# Apache Hadoop HDFS

**A distributed Java-based file system for storing large volumes of data.**

- HDFS and YARN form the data management layer of Apache Hadoop.
- YARN is the architectural center of Hadoop, the resource management framework that enables the enterprise to process data in multiple ways simultaneously—for batch, interactive and real-time data workloads on one shared dataset.
- YARN provides the resource management and HDFS provides the scalable, fault-tolerant, cost-efficient storage for big data.

# HDFS - Features

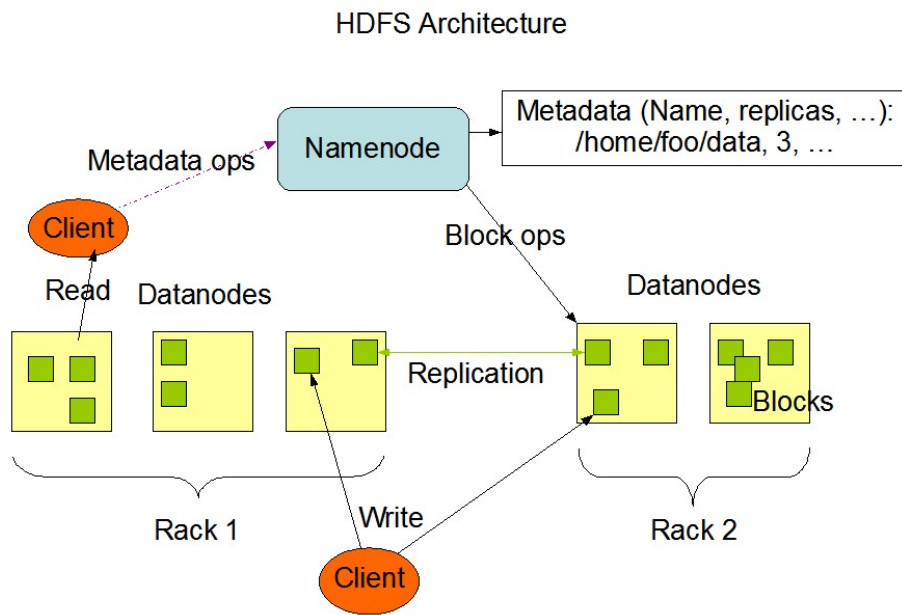| | |
|---|---|
| **Rack awareness** | Considers a node's physical location when allocating storage and scheduling tasks |
| **Minimal data motion** | Hadoop moves compute processes to the data on HDFS and not the other way around. Processing tasks can occur on the physical node where the data resides, which significantly reduces network I/O and provides very high aggregate bandwidth. |
| **Utilities** | Dynamically diagnose the health of the file system and rebalance the data on different nodes |
| **Rollback** | Allows operators to bring back the previous version of HDFS after an upgrade, in case of human or systemic errors |
| **Standby NameNode** | Provides redundancy and supports high availability (HA) |

# HDFS - Architecture



HDFS Architecture

- The file content is split into large blocks (typically 128 megabytes), and each block of the file is independently replicated at multiple DataNodes. The blocks are stored on the local file system on the DataNodes.
- The Namenode actively monitors the number of replicas of a block. When a replica of a block is lost due to a DataNode failure or disk failure, the NameNode creates another replica of the block. The NameNode maintains the namespace tree and the mapping of blocks to DataNodes, holding the entire namespace image in RAM.

# HDFS Azure Support: Azure Blob Storage and ADL

**The hadoop-azure** module provides support for integration with Azure Blob Storage. The built jar file, named hadoop-azure.jar, also declares transitive dependencies on the additional artifacts it requires, notably the Azure Storage SDK for Java.

**Features**

- Read and write data stored in an Azure Blob Storage account.
- Present a hierarchical file system view by implementing the standard Hadoop FileSystem interface.
- Supports configuration of multiple Azure Blob Storage accounts.
- Supports both page blobs (suitable for most use cases, such as MapReduce) and block blobs (suitable for continuous write use cases, such as an HBase write-ahead log).
- Reference file system paths using URLs using the wasb scheme.
- Also reference file system paths using URLs with the wasbs scheme for SSL encrypted access.
- Can act as a source of data in a MapReduce job, or a sink.
- Tested on both Linux and Windows.
- Tested at scale.

# HDFS Commands

User Commands
- classpath
- dfs
- fetchdt
- fsck
- getconf
- groups
- lsSnapshottableDir
- jmxget
- oev
- oiv
- oiv_legacy
- snapshotDiff
- version

Administration Commands
- balancer
- cacheadmin
- crypto
- datanode
- dfsadmin
- haadmin
- journalnode
- mover
- namenode
- nfs3
- portmap
- secondarynamenode
- storagepolicies
- zkfc