

# HBASE - OVERVIEW

[http://www.tutorialspoint.com/hbase/hbase\\_overview.htm](http://www.tutorialspoint.com/hbase/hbase_overview.htm)

Copyright © tutorialspoint.com

Since 1970, RDBMS is the solution for data storage and maintenance related problems. After the advent of big data, companies realized the benefit of processing big data and started opting for solutions like Hadoop.

Hadoop uses distributed file system for storing big data, and MapReduce to process it. Hadoop excels in storing and processing of huge data of various formats such as arbitrary, semi-, or even unstructured.

## Limitations of Hadoop

Hadoop can perform only batch processing, and data will be accessed only in a sequential manner. That means one has to search the entire dataset even for the simplest of jobs.

A huge dataset when processed results in another huge data set, which should also be processed sequentially. At this point, a new solution is needed to access any point of data in a single unit of time *randomaccess*.

## Hadoop Random Access Databases

Applications such as HBase, Cassandra, couchDB, Dynamo, and MongoDB are some of the databases that store huge amounts of data and access the data in a random manner.

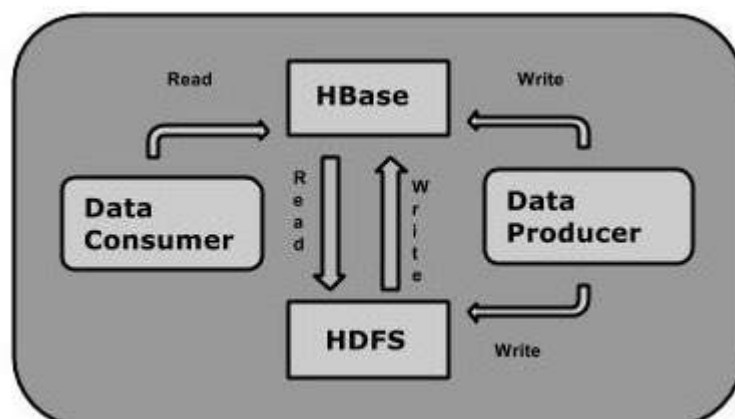
## What is HBase?

HBase is a distributed column-oriented database built on top of the Hadoop file system. It is an open-source project and is horizontally scalable.

HBase is a data model that is similar to Google's big table designed to provide quick random access to huge amounts of structured data. It leverages the fault tolerance provided by the Hadoop File System *HDFS*.

It is a part of the Hadoop ecosystem that provides random real-time read/write access to data in the Hadoop File System.

One can store the data in HDFS either directly or through HBase. Data consumer reads/accesses the data in HDFS randomly using HBase. HBase sits on top of the Hadoop File System and provides read and write access.



## HBase and HDFS

### HDFS

HDFS is a distributed file system suitable for storing large files.

### HBase

HBase is a database built on top of the HDFS.

HDFS does not support fast individual record lookups.

It provides high latency batch processing; no concept of batch processing.

It provides only sequential access of data.

HBase provides fast lookups for larger tables.

It provides low latency access to single rows from billions of records *Randomaccess*.

HBase internally uses Hash tables and provides random access, and it stores the data in indexed HDFS files for faster lookups.

## Storage Mechanism in HBase

HBase is a **column-oriented database** and the tables in it are sorted by row. The table schema defines only column families, which are the key value pairs. A table have multiple column families and each column family can have any number of columns. Subsequent column values are stored contiguously on the disk. Each cell value of the table has a timestamp. In short, in an HBase:

- Table is a collection of rows.
- Row is a collection of column families.
- Column family is a collection of columns.
- Column is a collection of key value pairs.

Given below is an example schema of table in HBase.

Rowid	Column Family	Column Family	Column Family	Column Family								
	col1	col2	col3	col1	col2	col3	col1	col2	col3	col1	col2	col3

1

2

3

## Column Oriented and Row Oriented

Column-oriented databases are those that store data tables as sections of columns of data, rather than as rows of data. Shortly, they will have column families.

### Row-Oriented Database

It is suitable for Online Transaction Process *OLTP*.

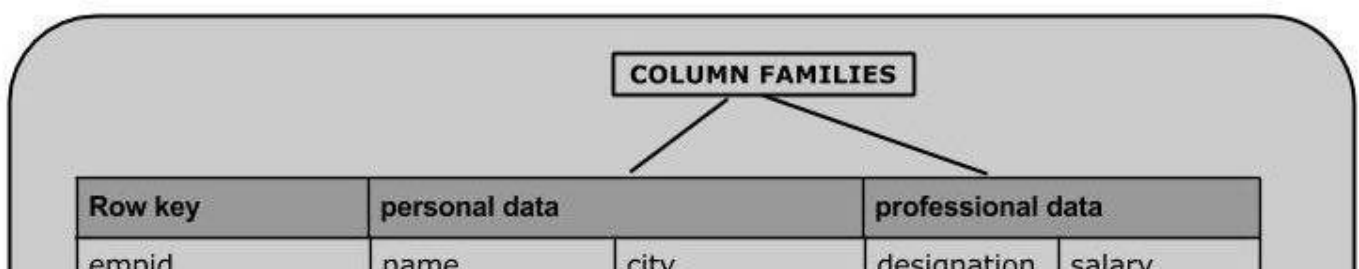
Such databases are designed for small number of rows and columns.

### Column-Oriented Database

It is suitable for Online Analytical Processing *OLAP*.

Column-oriented databases are designed for huge tables.

The following image shows column families in a column-oriented database:



1	raju	hyderabad	manager	50,000
2	ravi	chennai	sr.engineer	30,000
3	rajesh	delhi	jr.engineer	25,000

## HBase and RDBMS

### HBase

HBase is schema-less, it doesn't have the concept of fixed columns schema; defines only column families.

It is built for wide tables. HBase is horizontally scalable.

No transactions are there in HBase.

It has de-normalized data.

It is good for semi-structured as well as structured data.

### RDBMS

An RDBMS is governed by its schema, which describes the whole structure of tables.

It is thin and built for small tables. Hard to scale.

RDBMS is transactional.

It will have normalized data.

It is good for structured data.

## Features of HBase

- HBase is linearly scalable.
- It has automatic failure support.
- It provides consistent read and writes.
- It integrates with Hadoop, both as a source and a destination.
- It has easy java API for client.
- It provides data replication across clusters.

## Where to Use HBase

- Apache HBase is used to have random, real-time read/write access to Big Data.
- It hosts very large tables on top of clusters of commodity hardware.
- Apache HBase is a non-relational database modeled after Google's Bigtable. Bigtable acts up on Google File System, likewise Apache HBase works on top of Hadoop and HDFS.

## Applications of HBase

- It is used whenever there is a need to write heavy applications.
- HBase is used whenever we need to provide fast random access to available data.
- Companies such as Facebook, Twitter, Yahoo, and Adobe use HBase internally.

## HBase History

### Year

### Event

Nov 2006 Google released the paper on BigTable.

Feb 2007 Initial HBase prototype was created as a Hadoop contribution.  
Oct 2007 The first usable HBase along with Hadoop 0.15.0 was released.  
Jan 2008 HBase became the sub project of Hadoop.  
Oct 2008 HBase 0.18.1 was released.  
Jan 2009 HBase 0.19.0 was released.  
Sept 2009 HBase 0.20.0 was released.  
May 2010 HBase became Apache top-level project.

Loading [MathJax]/jax/output/HTML-CSS/jax.js