

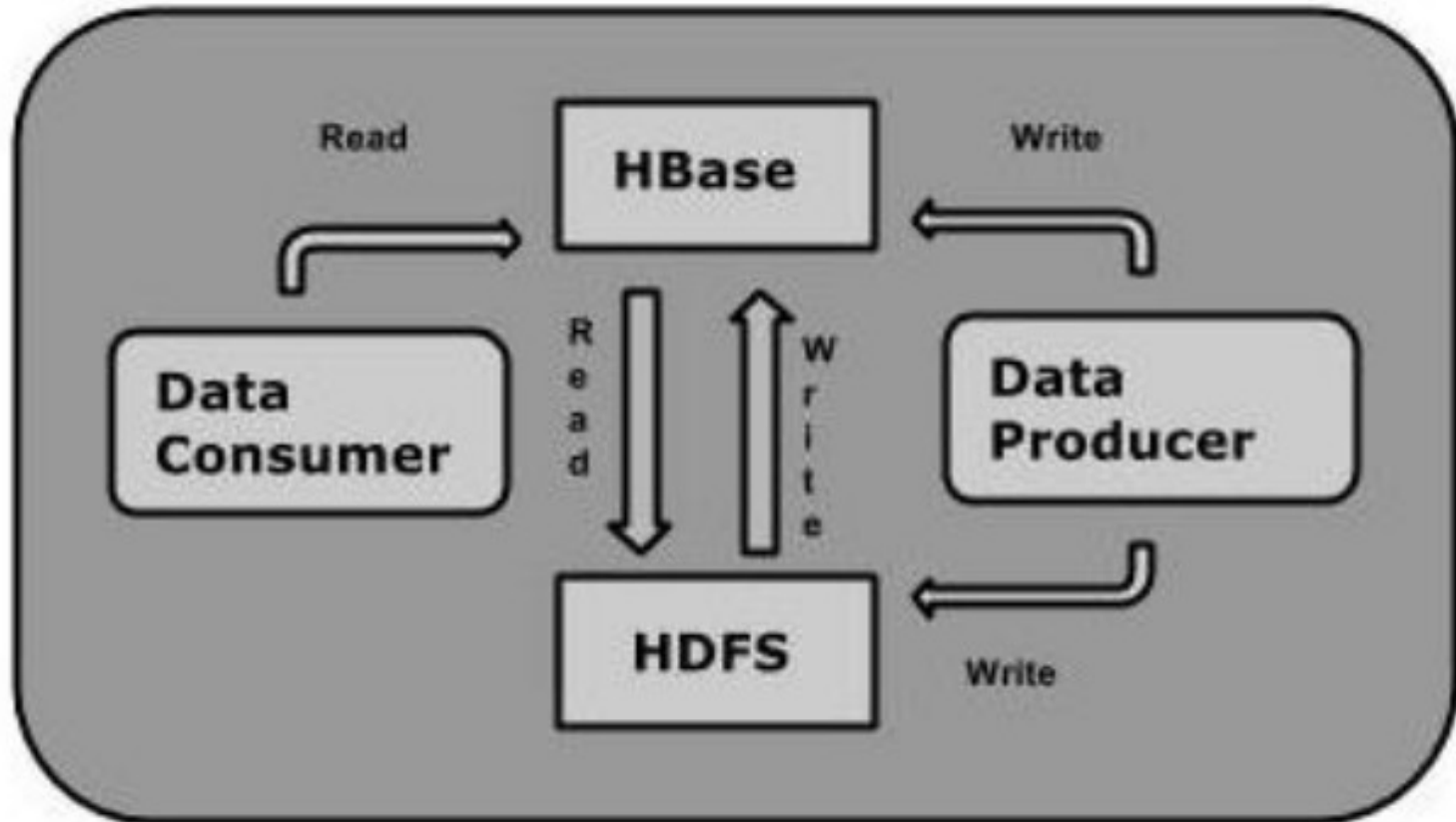
Apache HBase

Tushar B. Kute,
<http://tusharkute.com>

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



HDFS vs. HBase

HDFS	HBase
HDFS is a distributed file system suitable for storing large files.	HBase is a database built on top of the HDFS.
HDFS does not support fast individual record lookups.	HBase provides fast lookups for larger tables.
It provides high latency batch processing; no concept of batch processing.	It provides low latency access to single rows from billions of records (Random access).
It provides only sequential access of data.	HBase internally uses Hash tables and provides random access, and it stores the data in indexed HDFS files for faster lookups.

Storage Mechanism

- 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.

Storage Mechanism

Rowid	Column Family			Column Family			Column Family			Column Family		
	clo1	col 2	col 3	col 1	col 2	col 3	col1	col2	col3	col1	col2	col3
1												
2												
3												

Column Oriented vs. Row Oriented

Row-Oriented Database	Column-Oriented Database
It is suitable for Online Transaction Process (OLTP).	It is suitable for Online Analytical Processing (OLAP).
Such databases are designed for small number of rows and columns.	Column-oriented databases are designed for huge tables.

Example:

COLUMN FAMILIES

Row key	personal data		professional data	
empid	name	city	designation	salary
1	raju	hyderabad	manager	50,000
2	ravi	chennai	sr.engineer	30,000
3	rajesh	delhi	jr.engineer	25,000

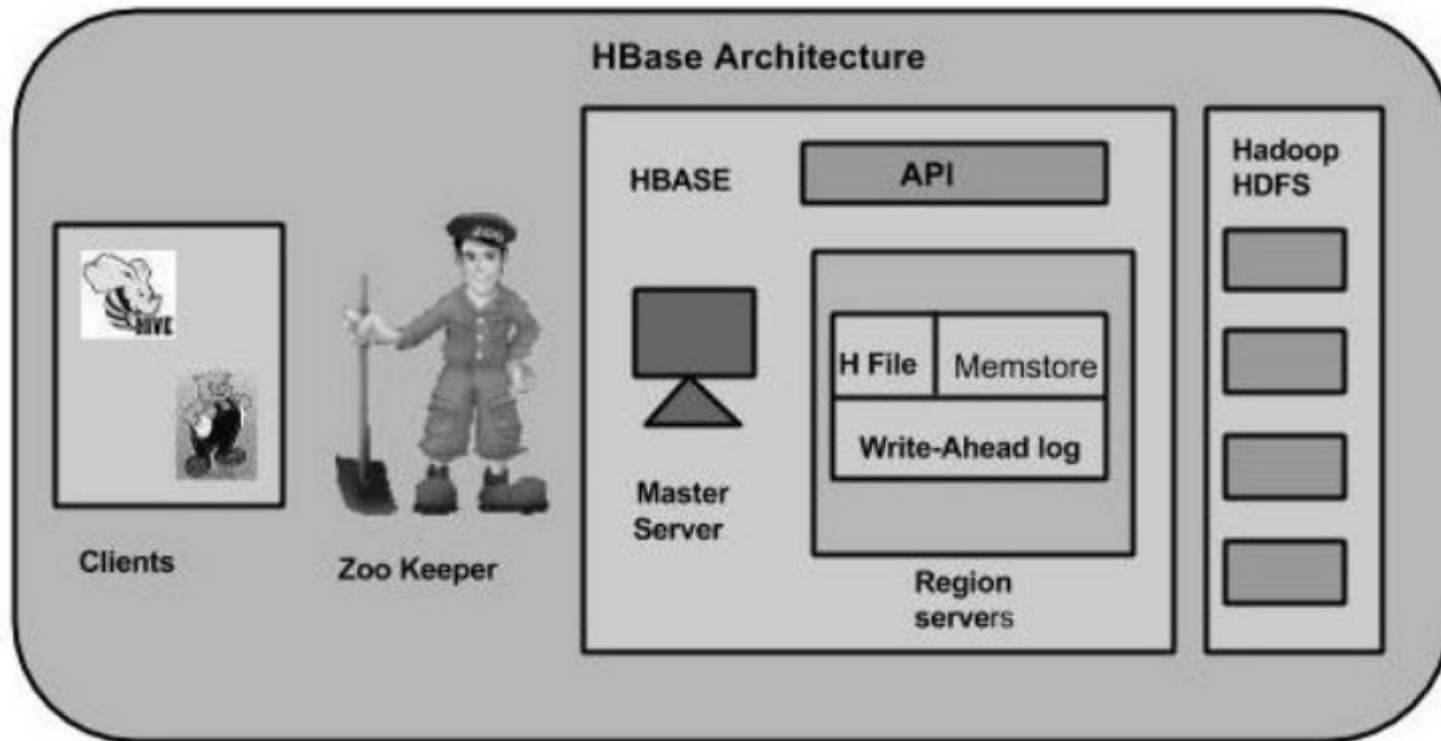
Features

- 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.

Uses

- 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.

Architecture



Master server

- Assigns regions to the region servers and takes the help of Apache ZooKeeper for this task.
- Handles load balancing of the regions across region servers. It unloads the busy servers and shifts the regions to less occupied servers.
- Maintains the state of the cluster by negotiating the load balancing.
- Is responsible for schema changes and other metadata operations such as creation of tables and column families.

Region server

- Regions
- Regions are nothing but tables that are split up and spread across the region servers.
- Region server
 - The region servers have regions that -
 - Communicate with the client and handle data-related operations. Handle read and write requests for all the regions under it.
 - Decide the size of the region by following the region size thresholds.

Zookeeper

- Zookeeper is an open-source project that provides services like maintaining configuration information, naming, providing distributed synchronization, etc.
- Zookeeper has ephemeral nodes representing different region servers. Master servers use these nodes to discover available servers.
- In addition to availability, the nodes are also used to track server failures or network partitions.
- Clients communicate with region servers via zookeeper.
- In pseudo and standalone modes, HBase itself will take care of zookeeper.

Hbase Shell

- HBase contains a shell using which you can communicate with HBase.
- Hbase uses the Hadoop File System to store its data.
- It will have a master server and region servers. The data storage will be in the form of regions (tables).
- These regions will be split up and stored in region servers.

General Commands

- **status:** Provides the status of HBase, for example, the number of servers.
- **version:** Provides the version of HBase being used.
- **table_help:** Provides help for table-reference commands.
- **whoami:** Provides information about the user.

DDL Commands

- create: Creates a table.
- list: Lists all the tables in HBase.
- disable: Disables a table.
- is_disabled: Verifies whether a table is disabled.
- enable: Enables a table.
- is_enabled: Verifies whether a table is enabled.
- describe: Provides the description of a table.
- alter: Alters a table.
- exists: Verifies whether a table exists.
- drop: Drops a table from HBase.

DML Commands

- put: Puts a cell value at a specified column in a specified row in a particular table.
- get: Fetches the contents of row or a cell.
- delete: Deletes a cell value in a table.
- deleteall: Deletes all the cells in a given row.
- scan: Scans and returns the table data.
- count: Counts and returns the number of rows in a table.
- truncate: Disables, drops, and recreates a specified table.

Create table

Row key	personal data	professional data

- create 'emp', 'personal data', 'professional data'

Describe table

- describe 'tablename'

Create data

- To create data in an HBase table, the following commands and methods are used:
 - put command,
 - add() method of Put class, and
 - put() method of HTable class.

Example:

COLUMN FAMILIES

Row key	personal data		professional data	
empid	name	city	designation	salary
1	raju	hyderabad	manager	50,000
2	ravi	chennai	sr.engineer	30,000
3	rajesh	delhi	jr.engineer	25,000

Example:

- `put 'emp','1','personal data:name','raju'`
- `put 'emp','1','personal data:city','mumbai'`
- `put 'emp','1','professional data:designation','manager'`
- `put 'emp','1','professional data:salary','50000'`
-

Display:

- scan 'tablename'

Update:

- `put 'emp',1 , 'personal:city','Delhi'`

Read Data:

- `get 'emp', '1'`
- `get 'emp', '1', {COLUMN=>'personal:name'}`

Read Data:

- delete 'emp', '1', 'personal data:city'

Count and truncate

- You can count the number of rows of a table using the count command.
count 'emp'
- This command disables drops and recreates a table.
truncate 'tablename'

Thank you

This presentation is created using LibreOffice Impress 5.3.2.2, can be used freely as per GNU General Public License



@mitu_skillologies



/mITuSkillologies



@mitu_group

Web Resources

<http://mitu.co.in>
<http://tusharkute.com>

Blogs

<http://digitallocha.blogspot.in>
<http://kyamputar.blogspot.in>

contact@mitu.co.in
tushar@tusharkute.com