

# ZooKeeper Administrator's Guide

## A Guide to Deployment and Administration

by

### Table of contents

1 Deployment.....	2
1.1 System Requirements.....	2
1.2 Clustered (Multi-Server) Setup.....	2
1.3 Single Server and Developer Setup.....	4
2 Administration.....	4
2.1 Configuration Parameters.....	4
2.2 Zookeeper Commands: The Four Letter Words.....	8
2.3 Monitoring.....	9
2.4 Data File Management.....	9
2.5 Things to Avoid.....	10
2.6 Best Practices.....	11

## 1. Deployment

This section contains information about deploying Zookeeper and covers these topics:

- [System Requirements](#)
- [Clustered \(Multi-Server\) Setup](#)
- [Single Server and Developer Setup](#)

The first two sections assume you are interested in installing Zookeeper in a production environment such as a datacenter. The final section covers situations in which you are setting up Zookeeper on a limited basis - for evaluation, testing, or development - but not in a production environment.

### 1.1. System Requirements

Zookeeper runs in Java, release 1.5 or greater, as group of hosts called a quorum. Three Zookeeper hosts per quorum is the minimum recommended quorum size. At Yahoo!, Zookeeper is usually deployed on dedicated RHEL boxes, with dual-core processors, 2GB of RAM, and 80GB IDE harddrives.

### 1.2. Clustered (Multi-Server) Setup

For reliable ZooKeeper service, you should deploy ZooKeeper in a cluster known as a *quorum*. As long as a majority of the quorum are up, the service will be available. Because Zookeeper requires a majority, it is best to use an odd number of machines. For example, with four machines ZooKeeper can only handle the failure of a single machine; if two machines fail, the remaining two machines do not constitute a majority. However, with five machines ZooKeeper can handle the failure of two machines.

Here are the steps to setting a server that will be part of a quorum. These steps should be performed on every host in the quorum:

1. Install the Java JDK:

```
$yinst -i jdk-1.6.0.00_3 -br test
```

2. Set the Java heap size. This is very important, to avoid swapping, which will seriously degrade Zookeeper performance. To determine the correct value, load tests, make sure you are well below the usage limit that would cause you to swap. Be conservative - use a maximum heap size of 3GB for a 4GB machine.
3. Install the Zookeeper Server Package:

```
$ yinst install -nostart zookeeper_server
```

4. Create a configuration file. This file can be called anything. Use the following settings as a starting point:

```
tickTime=2000 dataDir=/var/zookeeper/ clientPort=2181
initLimit=5 syncLimit=2 server.1=zoo1:2888
server.2=zoo2:2888 server.3=zoo3:2888
```

You can find the meanings of these and other configuration settings in the section [Configuration Parameters](#). A word though about a few here:

Every machine that is part of the ZooKeeper quorum should know about every other machine in the quorum. You accomplish this with the series of lines of the form **server.id=host:port**. The integers **host** and **port** are straightforward. You attribute the server id to each machine by creating a file named **myid**, one for each server, which resides in that server's data directory, as specified by the configuration file parameter **dataDir**. The **myid** file consists of a single line containing only the text of that machine's id. So **myid** of server 1 would contain the text "1" and nothing else. The id must be unique within the quorum.

5. If your configuration file is set up, you can start Zookeeper:

```
$ java -cp zookeeper-dev.jar:java/lib/log4j-1.2.15.jar:conf \
\ org.apache.zookeeper.server.quorum.QuorumPeerMain zoo.cfg
```

6. Test your deployment by connecting to the hosts:

- In Java, you can run the following command to execute simple operations:

```
$ java -cp zookeeper.jar:java/lib/log4j-1.2.15.jar:conf \
org.apache.zookeeper.ZooKeeperMain 127.0.0.1:2181
```

- In C, you can compile either the single threaded client or the multithreaded client: or in the **c** subdirectory in the Zookeeper sources. This compiles the single threaded client:

```
$ _make cli_st_
```

And this compiles the multithreaded client:

```
$ _make cli_mt_
```

Running either program gives you a shell in which to execute simple file-system-like operations. To connect to Zookeeper with the multithreaded client, for example, you would run:

```
$ cli_mt 127.0.0.1:2181
```

### 1.3. Single Server and Developer Setup

If you want to setup Zookeeper for development purposes, you will probably want to setup a single server instance of Zookeeper, and then install either the Java or C client-side libraries and bindings on your development machine.

The steps to setting up a single server instance are the similar to the above, except the configuration file is simpler. You can find the complete instructions in the [Installing and Running Zookeeper in Single Server Mode](#) section of the [Zookeeper Getting Started Guide](#).

For information on installing the client side libraries, refer to the [Bindings](#) section of the [Zookeeper Programmer's Guide](#).

## 2. Administration

This section contains information about running and maintaining ZooKeeper and covers these topics:

- [Configuration Parameters](#)
- [Zookeeper Commands: The Four Letter Words](#)
- [Data File Management](#)
- [Things to Avoid](#)
- [Best Practices](#)

### 2.1. Configuration Parameters

ZooKeeper's behavior is governed by the ZooKeeper configuration file. This file is designed so that the exact same file can be used by all the servers that make up a ZooKeeper server assuming the disk layouts are the same. If servers use different configuration files, care must be taken to ensure that the list of servers in all of the different configuration files match.

#### 2.1.1. Minimum Configuration

Here are the minimum configuration keywords that must be defined in the configuration file:

**clientPort**

the port to listen for client connections; that is, the port that clients attempt to connect to.

**dataDir**

the location where Zookeeper will store the in-memory database snapshots and, unless specified otherwise, the transaction log of updates to the database.

**Note:**

Be careful where you put the transaction log. A dedicated transaction log device is key to consistent good performance. Putting the log on a busy device will adversely effect performance.

## **tickTime**

the length of a single tick, which is the basic time unit used by ZooKeeper, as measured in milliseconds. It is used to regulate heartbeats, and timeouts. For example, the minimum session timeout will be two ticks.

### **2.1.2. Advanced Configuration**

The configuration settings in the section are optional. You can use them to further fine tune the behaviour of your Zookeeper servers. Some can also be set using Java system properties, generally of the form *zookeeper.keyword*. The exact system property, when available, is noted below.

#### **dataLogDir**

(No Java system property)

This option will direct the machine to write the transaction log to the **dataLogDir** rather than the **dataDir**. This allows a dedicated log device to be used, and helps avoid competition between logging and snapshots.

**Note:**

Having a dedicated log device has a large impact on throughput and stable latencies. It is highly recommended to dedicate a log device and set **dataLogDir** to point to a directory on that device, and then make sure to point **dataDir** to a directory *not* residing on that device.

#### **globalOutstandingLimit**

(Java system property: **zookeeper.globalOutstandingLimit**.)

Clients can submit requests faster than ZooKeeper can process them, especially if there are a lot of clients. To prevent ZooKeeper from running out of memory due to queued requests, ZooKeeper will throttle clients so that there is no more than **globalOutstandingLimit** outstanding requests in the system. The default limit is 1,000.

#### **preAllocSize**

(Java system property: **zookeeper.preAllocSize**)

To avoid seeks ZooKeeper allocates space in the transaction log file in blocks of `preAllocSize` kilobytes. The default block size is 64M. One reason for changing the size of the blocks is to reduce the block size if snapshots are taken more often. (Also, see **snapCount**).

### **snapCount**

(Java system property: **zookeeper.snapCount**)

Clients can submit requests faster than ZooKeeper can process them, especially if there are a lot of clients. To prevent ZooKeeper from running out of memory due to queued requests, ZooKeeper will throttle clients so that there is no more than `globalOutstandingLimit` outstanding requests in the system. The default limit is 1,000. ZooKeeper logs transactions to a transaction log. After `snapCount` transactions are written to a log file a snapshot is started and a new transaction log file is started. The default `snapCount` is 10,000.

### **traceFile**

(Java system property: **requestTraceFile**)

If this option is defined, requests will be logged to a trace file named `traceFile.year.month.day`. Use of this option provides useful debugging information, but will impact performance. (Note: The system property has no `zookeeper` prefix, and the configuration variable name is different from the system property. Yes - it's not consistent, and it's annoying.)

## **2.1.3. Cluster Options**

The options in this section are designed for use in quorums -- that is, when deploying clusters of servers.

### **electionAlg:**

(No Java system property)

Election implementation to use. A value of "0" corresponds to the original UDP-based version, "1" corresponds to the non-authenticated UDP-based version of fast leader election, "2" corresponds to the authenticated UDP-based version of fast leader election, and "3" corresponds to TCP-based version of fast leader election

### **electionPort**

(No Java system property)

Port used for leader election. It is only used when the election algorithm is not "0". When the election algorithm is "0" a UDP port with the same port number as the port listed in

the **server.num** option will be used.

### **initLimit**

(No Java system property)

Amount of time, in ticks (see [tickTime](#)), to allow followers to connect and sync to a leader. Increased this value as needed, if the amount of data managed by ZooKeeper is large.

### **leaderServes**

(Java system property: `zookeeper.leaderServes`)

Leader accepts client connections. Default value is "yes". The leader machine coordinates updates. For higher update throughput at the slight expense of read throughput the leader can be configured to not accept clients and focus on coordination. The default to this option is yes, which means that a leader will accept client connections.

#### **Note:**

Turning on leader selection is highly recommended when you have more than three Zookeeper servers in a quorum.

### **server.x=[hostname]:nnnn, etc**

(No Java system property)

servers making up the Zookeeper quorum. When the server starts up, it determines which server it is by looking for the file `myid` in the data directory. That file contains the server number, in ASCII, and it should match **x** in **server.x** in the left hand side of this setting.

The list of servers that make up ZooKeeper servers that is used by the clients must match the list of ZooKeeper servers that each ZooKeeper server has.

The port numbers **nnnn** in this setting are the *electionPort* numbers of the servers (as opposed to *clientPorts*). If you want to test multiple servers on a single machine, the individual choices of *electionPort* for each server can be defined in each server's config files using the line `electionPort=xxxx` to avoid clashes.

### **syncLimit**

(No Java system property)

Amount of time, in ticks (see [tickTime](#)), to allow followers to sync with ZooKeeper. If followers fall too far behind a leader, they will be dropped.

## **2.1.4. Unsafe Options**

The following options can be useful, but be careful when you use them. The risk of each is explained along with the explanation of what the variable does.

### **forceSync**

(Java system property: **zookeeper.forceSync**)

Requires updates to be synced to media of the transaction log before finishing processing the update. If this option is set to no, ZooKeeper will not require updates to be synced to the media.

### **jute.maxbuffer:**

(Java system property: **jute.maxbuffer**)

This option can only be set as a Java system property. There is no zookeeper prefix on it. It specifies the maximum size of the data that can be stored in a znode. The default is 0xfffff, or just under 1M. If this option is changed, the system property must be set on all servers and clients otherwise problems will arise. This is really a sanity check. ZooKeeper is designed to store data on the order of kilobytes in size.

### **skipACL**

(Java system property: **zookeeper.skipACL**)

Skips ACL checks. This results in a boost in throughput, but opens up full access to the data tree to everyone.

## **2.2. Zookeeper Commands: The Four Letter Words**

Zookeeper responds to a small set of commands. Each command is composed of four letters. You issue the commands to Zookeeper via telnet or nc, at the client port.

### **dump**

Lists the outstanding sessions and ephemeral nodes. This only works on the leader.

### **kill**

Shuts down the server. This must be issued from the machine the Zookeeper server is running on.

### **ruok**

Tests if server is running in a non-error state. The server will respond with imok if it is running. Otherwise it will not respond at all.

### **stat**

Lists statistics about performance and connected clients.



Here's an example of the **ruok** command:

```
$ echo ruok | nc 127.0.0.1 5111 imok
```

## 2.3. Monitoring

*[tbd]*

## 2.4. Data File Management

ZooKeeper stores its data in a data directory and its transaction log in a transaction log directory. By default these two directories are the same. The server can (and should) be configured to store the transaction log files in a separate directory than the data files. Throughput increases and latency decreases when transaction logs reside on a dedicated log devices.

### 2.4.1. The Data Directory

This directory has two files in it:

- `myid` - contains a single integer in human readable ASCII text that represents the server id.
- `snapshot.<zxid>` - holds the fuzzy snapshot of a data tree.

Each ZooKeeper server has a unique id. This id is used in two places: the `myid` file and the configuration file. The `myid` file identifies the server that corresponds to the given data directory. The configuration file lists the contact information for each server identified by its server id. When a ZooKeeper server instance starts, it reads its id from the `myid` file and then, using that id, reads from the configuration file, looking up the port on which it should listen.

The `snapshot` files stored in the data directory are fuzzy snapshots in the sense that during the time the ZooKeeper server is taking the snapshot, updates are occurring to the data tree. The suffix of the `snapshot` file names is the `zxid`, the ZooKeeper transaction id, of the last committed transaction at the start of the snapshot. Thus, the snapshot includes a subset of the updates to the data tree that occurred while the snapshot was in process. The snapshot, then, may not correspond to any data tree that actually existed, and for this reason we refer to it as a fuzzy snapshot. Still, ZooKeeper can recover using this snapshot because it takes advantage of the idempotent nature of its updates. By replaying the transaction log against fuzzy snapshots ZooKeeper gets the state of the system at the end of the log.

### 2.4.2. The Log Directory

The Log Directory contains the ZooKeeper transaction logs. Before any update takes place, ZooKeeper ensures that the transaction that represents the update is written to non-volatile storage. A new log file is started each time a snapshot is begun. The log file's suffix is the first zxid written to that log.

### 2.4.3. File Management

The format of snapshot and log files does not change between standalone ZooKeeper servers and different configurations of replicated ZooKeeper servers. Therefore, you can pull these files from a running replicated ZooKeeper server to a development machine with a stand-alone ZooKeeper server for trouble shooting.

Using older log and snapshot files, you can look at the previous state of ZooKeeper servers and even restore that state. The LogFormatter class allows an administrator to look at the transactions in a log.

The ZooKeeper server creates snapshot and log files, but never deletes them. The retention policy of the data and log files is implemented outside of the ZooKeeper server. The server itself only needs the latest complete fuzzy snapshot and the log files from the start of that snapshot. The PurgeTxnLog utility implements a simple retention policy that administrators can use.

## 2.5. Things to Avoid

Here are some common problems you can avoid by configuring ZooKeeper correctly:

### **inconsistent lists of servers**

The list of Zookeeper servers used by the clients must match the list of ZooKeeper servers that each ZooKeeper server has. Things work okay if the client list is a subset of the real list, but things will really act strange if clients have a list of ZooKeeper servers that are in different ZooKeeper clusters. Also, the server lists in each Zookeeper server configuration file should be consistent with one another.

### **incorrect placement of transaction log**

The most performance critical part of ZooKeeper is the transaction log. Zookeeper syncs transactions to media before it returns a response. A dedicated transaction log device is key to consistent good performance. Putting the log on a busy device will adversely effect performance. If you only have one storage device, put trace files on NFS and increase the snapshotCount; it doesn't eliminate the problem, but it should mitigate it.

### **incorrect Java heap size**

You should take special care to set your Java max heap size correctly. In particular, you should not create a situation in which Zookeeper swaps to disk. The disk is death to ZooKeeper. Everything is ordered, so if processing one request swaps the disk, all other queued requests will probably do the same. the disk. DON'T SWAP.

Be conservative in your estimates: if you have 4G of RAM, do not set the Java max heap size to 6G or even 4G. For example, it is more likely you would use a 3G heap for a 4G machine, as the operating system and the cache also need memory. The best and only recommend practice for estimating the heap size your system needs is to run load tests, and then make sure you are well below the usage limit that would cause the system to swap.

## **2.6. Best Practices**

For best results, take note of the following list of good Zookeeper practices. *[tbd...]*