

# Hadoop Distributed File System (HDFS)

## Overview

- Provides inexpensive and reliable storage for massive amounts of data.
- Designed for:
  - o Big files (100 MB to several TBs file sizes).
  - o Write once, read many times (Appending only).
  - o Running on commodity hardware.
- Hierarchical UNIX style file systems.
  - o UNIX style file ownership and permissions.

## HDFS main design principles

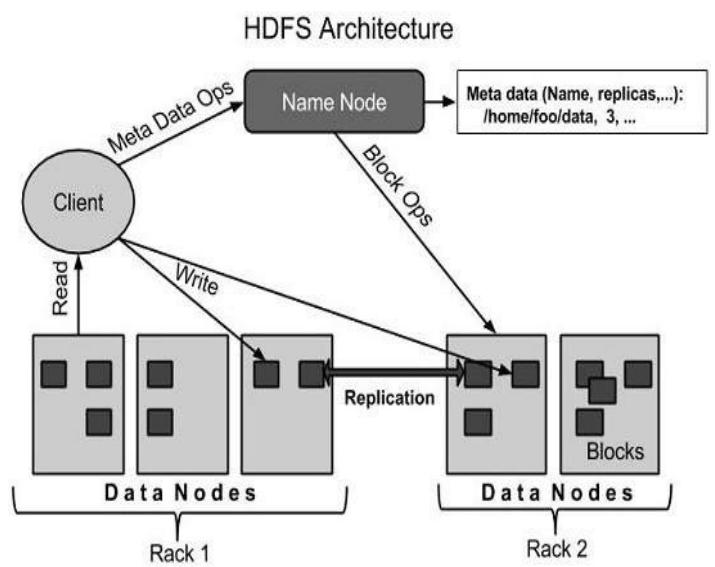
- I/O pattern:
  - o Append only → reduce synchronization.
- Data distribution:
  - o File is split in big chunks (64 MB).
    - Reduce metadata size.
    - Reduce network communication.
- Data replication:
  - o Each chunk is usually replicated in 3 different nodes.
- Fault tolerance:
  - o Data node: re-replication.
  - o Name node:
    - Secondary name node.
    - Standby, active name node.

## HDFS Architecture

- Master/slave architecture.
- HDFS master: Name node.
  - o Manage namespace and metadata.
  - o Monitor data nodes.
- HDFS slaves: Data nodes.
  - o Handle read/write the actual data (chunks).
  - o Chunks are local files in the local file systems.

## Functions of a Name Node

- Manage File System Namespace.
  - o Maps a file name to a set of blocks.
  - o Maps a block to the Data nodes where it resides.



- Cluster Configuration Management.
- Replication Engine for Blocks.

## Name Node metadata

- Metadata in memory.
  - o The entire metadata is in main memory.
  - o No demand paging of metadata.
- Types of metadata:
  - o List of files.
  - o List of Blocks for each file.
  - o List of Data nodes for each block.
  - o File attributes, e.g., creation time, replication factor.
- A Transaction Log:
  - o Records file creations, file deletions, etc.

## Data Node

- A Block Server:
  - o Stores data in the local file system (e.g., ext3).
  - o Stores metadata of a block (e.g., CRC).
  - o Serves data and metadata to Clients.
- Block Report:
  - o Periodically sends a report of all existing blocks to the Name node.
- Facilitates Pipelining of Data:
  - o Forwards data to other specified Data node.
- Heartbeat:
  - o Data nodes send heartbeat to the Name node.
    - Once every 3 seconds.
  - o Name node uses heartbeats to detect Data nodes failure.

## Data Replication

- Chunk replacement.
  - o Current Strategy.
    - One replica on local node.
    - Second replica on remote rack.
    - Third replica on the same remote rack.
    - Additional replicas are randomly placed.
  - o Clients read from nearest replicas.
- Name node detects Data node failures:
  - o Chooses new Data nodes for new replicas.
  - o Balances disk usage.
  - o Balances communication traffic to Data nodes.

## Data Rebalance

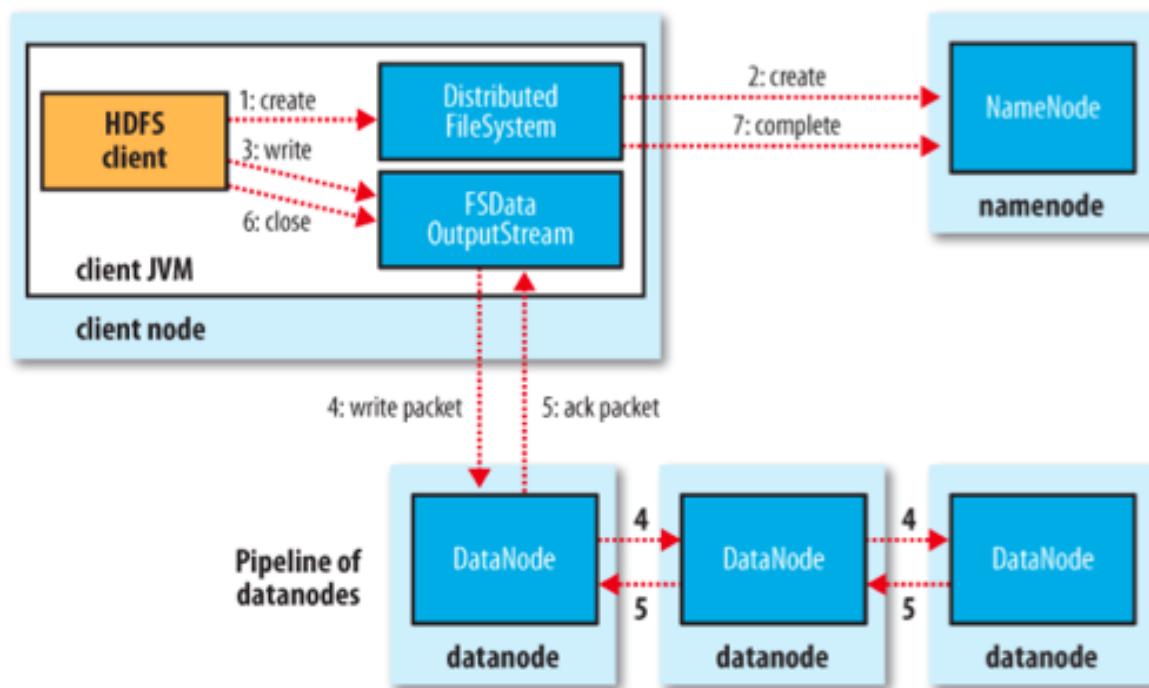
- Goal: % disk full on Data nodes should be similar.
  - o Usually run when new Data nodes are added.
  - o Cluster is online when Rebalancer is active.
  - o Rebalancer is throttled to avoid network congestion.
  - o Command line tool.

## Data Correctness

- Use Checksums to validate data.
  - o CRC32.
- File Creation:
  - o Client computes checksum per 512 bytes.
  - o Data node stores the checksum.
- File Access:
  - o Client retrieves the data and checksum from Data node.
  - o If validation fails, Client tries other replicas.

## Data pipelining

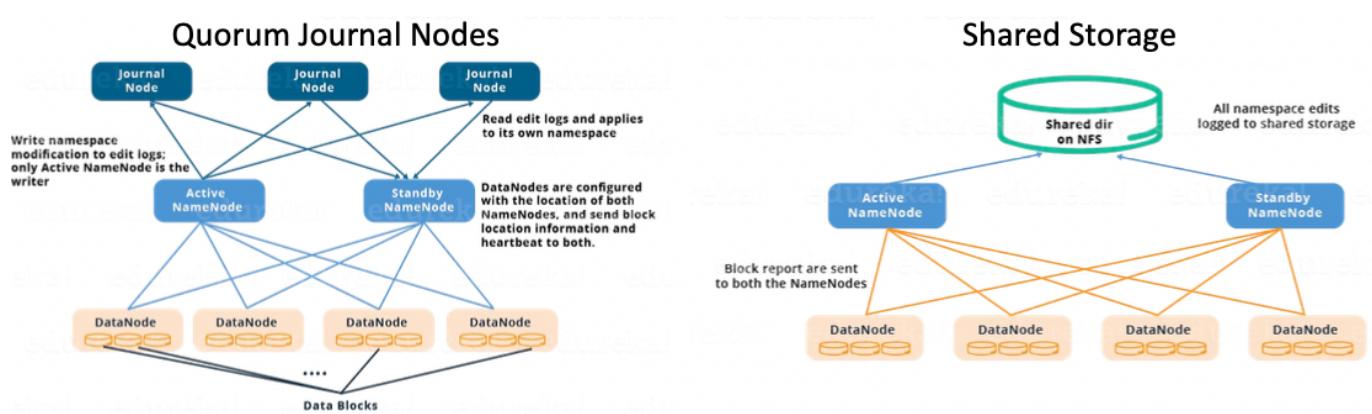
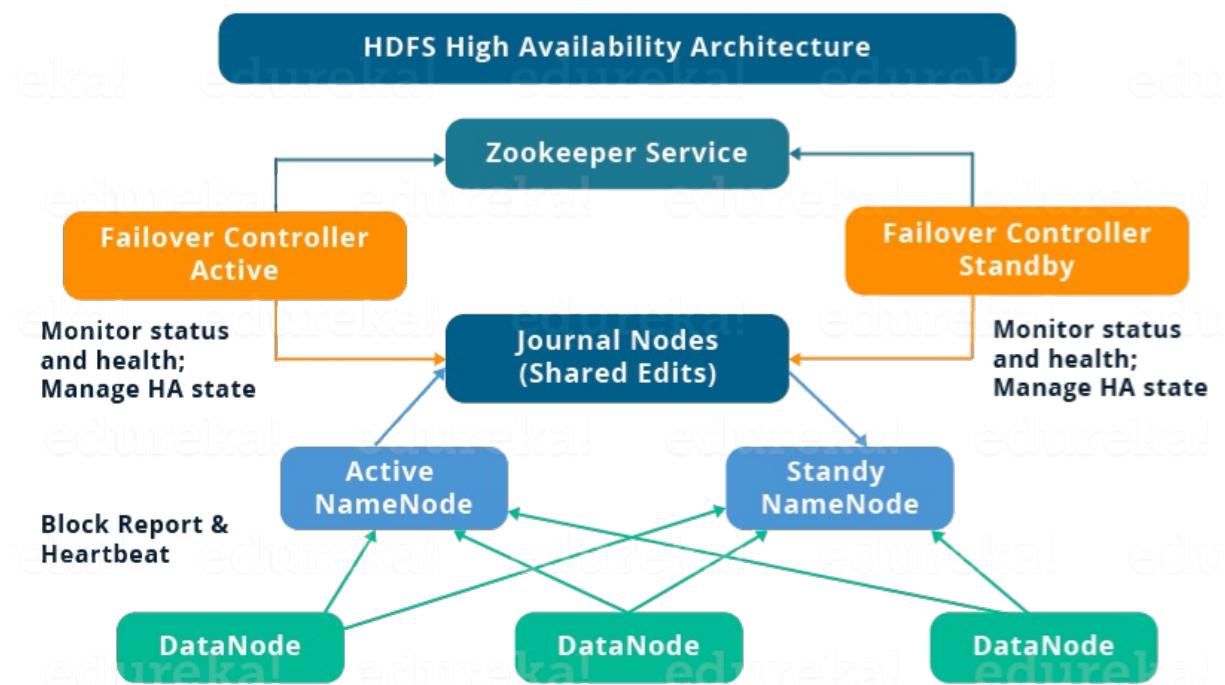
- Client retrieves a list of Data nodes on which to replace replicas of a block.
- Client writes block to the first Data node.
- The first Data node forwards the data to the next node in the Pipeline.
- When all replicas are written, the Client moves on to write the next block in file.



## Secondary Name node

- Name node is a single point of failure.
- Secondary Name node:
  - o Checkpointing latest copy of the Fs Image and the Transaction Log Files.
    - Copies Fs Image and Transaction Log from Name node to a temporary directory.
- When Name node restarted:
  - o Merges Fs Image and Transaction Log into a new Fs Image in temporary directory.
  - o Uploads new Fs Image to the Name node.
    - Transaction Log on Name node is purged.

## Name node high availability (HA)



## HDFS command-line interface

List Files	
hdfs dfs -ls /	List all the files/directories for the given hdfs destination path.
hdfs dfs -ls -d /hadoop	Directories are listed as plain files. In this case, this command will list the details of hadoop folder.
hdfs dfs -ls -h /data	Format file sizes in a human-readable fashion (eg 64.0m instead of 67108864).
hdfs dfs -ls -R /hadoop	Recursively list all files in hadoop directory and all subdirectories in hadoop directory.
hdfs dfs -ls /hadoop/dat*	List all the files matching the pattern. In this case, it will list all the files inside hadoop directory which starts with 'dat'.

Read/Write Files	
hdfs dfs -text /hadoop/derby.log	HDFS Command that takes a source file and outputs the file in text format on the terminal. The allowed formats are zip and TextRecordInputStream.
hdfs dfs -cat /hadoop/test	This command will display the content of the HDFS file test on your stdout .
hdfs dfs -appendToFile /home/ubuntu/test1 /hadoop/text2	Appends the content of a local file test1 to a hdfs file test2.

## Upload, download files

Upload/Download Files	
hdfs dfs -put /home/ubuntu/sample /hadoop	Copies the file from local file system to HDFS.
hdfs dfs -put -f /home/ubuntu/sample /hadoop	Copies the file from local file system to HDFS, and in case the local already exists in the given destination path, using -f option with put command will overwrite it.
hdfs dfs -put -l /home/ubuntu/sample /hadoop	Copies the file from local file system to HDFS. Allow DataNode to lazily persist the file to disk. Forces replication factor of 1.
hdfs dfs -put -p /home/ubuntu/sample /hadoop	Copies the file from local file system to HDFS. Passing -p preserves access and modification times, ownership and the mode.
hdfs dfs -get /newfile /home/ubuntu/	Copies the file from HDFS to local file system.
hdfs dfs -get -p /newfile /home/ubuntu/	Copies the file from HDFS to local file system. Passing -p preserves access and modification times, ownership and the mode.
hdfs dfs -get /hadoop/*.txt /home/ubuntu/	Copies all the files matching the pattern from local file system to HDFS.
hdfs dfs -copyFromLocal /home/ubuntu/sample /hadoop	Works similarly to the put command, except that the source is restricted to a local file reference.
hdfs dfs -copyToLocal /newfile /home/ubuntu/	Works similarly to the put command, except that the destination is restricted to a local file reference.
hdfs dfs -moveFromLocal /home/ubuntu/sample /hadoop	Works similarly to the put command, except that the source is deleted after it's copied.

## File management

File Management	
hdfs dfs -cp /hadoop/file1 /hadoop1	Copies file from source to destination on HDFS. In this case, copying file1 from hadoop directory to hadoop1 directory.
hdfs dfs -cp -p /hadoop/file1 /hadoop1	Copies file from source to destination on HDFS. Passing -p preserves access and modification times, ownership and the mode.
hdfs dfs -cp -f /hadoop/file1 /hadoop1	Copies file from source to destination on HDFS. Passing -f overwrites the destination if it already exists.
hdfs dfs -mv /hadoop/file1 /hadoop1	Move files that match the specified file pattern <src> to a destination <dst>. When moving multiple files, the destination must be a directory.
hdfs dfs -rm /hadoop/file1	Deletes the file (sends it to the trash).
hdfs dfs -rm -r /hadoop hdfs dfs -rm -R /hadoop hdfs dfs -rmr /hadoop	Deletes the directory and any content under it recursively.
hdfs dfs -rm -skipTrash /hadoop	The -skipTrash option will bypass trash, if enabled, and delete the specified file(s) immediately.
hdfs dfs -rm -f /hadoop	If the file does not exist, do not display a diagnostic message or modify the exit status to reflect an error.
hdfs dfs -rmdir /hadoop1	Delete a directory.
hdfs dfs -mkdir /hadoop2	Create a directory in specified HDFS location.
hdfs dfs -mkdir -f /hadoop2	Create a directory in specified HDFS location. This command does not fail even if the directory already exists.
hdfs dfs -touchz /hadoop3	Creates a file of zero length at <path> with current time as the timestamp of that <path>.

## Ownership and validation

Ownership and Validation	
hdfs dfs -checksum /hadoop/file1	Dump checksum information for files that match the file pattern <src> to stdout.
hdfs dfs -chmod 755 /hadoop/file1	Changes permissions of the file.
hdfs dfs -chmod -R 755 /hadoop	Changes permissions of the files recursively.
hdfs dfs -chown ubuntu:ubuntu /hadoop	Changes owner of the file. 1st ubuntu in the command is owner and 2nd one is group.
hdfs dfs -chown -R ubuntu:ubuntu /hadoop	Changes owner of the files recursively.
hdfs dfs -chgrp ubuntu /hadoop	Changes group association of the file.
hdfs dfs -chgrp -R ubuntu /hadoop	Changes group association of the files recursively.
Filesystem	
hdfs dfs -df /hadoop	Shows the capacity, free and used space of the filesystem.
hdfs dfs -df -h /hadoop	Shows the capacity, free and used space of the filesystem. -h parameter Formats the sizes of files in a human-readable fashion.
hdfs dfs -du /hadoop/file	Show the amount of space, in bytes, used by the files that match the specified file pattern.
hdfs dfs -du -s /hadoop/file	Rather than showing the size of each individual file that matches the pattern, shows the total (summary) size.
hdfs dfs -du -h /hadoop/file	Show the amount of space, in bytes, used by the files that match the specified file pattern. Formats the sizes of files in a human-readable fashion.

## Administration

Administration	
<b>hdfs balancer -threshold 30</b>	Runs a cluster balancing utility. Percentage of disk capacity. This overwrites the default threshold.
<b>hadoop version</b>	To check the version of Hadoop.
<b>hdfs fsck /</b>	It checks the health of the Hadoop file system.
<b>hdfs dfsadmin -safemode leave</b>	The command to turn off the safemode of NameNode.
<b>hdfs dfsadmin -refreshNodes</b>	Re-read the hosts and exclude files to update the set of Datanodes that are allowed to connect to the Namenode and those that should be decommissioned or recommissioned.
<b>hdfs namenode -format</b>	Formats the NameNode.

## HDFS Name Node UI

Hadoop	Overview	Datanodes	Datanode Volume Failures	Snapshot	Startup Progress	Utilities ▾
--------	----------	-----------	--------------------------	----------	------------------	-------------

### Overview 'hd01:8020' (active)

Started:	Thu Mar 14 11:01:37 +0700 2019
Version:	3.1.1.3.1.0.0-78, re4f82af51faec922b4804d0232a637422ec29e64
Compiled:	Thu Dec 06 20:34:00 +0700 2018 by jenkins from (HEAD detached at e4f82af)
Cluster ID:	CID-a1dea38d-b6cf-44e4-b5e8-3bde87ffad35
Block Pool ID:	BP-1412866890-10.10.137.41-1544787807355

### Summary

Security is off.  
Safemode is off.  
165,739 files and directories, 99,858 blocks (99,858 replicated blocks, 0 erasure coded block groups) = 265,597 total filesystem object(s).  
Heap Memory used 250.76 MB of 1011.25 MB Heap Memory. Max Heap Memory is 1011.25 MB.

### In operation

In operation								
Show	25	entries	Search:					
Node	Http Address	Last contact	Last Block Report	Capacity	Blocks	Block pool used	Version	
✓hd01:50010 (10.10.137.41:50010)	<a href="http://hd01:50075">http://hd01:50075</a>	1s	251m	296.83 GB	<div style="width: 30%; background-color: orange;"></div>	99858	166.56 GB (56.11%)	3.1.1.3.1.0.0-78
✓hd02:50010 (10.10.137.42:50010)	<a href="http://hd02:50075">http://hd02:50075</a>	0s	57m	296.83 GB	<div style="width: 30%; background-color: green;"></div>	99858	166.57 GB (56.12%)	3.1.1.3.1.0.0-78
✓hd03:50010 (10.10.137.43:50010)	<a href="http://hd03:50075">http://hd03:50075</a>	2s	197m	296.83 GB	<div style="width: 30%; background-color: green;"></div>	99858	166.57 GB (56.12%)	3.1.1.3.1.0.0-78

Showing 1 to 3 of 3 entries

Previous 1 Next

## Browse Directory

/									Go!			
Show 25 entries									Search:			
	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name				
<input type="checkbox"/>	drwxrwxrwt	yarn	hadoop	0 B	Jun 01 17:52	0	0 B	app-logs				
<input type="checkbox"/>	drwxr-xr-x	hdfs	hdfs	0 B	Dec 18 2018	0	0 B	apps				
<input type="checkbox"/>	drwxr-xr-x	yarn	hadoop	0 B	Dec 14 2018	0	0 B	ats				
<input type="checkbox"/>	drwxr-xr-x	hdfs	hdfs	0 B	Dec 14 2018	0	0 B	atsv2				

## Other HDFS interfaces

- Java API.
- Thrift API.
- Fuse.
- WebDAV.

## HDFS data format

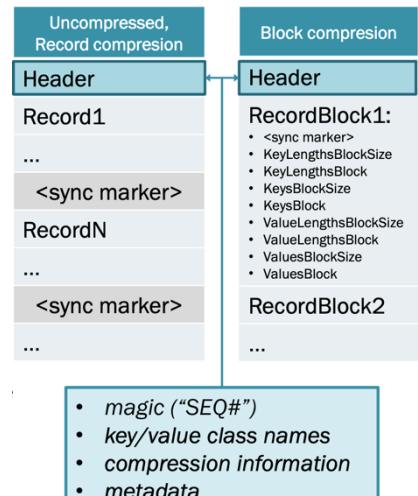
- Text.
- Sequence file.
- Avro.
- Parquet.
- Optimized Row Columnar (ORC).

## Text file

- CSV, TSV, JSON records.
- Convenient format to use to Exchange between applications or scripts.
- Human readable and parse able.
- Do not support block compression.
- Not as efficient to query.
- Good for the beginning, but not good enough for real life.

## Sequence file

- Provides a persistent data structure for binary key-value pairs.
- Commonly used to transfer data between Map Reduce jobs.
- Can be used as an archive to pack small files in Hadoop.
- Row-based.
- Compression.
- Splittable.
  - o Support splitting even when the data is compressed.



## Avro

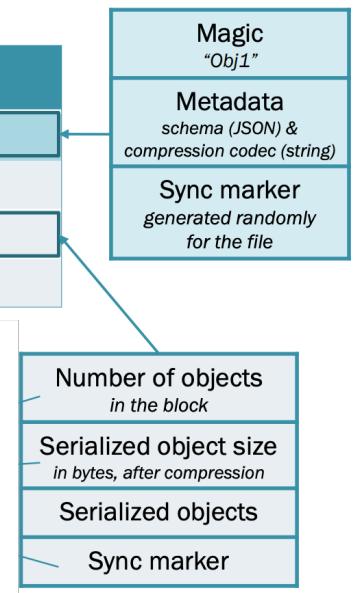
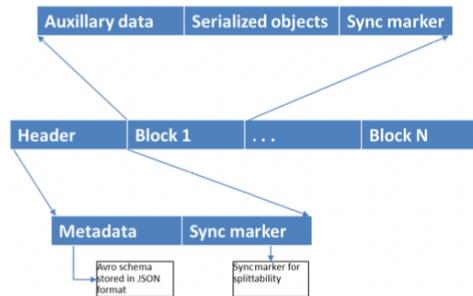
- Row-based.
- Supports (object) compression and splitting.
- Flexible data scheme.
  - o Scheme (JSON) included to the file.
- Data types:
  - o Primitive: null, Boolean, int, long, etc.
  - o Complex: records, arrays, maps, etc.
- Binary and JSON data serialization.
- Data corruption detection.

## Avro – File structure and example

Sample AVRO schema in JSON format

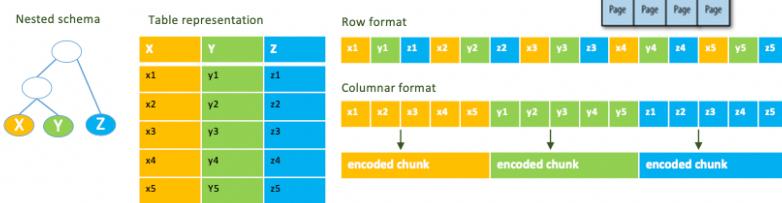
```
{
  "type" : "record",
  "name" : "tweets",
  "fields" : [ {
    "name" : "username",
    "type" : "string",
  }, {
    "name" : "tweet",
    "type" : "string",
  }, {
    "name" : "timestamp",
    "type" : "long",
  }],
  "doc:" : "schema for storing tweets"
}
```

Avro file structure

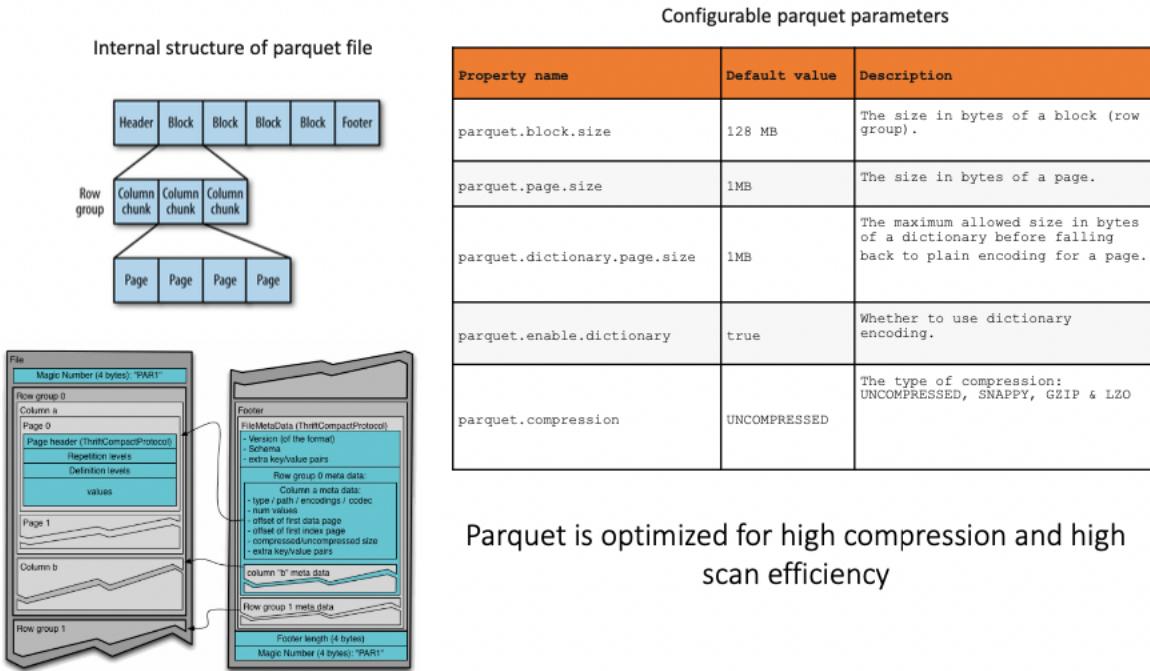


## Parquet

- Column-oriented binary file format
- Efficient in terms of disk I/O when specific columns need to be queried
- Supports (page) compression and splitting
- Supports nested columns (Dremel encoding)

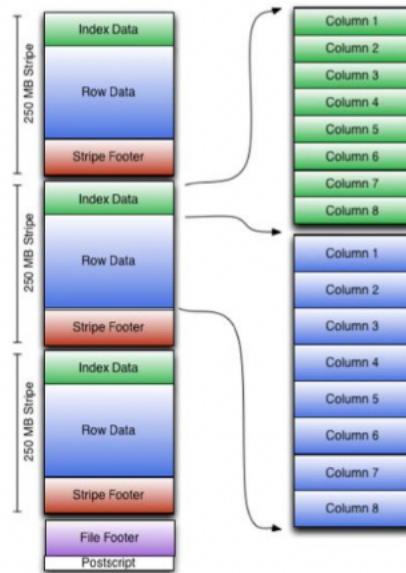


## Parquet file structure and configuration



## Optimized row columnar (ORC)

- **RCFile**
  - Every column is compressed individually within the row group
- **ORC File**
  - Block-mode compression
  - Data type support
  - Ordered data store (within one stripe)
- Stores collections of rows and within the collection the row data is stored in columnar format
- Introduces a lightweight indexing that enables skipping of irrelevant blocks of rows
- Splittable: allows parallel processing of row collections
- Indices with column-level aggregated values (min, max, sum and count)



## Table of Contents

<i>Overview</i> .....	1
<i>HDFS main design principles</i> .....	1
<i>HDFS Architecture</i> .....	1
<i>Functions of a Name Node</i> .....	1
<i>Name Node metadata</i> .....	2
<i>Data Node</i> .....	2
<i>Data Replication</i> .....	2
<i>Data Rebalance</i> .....	3
<i>Data Correctness</i> .....	3
<i>Data pipelining</i> .....	3
<i>Secondary Name node</i> .....	4
<i>Name node high availability (HA)</i> .....	4
<i>HDFS command-line interface</i> .....	5
<i>Upload, download files</i> .....	5
<i>File management</i> .....	6
<i>Ownership and validation</i> .....	6
<i>Administration</i> .....	7
<i>HDFS Name Node UI</i> .....	7
<i>Other HDFS interfaces</i> .....	8
<i>HDFS data format</i> .....	8
<i>Text file</i> .....	8
<i>Sequence file</i> .....	8
<i>Avro</i> .....	9
<i>Arvo – File structure and example</i> .....	9
<i>Parquet</i> .....	9
<i>Parquet file structure and configuration</i> .....	10
<i>Optimized row columnar (ORC)</i> .....	10