domains running on the same version of Informatica, you configure a new Blaze instance on the domain where you want to run the new Blaze instance.

In the "Advanced Configuration" section of `hadoopEnv.properties`, type the following property to designate a namespace for the Data Integration Service.

infagrid.cadi.namespace

Namespace for the Data Integration Service to use.

Configure the property as follows:

```
infagrid.cadi.namespace=<unique value>
```

For example,

```
infagrid.cadi.namespace=TestUser1_namespace
```

Configure Ports

Search for the following properties and enter port numbers that no other cluster processes use.

infagrid.blaze.console.jsfport

JSF port for the Blaze engine console.

Configure the property as follows:

```
infagrid.blaze.console.jsfport=<unique value>
```

For example,

```
infagrid.blaze.console.jsfport=9090
```

infagrid.blaze.console.httpport

HTTP port for the Blaze engine console.

Configure the property as follows:

```
infagrid.blaze.console.httpport=<unique value>
```

For example,

```
infagrid.blaze.console.httpport=9091
```

Configure Directory Paths

Search for the following properties and enter paths for the Blaze service logs and persistent logs.

infagrid.node.local.root.log.dir

Path for the Blaze service logs.

Note: This is the path that you configured in Step 1 as the local services log directory.

Configure the property as follows:

```
infagrid.node.local.root.log.dir=<directory path>
```

For example,

```
infagrid.node.local.root.log.dir=/opt/Informatica/blazeLogs
```

infacal.hadoop.logs.directory

Path in HDFS for the persistent Blaze logs.

Note: This is the path that you configured in Step 1 as the persistent log directory.

Configure the property as follows:

```
infacal.hadoop.logs.directory=<directory path>
```

For example,

```
infacal.hadoop.logs.directory=infacal.hadoop.logs.directory=/var/log/Hadoop-yarn/
apps/Informatica
```

Step 4. Create a Hadoop Connection

Create a Hadoop connection for the Blaze instance to use.

1. In Step 1 of the New Connection wizard, configure the Impersonation User Name property with the same impersonation user that you configured in Step 1, [Create a Blaze User Account](#).

The following image shows the Impersonation User Name property in the New Connection wizard:

The screenshot shows a wizard window titled "New Connection - Step 1 of 4". It contains several input fields for configuring a Hadoop connection. The "Impersonation User Name" field is highlighted with a red rectangular box and contains the text "blaze_user". Other fields include "Name" (Hadoop_connection_test), "ID" (Hadoop_connection_test), "Resource Manager Address" (hpwp:8032), "Default File System URI" (hdfs://hpwp:8020/), and "Temporary Table Compression Codec" (Snappy). The bottom of the window has buttons for "Test Connection", "< Back", "Next >", "Finish", and "Cancel".

2. In Step 3 of the New Connection wizard, configure the Temporary Working Directory on HDFS property with the path that you configured on the cluster in "[Create Blaze Engine Directories and Grant Permissions](#)."

The following image shows the Temporary Working Directory on HDFS property in the New Connection wizard:

New Connection - Step 3 of 4

Fields marked with an asterisk (*) are required.

Use this wizard to create a new connection.

Specify properties for Hadoop connection.

Blaze Service

Temporary Working Directory on HDFS *

/blaze/workdir

Blaze Service User Name

blaze_user

Minimum Port *

12300

Maximum Port *

12600

Yarn Queue Name

Blaze Service Custom Properties

?

Test Connection

< Back

Next >

Finish

Cancel

- Configure the Blaze Service User Name property with the same user name that you used to configure the Impersonation User in Step 1 of this topic.
- Configure the Minimum Port and Maximum Port properties with a port range for the connection.
You can supply a range within the port range of an existing Grid Manager, as long as ports are available when the mapping runs. The default range is 300 ports.

The following image shows the Minimum Port and Maximum Port properties in the New Connection wizard:

New Connection - Step 3 of 4

Fields marked with an asterisk (*) are required.

Use this wizard to create a new connection.

Specify properties for Hadoop connection.

Blaze Service

Temporary Working Directory on HDFS *

/blaze/workdir

Blaze Service User Name

blaze_user

Minimum Port *

12300

Maximum Port *

12600

Yarn Queue Name

Blaze Service Custom Properties

?

Test Connection

< Back

Next >

Finish

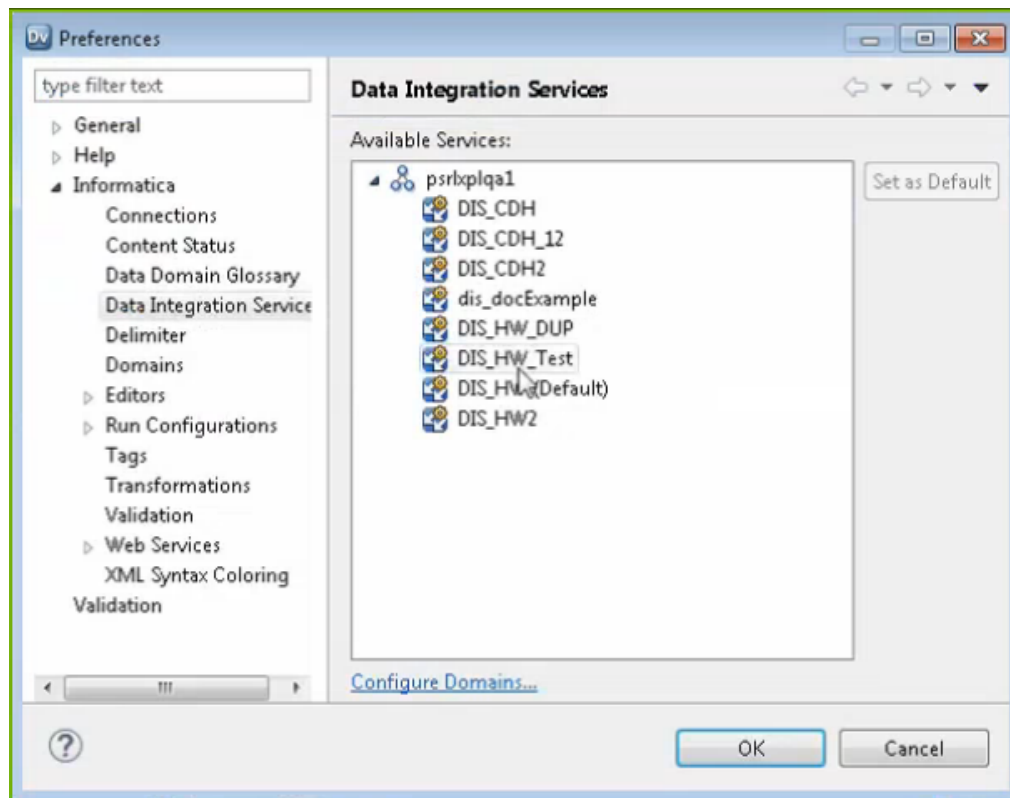
Cancel

Step 5. Set Mapping Preferences

Before you run the mapping in the Developer tool, configure the mapping to use the Data Integration Service and Hadoop connection you want to use to run the mapping.

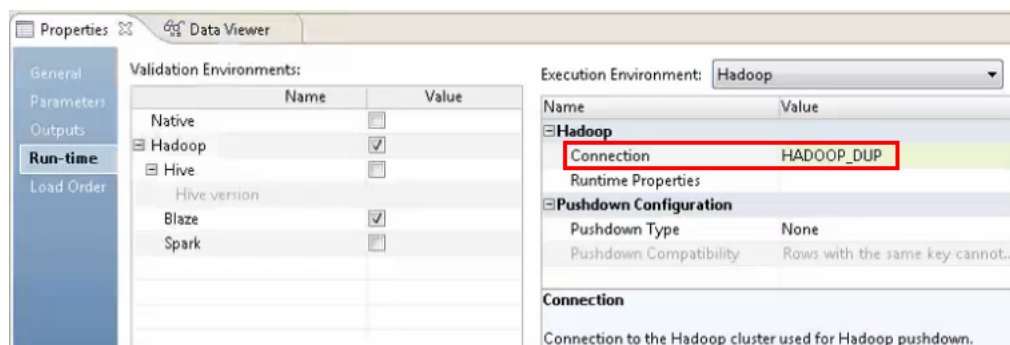
1. In the Developer tool, select **Mapping > Preferences**.
2. Expand the **Informatica** node, and then select **Data Integration Service**.

The following image shows the list of available services in the **Preferences** window:



3. Select the Data Integration Service that you want to use, and then click **OK**.
4. In the **Properties** tab of the mapping, select the **Run-time** sub-tab.
5. In the Execution Environment area, set the Connection property to the Hadoop connection that you created.

The following image shows the Connection property in the **Properties** tab:

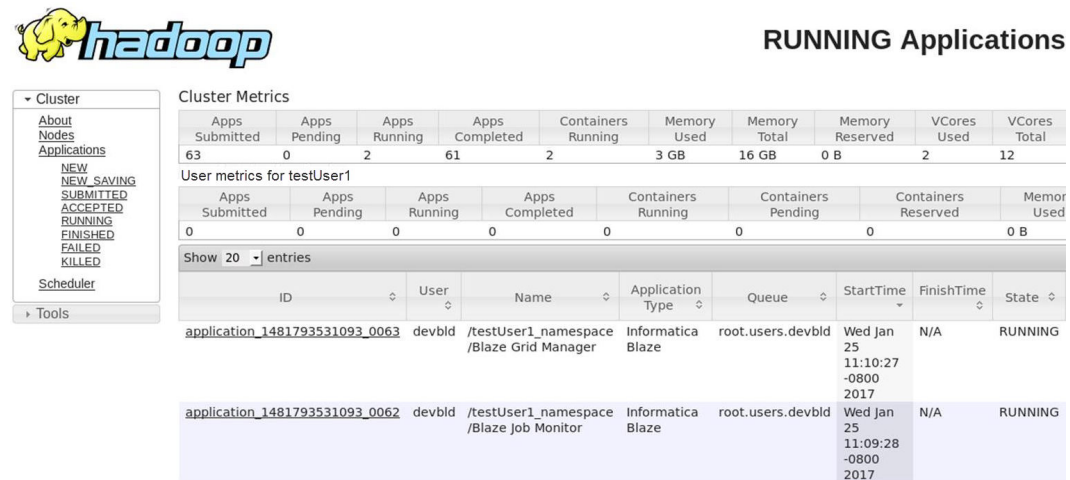


Result

The Data Integration Service creates a Grid Manager on the cluster the first time that it runs a mapping using the Blaze engine.

After you run the mapping, you can verify that the mapping used the Data Integration Service and new Grid Manager that you intended to use to run the mapping. Verify the resources that the mapping used by examining the Running Applications list in the Hadoop Resource Manager web interface. Look for applications that correspond to the namespace that you configured for the Blaze instance.

The following image shows applications with a name that includes the namespace, "testuser1_namespace," that you configured for the Grid Manager:



The screenshot displays the Hadoop Resource Manager web interface. On the left is a navigation menu with options: Cluster, About Nodes, Applications, NEW, NEW SAVING, SUBMITTED, ACCEPTED, RUNNING, FINISHED, FAILED, KILLED, Scheduler, and Tools. The 'Applications' section is selected. The main area is titled 'RUNNING Applications'. It features a 'Cluster Metrics' table and a 'User metrics for testUser1' table. Below these is a table of running applications.

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	VCores Used	VCores Total
63	0	2	61	2	3 GB	16 GB	0 B	2	12

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Containers Pending	Containers Reserved	Memory Used
0	0	0	0	0	0	0	0 B

ID	User	Name	Application Type	Queue	StartTime	FinishTime	State
application_1481793531093_0063	devbld	/testUser1_namespace /Blaze Grid Manager	Informatica Blaze	root.users.devbld	Wed Jan 25 11:10:27 -0800 2017	N/A	RUNNING
application_1481793531093_0062	devbld	/testUser1_namespace /Blaze Job Monitor	Informatica Blaze	root.users.devbld	Wed Jan 25 11:09:28 -0800 2017	N/A	RUNNING

After completion of the mapping run, the Grid Manager persists. The mapping uses the same Grid Manager whenever it runs with the unique combination of Data Integration Service and connection.

To use multiple connections that use the same Grid Manager, use an identical namespace in each connection to refer to the Blaze instance. Verify that each connection also uses identical values for the Blaze user name and queue name. If you use different values for the Blaze user name and queue name to connect to the same Blaze instance, the mapping fails.

INDEX

A

architecture
 Big Data Management [15](#)
 Hadoop environment [16](#)

B

big data
 access [12](#)
 application services [16](#)
 big data process [21](#)
 data lineage [14](#)
 repositories [16](#)
big data process
 collect your data [21](#)
Blaze engine
 transformation support [79](#)
 Blaze engine architecture [18](#)
 connection properties [24](#)
 mapping properties [96](#)
 monitoring [92–94](#)
 segment time [95](#)
 summary report [94–96, 98](#)
 tasklet execution time [96](#)
 tasklet information [98](#)
Blaze execution plan
 monitoring [91](#)
Blaze Job Monitor
 logging [99](#)

C

complex file sources
 in Hadoop environment [71](#)
component architecture
 clients and tools [15](#)
connections
 properties [24](#)
 HBase [23](#)
 HDFS [23](#)
 Hive [23](#)
 JDBC [23](#)
creating a column profile
 profiles [114](#)

D

Data Discovery
 description [13](#)
Data Integration Service grid [120](#)
data object profiles
 creating a single profile [112](#)

data object profiles (*continued*)
 enterprise discovery [113](#)
data types
 Hive [124](#)
 Hive complex data types [126](#)
 processing in a Hadoop environment [86](#)
 support [123](#)

E

enterprise discovery
 running in Informatica Analyst [115](#)
execution plan
 Spark engine [59](#)

F

flat file sources
 in Hadoop environment [68](#)
functions
 processing in a Hadoop environment [86](#)

G

grid
 Data Integration Service [120](#)
 description [119](#)
 optimization [120](#)

H

Hadoop [23](#)
Hadoop connections
 creating [46](#)
Hadoop environment
 transformation support [78](#)
 complex file sources [71](#)
 flat file limitations [68](#)
 flat file targets [73](#)
 Hive targets [74](#)
 logs [88](#)
 optimization [60](#)
 parameter usage [135](#)
 parameters [134](#)
 relational sources [71](#)
 Sqoop sources restrictions [72](#)
 Sqoop targets restrictions [77](#)
 valid sources [67](#)
Hadoop execution plan
 description, for mapping [48](#)
 overview [57](#)

- Hadoop mapping
 - run-time properties [49](#)
- hadoop utilities
 - Sqoop [17](#)
- HBase connections
 - properties [34](#)
- HDFS connections
 - creating [45](#)
 - properties [33](#)
- HDFS mappings
 - data extraction example [106](#)
 - description [106](#)
- high availability
 - description [121](#)
- Hive
 - target limitations [74](#)
- Hive connections
 - creating [45](#)
 - properties [35](#)
- Hive engine
 - data type processing [87](#)
 - function processing [87](#)
 - rules and guidelines [87](#)
 - transformation support [83](#)
 - Hive engine architecture [20](#)
 - Hive engine execution plan [60](#)
 - monitoring [103](#)
- Hive execution plan
 - monitoring [91](#)
- Hive mappings
 - description [107](#)
 - workflows [56](#)
- Hive pushdown
 - connection properties [24](#)
- Hive query
 - description, for mapping [48](#)
- Hive query plan
 - viewing, for mapping [60](#)
 - viewing, for profile [117](#)
- Hive script
 - description, for mapping [48](#)
- Hive sources
 - with Blaze engine [69](#)
 - with Informatica mappings [68](#)
- Hive targets
 - with Blaze engine [74](#)

I

- Informatica Big Data Management
 - overview [11](#)
- Informatica engine
 - Informatica engine execution plan [58](#)

J

- JDBC connections
 - properties [41](#)
 - Sqoop configuration [41](#)

L

- logging
 - mapping run on Hadoop [99](#)
 - Spark engine [102](#)

- logs
 - Blaze engine [98](#)
 - Hadoop environment [88](#)
 - Hive engine [104](#)
 - Spark engine [102](#)
- logs URL
 - YARN web user interface [102](#)

M

- mapping example
 - Hive [108](#)
 - Twitter [109](#)
- mapping execution plans
 - overview [57](#)
- mapping run on Hadoop
 - logging [90](#)
 - monitoring [91](#)
 - overview [48](#)
- MDM Big Data Relationship Management
 - description [14](#)
- Monitoring URL
 - Blaze and Spark jobs [89](#)

N

- native environment
 - high availability [121](#)
 - mappings [105](#)
 - optimization [119](#)
 - partitioning [120](#)

O

- optimization
 - compress temporary staging tables [62](#)
 - node labeling [64](#)
 - queuing [64](#)
 - scheduling [64](#)
 - truncate partitions [61](#)

P

- parameters
 - Hadoop environment [134](#)
- partitioning
 - description [120](#)
 - optimization [121](#)
- profile run on Blaze engine
 - Overview [110](#)
- profile run on Hadoop
 - monitoring [117](#)
- profile run on Hive
 - Overview [110](#)
- profiles
 - creating a column profile [114](#)

R

- rules and guidelines
 - Spark engine [86](#)
 - Hive engine [87](#)

run-time properties
Hadoop mapping [49](#)

S

social media mappings
description [108](#)

sources
in Hadoop environment [67](#)

Spark
execution plans [59](#)

Spark deploy mode
Hadoop connection properties [24](#)

Spark engine
data type processing [86](#)
function processing [86](#)
rules and guidelines [86](#)
transformation support [82](#)
connection properties [24](#)
monitoring [100](#)

Spark Event Log directory
Hadoop connection properties [24](#)

Spark execution parameters
Hadoop connection properties [24](#)

Spark HDFS staging directory
Hadoop connection properties [24](#)

Sqoop configuration
mapping properties [55](#)
profiling [111](#)

Sqoop connection arguments
-Dsqoop.connection.factories [43](#)
connect [43](#)
direct [43](#)
driver [43](#)

Sqoop connectivity
supported data types [126](#)

Sqoop data types
Aurora [126](#)
Greenplum [127](#)
IBM DB2 [127](#)
IBM DB2 for z/OS [127](#)
Microsoft SQL Server [128](#)
Netezza [128](#)
Oracle [129](#)
Teradata [129](#)
Teradata Data Types with TDCH Sqoop Specialized Connectors [130](#)

Sqoop mapping arguments
batch [55](#)
m [54](#)
num-mappers [54](#)

Sqoop mapping arguments (*continued*)
split-by [54](#)

Sqoop mappings
overview [53](#)
supported engines [53](#)

Sqoop sources
in Hadoop environment [72](#)

Sqoop targets
in Hadoop environment [77](#)

T

targets
flat files in Hadoop mapping [73](#)

TDCH connection factory
-Dsqoop.connection.factories [43](#)

third-party tools
hadoop cluster [17](#)

transformations
in Hadoop environment [78](#)
support on the Blaze engine [79](#)
support on the Spark engine [82](#)

U

utilities
hadoop cluster [17](#)

V

validation environments
Hadoop mapping [49](#)

Vibe Data Stream
description [14](#)

W

workflows
Hive mappings [56](#)

Y

YARN web user interface
description [89](#)