

# Apache Hadoop

From Wikipedia, the free encyclopedia

**Apache Hadoop** is a set of algorithms (an open-source software framework) for distributed storage and distributed processing of very large data sets (Big Data) on computer clusters built from commodity hardware. All the modules in Hadoop are designed with a fundamental assumption that hardware failures (of individual machines, or racks of machines) are commonplace and thus should be automatically handled in software by the framework.

The core of Apache Hadoop consists of a storage part (Hadoop Distributed File System (HDFS)) and a processing part (MapReduce). Hadoop splits files into large blocks (default 64MB or 128MB) and distributes the blocks amongst the nodes in the cluster. To process the data, Hadoop Map/Reduce transfers code (specifically Jar files) to nodes that have the required data, which the nodes then process in parallel. This approach takes advantage of data locality<sup>[3]</sup> to allow the data to be processed faster and more efficiently via distributed processing than by using a more conventional supercomputer architecture that relies on a parallel file system where computation and data are connected via high-speed networking.<sup>[4]</sup>

The base Apache Hadoop framework is composed of the following modules:

- *Hadoop Common* – contains libraries and utilities needed by other Hadoop modules;
- *Hadoop Distributed File System (HDFS)* – a distributed file-system that stores data on commodity machines, providing very high aggregate bandwidth across the cluster;
- *Hadoop YARN* – a resource-management platform responsible for managing compute resources in clusters and using them for scheduling of users' applications;<sup>[5][6]</sup> and
- *Hadoop MapReduce* – a programming model for large scale data processing.

Since 2012,<sup>[7]</sup> the term "Hadoop" often refers not just to the base modules above but also to the collection of additional software packages that can be installed on top of or alongside Hadoop, such as Apache Pig, Apache Hive, Apache HBase, Apache Spark, and others.<sup>[8][9]</sup>

Apache Hadoop's MapReduce and HDFS components originally derived respectively from Google's MapReduce and Google File System.

The Hadoop framework itself is mostly written in the Java programming language, with some native code in C and command line utilities written as shell-scripts. For end-users, though MapReduce Java code is common, any programming language can be used with "Hadoop Streaming" to implement the "map" and "reduce" parts of the user's program.<sup>[10]</sup> Other related projects expose other higher level user interfaces.

## Apache Hadoop



<b>Developer(s)</b>	Apache Software Foundation
<b>Initial release</b>	December 10, 2011 <sup>[1]</sup>
<b>Stable release</b>	2.6.0 / November 18, 2014 <sup>[2]</sup>
<b>Development status</b>	Active
<b>Written in</b>	Java
<b>Operating system</b>	Cross-platform
<b>Type</b>	Distributed file system
<b>License</b>	Apache License 2.0
<b>Website</b>	<a href="http://hadoop.apache.org">hadoop.apache.org</a> ( <a href="http://hadoop.apache.org/">http://hadoop.apache.org/</a> )

Prominent corporate users of Hadoop include Facebook and Yahoo. It can be deployed in traditional onsite datacenters as well as via the cloud; e.g., it is available on Microsoft Azure, Amazon Elastic Compute Cloud (EC2) and Amazon Simple Storage Service (S3) cloud services.

Apache Hadoop is a registered trademark of the Apache Software Foundation.

## Contents

- 1 History
- 2 Architecture
  - 2.1 File system
    - 2.1.1 Hadoop distributed file system
    - 2.1.2 Other file systems
  - 2.2 JobTracker and TaskTracker: the MapReduce engine
    - 2.2.1 Scheduling
      - 2.2.1.1 Fair scheduler
      - 2.2.1.2 Capacity scheduler
  - 2.3 Other applications
- 3 Prominent users
  - 3.1 Yahoo!
  - 3.2 Facebook
  - 3.3 Other users
- 4 Hadoop hosting in the Cloud
  - 4.1 Hadoop on Microsoft Azure
  - 4.2 Hadoop on Amazon EC2/S3 services
  - 4.3 Amazon Elastic MapReduce
- 5 Commercial support
  - 5.1 ASF's view on the use of "Hadoop" in product names
- 6 Papers
- 7 See also
- 8 References
- 9 Bibliography
- 10 External links

## History

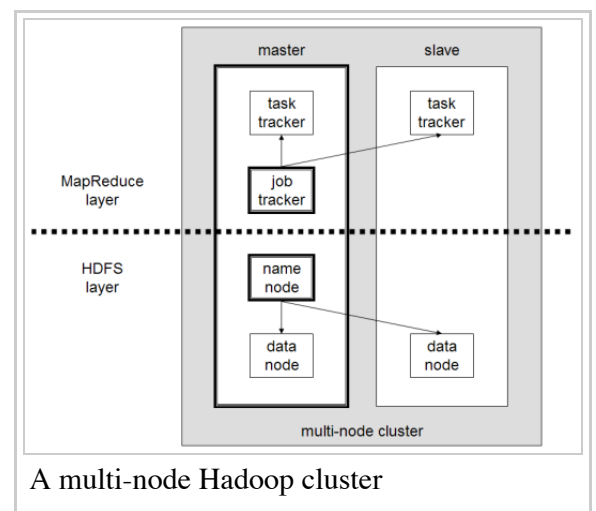
Hadoop was created by Doug Cutting and Mike Cafarella<sup>[11]</sup> in 2005. Cutting, who was working at Yahoo! at the time,<sup>[12]</sup> named it after his son's toy elephant.<sup>[13]</sup> It was originally developed to support distribution for the Nutch search engine project.<sup>[14]</sup>

# Architecture

Hadoop consists of the *Hadoop Common* package, which provides filesystem and OS level abstractions, a MapReduce engine (either MapReduce/MR1 or YARN/MR2)<sup>[15]</sup> and the Hadoop Distributed File System (HDFS). The Hadoop Common package contains the necessary Java ARchive (JAR) files and scripts needed to start Hadoop. The package also provides source code, documentation, and a contribution section that includes projects from the Hadoop Community.

For effective scheduling of work, every Hadoop-compatible file system should provide location awareness: the name of the rack (more precisely, of the network switch) where a worker node is. Hadoop applications can use this information to run work on the node where the data is, and, failing that, on the same rack/switch, reducing backbone traffic. HDFS uses this method when replicating data to try to keep different copies of the data on different racks. The goal is to reduce the impact of a rack power outage or switch failure, so that even if these events occur, the data may still be readable.<sup>[16]</sup>

A small Hadoop cluster includes a single master and multiple worker nodes. The master node consists of a JobTracker, TaskTracker, NameNode and DataNode. A slave or *worker node* acts as both a DataNode and TaskTracker, though it is possible to have data-only worker nodes and compute-only worker nodes. These are normally used only in nonstandard applications.<sup>[17]</sup> Hadoop requires Java Runtime Environment (JRE) 1.6 or higher. The standard startup and shutdown scripts require that Secure Shell (ssh) be set up between nodes in the cluster.<sup>[18]</sup>



In a larger cluster, the HDFS is managed through a dedicated NameNode server to host the file system index, and a secondary NameNode that can generate snapshots of the namenode's memory structures, thus preventing file-system corruption and reducing loss of data. Similarly, a standalone JobTracker server can manage job scheduling. In clusters where the Hadoop MapReduce engine is deployed against an alternate file system, the NameNode, secondary NameNode, and DataNode architecture of HDFS are replaced by the file-system-specific equivalents.

## File system

### Hadoop distributed file system

The **Hadoop distributed file system (HDFS)** is a distributed, scalable, and portable file-system written in Java for the Hadoop framework. A Hadoop cluster has nominally a single namenode plus a cluster of datanodes, although redundancy options are available for the namenode due to its criticality. Each datanode serves up blocks of data over the network using a block protocol specific to HDFS. The file system uses TCP/IP sockets for communication. Clients use remote procedure call (RPC) to communicate between each other.

HDFS stores large files (typically in the range of gigabytes to terabytes<sup>[19]</sup>) across multiple machines. It achieves reliability by replicating the data across multiple hosts, and hence theoretically does not require RAID storage on hosts (but to increase I/O performance some RAID configurations are still useful). With the default replication value, 3, data is stored on three nodes: two on the same rack, and one on a

different rack. Data nodes can talk to each other to rebalance data, to move copies around, and to keep the replication of data high. HDFS is not fully POSIX-compliant, because the requirements for a POSIX file-system differ from the target goals for a Hadoop application. The tradeoff of not having a fully POSIX-compliant file-system is increased performance for data throughput and support for non-POSIX operations such as Append.<sup>[20]</sup>

HDFS added the high-availability capabilities, as announced for release 2.0 in May 2012,<sup>[21]</sup> letting the main metadata server (the NameNode) fail over manually to a backup. The project has also started developing automatic fail-over.

The HDFS file system includes a so-called *secondary namenode*, a misleading name that some might incorrectly interpret as a backup namenode for when the primary namenode goes offline. In fact, the secondary namenode regularly connects with the primary namenode and builds snapshots of the primary namenode's directory information, which the system then saves to local or remote directories. These checkpointed images can be used to restart a failed primary namenode without having to replay the entire journal of file-system actions, then to edit the log to create an up-to-date directory structure. Because the namenode is the single point for storage and management of metadata, it can become a bottleneck for supporting a huge number of files, especially a large number of small files. HDFS Federation, a new addition, aims to tackle this problem to a certain extent by allowing multiple namespaces served by separate namenodes.

An advantage of using HDFS is data awareness between the job tracker and task tracker. The job tracker schedules map or reduce jobs to task trackers with an awareness of the data location. For example: if node A contains data (x,y,z) and node B contains data (a,b,c), the job tracker schedules node B to perform map or reduce tasks on (a,b,c) and node A would be scheduled to perform map or reduce tasks on (x,y,z). This reduces the amount of traffic that goes over the network and prevents unnecessary data transfer. When Hadoop is used with other file systems, this advantage is not always available. This can have a significant impact on job-completion times, which has been demonstrated when running data-intensive jobs.<sup>[22]</sup>

HDFS was designed for mostly immutable files<sup>[20]</sup> and may not be suitable for systems requiring concurrent write-operations.

HDFS can be mounted directly with a Filesystem in Userspace (FUSE) virtual file system on Linux and some other Unix systems.

File access can be achieved through the native Java API, the Thrift API to generate a client in the language of the users' choosing (C++, Java, Python, PHP, Ruby, Erlang, Perl, Haskell, C#, Cocoa, Smalltalk, and OCaml), the command-line interface, browsed through the HDFS-UI webapp over HTTP, or via 3rd-party network client libraries.<sup>[23]</sup>

## Other file systems

Hadoop works directly with any distributed file system that can be mounted by the underlying operating system simply by using a file:// URL; however, this comes at a price: the loss of locality. To reduce network traffic, Hadoop needs to know which servers are closest to the data; this is information that Hadoop-specific file system bridges can provide.

In May 2011, the list of supported file systems bundled with Apache Hadoop were:

- HDFS: Hadoop's own rack-aware file system.<sup>[24]</sup> This is designed to scale to tens of petabytes of storage and runs on top of the file systems of the underlying operating systems.
- FTP File system: this stores all its data on remotely accessible FTP servers.
- Amazon S3 file system. This is targeted at clusters hosted on the Amazon Elastic Compute Cloud server-on-demand infrastructure. There is no rack-awareness in this file system, as it is all remote.
- Windows Azure Storage Blobs (WASB) (<http://azure.microsoft.com/en-us/documentation/articles/hdinsight-use-blob-storage>) file system. WASB, an extension on top of HDFS, allows distributions of Hadoop to access data in Azure blob stores without moving the data permanently into the cluster.

A number of third-party file system bridges have also been written, none of which are currently in Hadoop distributions. However, some commercial distributions of Hadoop ship with an alternative filesystem as the default, -specifically IBM and MapR.

- In 2009 IBM discussed running Hadoop over the IBM General Parallel File System.<sup>[25]</sup> The source code was published in October 2009.<sup>[26]</sup>
- In April 2010, Parascle published the source code to run Hadoop against the Parascle file system.<sup>[27]</sup>
- In April 2010, Appistry released a Hadoop file system driver for use with its own CloudIQ Storage product.<sup>[28]</sup>
- In June 2010, HP discussed a location-aware IBRIX Fusion file system driver.<sup>[29]</sup>
- In May 2011, MapR Technologies, Inc. announced the availability of an alternative file system for Hadoop, which replaced the HDFS file system with a full random-access read/write file system.

## JobTracker and TaskTracker: the MapReduce engine

Above the file systems comes the MapReduce engine, which consists of one *JobTracker*, to which client applications submit MapReduce jobs. The JobTracker pushes work out to available *TaskTracker* nodes in the cluster, striving to keep the work as close to the data as possible. With a rack-aware file system, the JobTracker knows which node contains the data, and which other machines are nearby. If the work cannot be hosted on the actual node where the data resides, priority is given to nodes in the same rack. This reduces network traffic on the main backbone network. If a TaskTracker fails or times out, that part of the job is rescheduled. The TaskTracker on each node spawns off a separate Java Virtual Machine process to prevent the TaskTracker itself from failing if the running job crashes the JVM. A heartbeat is sent from the TaskTracker to the JobTracker every few minutes to check its status. The Job Tracker and TaskTracker status and information is exposed by Jetty and can be viewed from a web browser

Known limitations of this approach are:

- The allocation of work to TaskTrackers is very simple. Every TaskTracker has a number of available *slots* (such as "4 slots"). Every active map or reduce task takes up one slot. The Job Tracker allocates work to the tracker nearest to the data with an available slot. There is no

consideration of the current system load of the allocated machine, and hence its actual availability.

- If one TaskTracker is very slow, it can delay the entire MapReduce job – especially towards the end of a job, where everything can end up waiting for the slowest task. With speculative execution enabled, however, a single task can be executed on multiple slave nodes.

## Scheduling

By default Hadoop uses FIFO, and optionally 5 scheduling priorities to schedule jobs from a work queue.<sup>[30]</sup> In version 0.19 the job scheduler was refactored out of the JobTracker, while adding the ability to use an alternate scheduler (such as the *Fair scheduler* or the *Capacity scheduler*, described next).<sup>[31]</sup>

### Fair scheduler

The fair scheduler was developed by Facebook.<sup>[32]</sup> The goal of the fair scheduler is to provide fast response times for small jobs and QoS for production jobs. The fair scheduler has three basic concepts.<sup>[33]</sup>

1. Jobs are grouped into pools.
2. Each pool is assigned a guaranteed minimum share.
3. Excess capacity is split between jobs.

By default, jobs that are uncategorized go into a default pool. Pools have to specify the minimum number of map slots, reduce slots, and a limit on the number of running jobs.

### Capacity scheduler

The capacity scheduler was developed by Yahoo. The capacity scheduler supports several features that are similar to the fair scheduler.<sup>[34]</sup>

- Jobs are submitted into queues.
- Queues are allocated a fraction of the total resource capacity.
- Free resources are allocated to queues beyond their total capacity.
- Within a queue a job with a high level of priority has access to the queue's resources.

There is no preemption once a job is running.

## Other applications

The HDFS file system is not restricted to MapReduce jobs. It can be used for other applications, many of which are under development at Apache. The list includes the HBase database, the Apache Mahout machine learning system, and the Apache Hive Data Warehouse system. Hadoop can in theory be used for any sort of work that is batch-oriented rather than real-time, is very data-intensive, and benefits from parallel processing of data. It can also be used to complement a real-time system, such as lambda architecture.

As of October 2009, commercial applications of Hadoop<sup>[35]</sup> included:

- Log and/or clickstream analysis of various kinds
- Marketing analytics
- Machine learning and/or sophisticated data mining
- Image processing
- Processing of XML messages
- Web crawling and/or text processing
- General archiving, including of relational/tabular data, e.g. for compliance

## Prominent users

### Yahoo!

On February 19, 2008, Yahoo! Inc. launched what it claimed was the world's largest Hadoop production application. The Yahoo! Search Webmap is a Hadoop application that runs on a Linux cluster with more than 10,000 cores and produced data that was used in every Yahoo! web search query.<sup>[36]</sup>

There are multiple Hadoop clusters at Yahoo! and no HDFS file systems or MapReduce jobs are split across multiple datacenters. Every Hadoop cluster node bootstraps the Linux image, including the Hadoop distribution. Work that the clusters perform is known to include the index calculations for the Yahoo! search engine.

On June 10, 2009, Yahoo! made the source code of the version of Hadoop it runs in production available to the public.<sup>[37]</sup> Yahoo! contributes all the work it does on Hadoop to the open-source community. The company's developers also fix bugs, provide stability improvements internally, and release this patched source code so that other users may benefit from their effort.

### Facebook

In 2010 Facebook claimed that they had the largest Hadoop cluster in the world with 21 PB of storage.<sup>[38]</sup> On June 13, 2012 they announced the data had grown to 100 PB.<sup>[39]</sup> On November 8, 2012 they announced the data gathered in the warehouse grows by roughly half a PB per day.<sup>[40]</sup>

### Other users

As of 2013, Hadoop adoption is widespread. For example, more than half of the Fortune 50 use Hadoop.<sup>[41]</sup>

## Hadoop hosting in the Cloud

Hadoop can be deployed in a traditional onsite datacenter as well as in the cloud.<sup>[42]</sup> The cloud allows organizations to deploy Hadoop without hardware to acquire or specific setup expertise.<sup>[43]</sup> Vendors who currently have an offer for the cloud include Microsoft, Amazon, and Google.

### Hadoop on Microsoft Azure

Azure HDInsight<sup>[44]</sup> is a service that deploys Hadoop on Microsoft Azure. HDInsight uses a Windows-based Hadoop distribution that was jointly developed with Hortonworks and allows programming extensions with .NET (in addition to Java).<sup>[44]</sup> By deploying HDInsight in the cloud, organizations can spin up the number of nodes they want and only get charged for the compute and storage that is used.<sup>[44]</sup> Hortonworks implementations can also move data from the on-premises datacenter to the cloud for backup, development/test, and bursting scenarios.<sup>[44]</sup>

## Hadoop on Amazon EC2/S3 services

It is possible to run Hadoop on Amazon Elastic Compute Cloud (EC2) and Amazon Simple Storage Service (S3).<sup>[45]</sup> As an example The New York Times used 100 Amazon EC2 instances and a Hadoop application to process 4 TB of raw image TIFF data (stored in S3) into 11 million finished PDFs in the space of 24 hours at a computation cost of about \$240 (not including bandwidth).<sup>[46]</sup>

There is support for the S3 file system in Hadoop distributions, and the Hadoop team generates EC2 machine images after every release. From a pure performance perspective, Hadoop on S3/EC2 is inefficient, as the S3 file system is remote and delays returning from every write operation until the data is guaranteed not lost. This removes the locality advantages of Hadoop, which schedules work near data to save on network load.

## Amazon Elastic MapReduce

Elastic MapReduce (EMR)<sup>[47]</sup> was introduced by Amazon in April 2009. Provisioning of the Hadoop cluster, running and terminating jobs, and handling data transfer between EC2(VM) and S3(Object Storage) are automated by Elastic MapReduce. Apache Hive, which is built on top of Hadoop for providing data warehouse services, is also offered in Elastic MapReduce.<sup>[48]</sup>

Support for using Spot Instances<sup>[49]</sup> was later added in August 2011.<sup>[50]</sup> Elastic MapReduce is fault tolerant for slave failures,<sup>[51]</sup> and it is recommended to only run the Task Instance Group on spot instances to take advantage of the lower cost while maintaining availability.<sup>[52]</sup>

## Commercial support

A number of companies offer commercial implementations or support for Hadoop.<sup>[53]</sup>

## ASF's view on the use of "Hadoop" in product names

The Apache Software Foundation has stated that only software officially released by the Apache Hadoop Project can be called *Apache Hadoop* or *Distributions of Apache Hadoop*.<sup>[54]</sup> The naming of products and derivative works from other vendors and the term "compatible" are somewhat controversial within the Hadoop developer community.<sup>[55]</sup>

## Papers

Some papers influenced the birth and growth of Hadoop and big data processing. Here is a partial list:



- 2004 MapReduce: Simplified Data Processing on Large Clusters  
([https://www.usenix.org/legacy/publications/library/proceedings/osdi04/tech/full\\_papers/dean/dean\\_html/index.html](https://www.usenix.org/legacy/publications/library/proceedings/osdi04/tech/full_papers/dean/dean_html/index.html)) by Jeffrey Dean and Sanjay Ghemawat from Google. This paper inspired Doug Cutting to develop an open-source implementation of the Map-Reduce framework. He named it Hadoop (<http://hadoop.apache.org/>), after his son's toy elephant.
- 2005 From Databases to Dataspaces: A New Abstraction for Information Management  
(<http://www.eecs.berkeley.edu/~franklin/Papers/dataspaceSR.pdf>), the authors highlight the need for storage systems to accept all data formats and to provide APIs for data access that evolve based on the storage system's understanding of the data.
- 2006 Bigtable: A Distributed Storage System for Structured Data  
([http://static.googleusercontent.com/external\\_content/untrusted\\_dlcp/research.google.com/en/us/archive/bigtable-osdi06.pdf](http://static.googleusercontent.com/external_content/untrusted_dlcp/research.google.com/en/us/archive/bigtable-osdi06.pdf)) from Google.
- 2008 H-store: a high-performance, distributed main memory transaction processing system  
(<http://www.vldb.org/pvldb/1/1454211.pdf>)
- 2009 MAD Skills: New Analysis Practices for Big Data  
(<http://db.cs.berkeley.edu/jmh/papers/madskills-032009.pdf>)
- 2011 Apache Hadoop Goes Realtime at Facebook  
(<http://borthakur.com/ftp/RealtimeHadoopSigmod2011.pdf>)

## See also

- Apache Accumulo – Secure BigTable<sup>[56]</sup>
- Apache Bigtop - Packaging and interoperability testing of Hadoop-related projects
- Apache Cassandra – A column-oriented database that supports access from Hadoop
- Apache CouchDB is a database that uses JSON for documents, JavaScript for MapReduce queries, and regular HTTP for an API
- Apache Mahout – Machine Learning algorithms implemented on Hadoop
- Big data
- Cask (company)
- Cloud computing
- Data Intensive Computing
- Datameer Analytics Solution (DAS) – data source integration, storage, analytics engine and visualization
- Druid (open-source data store) - Provides a native indexing service for ingesting from HDFS.
- HBase – BigTable-model database
- HPCC – LexisNexis Risk Solutions High Performance Computing Cluster
- Hortonworks - Open source, Hortonworks Data Platform (HDP) provides Hadoop designed for enterprise data processing
- Hypertable – HBase alternative

- MapReduce – Hadoop's fundamental data filtering algorithm
- Nutch – An effort to build an open source search engine based on Lucene and Hadoop, also created by Doug Cutting
- Pentaho – Open source data integration (Kettle), analytics, reporting, visualization and predictive analytics directly from Hadoop nodes
- Pivotal HD - Apache Hadoop distribution enhanced to support enterprise Big Data analytics. Industry's first native massively parallel processing (MPP) SQL database on Hadoop.
- Qubole - a cloud-based Big Data as a service developer
- RapidMiner Radoop – In-Hadoop big data analytics providing a set of algorithms for doing scalable data transformations, advanced analytics, and predictive modeling
- Sector/Sphere – Open source distributed storage and processing
- Simple Linux Utility for Resource Management
- Talend – An open source integration software

## References

1. ^ "Hadoop Releases" (<http://hadoop.apache.org/releases.html#27+December%2C+2011%3A+release+1.0.0+available>). *apache.org*. Apache Software Foundation. Retrieved 2014-12-06.
2. ^ "Hadoop Releases" (<http://hadoop.apache.org/docs/r2.6.0/hadoop-project-dist/hadoop-common/releasenotes.html>). Hadoop.apache.org. Retrieved 2014-12-01.
3. ^ "What is the Hadoop Distributed File System (HDFS)?" (<http://www-01.ibm.com/software/data/infosphere/hadoop/hdfs/>). *ibm.com*. IBM. Retrieved 2014-10-30.
4. ^ Malak, Michael (2014-09-19). "Data Locality: HPC vs. Hadoop vs. Spark" (<http://www.datascienceassn.org/content/data-locality-hpc-vs-hadoop-vs-spark>). *datascienceassn.org*. Data Science Association. Retrieved 2014-10-30.
5. ^ "Resource (Apache Hadoop Main 2.5.1 API)" ([http://hadoop.apache.org/docs/r2.5.1/api/org/apache/hadoop/yarn/api/records/Resource.html#newInstance\(int\)](http://hadoop.apache.org/docs/r2.5.1/api/org/apache/hadoop/yarn/api/records/Resource.html#newInstance(int))). *apache.org*. Apache Software Foundation. 2014-09-12. Retrieved 2014-09-30.
6. ^ Murthy, Arun (2012-08-15). "Apache Hadoop YARN – Concepts and Applications" (<http://hortonworks.com/blog/apache-hadoop-yarn-concepts-and-applications/>). *hortonworks.com*. Hortonworks. Retrieved 2014-09-30.
7. ^ "Continuuity Raises \$10 Million Series A Round to Ignite Big Data Application Development Within the Hadoop Ecosystem" (<http://finance.yahoo.com/news/continuuity-raises-10-million-series-120500471.html>). *finance.yahoo.com*. Marketwired. 2012-11-14. Retrieved 2014-10-30.
8. ^ "Hadoop-related projects at" (<http://hadoop.apache.org/>). Hadoop.apache.org. Retrieved 2013-10-17.
9. ^ Roman, Javi. "The Hadoop Ecosystem Table" (<http://hadoopecosystemtable.github.io/>). *github.com*. Retrieved 2014-12-06.
10. ^ "[nlpatumd] Adventures with Hadoop and Perl" (<http://www.mail-archive.com/nlpatumd@yahoogroups.com/msg00570.html>). Mail-archive.com. 2010-05-02. Retrieved 2013-04-05.

11. ^ "Michael J. Cafarella" (<http://web.eecs.umich.edu/~michjc/bio.html>). Web.eecs.umich.edu. Retrieved 2013-04-05.
12. ^ Hadoop creator goes to Cloudera (<http://www.sdtimes.com/blog/post/2009/08/10/Hadoop-creator-goes-to-Cloudera.aspx>)
13. ^ Ashlee Vance (2009-03-17). "Hadoop, a Free Software Program, Finds Uses Beyond Search" (<http://www.nytimes.com/2009/03/17/technology/business-computing/17cloud.html>). *The New York Times*. Archived (<http://web.archive.org/web/20100211022503/http://www.nytimes.com/2009/03/17/technology/business-computing/17cloud.html?>) from the original on 11 February 2010. Retrieved 2010-01-20.
14. ^ "Hadoop contains the distributed computing platform that was formerly a part of Nutch. This includes the Hadoop Distributed Filesystem (HDFS) and an implementation of MapReduce." About Hadoop (<http://hadoop.apache.org/core/>)
15. ^ Harsh Chouraria (21 October 2012). "MR2 and YARN Briefly Explained" (<http://blog.cloudera.com/blog/2012/10/mr2-and-yarn-briefly-explained/>). *cloudera.com*. Cloudera. Retrieved 23 October 2013.
16. ^ "HDFS User Guide" (<http://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-hdfs/HdfsUserGuide.html>). Hadoop.apache.org. Retrieved 2014-09-04.
17. ^ "Running Hadoop on Ubuntu Linux (Multi-Node Cluster)" (<http://www.michael-noll.com/tutorials/running-hadoop-on-ubuntu-linux-multi-node-cluster/>).
18. ^ "Running Hadoop on Ubuntu Linux (Single-Node Cluster)" (<http://www.michael-noll.com/tutorials/running-hadoop-on-ubuntu-linux-single-node-cluster/#prerequisites>). Retrieved 6 June 2013.
19. ^ "HDFS Architecture" ([http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsDesign.html#Large\\_Data\\_Sets](http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsDesign.html#Large_Data_Sets)). Retrieved 1 September 2013.
20. ^ <sup>a</sup> <sup>b</sup> Yaniv Pessach (2013). "Distributed Storage" ([http://openlibrary.org/books/OL25423189M/Distributed\\_Storage\\_Concepts\\_Algorithms\\_and\\_Implementations](http://openlibrary.org/books/OL25423189M/Distributed_Storage_Concepts_Algorithms_and_Implementations)) (Distributed Storage: Concepts, Algorithms, and Implementations ed.). Amazon.com
21. ^ "Version 2.0 provides for manual failover and they are working on automatic failover:" (<https://hadoop.apache.org/releases.html#23+May%2C+2012%3A+Release+2.0.0-alpha+available>). Hadoop.apache.org. Retrieved 30 July 2013.
22. ^ "Improving MapReduce performance through data placement in heterogeneous Hadoop Clusters" (<http://www.eng.auburn.edu/~xqin/pubs/hcw10.pdf>) (PDF). Eng.auburn.ed. April 2010.
23. ^ "Mounting HDFS" (<https://wiki.apache.org/hadoop/MountableHDFS>). Retrieved May 2014.
24. ^ "HDFS Users Guide – Rack Awareness" ([http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsUserGuide.html#Rack\\_Awareness](http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsUserGuide.html#Rack_Awareness)). Hadoop.apache.org. Retrieved 2013-10-17.
25. ^ "Cloud analytics: Do we really need to reinvent the storage stack?" ([http://www.usenix.org/events/hotcloud09/tech/full\\_papers/ananthanarayanan.pdf](http://www.usenix.org/events/hotcloud09/tech/full_papers/ananthanarayanan.pdf)). IBM. June 2009.
26. ^ "HADOOP-6330: Integrating IBM General Parallel File System implementation of Hadoop Filesystem interface" (<https://issues.apache.org/jira/browse/HADOOP-6330>). IBM. 2009-10-23.
27. ^ "HADOOP-6704: add support for Parascle filesystem" (<https://issues.apache.org/jira/browse/HADOOP-6704>). Parascle. 2010-04-14.
28. ^ "HDFS with CloudIQ Storage" (<http://resources.appistry.com/news-and-events/press/06072010-appistry-cloudiq-storage-now-generally-available>). Appistry,Inc. 2010-07-06.
29. ^ "High Availability Hadoop" ([http://www.slideshare.net/steve\\_l/high-availability-hadoop](http://www.slideshare.net/steve_l/high-availability-hadoop)). HP. 2010-06-09.

30. ^ job ([http://hadoop.apache.org/common/docs/current/commands\\_manual.html#job](http://hadoop.apache.org/common/docs/current/commands_manual.html#job))
31. ^ "Refactor the scheduler out of the JobTracker" (<https://issues.apache.org/jira/browse/HADOOP-3412>). *Hadoop Common*. Apache Software Foundation. Retrieved 9 June 2012.
32. ^ M. Tim Jones (6 December 2011). "Scheduling in Hadoop" (<http://www.ibm.com/developerworks/library/os-hadoop-scheduling/>). *ibm.com*. IBM. Retrieved 20 November 2013.
33. ^ [1] ([https://svn.apache.org/repos/asf/hadoop/common/branches/MAPREDUCE-233/src/contrib/fairscheduler/designdoc/fair\\_scheduler\\_design\\_doc.pdf](https://svn.apache.org/repos/asf/hadoop/common/branches/MAPREDUCE-233/src/contrib/fairscheduler/designdoc/fair_scheduler_design_doc.pdf)) Hadoop Fair Scheduler Design Document
34. ^ [2] ([http://hadoop.apache.org/docs/stable1/capacity\\_scheduler.html](http://hadoop.apache.org/docs/stable1/capacity_scheduler.html)) Capacity Scheduler Guide
35. ^ October 10, 2009 (2009-10-10). " "How 30+ enterprises are using Hadoop", in DBMS2" (<http://www.dbms2.com/2009/10/10/enterprises-using-hadoop/>). Dbms2.com. Retrieved 2013-10-17.
36. ^ Yahoo! Launches World's Largest Hadoop Production Application (<https://developer.yahoo.com/blogs/hadoop/yahoo-launches-world-largest-hadoop-production-application-398.html>)
37. ^ "Hadoop and Distributed Computing at Yahoo!" (<http://developer.yahoo.com/hadoop/>). Yahoo!. 2011-04-20. Retrieved 2013-10-17.
38. ^ "HDFS: Facebook has the world's largest Hadoop cluster!" (<http://hadoopblog.blogspot.com/2010/05/facebook-has-worlds-largest-hadoop.html>). Hadoopblog.blogspot.com. 2010-05-09. Retrieved 2012-05-23.
39. ^ "Under the Hood: Hadoop Distributed File system reliability with Namenode and Avatarnode" (<http://www.facebook.com/notes/facebook-engineering/under-the-hood-hadoop-distributed-filesystem-reliability-with-namenode-and-avata/10150888759153920>). Facebook. Retrieved 2012-09-13.
40. ^ "Under the Hood: Scheduling MapReduce jobs more efficiently with Corona" (<https://www.facebook.com/notes/facebook-engineering/under-the-hood-scheduling-mapreduce-jobs-more-efficiently-with-corona/10151142560538920>). Facebook. Retrieved 2012-11-09.
41. ^ "Altior's AltraSTAR – Hadoop Storage Accelerator and Optimizer Now Certified on CDH4 (Cloudera's Distribution Including Apache Hadoop Version 4)" (<http://www.prnewswire.com/news-releases/altiors-altrastar---hadoop-storage-accelerator-and-optimizer-now-certified-on-cdh4-clouderas-distribution-including-apache-hadoop-version-4-183906141.html>) (Press release). Eatontown, New Jersey: Altior Inc. 2012-12-18. Retrieved 2013-10-30.
42. ^ "What is Hadoop?" (<http://azure.microsoft.com/en-us/solutions/hadoop/>).
43. ^ "Hadoop" (<http://azure.microsoft.com/en-us/solutions/hadoop/>). Azure.microsoft.com. Retrieved 2014-07-22.
44. ^ <sup>a</sup> <sup>b</sup> <sup>c</sup> <sup>d</sup> "HDInsight | Cloud Hadoop" (<http://azure.microsoft.com/en-us/services/hdinsight/>). Azure.microsoft.com. Retrieved 2014-07-22.
45. ^ Varia, Jinesh (@jinman). "Taking Massive Distributed Computing to the Common Man – Hadoop on Amazon EC2/S3" (<http://aws.typepad.com/aws/2008/02/taking-massive.html>). *Amazon Web Services Blog*. Amazon.com. Retrieved 9 June 2012.
46. ^ Gottfrid, Derek (November 1, 2007). "Self-service, Prorated Super Computing Fun!" (<http://open.blogs.nytimes.com/2007/11/01/self-service-prorated-super-computing-fun/?scp=1&sq=self%20service%20prorated&st=cse>). *The New York Times*. Retrieved May 4, 2010.
47. ^ "AWS | Amazon Elastic MapReduce (EMR) | Hadoop MapReduce in the Cloud" (<http://aws.amazon.com/elasticmapreduce/>). Aws.amazon.com. Retrieved 2014-07-22.

48. ^ "Amazon Elastic MapReduce Developer Guide" (<http://s3.amazonaws.com/awsdocs/ElasticMapReduce/latest/emr-dg.pdf>) (PDF). Retrieved 2013-10-17.
49. ^ "Amazon EC2 Spot Instances" (<http://aws.amazon.com/ec2/spot-instances/>). Aws.amazon.com. Retrieved 2014-07-22.
50. ^ "Amazon Elastic MapReduce Now Supports Spot Instances" (<http://aws.amazon.com/about-aws/whats-new/2011/08/18/amazon-elastic-mapreduce-now-supports-spot-instances/>). Amazon.com. 2011-08-18. Retrieved 2013-10-17.
51. ^ "Amazon Elastic MapReduce FAQs" (<http://aws.amazon.com/elasticmapreduce/faqs/#cluster-10>). Amazon.com. Retrieved 2013-10-17.
52. ^ Using Spot Instances with EMR (<https://www.youtube.com/watch?v=66rfnFA0jpM>) on YouTube
53. ^ "Why the Pace of Hadoop Innovation Has to Pick Up" (<http://gigaom.com/cloud/why-we-need-more-hadoop-innovation/>). Gigaom.com. 2011-04-25. Retrieved 2013-10-17.
54. ^ "Defining Hadoop" (<http://wiki.apache.org/hadoop/Defining%20Hadoop>). Wiki.apache.org. 2013-03-30. Retrieved 2013-10-17.
55. ^ "Defining Hadoop Compatibility: revisited" ([http://mail-archives.apache.org/mod\\_mbox/hadoop-general/201105.mbox/%3C4DC91392.2010308@apache.org%3E](http://mail-archives.apache.org/mod_mbox/hadoop-general/201105.mbox/%3C4DC91392.2010308@apache.org%3E)). Mail-archives.apache.org. 2011-05-10. Retrieved 2013-10-17.
56. ^ "Apache Accumulo User Manual: Security" ([https://accumulo.apache.org/1.4/user\\_manual/Security.html](https://accumulo.apache.org/1.4/user_manual/Security.html)). *apache.org*. Apache Software Foundation. Retrieved 2014-12-03.

## Bibliography

- Lam, Chuck (July 28, 2010). *Hadoop in Action* (1st ed.). Manning Publications. p. 325. ISBN 1-935182-19-6.
- Venner, Jason (June 22, 2009). *Pro Hadoop* (<http://www.apress.com/book/view/1430219424>) (1st ed.). Apress. p. 440. ISBN 1-4302-1942-4.
- White, Tom (June 16, 2009). *Hadoop: The Definitive Guide* (<http://oreilly.com/catalog/9780596521974>) (1st ed.). O'Reilly Media. p. 524. ISBN 0-596-52197-9.

## External links

- Official Hadoop Homepage (<http://hadoop.apache.org>)
- Official Hadoop Wiki (<http://wiki.apache.org/hadoop/>)
- Introducing Apache Hadoop: The Modern Data Operating System (<http://www.stanford.edu/class/ee380/Abstracts/111116.html>) — lecture given at Stanford University by Co-Founder and CTO of Cloudera, Amr Awadallah (video archive (<http://ee380.stanford.edu/cgi-bin/videologger.php?target=111116-ee380-300.asx>)) (YouTube (<http://www.youtube.com/watch?v=d2xeNpfzsYI>)))
- Hadoop with Philip Zeyliger, Software Engineering Radio, IEEE Computer Society, March 8 2010 (<http://www.se-radio.net/2010/03/episode-157-hadoop-with-philip-zeyliger/>)

Retrieved from "http://en.wikipedia.org/w/index.php?title=Apache\_Hadoop&oldid=644013236"

Categories: [Hadoop](#) | [Apache Software Foundation](#) | [Software using the Apache license](#)  
| [Free software programmed in Java \(programming language\)](#) | [Free system software](#)  
| [Distributed file systems](#) | [Cloud infrastructure](#) | [Free software for cloud computing](#)

---

- This page was last modified on 24 January 2015, at 21:40.
- Text is available under the Creative Commons Attribution-ShareAlike License; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy. Wikipedia® is a registered trademark of the Wikimedia Foundation, Inc., a non-profit organization.