



University of  
South Australia

## HDFS in detail

1



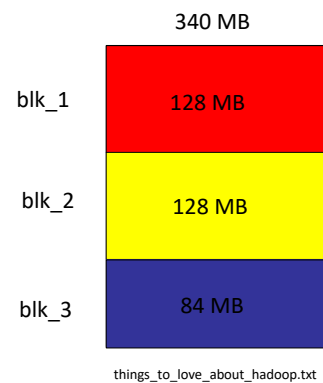
## Block size

128 MB vs 64 MB?

What are the pros and cons of having small or large blocks?

**Too large** – processing time increases for each block

**Too small** – Too many data blocks, lots of metadata, seeking blocks takes time



2



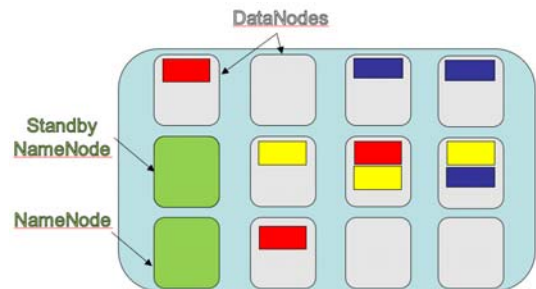
## Replication factor

Hadoop's default replication factor is 3

What are the pros and cons of having small or large replication factors?

**Larger**– Safer, requires more storage, almost always unnecessary.

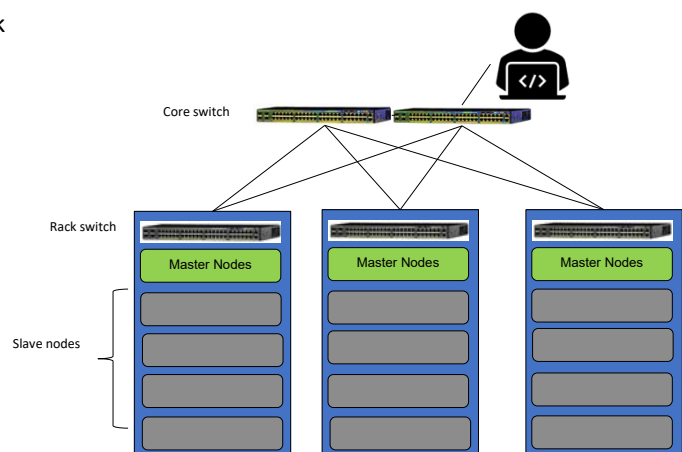
**Smaller**– Faster, requires less space, but dangerous



## Data centre racks

A 'rack' in a data centre is a physical framework used to store a set of computing equipment.

In our case the racks would consist of nodes in our cluster. That is, the machines used for our distributed storage and processing.



Yahoo's massive (!!) Hadoop cluster



## Rack awareness

Hi, I'm the client. I want to store blk\_1 of file01.txt



## Rack awareness

Hi, I'm the client. I want to store blk\_1 of file01.txt

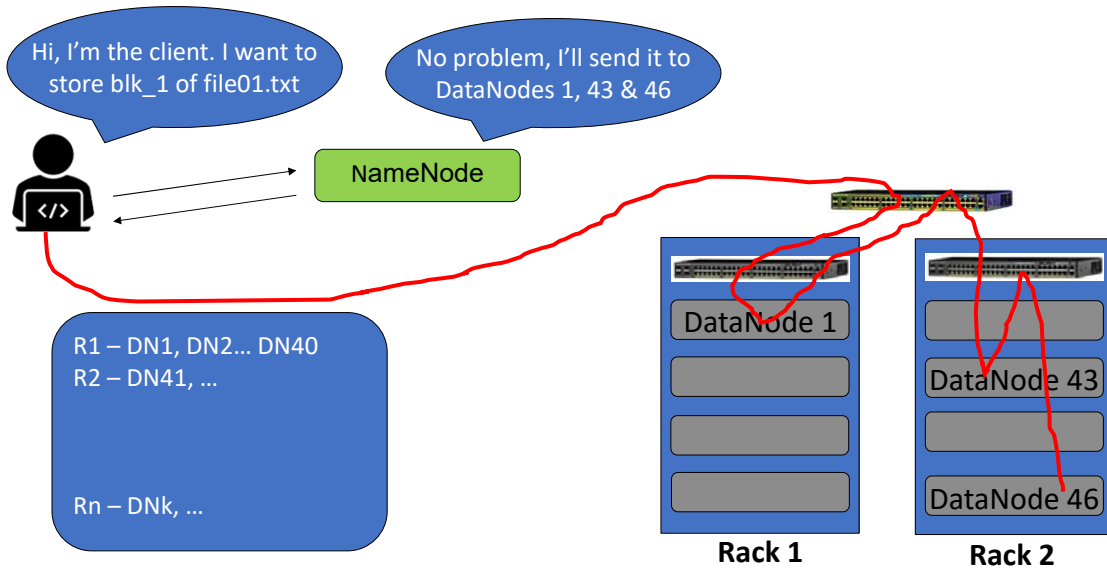
No problem, I'll send it to DataNodes 1, 43 & 46



NameNode

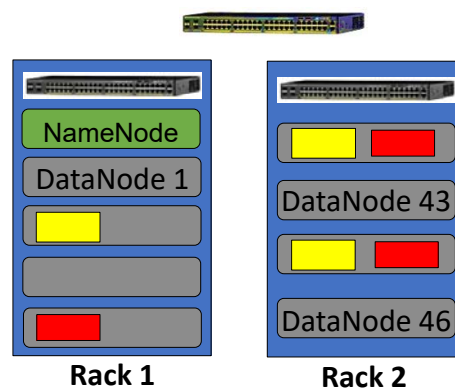


## Rack awareness



## Heartbeat

NameNode → DataNode



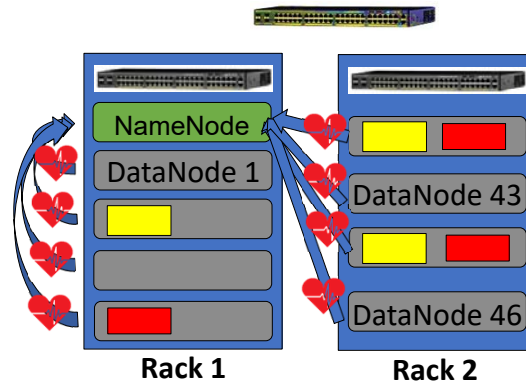


# Heartbeat

NameNode ← DataNode

By default a pulse is sent from the nameNode to the dataNode every three seconds

Once again... that value can be adjusted depending on context



# How many parameters can we change?

<code>dfs.namenode.backup.address</code>	<code>0.0.0.0:50100</code>	The backup node server address and port. If the port is 0 then the server will start on a free port.
<code>dfs.namenode.backup.tmp.address</code>	<code>0.0.0.0:50103</code>	The backup node tmp server address and port. If the port is 0 then the server will start on a free port.
<code>dfs.namenode.replication.considerLoad</code>	<code>true</code>	Decide if chooseTarget considers the target's load or not.
<code>dfs.default.chunk.view.size</code>	<code>32768</code>	The number of bytes to view for a file on the browser.
<code>dfs.datanode.du.reserved</code>	<code>0</code>	Reserved space in bytes per volume. Always leave this much space free for non-DFS use.
<code>dfs.namenode.name.dir</code>	<code>file:/5(hadoop.tmp.dir)/dfs/name</code>	Determines where on the local filesystem the DFS name node should store the name table(snapshots). If this is a comma-delimited list of directories, the first directory is used for the primary namespace, and the others are used for redundancy.
<code>dfs.namenode.name.dir.restore</code>	<code>false</code>	Set to true to enable NameNode to attempt recovering a previously failed <code>dfs.namenode.name.dir</code> . When enabled, it will attempt to recover the namespace during checkpoint.
<code>dfs.namenode.fs-limits.max-component-length</code>	<code>255</code>	Defines the maximum number of bytes in UTF-8 encoding in each component of a path. A value of 0 will disable this limit.
<code>dfs.namenode.fs-limits.max-directory-items</code>	<code>1048576</code>	Defines the maximum number of items that a directory may contain. Cannot set the property to a value less than 1.
<code>dfs.namenode.fs-limits.min-block-size</code>	<code>1048576</code>	Minimum block size in bytes, enforced by the NameNode at create time. This prevents the accidental creation of files which can degrade performance.
<code>dfs.namenode.fs-limits.max-blocks-per-file</code>	<code>1048576</code>	Maximum number of blocks per file, enforced by the NameNode on write. This prevents the creation of extremely large files.
<code>dfs.namenode.edits.dir</code>	<code>\$(dfs.namenode.name.dir)</code>	Determines where on the local filesystem the DFS name node should store the transaction (edits) file. If this is a comma-delimited list of directories, the first directory is used for the primary namespace, and the others are used for redundancy.
<code>dfs.namenode.shared.edits.dir</code>	<code>\$(dfs.namenode.name.dir)</code>	A directory on shared storage between the multiple namenodes in an HA cluster. This directory will be written by the primary namenode and read by the standbys. The namespace is synchronized. This directory does not need to be listed in <code>dfs.namenode.edits.dir</code> above. It should be listed in <code>dfs.namenode.name.dir</code> .
<code>dfs.namenode.edits.journal.plugin</code>	<code>org.apache.hadoop.hdfs.qjournal.client.QuorumJournalManager</code>	The plugin to use for the Quorum Journal Manager.
<code>dfs.permissions.enabled</code>	<code>true</code>	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.
<code>dfs.permissions.supergroup</code>	<code>supergroup</code>	The name of the group of super users.
<code>dfs.namenode.acls.enabled</code>	<code>false</code>	Set to true to enable support for HDFS ACLs (Access Control Lists). By default, ACLs are disabled. When ACLs are disabled, the NameNode rejects all RPCs related to setting or getting ACLs.
<code>dfs.namenode.lazypersist.file.scrub.interval.sec</code>	<code>300</code>	The NameNode periodically scans the namespace for LazyPersist files with missing blocks and unlinks them from the namespace. This configuration key controls the interval between successive scans. Set it to a negative value to disable this behavior.
<code>dfs.block.access.token.enable</code>	<code>false</code>	If "true", access tokens are used as capabilities for accessing datanodes. If "false", no access tokens are checked on accessing datanodes.
<code>dfs.block.access.key.update.interval</code>	<code>600</code>	Interval in minutes at which namenode updates its access keys.
<code>dfs.block.access.token.lifetime</code>	<code>600</code>	The lifetime of access tokens in minutes.
<code>dfs.datanode.data.dir</code>	<code>file:/5(hadoop.tmp.dir)/dfs/data</code>	Determines where on the local filesystem the 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.
<code>dfs.datanode.data.dir.perm</code>	<code>700</code>	Permissions for the directories on the local filesystem where the DFS data node store its blocks. The permissions can be either octal or symbolic.
<code>dfs.replication</code>	<code>3</code>	The 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.
<code>dfs.replication.max</code>	<code>512</code>	Maximal block replication.
<code>dfs.namenode.replication.max</code>	<code>512</code>	Maximal block replication.
<code>dfs.blocksize</code>	<code>134217728</code>	The default block size for new files, in bytes. You can use the following suffix (case insensitive): k(kilo), m(mega), g(giga), t(tera), p(peta), e(exa) to specify the size (such as 1KB, 512m, 1g, etc.). Or provide complete size in bytes (such as 134217728 for 128 MB).
<code>dfs.client.block.write.retry</code>	<code>3</code>	The number of retries for writing blocks to the data nodes, before we signal failure to the application.
<code>dfs.client.block.write.replace-datanode-on-failure.enable</code>	<code>true</code>	If there is a datanode network failure in the write pipeline, DFSClient will try to remove the failed datanode from the pipeline and then continue writing with the remaining datanodes. As a result, the number of datanodes in the pipeline is decreased. The feature is to add new datanodes to the pipeline. This is a site-wide property to enable/disable the feature. When the cluster size is extremely small, e.g. 3 nodes or less, cluster administrators may want to set the policy to NEVER in the default configuration file or disable this feature. Otherwise, users may experience an unusually high rate of pipeline failures since it is impossible to find new datanodes for replacement. See also <code>dfs.client.block.write.replace-datanode-on-failure.policy</code> .

There's roughly 8 more tables of this size, full of values that we can change

**WARNING**

This material has been reproduced and communicated to you by or on behalf of the **University of South Australia** in accordance with section 113P of the *Copyright Act 1968* (**Act**).

The material in this communication may be subject to copyright under the Act. Any further reproduction or communication of this material by you may be the subject of copyright protection under the Act.

**Do not remove this notice**