

Oracle® Big Data Discovery

Installation Guide

Version 1.4.0 • Revision B • January 2017

Copyright and disclaimer

Copyright © 2015, 2017, Oracle and/or its affiliates. All rights reserved.

Oracle and Java are registered trademarks of Oracle and/or its affiliates. Other names may be trademarks of their respective owners. UNIX is a registered trademark of The Open Group.

This software and related documentation are provided under a license agreement containing restrictions on use and disclosure and are protected by intellectual property laws. Except as expressly permitted in your license agreement or allowed by law, you may not use, copy, reproduce, translate, broadcast, modify, license, transmit, distribute, exhibit, perform, publish or display any part, in any form, or by any means. Reverse engineering, disassembly, or decompilation of this software, unless required by law for interoperability, is prohibited.

The information contained herein is subject to change without notice and is not warranted to be error-free. If you find any errors, please report them to us in writing.

If this is software or related documentation that is delivered to the U.S. Government or anyone licensing it on behalf of the U.S. Government, the following notice is applicable:

U.S. GOVERNMENT END USERS: Oracle programs, including any operating system, integrated software, any programs installed on the hardware, and/or documentation, delivered to U.S. Government end users are "commercial computer software" pursuant to the applicable Federal Acquisition Regulation and agency-specific supplemental regulations. As such, use, duplication, disclosure, modification, and adaptation of the programs, including any operating system, integrated software, any programs installed on the hardware, and/or documentation, shall be subject to license terms and license restrictions applicable to the programs. No other rights are granted to the U.S. Government.

This software or hardware is developed for general use in a variety of information management applications. It is not developed or intended for use in any inherently dangerous applications, including applications that may create a risk of personal injury. If you use this software or hardware in dangerous applications, then you shall be responsible to take all appropriate fail-safe, backup, redundancy, and other measures to ensure its safe use. Oracle Corporation and its affiliates disclaim any liability for any damages caused by use of this software or hardware in dangerous applications.

This software or hardware and documentation may provide access to or information on content, products and services from third parties. Oracle Corporation and its affiliates are not responsible for and expressly disclaim all warranties of any kind with respect to third-party content, products, and services. Oracle Corporation and its affiliates will not be responsible for any loss, costs, or damages incurred due to your access to or use of third-party content, products, or services.

Table of Contents

Copyright and disclaimer	2
Preface	6
About this guide	6
Audience	6
Conventions	6
Contacting Oracle Customer Support	7
 Part I: Before You Install	
Chapter 1: Introduction	9
Big Data Discovery overview	9
Studio	9
Data Processing	10
The Dgraph	10
Integration with Hadoop	11
Integration with WebLogic	11
Integration with Jetty	11
Cluster configurations and diagrams	12
A note about component names	14
Chapter 2: Prerequisites	15
Supported platforms	15
Hardware requirements	20
Memory requirements	20
Disk space requirements	21
Network requirements	22
Supported operating systems	22
Required Linux utilities	22
Installing the required Perl modules	23
OS user requirements	24
Enabling passwordless SSH	24
Hadoop requirements	25
YARN setting changes	27
Required Hadoop client libraries	27
Required HDP JARs	28
MapR-specific requirements	29
Updating the YARN ResourceManager configuration	29
Applying the MapR patches	30
JDK requirements	31
Security options	31

Kerberos	31
Sentry	32
TLS/SSL	33
HDFS data at rest encryption	34
Other security options.....	35
Component database requirements	35
Dgraph database requirements.....	36
HDFS	36
Setting up cgroups.....	37
Installing the HDFS NFS Gateway service	38
NFS	38
Increasing the numbers of open file descriptors and processes	38
Studio database requirements	39
Workflow Manager Service database requirements.....	40
Sample commands for production databases	41
Supported Web browsers	41
Screen resolution requirements	42
Studio support for iPad	42

Part II: Installing Big Data Discovery

Chapter 3: Prerequisite checklist.....	44
Chapter 4: QuickStart Installation.....	49
Installing BDD with quickstart	49
Chapter 5: Single-Node Installation	51
Installing BDD on a single node.....	51
Configuring a single-node installation	52
Chapter 6: Cluster Installation	58
Installation overview.....	58
Setting up the install machine	60
Downloading the BDD media pack.....	60
Downloading a WebLogic Server patch	61
Configuring BDD	62
Required settings	63
Running the prerequisite checker	70
Installing BDD on a cluster.....	71
Chapter 7: Troubleshooting a Failed Installation	73
Failed ZooKeeper check.....	73
Failure to download the Hadoop client libraries	73
Failure to generate the Hadoop fat JAR	74
Rerunning the installer	74

Part III: After You Install

Chapter 8: Post-Installation Tasks	77
Verifying your installation	77
Verifying your cluster's health	77
Verifying Data Processing	78
Navigating the BDD directory structure	78
Configuring load balancing	82
Configuring load balancing for Studio	82
Configuring load balancing for the Transform Service	83
Updating the DP CLI whitelist and blacklist	83
Signing in to Studio as an administrator	84
Backing up your cluster	84
Replacing certificates	84
Increasing Linux file descriptors	85
Customizing the WebLogic JVM heap size	85
Configuring Studio database caching	85
Customizing Studio database caching	86
Disabling Studio database caching	87
Clearing the Studio database cache	87
Chapter 9: Using Studio with a Reverse Proxy	88
About reverse proxies	88
Types of reverse proxies	88
Example sequence for a reverse proxy request	89
Recommendations for reverse proxy configuration	89
Preserving HTTP 1.1 Host: headers	90
Enabling the Apache ProxyPreserveHost directive	90
Reverse proxy configuration options for Studio	91
Simple Studio reverse proxy configuration	91
Studio reverse proxy configuration without preserving Host: headers	91
Configuring Studio to support an SSL-enabled reverse proxy	92

Part IV: Uninstalling Big Data Discovery

Chapter 10: Uninstalling BDD	94
-------------------------------------	-----------

Appendix A: Optional and Internal BDD Properties

Optional settings	95
Internal settings	101

Preface

Oracle Big Data Discovery is a set of end-to-end visual analytic capabilities that leverage the power of Apache Spark to turn raw data into business insight in minutes, without the need to learn specialist big data tools or rely only on highly skilled resources. The visual user interface empowers business analysts to find, explore, transform, blend and analyze big data, and then easily share results.

About this guide

This guide describes how to configure and install Oracle Big Data Discovery. It also describes tasks that should be performed immediately after installing, as well as instructions for uninstalling the product.

This guide relates specifically to Big Data Discovery version 1.4.0. The most up-to-date version of this document is available on the <http://www.oracle.com/technetwork/index.html>.



Note: This guide does *not* describe how to install Big Data Discovery on the Oracle Big Data Appliance. If you want to install on the Big Data Appliance, see the *Oracle Big Data Appliance Owner's Guide Release 4 (4.x)* and the corresponding MOS note.

Audience

This guide addresses administrators and engineers who need to install and deploy Big Data Discovery within their existing Hadoop environment.

Conventions

The following conventions are used in this document.

Typographic conventions

The following table describes the typographic conventions used in this document.

Typeface	Meaning
User Interface Elements	This formatting is used for graphical user interface elements such as pages, dialog boxes, buttons, and fields.
Code Sample	This formatting is used for sample code segments within a paragraph.
<i>Variable</i>	This formatting is used for variable values. For variables within a code sample, the formatting is <i>Variable</i> .
File Path	This formatting is used for file names and paths.

Path variable conventions

This table describes the path variable conventions used in this document.

Path variable	Meaning
\$ORACLE_HOME	Indicates the absolute path to your Oracle Middleware home directory, where BDD and WebLogic Server are installed.
\$BDD_HOME	Indicates the absolute path to your Oracle Big Data Discovery home directory, \$ORACLE_HOME/BDD-<version>.
\$DOMAIN_HOME	Indicates the absolute path to your WebLogic domain home directory. For example, if your domain is named bdd-<version>_domain, then \$DOMAIN_HOME is \$ORACLE_HOME/user_projects/domains/bdd-<version>_domain.
\$DGRAPH_HOME	Indicates the absolute path to your Dgraph home directory, \$BDD_HOME/dgraph.

Contacting Oracle Customer Support

Oracle customers that have purchased support have access to electronic support through My Oracle Support. This includes important information regarding Oracle software, implementation questions, product and solution help, as well as overall news and updates from Oracle.

You can contact Oracle Customer Support through Oracle's Support portal, My Oracle Support at <https://support.oracle.com>.

Part I

Before You Install



Chapter 1

Introduction

The following sections describe Oracle Big Data Discovery (BDD) and how it integrates with other software products. They also describe some of the different cluster configurations BDD supports.

[*Big Data Discovery overview*](#)

[*Integration with Hadoop*](#)

[*Integration with WebLogic*](#)

[*Integration with Jetty*](#)

[*Cluster configurations and diagrams*](#)

[*A note about component names*](#)

Big Data Discovery overview

Oracle Big Data Discover is made up of a number of distinct components: three main ones and a number of others that work closely with them.

[*Studio*](#)

[*Data Processing*](#)

[*The Dgraph*](#)

Studio

Studio is BDD's front-end web application. It provides tools that you can use to create and manage data sets and projects, as well as administrator tools for managing end user access and other settings.

Studio is a Java-based application. It runs in a Java container provided by the WebLogic Server, which is automatically installed with BDD.

Transform Service

The Transform Service processes end user-defined changes to data sets, called *transformations*, on behalf of Studio. It enables users to preview the effects their transformations will have on their data before saving them.

The Transform Service is a RESTful web application that runs inside a Jetty container. Like WebLogic Server, Jetty is automatically installed with BDD.

Data Processing

Data Processing collectively refers to a set of processes and jobs that discover, sample, profile, and enrich source data.

Many of these processes run within Hadoop, so the Data Processing libraries are installed on Hadoop nodes.

Workflow Manager Service

The Workflow Manager Service launches and manages all Data Processing jobs on behalf of the other BDD components. It runs inside its own Jetty container, like the Transform Service, and can only be installed on a single node.

Data Processing CLI

The Data Processing Command Line Interface (DP CLI) provides a way to manually launch Data Processing jobs and invoke the Hive Table Detector (see below). It can also be configured to run automatically as a cron job.

The DP CLI is installed on all Studio and Dgraph nodes. It can later be moved to any node that has access to the BDD cluster.

Hive Table Detector

The Hive Table Detector is a Data Processing component that monitors the Hive database for new and deleted tables, and launches Data Processing workflows in response.

The Hive Table Detector is invoked by the CLI, either manually by the Hive administrator or via the CLI cron job. If you enable the CLI to run as a cron job, the Hive Table Detector runs at each invocation of the cron job.

The Dgraph

The Dgraph indexes the data sets produced by Data Processing and stores them in databases. It also responds in real time to Studio user queries for the indexed data, which are routed to it by the Dgraph Gateway.

Dgraph Gateway

The Dgraph Gateway is a Java-based interface that routes requests to the Dgraph instances and provides caching and business logic. It also handles cluster services for the Dgraph instances by leveraging Apache ZooKeeper, which is part of the Hadoop ecosystem.

The Dgraph Gateway runs inside WebLogic Server, along with Studio.

Dgraph HDFS Agent

The Dgraph HDFS Agent acts as a data transport layer between the Dgraph and the Hadoop Distributed File System (HDFS). It exports records to HDFS on behalf of the Dgraph, and imports them from HDFS during data ingest operations.

The HDFS Agent is automatically installed on the same nodes as the Dgraph.

Integration with Hadoop

BDD runs on top of an existing Hadoop cluster, which provides a number of components and tools that BDD requires to process and manage data. For example, the source data you load into BDD is stored in HDFS and processed by Spark on YARN.

BDD supports the following Hadoop distributions:

- Cloudera Distribution for Hadoop (CDH)
- Hortonworks Data Platform (HDP)
- MapR Converged Data Platform (MapR)

You must have one of these installed on your cluster before installing BDD, as the configuration of your Hadoop cluster determines where many of the BDD components will be installed. For supported versions and a list of required Hadoop components, see [Hadoop requirements on page 25](#).

Integration with WebLogic

WebLogic Server provides a J2EE container for hosting and managing Studio and the Dgraph Gateway, which are J2EE applications. Additionally, WebLogic's Admin Server plays an important role in installing and administering BDD.

WebLogic Server 12c (12.1.3) is included in the BDD media pack and automatically installed on all nodes that will host Studio and the Dgraph Gateway.



Note: BDD does not currently support integration with an existing WebLogic installation. You must use the version included with the BDD packages.

The WebLogic Admin Server serves as a central point of control for your BDD cluster. Before installing, you select a node to be the Admin Server and perform the entire installation from it. After installation, you can perform script-based administrative tasks, such as starting individual components and updating the cluster configuration, from this node.

You can also use the WebLogic Administration Console and WLST (WebLogic Server Scripting Tool) for starting and stopping the Managed Servers that host Studio and the Dgraph Gateway.

Integration with Jetty

Jetty provides open-source `javax.servlet` containers for hosting the Transform Service and the Workflow Manager Service.

Jetty 9 is included in the BDD media pack and automatically installed on all nodes that will host the Transform Service and Workflow Manager Service.

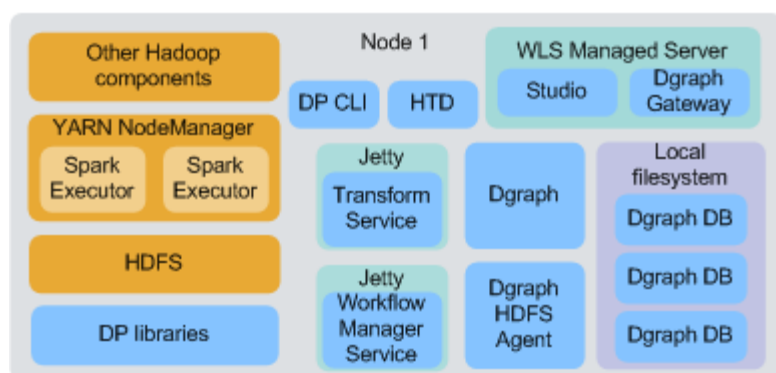
Cluster configurations and diagrams

BDD supports many different cluster configurations. The following sections describe three suitable for demonstration, development, and production environments, and their possible variations.

Note that you aren't limited to these examples and can install in any configuration suited to your resources and data processing needs.

Single-node demo environment

You can install BDD in a demo environment running on a single physical or virtual machine. This configuration can only handle a limited amount of data, so it is recommended solely for demonstrating the product's functionality with a small sample database.

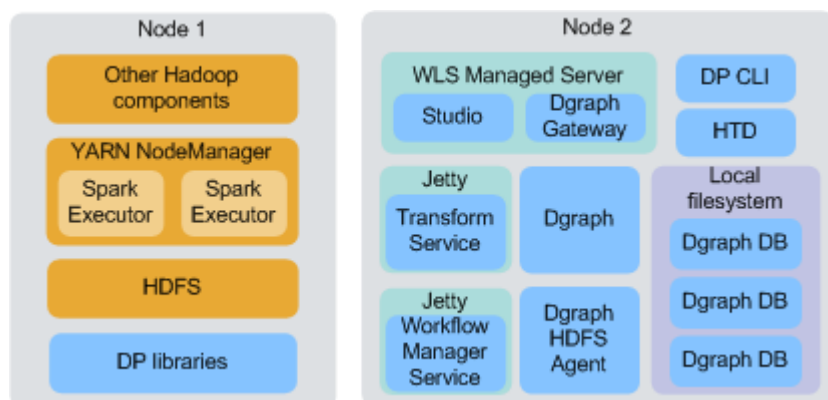


In a single-node deployment, all BDD and required Hadoop components are hosted on the same node, and the Dgraph databases are stored on the local filesystem.

For single-node installations, BDD provides a QuickStart option that enables you to install quickly with default configuration. For more information, see [QuickStart Installation on page 48](#).

Two-node development environment

You can install BDD in a development environment running on two nodes. This configuration can handle a slightly larger database than a single-node deployment, but still has limited processing capacity and doesn't provide high availability for any BDD components.

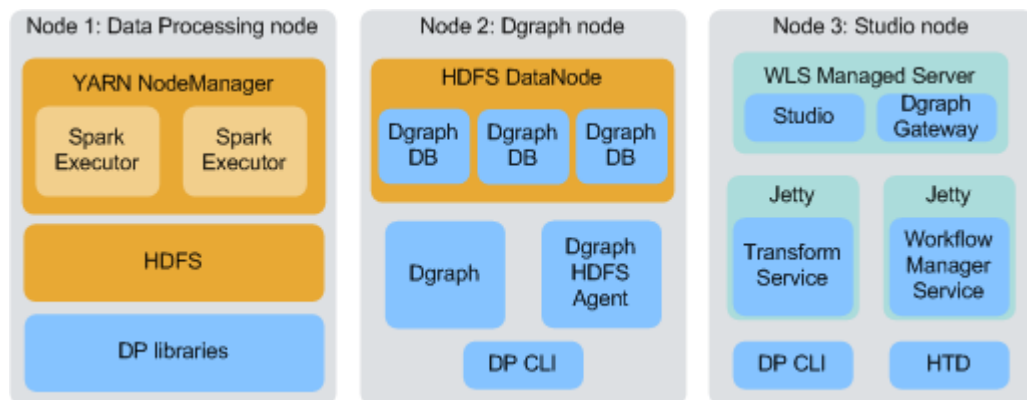


In a two-node configuration, all Hadoop components and the Data Processing libraries are hosted on one node, while the remaining BDD components are hosted on the other. Although the diagram above shows the Dgraph databases stored on the local file system, they could also be stored in HDFS on Node 1.

Multi-node production environment

For a production environment, BDD should be installed on a multi-node cluster. The size of your cluster depends on the amount of data you plan on processing and the number of end users expected to be querying that data at any given time; however, a minimum of six nodes ensures high availability for all components.

A typical BDD cluster will consist of nodes that look similar to the following:



- Node 1 is running the Data Processing libraries, along with the YARN NodeManager service, Spark on YARN, and HDFS, all of which Data Processing requires to function. The number of Data Processing nodes your cluster should contain depends on the amount of data you have and its size, although a minimum of three ensures high availability.
- Node 2 is running the Dgraph, the Dgraph HDFS Agent, the DP CLI, and the HDFS DataNode service. The Dgraph databases are stored in HDFS, which is recommended for production environments. (Note that they could also be stored on an NFS, in which case the DataNode service wouldn't be required.) A typical cluster would contain two or more Dgraph nodes.
- Node 3 is running Studio and the Dgraph Gateway inside a WebLogic Managed Server container; the Transform Service and the Workflow Manager service, each inside a Jetty container; the DP CLI; and the Hive Table Detector. A typical cluster would contain one or more Studio nodes, depending on the number of users making concurrent queries. Note that in a cluster with multiple Studio nodes, the Workflow Manager Service and Hive Table Detector would each only be installed on one of them. Additionally, one Studio node in the cluster must serve as the Admin Server.

Co-locating components

One way to configure your BDD cluster is to co-locate different components on the same nodes. This is a more efficient use of your hardware, since you don't have to devote an entire node to any specific component.

Be aware, however, that the co-located components will compete for memory and other resources, which can have a negative impact on performance. The decision to host different components on the same nodes depends on your site's production requirements and your hardware's capacity.

Most combinations of BDD and Hadoop components will work. If you choose to do this, however, keep the following in mind:

- You shouldn't co-locate the Dgraph with Hadoop components, other than the HDFS DataNode service. In particular, you shouldn't host it on the same nodes as Spark, as both require a lot of memory. If you have to do this, you should use cgroups ensure each has access to sufficient resources. For more information, see [Setting up cgroups on page 37](#).
- Similarly, you shouldn't co-locate the Dgraph with the Transform Service, which also requires a lot of memory.
- While you can co-locate Managed Servers with either the Dgraph or any Hadoop components, you should limit the amount of memory WebLogic Server can consume to ensure the other components have access to the resources they require.

A note about component names

Some of the installation files and scripts may contain references to the Endeca Server, which is a legacy name for the Dgraph Gateway. This document refers to the component as the Dgraph Gateway, and notes any discrepancies to avoid confusion.



Chapter 2

Prerequisites

The following sections describe the hardware and software requirements your environment must meet before you can install BDD.

[*Supported platforms*](#)

[*Hardware requirements*](#)

[*Memory requirements*](#)

[*Disk space requirements*](#)

[*Network requirements*](#)

[*Supported operating systems*](#)

[*Required Linux utilities*](#)

[*OS user requirements*](#)

[*Hadoop requirements*](#)

[*JDK requirements*](#)

[*Security options*](#)

[*Component database requirements*](#)

[*Supported Web browsers*](#)

[*Screen resolution requirements*](#)

[*Studio support for iPad*](#)

Supported platforms

The following tables list the platforms and versions supported in each BDD release.

Note that this is not an exhaustive list of BDD's requirements. Be sure to read through the rest of this chapter before installing for more information about the components and configuration changes BDD requires.

Supported Hadoop distributions

Big Data Discovery version	Hadoop distribution	Supported version(s)
1.0	Cloudera Distribution for Hadoop	5.3.0

Big Data Discovery version	Hadoop distribution	Supported version(s)
1.1.x	Cloudera Distribution for Hadoop Hortonworks Data Platform	5.3.x, 5.4.x, 5.5.2 2.2.4-2.3.x
1.2.0	Cloudera Distribution for Hadoop Hortonworks Data Platform	5.5.2+ 2.3.4.17-5
1.2.2	Cloudera Distribution for Hadoop Hortonworks Data Platform	5.5.x (min. 5.5.2), 5.6, 5.7.1 2.3.4.17-5, 2.4.x (min. 2.4.2)
1.3.x	Cloudera Distribution for Hadoop Hortonworks Data Platform MapR Converged Data Platform	5.5.x (min. 5.5.2), 5.6, 5.7.x (min. 5.7.1), 5.8, 5.9 2.3.4.17-5, 2.4.x (min. 2.4.2) 5.1
1.4.0	Cloudera Distribution for Hadoop Hortonworks Data Platform MapR Converged Data Platform	5.7.x (min. 5.7.1), 5.8.x, 5.9.x, 5.10.x, 5.11.x 2.4.x (min. 2.4.2) 5.1+

Supported Big Data Appliance versions



Note: This guide does *not* describe how to install Big Data Discovery on the Oracle Big Data Appliance. If you want to install on the Big Data Appliance, see the *Oracle Big Data Appliance Owner's Guide Release 4 (4.x)* and the corresponding MOS note.

Big Data Discovery version	Supported Big Data Appliance version(s)
1.0	N/A
1.1.x	4.3, 4.4
1.2.0	4.4
1.2.2	4.4, 4.5
1.3.x	4.5, 4.6, 4.7
1.4.0	4.7, 4.8, 4.9

Supported operating systems

Big Data Discovery version	Operating system	Supported version(s)
1.0	Oracle Enterprise Linux	6
	Red Hat Enterprise Linux	6
1.1.x	Oracle Enterprise Linux	6.4+
	Red Hat Enterprise Linux	6.4+
1.2.0	Oracle Enterprise Linux	6.4+, 7.1
	Red Hat Enterprise Linux	6.4+, 7.1
1.2.2	Oracle Enterprise Linux	6.4+, 7.1
	Red Hat Enterprise Linux	6.4+, 7.1
1.3.0	Oracle Enterprise Linux	6.4+, 7.1
	Red Hat Enterprise Linux	6.4+, 7.1
1.4.0	Oracle Enterprise Linux	6.4+, 7.1, 7.2
	Red Hat Enterprise Linux	6.4+, 7.1, 7.2

Supported application servers

Big Data Discovery version	Application server	Supported version(s)
1.0	Oracle WebLogic Server	12c 12.1.3
1.1.x	Oracle WebLogic Server	12c 12.1.3
1.2.0	Oracle WebLogic Server	12c 12.1.3
1.2.2	Oracle WebLogic Server	12c 12.1.3
1.3.0	Oracle WebLogic Server	12c 12.1.3
1.4.0	Oracle WebLogic Server	12c 12.1.3

Supported JDK versions

Big Data Discovery version	Supported JDK version(s)
1.0	HotSpot JDK 7U67+ x64

Big Data Discovery version	Supported JDK version(s)
1.1.x	HotSpot JDK 7u67+ x64 HotSpot JDK 8u45+ x64
1.2.0	HotSpot JDK 7u67+ x64 HotSpot JDK 8u45+ x64
1.2.2	HotSpot JDK 7u67+ x64 HotSpot JDK 8u45+ x64
1.3.0	HotSpot JDK 7u67+ x64 HotSpot JDK 8u45+ x64
1.4.0	HotSpot JDK 7u67+ x64 HotSpot JDK 8u45+ x64

Supported Studio database servers

Big Data Discovery version	Database server	Supported version(s)
1.0	Oracle MySQL Hypersonic (non-prod environments, only)	11g, 12c 12.1.0.1.0+ 5.5.3+ N/A
1.1.x	Oracle MySQL Hypersonic (non-prod environments, only)	11g, 12c 12.1.0.1.0+ 5.5.3+ N/A
1.2.0	Oracle MySQL Hypersonic (non-prod environments, only)	11g, 12c 12.1.0.1.0+ 5.5.3+ N/A
1.2.2	Oracle MySQL Hypersonic (non-prod environments, only)	11g, 12c 12.1.0.1.0+ 5.5.3+ N/A
1.3.0	Oracle MySQL Hypersonic (non-prod environments, only)	11g, 12c 12.1.0.1.0+ 5.5.3+ N/A

Big Data Discovery version	Database server	Supported version(s)
1.4.0	Oracle MySQL Hypersonic (non-prod environments, only)	11g, 12c 12.1.0.1.0+ 5.5.3+ N/A

Supported Workflow Manager database servers

Big Data Discovery version	Database server	Supported version(s)
1.4.0	Oracle MySQL	11g, 12 c 12.1.0.1.0+ 5.5.3+

Supported browsers

Big Data Discovery version	Supported browsers
1.0	Internet Explorer 10, 11 Firefox ESR Chrome for Business Safari Mobile 7.x
1.1.x	Internet Explorer 10, 11 Firefox ESR Chrome for Business Safari Mobile 8.x
1.2.0	Internet Explorer 11 Firefox ESR Chrome for Business Safari Mobile 9.x
1.2.2	Internet Explorer 11 Firefox ESR Chrome for Business Safari Mobile 9.x

Big Data Discovery version	Supported browsers
1.3.0	Internet Explorer 11 Firefox ESR Chrome for Business Safari Mobile 9.x
1.4.0	Internet Explorer 11 Microsoft Edge Firefox ESR Chrome for Business Safari Mobile 9.x

Hardware requirements

The hardware requirements for your BDD installation depend on the amount of data you will process. Oracle recommends the following minimum requirements:

- x86_64 dual-core CPU for Dgraph nodes
- x86_64 quad-core CPU for WebLogic Managed Servers, which will run Studio and the Dgraph Gateway

In this guide, the term "x64" refers to any processor compatible with the AMD64/EM64T architecture. You might need to upgrade your hardware, depending on the data you are processing. All run-time code must fit entirely in RAM. Likewise, hard disk capacity must be sufficient based on the size of your data set. Please contact your Oracle representative if you need more information on sizing your hardware.

Note that Oracle recommends turning off hyper-threading for Dgraph nodes. Because of the way the Dgraph works, hyper-threading is actually detrimental to cache performance.

Memory requirements

The amount of RAM your system requires depends on the amount of data you plan on processing.

The following table lists the minimum amounts of RAM required to install BDD on each type of node.



Important: Be aware that these are the amounts required by the product itself and don't account for storing or processing data—full-scale installations will require more. You should work with your Oracle representative to determine an appropriate amount for your processing needs before installing.

Type of node	Requirements
WebLogic	<p>16GB</p> <p>This breaks down into 5GB for WebLogic Server and 11GB for the Transform Service.</p> <p>Note that installing the Transform Service on WebLogic nodes is recommended, but not required. If you decide to host it on a different type of node, verify that it has enough RAM.</p>
Dgraph	<p>5GB</p> <p>If you're planning on storing your databases on HDFS, your Dgraph nodes should have 5GB of RAM plus the amount required by HDFS and any other Hadoop components running on them. For more information, see Dgraph database requirements on page 36.</p>
Data Processing (YARN cluster)	<p>16GB</p> <p>Note that this is for the entire YARN cluster combined, not per node.</p>

Disk space requirements

You must ensure that each node contains enough space to install BDD.

The product has the following *minimum* space requirements:

- 30GB in the `ORACLE_HOME` directory on all BDD nodes. You will define the location of this directory in BDD's configuration file before installing.
- 20GB in the `TEMP_FOLDER_PATH` directory on all BDD nodes. You will define the location of this directory in BDD's configuration file before installing.
- 10GB in the `INSTALLER_PATH` directory on the install machine. You will define the location of this directory in BDD's configuration file before installing.
- 512MB swap space on the install machine and all Managed Servers. If these nodes don't meet this requirement, be sure to set the `WLS_NO_SWAP` property in BDD's configuration file to `TRUE`.
- 39GB virtual memory on all Transform Service nodes.



Important: Be aware that these are the amounts required by the product itself and don't account for storing or processing data—full-scale installations will require more. You should work with your Oracle representative to determine an appropriate amount of space for your processing needs before installing.

Network requirements

The hostname of each BDD machine must be externally-resolvable and accessible using the machine's IP address. Oracle recommends using only Fully Qualified Domain Names (FQDNs).

Supported operating systems

BDD supports the following operating systems:

- Oracle Enterprise Linux 6.4+, 7.1, 7.2 x64
- Red Hat Enterprise Linux 6.4+, 7.1, 7.2 x64

One of these must be installed on all nodes in the BDD cluster, including Hadoop nodes.

Required Linux utilities

The BDD installer requires several Linux utilities.

The following must be present in the `/bin` directory:

```
basename  
cat  
chgrp  
chown  
date  
dd  
df  
mkdir  
more  
rm  
sed  
tar  
true
```

The following must be present in the `/usr/bin` directory:

```
awk  
cksum  
cut  
dirname  
expr  
gzip  
head  
id  
netcat  
perl (see below)  
printf  
sudo (Note: This is the default version on OEL 6.x.)  
tail  
tr  
unzip  
wc  
which
```

In addition to these, BDD requires the following:

- Perl 5.10+ with multithreading. This must be set as the default version on all BDD nodes. Additionally, the install machine requires some specific Perl modules; see [Installing the required Perl modules on page 23](#) for instructions on installing them.

- The default `umask` set to 022 on all BDD nodes, including Hadoop nodes.
- `curl 7.19.7+` with support for the `--tlsv1.2` and `--negotiate` options. This must be installed on all nodes that will host Studio.
- Network Security Services (NSS) 3.16.1+ on all nodes that will host Studio.
- `nss-devel` on all nodes that will host Studio. This contains the `nss-config` command, which must be installed in `/usr/bin`.

`nss-devel` is included in Linux 6.7 and higher, but needs to be installed manually on older versions. To see if it's installed, run:

```
sudo rpm -q nss-devel
```

If `nss-devel` is installed, the above command should return its version number. You should also verify that `nss-config` is available in `/usr/bin`.

If you don't have `nss-devel`, install it by running:

```
sudo yum install nss-devel
```

`nss-config` will be installed in `/usr/bin` by default.

- `tty` disabled for `sudo`. If it's currently enabled, comment out the line `Defaults requiretty` in `/etc/sudoers` on all nodes:

```
#Defaults requiretty
```

- Apache Ant 1.7.1+ installed and added to the `PATH` on all nodes, including Hadoop nodes.

Installing the required Perl modules

Installing the required Perl modules

The `Mail::Address` and `XML::Parser` Perl modules are required on the install machine.

You only need to perform this procedure on the install machine. These modules aren't required on any other nodes.

To install the required Perl modules:

1. Install `Mail::Address`:

① Download `Mail::Address` from <http://pkgs.fedoraproject.org/repo/pkgs/perl-MailTools/MailTools-2.14.tar.gz/813ae849683367bb75e6be89e4e8cc46/MailTools-2.14.tar.gz>.

② Extract `MailTools-2.14.tar.gz`:

```
tar -xvf MailTools-2.14.tar.gz
```

This creates a directory called `/MailTools-2.14`.

③ Go to `/MailTools-2.14` and run the following commands to install the module:

```
perl Makefile.PL
make
make test
sudo make install
```

2. Install XML::Parser:

- ① Download XML::Parser from <http://search.cpan.org/CPAN/authors/id/T/TO/TODDR/XML-Parser-2.44.tar.gz>.

- ② Extract XML-Parser-2.44.tar.gz:

```
tar -xvf XML-Parser-2.44.tar.gz
```

This creates a directory called /XML-Parser-2.44.

- ③ Go to /XML-Parser-2.44 and run the following commands to install the module:

```
perl Makefile.PL
make
make test
sudo make install
```

OS user requirements

The entire installation must be performed by a single OS user, called the `bdd` user. After installing, this user will run all BDD processes.

You must create this user or select an existing one to fill this role before installing. Although this document refers to it as the `bdd` user, its name is arbitrary.

The user you choose must meet the following requirements:

- It can't be the root user.
- Its UID must be the same on all nodes in the cluster, including Hadoop nodes.
- It must have passwordless `sudo` enabled on all nodes in the cluster, including Hadoop nodes.
- It must have passwordless SSH enabled on all nodes in the cluster, including Hadoop nodes, so that it can log into each node from the install machine. For instructions on enabling this, see [Enabling passwordless SSH on page 24](#).
- It must have bash set as its default shell on all nodes in the cluster, including Hadoop nodes.
- It must have permission to create the directory BDD will be installed in on all nodes in the cluster, including Hadoop nodes. This directory is defined by the `ORACLE_HOME` property in the BDD configuration file.

If your databases are located on HDFS, the `bdd` user has additional requirements. These are described in [Dgraph database requirements on page 36](#).

[Enabling passwordless SSH](#)

Enabling passwordless SSH

You must enable passwordless SSH on all nodes in the cluster for the `bdd` user.

To enable passwordless SSH for the `bdd` user:

1. Generate SSH keys on all nodes in the cluster, including Hadoop nodes.
2. Copy the keys to the install machine to create `known_hosts` and `authorized_keys` files.

3. Copy the `known_hosts` and `authorized_keys` files to all servers in the cluster.

Hadoop requirements

BDD supports the following Hadoop distributions:

- Cloudera Distribution for Hadoop (CDH) 5.7.x (min. 5.7.1), 5.8.x, 5.9.x, 5.10.x, 5.11.x. Enterprise edition is recommended.
- Hortonworks Data Platform (HDP) 2.4.x (min. 2.4.2)
- MapR Converged Data Platform (MapR) 5.1+



You must have one of these installed before installing BDD. Note that you can't connect BDD to more than one Hadoop cluster.



Note: You can switch to a different version of your Hadoop distribution after installing BDD, if necessary. See the *Administrator's Guide* for more information.

BDD doesn't require all of the components each distribution provides, and the components it does require don't need to be installed on all BDD nodes. The following table lists the required Hadoop components and the node(s) they must be installed on. If you're installing on a single machine, it must be running all required components.

Component	Description
Cluster manager	<p>Your cluster manager depends on your Hadoop distribution:</p> <ul style="list-style-type: none"> • CDH: Cloudera Manager • HDP: Ambari • MapR: MapR Control System (MCS) <p>The installer uses a RESTful API to query your Hadoop cluster manager for information about your Hadoop nodes, such as their hostnames and port numbers. Post-install, the <code>bdd-admin</code> script will query it for similar information when performing administrative tasks.</p> <p>Your cluster manager must be installed on at least one node in your Hadoop cluster, although it doesn't have to be on any that will host BDD.</p>
ZooKeeper	<p>BDD uses ZooKeeper to manage the Dgraph instances and ensure high availability of Dgraph query processing. ZooKeeper must be installed on at least one node in your Hadoop cluster, although to ensure high availability, it should be on three or more. These don't have to be BDD nodes, although each Managed Server must be able to connect to at least one of them.</p>

Component	Description
HDFS/MapR-FS	<p>The tables that contain your source data are stored in HDFS. It must be installed on all nodes that will run Data Processing. Additionally, if you choose to store your Dgraph databases on HDFS, the HDFS DataNode service must be installed on all Dgraph nodes.</p> <p> Note: MapR uses the MapR File System (MapR-FS) instead of standard HDFS. For simplicity, this document typically refers only to HDFS. Any requirements specific to MapR-FS will be called out explicitly.</p>
YARN	The YARN NodeManager service runs all Data Processing jobs. YARN must be installed on all nodes that will run Data Processing.
Spark on YARN	<p>BDD uses Spark on YARN to run all Data Processing jobs. Spark on YARN must be installed on all nodes that will run Data Processing.</p> <p>Note that BDD requires Spark 1.6+. Verify the version you have and upgrade it, if necessary.</p>
Hive	All of your data is stored in Hive tables within HDFS. When BDD discovers a new or modified Hive table, it launches a Data Processing workflow for that table.
HCatalog	The Hive Table Detector monitors HCatalog for new and deleted tables that require processing. HCatalog must be installed on at least one node in your Hadoop cluster, although it doesn't have to be one that will host BDD.
Hue	<p>You can use Hue to load your source data into Hive and to view data exported from Studio. Hue must be installed on at least one node in your Hadoop cluster, although it doesn't have to be one that will host BDD.</p> <p> Note: HDP doesn't include Hue. If you have HDP, you must install Hue separately and set the <code>HUE_URI</code> property in BDD's configuration file. You can also use the <code>bdd-admin</code> script to update this property after installation, if necessary. For more information, see the <i>Administrator's Guide</i>.</p>

To reiterate, Data Processing will automatically be installed on nodes running the following:

- YARN
- Spark on YARN
- HDFS

You must also make a few changes within your Hadoop cluster to ensure that BDD can communicate with your Hadoop nodes. These changes are described below.

[YARN setting changes](#)

[Required Hadoop client libraries](#)

[Required HDP JARs](#)

[MapR-specific requirements](#)

YARN setting changes

To ensure that each node in your YARN cluster has access to sufficient resources during processing, you need to update the following YARN-specific Hadoop properties.

You can access these properties from your Hadoop cluster manager (Cloudera Manager, Ambari, or MCS). If you need help locating any of them, refer to your distribution's documentation.

Property	Description
<code>yarn.nodemanager.resource.memory-mb</code>	The total amount of memory that YARN can use on a given node. This should be at least 16GB, although you might need to set it higher depending on the amount of data you plan on processing.
<code>yarn.scheduler.maximum-allocation-vcores</code>	The maximum number of virtual CPU cores allocated to each YARN container per request. If your Hadoop cluster contains only one YARN worker node, this should be less than or equal to half of that node's cores. If it contains multiple YARN worker nodes, this should be less than or equal to each node's total number of cores.
<code>yarn.scheduler.maximum-allocation-mb</code>	The maximum amount of RAM allocated to each YARN container per request. This should be at least 16GB. Additionally: <ul style="list-style-type: none"> • If your Hadoop cluster contains only one YARN node, this should be less than or equal to half of that node's RAM. • If your Hadoop cluster contains multiple YARN nodes, this should be less than or equal to each node's total amount of RAM.
<code>yarn.scheduler.capacity.maximum-applications</code>	The maximum number of concurrently-running jobs allowed on each node. This can be between 2 and 8. Note that setting this value higher could cause jobs submitted at the same time to hang indefinitely.

Required Hadoop client libraries

BDD requires a number of client libraries to interact with Hadoop. When the installer runs, it adds these libraries to a single JAR, called the Hadoop fat JAR, which is distributed to all BDD nodes.

How you obtain the client libraries depends on your Hadoop distribution:

- **CDH:** The installer will download the required libraries automatically. Note that this requires an internet connection on the install machine. If the script can't download all of the client libraries, it will fail and you will have to download them manually. See [Failure to download the Hadoop client libraries on page 73](#) for more information.
- **HDP:** Locate the following directories on your Hadoop nodes and copy them to the install machine. Note that they might not all be on the same node.
 - `/usr/hdp/<version>/hive/lib/`

- /usr/hdp/<version>/spark/lib/
 - /usr/hdp/<version>/hadoop/
 - /usr/hdp/<version>/hadoop/lib/
 - /usr/hdp/<version>/hadoop-hdfs/
 - /usr/hdp/<version>/hadoop-hdfs/lib/
 - /usr/hdp/<version>/hadoop-yarn/
 - /usr/hdp/<version>/hadoop-yarn/lib/
 - /usr/hdp/<version>/hadoop-mapreduce/
 - /usr/hdp/<version>/hadoop-mapreduce/lib/
- **MapR:** Locate the following directories on your Hadoop nodes and copy them to the install machine. Note that they might not all be on the same node.
- /opt/mapr/spark/spark-1.6.1/lib
 - /opt/mapr/hive/hive-1.2/lib
 - /opt/mapr/zookeeper/zookeeper-3.4.5
 - /opt/mapr/zookeeper/zookeeper-3.4.5/lib
 - /opt/mapr/hadoop/hadoop-2.7.0/share/hadoop/common
 - /opt/mapr/hadoop/hadoop-2.7.0/share/hadoop/common/lib
 - /opt/mapr/hadoop/hadoop-2.7.0/share/hadoop/hdfs
 - /opt/mapr/hadoop/hadoop-2.7.0/share/hadoop/hdfs/lib
 - /opt/mapr/hadoop/hadoop-2.7.0/share/hadoop/mapreduce
 - /opt/mapr/hadoop/hadoop-2.7.0/share/hadoop/mapreduce/lib
 - /opt/mapr/hadoop/hadoop-2.7.0/share/hadoop/tools/lib
 - /opt/mapr/hadoop/hadoop-2.7.0/share/hadoop/yarn
 - /opt/mapr/hadoop/hadoop-2.7.0/share/hadoop/yarn/lib

Required HDP JARs

If you have HDP, make sure that the following JAR files are present on all of your Hadoop nodes.

Note that this isn't required if you have CDH or MapR.

- /usr/hdp/<version>/hive/lib/hive-metastore.jar
- /usr/hdp/<version>/spark/lib/spark-assembly-1.2.1.2.3.X-hadoop2.6.0.2.3.X.jar
- /usr/hdp/<version>/hive/lib/hive-exec.jar

If any are missing, copy them over from one of your Hive or Spark nodes.

MapR-specific requirements

If you have MapR, your system must meet the following additional requirements.

- The MapR Client must be installed and added to the `$PATH` on all *non-MapR* nodes that will host the Dgraph, Studio, and the Transform Service (if different from Studio nodes). Note that the Client *isn't* required on these nodes if they host any MapR processes. For instructions on installing the MapR Client, see [Installing the MapR Client](#) in MapR's documentation.
- Pluggable authentication modules (PAMs) must be disabled for the installation.
- The `yarn.resourcemanager.hostname` property in `yarn-site.xml` must be set to the fully-qualified domain name (FQDN) of your YARN ResourceManager. For instructions on updating this property, see [Updating the YARN ResourceManager configuration on page 29](#).
- The directories `/user/HDFS_DP_USER_DIR/<bdd>` and `/user/HDFS_DP_USER_DIR/edp/data` must be either nonexistent or mounted with a volume. `HDFS_DP_USER_DIR` is defined in BDD's configuration file, and `<bdd>` is the name of the `bdd` user.
- The directories `/opt/mapr/zkdata` and `/opt/mapr/zookeeper/zookeeper-3.4.5/logs` must have their permissions set to `755`.
- If you want to store your Dgraph databases on MapR-FS, the directory defined by `DGRAPH_INDEX_DIR` in BDD's configuration file must be either nonexistent or mounted with a volume. Additionally, the MapR NFS service must be installed on all nodes that will host the Dgraph. For more information, see [HDFS on page 36](#).
- The required Spark, ZooKeeper, and Hive patches must be installed as described in [Applying the MapR patches on page 30](#).

Updating the YARN ResourceManager configuration

If you have MapR, you must set the `yarn.resourcemanager.hostname` property in `yarn-site.xml` to the fully-qualified domain name (FQDN) of your YARN ResourceManager.

Note that this isn't required if you have CDH or HDP.

The property is set to `0.0.0.0` by default. To update it, run the following command on the machine hosting MCS:

```
/opt/mapr/server/configure.sh -C <cldb_host>[:<cldb_port>][,<cldb_host>[:<cldb_port>]...]
-Z <zk_host>[:<zk_port>][,<zk_host>[:<zk_port>]...] [-RM <rm_host>] [-HS <hs_host>] [-L <logfile>]
[-N <cluster_name>]
```

Where:

- `<cldb_host>` and `<cldb_port>` are the FQDNs and ports of your container location database (CLDB) nodes
- `<zk_host>` and `<zk_port>` are the FQDNs and ports of your ZooKeeper nodes
- `<rm_host>` is the FQDN of your ResourceManager
- `<hs_host>` is the FQDN of your HistoryServer
- `<logfile>` is the log file `configure.sh` will write to
- `<cluster_name>` is the name of your MapR cluster

For more information on updating node configuration, see [configure.sh](#) in MapR's documentation.

Applying the MapR patches

If you have MapR, you must apply three sets of patches to your Hadoop cluster before installing BDD.

Note that this isn't required if you have CDH or HDP.

The patches are required to upgrade the versions of Spark, ZooKeeper, and Hive you have installed. Otherwise, BDD won't be able to work with them.

To apply the patches:

1. To apply the Spark patches, do the following on each Spark node:
 - (a) Download the following patches from <http://archive.mapr.com/releases/ecosystem-5.x/redhat/>:
 - `mapr-spark-master-1.6.1.201605311547-1.noarch.rpm`
 - `mapr-spark-1.6.1.201605311547-1.noarch.rpm`
 - `mapr-spark-historyserver-1.6.1.201605311547-1.noarch.rpm`
 - (b) Go to the directory you put the patches in and install each by running:

```
rpm -ivh <patch>
```

If the patches succeeded, your Spark nodes should contain the directory `/opt/mapr/spark/spark-1.6.1/`.

2. To apply the ZooKeeper patch, do the following on each ZooKeeper node:
 - (a) Download the following patch from <http://package.mapr.com/patches/releases/v5.1.0/redhat/>:
 - `mapr-patch-5.1.0.37549.GA-38290.x86_64.rpm`
 - (b) Apply the patch according to the instructions in MapR's [Patch Installation Guide](#).
 - (c) Restart ZooKeeper by running:

```
sudo service mapr-zookeeper restart
```

- (d) Verify that the patch succeeded by running:

```
echo status|nc <hostname> 5181|grep "Zookeeper version"
```

Where `<hostname>` is the hostname of the current ZooKeeper node.

The output should report ZooKeeper's current version as 1604, and not 1503:

```
Zookeeper version: 3.4.5-mapr-1604--1, built on 05/18/2016 14:50 GMT
```

3. To apply the Hive patches:
 - (a) Download the following patches from <http://archive.mapr.com/releases/ecosystem-5.x/redhat/> and copy them to each Hive node:
 - `mapr-hive-1.2.201606020917-1.noarch.rpm`
 - `mapr-hivemetastore-1.2.201606020917-1.noarch.rpm`
 - `mapr-hiveserver2-1.2.201606020917-1.noarch.rpm`
 - `mapr-hivewebhcat-1.2.201606020917-1.noarch.rpm`
 - (b) On each Hive node, go to the directory you put the patches in and install them by running:

```
rpm -Uvh <patch>
```

- (c) Go to MCS and restart the HiveServer 2, Hivemeta, and WebHcat services.
4. Update your MapR cluster's configuration by running the following command:

```
/opt/mapr/server/configure.sh -R
```

JDK requirements

BDD requires one of the following JDK versions installed in the same location on all nodes. If one of these is installed on your Hadoop nodes, you can copy it to your BDD nodes.

- [JDK 7u67+ x64](#)
- [JDK 8u45+ x64](#)

BDD requires a JDK that includes the HotSpot JVM, which must support the MD5 algorithm. These requirements will be met by any version you download using the links above, as long as you *don't* select a version from the JRockit Family.

Also, be sure to set the `$JAVA_HOME` environment variable on all nodes. If you have multiple versions of the JDK installed, be sure that this points to the correct one. If the path is set to or contains a symlink, the symlink must be identical on all other nodes.

Security options

The following sections describe methods for securing your BDD cluster.

Additional information on BDD security is available in the *Security Guide*.

[Kerberos](#)

[Sentry](#)

[TLS/SSL](#)

[HDFS data at rest encryption](#)

[Other security options](#)

Kerberos

The Kerberos network authentication protocol enables client/server applications to identify one another in a secure manner, even when communicating over an unsecured network.

In Kerberos terminology, individual applications are called *principals*. Each principal has a *keytab file*, which contains its *key*, or password. Keytab files enable principals to authenticate automatically, without human interaction. When one principal wants to communicate with another, it uses its keytab file to obtain a *ticket*. It then uses its ticket to gain access to the other principal.

Because Kerberos authentication uses strong encryption, it can work over unsecured networks. Additionally, tickets can be configured to expire after a set period of time to minimize risk should they become compromised.

You can configure BDD to use Kerberos authentication for its communications with Hadoop. This is required if Kerberos is already enabled in your Hadoop cluster, and strongly recommended for production environments in general. BDD supports integration with Kerberos 5+.

This procedure assumes you already have Kerberos installed on your system and configured for your Hadoop cluster.

To enable Kerberos:

1. Create the following directories in HDFS:

- `/user/<bdd user>`, where `<bdd user>` is the name of the bdd user.
- `/user/<HDFS_DP_USER_DIR>`, where `<HDFS_DP_USER_DIR>` is the value of `HDFS_DP_USER_DIR` in BDD's configuration file.

The owner of both directories must be the `bdd` user. Their group must be the HDFS super users group, which is defined by the `dfs.permissions.supergroup` configuration parameter. The default value is `supergroup`.

2. Add the `bdd` user to the `hive` group.
3. Add the `bdd` user to the `hdfs` group on all BDD nodes.
4. Create a BDD principal.

The primary component must be the name of the `bdd` user. The realm must be your default realm.

5. Generate a keytab file for the BDD principal and copy it to the install machine.

The name and location of this file are arbitrary. The installer will rename it `bdd.keytab` and copy it to all BDD nodes.

6. Copy the `krb5.conf` file from one of your Hadoop nodes to the install machine.

The location you put it in is arbitrary. The installer will copy it to `/etc` on all BDD nodes.

7. Install the `kinit` and `kdestroy` utilities on all BDD nodes.

These are required to enable ticket expiration.

8. If you have HDP, set the `hadoop.proxyuser.hive.groups` property in `core-site.xml` to `*`.

You can do this in Ambari.

You must also set the Kerberos-related properties in BDD's configuration file. For more information, see [Configuring BDD on page 62](#).

Sentry

Sentry provides role-based authorization in Hadoop clusters. Among other things, it can be used to restrict access to Hive data at a granular level.

Oracle strongly recommends using Sentry to protect your data from outside users. If you already have it set up in your Hadoop cluster, you must do a few things to enable BDD to work with it.



Note: The first two steps in this procedure are also required to enable Kerberos. If you've already done them, you can skip them.

To enable Sentry:

1. If you haven't already, create the following directories in HDFS:
 - `/user/<bdd user>`, where `<bdd user>` is the name of the bdd user.
 - `/user/<HDFS_DP_USER_DIR>`, where `<HDFS_DP_USER_DIR>` is the value of `HDFS_DP_USER_DIR` in BDD's configuration file.

The owner of both directories must be the `bdd` user. Their group must be the HDFS super users group, which is defined by the `dfs.permissions.supergroup` configuration parameter. The default value is `supergroup`.

2. If you haven't already, add the `bdd` user to the `hive` group.
3. Create a new role for BDD:

```
create role <BDD_role>;
grant all on server server1 to role <BDD_role>;
show grant role <BDD_role>;
grant role <BDD_role> to group hive;
```

TLS/SSL

BDD can be installed on Hadoop clusters secured with TLS/SSL.

TLS/SSL can be configured for specific Hadoop services to encrypt communication between them. If you have it enabled in Hadoop, you can enable it for BDD to encrypt its communications with your Hadoop cluster.

If your Hadoop cluster has TLS/SSL enabled, verify that your system meets the following requirements:

- Kerberos is enabled for both Hadoop and BDD. Note that this isn't required, but is strongly recommended. For more information, see [Kerberos on page 31](#).
- TLS/SSL is enabled in your Hadoop cluster for the HDFS, YARN, Hive, and/or Key Management Server (KMS) services.
- The KMS service is installed in your Hadoop cluster. You should have already done this as part of enabling TLS/SSL.

To enable BDD to run on a Hadoop cluster secured with TLS/SSL:

1. Export the public key certificates for all nodes running TLS/SSL-enabled HDFS, YARN, Hive, and/or KMS.

You can do this with the following command:

```
keytool -exportcert -alias <alias> -keystore <keystore_filename> -file <export_filename>
```

Where:

- `<alias>` is the certificate's alias.
 - `<keystore_filename>` is the absolute path to your keystore file. You can find this in Cloudera Manager, Ambari, or MCS.
 - `<export_filename>` is the name of the file you want to export the keystore to.
2. Copy the exported certificates to a single directory on the install machine.

The location of this directory is arbitrary, as you will define it in BDD's configuration file before installing. Don't remove this directory after installing, as you will use it if you have to update the certificates.

3. Verify that the password for `$JAVA_HOME/jre/lib/security/cacerts` is set to the default, `changeit`.

This is required by the installer. If it has been changed, be sure to set it back to the default.

When the installer runs, it imports the certificates to the custom truststore file, then copies the truststore to `$BDD_HOME/common/security/cacerts` on all BDD nodes.

HDFS data at rest encryption

HDFS data at rest encryption allows data to be stored in encrypted HDFS directories called *encryption zones*. All files within an encryption zone are transparently encrypted and decrypted on the client side, meaning decrypted data is never stored in HDFS.

If HDFS data at rest encryption is enabled in your Hadoop cluster, you must enable it for BDD. Before doing this, verify that your system meets the following requirements:

- The key trustee KMS and key trustee server are installed and configured in your Hadoop cluster. You should have already done this as part of enabling HDFS data at rest encryption.
- Kerberos is enabled for both Hadoop and BDD. Note that this isn't required, but is strongly recommended. For more information, see [Kerberos on page 31](#).
- TLS/SSL is enabled for both Hadoop and BDD. Note that this isn't required, but is strongly recommended. For more information, see [TLS/SSL on page 33](#).

To enable HDFS data at rest encryption for BDD:

1. Create an encryption zone in HDFS for your BDD files.
For instructions, refer to the documentation for your Hadoop distribution.
2. Grant the `bdd` user the `GENERATE_EEK` and `DECRYPT_EEK` privileges for the encryption and decryption keys.

You can do this in Cloudera Manager, Ambari, or MCS by adding the following properties to the KMS service's `kms-acls.xml` file. If you need help locating them, refer to your distribution's documentation.

```
<property>
  <name>key.acl.bdd_key.DECRYPT_EEK</name>
  <value>bdd,hdfs supergroup</value>
  <description>
    ACL for DECRYPT_EEK operations on key 'bdd_key'.
  </description>
</property>
<property>
  <name>key.acl.bdd_key.GENERATE_EEK</name>
  <value>bdd supergroup</value>
  <description>
    ACL for GENERATE_EEK operations on key 'bdd_key'.
  </description>
</property>
```

Be sure to replace `bdd` in the above code with the name of the `bdd` user and `supergroup` with the name of the HDFS super users group, which is defined by the `dfs.permissions.supergroup` configuration parameter.

Also note that the `hdfs` user is included in the value of the `DECRYPT_EEK` property. This is required if you're storing your Dgraph databases on HDFS, but can be omitted otherwise. For more information, see [Installing the HDFS NFS Gateway service on page 38](#).

Other security options

You can further protect BDD by installing it behind a firewall and enabling TLS/SSL on Studio's outward-facing ports.

Firewalls

Oracle recommends using a firewall to protect your network and BDD cluster from external entities. A firewall limits traffic into and out of your network, creating a secure barrier around it. It can consist of a combination of software and hardware, including routers and dedicated gateway machines.

There are multiple types of firewalls, so be sure to choose one that suits your resources and needs. One option is to use a reverse proxy server as part of your firewall, which you can configure after installing BDD. For instructions, see [Using Studio with a Reverse Proxy on page 87](#).

TLS/SSL in Studio

You can enable TLS/SSL on Studio's outward-facing ports in one or both of the following ways:

- Enable encryption through WebLogic Server. You can do this by setting `WLS_SECURE_MODE` to `TRUE` in BDD's configuration file.

This method activates WebLogic's default demo keystores, which you should replace with your own certificates after deployment. For more information, see [Replacing certificates on page 84](#).

- Set up a reverse-proxy server. For instructions on how to do this, see [About reverse proxies on page 88](#).

Be aware that these methods don't enable encryption on the inward-facing port on which the Dgraph Gateway listens for requests from Studio.

Component database requirements

The Dgraph, Studio, and the Workflow Manager Service all require databases. The following sections describe the requirements for each.

[Dgraph database requirements](#)

[Studio database requirements](#)

[Workflow Manager Service database requirements](#)

[Sample commands for production databases](#)

Dgraph database requirements

The Dgraph stores the data sets it queries in a set of databases. For high availability, these can be stored on HDFS/MapR-FS or a shared NFS. They can also be stored on the local disk for a non-HA option.

The location you choose determines the database requirements, as well as where the Dgraph will be installed and its behavior. For more information, see:

- [HDFS on page 36](#)
- [NFS on page 38](#)



Note: You can install with pre-existing BDD-formatted databases if you have any you want to use. To do this, put them in the directory you want to store your databases in and point BDD's configuration file to it. For more information, see [Configuring BDD on page 62](#).

Regardless of where you put your Dgraph databases, you must increase the allowed numbers of open file descriptors and processes on all nodes in the cluster, including Hadoop nodes. For more information, see [Increasing the numbers of open file descriptors and processes on page 38](#).

HDFS

Storing your databases on HDFS provides increased high availability for the Dgraph—the contents of the databases are distributed across multiple nodes, so the Dgraph can continue to process queries if a node goes down. It also increases the amount of data your databases can contain.



Note: This information also applies to MapR-FS.

To store your databases on HDFS, your system must meet the following requirements:

- The HDFS DataNode service must be running on all nodes that will host the Dgraph. For best performance, this should be the only Hadoop service running on your Dgraph nodes. In particular, the Dgraph shouldn't be co-located with Spark, as both services require a lot of resources.

If you have to co-locate the Dgraph with Spark or any other Hadoop services, you should use cgroups to isolate resources for it. For more information, see [Setting up cgroups on page 37](#).

- For best performance, configure short-circuit reads in HDFS. This enables the Dgraph to access the local database files directly, rather than using the DataNode's network sockets to transfer the data. For instructions, refer to the documentation for your Hadoop distribution.
- The `bdd` user must have **read** and **write** permissions for the HDFS directory where the databases will be stored. Be sure to set this on all Dgraph nodes.
- If you have HDFS data at rest encryption enabled in Hadoop, you must store your databases in an encryption zone. For more information, see [HDFS data at rest encryption on page 34](#).
- If you decide to not use the default HDFS mount point (the local directory where the Dgraph mounts the HDFS root directory), make sure the one you use is empty and has **read**, **write**, and **execute** permissions for the `bdd` user. This must be set on all Dgraph nodes.
- Be sure to set the `DGRAPH_USE_MOUNT_HDFS` property in BDD's configuration file to `TRUE`.
- To enable the Dgraph to access its databases in HDFS, you must install the HDFS NFS Gateway (called MapR NFS in MapR). For more information, see [Installing the HDFS NFS Gateway service on page 38](#).

Setting up cgroups

Control groups, or cgroups, are a Linux kernel feature that enable you to allocate resources like CPU time and system memory to specific processes or groups of processes. If you need to host the Dgraph on nodes running Spark, you should use cgroups to ensure sufficient resources are available to it.



Note: Installing the Dgraph on Spark nodes is not recommended and should only be done if absolutely necessary.

To do this, you enable cgroups in Hadoop and create one for YARN that limits the amounts of CPU and memory it can consume. You then create a separate cgroup for the Dgraph.

To set up cgroups:

- 1 If your system doesn't currently have the `libcgroup` package, install it as root.
This creates `/etc/cgconfig.conf`, which is used to configure cgroups.

- 2 Enable the `cgconfig` service to run automatically:

```
chkconfig cgconfig on
```

- 3 Create a cgroup for YARN. You must do this within Hadoop. For instructions, refer to the documentation for your Hadoop distribution.

The YARN cgroup should limit the amounts of CPU and memory allocated to all YARN containers. The appropriate limits to set depend on your system and the amount of data you will process. At a minimum, you should reserve the following for the Dgraph:

- 5GB of RAM
- 2 CPU cores

The number of CPU cores YARN is allowed to use must be specified as a percentage. For example, on a quad-core machine, YARN should only get two cores, or 50%. On an eight-core machine, YARN could get up to six of them, or 75%. When setting this amount, remember that allocating more cores to the Dgraph will boost its performance.

- 4 Create a cgroup for the Dgraph by adding the following to `cgconfig.conf`:

```
# Create a Dgraph cgroup named "dgraph"
group dgraph {
# Specify which users can edit this group
  perm {
    admin {
      uid = $BDD_USER;
    }
    # Specify which users can add tasks for this group
    task {
      uid = $BDD_USER;
    }
  }
# Set the memory and swap limits for this group
  memory {
    # Sets memory limit to 10GB
    memory.limit_in_bytes = 10000000000;

    # Sets memory + swap limit to 12GB
    memory.memsw.limit_in_bytes = 12000000000;
  }
}
```

Where `$BDD_USER` is the name of the `bdduser`.



Important: The values given for `memory.limit_in_bytes` and `memory.memsw.limit_in_bytes` above are the *absolute minimum* requirements. You should use higher values, if possible.

- 5 Restart `cfconfig` to enable your changes.

Installing the HDFS NFS Gateway service

If you want to store your Dgraph databases on HDFS, you must install the HDFS NFS Gateway service (called the MapR NFS service in MapR).

The NFS Gateway service enables client applications to mount HDFS as part of the local file system. Clients can then search for, read from, and write to HDFS files as if they were stored locally. In the context of BDD, the NFS Gateway allows the Dgraph to access its databases when they're stored in HDFS.

To enable this for BDD, the NFS Gateway service must be installed on all Dgraph nodes. For instructions on installing it, refer to the documentation for your Hadoop distribution.

The NFS Gateway service must be running when you install BDD. The installer will automatically detect it at runtime and add the following properties to BDD's configuration file:

```
NFS_GATEWAY_SERVERS=<list of NFS Gateway nodes>
DGRAPH_USE_NFS_MOUNT=TRUE
```

After installing, the Dgraph will mount HDFS via the NFS Gateway when it starts.

NFS

If you don't want to store your databases on HDFS, you can keep them on an NFS.

NFS (Network File System) is a distributed file system that enables clients to access data stored on a separate machine over the network. Storing your Dgraph databases on one ensures that all Dgraph instances will be able to access them.

If you want to use NFS, be sure that the NFS server is properly set up and that all Dgraph nodes have read and write access to it. You should also ensure that it contains enough storage space for the amount of data you plan on processing. You may want to keep your databases on a separate partition from system files and any other data on the NFS.

Increasing the numbers of open file descriptors and processes

Regardless of where you put your Dgraph databases, you must increase the maximum numbers of open file descriptors and processes, or the Dgraph may crash during processing.

The number of open file descriptors should have hard and soft limits of 65536, at a minimum. The number of open processes should have a soft limit of 65536 and an unlimited hard limit.

To set these, do the following on *each node in your cluster* (including Hadoop nodes):

- 1 Create a process limit configuration file for the `bdd` user named `/etc/security/limits.d/<bdd>.conf`, where `<bdd>` is the name of the `bdduser`.
- 2 Open `<bdd>.conf` and add the following:

```
<bdd> soft nofile 65536
```

<bdd>	hard	nofile	65536
<bdd>	soft	nproc	65536
<bdd>	hard	nproc	unlimited

Where <bdd> is the name of the bdd user.

- 3 Save and close the file.
- 4 Log out and then log back in so that your changes will take effect.
- 5 Run the following to verify your changes:

```
ulimit -n
```

The above command should output 65536.

Studio database requirements

Studio requires a relational database to store configuration and state, including component configuration, user permissions, and system settings. If you install with multiple Studio instances, all of them must be connected to the same database.

Studio supports the following database types:

- Oracle 11g
- Oracle 12c 12.1.0.1.0+
- MySQL 5.5.3+

If you're installing BDD in a production environment, you must create the following:

- A database of one of the types listed above.
- A database username and password.
- An empty schema. The name of this is arbitrary.

Note that BDD doesn't currently support database migration. If you decide to switch to a different type of database later on, you must reinstall BDD with a new database instance. If you're installing BDD in a non-production environment with the QuickStart option, you *must* use a MySQL database named `studio`. For more information, see [QuickStart Installation on page 48](#).

You can optionally use a clustered database configuration. For clustering, Oracle 11g uses RAC and MySQL has MySQL Cluster. Refer to the documentation for your database system for details on setting up a clustered configuration.

Additionally:

- You must install the database client on the install machine. For MySQL, this should be MySQL client. For Oracle databases, this should be Oracle Database Client, installed with a type of Administrator. Note that the Instant Client is not supported.
- If you have a MySQL database, you must set UTF-8 as the default character set.
- If you have an Oracle database, you must set the `ORACLE_HOME` environment variable to the directory one level above the `/bin` directory that the `sqlplus` executable is located in. For example, if the `sqlplus` executable is located in `/u01/app/oracle/product/11/2/0/dbhome/bin`, you should set `ORACLE_HOME` to `/u01/app/oracle/product/11/2/0/dbhome`. Note that this is different from the `ORACLE_HOME` property in BDD's configuration file.

Sample commands for creating Oracle and MySQL database users and schemas are available in [Sample commands for production databases on page 41](#).

Studio database requirements in demo environments

In demo environments, Studio supports Hypersonic (HSQL) databases in addition to the types listed above. Hypersonic is an embedded database that runs inside the JVM. It is useful for getting Studio up and running quickly, but can't be used in a production environment due to performance issues and its inability to support multiple Studio nodes.



Note: The Connector Service and the Component Registry *don't* support Hypersonic databases, even in demo environments.

If you want to use a Hypersonic database, the installer will create it for you. You can enable this in BDD's configuration file.



Important: If you install in a demo environment with a Hypersonic database and later decide to scale up to a production environment, you must reinstall BDD with one of the supported MySQL or Oracle databases listed above.

Workflow Manager Service database requirements

The Workflow Manager Service requires a relational database to store state information.

Like Studio, the Workflow Manager Service supports the following types of databases:

- Oracle 11g
- Oracle 12c 12.1.0.1.0+
- MySQL 5.5.3+

You must create the following for the Workflow Manager Service:

- A database of one of the types listed above. Note that this must be separate from the Studio database.
- A database username and password.
- An empty schema. The name of this is arbitrary.

Note that BDD doesn't currently support database migration. If you decide to switch to a different type of database later on, you must reinstall BDD with a new database instance. If you're installing BDD in a non-production environment with the QuickStart option, you *must* use a MySQL database named `workflow`. For more information, see [QuickStart Installation on page 48](#).

Additionally:

- You must install the database client on the install machine. For MySQL, this should be MySQL client. For Oracle databases, this should be Oracle Database Client, installed with a type of Administrator. Note that the Instant Client is not supported.
- If you have a MySQL database, you must set UTF-8 as the default character set.
- If you have an Oracle database, you must set the `ORACLE_HOME` environment variable to the directory one level above the `/bin` directory that the `sqlplus` executable is located in. For example, if the `sqlplus` executable is located in `/u01/app/oracle/product/11/2/0/dbhome/bin`, you should set

ORACLE_HOME to /u01/app/oracle/product/11/2/0/dbhome. Note that this is different from the ORACLE_HOME property in BDD's configuration file.

Sample commands for creating Oracle and MySQL database users and schemas are available in [Sample commands for production databases on page 41](#).

Sample commands for production databases

Below are sample commands you can use to create users and schemas for Oracle and MySQL databases. You are not required to use these exact commands when setting up your component databases—these are just examples to help get you started.

Oracle database

You can use the following commands to create a user and schema for an Oracle 11g or 12c database.

```
CREATE USER <username> PROFILE "DEFAULT" IDENTIFIED BY <password> DEFAULT TABLESPACE "USERS"
TEMPORARY TABLESPACE "TEMP" ACCOUNT UNLOCK;
GRANT CREATE PROCEDURE TO <username>;
GRANT CREATE SESSION TO <username>;
GRANT CREATE SYNONYM TO <username>;
GRANT CREATE TABLE TO <username>;
GRANT CREATE VIEW TO <username>;
GRANT UNLIMITED TABLESPACE TO <username>;
GRANT CONNECT TO <username>;
GRANT RESOURCE TO <username>;
```

MySQL database

You can use the following commands to create a user and schema for a MySQL database.



Note: MySQL databases must use UTF-8 as the default character encoding.

```
create user '<username>'@'%' identified by '<password>';
create database <database name> default character set utf8 default collate utf8_general_ci;
grant all on <database name>.* to '<username>'@'%' identified by '<password>' with grant option;
flush privileges;
```

Supported Web browsers

Studio supports the following Web browsers:

- Internet Explorer 11 (compatibility mode is not supported)
- Microsoft Edge
- Firefox ESR
- Chrome for Business
- Safari 9+ (for mobile)

Screen resolution requirements

BDD has the following screen resolution requirements:

- Minimum: 1366x768
- Recommended: 1920x1080

Studio support for iPad

You can use the Safari Web browser on an iPad running iOS 7+ to sign in to Studio and view projects. You cannot use an iPad to create, configure, or export projects.

While the iPad can support most component functions, the component export option is disabled.

Part II

Installing Big Data Discovery



Chapter 3

Prerequisite checklist

Before installing, run through the following checklist to verify you've satisfied all prerequisites.

For more information on each prerequisite, refer to the relevant section in [Prerequisites on page 14](#).




Note: BDD includes a script called the *prerequisite checker* that verifies whether your system meets all install requirements. You can run this script after you update BDD's configuration file. For more information, see [Running the prerequisite checker on page 70](#).

Prerequisite	Description
Hardware	<p>Minimum requirements:</p> <ul style="list-style-type: none">• WebLogic nodes: quad-core CPU• Dgraph nodes: dual-core CPU <p>Note that these are the minimum amounts required to install BDD. A full-scale installation will require more.</p>
Memory	<p>Minimum requirements:</p> <ul style="list-style-type: none">• Managed Servers: 16GB (5GB for WebLogic Server and 11GB for the Transform Service)• Dgraph nodes: 5GB (excluding requirements for HDFS, if applicable)• YARN cluster: 16GB (combined) <p>Note that these are the minimum amounts required to install BDD. A full-scale installation will require more.</p>
Disk space	<p>Minimum requirements:</p> <ul style="list-style-type: none">• 30GB in <code>ORACLE_HOME</code> on all BDD nodes• 20GB in <code>TEMP_FOLDER_PATH</code> on all BDD nodes• 10GB in <code>INSTALLER_PATH</code> on the install machine• 512MB swap space on the install machine and all Managed Servers• 39GB virtual memory on all Transform Service nodes <p>Note that these are the minimum amounts required to install BDD. A full-scale installation will require more.</p>
Network	<p>The hostname of each BDD machine can be externally resolved and accessed using the machine's IP address.</p>

Prerequisite	Description																																				
Operating system	<ul style="list-style-type: none">• OEL 6.4+, 7.1• RHEL 6.4+, 7.1																																				
Linux utilities	<ul style="list-style-type: none">• /bin:<table><tr><td>basename</td><td>date</td><td>more</td><td>true</td></tr><tr><td>cat</td><td>dd</td><td>rm</td><td></td></tr><tr><td>chgrp</td><td>df</td><td>sed</td><td></td></tr><tr><td>chown</td><td>mkdir</td><td>tar</td><td></td></tr></table>• /usr/bin:<table><tr><td>awk</td><td>expr</td><td>netcat</td><td>tail</td><td>which</td></tr><tr><td>cksum</td><td>gzip</td><td>perl</td><td>tr</td><td></td></tr><tr><td>cut</td><td>head</td><td>printf</td><td>unzip</td><td></td></tr><tr><td>dirname</td><td>id</td><td>sudo</td><td>wc</td><td></td></tr></table>• Perl 5.10+ with multithreading• The Mail::Address and XML::Parser Perl modules• The default umask set to 022• curl 7.19.7+ (with support for --tlsv1.2 and --negotiate) on all nodes that will host Studio• Network Security Services (NSS) 3.16.1+ and nss-devel on all nodes that will host Studio• tty disabled for sudo• Apache Ant 1.7.1+ installed and added to the PATH on all nodes, including Hadoop nodes	basename	date	more	true	cat	dd	rm		chgrp	df	sed		chown	mkdir	tar		awk	expr	netcat	tail	which	cksum	gzip	perl	tr		cut	head	printf	unzip		dirname	id	sudo	wc	
basename	date	more	true																																		
cat	dd	rm																																			
chgrp	df	sed																																			
chown	mkdir	tar																																			
awk	expr	netcat	tail	which																																	
cksum	gzip	perl	tr																																		
cut	head	printf	unzip																																		
dirname	id	sudo	wc																																		
OS user	<p>The following are set for the bdd user:</p> <ul style="list-style-type: none">• Passwordless sudo and SSH on all nodes, including Hadoop nodes• Passwordless SSH on all nodes, including Hadoop nodes• Bash set as the default shell• Permission to create the ORACLE_HOME directory on all nodes																																				

Prerequisite	Description
Hadoop	<ul style="list-style-type: none"> Distributions: <ul style="list-style-type: none"> CDH 5.7.x (min. 5.7.1), 5.8.x, 5.9.x, 5.10.x, 5.11.x HDP 2.4.x (min. 2.4.2) MapR 5.1+ Components: <ul style="list-style-type: none"> Cluster manager: Cloudera Manager, Ambari, or MCS ZooKeeper HDFS HCatalog Hive Spark on YARN Hue YARN Spark on YARN, YARN, and HDFS are on all Data Processing nodes YARN configuration has been updated
HDP-specific requirements	<ul style="list-style-type: none"> The required client libraries are on the install machine The <code>hive-metastore</code>, <code>spark-assembly</code>, and <code>hive-exec</code> JARs are on all Hadoop nodes
MapR-specific requirements	<ul style="list-style-type: none"> The MapR Client is installed on all non-MapR nodes that will host the Dgraph, Studio, and the Transform Service PAMs are disabled The YARN Resource Manager IP is configured correctly on the machine hosting MCS The directories <code>/user/HDFS_DP_USER_DIR/</code> and <code>/user/HDFS_DP_USER_DIR/edp/data</code> are either nonexistent or mounted with a volume The permissions for the <code>/opt/mapr/zkdata</code> and <code>/opt/mapr/zookeeper/zookeeper-3.4.5/logs</code> directories are set to 755 The required Spark, ZooKeeper, and Hive patches have been applied

Prerequisite	Description
JDK	<ul style="list-style-type: none"> • JDK 7u67+ • JDK 8u45+ • The installed JDK contains the HotSpot JVM, which supports MD5 • <code>\$JAVA_HOME</code> set on all nodes
Kerberos	<ul style="list-style-type: none"> • <code>/user/<bdd_user></code> and <code>/user/<HDFS_DP_USER_DIR></code> created in HDFS • bdd user is a member of the <code>hive</code> and <code>hdfs</code> groups • bdd principal and keytab file have been generated • bdd keytab file and <code>krb5.conf</code> are on the install machine • <code>kinit</code> and <code>kdestroy</code> are installed on BDD nodes • <code>core-site.xml</code> has been updated (HDP only)
Sentry	<ul style="list-style-type: none"> • <code>/user/<bdd_user></code> and <code>/user/<HDFS_DP_USER_DIR></code> in HDFS • bdd user is a member of the <code>hive</code> group • BDD role
TLS/SSL	<ul style="list-style-type: none"> • Kerberos enabled for BDD and Hadoop • KMS is installed and configured • TLS/SSL enabled in Hadoop for HDFS, YARN, Hive, and/or KMS • The public key certificates for all TLS/SSL enabled services (HDFS, YARN, Hive, and/or KMS) have been exported and copied to the install machine • The password for <code>cacerts</code> is set to the default (<code>chageit</code>)
HDFS data at rest encryption	<ul style="list-style-type: none"> • Kerberos and TLS/SSL enabled for BDD and Hadoop • The key trustee KMS and key trustee server installed and configured in Hadoop • HDFS data at rest encryption enabled in Hadoop • A BDD encryption zone has been created in HDFS • The bdd user has <code>GENERATE_EEK</code> and <code>DECRYPT_EEK</code> privileges for the encryption and decryption keys

Prerequisite	Description
Dgraph databases	<ul style="list-style-type: none"> If stored on HDFS: <ul style="list-style-type: none"> The HDFS DataNode service is on all Dgraph nodes cgroups are set up, if necessary (Optional) Short-circuit reads are enabled in HDFS The <code>bdd</code> user has read and write permissions to the databases directory in HDFS If using a non-default mount point, it's empty and the <code>bdd</code> user has read, write, and execute permissions for it You installed the HDFS NFS Gateway service If stored on an NFS: <ul style="list-style-type: none"> The NFS is set up All Dgraph nodes can write to it The number of open file descriptors is set to 65536 on all Dgraph nodes
Studio database	<p>The following have been created:</p> <ul style="list-style-type: none"> One of the following databases: <ul style="list-style-type: none"> Oracle 11g Oracle 12c 12.1.0.1.0+ MySQL 5.5.3+ A database username and password An empty schema <p> Note: You can also configure the installer to create an HSQL database for you, although this isn't supported for production environments.</p>
Web browser	<ul style="list-style-type: none"> Firefox ESR Internet Explorer 11 (compatibility mode not supported) Chrome for Business Safari 9+ (for mobile)



Chapter 4

QuickStart Installation

The BDD installer includes a `quickstart` option, which installs the software on a single machine with default configuration suitable for a demo environment. You can use `quickstart` to install BDD quickly and easily, without having to worry about setting it up yourself.



Important: Single-node installations can only be used for demo purposes; you can't host a production environment on a single machine. If you want to install BDD in a production environment, see [Cluster Installation on page 57](#).

Before you can install BDD with `quickstart`, you must satisfy all of the prerequisites described in [Prerequisites on page 14](#), with a few exceptions:

- You must use CDH. HDP and MapR aren't supported.
- You must have MySQL databases for Studio and the Workflow Manager Service. These must be named `studio` and `workflow`, respectively.
- You can't have Kerberos installed.
- You can't have TLS/SSL or HDFS data at rest encryption enabled in Hadoop.
- You can't use any existing Dgraph databases.



Note: If you want to install BDD on a single machine but need more control and flexibility than `quickstart` offers, see [Single-Node Installation on page 50](#).

[Installing BDD with quickstart](#)

Installing BDD with quickstart

Once you've satisfied all of BDD's prerequisites, you can download and install the software.

Before installing, verify that:

- CDH is installed.
- You satisfied all requirements described in [Prerequisites on page 14](#).
- The `bdd` user meets all requirements described in [OS user requirements on page 24](#).
- You set up MySQL databases (including usernames, passwords, and schemas) for Studio and the Workflow Manager Service.
- The following Hadoop components are running:
 - Cloudera Manager
 - ZooKeeper
 - HDFS

- Hive
- Spark on YARN
- YARN
- Hue

To install BDD with `quickstart`:

1. On your machine, create a new directory or choose an existing one to be the installation source directory.
This directory must contain at least 10GB of free space.
2. Within the installation source directory, create a new directory named `packages`.
3. Download the BDD media pack from the [Oracle Software Delivery Cloud](#).
Be sure to download all packages in the media pack. Make a note of each file's part number, as you will need this to identify it later.
4. Move the BDD installer, BDD binary, and WebLogic Server packages from the download location to the `packages` directory.
5. Rename the first BDD binary package `bdd1.zip` and the second `bdd2.zip`.
This ensures that the installer will recognize them.
6. Extract the WebLogic Server package.
This creates a file called `fmw_12.1.3.0.0_wls.jar`, which contains the WebLogic Server installer.
7. Navigate back to the installation source directory and extract the BDD installer package:

```
unzip packages/<BDD_installer_package>.zip
```


This creates a new directory called `installer`, which contains the install script and other files it requires.
8. Go to the `installer` directory and run:

```
./setup.sh --quickstart
```
9. Enter the following when prompted:
 - The username and password for Cloudera Manager.
 - A username and password for the WebLogic Server admin. The password must contain at least 8 characters, one of which must be a number, and can't begin with a number.
 - The username and password for the Studio and Workflow Manager Service databases.
 - The password for the Studio admin. This must contain at least 8 characters, one of which must be a non-alphanumeric character.

If the script succeeded, BDD is now installed under the current directory and ready for you to begin working with it. See [Post-Installation Tasks on page 76](#) to learn more about your installation and how to verify it.

If the script failed, see [Troubleshooting a Failed Installation on page 72](#).



Chapter 5

Single-Node Installation

If you want to demo BDD before committing to a full-cluster installation, you can install it on a single node. This gives you the chance to learn more about the software and see how it performs on a smaller scale. The following sections describe how to get BDD running on your machine quickly and easily.



Important: Single-node installations can only be used for demo purposes; you can't host a production environment on a single machine. If you want to install BDD in a production environment, see [Cluster Installation on page 57](#).

[Installing BDD on a single node](#)

[Configuring a single-node installation](#)

Installing BDD on a single node

Once you've satisfied all of BDD's prerequisites, you can download and install the software.

Before installing, verify that:

- You satisfied all requirements described in [Prerequisites on page 14](#).
- The `bdd` user meets the requirements described in [OS user requirements on page 24](#).
- The Studio and Workflow Manager Service databases (including their usernames, passwords, and schemas) are set up.
- The following Hadoop components are running:
 - Cloudera Manager/Ambari/MCS
 - ZooKeeper
 - HDFS
 - Hive
 - Spark on YARN
 - YARN
 - Hue

To install BDD:

1. On your machine, create a new directory or choose an existing one to be the installation source directory.
This directory must contain at least 10GB of free space.
2. Within the installation source directory, create a new directory named `packages`.

3. Download the BDD media pack from the [Oracle Software Delivery Cloud](#).
Be sure to download all packages in the media pack. Make a note of each file's part number, as you will need this to identify it later.
4. Move the BDD installer, BDD binary, and WebLogic Server packages from the download location to the `packages` directory.
5. Rename the first BDD binary package `bdd1.zip` and the second `bdd2.zip`.
This ensures that the installer will recognize them.
6. Extract the WebLogic Server package.
This creates a file called `fmw_12.1.3.0.0_wls.jar`, which contains the WebLogic Server installer.
7. Navigate back to the installation source directory and extract the BDD installer package:

```
unzip packages/<BDD_installer_package>.zip
```


This creates a new directory called `installer`, which contains the install script and other files it requires.
8. Open BDD's configuration file, `bdd.conf`, in a text editor and update the Required Settings section.
See [Configuring a single-node installation on page 52](#) for instructions.
9. Run the prerequisite checker to verify whether your system meets all install requirements.
See [Running the prerequisite checker on page 70](#) for instructions.
10. Go to the `installer` directory and run:

```
./setup.sh
```
11. Enter the following when prompted:
 - The username and password for your cluster manager.
 - A username and password for the WebLogic Server admin. The password must contain at least 8 characters, one of which must be a number, and can't begin with a number.
 - The username and password for the database.
 - The password for the Studio admin. This must contain at least 8 characters, one of which must be a non-alphanumeric character.

If the script succeeded, BDD is now installed on your machine and ready for you to begin working with it. See [Post-Installation Tasks on page 76](#) to learn more about your installation and how to verify it.

If the script failed, see [Troubleshooting a Failed Installation on page 72](#).

Configuring a single-node installation

The table below describes the properties you should set for a single-node installation. You can modify `bdd.conf` in any text editor.


Keep the following in mind when editing the file:

- The accepted values for some properties are case-sensitive and must be entered exactly as they appear in the table below.
- All hostnames must be Fully Qualified Domain Names (FQDNs).

- Each port setting must have a unique value.
- Some of the directories defined in `bdd.conf` have location requirements. These are specified below.

Configuration property	Description
ORACLE_HOME	<p>The path to the directory BDD will be installed in. This must not exist and the system must contain at least 30GB of free space to create this directory. Additionally, its parent directories' permissions must be set to either 755 or 775.</p> <p>Note that this setting is different from the <code>ORACLE_HOME</code> environment variable required by the Studio database.</p>
ORACLE_INV_PTR	<p>The absolute path to the Oracle inventory pointer file, which the installer will create when it runs. This can't be located in the <code>ORACLE_HOME</code> directory.</p> <p>If you have any other Oracle software products installed, this file will already exist. Update this property to point to it.</p>
INSTALLER_PATH	The absolute path to the installation source directory.
DGRAPH_INDEX_DIR	<p>The absolute path to the Dgraph databases. This directory shouldn't be located under <code>ORACLE_HOME</code>, or it will be deleted.</p> <p>The script will create this directory if it doesn't currently exist. If you're installing with existing databases, set this property to their parent directory.</p>
HADOOP_UI_HOST	The hostname of the machine running your Hadoop manager (Cloudera Manager, Ambari, or MCS). Set this to your machine's hostname.

Configuration property	Description
STUDIO_JDBC_URL	<p>The JDBC URL for your Studio database, which Studio requires to connect to it.</p> <p>There are three templates for this property. Copy the template that corresponds to your database type to STUDIO_JDBC_URL and update the URL to point to your database.</p> <ul style="list-style-type: none"> If you have a MySQL database, use the first template and update the URL as follows: <pre>jdbc:mysql://<database hostname>:<port number> /<database name>?useUnicode=true&characterEncoding=UTF-8&useFastDateParsing=false</pre> If you have an Oracle database, use the first template and update the URL as follows: <pre>jdbc:oracle:thin: @<database hostname>:<port number>:<database SID></pre> If you're not installing on a production environment and want the installer to create a Hypersonic database for you, use the third template. The script will create the database for you in the location defined by the URL.
WORKFLOW_MANAGER_JDBC_URL	<p>The JDBC URL for the Workflow Manager Service database.</p> <p>There are two templates for this property. Copy the template that corresponds to your database type to WORKFLOW_MANAGER_JDBC_URL and update the URL to point to your database.</p> <ul style="list-style-type: none"> For MySQL databases, use the first template and update the URL as follows: <pre>jdbc:mysql://<database hostname>:<port number> /<database name>?useUnicode=true&characterEncoding=UTF-8&useFastDateParsing=false</pre> For Oracle databases, use the first template and update the URL as follows: <pre>jdbc:oracle:thin: @<database hostname>:<port number>:<database SID></pre> <p>Note that BDD doesn't currently support database migration. After deployment, the only ways to change to a different database are to reconfigure the database itself or reinstall BDD.</p>
INSTALL_TYPE	<p>Determines the installation type according to your hardware and Hadoop distribution. Set this to one of the following:</p> <ul style="list-style-type: none"> CDH HW MAPR

Configuration property	Description
JAVA_HOME	<p>The absolute path to the JDK install directory. This should have the same value as the <code>\$JAVA_HOME</code> environment variable.</p> <p>If you have multiple versions of the JDK installed, be sure that this points to the correct one.</p>
TEMP_FOLDER_PATH	The temporary directory used by the installer. This must exist and contain at least 20GB of free space.
HADOOP_UI_PORT	The port number for the Hadoop manager.
HADOOP_UI_CLUSTER_NAME	The name of your Hadoop cluster, which is listed in the manager. Be sure to replace any spaces in the cluster name with <code>%20</code> .
HUE_URI	The hostname and port for Hue, in the format <code><hostname>:<port></code> . This property is only required for HDP.
HADOOP_CLIENT_LIB_PATHS	<p>A comma-separated list of the absolute paths to the Hadoop client libraries.</p> <p> Note: You only need to set this property before installing if you have HDP or MapR. For CDH, the installer will download the required libraries and set this property automatically. This requires an internet connection. If the script is unable to download the libraries, it will fail; see Failure to download the Hadoop client libraries on page 73 for instructions on solving this issue.</p> <p>To set this property, copy the template that corresponds to your Hadoop distribution to <code>HADOOP_CLIENT_LIB_PATHS</code> and update the paths to point to the libraries you copied to the install machine. Be sure to replace all instances of <code><UNZIPPED_XXX_BASE></code> with the absolute path to the correct library.</p> <p>Don't change the order of the paths in the list as they <i>must</i> be specified as they appear.</p>
HADOOP_CERTIFICATES_PATH	<p>Only required for Hadoop clusters with TLS/SSL enabled. The absolute path to the directory on the install machine where you put the certificates for HDFS, YARN, Hive, and the KMS.</p> <p>Don't remove this directory after installing, as you will use it if you have to update the certificates.</p>
ENABLE_KERBEROS	Enables Kerberos. If you have Kerberos 5+ installed, set this value to <code>TRUE</code> ; if not, set it to <code>FALSE</code> .
KERBEROS_PRINCIPAL	<p>The name of the BDD principal. This should include the name of your domain; for example, <code>bdd-service@EXAMPLE.COM</code>.</p> <p>This property is only required if <code>ENABLE_KERBEROS</code> is set to <code>TRUE</code>.</p>

Configuration property	Description
KERBEROS_KEYTAB_PATH	The absolute path to the BDD <code>keytab</code> file. This property is only required if <code>ENABLE_KERBEROS</code> is set to <code>TRUE</code> .
KRB5_CONF_PATH	The absolute path to the <code>krb5.conf</code> file. This property is only required if <code>ENABLE_KERBEROS</code> is set to <code>TRUE</code> .
ADMIN_SERVER	The hostname of the WebLogic Admin Server. This will default to your machine's hostname, so you don't need to set it.
MANAGED_SERVERS	The hostname of the Managed Server. Leave this set to <code>\${ADMIN_SERVER}</code> .
DGRAPH_SERVERS	The Dgraph hostname. Leave this set to <code>\${ADMIN_SERVER}</code> .
DGRAPH_THREADS	The number of threads the Dgraph starts with. This will default to the number of cores your machine has minus 2, so you don't need to set it.
DGRAPH_CACHE	The size of the Dgraph cache, in MB. This will default to either 50% of your RAM or the total amount of free memory minus 2GB (whichever is larger), so you don't need to set it.
ZOOKEEPER_INDEX	The index of the Dgraph cluster in the ZooKeeper ensemble, which ZooKeeper uses to identify it.
HDFS_DP_USER_DIR	The location within the HDFS <code>/user</code> directory that stores the sample files created when Studio users export data. The installer will create this directory if it doesn't already exist. The name of this directory can't include spaces or slashes (/).
YARN_QUEUE	The YARN queue Data Processing jobs are submitted to.
HIVE_DATABASE_NAME	The name of the Hive database that stores the source data for Studio data sets.

Configuration property	Description
SPARK_ON_YARN_JAR	<p>The absolute path to the Spark on YARN JAR on your Hadoop nodes. This will be added to the CLI classpath.</p> <p>There are two templates for this property. Copy the value of the template that corresponds to your Hadoop distribution to <code>SPARK_ON_YARN_JAR</code> and update its value as follows:</p> <ul style="list-style-type: none"> For CDH, use the first template. This should be the absolute path to <code>spark-assembly.jar</code>. For HDP, use the second template. This should be the absolute paths to <code>hive-metastore.jar</code>, <code>hive-exec.jar</code> and <code>spark-assembly.jar</code>, separated by a colon: <pre><path/to/hive-metastore.jar>:<path/to/hive-exec.jar>:<path/to/spark-assembly.jar></pre> For MapR, use the third template. This should be the absolute path to <code>spark-assembly-1.5.2-mapr-1602-hadoop2.7.0-mapr-1602.jar</code>.
TRANSFORM_SERVICE_SERVERS	A comma-separated list of the Transform Service nodes. For best performance, these should all be Managed Servers. In particular, they shouldn't be Dgraph nodes, as both the Dgraph and the Transform Service require a lot of memory.
TRANSFORM_SERVICE_PORT	The port the Transform Service listens on for requests from Studio.
ENABLE_CLUSTERING_SERVICE	For use by Oracle Support only. Leave this property set to <code>FALSE</code> .
CLUSTERING_SERVICE_SERVERS	For use by Oracle Support only. Don't modify this property.
CLUSTERING_SERVICE_PORT	For use by Oracle Support only. Don't modify this property.
WORKFLOW_MANAGER_SERVERS	The Workflow Manager Service node.
WORKFLOW_MANAGER_PORT	The port the Workflow Manager Service listens on for data processing requests.



Chapter 6

Cluster Installation

The following sections describe how to install BDD on multiple nodes, and provide tips on troubleshooting a failed installation.

Installation overview

Setting up the install machine

Downloading the BDD media pack

Downloading a WebLogic Server patch

Configuring BDD

Running the prerequisite checker

Installing BDD on a cluster

Installation overview

You install BDD by running a single script, which installs all of its components at once. When the script completes, your cluster will be running and ready to use.

The installer is contained in one of the BDD installation packages, which you will download to the install machine. The same package also contains BDD's configuration file and a second script that verifies whether your system meets all prerequisites.

The following sections describe the installation process, from preparing the install machine to running the installer.

Silent installation

Normally, the BDD installer prompts for the following information at runtime:

- The username and password for your cluster manager (Cloudera Manager, Ambari, or MCS), which the script uses to query your cluster manager for information related to your Hadoop cluster.
- The username and password for the WebLogic Server admin. The script will create this user when it deploys WebLogic.
- The JDBC usernames and passwords for the Studio and Workflow Manager Service databases.
- The username and password for the Studio admin.
- The absolute path to the location of the installation packages.

You can avoid entering this information manually by running the installer in silent mode. To do this, set the following environment variables before installing. The installer will check for them when it runs and execute

silently if it finds them. Note that many of these will be useful after installing, as they are used by BDD's administration script.

Environment variable	Value
BDD_HADOOP_UI_USERNAME	The username for your cluster manager (Cloudera Manager, Ambari, or MCS).
BDD_HADOOP_UI_PASSWORD	The password for your cluster manager.
BDD_WLS_USERNAME	The username for the WebLogic Server administrator.
BDD_WLS_PASSWORD	The password for the WebLogic Server administrator. This must contain at least 8 characters, one of which must be a number, and cannot start with a number.
BDD_STUDIO_JDBC_USERNAME	The username for the Studio database.
BDD_STUDIO_JDBC_PASSWORD	The password for the Studio database.
BDD_WORKFLOW_MANAGER_JDBC_USERNAME	The username for the Workflow Manager Service database.
BDD_WORKFLOW_MANAGER_JDBC_PASSWORD	The password for the Workflow Manager Service database.
BDD_STUDIO_ADMIN_USERNAME	<p>The email address of the Studio admin, which will be their username. This must be a full email address and can't begin with <code>root@</code> or <code>postmaster@</code>.</p> <p>The installer will automatically populate this value to the <code>STUDIO_ADMIN_EMAIL_ADDERESS</code> property in <code>bdd.conf</code>, overwriting any existing value. If you set <code>STUDIO_ADMIN_EMAIL_ADDERESS</code> instead of this environment variable, the installer will still execute silently.</p>
BDD_STUDIO_ADMIN_PASSWORD	<p>The password for the Studio admin. This must contain at least 8 characters, one of which must be a non-alphanumeric character.</p> <p>The Studio admin will be asked to reset their password the first time they log in if you set the <code>STUDIO_ADMIN_PASSWORD_RESET_REQUIRED</code> property to <code>TRUE</code>.</p>
INSTALLER_PATH	The absolute path to the location of the installation packages. This is only required if you don't set the <code>INSTALLER_PATH</code> property in <code>bdd.conf</code> .

Setting up the install machine

The first step in the installation process is to set up the install machine.

To set up the install machine:

1. Select one machine in your cluster to be the install machine.

This can be any machine in your cluster that has the following:

- A supported operating system and JDK
- Perl 5.10+ with multithreading
- The `Mail::Address` and `XML::Parser` Perl modules
- Passwordless `sudo` and SSH enabled for the `bdduser`
- Bash set as the default shell for the `bdd` user

2. Choose an existing directory or create a new one to be the installation source directory.

You'll perform the entire installation process from this directory. Its name and location are arbitrary and it must contain at least 10GB of free space.

3. Within the installation source directory, create a new directory named `packages`.

Next, download the BDD media pack.

Downloading the BDD media pack

After you set up the install machine, you can download the BDD media pack from the Oracle Software Delivery Cloud.

To download the media pack:

1. Go to the [Oracle Software Delivery Cloud](#) and sign in.
2. Accept the Export Restrictions.
3. Check **Programs** if it isn't already.
4. In the **Product** text box, enter `Oracle Big Data Discovery`.
5. Click **Select Platform** and check **Linux x86-64**.
Oracle Big Data Discovery displays in the **Selected Products** table.
6. Click **Continue**.
7. Verify that **Available Release** and **Oracle Big Data Discovery 1.4.x.x.x for Linux x86-64** are both checked, then click **Continue**.
8. Accept the Oracle Standard Terms and Restrictions and click **Continue**.
9. In the **File Download** popup, click **Download All**.

This downloads the following packages to your machine:

- **First of two parts of the Oracle Big Data Discovery binary**
- **Second of two parts of the Oracle Big Data Discovery binary**
- **Installer for Oracle Big Data Discovery**

- **SDK for Oracle Big Data Discovery**
- **Documentation for Oracle Big Data Discovery**
- **Oracle Fusion Middleware 12c (12.1.3.0.0) WebLogic Server and Coherence**

You should also make a note of each file's part number, as you will need this information to identify it.

10. Move the BDD installer, BDD binary, and WebLogic Server packages from the download location to the `packages` directory.

11. Rename the first BDD binary package `bdd1.zip` and the second `bdd2.zip`.

This ensures that the installer will recognize them.

12. Extract the WebLogic Server package.

This creates a file called `fmw_12.1.3.0.0_wls.jar`, which contains the WebLogic Server installer.

13. Navigate back to the installation source directory and extract the installer package:

```
unzip packages/<installer_package>.zip
```

This creates a new directory within the installation source directory called `installer`, which contains the installer, `bdd.conf`, and other files required by the installer.

Next, you can download a WebLogic Server patch for the installer to apply. If you don't want to patch WebLogic Server, you should configure your BDD installation.

Downloading a WebLogic Server patch

You can optionally download a WebLogic Server patch for the installer to apply when it runs.

You can only apply one patch when installing. If it, the installer will remove it and continue running.

For more information on patching WebLogic Server, see [Oracle Fusion Middleware Patching with OPatch](#).

To download a WebLogic Server patch:

1. Within the installation source directory, create a new directory called `WLSPatches`.
Don't change the name of this directory or the installer won't recognize it.
2. Go to [My Oracle Support](#) and log in.
3. On the **Patches & Updates** tab, find and download the patch you want to apply.
4. Move all ZIP files associated with the patch to `WLSPatches/`.
Don't extract the files. The installer will do this when it runs.

Next, you should configure your BDD installation.

Configuring BDD

Before installing, you must update BDD's configuration file, `bdd.conf`, which is located in the `/<installation_src_dir>/installer` directory.

`bdd.conf` defines the configuration of your BDD cluster and provides the installer with parameters it requires to run. **Updating it is the most important step of the installation process.** If you don't modify it, or if you modify it incorrectly, the installer could fail or your cluster could be configured differently than you intended.

You can edit the file in any text editor. Be sure to save your changes before closing.

The installer validates `bdd.conf` at runtime and fails if it contains any invalid values. To avoid this, keep the following in mind when updating it:

- The accepted values for some properties are case-sensitive and must be entered exactly as they appear in this document.
- All hostnames must be Fully Qualified Domain Names (FQDNs).
- Any symlinks in paths must be identical on all nodes. If any are different or don't exist, the installation may fail.
- Each port setting must have a unique value.
- Some of the directories defined in `bdd.conf` have location requirements. These are specified in this document.

`bdd.conf` is divided into three parts:

- **Required settings:** You must update these properties with information specific to your system and installation, or the installer may fail. See [Required settings on page 63](#).
- **Optional settings:** You can update these settings if you want to further customize your installation, but the defaults will work for most. See [Optional settings on page 95](#).
- **Internal settings:** These are intended for use by Oracle Support, only. Don't edit these unless instructed to do so by a support representative. See [Internal settings on page 101](#).

[Required settings](#)

Required settings

The first part of `bdd.conf` contains required settings. You must update these with information specific to your system, or the installer could fail.

Must Set

This section contains blank settings that you must provide values for. If you don't set these, the installation will fail.

Configuration property	Description
ORACLE_HOME	<p>The path to the BDD root directory, where BDD will be installed on each node in the cluster. This directory must not exist. To ensure that the installer will be able to create it, its parent directories' permissions must be set to either 755 or 775, and there must be at least 30GB of space available on each BDD node.</p> <p>Note that this is different from the <code>ORACLE_HOME</code> environment variable required by the Studio database.</p>
ORACLE_INV_PTR	<p>The absolute path to the Oracle inventory pointer file, which the installer will create. This file can't be located in the <code>ORACLE_HOME</code> directory.</p> <p>If you have any other Oracle software products installed, this file will already exist. Update this property to point to it.</p>
INSTALLER_PATH	<p>Optional. The absolute path to the installation source directory. This must contain at least 10GB of free space.</p> <p>If you don't set this property, you can either set the <code>INSTALLER_PATH</code> environment variable or specify the path at runtime. For more information, see Installation overview on page 58.</p>
DGRAPH_INDEX_DIR	<p>The absolute path to the Dgraph databases. This directory shouldn't be located under <code>ORACLE_HOME</code>, or it will be deleted.</p> <p>The script will create this directory if it doesn't currently exist. If you're installing with existing databases, set this property to their parent directory.</p> <p>If you have HDFS data at rest encryption enabled in Hadoop and you want to store your databases on HDFS, be sure that this is in an encryption zone.</p>
HADOOP_UI_HOST	<p>The name of the server hosting your Hadoop manager (Cloudera Manager, Ambari, or MCS).</p>

Configuration property	Description
STUDIO_JDBC_URL	<p>The JDBC URL for the Studio database.</p> <p>There are three templates for this property. Copy the template that corresponds to your database type to <code>STUDIO_JDBC_URL</code> and update the URL to point to your database.</p> <ul style="list-style-type: none"> For MySQL databases, use the first template and update the URL as follows: <pre>jdbc:mysql://<database hostname>:<port number> /<database name>?useUnicode=true&characterEncoding =UTF-8&useFastDateParsing=false</pre> For Oracle databases, use the first template and update the URL as follows: <pre>jdbc:oracle:thin: @<database hostname>:<port number>:<database SID></pre> If you're not installing on a production environment and want the installer to create a Hypersonic database for you, use the third template. The script will create the database for you in the location defined by the URL. <p>If you're installing on more than one machine, be sure to use the database host's FQDN and not <code>localhost</code>.</p>
WORKFLOW_MANAGER_JDBC_URL	<p>The JDBC URL for the Workflow Manager Service database.</p> <p>There are two templates for this property. Copy the template that corresponds to your database type to <code>WORKFLOW_MANAGER_JDBC_URL</code> and update the URL to point to your database.</p> <ul style="list-style-type: none"> For MySQL databases, use the first template and update the URL as follows: <pre>jdbc:mysql://<database hostname>:<port number> /<database name>?useUnicode=true&characterEncoding =UTF-8&useFastDateParsing=false</pre> For Oracle databases, use the first template and update the URL as follows: <pre>jdbc:oracle:thin: @<database hostname>:<port number>:<database SID></pre> <p>If you're installing on more than one machine, be sure to use the database host's FQDN and not <code>localhost</code>.</p>

General


This section configures settings relevant to all components and the installation process itself.

Configuration property	Description
INSTALL_TYPE	<p>Determines the installation type according to your hardware and Hadoop distribution. Set this to one of the following:</p> <ul style="list-style-type: none"> • CDH • HW • MAPR <p>This document doesn't cover Oracle Big Data Appliance (BDA) or Oracle Public Cloud (OPC) installations. If you want to install on the Big Data Appliance, see the <i>Oracle Big Data Appliance Owner's Guide Release 4 (4.x)</i> and any corresponding MOS notes.</p>
JAVA_HOME	<p>The absolute path to the JDK install directory. This must be the same on all BDD servers and should have the same value as the <code>\$JAVA_HOME</code> environment variable.</p> <p>If you have multiple versions of the JDK installed, be sure that this points to the correct one.</p>
TEMP_FOLDER_PATH	<p>The temporary directory used on each node during the installation. This directory must exist on all BDD nodes and must contain at least 20GB of free space.</p>

CDH/HDP/MapR

This section contains properties related to Hadoop. The installer uses these properties to query the Hadoop cluster manager (Cloudera Manager, Ambari, or MCS) for information about the Hadoop components, such as the URIs and names of their host servers.

Configuration property	Description and possible settings
HADOOP_UI_PORT	The port number of the server running the Hadoop cluster manager.
HADOOP_UI_CLUSTER_NAME	The name of your Hadoop cluster, which is listed in the cluster manager. Be sure to replace any spaces in the cluster name with <code>%20</code> .
HUE_URI	HDP only. The hostname and port of the node running Hue, in the format <code><hostname>:<port></code> .

Configuration property	Description and possible settings
HADOOP_CLIENT_LIB_PATHS	<p>A comma-separated list of the absolute paths to the Hadoop client libraries.</p> <p> Note: You only need to set this property before installing if you have HDP or MapR. For CDH, the installer will download the required libraries and set this property automatically. Note that this requires an internet connection. If the script is unable to download the libraries, it will fail; see Failure to download the Hadoop client libraries on page 73 for instructions on solving this issue.</p> <p>To set this property, copy the template for your Hadoop distribution to <code>HADOOP_CLIENT_LIB_PATHS</code> and update the paths to point to the client libraries you copied to the install machine. Be sure to replace all instances of <code><UNZIPPED_XXX_BASE></code> with the absolute path to the correct library.</p> <p>Don't change the order of the paths in the list as they <i>must</i> be specified as they appear.</p>
HADOOP_CERTIFICATES_PATH	<p>Only required for Hadoop clusters with TLS/SSL enabled. The absolute path to the directory on the install machine where you put the certificates for HDFS, YARN, Hive, and the KMS.</p> <p>Don't remove this directory after installing, as you will use it if you have to update the certificates.</p>

Kerberos

This section configures Kerberos for BDD. Only modify these properties if you want to enable Kerberos.

Configuration property	Description and possible settings
ENABLE_KERBEROS	Enables Kerberos in the BDD cluster. If Kerberos is installed on your cluster and you want BDD to integrate with it, set this value to <code>TRUE</code> ; if not, set it to <code>FALSE</code> .
KERBEROS_PRINCIPAL	<p>The name of the BDD principal. This should include the name of your domain; for example, <code>bdd-service@EXAMPLE.COM</code>.</p> <p>This property is only required if <code>ENABLE_KERBEROS</code> is set to <code>TRUE</code>.</p>
KERBEROS_KEYTAB_PATH	<p>The absolute path to the BDD <code>keytab</code> file on the install machine.</p> <p>The installer will rename this to <code>bdd.keytab</code> and copy it to <code>\$BDD_HOME/common/kerberos/</code> on all BDD nodes.</p> <p>This property is only required if <code>ENABLE_KERBEROS</code> is set to <code>TRUE</code>.</p>

Configuration property	Description and possible settings
KRB5_CONF_PATH	<p>The absolute path to the <code>krb5.conf</code> file on the install machine. The installer will copy this to <code>/etc</code> on all BDD nodes.</p> <p>This property is only required if <code>ENABLE_KERBEROS</code> is set to <code>TRUE</code>.</p>

WebLogic (BDD Server)

This section configures the WebLogic Server, including the Admin Server and all Managed Servers.

Configuration property	Description and possible settings
ADMIN_SERVER	<p>The hostname of the install machine, which will become the Admin Server.</p> <p>If you leave this blank, it will default to the hostname of the machine you're on.</p>
MANAGED_SERVERS	<p>A comma-separated list of the Managed Server hostnames (the servers that will run WebLogic, Studio, and the Dgraph Gateway). This list must include the Admin Server and can't contain duplicate values.</p> <p>If you define more than one Managed Server, you must set up a load balancer in front of them after installing. For more information, see Configuring load balancing for Studio on page 82.</p>

Dgraph and HDFS Agent

This section configures the Dgraph and the HDFS Agent.

Configuration property	Description and possible settings
DGRAPH_SERVERS	<p>A comma-separated list of the hostnames of the nodes that will run the Dgraph and the Dgraph HDFS Agent.</p> <p>This list can't contain duplicate values. If you plan on storing your databases on HDFS, these must be HDFS DataNodes. For best performance, there shouldn't be any other Hadoop services running on these nodes, especially Spark.</p>

Configuration property	Description and possible settings
DGRAPH_THREADS	<p>The number of threads the Dgraph starts with. This should be at least 2. The exact number depends on the other services running on the machine:</p> <ul style="list-style-type: none"> For machines running only the Dgraph, the number of threads should be equal to the number of cores on the machine. For machines running the Dgraph and other BDD components, the number of threads should be the number of cores minus 2. For example, a quad-core machine should have 2 threads. For HDFS nodes running the Dgraph, the number of threads should be the number of CPU cores minus the number required for the Hadoop services. For example, a quad-core machine running Hadoop services that require 2 cores should have 2 threads. <p>If you leave this property blank, it will default to the number of CPU cores minus 2.</p> <p>Be sure that the number you use is in compliance with the licensing agreement.</p>
DGRAPH_CACHE	<p>The size of the Dgraph cache, in MB. Only specify the number; don't include MB.</p> <p>If you leave this property blank, it will default to either 50% of the node's available RAM or the total mount of free memory minus 2GB (whichever is larger).</p> <p>Oracle recommends allocating at least 50% of the node's available RAM to the Dgraph cache. If you later find that queries are getting cancelled because there isn't enough available memory to process them, experiment with gradually decreasing this amount.</p>
ZOOKEEPER_INDEX	<p>The index of the Dgraph cluster in the ZooKeeper ensemble, which ZooKeeper uses to identify it.</p>

Data Processing

This section configures Data Processing and the Hive Table Detector.

Configuration property	Description and possible settings
HDFS_DP_USER_DIR	<p>The location within the HDFS <code>/user</code> directory that stores the sample files created when Studio users export data. The name of this directory must not include spaces or slashes (/). The installer will create it if it doesn't already exist.</p> <p>If you have MapR and want to use an existing directory, it must be mounted with a volume.</p>
YARN_QUEUE	<p>The YARN queue Data Processing jobs are submitted to.</p>

Configuration property	Description and possible settings
HIVE_DATABASE_NAME	<p>The name of the Hive database that stores the source data for Studio data sets.</p> <p>The default value is <code>default</code>. This is the same as the default value of <code>DETECTOR_HIVE_DATABASE</code>, which is used by the Hive Table Detector. It is possible to use different databases for these properties, but it is recommended that you start with one for a first time installation.</p>
SPARK_ON_YARN_JAR	<p>The absolute path to the Spark on YARN JAR on your Hadoop nodes. This will be added to the CLI classpath.</p> <p>There are two templates for this property. Copy the value of the template that corresponds to your Hadoop distribution to <code>SPARK_ON_YARN_JAR</code> and update its value as follows:</p> <ul style="list-style-type: none"> If you have CDH, use the first template. This should be the absolute path to <code>spark-assembly.jar</code>. For HDP, use the second template. This should be the absolute paths to <code>hive-metastore.jar</code>, <code>hive-exec.jar</code> and <code>spark-assembly.jar</code>, separated by a colon: <pre><path/to/hive-metastore.jar>:<path/to/hive-exec.jar>:<path/to/spark-assembly.jar></pre> If you have MapR, use the third template. This should be the absolute path to <code>spark-assembly-1.5.2-mapr-1602-hadoop2.7.0-mapr-1602.jar</code>. <p>This JAR must be located in the same location on all Hadoop nodes.</p>

Micro Service

This section configures the Transform Service.

Configuration property	Description and possible settings
TRANSFORM_SERVICE_SERVERS	<p>A comma-separated list of the Transform Service nodes.</p> <p>For best performance, these should all be Managed Servers. In particular, they shouldn't be Dgraph nodes, as both the Dgraph and the Transform Service require a lot of memory.</p> <p>If you define multiple Transform Service nodes, you must set up a load balancer in front of them after installing. For instructions, see Configuring load balancing for the Transform Service on page 83.</p>
TRANSFORM_SERVICE_PORT	The port the Transform Service listens on for requests from Studio.
ENABLE_CLUSTERING_SERVICE	For use by Oracle Support only. Leave this property set to <code>FALSE</code> .

Configuration property	Description and possible settings
CLUSTERING_SERVICE_SERVERS	For use by Oracle Support only. Don't modify this property.
CLUSTERING_SERVICE_PORT	For use by Oracle Support only. Don't modify this property.
WORKFLOW_MANAGER_SERVERS	The Workflow Manager Service node. Note that you can only define one.
WORKFLOW_MANAGER_PORT	The port the Workflow Manager Service listens on for data processing requests.

Running the prerequisite checker

After you update `bdd.conf`, you should run the prerequisite checker.

This script checks your system to make sure it meets each requirement and verifies that `bdd.conf` has been properly updated. It outputs the results to an HTML file, which you can view in your browser.

When the script runs, it prompts you for the username and password for your Hadoop cluster manager and Studio database. You can avoid this by setting the following environment variables beforehand. Note that these are different from the environment variables used by the installer.

- `PREREQ_HADOOP_USERNAME`
- `PREREQ_HADOOP_PASSWORD`
- `PREREQ_STUDIO_DATABASE_USERNAME`
- `PREREQ_STUDIO_DATABASE_PASSWORD`

To run the prerequisite checker:

1. On the install machine, open a new terminal window and go to `<install_source_dir>/installer/linux/utils/prerequisite_validation/`.
2. Run the following command:

```
python prerequisite_validation.py <path_to_bdd.conf>
```

Where `<path_to_bdd.conf>` is the absolute path to `bdd.conf`.
3. Enter the username and password for your Hadoop cluster manager and Studio database, if prompted.
4. When the script completes, go to the timestamped output directory and open `test_report.html` in a browser.

The report lists all BDD requirements and whether each passed, failed, or was ignored. Ignored requirements aren't applicable to your system.

If everything passed, you're ready to install BDD. If any requirement failed, update your system or `bdd.conf` accordingly and rerun the prerequisite checker.

Installing BDD on a cluster

After you update `bdd.conf` and verify that you satisfied all prerequisites, you can install BDD.

Before running the installer, verify that all of BDD's prerequisites have been satisfied. Specifically, make sure that:

- You satisfied all requirements described in [Prerequisites on page 14](#).
- The `bdd` user meets the requirements described in [OS user requirements on page 24](#).
- You are working on the install machine, which is properly set up.
- The Studio and Workflow Manager Service databases (including their usernames, passwords, and schemas) are set up.
- If you are installing with existing Dgraph databases, the files are on either HDFS or the NFS and `DRAPH_INDEX_DIR` points to the correct location.
- If you want to run the script in silent mode, you set the environment variables described in [Installation overview on page 58](#).
- `bdd.conf` is available and properly configured.
- The following Hadoop components are running:
 - Cloudera Manager/Ambari/MCS
 - ZooKeeper
 - HDFS
 - Hive
 - Spark on YARN
 - YARN
 - Hue
 - NFS Gateway (if required)

To install BDD:

- 1 On the install machine, open a new terminal window and go to the `/installer` directory.
- 2 Run the installer:

```
./setup.sh
```

- 3 If you are not running the script in silent mode, enter the following information when prompted:
 - The username and password for the cluster manager.
 - A username and password for the WebLogic Server admin. The password must contain at least 8 characters, including at least 1 number, and can't begin with a number.
 - The username and password for the Studio database.
 - The password for the Studio admin. This must contain at least 8 characters, including at least 1 non-alphanumeric character.
 - The absolute path to the installation source directory, if you didn't set `INSTALLER_PATH` in `bdd.conf`.

If the script succeeds, BDD will be fully installed and running. See [Post-Installation Tasks on page 76](#) to learn more about your installation and how to verify it.

If the script fails, see [Troubleshooting a Failed Installation on page 72](#).



Chapter 7

Troubleshooting a Failed Installation

If the installer fails, you can use its console output and log files to determine why.

The installer's console output specifies the steps it performed and whether each passed or failed. For failed steps, the output indicates the cause of the failure. If a step failed on one or more specific servers, the output will also list the hostnames of those servers. For example:

```
[Installer] Error! Fail to copy Data Processing package to servers: <hostname1, hostname2>
```

You can then check the log files on those servers for more information about the failure. The installer's log files are located on each server in the directory defined by `TEMP_FOLDER_PATH`.

Once you determine what caused the failure, you can fix it and rerun the installer.

[Failed ZooKeeper check](#)

[Failure to download the Hadoop client libraries](#)

[Failure to generate the Hadoop fat JAR](#)

[Rerunning the installer](#)

Failed ZooKeeper check

The installer will fail if it can't connect to the ZooKeeper. This can occur if the ZooKeeper crashes during the installation.

If this happens, you will receive an error similar to the following:

```
Checking Zookpeers...Exception in thread "main" org.apache.zookeeper ...  
Fail! Error executing zookeeper-client on jdoe.example.com. Return code 1.
```

To fix this problem, try rerunning the installer according to the instructions in [Rerunning the installer on page 74](#). If it continues to fail, check if ZooKeeper is completely down and restart it if it is.

Failure to download the Hadoop client libraries

If you have CDH, the installer will fail if it can't download the required Hadoop client libraries. This can occur if you don't have an internet connection, or if some of the libraries are missing or incomplete.

If this occurs, you'll receive an error similar to the following:

```
Error! Cannot download <client_library_package>
```

To fix this problem:

1. On the install machine, download the following packages from <http://archive-primary.cloudera.com/cdh5/cdh/5/> and extract them:



Note: It is recommended that you use a browser other than Chrome for this.

- `spark-<spark_version>.cdh.<cdh_version>.tar.gz`
- `hive-<hive_version>.cdh.<cdh_version>.tar.gz`
- `hadoop-<hadoop_version>.cdh.<cdh_version>.tar.gz`
- `avro-<avro_version>.cdh.<cdh_version>.tar.gz`

The location you extract them to is arbitrary.

2. Open `bdd.conf` in a text editor and locate the `HADOOP_CLIENT_LIB_PATHS` property. Note that there are three templates below this property.
3. Copy and paste the value of the first template to `HADOOP_CLIENT_LIB_PATHS` and replace each instance of `$UNZIPPED_XXX_BASE` with the absolute path to that library's location on the install machine.
4. Rerun the installer.

For instructions on rerunning the installer, see [Rerunning the installer on page 74](#).

Failure to generate the Hadoop fat JAR

If you have HDP, the installer will fail if it's unable to generate the Hadoop fat JAR. This can occur if it can't find the `ojdbc6.jar` file.

To fix this problem:

1. On the install machine, create a directory called `/usr/share/java`.
2. Download `ojdbc6.jar` from <http://www.oracle.com/technetwork/apps-tech/jdbc-112010-090769.html> and copy it to `/usr/share/java`.
3. Rerun the installer.

For instructions on rerunning the installer, see [Rerunning the installer on page 74](#).

Rerunning the installer

After you have fixed the errors that caused the installer to fail, you can reinstall BDD.

To rerun the installer:

1. On the install machine, go to `$BDD_HOME/uninstall/` and run:

```
./uninstall.sh [--silent]
```

This removes many of the files created the last time you ran the installer and cleans up your environment. The `--silent` option runs the script in silent mode, which enables you to skip the confirmation step.

- 2 If you're not running the script in silent mode, enter `yes` or `y` when asked if you're sure you want to uninstall BDD.
- 3 If the installer was previously run by a different Linux user, delete the `TEMP_FOLDER_PATH` directory from all nodes.
- 4 Clean up any existing tables in the Studio and Workflow Manager Service databases.
- 5 Rerun the installer.

The installer removes any files created the last time it ran and runs again on the clean system.

Part III

After You Install



Chapter 8

Post-Installation Tasks

The following sections describe tasks you can perform after you install BDD, such as verifying your installation and increasing Linux file descriptors.

[*Verifying your installation*](#)

[*Navigating the BDD directory structure*](#)

[*Configuring load balancing*](#)

[*Updating the DP CLI whitelist and blacklist*](#)

[*Signing in to Studio as an administrator*](#)

[*Backing up your cluster*](#)

[*Replacing certificates*](#)

[*Increasing Linux file descriptors*](#)

[*Customizing the WebLogic JVM heap size*](#)

[*Configuring Studio database caching*](#)

Verifying your installation

Once the installer completes, you can verify that each of the major BDD components were installed properly and are running.

[*Verifying your cluster's health*](#)

[*Verifying Data Processing*](#)

Verifying your cluster's health

Use the `bdd-admin` script to verify the overall health of your cluster.

More information on the `bdd-admin` script is available in the *Administrator's Guide*.

To verify the deployed components:

1. On the Admin Server, open a new terminal window and navigate to the `$BDD_HOME/BDD_manager/bin` directory.
2. Run the following:

```
./bdd-admin.sh status --health-check
```

If your cluster is healthy, the script's output should be similar to the following:

```
[2015/06/19 04:18:55 -0700] [Admin Server] Checking health of BDD cluster...
[2015/06/19 04:20:39 -0700] [web009.us.example.com] Check BDD functionality.....Pass!
[2015/06/19 04:20:39 -0700] [web009.us.example.com] Check Hive Data Detector health.....Hive Data Detector
has previously run
[2015/06/19 04:20:39 -0700] [Admin Server] Successfully checked statuses.
```

Verifying Data Processing

To verify that Data Processing is running, you must launch a Data Processing workflow. You can do this in two ways:

- Use the CLI to launch a Data Processing workflow. For more information, see the *Data Processing Guide*.
- Create a data set in Studio. For more information, see the *Studio User's Guide*.



Note: If you use the CLI to verify Data Processing, you must first add the table(s) you want processed to the CLI whitelist. For more information, see [Updating the DP CLI whitelist and blacklist on page 83](#).

Navigating the BDD directory structure


Your BDD installation consists of two main directories: `$BDD_HOME` and `$DOMAIN_HOME`.

`$BDD_HOME`

`$BDD_HOME` is the root directory of your BDD installation. Its default path is:

```
$ORACLE_HOME/BDD-<version>
```

`$BDD_HOME` contains the following subdirectories.

Directory name	Description
/BDD_manager	<p>Directories related to the <code>bdd-admin</code> script:</p> <ul style="list-style-type: none"> • <code>/bin</code>: The <code>bdd-admin</code> script, which you can use to administer your cluster from the command line. • <code>/commands</code>: Scripts invoked by <code>bdd-admin</code>. • <code>/conf</code>: Contains <code>bdd.conf</code>. • <code>/lib</code>: Additional files required by <code>bdd-admin</code>. • <code>/log</code>: The <code>bdd-admin</code> log files. • <code>version.txt</code>: Version information for <code>bdd-admin</code>. <p>More information on the <code>bdd-admin</code> script is available in the <i>Administrator's Guide</i>.</p> <p> Note: Although the <code>bdd-admin</code> script can only be run from the Admin Server, this directory is created on all nodes BDD is installed on because it's required for updating cluster configuration post-installation.</p>

Directory name	Description
/bdd-shell	Files related to the optional BDD Shell component. For more information, see the <i>BDD Shell Guide</i> .
/clusteringervice	For use by Oracle Support, only. Files and directories related to the Cluster Analysis service.
/common	Files and directories required by all BDD components: <ul style="list-style-type: none">• /edp: Libraries and OLT files required by Data Processing.• /hadoop: The Hadoop fat JAR generated from the client libraries, and other Hadoop configuration files required by BDD.• /security/cacerts: Only available when BDD is installed on secure Hadoop clusters. Contains the certificates for HDFS, YARN, Hive, and the KMS services.• /templates: Additional JARs required by BDD components.
/csfmanagerservice	Install location of the Credential Store Framework (CSF) management utilities, which store the credentials the Transform Service and Workflow Manager Service used to connect to other components.
/dataprocessing/edp_cli	The DP CLI and related files.

Directory name	Description
/dgraph	<p>Files and directories related to the Dgraph, including:</p> <ul style="list-style-type: none"> • /bin: Scripts for administering the Dgraph. • /bin/trace_logs: The Dgraph Tracing Utility logs. • /bin/zk_session: ZooKeeper session information. • /conf: Stylesheets for Dgraph statistics pages and schemas for Dgraph queries and responses. • /dgraph-hdfs-agent: Scripts for administering the HDFS Agent and its libraries. • /doc: Schemas for communications between the Dgraph and other services. • /hdfs_root: The mount point for the HDFS root directory, which enables the Dgraph to access the databases. This is only used if your databases are on HDFS. • /lib and /lib64: Dgraph libraries. • /msg: Localized messages for EQL queries. • /olt: Files related to the OLT. • /ssl: File for configuring SSL. • version.txt: Contains version information for the Dgraph and HDFS Agent components. • /xquery: XQuery documents for communications between the Dgraph and other services.
/jetty	The Jetty installation location.
/logs	BDD log files.
/microservices	The Jetty and OPSS installation packages.
/opss and /opss_standalone	Install locations of the Oracle Platform Security Services application, which provides the CSF required by the Transform Service and Workflow Manager Service.
/server	<p>Files and directories related to the Dgraph Gateway, including:</p> <ul style="list-style-type: none"> • /common: JARs required by the Dgraph Gateway. • /endeca-server: EAR file for the Dgraph Gateway application. • README_BDD.txt: The BDD release notes. • version.txt: Contains version information for the Dgraph Gateway component.

Directory name	Description
/studio	Contains the EAR file for the Studio application and a version file for Studio.
/transformservice	Scripts and other resources required by the Transform Service.
/uninstall	The uninstall script and required utilities.
version.txt	Version information for your BDD installation.
/workflowmanager	Files and directories related to the Workflow Manager Service, including: <ul style="list-style-type: none"> • /dp/config: Workflow Manager Service configuration files. • /logs: Workflow Manager Service logs.

\$DOMAIN_HOME

\$DOMAIN_HOME is the root directory of Studio, the Dgraph Gateway, and your WebLogic domain. Its default path is:

```
$ORACLE_HOME/user_projects/domains/bdd-<version>_domain
```

\$DOMAIN_HOME contains the following subdirectories.

Directory name	Description
/autodeploy	Provides a way to quickly deploy applications to a development server. You can place J2EE applications in this directory; these will be automatically deployed to the WebLogic Server when it is started in development mode.
/bin	Scripts for migrating servers and services, setting up domain and startup environments, and starting and stopping the WebLogic Server and other components.
/config	Data sources and configuration files for Studio and the Dgraph Gateway.
/console-ext	Console extensions. This directory is only used on the Admin Server.
edit.lock	Ensures can only edit the domain's configuration one at a time. Don't edit this file.
fileRealm.properties	Configuration file for the file realm.
/init-info	Schemas used by the Dgraph Gateway.

Directory name	Description
/lib	The domain library. JAR files placed in this directory will be dynamically added to the end of the Dgraph Gateway's classpath when the Dgraph Gateway is started. You use this directory to add application libraries to the Dgraph Gateway's classpath.
/nodemanager	Files used by the Node Manager. <code>nodemanager.domains</code> lists the locations of directories created by the configuration wizard, and <code>nodemanager.properties</code> configures the Node Manager.
/pending	Stores pending configuration changes.
/security	Files related to domain security.
/servers	Log files and security information for each server in the cluster.
startWebLogic.sh	Script for starting the WebLogic Server.
/tmp	Temporary directory.

Configuring load balancing

Studio and the Transform Service require load balancing when installed on multiple nodes.

A *load balancer* distributes client requests to individual nodes within a cluster. It improves the speed and efficiency of the cluster by ensuring individual nodes aren't overwhelmed with requests while others remain idle.

The following sections describe how to configure load balancing for Studio and the Transform Service.

[Configuring load balancing for Studio](#)

[Configuring load balancing for the Transform Service](#)

Configuring load balancing for Studio

If you installed Studio on multiple nodes, you need to set up a load balancer in front of them to ensure that user requests are always routed to available nodes.



Note: A load balancer isn't required if Studio is only installed on one node.

There are many load balancing options available. Oracle recommends an external HTTP load balancer, but you can use whatever option is best suited to your needs and available resources. Just be sure the option you choose uses *session affinity* (also called sticky sessions).

Session affinity forces all requests from a given session to be routed to the same node, resulting in one session token. Without this, requests from a single session could be handled by multiple nodes, which would create multiple session tokens.

Configuring load balancing for the Transform Service

If you installed the Transform Service on multiple nodes, you need to set up a load balancer in front of them.



Note: A load balancer isn't required if the Transform Service is installed on one node.

There are many load balancing options available. Be sure to choose one that:

- Uses session affinity, or "sticky sessions". For more information, see [Configuring load balancing for Studio on page 82](#).
- Can assign a virtual IP address to the Transform Service cluster. This is required for Studio to communicate with the cluster; without it, Studio will only send requests to the first Transform Service instance.

To configure load balancing for the Transform Service:

1. Set up the load balancer and configure a virtual IP address for the Transform Service cluster.
2. On all Studio nodes, open `$DOMAIN_HOME/config/studio/portal-ext.properties` and change the hostname portion of `bdd.microservice.transformservice.url` to the virtual IP for the Transform Service cluster.

Don't change the port number or anything after it. The new value should be similar to `http://<virtual_IP>:7203/bdd.transformservice/v1`.

Additionally, don't change the value of `TRANSFORM_SERVICE_NODES` in `bdd.conf`, as it's used by other BDD components to locate the Transform Service.

Updating the DP CLI whitelist and blacklist

In order to create data sets from existing Hive tables, you must update the DP CLI white- and blacklists that define which tables are processed by Data Processing.

The DP CLI whitelist specifies which Hive tables should be processed. Tables not included in this list are ignored by the Hive Table Detector and any Data Processing workflows invoked by the DP CLI. Similarly, the blacklist specifies the Hive tables that should not be processed. You can use one or both of these lists to control which of your Hive tables are processed and which are not.

Once you have updated the whitelist and/or blacklist as needed, you can either wait for the Hive Table Detector to process your tables automatically or use the DP CLI to start a Data Processing workflow immediately.

For information on the DP CLI white- and blacklists, see the *Data Processing Guide*.

Signing in to Studio as an administrator

After you complete the BDD installation and deployment, you can sign in to Studio as an administrator, begin to create new users, explore data sets, re-configure Studio settings as necessary, and so on.

To sign in to Studio as an administrator:

1. Ensure the WebLogic Server on the Admin Server node is running.
(This is the WebLogic instance running Studio.)
2. Open a Web browser and load Studio.
By default, the URL is `http://<Admin Server Name>:7003/bdd`.
3. Specify the admin username and password set during the installation and click **Sign In**.
If the admin username wasn't set, log in with `admin@oracle.com`.
4. Reset the password, if prompted.
The new password must contain:
 - At least 8 characters
 - At least one non-alphabetic character

Now you can add additional Studio users. There are several ways to add new Studio Users:

- Integrate Studio with an Oracle Single Sign On (SSO) system. For details, see the *Administrator's Guide*.
- Integrate Studio with an LDAP system. For details, see the *Administrator's Guide*.
- Or, while you are signed in as an administrator, you can create users manually in Studio from the **Control Panel>Users** page.

Backing up your cluster

Oracle recommends that you back up your BDD cluster immediately after deployment.

You can do this with the `bdd-admin` script. For more information, see the *Administrator's Guide*.

Replacing certificates

Enabling SSL for Studio activates WebLogic Server's default Demo Identity and Demo Trust Keystores. As their names suggest, these keystores are untrusted and meant for demo purposes only. After deployment, you should replace them with your own certificates.

More information on WebLogic's demo keystores is available in section [Configure keystores](#) of WebLogic's *Administration Console Online Help*.

Increasing Linux file descriptors

You should increase the number of file descriptors from the 1024 default.

Having a higher number of file descriptors ensures that the WebLogic Server can open sockets under high load and not abort requests coming in from clients.



Note: On Dgraph nodes, the recommended number of open file descriptors is 65536. For more information, see [Increasing the numbers of open file descriptors and processes on page 38](#).

To increase the number of file descriptors on Linux:

- 1 Edit the `/etc/security/limits.conf` file.
- 2 Modify the **nofile** limit so that **soft** is 4096 and **hard** is 8192. Either edit existing lines or add these two lines to the file:

```
*      soft      nofile      4096
*      hard      nofile      8192
```

The "*" character is a wildcard that identifies all users.

Customizing the WebLogic JVM heap size

You can change the default JVM heap size to fit the needs of your deployment.

The default JVM heap size for WebLogic is 3GB. The size is set in the `setDomainEnv.sh` file, which is in the `$DOMAIN_HOME/bin` directory. The heap size is set with the `-Xmx` option.

To change the WebLogic JVM heap size:

- 1 Open the `setDomainEnv` file in a text editor.
- 2 Search for this comment line:

```
# IF USER_MEM_ARGS the environment variable is set, use it to override ALL MEM_ARGS values
```

- 3 Add the following line immediately after the comment line:

```
export USER_MEM_ARGS="-Xms128m -Xmx3072m ${MEM_DEV_ARGS} ${MEM_MAX_PERM_SIZE}"
```

- 4 Save and close the file.
- 5 Re-start WebLogic Server.

Configuring Studio database caching

All Studio instances are automatically configured to use synchronized database caching, so that information cached on one instance is available to the others.

Studio uses Ehcache (www.ehcache.org), which uses RMI (Remote Method Invocation) multicast to notify each instance when the cache has been updated.

Although the default caching configuration will work in most cases, you may want to customize it. You might also want to disable it entirely, depending on your environment.

[Customizing Studio database caching](#)

[Disabling Studio database caching](#)

[Clearing the Studio database cache](#)

Customizing Studio database caching

You can customize Studio's database cache configuration, if needed.

The most likely change you'd want to make would be to update the IP address and port number at the top of each configuration file:

```
<cacheManagerPeerProviderFactory
  class="net.sf.ehcache.distribution.RMICacheManagerPeerProviderFactory"
  properties="peerDiscovery=automatic,multicastGroupAddress=230.0.0.1,multicastGroupPort
=4446,timeToLive=1"
  propertySeparator=","
/>
```

Note that any changes you make must be made on all Studio nodes.

To customize Studio's database caching:

- 1 Extract the default files from the `ehcache` directory in `portal-impl.jar`.

The file is in the `WEB-INF/lib` directory, which is located in `endeca-portal.war`, which is in `bdd-studio.ear`.

- 2 Update the files as needed.

To ensure that Studio uses the correct files, you may want to rename the customized files to something like:

- `hibernate-clustered-custom.xml`
- `liferay-multi-vm-clustered-custom.xml`

- 3 Deploy the customized files:

(a) Undeploy `bdd-studio.ear`.

Use the appropriate method to undeploy the file based on whether you auto-deployed the `.ear` file or installed it.

(b) Update `bdd-studio.ear` to add a subdirectory `APP-INF/classes/ehcache/` that contains the customized XML files.

(c) Redeploy the updated `.ear` file.

- 4 If needed, update `portal-ext.properties` to reflect the customized file names:

```
net.sf.ehcache.configurationResourceName=/ehcache/hibernate-clustered-custom.xml
ehcache.multi.vm.config.location=/ehcache/liferay-multi-vm-clustered-custom.xml
```

Disabling Studio database caching

Database caching is enabled for Studio by default. This provides better network efficiency for most clusters, but can in some cases cause issues in Studio.

You will likely want to disable database caching if you installed or plan on installing Studio on multiple nodes *and* either of the following is true:

- Your network or host environment doesn't support multicast UDP traffic. This is sometimes true of VM environments.
- Your Studio nodes are on separate LANs that don't use multicast routing.

To disable database caching for Studio:

- 1 Before installing, set `STUDIO_JDBC_CACHE` to `FALSE` in `bdd.conf`.

You can also do this after installing. For instructions on updating BDD's configuration post-install, see the *Administrator's Guide*.

- 2 After installing, open `$DOMAIN_HOME/bin/setUserOverrides.sh` on each Studio node and add the following argument to `JAVA_OPTIONS`, before the final quotation mark:

```
-Dnet.sf.ehcache.disabled=true
```

- 3 Restart each Studio node.

Clearing the Studio database cache

As part of troubleshooting issues with Studio, you can clear the cache for either a single Studio instance or the entire Studio cluster.

To clear the Studio cache:

- 1 Click the **Configuration Options** icon, then click **Control Panel**.
- 2 Click **Server > Server Administration**.
- 3 In the **Actions** tab at the bottom of the page:
 - To clear the cache for the current instance only, click the **Execute** button next to **Clear content cached by this VM**.
 - To clear the cache for the entire Studio cluster, click the **Execute** button next to **Clear content cached across the cluster**.



Chapter 9

Using Studio with a Reverse Proxy

Studio can be configured to use a reverse proxy.

[About reverse proxies](#)

[Types of reverse proxies](#)

[Example sequence for a reverse proxy request](#)

[Recommendations for reverse proxy configuration](#)

[Reverse proxy configuration options for Studio](#)

About reverse proxies

A reverse proxy provides a more secure way for users to get access to application servers by retrieving resources on behalf of a client from one or more servers and returning them to the client as though they came from the server itself.

A reverse proxy is located between the client and the proxied server(s). Clients access content through the proxy server. The reverse proxy server assumes the public hostname of the proxied server. The hostname(s) of the actual/proxied servers are often internal and unknown to the client browser.

Some common reasons for implementing a reverse proxy include:

- Security or firewalling
- SSL termination
- Load balancing and failover
- Resource caching/acceleration
- URL partitioning

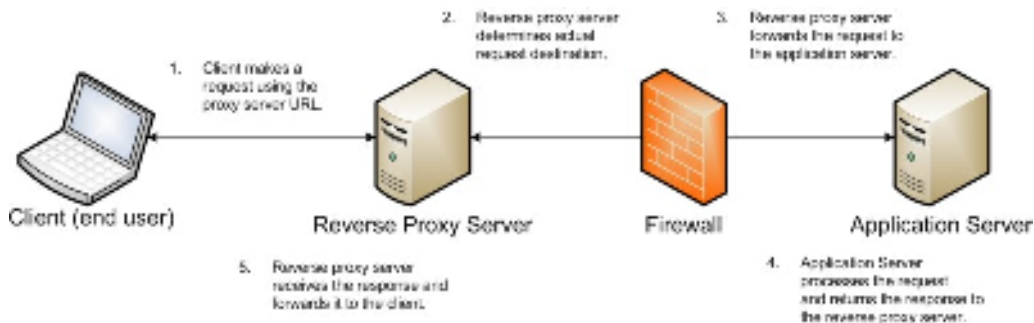
Types of reverse proxies

Reverse proxies may be either devices/appliances or specially configured web servers.

A very popular software-based reverse proxy is the Apache HTTP Server configured with the `mod_proxy` module. Many commercial web servers and reverse proxy solutions are built on top of Apache HTTP Server, including Oracle HTTP Server.

Example sequence for a reverse proxy request

Here is an example of the typical sequence for a request processed using a reverse proxy server.



1. The client makes a request to the public URL.

For this example, for a Studio project, the request URL might be something like `http://mybdd/bdd/web/myproject`, using the default port 80.

The hostname resolves to the address of the reverse proxy server. The reverse proxy is listening on this address and receives the request.

2. The reverse proxy server analyzes the URL to determine where the request needs to be proxied to.

A reverse proxy might use any part of the URL to route the request, such as the protocol, host, port, path, or query-string. Typically the path is the main data used for routing.

The reverse proxy configuration rules determine the outbound URL to send the request to. This destination is usually the end server responsible for serving the content. The reverse proxy server may also rewrite parts of the request. For example, it may change or make additions to path segments.

Reverse proxies can also add standard or custom headers to the request.

For example, the URL `http://mybdd/web/myproject` might be proxied to `http://bddserver1:8080/bdd/web/myproject`. In this case:

- The hostname of the target server is `bddserver1`
- The port is changed to 8080
- The context path `/bdd/` is added

3. The reverse proxy server sends the request to the target server.
4. The target server sends the response to the reverse proxy server.
5. The reverse proxy server reads the request and returns it to the client.

Recommendations for reverse proxy configuration

Here are some general configuration recommendations for setting up a reverse proxy.

[Preserving HTTP 1.1 Host: headers](#)

[Enabling the Apache ProxyPreserveHost directive](#)

Preserving HTTP 1.1 Host: headers

HTTP 1.1 requests often include a `Host:` header, which contains the hostname from the client request. This is because a server may use a single IP address or interface to accept requests for multiple DNS hostnames.

The `Host:` header identifies the server requested by the client. When a reverse proxy proxies an HTTP 1.1 request between a client and a target server, when it makes the request, it must add the `Host:` header to the outbound request. The `Host:` header it sends to the target server should be the same as the `Host:` header it received from the client. It should not be the `Host:` header that would be sent if accessing the target server directly.

When the application server needs to create an absolute, fully-qualified URL, such as for a redirect URL or an absolute path to an image or CSS file, it must provide the correct hostname to the client to use in a subsequent request.

For example, a Java application server sends a client-side redirect to a browser (HTTP 302 Moved). It uses the `ServletRequest.getServerName()` method to fetch the hostname in the request, then constructs a `Host:` header.

The URL sent by the client is `http://mystudio/web/myapp`. The actual internal target URL generated by the reverse proxy will be `http://studioserver1:8080/bdd/web/myapp`.

If there is no specific configuration for the target server, then if the reverse proxy retains the `Host:` header, the header is:

```
Host: http://mystudio
```

If the reverse proxy does not retain the `Host:` header, the result is:

```
Host: http://studioserver1:8080
```

In the latter case, where the header uses the actual target server hostname, the client may not have access to `studioserver1`, or may not be able to resolve the hostname. It also will bypass the reverse proxy on the next request, which may cause security issues.

If the `Host:` header cannot be relied on as correct for the client, then it must be configured specifically for the web or application server, so that it can render correct absolute URLs.

Most reverse proxy solutions should have a configuration option to allow the `Host:` header to be preserved.

Enabling the Apache ProxyPreserveHost directive

The `ProxyPreserveHost` directive is used to instruct Apache `mod_proxy`, when acting as a reverse proxy, to preserve and retain the original `Host:` header from the client browser when constructing the proxied request to send to the target server.

The default setting for this configuration directive is `Off`, indicating to not preserve the `Host:` header and instead generate a `Host:` header based on the target server's hostname.

Because this is often not what is wanted, you should add the `ProxyPreserveHost On` directive to the Apache HTTPD configuration, either in `httpd.conf` or related/equivalent configuration files.

Reverse proxy configuration options for Studio

Here are some options for configuring reverse proxy for Studio.

[Simple Studio reverse proxy configuration](#)

[Studio reverse proxy configuration without preserving Host: headers](#)

[Configuring Studio to support an SSL-enabled reverse proxy](#)

Simple Studio reverse proxy configuration

Here is a brief overview of a simple reverse proxy configuration for Studio. The configuration preserves the `Host:` header, and does not use SSL or path remapping. Studio only supports matching context paths.

In this simple configuration:

- A reverse proxy server is in front of a single Studio application server.
- The reverse proxy server is configured to preserve the `Host:` header.
- The context paths match.
- Neither the reverse proxy nor the application server is configured for SSL.

With this setup, you should be able to access Studio correctly using the reverse proxy without additional configuration.

Studio reverse proxy configuration without preserving Host: headers

If a reverse proxy used by Studio does not preserve the `Host:` header, and instead makes a request with a `Host:` header referring to the target application server, Studio and the application server receive an incorrect hostname. This causes Studio to generate absolute URLs that refer to the proxied application server instead of to the reverse proxy server.

If the reverse proxy cannot be configured to preserve the `Host:` header, you must configure a fixed hostname and port. To do this, you can either:

- Configure the application server to have a fixed hostname and port
- Use `portal-ext.properties` to configure Studio with a fixed hostname and port

Configuring a fixed hostname for the application server

In WebLogic, set up a virtual host with the fixed hostname and port.

Configuring Studio with a fixed hostname

To configure Studio with a fixed hostname and port, add the following properties to `portal-ext.properties`:

```
web.server.host=<reverseProxyHostName>
web.server.http.port=<reverseProxyPort>
```

Configuring Studio to support an SSL-enabled reverse proxy

If Studio is installed behind a reverse proxy that has SSL capabilities, and the client SSL is terminated on the reverse proxy, you must configure Studio to set the preferred protocol to HTTPS, and provide the host and port for the reverse proxy server.

To do this, add the following settings to `portal-ext.properties`:

```
web.server.protocol=https
web.server.host=<reverseProxyHostName>
web.server.https.port=<reverseProxyPort>
```

Where:

- *reverseProxyHostName* is the host name of the reverse proxy server.
- *reverseProxyPort* is the port number for the reverse proxy server.

Part IV

Uninstalling Big Data Discovery



Chapter 10

Uninstalling BDD

You uninstall BDD by running the `uninstall.sh` script from the Admin Server.

The script removes all BDD data from your system, except for the following:

- The empty BDD directories. For example, the script removes everything inside of `$ORACLE_HOME`, but leaves the directory itself. You can remove these manually when the script finishes running, although this isn't required if you're going to reinstall.
- The Dgraph databases. If you plan on reinstalling BDD, you can leave them where they are and reuse them.
- The sample files created by Data Processing.
- The `/oraInventory` directory and the `oraInst.loc` file.

Additionally, if you have MapR and moved your Dgraph databases to MapR-FS after installing, the uninstaller won't remove the mount point you created. This must be removed manually.

Note that if you upgraded BDD at any point, the script will remove all remaining files from the previous versions. You should back these up before uninstalling, if necessary.

To uninstall BDD:

1. On the Admin Server, open a command prompt and go to `$BDD_HOME/uninstall`.
2. Run the uninstallation script:

```
./uninstall.sh [--silent]
```

The `[--silent]` option runs the script in silent mode, which enables you to skip the following confirmation step.

3. Enter `yes` or `y` when asked if you're sure you want to uninstall BDD.



Appendix A

Optional and Internal BDD Properties

The following sections describe the optional and internal properties in `bdd.conf`.

[Optional settings](#)

[Internal settings](#)

Optional settings

The second part of `bdd.conf` contains optional properties. You can update these if you want, but the default values will work for most installations.

General


This section configures settings relevant to all components and the installation process itself.

Configuration property	Description
FORCE	Determines whether the installer removes files and directories left over from previous installations. Use <code>FALSE</code> if this is your first time installing BDD. Use <code>TRUE</code> if you're reinstalling after either a failed installation or an uninstallation. Note that this property only accepts UPPERCASE values.
ENABLE_AUTOSTART	Determines whether the BDD components restart automatically after their servers are rebooted. When set to <code>FALSE</code> , all components must be restarted manually. Note that this property only accepts UPPERCASE values.
BACKUP_LOCAL_TEMP_FOLDER_PATH	The absolute path to the default temporary folder on the Admin Server used during backup and restore operations. This can be overridden on a case-by-case basis by the <code>bdd-admin</code> script.
BACKUP_HDFS_TEMP_FOLDER_PATH	The absolute path to the default temporary folder on HDFS used during backup and restore operations. This can be overridden on a case-by-case basis by the <code>bdd-admin</code> script.

WebLogic (BDD Server)

This section configures WebLogic Server, including the Admin Server and all Managed Servers. It doesn't configure Studio or the Dgraph Gateway.


Configuration property	Description and possible settings
WLS_START_MODE	<p>Defines the mode WebLogic Server starts in:</p> <ul style="list-style-type: none"> <code>prod</code>: Starts WebLogic in production mode, which requires a username and password when it starts. Use this if you're installing on a production environment, as its more secure. <code>dev</code>: Starts WebLogic in development mode, which doesn't require a username or password. The installer will still prompt you for a username and password at runtime, but these will not be required when starting WebLogic Server. <p>Note that this property only accepts lowercase values.</p>
WLS_NO_SWAP	<p>Determines whether the installer checks for the required amount of free swap space (512MB) on the Admin Server and all Managed Servers before installing WebLogic Server.</p> <p>Use <code>TRUE</code> (no swap space check) if you're installing WebLogic Server on nodes that don't meet the swap space requirement.</p> <p>For more information, see Disk space requirements on page 21.</p>
WEBLOGIC_DOMAIN_NAME	The name of the WebLogic domain, which Studio and the Dgraph Gateway run in. This is automatically created by the installer.
ADMIN_SERVER_PORT	The Admin Server's port number. This number must be unique.
MANAGED_SERVER_PORT	<p>The port used by the Managed Server (i.e., Studio). This number must be unique.</p> <p>This property is still required if you're installing on a single server.</p>
WLS_SECURE_MODE	<p>Toggles SSL for Studio's outward-facing ports.</p> <p>When set to <code>TRUE</code>, the Studio instances on the Admin Server and the Managed Servers listen for requests on the <code>ADMIN_SERVER_SECURE_PORT</code> and <code>MANAGED_SERVER_SECURE_PORT</code>, respectively.</p> <p>Note that this property doesn't enable SSL for any other BDD components.</p>
ADMIN_SERVER_SECURE_PORT	<p>The secure port on the Admin Server that Studio listens on when <code>WLS_SECURE_MODE</code> is set to <code>TRUE</code>.</p> <p>Note that when SSL is enabled, Studio still listens on the un-secure <code>ADMIN_SERVER_PORT</code> for requests from the Dgraph Gateway.</p>

Configuration property	Description and possible settings
MANAGED_SERVER_SECURE_PORT	<p>The secure port on the Managed Server that Studio listens on when <code>WLS_SECURE_MODE</code> is set to <code>TRUE</code>.</p> <p>Note that when SSL is enabled, Studio still listens on the un-secure <code>MANAGED_SERVER_PORT</code> for requests from the Dgraph Gateway.</p>
ENDECA_SERVER_LOG_LEVEL	<p>The log level used by the Dgraph Gateway:</p> <ul style="list-style-type: none"> • <code>INCIDENT_ERROR</code> • <code>ERROR</code> • <code>WARNING</code> • <code>NOTIFICATION</code> • <code>TRACE</code> <p>More information on Dgraph Gateway log levels is available in the <i>Administrator's Guide</i>.</p>
SERVER_TIMEOUT	The timeout value (in milliseconds) used when responding to requests sent to all Dgraph Gateway web services except the Data Ingest Web Service. A value of 0 means there is no timeout.
SERVER_INGEST_TIMEOUT	The timeout value (in milliseconds) used when responding to requests sent to the Data Ingest Web Service. A value of 0 means there is no timeout.
SERVER_HEALTHCHECK_TIMEOUT	The timeout value (in milliseconds) used when checking data source availability when connections are initialized. A value of 0 means there is no timeout.
STUDIO_JDBC_CACHE	<p>Enables/disables database caching for Studio.</p> <p>You may want to set this to <code>FALSE</code>, depending on your environment. For more information, see Disabling Studio database caching on page 87.</p>
STUDIO_ADMIN_SCREEN_NAME	The Studio admin's screen name. This can only contain alphanumeric characters, periods (.), and hyphens (-).
STUDIO_ADMIN_EMAIL_ADDRESS	<p>The Studio admin's email address, which will be their username. This must be a full email address and can't begin with <code>root@</code> or <code>postmaster@</code>.</p> <p> Note: If you set the <code>BDD_STUDIO_ADMIN_USERNAME</code> environment variable for a silent installation, you don't need to set this property. If you do, the installer will overwrite this value with the value of <code>BDD_STUDIO_ADMIN_USERNAME</code>.</p>
STUDIO_ADMIN_PASSWORD_RESET_REQUIRED	Determines whether the Studio admin is asked to reset their password the first time they log in.

Configuration property	Description and possible settings
STUDIO_ADMIN_FIRST_NAME	The Studio admin's first name.
STUDIO_ADMIN_MIDDLE_NAME	The Studio admin's middle name.
STUDIO_ADMIN_LAST_NAME	The Studio admin's last name.

Dgraph and HDFS Agent

This section configures the Dgraph and the HDFS Agent.

Configuration property	Description and possible settings
DGRAPH_WS_PORT	The port the Dgraph listens on for requests.
DGRAPH_BULKLOAD_PORT	The port that the Dgraph listens on for bulk load ingest requests.
DGRAPH_OUT_FILE	The path to the Dgraph's stdout/stderr file.
DGRAPH_LOG_LEVEL	<p>Defines the log levels for the Dgraph's out log subsystems. This must be formatted as:</p> <pre>subsystem1 level1 subsystem2,subsystem3 level2 subsystemN levelN</pre> <p>For example:</p> <pre>DGRAPH_LOG_LEVEL =bulk_ingest WARNING cluster ERROR dgraph, eql, eve INCIDENT_ERROR</pre> <p>You can include as many subsystems as you want. Unspecified subsystems and unsupported/improperly formatted values default to NOTIFICATION.</p> <p>For more information on the Dgraph's out log subsystems and their supported levels, see the <i>Administrator's Guide</i>.</p>
DGRAPH_ADDITIONAL_ARG	<p> Note: This property is only intended for use by Oracle Support. Don't provide a value for this property when installing BDD.</p> <p>Defines one or more flags to start the Dgraph with. More information on Dgraph flags is available in the <i>Administrator's Guide</i>.</p>
DGRAPH_USE_MOUNT_HDFS	Specifies whether the Dgraph databases are stored on HDFS. When set to <code>TRUE</code> , the Dgraph runs on Hadoop DataNodes and mounts HDFS when it starts.

Configuration property	Description and possible settings
DGRAPH_HDFS_MOUNT_DIR	<p>The absolute path to the local directory where the Dgraph mounts the HDFS root directory.</p> <p>Use a nonexistent directory when installing. If this location changes after installing, the new location must be empty and have read, write, and execute permissions for the <code>bdd</code> user.</p> <p>This setting is only required if <code>DGRAPH_USE_MOUNT_HDFS</code> is set to <code>TRUE</code>.</p>
DGRAPH_ENABLE_MPP	For use by Oracle Support only. Don't modify this property.
DGRAPH_MPP_PORT	For use by Oracle Support only. Don't modify this property.
KERBEROS_TICKET_REFRESH_INTERVAL	<p>The interval (in minutes) at which the Dgraph's Kerberos ticket is refreshed. For example, if set to <code>60</code>, the Dgraph's ticket would be refreshed every 60 minutes, or every hour.</p> <p>This setting is only required if <code>DGRAPH_USE_MOUNT_HDFS</code> and <code>ENABLE_KERBEROS</code> are set to <code>TRUE</code>.</p>
KERBEROS_TICKET_LIFETIME	<p>The amount of time that the Dgraph's Kerberos ticket is valid. This should be given as a number followed by a supported unit of time: <code>s</code>, <code>m</code>, <code>h</code>, or <code>d</code>. For example, <code>10h</code> (10 hours), or <code>10m</code> (10 minutes).</p> <p>This setting is only required if <code>DGRAPH_USE_MOUNT_HDFS</code> and <code>ENABLE_KERBEROS</code> are set to <code>TRUE</code>.</p>
DGRAPH_ENABLE_CGROUP	<p>Enables cgroups for the Dgraph. This must be set to <code>TRUE</code> if you created a cgroup for the Dgraph.</p> <p>If set to <code>TRUE</code>, <code>DGRAPH_CGROUP_NAME</code> must also be set.</p>
DGRAPH_CGROUP_NAME	The name of the cgroup that controls the Dgraph. This is required if <code>DGRAPH_ENABLE_CGROUP</code> is set to <code>TRUE</code> . You must create this before installing; for more information, see Setting up cgroups on page 37 .
AGENT_PORT	The port that the HDFS Agent listens on for HTTP requests.
AGENT_EXPORT_PORT	The port that the HDFS Agent listens on for requests from the Dgraph.
AGENT_OUT_FILE	The path to the HDFS Agent's stdout/stderr file.

Data Processing

This section configures Data Processing and the Hive Table Detector.

Configuration property	Description and possible settings
ENABLE_HIVE_TABLE_DETECTOR	<p>Enables the DP CLI to automatically run the Hive Table Detector according to the schedule defined by the subsequent properties.</p> <p>When set to <code>TRUE</code>, the Hive Table Detector runs automatically on the <code>DETECTOR_SERVER</code>. By default, it does the following when it runs:</p> <ul style="list-style-type: none"> Provisions any new Hive table in the "default" database, if that table passes the whitelist and blacklist. Deletes any BDD data sets that don't have corresponding source Hive tables. This is an action that you can't prevent. <p>When set to <code>FALSE</code>, the Hive Table Detector doesn't run.</p>
DETECTOR_SERVER	The hostname of the server the Hive Table Detector runs on. This must be one of the WebLogic Managed Servers.
DETECTOR_HIVE_DATABASE	<p>The name of the Hive database that the Hive Table Detector monitors.</p> <p>The default value is <code>default</code>. This is the same as the default value of <code>HIVE_DATABASE_NAME</code>, which is used by Studio and the CLI. You can use a different database for each these properties, but Oracle recommends you start with one for a first time installation.</p> <p>This value can't contain semicolons (;).</p>
DETECTOR_MAXIMUM_WAIT_TIME	The maximum amount of time (in seconds) that the Hive Table Detector waits before submitting update jobs.
DETECTOR_SCHEDULE	The cron schedule that specifies how often the Hive Table Detector runs. The default value is <code>0 0 * * *</code> , which sets the Hive Table Detector to run at midnight every day of every month.
ENABLE_ENRICHMENTS	<p>Enables the following data enrichment modules to run during the sampling phase of data processing: Language Detection, Term Extraction, Geocoding Address, Geocoding IP, and Reverse Geotagger.</p> <p>When set to <code>true</code>, all of the data enrichments run. When set to <code>false</code>, none of them run.</p> <p>For more information on data enrichments, see the <i>Data Processing Guide</i>.</p>

Configuration property	Description and possible settings
MAX_RECORDS	<p>The maximum number of records included in a data set. For example, if a Hive table has 1,000,000 records, you could restrict the total number of sampled records to 100,000.</p> <p>Note that the actual number of records in each data set may be slightly higher or less than this value.</p>
SANDBOX_PATH	The path to the HDFS directory where the sample files created when Studio users export data are stored.
LANGUAGE	<p>Specifies either a supported ISO-639 language code (<code>en</code>, <code>de</code>, <code>fr</code>, etc.) or a value of <code>unknown</code> to set the language property for all attributes in the data set. This controls whether Oracle Language Technology (OLT) libraries are invoked during indexing.</p> <p>A language code requires more processing but produces better processing and indexing results by using the OLT libraries for the specified language. If the value is <code>unknown</code>, the processing time is faster but the processing and indexing results are more generic and OLT is not invoked.</p> <p>For a complete list of the languages BDD supports, see the <i>Data Processing Guide</i>.</p>
DP_ADDITIONAL_JARS	<p>A colon-separated list of the absolute paths to additional JARs, such as custom SerDe JARs, used during data processing. These are added to the CLI classpath.</p> <p>Note that you must manually copy each SerDe JAR to the same location on all cluster nodes before installing.</p>

Internal settings

The third part of `bdd.conf` contains internal settings either required by the installer or intended for use by Oracle Support. Note that the installer will automatically add properties to this section when it runs.



Note: Don't modify any properties in this part unless instructed to by Oracle Support.

Configuration property	Description
DP_POOL_SIZE	The maximum number of concurrent calls Studio can make to Data Processing.
DP_TASK_QUEUE_SIZE	The maximum number of jobs Studio can add to the Data Processing queue.

Configuration property	Description
MAX_INPUT_SPLIT_SIZE	<p>The maximum partition size used for Spark inputs, in MB. This controls the size of the blocks of data handled by Data Processing jobs.</p> <p>Partition size directly affects Data Processing performance. When partitions are smaller, more jobs run in parallel and cluster resources are used more efficiently. This improves both speed and stability.</p> <p>The default value is 32. This amount should be sufficient for most clusters, with a few exceptions:</p> <ul style="list-style-type: none"> • If your Hadoop cluster has a very large processing capacity and most of your data sets are small (around 1GB), you can decrease this value. • In rare cases, when data enrichments are enabled, the enriched data set in a partition can become too large for its YARN container to handle. If this occurs, you can decrease this value to reduce the amount of memory each partition requires. <p>Note that this property overrides the HDFS block size used in Hadoop.</p>
SPARK_DYNAMIC_ALLOCATION	<p>Determines whether Data Processing dynamically computes the resources allocated to the Spark executors during processing. This value should always be set to <code>true</code>.</p> <p><code>false</code> is only intended for use by Oracle Support. When set, Data Processing allocates Spark resources according to the static configuration defined by the following properties:</p> <ul style="list-style-type: none"> • SPARK_DRIVER_CORES • SPARK_DRIVER_MEMORY • SPARK_EXECUTORS • SPARK_EXECUTOR_CORES • SPARK_EXECUTOR_MEMORY
SPARK_DRIVER_CORES	The number of cores used by the Spark job driver.
SPARK_DRIVER_MEMORY	The maximum memory heap size for the Spark job driver. This must be in the same format as JVM memory settings; for example, 512m or 2g.
SPARK_EXECUTORS	The total number of Spark executors to launch.
SPARK_EXECUTOR_CORES	The number of cores for each Spark executor.

Configuration property	Description
SPARK_EXECUTOR_MEMORY	The maximum memory heap size for each Spark executor. This must be in the same format as JVM memory settings; for example, 512M or 2g.
RECORD_SEARCH_THRESHOLD	The minimum number of characters the average value of a String attribute must contain to be record searchable.
VALUE_SEARCH_THRESHOLD	The minimum number of characters the average value of a String attribute must contain to be value searchable.
BDD_VERSION	The version of BDD. This property is intended for use by Oracle Support and shouldn't be changed.
BDD_RELEASE_VERSION	The BDD hotfix or patch version. This property is intended for use by Oracle Support and shouldn't be changed.

Index

A

Admin Server, about 11

B

backups 84

bdd.conf

internal settings 101

optional settings 95

overview 62

required settings 63

Big Data Discovery

cluster configuration 12

integration with Hadoop 11

integration with WebLogic 11

overview 9

C

cgroups 37

co-locating BDD components 13

configuration

internal settings 101

optional settings 95

required settings 63

D

Data Processing, about 10

Dgraph, about 10

Dgraph Gateway, about 10

Dgraph HDFS Agent, about 10

Dgraph requirements

about 36

file descriptors and processes 38

HDFS 36

NFS Gateway 38

directory structure

\$BDD_HOME 78

\$DOMAIN_HOME 81

DP CLI

about 10

whitelist and blacklist, updating 83

E

Endeca Server 14

F

file descriptors, increasing 85

H

Hadoop, about 11

Hadoop requirements

client libraries 27

distributions and components 25

HDP JARs 28

YARN setting changes 27

Hive Table Detector, about 10

I

installation

configuration 62

installing 71

install machine, setting up 60

media pack, downloading 60

overview 58

prerequisite checker, running 70

silent mode 58

troubleshooting 73

WebLogic Server patch, downloading 61

iPad, using to view projects 42

J

Jetty 11

JVM heap size, setting 85

K

Kerberos 31

L

load balancing

overview 82

Studio 82

Transform Service 83

M

MapR

configuration 29

patches 30

special requirements 29

P

prerequisite checker, running 70

prerequisite checklist 44

prerequisites

authentication 31

authorization 32

bdd user 24

- bdd user, enabling passwordless SSH 24
- Dgraph databases 36
- encryption 33
- Hadoop client libraries 27
- Hadoop requirements 25
- hardware 20
- HDFS encryption 34
- JDK 31
- memory 20
- network 22
- operating system 22
- Perl modules, installing 23
- physical memory and disk space 21
- screen resolution 42
- Studio database 39
- Studio database commands 41
- supported browsers 41
- supported platforms 15
- YARN setting changes 27

Q

- quickstart installation 49, 50

R

- reverse proxy, using with Studio 88

S

- security
 - firewalls 35
 - Hadoop encryption 33
 - HDFS encryption 34
 - Kerberos 31
 - replacing certificates 84
 - reverse proxy 88
 - Sentry 32
 - Studio encryption 35
- Sentry 32
- silent installation 58
- single-node installation
 - about 51
 - configuring 52
 - installing 51
- Studio
 - about 9
 - database, creating 41
 - disabling 87
 - projects, viewing on iPad 42
 - signing in 84
- Studio database caching
 - clearing cache 87
 - customizing 86

- overview 85
- supported platforms 15
- system requirements
 - authentication 31
 - authorization 32
 - bdd user, enabling passwordless SSH 24
 - Dgraph databases 36
 - encryption 33
 - Hadoop client libraries 27
 - Hadoop requirements 25
 - hardware 20
 - HDFS encryption 34
 - JDK 31
 - Linux utilities 22
 - memory 20
 - operating system 22
 - Perl modules, installing 23
 - physical memory and disk space 21
 - screen resolution 42
 - Studio database 39
 - Studio database commands 41
 - supported browsers 41
 - supported platforms 15
 - YARN setting changes 27

T

- Transform Service, about 9
- troubleshooting
 - about 73
 - failed ZooKeeper check 73
 - failure to download Hadoop client libraries 73
 - failure to generate Hadoop fat JAR 74
 - installer, rerunning 74

U

- uninstalling 94

V

- verification
 - Data Processing 78
 - deployed components 77

W

- WebLogic Server
 - about 11
 - patches, downloading 61
 - setting JVM heap size 85
- Workflow Manager Service
 - about 10
 - database requirements 40