Dell | Apache Hadoop Solution Hadoop Barclamps User Guide Crowbar Version 1.2



Notes, Cautions, and Warnings



NOTE: A NOTE indicates important information that helps you make better use of your computer



CAUTION: A CAUTION indicates potential damage to hardware or loss of data if instructions are not followed.



WARNING: A WARNING indicates a potential for property damage, personal injury, or death.

Information in this document is subject to change without notice.

© 2011 Dell Inc. All rights reserved.

Reproduction of these materials is allowed under the Apache 2 license.

Information in this document is subject to change without notice.

© 2011 Dell Inc. All rights reserved.

Reproduction of these materials in any manner whatsoever without the written permission of Dell Inc. is strictly forbidden.

Trademarks used in this text: DellTM, the DELL logo, CloudEraTM, NagiosTM, GangliaTM, Opescode ChefTM, OpenStackTM, Canonical UbuntuTM, VmWareTM, Dell PrecisionTM, OptiPlexTM, LatitudeTM, PowerEdgeTM, PowerVaultTM, PowerConnectTM, OpenManageTM, EqualLogicTM, KACETM, FlexAddressTM and VostroTM are trademarks of Dell Inc. Intel®, Pentium®, Xeon®, CoreTM and Celeron® are registered trademarks of Intel Corporation in the U.S. and other countries. AMD® is a registered trademark and AMD OpteronTM, AMD PhenomTM, and AMD SempronTM are trademarks of Advanced Micro Devices, Inc. Microsoft®, Windows®, Windows Server®, MS-DOS® and Windows Vista® are either trademarks or registered trademarks of Microsoft Corporation in the United States and/or other countries. Red Hat Enterprise Linux® and Enterprise Linux® are registered trademarks of Red Hat, Inc. in the United States and/or other countries. Novell® is a registered trademark of Novell Inc. in the United States and other countries. Oracle® is a registered trademark of Oracle Corporation and/or its affiliates. Citrix®, Xen®, XenServer® and XenMotion® are either registered trademarks or trademarks of VMWare, Inc. in the United States and/or other countries. VMware®, Virtual SMP®, vMotion®, vCenter®, and vSphere® are registered trademarks or trademarks of VMWare, Inc. in the United States or other countries.

Other trademarks and trade names may be used in this publication to refer to either the entities claiming the marks and names or their products. Dell Inc. disclaims any proprietary interest in trademarks and trade names other than its own.

October 2011

Contents

INTRODUCTION	
CONCEPTS	2
CLOUDERA APACHE DISTRIBUTION	
Opscode Chef Server	
Dell Specific Options	
THE CROWBAR SOLUTION	
Architecture	
Network Setup	
Managing Growth	
Rack	
Pod	
Cluster	
DEFAULT NETWORKS	
Layout	
IP Addressing	
Rack Awareness	
USER INTERFACE	
BARCLAMPS	9
Hadoop Barclamps	
BARCLAMP DETAILS	10
HADOOP BARCLAMP	10
Hadoop Overview	10
HDFS Overview	
Cluster Deployment Topology	
Hadoop Barclamp Parameters	
PIG BARCLAMP	33
Barclamp Parameters	32
HIVE BARCLAMP	32
Barclamp Parameters	
SQOOP BARCLAMP	36
Barclamp Parameters	36
ZOOKEEPER BARCLAMP	36
Barclamp Parameters	37
SUPPORT	37
CLOUDERA SUPPORT	37

Introduction

This document provides instructions you to use when deploying **Hadoop** components on **Crowbar 1.2.** This guide is for use with the Crowbar Users Guide, it is *not* a stand alone document.

Concepts

The purpose of this guide is to explain the user interface of Crowbar. Use the Crowbar User Guide for assistance with installing Crowbar and configuring the target system.



Note: Concepts beyond the scope of this guide are introduced as needed in notes and references to other documentation.

Cloudera Apache Distribution

The focus of this guide is the use of Crowbar, *not* Hadoop. While Crowbar includes substantial components to assist in the deployment of Hadoop, its operational aspects are completely independent of Hadoop. For more information about Hadoop, please visit http://hadoop.apache.org/.

This guide provides this additional information about Cloudera as notes flagged with the Cloudera logo. For detailed operational support for Hadoop, we suggest visiting the Coudera documentation web site at http://www.cloudera.com.

Opscode Chef Server

Crowbar makes extensive use of Opscode Chef Server, http://opscode.com. To explain Crowbar actions, you should understand the underlying Chef implementation.



To use Crowbar, it is not necessary to log into the Chef Server; consequently, use of the Chef UI is not covered in this guide. Supplemental information about Chef is included.

This guide provides this additional Chef information as notes flagged with the Opscode logo.

Dell Specific Options

The Dell EULA version of Crowbar provides additional functionality and color pallets than the open source version. When divergences are relevant, they are identified.



To perform some configuration options and provide some integration, we use libraries that cannot be distributed using open source.

Crowbar is not limited to managing Dell servers and components. Due to driver requirements, some barclamps, for example: BIOS & RAID, must be targeted to specific hardware; however, those barclamps are not required for system configuration.

The Crowbar Solution

The Crowbar solution details are covered in the Crowbar Users Guide.

Architecture

For Hadoop and eco-system components (Pig, Hive, Sqoop, and Zookeeper), employ Crowbar tools to construct a starting proposal and then edit it to fit the specific requirements for your environment. Once the proposal is ready, you commit the proposal.

Network Setup

The network configuration assumes a flat L2 wiring – all network connections should be accessible at that layer. Where isolation between different logical networks is required, VLANs are used.

Managing Growth

The Dell | Cloudera Reference Architecture is organized into three components, for sizing as the Hadoop environment grows. From smallest to largest, they are:

- Rack
- Pod
- Cluster

Each has specific characteristics and sizing considerations. You can scale the environment by adding additional capacity as needed, without the need to replace any existing components.

Rack

A rack is the smallest component in a Hadoop environment, and consists of all of the power, network cabling, and two Ethernet switches required to support up to 20 data nodes. These nodes should utilize their own power connectivity and data center space – separate from other racks – and be treated as a fault zone.

Pod

A pod is a single set of stacked Ethernet switches. For the Dell | Cloudera Reference Architecture, both the maximum and minimum are six. A pod consists of the administration and operation infrastructure to support three racks.

Cluster

A cluster is a set of greater than one pod, up to a maximum of 12 pods. A cluster is a set of Hadoop nodes that share the same Network Node and management tools for operating the Hadoop environment.



Note: Please see the Dell | Cloudera Solution Reference Architecture Guide for more information.

Default Networks

The default networks are presented in the following table.

Table 1-1: Default Networks

Usage	Description	Default reserved vLAN tag	Tagged
Admin/Internal vLAN	Used for administrative functions such as Crowbar node installation, TFTP booting, DHCP assignments, KVM, system logs, backups, and other monitoring. There is only one vLAN set up for this function and it is spanned across the entire network.	100	Not tagged
BMC vLAN	Used for connecting to the BMC of each node.	100	Not tagged
Storage vLAN	Used by the Swift storage system for replication of data between machines, monitoring of data integrity, and other storage specific functions (802.1q Tagged).	200	Tagged
Edge vLANs	Used for connections to devices external to the OpenStack Cloud infrastructure; these include externally visible services such as load balancers and web servers. Use one or many of these networks, dependent on the need to segregate traffic among groups of servers (802.1q Tagged).	300	Tagged



Note: The admin and BMC networks are expected to be in the same L2 network.

Layout

Due to the nature of Crowbar's network layout, addresses are assigned to a whole network based upon interface, Network Type (Production, Management, and External) and teaming type.

Table 1-2: Master/Secondary (Admin) Name Nodes Network Connections

Interface	Network Type	Teaming Type
BMC	Management LAN	Single
LOM1	Production LAN	Teamed
LOM2	Production LAN	Teamed
Eth1	Production LAN	Teamed
Eth2	Management LAN	Single

Table 1-3: Edge Nodes Network Connections

Interface	Network Type	Teaming Type
BMC	Management LAN	Single
LOM1	Production LAN	Teamed 1

LOM2	Production LAN	Teamed 1
Eth1	External LAN	Teamed 2
Eth2	External LAN	Teamed 2

Table 1-4: Slave Nodes Network Connections

Interface	Network Type	Teaming Type
BMC	Management LAN	Single
LOM1	Production LAN	Teamed 1
LOM2	Production LAN	Teamed 1

IP Addressing

The IP address can be assigned in this fashion, using large subnets to support many machines on the production network. The Management network is a Class C network with 254 IP addresses; the Production network is what is known as a /23 with 512 IP addresses. In each network, the first 10 IP addresses are reserved for switches, routers, and firewalls.



Note: Each network's ".1" address is reserved for the network gateway.

Table 1-5: IP Addressing Schema

LAN	Network	Subnet	Gateway	Reserved
Management LAN	172.16.0.0	255.255.255.0	172.16.0.1	0.1 - 0.10
Production LAN	172.16.2.0	255.255.254.0	172.16.2.1	2.1-2.20
Name Nodes]	DHCP Allocated		
Slave Nodes		DHCP Allocated	·	•
External LAN	Т	GBD by Customer		

Rack Awareness

With the network set up using Top of Rack (ToR) switches, Rack Awareness can be programmed by using the Chef information about which switch the LOM1 is plugged into. A simple script has been added to the Hadoop configuration to pull the information out of Chef, and then use it for Rack Awareness.

Table 1-6: Pod 1 IP Addressing Layout

Network: 172.16.0.0 Netmask: 255.255.252.0

Multicast: 172.16.0.0 Broadcast 172.16.3.255

Pod	Rack Number	Network	Server Type	IP Range	Subnet Mask	Gateway
1	1	Production	Slave	172.16.0.1-42	255.255.252.0	172.16.0.1
1	2	Production	Slave	172.16.1. 1-42	255.255.252.0	172.16.0.1
1	3	Production	Slave	172.16.2. 1-42	255.255.252.0	172.16.0.1
1		Production	Master Name	172.16.3.1-19	255.255.252.0	172.16.0.1
1		Production	Secondary Name	172.16.3.20-30	255.255.252.0	172.16.0.1
1		Production	Edge Node	172.16.3.41-50	255.255.252.0	172.16.0.1
1	1	BMC	Slave	172.16.0.200-242	255.255.252.0	172.16.0.1
1	2	BMC	Slave	172.16.1.200-242	255.255.252.0	172.16.0.1
1	3	BMC	Slave	172.16.2.200.242	255.255.252.0	172.16.0.1
1		BMC	Master Name	172.16.3.201-219	255.255.252.0	172.16.0.1
1		BMC	Secondary Name	172.16.3.220-230	255.255.252.0	172.16.0.1
1		BMC	Edge Node	172.16.3.231-250	255.255.252.0	172.16.0.1

Table 1-7: Pod 2 IP Addressing Layout

Network: 172.16.0.0 Netmask: 255.255.252.0

Multicast: 172.16.0.0 Broadcast 172.16.3.255

Pod	Rack Number	Network	Server Type	IP Range	Subnet Mask	Gateway
2	1	Production	Slave	172.16.4.1-42	255.255.252.0	172.16.4.1
2	2	Production	Slave	172.16.5. 1-42	255.255.252.0	172.16.4.1
2	3	Production	Slave	172.16.6. 1-42	255.255.252.0	172.16.4.1
2		Production	Master Name	172.16.7.1-19	255.255.252.0	172.16.4.1
2		Production	Secondary Name	172.16.7.20-30	255.255.252.0	172.16.4.1
2		Production	Edge Node	172.16.7.41-50	255.255.252.0	172.16.4.1
2	1	BMC	Slave	172.16.4.200-242	255.255.252.0	172.16.4.1
2	2	BMC	Slave	172.16.5.200-242	255.255.252.0	172.16.4.1
2	3	BMC	Slave	172.16.6.200.242	255.255.252.0	172.16.4.1
2		BMC	Master Name	172.16.7.201-219	255.255.252.0	172.16.4.1
2		BMC	Secondary Name	172.16.7.220-230	255.255.252.0	172.16.4.1
2		BMC	Edge Node	172.16.7.231-250	255.255.252.0	172.16.4.1
2		External	Edge Node	TBD by Customer	TBD	TBD

For more information about Hadoop, please visit http://hadoop.apache.org/.

User Interface

Crowbar is delivered as a Web application available on the admin node using HTTP on port 3000. By default, you can access it using http://192.168.124.10:3000 Additionally, the default installation contains an implementation of Hadoop specific components (see table below).



Note: Nagios, Ganglia and Chef can be accessed directly from a web browser or via selecting one of the links on the Dashboard in the Crowbar UT.

Table 2-1: Service URLs

Service	URL	Credentials
Hadoop Jobtracker UI (master name node)	http://192.168.124.81:50070	

Note: Crowbar has been tested on the following browsers: FireFox 3.5+, FireFox 4.0, Internet Explorer 7, and Safari 5. A minimum screen resolution of 1024x768 is recommended.

The IP address (192.168.124.10) is the default address. Replace it with the address assigned to the Admin node.

Crowbar has two primary concepts for users Nodes and Barclamps. Before talking about the UI, it's important to understand how they are used by Crowbar.

Barclamps

Cloudera The Hadoop barclamp is *not* currently enabled to allow multiple proposals.

Hadoop Barclamps

Table 4-1: Barclamps

Barclamp	Function / Comments	Role
Hadoop	Common libraries and utilities that provides the basic Hadoop runtime environment (HDFS/Map Reduce). A set of components and interfaces which implements a distributed filesystem and provides general I/O access for the hadoop framework (serialization, Java RPC and persistent data storage).	Optional
Hive	Data warehouse that infrastructure provides SQL based data summarization and ad hoc querying.	Optional
Pig	Platform for analyzing large data sets that consists of a high-level language for expressing data algorithms.	Optional
Sqoop	SQL based command-line tool to assist with HDFS data import/export (SQL-to-Hadoop).	Optional

Barclamp	Function / Comments	Role
ZooKeeper	High-performance coordination service for distributed applications. ZooKeeper provides primitives such as distributed locks which can be used for building large scale distributed processing applications.	Optional

Details about individual barclamps are including in the Barclamps section below.

Barclamp Details

Hadoop Barclamp

The Hadoop software library is a framework that allows for the distributed processing of large data sets across clusters of computers using a simple programmatic driven processing model. Hadoop is designed to scale up from single servers to thousands of machines, each offering local computation and storage.

Rather than rely on hardware to deliver high-availability, the Hadoop library itself is designed to detect and handle failures at the application layer, so delivering a highly-available service on top of a cluster of computers, each of which may be prone to failures.

Hadoop is ideal for organizations with a growing need to store and process massive application datasets. It enables applications to work with thousands of nodes and petabytes of data. This Crowbar barclamp provides the ability to deploy and maintain Hadoop cluster Admin, Master, Slave and Edge nodes. It also provides the capability to configure and deliver Hadoop HDFS and MapReduce components.

Hadoop Overview

- **Hadoop Core**: The common libraries and utilities that provide the basic Hadoop runtime environment. A set of components and interfaces which implement a distributed filesystem and provide general I/O access for the Hadoop framework (serialization, Java RPC and persistent data storage).
- Hadoop Distributed File System (HDFS): A distributed file system that provides high-throughput access to application data.
- Hadoop MapReduce: A software framework for distributed processing of large data sets on compute clusters.

HDFS Overview

HDFS is a core component of the Hadoop framework and it is the underlining Hadoop virtual file system.

HDFS has the underlining concepts of three node classes:

- Master name node which is responsible for managing the file system metadata and transactions.
- Secondary name node which is responsible for checkpointing the name node's persistent state.
- Slave data nodes which are responsible for actually storing the file data.

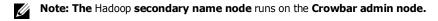
HDFS stores files as a series of blocks, each of which is by default 64MB in size. A block is the unit of storage for data nodes. Data nodes store and retrieve blocks, and have no concept of the files that are composed of these blocks.

- Master name node The master name node, is responsible for managing the filesystem metadata and data node mappings. The master name node holds that mapping from files to blocks, which it stores in memory as well as in a persistent metadata store on disk (e.g image file and edit log). The mapping between blocks and the data nodes they reside on is not stored persistently. Instead, it is stored in the name node's memory, and is built up from the periodic block reports that data nodes send to the name node. This is the primary metadata store for the cluster.
- Secondary name node The secondary name is a checkpointing mechanism which can take over the primary name node's functional aspects for this particular operation. During system operation, the name node maintains two ondisk data structures to represent the filesystem state (image file and an edit log). The image file is a checkpoint of the filesystem metadata at a point in time and the edit log is a transactional redo log of every filesystem metadata mutation since the image file was created. Incoming changes to the filesystem metadata (such as creating a new file) are written to the edit log2. When the name node starts, it reconstructs the current state by replaying the edit log. To ensure that the log doesn't grow without bound, at periodic intervals the edit log is rolled, and a new checkpoint is created by applying the old edit log to the image. This process is performed by the secondary name node daemon, often on a separate machine to the primary since creating a checkpoint has similar memory requirements to the name node itself. A side effect of the checkpointing mechanism is that the secondary holds an out-of-date copy of the primary's persistent state, which, in extreme cases, can be used to recover the filesystem's state. Blocks are stored on the underlying filesystem of the data node, as opposed to the data node managing their own storage, as native kernel level filesystems do.
- Slave nodes Slave nodes are the distributed collection points for data storage. Functioning data nodes send heartbeats to the name node every 3 seconds. This mechanism forms the communication channel between data node and name node. Occasionally, the name node will piggyback a command to a data node on the heartbeat response, for instance, "send a copy of block e to data node b". One of the first things that a data node does upon startup is send a block report to the name node, and this allows the name node to rapidly form a picture of the block distribution across the cluster.

Cluster Deployment Topology

The Crowbar Hadoop barclamp framework has expanded the concept of node deployment beyond HDFS in order to introduce the notion of a cloud edge node. The cloud edge node sits on the cloud boundary and provides the underlining interface between the data/processing capacity within the Hadoop cluster and the data consumer/end user environment. The addition of the cloud edge node serves to off-load external transactional processing requests from the data nodes and provide an additional level of security between the private cloud and the outside world.

- Master and secondary name nodes Runs all the basic services needed to manage the HDFS data storage and MapReduce task distribution and tracking.
- Slave node Runs all the services required to store blocks of data on the local hard drives and execute processing tasks against that data.
- Edge Node Provides the interface between a data and processing capacity available in the Hadoop cluster and a user of that capacity. Most of the Hadoop eco-system sub-components run on the edge node.
- Admin Node Provides cluster deployment/management capabilities and is used to deploy Hadoop to all the nodes in the cluster (The Crowbar administration node).



The typical deployment process is:

- 1. Deploy the core components: HDFS, MapReduce on the NameNodes and DataNodes.
- 2. Bring up the cluster.
- 3. Deploy the EdgeNode.
- 4. Deploy the eco-system sub-components within the cluster ZooKeeper on a slave node or Sqoop, Hive, Pig on the EdgeNode.

There may be cases when the customer may choose to deploy the add-on services on slave nodes or even the name nodes. Also when the cluster grows beyond a certain size the customer may need to run the name node (the HDFS manager) daemon and the JobTracker (the MapReduce manager) on different machines. In that case the customer needs to be able to terminate/uninstall the JobTracker daemon on the original name node and bring it up on the new JobTracker machine.

Eco-system sub-components need to be able to scale independently of the cluster configuration and/or capacity. For example, there may be cases when the data transfer capacity between the Hadoop and data warehouse (i.e. Aster Data) may exceed the max capacity of a single edge node. Adding a second edge node may be a viable alternative.

The design of the Hadoop add-on services need to separate the core Hadoop components (HDFS, MapReduce) from the add-on services and allow the customer to manipulate and deploy the services configuration that makes sense in his environment regardless of the size or topology of the actual Hadoop cluster.

Hadoop Barclamp Parameters

Table 4-21: Hadoop Barclamp configuration Parameters

Name	Description	Required	Default
fs_checkpoint_dir	Determines where on the local filesystem the DFS secondary name node should store the temporary images to merge. If this is a comma-delimited list of directories then the image is replicated in all of the directories for redundancy.	true	\${hadoop.tmp.dir}/dfs/namesecond ary
fs_checkpoint_edits_dir	Determines where on the local filesystem the DFS secondary name node should store the temporary edits to merge. If this is a comma-delimited list of directories then teh edits is replicated in all of the directories for redundancy. Default value is same as fs.checkpoint.dir.	true	\${fs.checkpoint.dir}
fs_checkpoint_period	The number of seconds between two periodic checkpoints.	true	3600
fs_checkpoint_size	The size of the current edit log (in bytes) that triggers a periodic checkpoint even if the fs.checkpoint.period hasn't expired.	true	67108864
fs_default_name	The name of the default filesystem. A URI whose scheme and authority determine the FileSystem implementation. The uri's scheme determines the config property (fs.SCHEME.impl) naming the FileSystem implementation class. The uri's authority is used to determine the host, port, etc. for a filesystem.	true	file:///
fs_file_impl	The FileSystem for file: uris.	true	org.apache.hadoop.fs.LocalFileSyst em
fs_ftp_impl	The FileSystem for ftp: uris.	true	org.apache.hadoop.fs.ftp.FTPFileS ystem

Name	Description	Required	Default
fs_har_impl	The FileSystem for Hadoop archives.	true	org.apache.hadoop.fs.HarFileSyste
fs_har_impl_disable_cache	Don't cache 'har' filesystem instances.	true	true
fs_hdfs_impl	The FileSystem for hdfs: uris.	true	org.apache.hadoop.hdfs.Distributed FileSystem
fs_hftp_impl		true	org.apache.hadoop.hdfs.HftpFileSy stem
fs_hsftp_impl		true	org.apache.hadoop.hdfs.HsftpFileS ystem
fs_kfs_impl	The FileSystem for kfs: uris.	true	org.apache.hadoop.fs.kfs.KosmosFi leSystem
fs_ramfs_impl	The FileSystem for ramfs: uris.	true	org.apache.hadoop.fs.InMemoryFil eSystem
fs_s3_block_size	Block size to use when writing files to S3.	true	67108864
fs_s3_buffer_dir	Determines where on the local filesystem the S3 filesystem should store files before sending them to S3 (or after retrieving them from S3).	true	\${hadoop.tmp.dir}/s3
fs_s3_impl	The FileSystem for s3: uris.	true	org.apache.hadoop.fs.s3.S3FileSyst
fs_s3_maxRetries	The maximum number of retries for reading or writing files to S3, before we signal failure to the application.	true	4
fs_s3_sleepTimeSeconds	The number of seconds to sleep between each S3 retry.	true	10
fs_s3n_impl	The FileSystem for s3n: (Native S3) uris.	true	org.apache.hadoop.fs.s3native.Nati veS3FileSystem
fs_trash_interval	Number of minutes between trash checkpoints. If zero, the trash feature is disabled.	false	
hadoop_http_filter_initializer s	A comma separated list of class names. Each class in the list must extend org.apache.hadoop.http.FilterInitializer. The corresponding Filter will be initialized. Then, the Filter will be applied to all user facing jsp and servlet web pages. The ordering of the list defines the ordering of the filters.	false	
hadoop_logfile_count	The max number of log files.	true	10
hadoop_logfile_size	The max size of each log file.	true	10000000
hadoop_native_lib	Should native hadoop libraries, if present, be used.	true	true
hadoop_rpc_socket_factory_ class_ClientProtocol	SocketFactory to use to connect to a DFS. If null or empty, use hadoop.rpc.socket.class.default. This socket factory is also used by DFSClient to create sockets to DataNodes.	false	
hadoop_rpc_socket_factory_	Default SocketFactory to use. This	true	org.apache.hadoop.net.StandardSoc

Name	Description	Required	Default
class_default	parameter is expected to be formatted as "package.FactoryClassName".		ketFactory
hadoop_security_authenticati on	Possible values are simple (no authentication), and kerberos.	true	simple
hadoop_security_authorizatio	Is service-level authorization enabled?	true	false
hadoop_security_group_map ping	Class for user to group mapping (get groups for a given user).	true	org.apache.hadoop.security.ShellBa sedUnixGroupsMapping
hadoop_security_uid_cache_ secs	NativeIO maintains a cache from UID to UserName. This is the timeout for an entry in that cache.	true	14400
hadoop_socks_server	Address (host:port) of the SOCKS server to be used by the SocksSocketFactory.	false	
hadoop_tmp_dir	A base for other temporary directories.	true	/tmp/hadoop-\${user.name}
hadoop_util_hash_type	The default implementation of Hash. Currently this can take one of the two values: 'murmur' to select MurmurHash and 'jenkins' to select JenkinsHash.	true	murmur
io_bytes_per_checksum	The number of bytes per checksum. Must not be larger than io.file.buffer.size.	true	512
io_compression_codecs	A list of the compression codec classes that can be used for compression/decompression.	true	org.apache.hadoop.io.compress.Def aultCodec,org.apache.hadoop.io.co mpress.GzipCodec,org.apache.hado op.io.compress.BZip2Codec
io_file_buffer_size	The size of buffer for use in sequence files. The size of this buffer should probably be a multiple of hardware page size (4096 on Intel x86), and it determines how much data is buffered during read and write operations.	true	4096
io_mapfile_bloom_error_rate	The rate of false positives in BloomFilter - s used in BloomMapFile. As this value decreases, the size of BloomFilter -s increases exponentially. This value is the probability of encountering false positives (default is 0.5%).	true	0.005
io_mapfile_bloom_size	The size of BloomFilter -s used in BloomMapFile . Each time this many keys is appended the next BloomFilter will be created (inside a DynamicBloomFilter). Larger values minimize the number of filters, which slightly increases the performance, but may waste too much space if the total number of keys is usually much smaller than this number.	true	1048576
io_seqfile_compress_blocksi ze	The minimum block size for compression in block compressed SequenceFiles .	true	1000000
io_seqfile_lazydecompress	Should values of block-compressed SequenceFiles be decompressed only when necessary.	true	true

Name	Description	Required	Default
io_seqfile_sorter_recordlimit	The limit on number of records to be kept in memory in a spill in SequenceFiles .Sorter.	true	1000000
io_serializations	A list of serialization classes that can be used for obtaining serializers and deserializers.	true	org.apache.hadoop.io.serializer.Wri tableSerialization
io_skip_checksum_errors	If true, when a checksum error is encountered while reading a sequence file, entries are skipped, instead of throwing an exception.	true	false
ipc_client_connect_max_retri	Indicates the number of retries a client will make to establish a server connection.	true	10
ipc_client_connection_maxid letime	The maximum time in msec after which a client will bring down the connection to the server.	true	10000
ipc_client_idlethreshold	Defines the threshold number of connections after which connections will be inspected for idleness.	true	4000
ipc_client_kill_max	Defines the maximum number of clients to disconnect in one go.	true	10
ipc_client_tcpnodelay	Turn on/off Nagle's algorithm for the TCP socket connection on the client. Setting to true disables the algorithm and may decrease latency with a cost of more/smaller packets.	true	false
ipc_server_listen_queue_size	Indicates the length of the listen queue for servers accepting client connections.	true	128
ipc_server_tcpnodelay	Turn on/off Nagle's algorithm for the TCP socket connection on the server. Setting to true disables the algorithm and may decrease latency with a cost of more/smaller packets.	true	false
local_cache_size	The limit on the size of cache you want to keep, set by default to 10GB. This will act as a soft limit on the cache directory for out of band data.	true	10737418240
topology_node_switch_mapp ing_impl	The default implementation of the DNSToSwitchMapping . It invokes a script specified in topology.script.file.name to resolve node names. If the value for topology.script.file.name is not set, the default value of DEFAULT_RACK is returned for all node names.	true	org.apache.hadoop.net.ScriptBased Mapping
topology_script_file_name	The script name that should be invoked to resolve DNS names to NetworkTopology names. Example: the script would take host.foobar as an argument, and return /rack1 as the output.	false	
topology_script_number_arg s	The max number of args that the script configured with topology.script.file.name	true	100

Name	Description	Required	Default
	should be run with. Each arg is an IP address.		
webinterface_private_actions	If set to true, the web interfaces of JT and NN may contain actions, such as kill job, delete file, etc., that should not be exposed to public. Enable this option if the interfaces are only reachable by those who have the right authorization.	true	false

Table 4-22. Hadoop Barclamp Environment Parameters

Name	Description	Required	Default
hadoop_datanode_opts	Command line configuration options for the data nodes.	false	
hadoop_heapsize	The maximum amount of heapsize to use, in MB (e.g. 1000MB). This is used to configure the heap size for the hadoop daemon. The default, value is 1000.	true	1000
hadoop_jobtracker_opts	Command line configuration options for the jobtracker.	false	
hadoop_log_dir	The directory where the daemons log files are stored. They are automatically created if they don't already exist.	true	/var/log/hadoop
hadoop_namenode_opts	Command line configuration options for the primary name node.	false	
hadoop_secondarynamenode_opts	Command line configuration options for the secondary name node.	false	
hadoop_tasktracker_opts	Command line configuration options for the tasktracker.	false	

Table 4-23: Hadoop Barclamp HDFS Parameters

Name	Description	Required	Default
dfs_access_time_precision	The access time for HDFS file is precise up to this value. The default value is 1 hour. Setting a value of 0 disables access times for HDFS.	true	3600000
dfs_balance_bandwidthPerSe c	Specifies the maximum amount of bandwidth that each DataNode can utilize for the balancing purpose in term of the number of bytes per second.	true	1048576
dfs_block_access_key_updat	Interval in minutes at which NameNode	true	600

Name	Description	Required	Default
e_interval	updates its access keys.		
dfs_block_access_token_ena ble	If "true", access tokens are used as capabilities for accessing DataNodes. If "false", no access tokens are checked on accessing DataNodes.	true	false
dfs_block_access_token_lifet ime	The lifetime of access tokens in minutes.	true	600
dfs_block_size	The default block size for new files.	true	67108864
dfs_blockreport_initialDelay	Delay for first block report in seconds.	false	
dfs_blockreport_intervalMse	Determines block reporting interval in milliseconds.	true	3600000
dfs_client_block_write_retrie s	The number of retries for writing blocks to the data nodes, before we signal failure to the application.	true	3
dfs_data_dir	Determines where on the local filesystem an DFS data node should store its blocks. If this is a comma-delimited list of directories, then data will be stored in all named directories, typically on different devices. Directories that do not exist are ignored.	true	\${hadoop.tmp.dir}/dfs/data
dfs_datanode_address	The address where the DataNode server will listen to. If the port is 0 then the server will start on a free port.	true	0.0.0.0:50010
dfs_datanode_data_dir_perm	Permissions for the directories on on the local filesystem where the DFS data node store its blocks. The permissions can either be octal or symbolic.	true	755
dfs_datanode_dns_interface	The name of the Network Interface from which a data node should report its IP address.	true	default
dfs_datanode_dns_nameserv er	The host name or IP address of the name server (DNS) which a DataNode should use to determine the host name used by the NameNode for communication and display purposes.	true	default
dfs_datanode_du_reserved	Reserved space in bytes per volume. Always leave this much space free for non dfs use.	false	
dfs_datanode_failed_volume s_tolerated	The number of volumes that are allowed to fail before a DataNode stops offering service. By default any volume failure will cause a DataNode to shut down.	false	
dfs_datanode_handler_count	The number of server threads for the DataNode.	true	3
dfs_datanode_http_address	The DataNode HTTP server address and port. If the port is 0 then the server will start on a free port.	true	0.0.0.0:50075
dfs_datanode_https_address	The DataNode HTTPS server address and	true	0.0.0.0:50475

Name	Description	Required	Default
	port.		
dfs_datanode_ipc_address	The DataNode IPC server address and port. If the port is 0 then the server will start on a free port.	true	0.0.0.0:50020
dfs_default_chunk_view_size	The number of bytes to view for a file on the browser.	true	32768
dfs_df_interval	Disk usage statistics refresh interval in msec.	true	60000
dfs_heartbeat_interval	Determines DataNode heartbeat interval in seconds.	true	3
dfs_hosts	Names a file that contains a list of hosts that are permitted to connect to the NameNode. The full pathname of the file must be specified. If the value is empty, all hosts are permitted.	false	
dfs_hosts_exclude	Names a file that contains a list of hosts that are not permitted to connect to the NameNode. The full pathname of the file must be specified. If the value is empty, no hosts are excluded.	false	
dfs_http_address	The address and the base port where the dfs NameNode web UI will listen on. If the port is 0 then the server will start on a free port.	true	0.0.0.0:50070
dfs_https_address		true	0.0.0.0:50470
dfs_https_client_keystore_res ource	Resource file from which SSL client keystore information will be extracted.	true	ssl-client.xml
dfs_https_enable	Decide if HTTPS (SSL) is supported on HDFS.	true	false
dfs_https_need_client_auth	Whether SSL client certificate authentication is required.	true	false
dfs_https_server_keystore_re source	Resource file from which SSL server keystore information will be extracted.	true	ssl-server.xml
dfs_max_objects	The maximum number of files, directories and blocks DFS supports. A value of zero indicates no limit to the number of objects that DFS supports.	false	
dfs_name_dir	Determines where on the local filesystem the DFS name node should store the name table(fsimage). If this is a commadelimited list of directories then the name table is replicated in all of the directories, for redundancy.	true	\${hadoop.tmp.dir}/dfs/name
dfs_name_edits_dir	Determines where on the local filesystem the DFS name node should store the transaction (edits) file. If this is a commadelimited list of directories then the transaction file is replicated in all of the directories, for redundancy. Default value is same as dfs.name.dir.	true	\${dfs.name.dir}

Name	Description	Required	Default
dfs_namenode_decommissio n_interval	Namenode periodicity in seconds to check if decommission is complete.	true	30
dfs_namenode_decommissio n_nodes_per_interval	The number of nodes NameNode checks if decommission is complete in each dfs.namenode.decommission.interval.	true	5
dfs_namenode_delegation_ke y_update_interval	The update interval for master key for delegation tokens in the NameNode in milliseconds.	true	86400000
dfs_namenode_delegation_to ken_max_lifetime	The maximum lifetime in milliseconds for which a delegation token is valid.	true	604800000
dfs_namenode_delegation_to ken_renew_interval	The renewal interval for delegation token in milliseconds.	true	86400000
dfs_namenode_handler_coun t	The number of server threads for the NameNode.	true	10
dfs_namenode_logging_level	The logging level for dfs namenode. Other values are "dir"(trac e namespace mutations), "block"(trace block under/over replications and block creations/deletions), or "all".	true	info
dfs_permissions	If "true", enable permission checking in HDFS. If "false", permission checking is turned off, but all other behavior is unchanged. Switching from one parameter value to the other does not change the mode, owner or group of files or directories.	true	true
dfs_permissions_supergroup	The name of the group of super-users.	true	supergroup
dfs_replication	Default block replication. The actual number of replications can be specified when the file is created. The default is used if replication is not specified in create time.	true	3
dfs_replication_considerLoad	Decide if chooseTarget considers the target's load or not.	true	true
dfs_replication_interval	The periodicity in seconds with which the NameNode computes replication work for DataNodes.	true	3
dfs_replication_max	Maximal block replication.	true	512
dfs_replication_min	Minimal block replication.	true	1
dfs_safemode_extension	Determines extension of safe mode in milliseconds after the threshold level is reached.	true	30000
dfs_safemode_threshold_pct	Specifies the percentage of blocks that should satisfy the minimal replication requirement defined by dfs.replication.min. Values less than or equal to 0 mean not to start in safe mode. Values greater than 1 will make safe mode permanent.	true	0.999f

Name	Description	Required	Default
dfs_secondary_http_address	The secondary NameNode HTTP server address and port. If the port is 0 then the server will start on a free port.	true	0.0.0.0:50090
dfs_support_append	Does HDFS allow appends to files? This is currently set to false because there are bugs in the "append code" and is not supported in any prodction cluster.	true	false
dfs_web_ugi	The user account used by the web interface. Syntax: USERNAME,GROUP1,GROUP2,	true	webuser,webgroup

Table 4-24: Hadoop Barclamp Map/Reduce Parameters

Name	Description	Required	Default
hadoop_job_history_location	If job tracker is static the history files are stored in this single well known place. If No value is set here, by default, it is in the local file system at \${hadoop.log.dir}/history.	false	
hadoop_job_history_user_loc ation	User can specify a location to store the history files of a particular job. If nothing is specified, the logs are stored in output directory. The files are stored in "_logs/history/" in the directory. User can stop logging by giving the value "none".	false	
hadoop_rpc_socket_factory_ class_JobSubmissionProtocol	SocketFactory to use to connect to a Map/Reduce master (JobTracker). If null or empty, then use hadoop.rpc.socket.class.default.	false	
io_map_index_skip	Number of index entries to skip between each entry. Zero by default. Setting this to values larger than zero can facilitate opening large map files using less memory.	false	
io_sort_factor	The number of streams to merge at once while sorting files. This determines the number of open file handles.	true	10
io_sort_mb	The total amount of buffer memory to use while sorting files, in megabytes. By default, gives each merge stream 1MB, which should minimize seeks.	true	100
io_sort_record_percent	The percentage of io.sort.mb dedicated to tracking record boundaries. Let this value be r, io.sort.mb be x. The maximum number of records collected before the collection thread must block is equal to (r * x) / 4.	true	0.05
io_sort_spill_percent	The soft limit in either the buffer or record	true	0.80

Name	Description	Required	Default
	collection buffers. Once reached, a thread will begin to spill the contents to disk in the background. Note that this does not imply any chunking of data to the spill. A value less than 0.5 is not recommended.		
job_end_retry_attempts	Indicates how many times Hadoop should attempt to contact the notification URL.	false	
job_end_retry_interval	Indicates time in milliseconds between notification URL retry calls.	true	30000
jobclient_output_filter	The filter for controlling the output of the task's userlogs sent to the console of the JobClient . The permissible options are: NONE, KILLED, FAILED, SUCCEEDED, and ALL.	true	FAILED
keep_failed_task_files	Controls whether the files for failed tasks be kept. This should only be used on jobs that are failing, because the storage is never reclaimed. It also prevents the map outputs from being erased from the reduce directory as they are consumed.	true	false
map_sort_class	The default sort class for sorting keys.	true	org.apache.hadoop.util.QuickSort
mapred_acls_enabled	Specifies whether ACLs should be checked for authorization of users for doing various queue and job level operations. ACLs are disabled by default. If enabled, access control checks are made by JobTracker and TaskTracker when requests are made by users for queue operations like submit job to a queue and kill a job in the queue and job operations like viewing the job-details (See mapreduce.job.acl-view-job) or for modifying the job (See mapreduce.job.acl-modify-job) using Map/Reduce APIs, RPCs or via the console and web user interfaces.	true	false
mapred_child_env	User added environment variables for the task tracker child processes. Example: 1) A=foo This will set the env variable A to foo 2) B=\$B:c This is inherit tasktracker's B env variable.	false	
mapred_child_java_opts	Java opts for the task tracker child processes. The following symbol, if present, will be interpolated: @taskid@ is replaced by current TaskID . Any other occurrences of '@' will go unchanged. For example, to enable verbose gc logging to a file named for the taskid in /tmp and to set the heap maximum to be a gigabyte, pass a 'value' of: -Xmx1024m -verbose:gc - Xloggc:/tmp/@taskid@.gc The configuration variable mapred.child.ulimit can be used to control the maximum virtual memory of the child processes.	true	-Xmx200m

Name	Description	Required	Default
mapred_child_tmp	To set the value of tmp directory for map and reduce tasks. If the value is an absolute path, it is directly assigned. Otherwise, it is prepended with task's working directory. The java tasks are executed with option - Djava.io.tmpdir='the absolute path of the tmp dir'. Pipes and streaming are set with environment variable, TMPDIR='the absolute path of the tmp dir'.	true	./tmp
mapred_child_ulimit	The maximum virtual memory, in KB, of a process launched by the Map-Reduce framework. This can be used to control both the Mapper/Reducer tasks and applications using Hadoop Pipes, Hadoop Streaming etc. By default it is left unspecified to let cluster admins control it via limits.conf and other such relevant mechanisms. Note: mapred.child.ulimit must be greater than or equal to the -Xmx passed to JavaVM, else the VM might not start.	false	
mapred_cluster_map_memor y_mb	The size, in terms of virtual memory, of a single map slot in the Map-Reduce framework, used by the scheduler. A job can ask for multiple slots for a single map task via mapred.job.map.memory.mb, upto the limit specified by mapred.cluster.max.map.memory.mb, if the scheduler supports the feature. The value of -1 indicates that this feature is turned off.	true	-1
mapred_cluster_max_map_m emory_mb	The maximum size, in terms of virtual memory, of a single map task launched by the Map-Reduce framework, used by the scheduler. A job can ask for multiple slots for a single map task via mapred.job.map.memory.mb, upto the limit specified by mapred.cluster.max.map.memory.mb, if the scheduler supports the feature. The value of -1 indicates that this feature is turned off.	true	-1
mapred_cluster_max_reduce _memory_mb	The maximum size, in terms of virtual memory, of a single reduce task launched by the Map-Reduce framework, used by the scheduler. A job can ask for multiple slots for a single reduce task via mapred.job.reduce.memory.mb, upto the limit specified by mapred.cluster.max.reduce.memory.mb, if the scheduler supports the feature. The value of -1 indicates that this feature is turned off.	true	-1

Name	Description	Required	Default
mapred_cluster_reduce_mem ory_mb	The size, in terms of virtual memory, of a single reduce slot in the Map-Reduce framework, used by the scheduler. A job can ask for multiple slots for a single reduce task via mapred.job.reduce.memory.mb, upto the limit specified by mapred.cluster.max.reduce.memory.mb, if the scheduler supports the feature. The value of -1 indicates that this feature is turned off.	true	-1
mapred_compress_map_outp ut	Controls whether the outputs of the maps be compressed before being sent across the network. Uses SequenceFile compression.	true	false
mapred_healthChecker_inter val	Frequency of the node health script to be run, in milliseconds.	true	60000
mapred_healthChecker_scrip t_args	List of arguments which are to be passed to node health script when it is being launched comma separated.	false	
mapred_healthChecker_scrip t_path	Absolute path to the script which is periodically run by the node health monitoring service to determine if the node is healthy or not. If the value of this key is empty or the file does not exist in the location configured here, the node health monitoring service is not started.	false	
mapred_healthChecker_scrip t_timeout	Time after node health script should be killed if unresponsive and considered that the script has failed.	true	600000
mapred_heartbeats_in_secon d	Expert: Approximate number of heart-beats that could arrive at JobTracker in a second. Assuming each RPC can be processed in 10msec, the default value is made 100 RPCs in a second.	true	100
mapred_hosts	Names a file that contains the list of nodes that may connect to the jobtracker. If the value is empty, all hosts are permitted.	false	
mapred_hosts_exclude	Names a file that contains the list of hosts that should be excluded by the jobtracker. If the value is empty, no hosts are excluded.	false	
mapred_inmem_merge_thres hold mapred_job_map_memory_	The threshold, in terms of the number of files, for the in-memory merge process. When we accumulate threshold number of files we initiate the in-memory merge and spill to disk. A value of 0 or less than 0 indicates we want to DON'T have any threshold and instead depend only on the ramfs's memory consumption to trigger the merge. The size, in terms of virtual memory, of a	true	-1

Name	Description	Required	Default
mb	single map task for the job. A job can ask for multiple slots for a single map task, rounded up to the next multiple of mapred.cluster.map.memory.mb and upto the limit specified by mapred.cluster.max.map.memory.mb, if the scheduler supports the feature. The value of -1 indicates that this feature is turned off iff mapred.cluster.map.memory.mb is also turned off (-1).		
mapred_job_queue_name	Queue to which a job is submitted. This must match one of the queues defined in mapred.queue.names for the system. Also, the ACL setup for the queue must allow the current user to submit a job to the queue. Before specifying a queue, ensure that the system is configured with the queue, and access is allowed for submitting jobs to the queue.	true	default
mapred_job_reduce_input_b uffer_percent	The percentage of memory- relative to the maximum heap size- to retain map outputs during the reduce. When the shuffle is concluded, any remaining map outputs in memory must consume less than this threshold before the reduce can begin.	true	0.0
mapred_job_reduce_memory _mb	The size, in terms of virtual memory, of a single reduce task for the job. A job can ask for multiple slots for a single map task, rounded up to the next multiple of mapred.cluster.reduce.memory.mb and upto the limit specified by mapred.cluster.max.reduce.memory.mb, if the scheduler supports the feature. The value of -1 indicates that this feature is turned off iff mapred.cluster.reduce.memory.mb is also turned off (-1).	true	-1
mapred_job_reuse_jvm_num _tasks	How many tasks to run per JVM. If set to - 1, there is no limit.	true	1
mapred_job_shuffle_input_b uffer_percent	The percentage of memory to be allocated from the maximum heap size to storing map outputs during the shuffle.	true	0.70
mapred_job_shuffle_merge_ percent	The usage threshold at which an in- memory merge will be initiated, expressed as a percentage of the total memory allocated to storing in-memory map outputs, as defined by mapred.job.shuffle.input.buffer.percent.	true	0.66
mapred_job_tracker	The host and port that the MapReduce job tracker runs at. If "local", then jobs are run in-process as a single map and reduce task.	true	local
mapred_job_tracker_handler	The number of server threads for the	true	10

Name	Description	Required	Default
_count	JobTracker . This should be roughly 4% of the number of tasktracker nodes.		
mapred_job_tracker_history_ completed_location	The completed job history files are stored at this single well known location. If nothing is specified, the files are stored at \${hadoop.job.history.location}/done.	false	
mapred_job_tracker_http_ad dress	The job tracker HTTP server address and port the server will listen on. If the port is 0 then the server will start on a free port.	true	0.0.0.0:50030
mapred_job_tracker_jobhisto ry_lru_cache_size	The number of job history files loaded in memory. The jobs are loaded when they are first accessed. The cache is cleared based on LRU.	true	5
mapred_job_tracker_persist_j obstatus_active	Indicates if persistency of job status information is active or not.	true	false
mapred_job_tracker_persist_j obstatus_dir	The directory where the job status information is persisted in a file system to be available after it drops of the memory queue and between jobtracker restarts.	true	/jobtracker/jobsInfo
mapred_job_tracker_persist_j obstatus_hours	The number of hours job status information is persisted in DFS. The job status information will be available after it drops of the memory queue and between jobtracker restarts. With a zero value the job status information is not persisted at all in DFS.	false	
mapred_job_tracker_retiredjo bs_cache_size	The number of retired job status to keep in the cache.	true	1000
mapred_jobtracker_blacklist_fault_bucket_width	The width (in minutes) of each bucket in the tasktracker fault timeout window. Each bucket is reused in a circular manner after a full timeout-window interval (defined by mapred.jobtracker.blacklist.fault-timeout-window).	true	15
mapred_jobtracker_blacklist_ fault_timeout_window	The timeout (in minutes) after which perjob tasktracker faults are forgiven. The window is logically a circular buffer of time-interval buckets whose width is defined by mapred.jobtracker.blacklist.fault-bucketwidth; when the "now" pointer moves across a bucket boundary, the previous contents (faults) of the new bucket are cleared. In other words, the timeout's granularity is determined by the bucket width.	true	180
mapred_jobtracker_complete userjobs_maximum	The maximum number of complete jobs per user to keep around before delegating them to the job history.	true	100
mapred_jobtracker_job_histo ry_block_size	The block size of the job history file. Since the job recovery uses job history, it is	true	3145728

Name	Description	Required	Default
	important to dump job history to disk as soon as possible. Note that this is an expert level parameter. The default value is set to 3 MB.		
mapred_jobtracker_maxtasks _per_job	The maximum number of tasks for a single job. A value of -1 indicates that there is no maximum.	true	-1
mapred_jobtracker_restart_re cover	"true" to enable (job) recovery upon restart, "false" to start afresh.	true	false
mapred_jobtracker_taskSche duler	The class responsible for scheduling the tasks.	true	org.apache.hadoop.mapred.JobQue ueTaskScheduler
mapred_jobtracker_taskSche duler_maxRunningTasksPerJ ob	The maximum number of running tasks for a job before it gets preempted. No limits if undefined.	false	
mapred_line_input_format_li nespermap	Number of lines per split in NLineInputFormat .	true	1
mapred_local_dir	The local directory where MapReduce stores intermediate data files. May be a comma-separated list of directories on different devices in order to spread disk i/o. Directories that do not exist are ignored.	true	\${hadoop.tmp.dir}/mapred/local
mapred_local_dir_minspacek ill	If the space in mapred.local.dir drops under this, do not ask more tasks until all the current ones have finished and cleaned up. Also, to save the rest of the tasks we have running, kill one of them, to clean up some space. Start with the reduce tasks, then go with the ones that have finished the least. Value in bytes.	false	
mapred_local_dir_minspaces tart	If the space in mapred.local.dir drops under this, do not ask for more tasks. Value in bytes.	false	
mapred_map_max_attempts	Expert: The maximum number of attempts per map task. In other words, framework will try to execute a map task these many number of times before giving up on it.	true	4
mapred_map_output_compre ssion_codec	If the map outputs are compressed, how should they be compressed?.	true	org.apache.hadoop.io.compress.Def aultCodec
mapred_map_tasks	The default number of map tasks per job. Ignored when mapred.job.tracker is "local".	true	2
mapred_map_tasks_speculati ve_execution	If true, then multiple instances of some map tasks may be executed in parallel.	true	true
mapred_max_tracker_blackli sts	The number of blacklists for a tasktracker by various jobs after which the tasktracker will be marked as potentially faulty and is a candidate for graylisting across all jobs. (Unlike blacklisting, this is advisory; the tracker remains active. However, it is reported as graylisted in the web UI, with	true	4

Name	Description	Required	Default
	the expectation that chronically graylisted trackers will be manually decommissioned.) This value is tied to mapred.jobtracker.blacklist.fault-timeout-window; faults older than the window width are forgiven, so the tracker will recover from transient problems. It will also become healthy after a restart.		
mapred_max_tracker_failure s	The number of task-failures on a tasktracker of a given job after which new tasks of that job aren't assigned to it.	true	4
mapred_merge_recordsBefor eProgress	The number of records to process during a merge before sending a progress notification to the TaskTracker.	true	10000
mapred_min_split_size	The minimum size chunk that map input should be split into. Note that some file formats may have minimum split sizes that take priority over this setting.	false	
mapred_output_compress	Controls whether the job outputs will be compressed.	true	false
mapred_output_compression _codec	If the job outputs are compressed, how should they be compressed?	true	org.apache.hadoop.io.compress.Def aultCodec
mapred_output_compression _type	If the job outputs are to compressed as SequenceFiles, how should they be compressed? Should be one of NONE, RECORD or BLOCK.	true	RECORD
mapred_queue_default_state	This values defines the state, default queue is in. the values can be either "STOPPED" or "RUNNING" This value can be changed at runtime.	true	RUNNING
mapred_queue_names	Comma separated list of queues configured for this jobtracker. Jobs are added to queues and schedulers can configure different scheduling properties for the various queues. To configure a property for a queue, the name of the queue must match the name specified in this value. Queue properties that are common to all schedulers are configured here with the naming convention, mapred.queue.\$QUEUE-NAME.\$PROPERTY-NAME, for e.g. mapred.queue.default.submit-job-acl. The number of queues configured in this parameter could depend on the type of scheduler being used, as specified in mapred.jobtracker.taskScheduler. For example, the JobQueueTaskScheduler supports only a single queue, which is the default configured here. Before adding more queues, ensure that the scheduler you've configured supports multiple queues.	true	default

Name	Description	Required	Default
mapred_reduce_copy_backof f	The maximum amount of time (in seconds) a reducer spends on fetching one map output before declaring it as failed.	true	300
mapred_reduce_max_attempt s	Expert: The maximum number of attempts per reduce task. In other words, framework will try to execute a reduce task these many number of times before giving up on it.	true	4
mapred_reduce_parallel_copi es	The default number of parallel transfers run by reduce during the copy(shuffle) phase.	true	5
mapred_reduce_slowstart_co mpleted_maps	Fraction of the number of maps in the job which should be complete before reduces are scheduled for the job.	true	0.05
mapred_reduce_tasks	The default number of reduce tasks per job. Typically set to 99% of the cluster's reduce capacity, so that if a node fails the reduces can still be executed in a single wave. Ignored when mapred.job.tracker is "local".	true	1
mapred_reduce_tasks_specul ative_execution	If true, then multiple instances of some reduce tasks may be executed in parallel.	true	true
mapred_skip_attempts_to_sta rt_skipping	The number of Task attempts AFTER which skip mode will be kicked off. When skip mode is kicked off, the tasks reports the range of records which it will process next, to the TaskTracker. So that on failures, TT knows which ones are possibly the bad records. On further executions, those are skipped.	true	2
mapred_skip_map_auto_incr _proc_count	The flag which if set to true, SkipBadRecords .COUNTER_MAP_PRO CESSED_RECORDS is incremented by MapRunner after invoking the map function. This value must be set to false for applications which process the records asynchronously or buffer the input records. For example streaming. In such cases applications should increment this counter on their own.	true	true
mapred_skip_map_max_skip _records	The number of acceptable skip records surrounding the bad record PER bad record in mapper. The number includes the bad record as well. To turn the feature of detection/skipping of bad records off, set the value to 0. The framework tries to narrow down the skipped range by retrying until this threshold is met OR all attempts get exhausted for this task. Set the value to Long.MAX_VALUE to indicate that framework need not try to narrow down. Whatever records (depends on application)	false	

Name	Description	Required	Default
	get skipped are acceptable.		
mapred_skip_out_dir	If no value is specified here, the skipped records are written to the output directory at _logs/skip. User can stop writing skipped records by giving the value "none".	false	
mapred_skip_reduce_auto_in cr_proc_count	The flag which if set to true, SkipBadRecords .COUNTER_REDUCE_ PROCESSED_GROUPS is incremented by framework after invoking the reduce function. This value must be set to false for applications which process the records asynchronously or buffer the input records. For example streaming. In such cases applications should increment this counter on their own.	true	true
mapred_skip_reduce_max_sk ip_groups	The number of acceptable skip groups surrounding the bad group PER bad group in reducer. The number includes the bad group as well. To turn the feature of detection/skipping of bad groups off, set the value to 0. The framework tries to narrow down the skipped range by retrying until this threshold is met OR all attempts get exhausted for this task. Set the value to Long.MAX_VALUE to indicate that framework need not try to narrow down. Whatever groups(depends on application) get skipped are acceptable.	false	
mapred_submit_replication	The replication level for submitted job files. This should be around the square root of the number of nodes.	true	10
mapred_system_dir	The directory where MapReduce stores control files.	true	\${hadoop.tmp.dir}/mapred/system
mapred_task_cache_levels	This is the max level of the task cache. For example, if the level is 2, the tasks cached are at the host level and at the rack level.	true	2
mapred_task_profile	To set whether the system should collect profiler information for some of the tasks in this job? The information is stored in the user log directory. The value is "true" if task profiling is enabled.	true	false
mapred_task_profile_maps	To set the ranges of map tasks to profile. mapred.task.profile has to be set to true for the value to be accounted.	true	0-2
mapred_task_profile_reduces	To set the ranges of reduce tasks to profile. mapred.task.profile has to be set to true for the value to be accounted.	true	0-2
mapred_task_timeout	The number of milliseconds before a task will be terminated, if it neither reads an input; writes an output; nor updates its status string.	true	600000

Name	Description	Required	Default
mapred_task_tracker_http_ad dress	The task tracker HTTP server address and port. If the port is 0 then the server will start on a free port.	true	0.0.0.0:50060
mapred_task_tracker_report_ address	The interface and port that task tracker server listens on. Since it is only connected to by the tasks, it uses the local interface. EXPERT ONLY. Should only be changed if your host does not have the loopback interface.	true	127.0.0.1:0
mapred_task_tracker_task_co ntroller	TaskController which is used to launch and manage task execution.	true	org.apache.hadoop.mapred.Default TaskController
mapred_tasktracker_dns_inte rface	The name of the Network Interface from which a task tracker should report its IP address.	true	default
mapred_tasktracker_dns_na meserver	The host name or IP address of the name server (DNS) which a TaskTracker should use to determine the host name used by the JobTracker for communication and display purposes.	true	default
mapred_tasktracker_expiry_i nterval	Expert: The time-interval, in milliseconds, after which a tasktracker is declared 'lost' if it doesn't send heartbeats.	true	600000
mapred_tasktracker_indexcac he_mb	The maximum memory that a task tracker allows for the index cache that is used when serving map outputs to reducers.	true	10
mapred_tasktracker_map_tas ks_maximum	The maximum number of map tasks that will be run simultaneously by a task tracker.	true	2
mapred_tasktracker_memory _calculator_plugin	Name of the class whose instance will be used to query memory information on the tasktracker. The class must be an instance of org.apache.hadoop.util.MemoryCalculator Plugin. If the value is null, the tasktracker attempts to use a class appropriate to the platform. Currently, the only platform supported is Linux.	false	
mapred_tasktracker_reduce_t asks_maximum	The maximum number of reduce tasks that will be run simultaneously by a task tracker.	true	2
mapred_tasktracker_taskme morymanager_monitoring_in terval	The interval, in milliseconds, for which the tasktracker waits between two cycles of monitoring its tasks' memory usage. Used only if tasks' memory management is enabled via mapred.tasktracker.tasks.maxmemory.	true	5000
mapred_tasktracker_tasks_sl eeptime_before_sigkill	The time, in milliseconds, the tasktracker waits for sending a SIGKILL to a process, after it has been sent a SIGTERM.	true	5000
mapred_temp_dir	A shared directory for temporary files.	true	\${hadoop.tmp.dir}/mapred/temp

Name	Description	Required	Default
mapred_user_jobconf_limit	The maximum allowed size of the user jobconf. The default is set to 5 MB.	true	5242880
mapred_userlog_limit_kb	The maximum size of user-logs of each task in KB. 0 disables the cap.	false	
mapred_userlog_retain_hours	The maximum time, in hours, for which the user-logs are to be retained after the job completion.	true	24
mapreduce_job_acl_modify_ job	Job specific access-control list for 'modifying' the job. It is only used if authorization is enabled in Map/Reduce by setting the configuration property mapred.acls.enabled to true. This specifies the list of users and/or groups who can do modification operations on the job. For specifying a list of users and groups the format to use is "user1,user2 group1,group". If set to '*', it allows all users/groups to modify this job. If set to '(i.e. space), it allows none. This configuration is used to guard all the modifications with respect to this job and takes care of all the following operations: o killing this job o killing a task of this job, failing a task of this job o setting the priority of this job Each of these operations are also protected by the per-queue level ACL "acl-administer-jobs" configured via mapred-queues.xml. So a caller should have the authorization to satisfy either the queue-level ACL or the job-level ACL. Irrespective of this ACL configuration, job-owner, the user who started the cluster, cluster administrators configured via mapreduce.cluster.administrators and queue administrators of the queue to which this job is submitted to configured via mapred.queue.queue-name.acl-administer-jobs in mapred-queue-acls.xml can do all the modification operations on a job. By default, nobody else besides job-owner, the user who started the cluster administrators and queue administrators can perform modification operations on a job.	false	
mapreduce_job_acl_view_job	Job specific access-control list for 'viewing' the job. It is only used if authorization is enabled in Map/Reduce by setting the configuration property mapred.acls.enabled to true. This specifies the list of users and/or groups who can view private details about the job. For specifying a list of users and groups the format to use is "user1,user2 group1,group". If set to '*', it allows all	false	

Name	Description	Required	Default
	users/groups to modify this job. If set to ' '(i.e. space), it allows none. This configuration is used to guard some of the job-views and at present only protects APIs that can return possibly sensitive information of the job-owner like o job- level counters o task-level counters o tasks' diagnostic information o task-logs displayed on the TaskTracker web-UI and o job.xml showed by the JobTracker 's web-UI Every other piece of information of jobs is still accessible by any other user, for e.g., JobStatus, JobProfile, list of jobs in the queue, etc. Irrespective of this ACL configuration, job-owner, the user who started the cluster, cluster administrators configured via mapreduce.cluster.administrators and queue administrators of the queue to which this job is submitted to configured via mapred.queue.queue-name.acl-administer- jobs in mapred-queue-acls.xml can do all the view operations on a job. By default, nobody else besides job-owner, the user who started the cluster, cluster administrators and queue administrators can perform view operations on a job.		
mapreduce_job_complete_ca ncel_delegation_tokens	If false - do not unregister/cancel delegation tokens from renewal, because same tokens may be used by spawned jobs.	true	true
mapreduce_job_counters_lim it	Limit on the number of counters allowed per job.	true	120
mapreduce_job_split_metainf o_maxsize	The maximum permissible size of the split metainfo file. The JobTracker won't attempt to read split metainfo files bigger than the configured value. No limits if set to -1.	true	10000000
mapreduce_jobtracker_stagin g_root_dir	The root of the staging area for users' job files In practice, this should be the directory where users' home directories are located (usually /user).	true	\${hadoop.tmp.dir}/mapred/staging
mapreduce_reduce_input_lim it	The limit on the input size of the reduce. If the estimated input size of the reduce is greater than this value, job is failed. A value of -1 means that there is no limit set.	true	-1
mapreduce_reduce_shuffle_c onnect_timeout	Expert: The maximum amount of time (in milli seconds) a reduce task spends in trying to connect to a tasktracker for getting map output.	true	180000
mapreduce_reduce_shuffle_r ead_timeout	Expert: The maximum amount of time (in milliseconds) a reduce task waits for map output data to be available for reading after obtaining connection.	true	180000

Name	Description	Required	Default
mapreduce_tasktracker_grou p	Expert: Group to which TaskTracker belongs. If LinuxTaskController is configured via mapreduce.tasktracker.taskcontroller, the group owner of the task-controller binary should be same as this group.	false	
mapreduce_tasktracker_outof band_heartbeat	Expert: Set this to true to let the tasktracker send an out-of-band heartbeat on task-completion for better latency.	true	false
tasktracker_http_threads	The number of worker threads that for the HTTP server. This is used for map output fetching.	true	40

Table 4-27: Hadoop Barclamp Fair Scheduler Parameters

Name	Description	Required	Default
default_min_share_preempti on_timeout	Sets the default minimum share preemption timeout for any pools where it is not specified.	true	600
default_pool_scheduling_mo de	Sets the default scheduling mode (fair or fifo) for pools whose mode is not specified.	true	fair
fair_share_preemption_timeo ut	Sets the preemption timeout used when jobs are below half their fair share.	true	600
pool_max_jobs_default	Sets the default running job limit for any pools whose limit is not specified.	true	20
user_max_jobs_default	Sets the default running job limit for any users whose limit is not specified.	true	10

Pig Barclamp

Apache Pig is a platform for analyzing large data sets that consists of a high-level language for expressing data analysis programs, coupled with infrastructure for evaluating these programs. The salient property of Pig programs is that their structure is amenable to substantial parallelization, which in turns enables them to handle very large data sets.

Pig's infrastructure layer consists of a compiler that produces sequences of MapReduce programs, for which large-scale parallel implementations already exist (e.g., the Hadoop subproject). Pig's language layer currently consists of a textual language called Pig Latin, which has the following key properties:

- Ease of programming. It is trivial to achieve parallel execution of simple, "embarrassingly parallel" data analysis tasks. Complex tasks comprised of multiple interrelated data transformations are explicitly encoded as data flow sequences, making them easy to write, understand, and maintain.
- Optimization opportunities. The way in which tasks are encoded permits the system to optimize their execution automatically, allowing the user to focus on semantics rather than efficiency.
- Extensibility. Users can create their own functions to do special-purpose processing.

Barclamp Parameters

Table 4-34: Pig Barclamp Parameters

Name	Description	Required	Default
java_home	JAVA_HOME environment variable.	true	/usr/java/jdk1.6.0_27/jre
log4jconf	log4jconf log4j configuration file.	true	./conf/log4j.properties
brief	brief logging - no timestamps.	true	false
cluster	Clustername, name of the hadoop jobtracker. If no port is defined port 50020 will be used.	false	
debug_level	Debug level, INFO is default.	true	INFO
file	A file that contains pig script.	false	
jar	Load jarfile, colon separated.	false	
verbose	Verbose print all log messages to screen (default to print only INFO and above to screen).	true	false
exectype	Exectype local or mapreduce - mapreduce is default.	true	mapreduce
ssh_gateway	HOD gateway property.	false	
hod_expect_root	HOD expect root property.	false	
hod_expect_uselatest	HOD use latest root property.	false	
hod_command	HOD command root property.	false	
hod_config_dir	HOD config directory property.	false	
hod_param	HOD param property.	false	
pig_spill_size_threshold	Do not spill temp files smaller than this size (bytes).	true	5000000
pig_spill_gc_activation_size	EXPERIMENT: Activate garbage collection when spilling a file bigger than this size (bytes). This should help reduce the number of files being spilled.	true	40000000
log_file	Log file location.	false	

Hive Barclamp

Hive is a data warehouse system for Hadoop that facilitates easy data summarization, ad-hoc queries, and the analysis of large datasets stored in Hadoop compatible file systems. Hive provides a mechanism to project structure onto this data and query the data using a SQL-like language called HiveQL. This language also allows traditional map/reduce programmers to plug in their custom mappers and reducers when it is inconvenient or inefficient to express this logic in HiveQL.

Barclamp Parameters

Table 4-41: Hive Barclamp Parameters

Name	Description	Required	Default	
hive_exec_scratchdir	Scratch space for Hive jobs.	true	/tmp/hive-\${user.name}	
hive_metastore_local	Controls whether to connect to remove metastore server or open a new metastore server in Hive Client JVM.	true	true	
javax_jdo_option_Connectio nURL	JDBC connect string for a JDBC metastore.	true	jdbc:derby:;databaseName=metasto re_db;create=true	
javax_jdo_option_Connectio nDriverName	Driver class name for a JDBC metastore.	true	org.apache.derby.jdbc.EmbeddedD river	
hive_metastore_metadb_dir	The location of filestore metadata base dir.	true	file:///var/metastore/metadb/	
hive_metastore_uris	Comma separated list of URIs of metastore servers. The first server that can be connected to will be used.	true	file:///var/metastore/metadb/	
hive_metastore_warehouse_d ir	The location of the default database for the warehouse.	true	/user/hive/warehouse	
hive_metastore_connect_retri	Number of retries while opening a connection to metastore.	true	5	
hive_metastore_rawstore_im pl	Name of the class that implements org.apache.hadoop.hive.metastore.rawstore interface. This class is used to store and retrieval of raw metadata objects such as table, database.	true	org.apache.hadoop.hive.metastore. ObjectStore	
hive_default_fileformat	Default file format for CREATE TABLE statement. Options are TextFile and SequenceFile.	true	TextFile	
hive_map_aggr	Whether to use map-side aggregation in Hive Group By queries.	true	false	
hive_join_emit_interval	How many rows in the right-most join operand Hive should buffer before emitting the join result.	true	1000	
hive_exec_script_maxerrsize	Maximum number of bytes a script is allowed to emit to standard error (per mapreduce task). This prevents runaway scripts from filling logs partitions to capacity.	true	100000	
hive_exec_compress_output	Controls whether the final outputs of a query (to a local/hdfs file or a hive table) is compressed. The compression codec and other options are determined from hadoop config variables mapred.output.compress.	true	false	
hive_exec_compress_interme diate	Controls whether intermediate files produced by hive between multiple mapreduce jobs are compressed. The compression codec and other options are determined from hadoop config variables mapred.output.compress.	true	false	

Sqoop Barclamp

SQL based command-line tool to assist with HDFS data import/export (SQL-to-Hadoop). Sqoop is a tool designed to transfer data between Hadoop and relational databases. You can use Sqoop to import data from a relational database management system (RDBMS) such as MySQL or Oracle into the Hadoop Distributed File System (HDFS), transform the data in Hadoop MapReduce, and then export the data back into an RDBMS.

Sqoop automates most of this process, relying on the database to describe the schema for the data to be imported. Sqoop uses MapReduce to import and export the data, which provides parallel operation as well as fault tolerance.

Barclamp Parameters

Table 4-48: Sqoop Barclamp Parameters

Name	Description		Default
sqoop_connection_factories	A comma-delimited list of ManagerFactory implementations which are consulted, in order, to instantiate ConnManager instances used to drive connections to databases.	false	
sqoop_tool_plugins	A comma-delimited list of ToolPlugin implementations which are consulted, in order, to register SqoopTool instances which allow third-party tools to be used.	false	
sqoop_metastore_client_enab le_autoconnect	If true, Sqoop will connect to a local metastore for job management when no other metastore arguments are provided.	true	false
sqoop_metastore_client_auto connect_url	The connect string to use when connecting to a job-management metastore. If unspecified, uses ~/.sqoop/. You can specify a different path here.	false	
sqoop_metastore_client_auto connect_username	The username to bind to the metastore.	false	
sqoop_metastore_client_auto connect_password	The password to bind to the metastore.	false	
sqoop_metastore_client_reco rd_password	If true, allow saved passwords in the metastore.	false	
sqoop_metastore_server_loca tion	Path to the shared metastore database files. If this is not set, it will be placed in ~/.sqoop/.	false	
sqoop_metastore_server_port	Port that this metastore should listen on.	false	

ZooKeeper Barclamp

ZooKeeper is a centralized service for maintaining configuration information, naming, providing distributed synchronization, and providing group services. All of these kinds of services are used in some form or another by distributed applications. Each time they are implemented there is a lot of work that goes into fixing the bugs and race conditions that are inevitable. Because of the difficulty of implementing these kinds of services, applications initially usually skimp on them, which makes them brittle in the presence of change, and difficult to manage. Even when done correctly, different implementations of these services lead to management complexity when the applications are deployed.

ZooKeeper aims to distill the essence of these different services into a very simple interface to a centralized coordination service. The service itself is distributed and highly reliable. Consensus, group management, and presence protocols will be implemented by the service so that the applications do not need to implement them on their own. Application specific uses of these will consist of a mixture of specific components of ZooKeeper and application specific conventions. ZooKeeper recipes show how this simple service can be used to build much more powerful abstractions.

Barclamp Parameters

Table 4-55: ZooKeeper Barclamp Parameters

Name	Description	Required	Default
cluster_name	Provides the ability separate zookeeper services in logical groups.	true	default
tick_time	The number of milliseconds of each tick.	true	2000
init_limit	The number of ticks that the initial synchronization phase can take.	true	10
sync_limit	The number of ticks that can pass between sending a request and getting an acknowledgement.	true	5
client_port	Port at which the clients will connect.	true	2181
peer_port	Server peer port.	true	2888
leader_port	Server leader port.	true	3888
data_dir	Directory where the Zookeeper snapshot is stored.	true	/var/zookeeper
jvm_flags	Increase the heapsize of the ZooKeeper -Server instance to 4GB.	true	-Dzookeeper.log.threshold=INFO - Xmx4G
data_log_dir	Directory where the data log is stored.	false	/var/log/zookeeper

Support

Cloudera Support

To obtain support for Hadoop:

• Open a request at Cloudera's support portal. http://www.cloudera.com/hadoop-support/

Printed in USA

www.dell.com | support.dell.com