

An In-Depth Look at the HBase Architecture



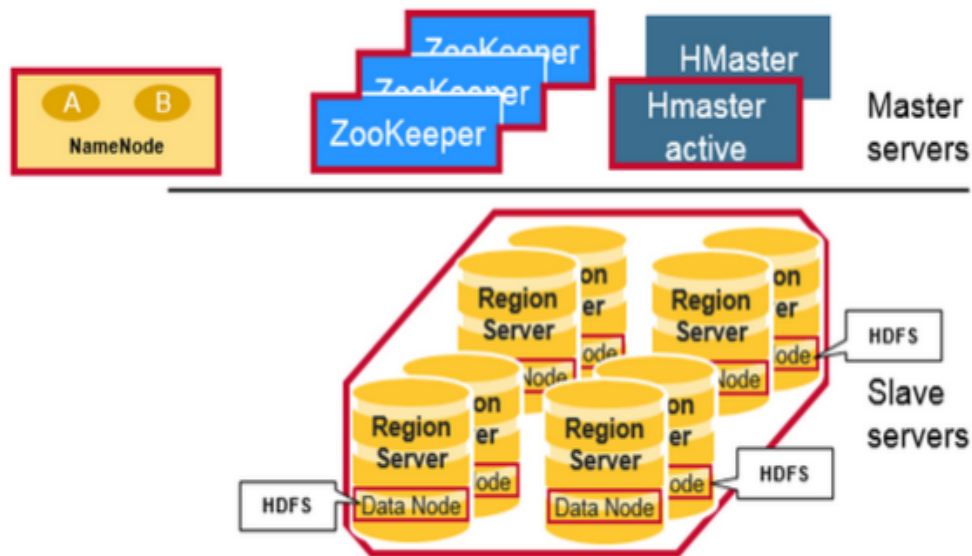
In this blog post, I'll give you an in-depth look at the HBase architecture and its main benefits over NoSQL data store solutions. Be sure and read the first blog post in this series, titled ["HBase and MapR-DB: Designed for Distribution, Scale, and Speed."](#)

HBase Architectural Components

Physically, HBase is composed of three types of servers in a master slave type of architecture. Region servers serve data for reads and writes. When accessing data, clients communicate with HBase RegionServers directly. Region assignment, DDL (create, delete tables) operations are handled by the HBase Master process. Zookeeper, which is part of HDFS, maintains a live cluster state.

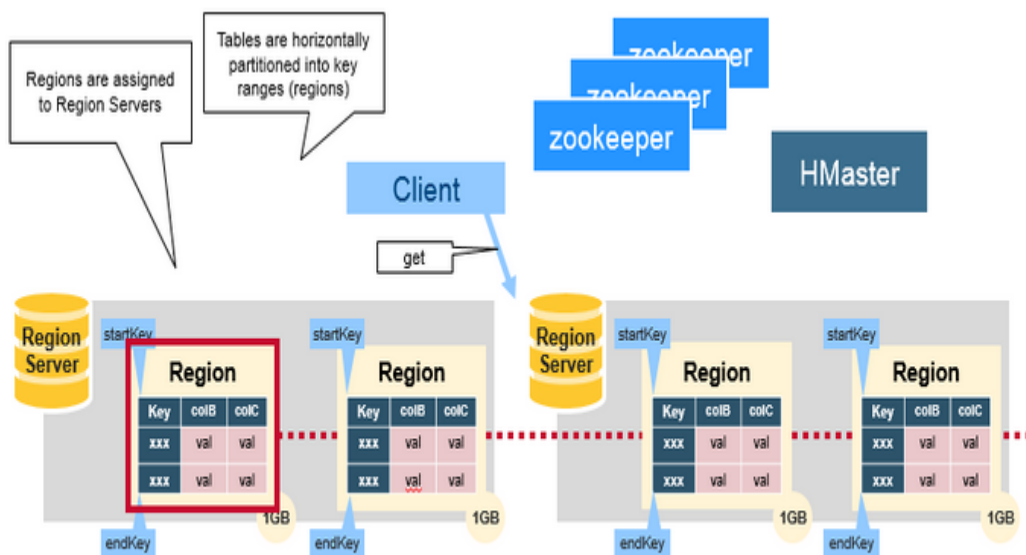
The Hadoop DataNode stores the data that the Region Server is managing. All HBase data is stored in HDFS files. Region Servers are collocated with the HDFS DataNodes, which enable data locality (putting the data close to where it is needed) for the data served by the RegionServers. HBase data is local when it is written, but when a region is moved, it is not local until compaction.

The NameNode maintains metadata information for all the physical data blocks that comprise the files.



Regions

HBase Tables are divided horizontally by row key range into “Regions.” A region contains all rows in the table between the region’s start key and end key. Regions are assigned to the nodes in the cluster, called “Region Servers,” and these serve data for reads and writes. A region server can serve about 1,000 regions.

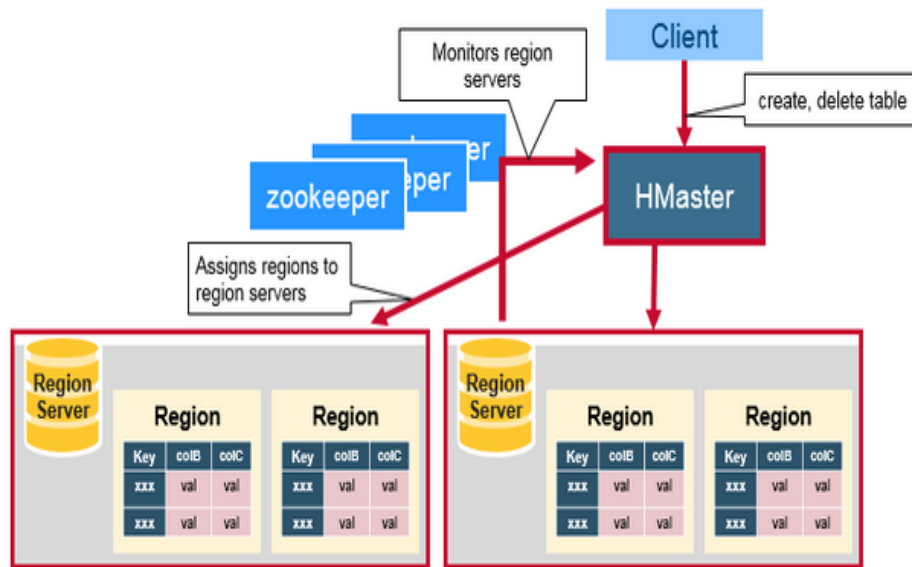


HBase HMaster

Region assignment, DDL (create, delete tables) operations are handled by the HBase Master.

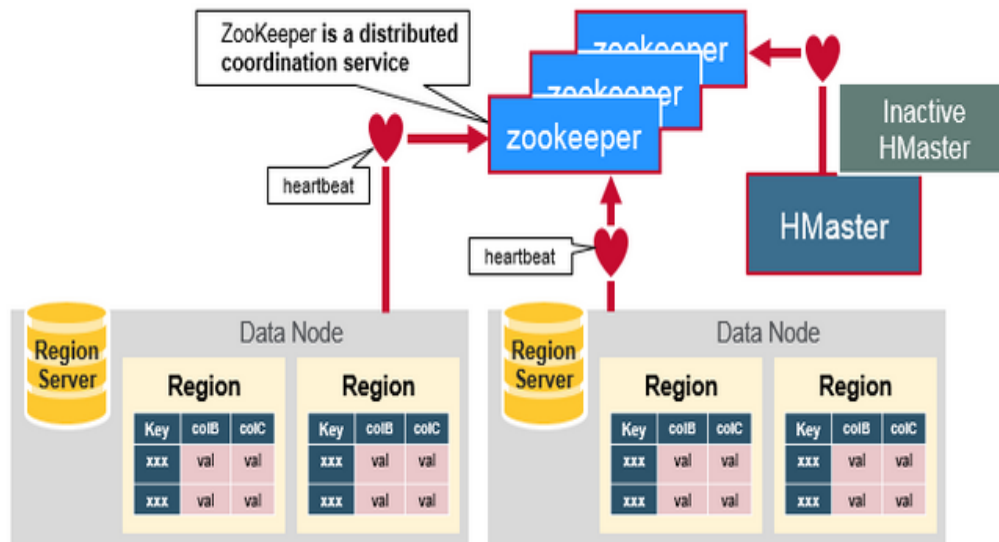
A master is responsible for:

- Coordinating the region servers
 - Assigning regions on startup , re-assigning regions for recovery or load balancing
 - Monitoring all RegionServer instances in the cluster (listens for notifications from zookeeper)
- Admin functions
 - Interface for creating, deleting, updating tables



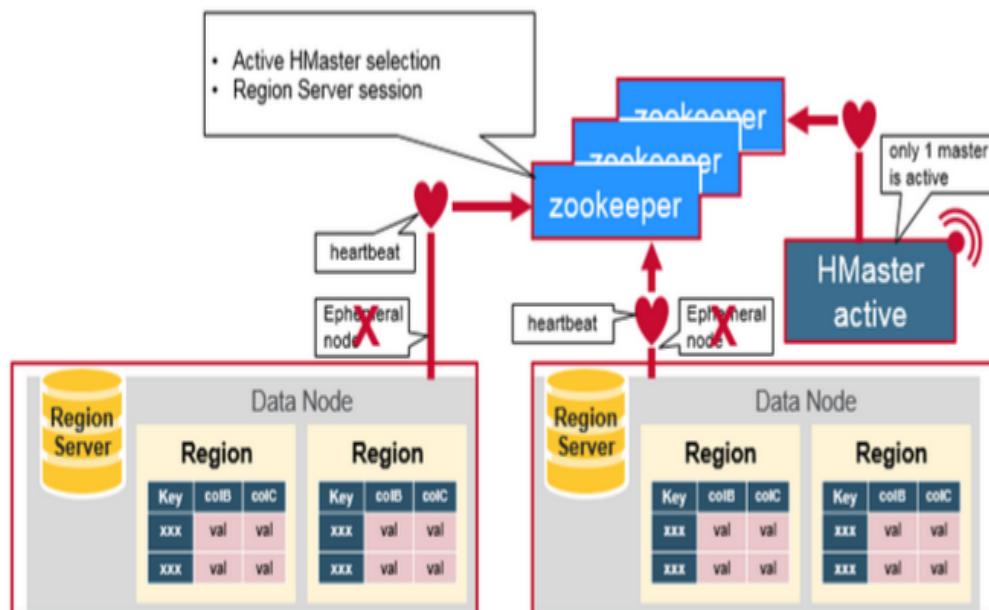
ZooKeeper: The Coordinator

HBase uses ZooKeeper as a distributed coordination service to maintain server state in the cluster. Zookeeper maintains which servers are alive and available, and provides server failure notification. Zookeeper uses consensus to guarantee common shared state. Note that there should be three or five machines for consensus.



How the Components Work Together

Zookeeper is used to coordinate shared state information for members of distributed systems. Region servers and the active HMaster connect with a session to ZooKeeper. The ZooKeeper maintains ephemeral nodes for active sessions via heartbeats.



Each Region Server creates an ephemeral node. The HMaster monitors these nodes to discover available region servers, and it also monitors these nodes for server failures. HMaster tries to create an ephemeral node. Zookeeper determines the first one and uses it to make sure that only one master is active. The active HMaster sends heartbeats to Zookeeper, and the inactive HMaster listens for notifications of the active HMaster failure.

If a region server or the active HMaster fails to send a heartbeat, the session is expired and the corresponding ephemeral node is deleted. Listeners for updates will be notified of the deleted nodes. The

active HMaster listens for region servers, and will recover region servers on failure. The Inactive HMaster listens for active HMaster failure, and if an active HMaster fails, the inactive HMaster becomes active.

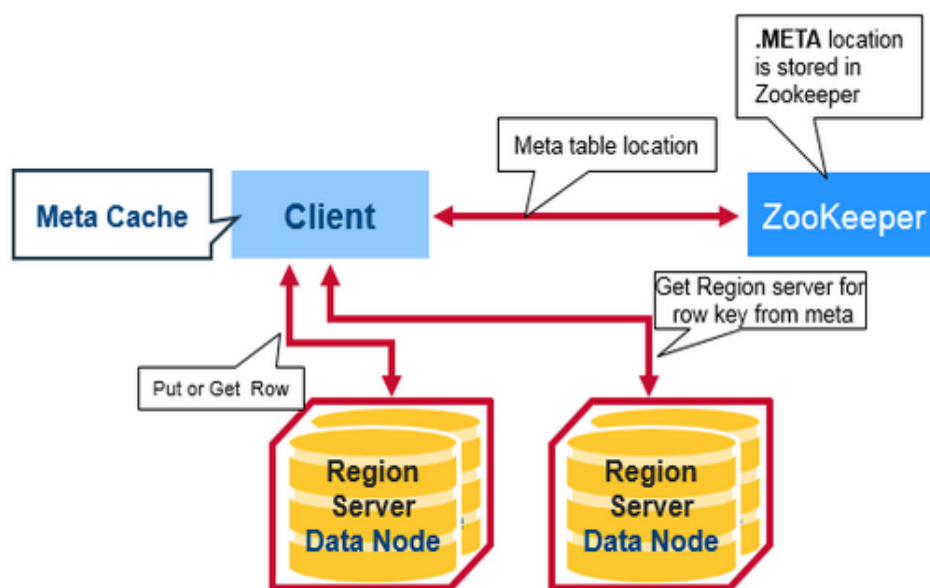
HBase First Read or Write

There is a special HBase Catalog table called the META table, which holds the location of the regions in the cluster. ZooKeeper stores the location of the META table.

This is what happens the first time a client reads or writes to HBase:

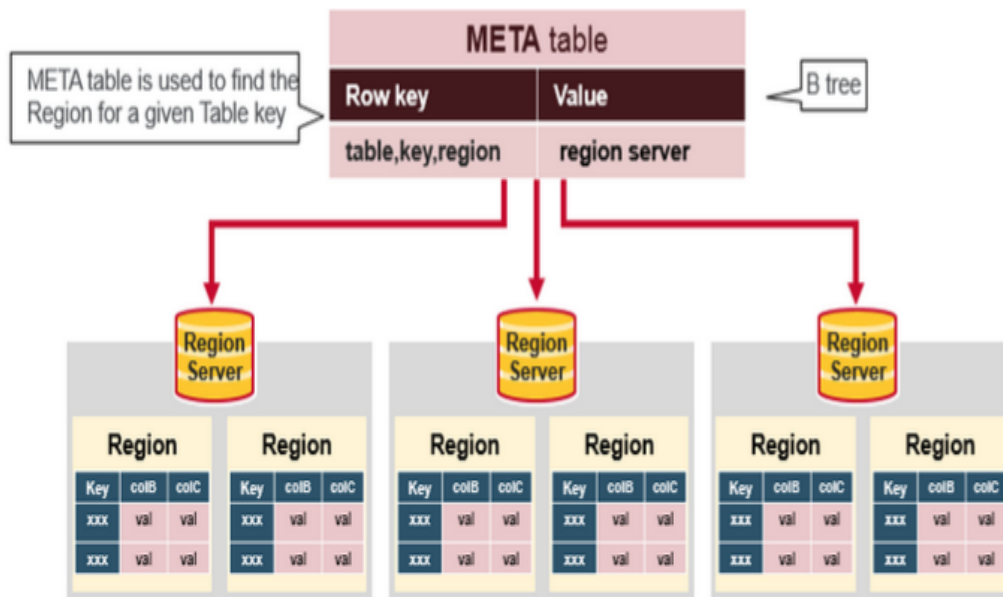
1. The client gets the Region server that hosts the META table from ZooKeeper.
2. The client will query the .META. server to get the region server corresponding to the row key it wants to access. The client caches this information along with the META table location.
3. It will get the Row from the corresponding Region Server.

For future reads, the client uses the cache to retrieve the META location and previously read row keys. Over time, it does not need to query the META table, unless there is a miss because a region has moved; then it will re-query and update the cache.



HBase Meta Table

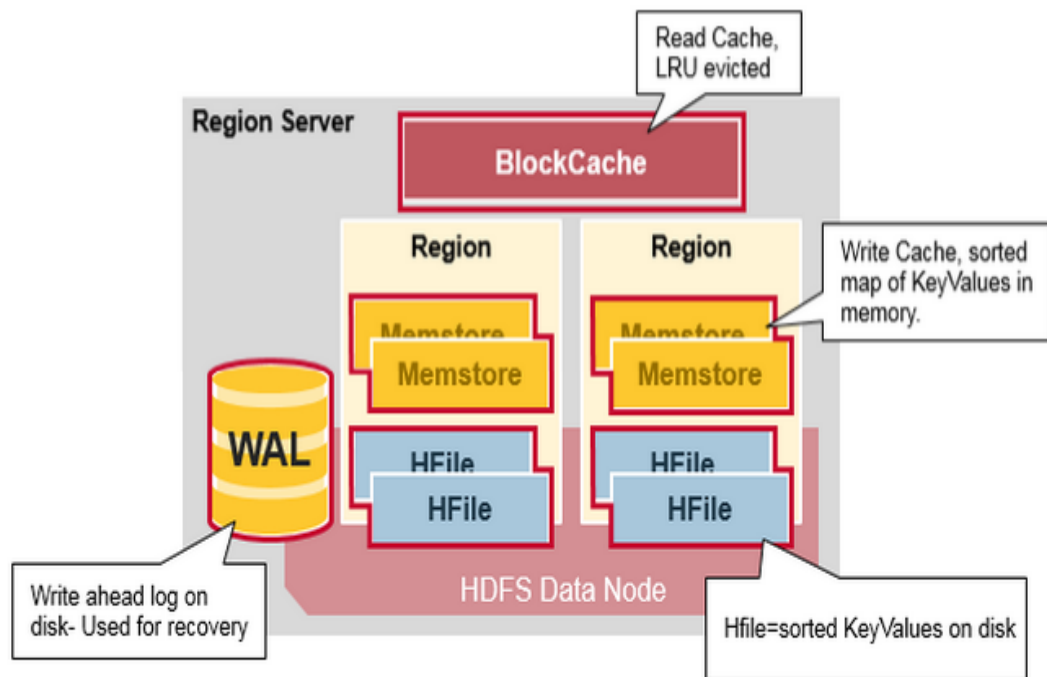
- This META table is an HBase table that keeps a list of all regions in the system.
- The .META. table is like a b tree.
- The .META. table structure is as follows:
 - Key: region start key, region id
 - Values: RegionServer



Region Server Components

A Region Server runs on an HDFS data node and has the following components:

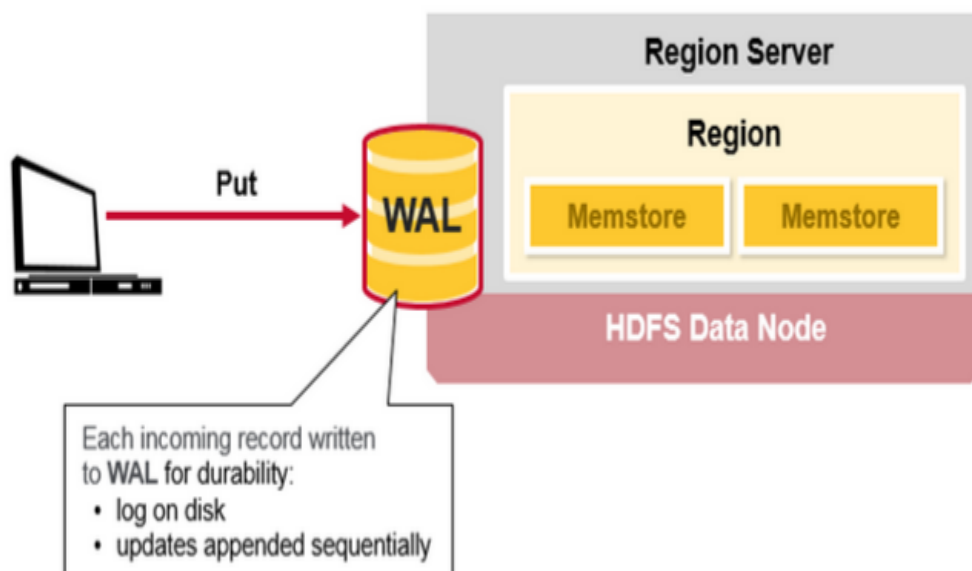
- WAL: Write Ahead Log is a file on the distributed file system. The WAL is used to store new data that hasn't yet been persisted to permanent storage; it is used for recovery in the case of failure.
- BlockCache: is the read cache. It stores frequently read data in memory. Least Recently Used data is evicted when full.
- MemStore: is the write cache. It stores new data which has not yet been written to disk. It is sorted before writing to disk. There is one MemStore per column family per region.
- Hfiles store the rows as sorted KeyValues on disk.



HBase Write Steps (1)

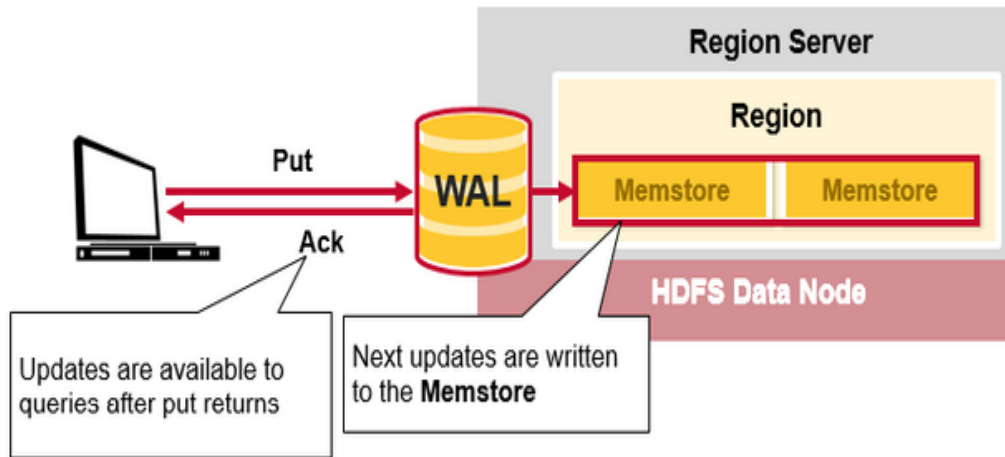
When the client issues a Put request, the first step is to write the data to the write-ahead log, the WAL:

- Edits are appended to the end of the WAL file that is stored on disk.
- The WAL is used to recover not-yet-persisted data in case a server crashes.



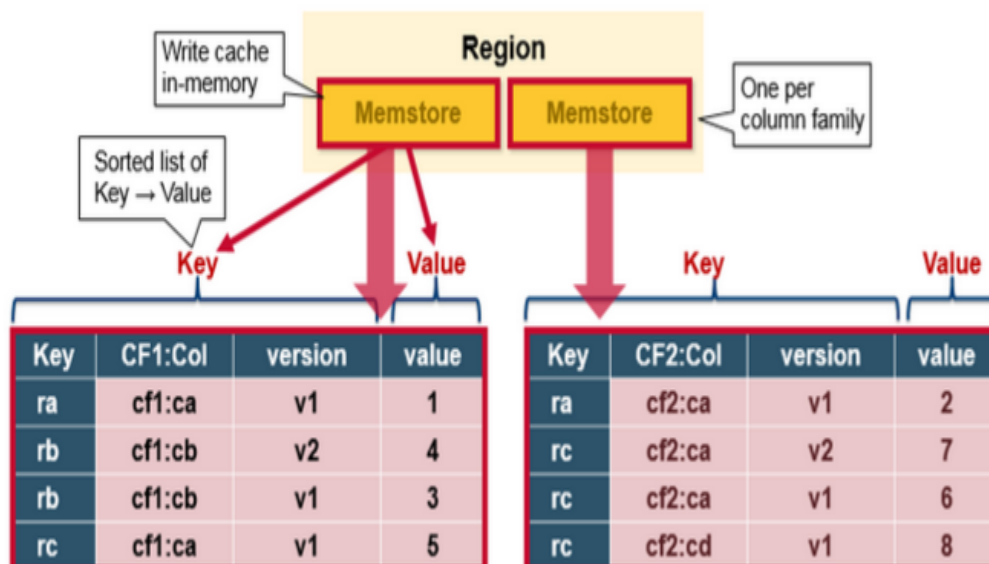
HBase Write Steps (2)

Once the data is written to the WAL, it is placed in the MemStore. Then, the put request acknowledgement returns to the client.



HBase MemStore

The MemStore stores updates in memory as sorted KeyValues, the same as it would be stored in an HFile. There is one MemStore per column family. The updates are sorted per column family.

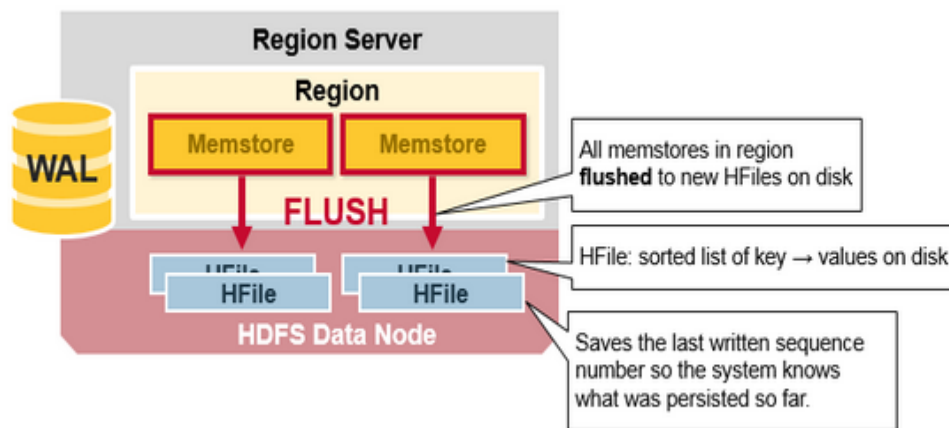


HBase Region Flush

When the MemStore accumulates enough data, the entire sorted set is written to a new HFile in HDFS. HBase uses multiple HFiles per column family, which contain the actual cells, or KeyValue instances. These files are created over time as KeyValue edits sorted in the MemStores are flushed as files to disk.

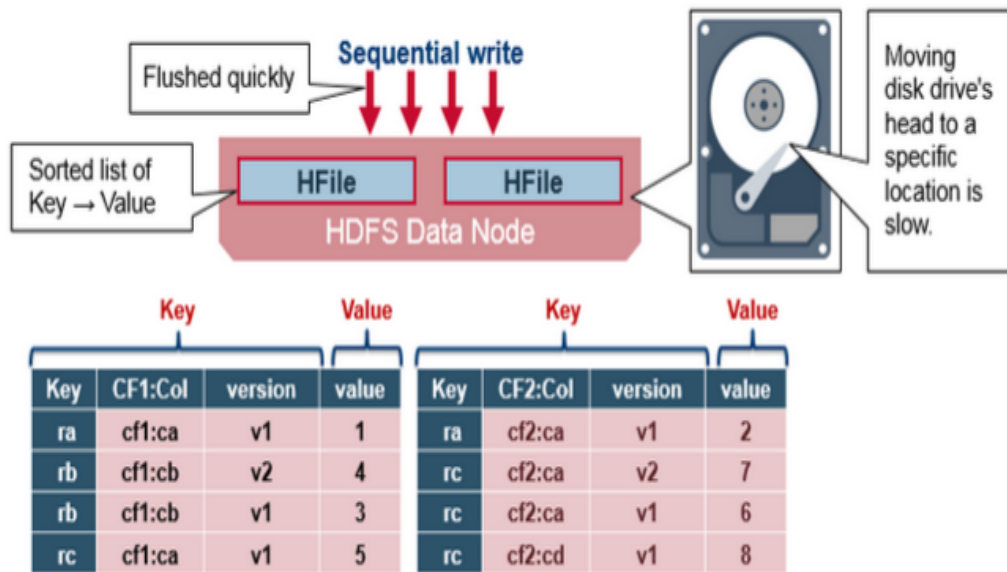
Note that this is one reason why there is a limit to the number of column families in HBase. There is one MemStore per CF; when one is full, they all flush. It also saves the last written sequence number so the system knows what was persisted so far.

The highest sequence number is stored as a meta field in each HFile, to reflect where persisting has ended and where to continue. On region startup, the sequence number is read, and the highest is used as the sequence number for new edits.



HBase HFile

Data is stored in an HFile which contains sorted key/values. When the MemStore accumulates enough data, the entire sorted KeyValue set is written to a new HFile in HDFS. This is a sequential write. It is very fast, as it avoids moving the disk drive head.

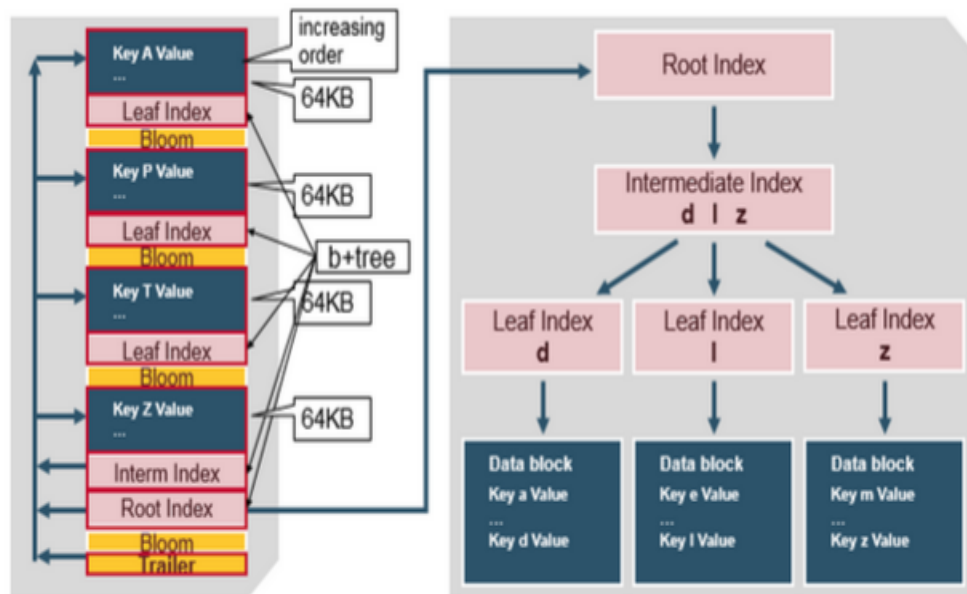


HBase HFile Structure

An HFile contains a multi-layered index which allows HBase to seek to the data without having to read the whole file. The multi-level index is like a b+tree:

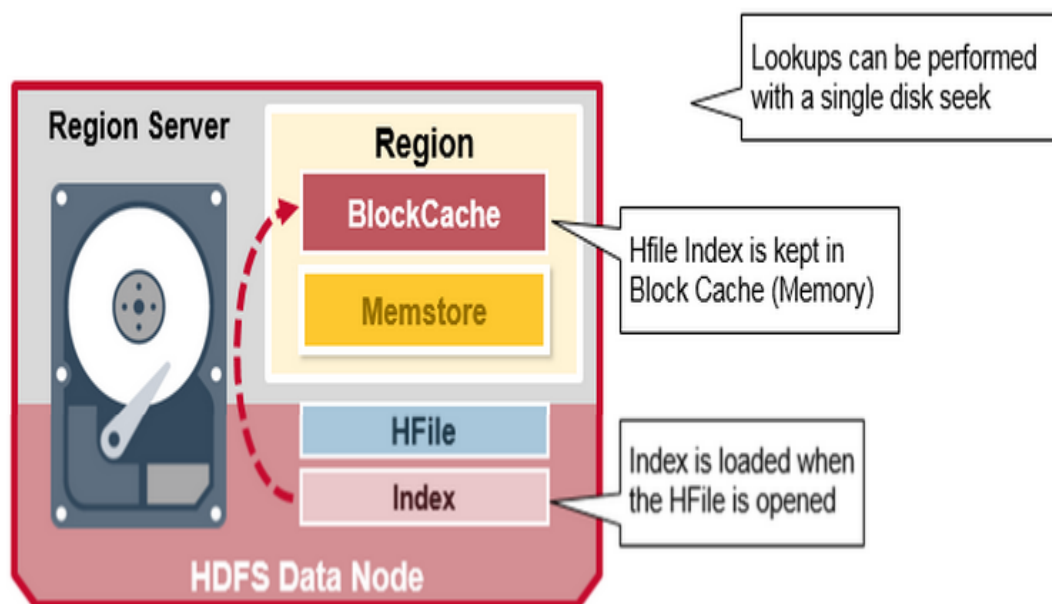
- Key value pairs are stored in increasing order
- Indexes point by row key to the key value data in 64KB “blocks”
- Each block has its own leaf-index
- The last key of each block is put in the intermediate index
- The root index points to the intermediate index

The trailer points to the meta blocks, and is written at the end of persisting the data to the file. The trailer also has information like bloom filters and time range info. Bloom filters help to skip files that do not contain a certain row key. The time range info is useful for skipping the file if it is not in the time range the read is looking for.



HFile Index

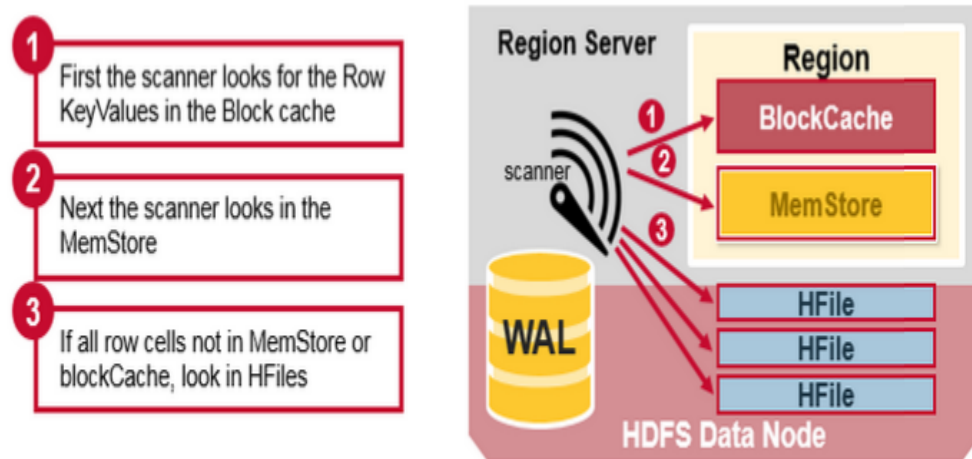
The index, which we just discussed, is loaded when the HFile is opened and kept in memory. This allows lookups to be performed with a single disk seek.



HBase Read Merge

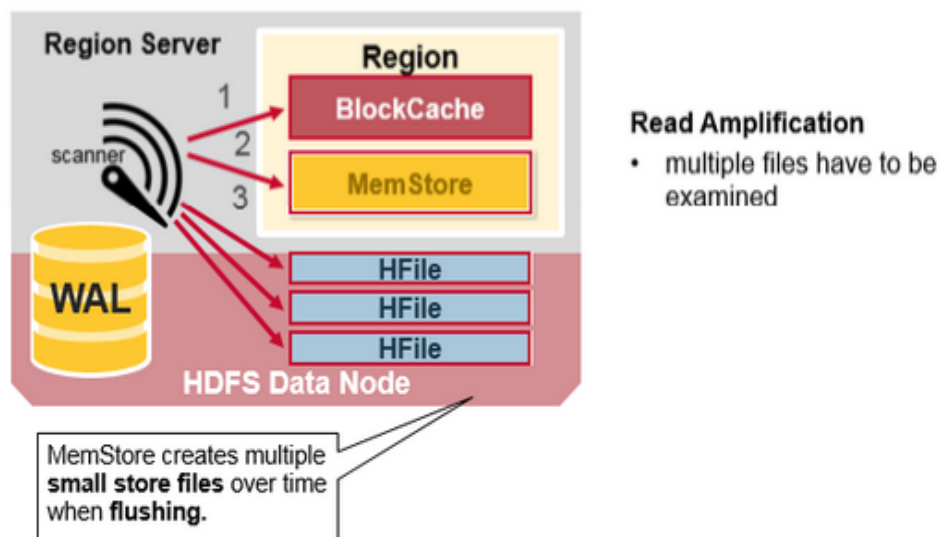
We have seen that the Key Value cells corresponding to one row can be in multiple places, row cells already persisted are in Hfiles, recently updated cells are in the MemStore, and recently read cells are in the Block cache. So when you read a row, how does the system get the corresponding cells to return? A Read merges Key Values from the block cache, MemStore, and HFiles in the following steps:

1. First, the scanner looks for the Row cells in the Block cache - the read cache. Recently Read Key Values are cached here, and Least Recently Used are evicted when memory is needed.
2. Next, the scanner looks in the MemStore, the write cache in memory containing the most recent writes.
3. If the scanner does not find all of the row cells in the MemStore and Block Cache, then HBase will use the Block Cache indexes and bloom filters to load HFiles into memory, which may contain the target row cells.



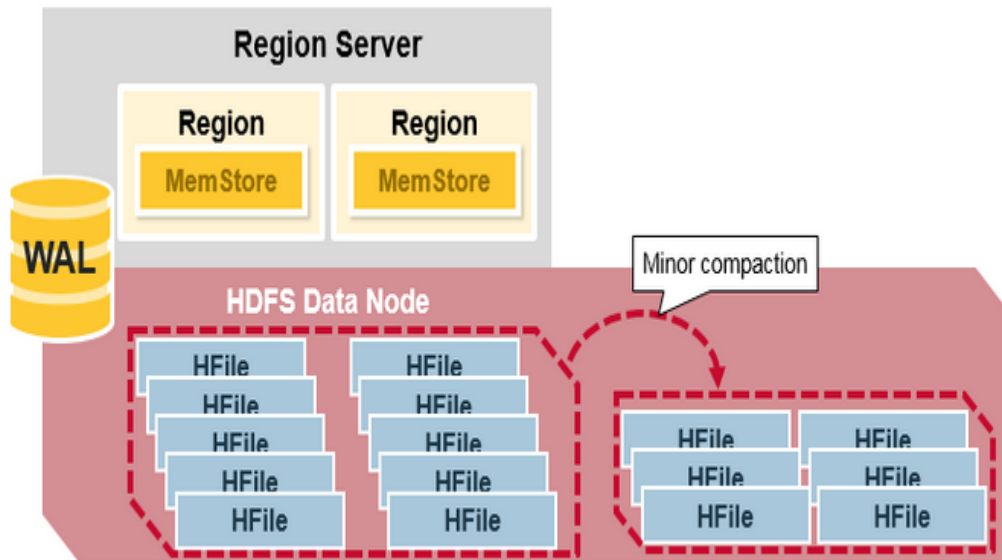
HBase Read Merge

As discussed earlier, there may be many HFiles per MemStore, which means for a read, multiple files may have to be examined, which can affect the performance. This is called read amplification.



HBase Minor Compaction

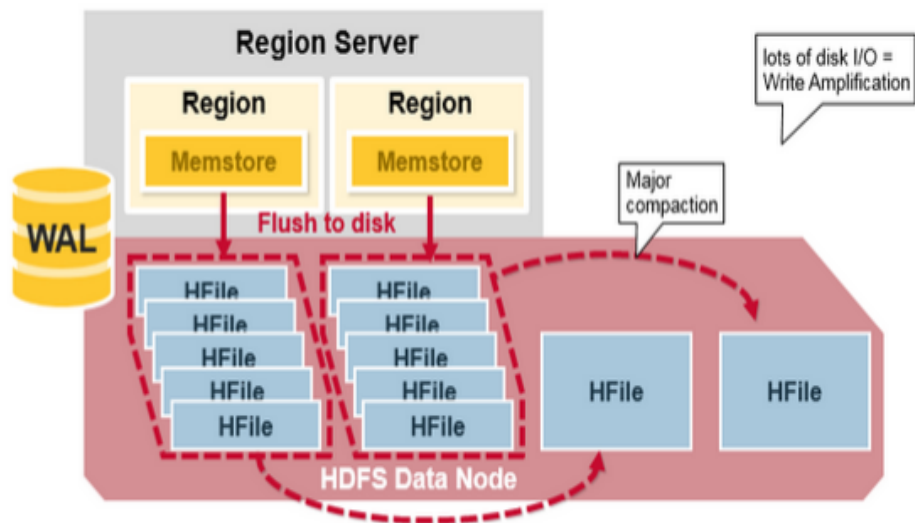
HBase will automatically pick some smaller HFiles and rewrite them into fewer bigger Hfiles. This process is called minor compaction. Minor compaction reduces the number of storage files by rewriting smaller files into fewer but larger ones, performing a merge sort.



HBase Major Compaction

Major compaction merges and rewrites all the HFiles in a region to one HFile per column family, and in the process, drops deleted or expired cells. This improves read performance; however, since major compaction rewrites all of the files, lots of disk I/O and network traffic might occur during the process. This is called write amplification.

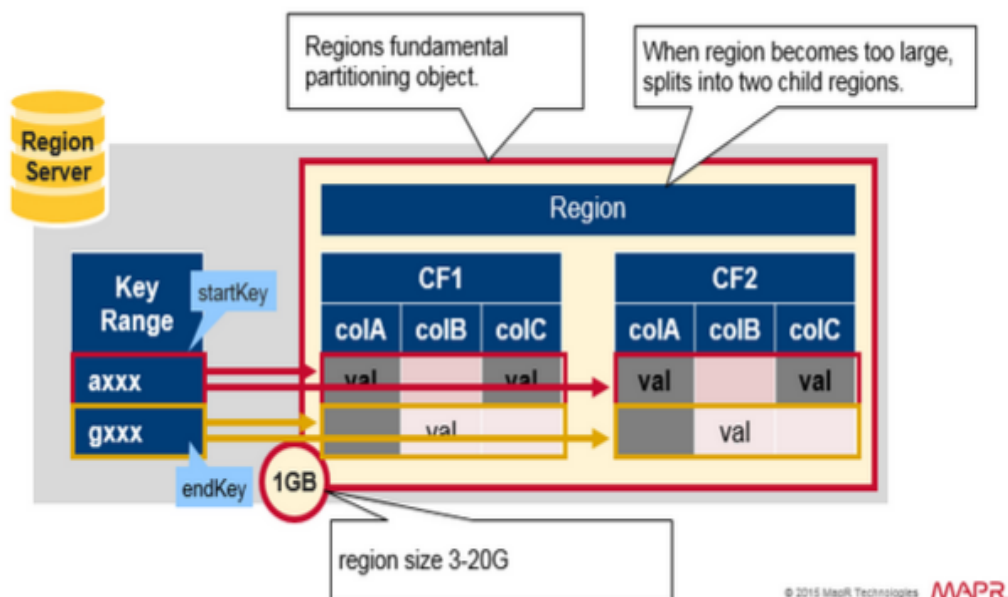
Major compactions can be scheduled to run automatically. Due to write amplification, major compactions are usually scheduled for weekends or evenings. Note that MapR-DB has made improvements and does not need to do compactions. A major compaction also makes any data files that were remote, due to server failure or load balancing, local to the region server.



Region = Contiguous Keys

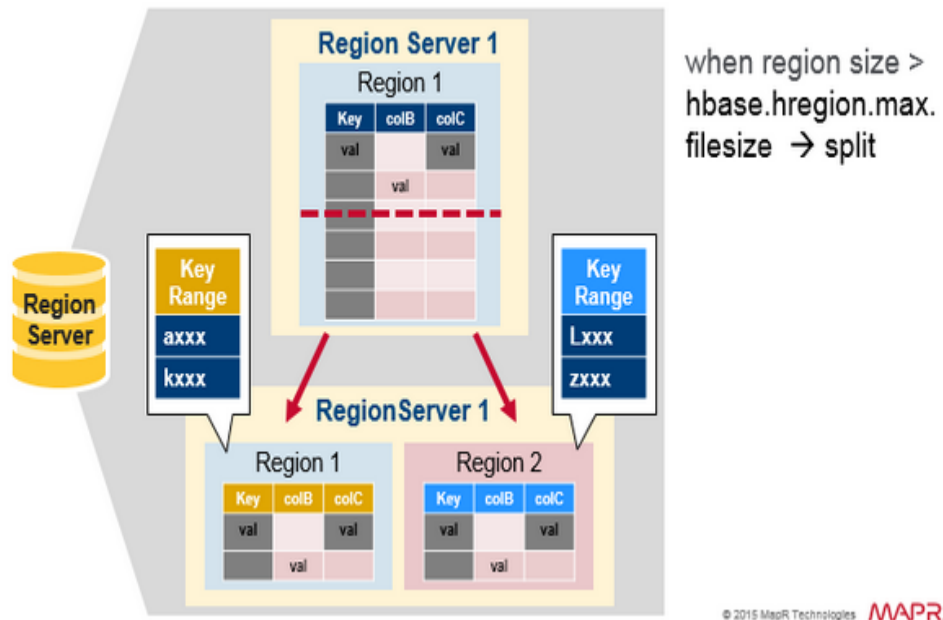
Let's do a quick review of regions:

- A table can be divided horizontally into one or more regions. A region contains a contiguous, sorted range of rows between a start key and an end key
- Each region is 1GB in size (default)
- A region of a table is served to the client by a RegionServer
- A region server can serve about 1,000 regions (which may belong to the same table or different tables)



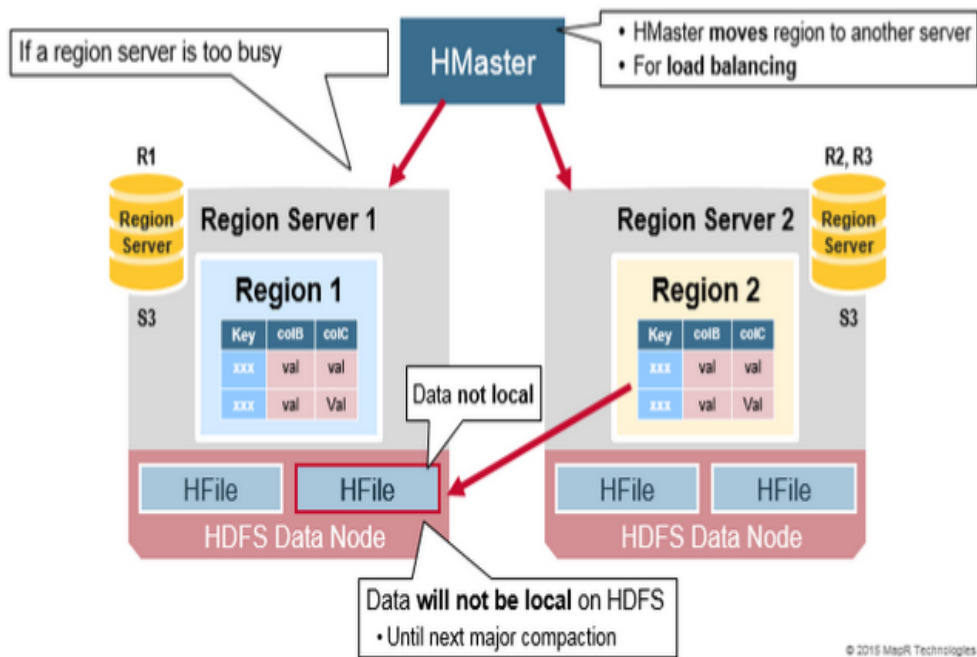
Region Split

Initially there is one region per table. When a region grows too large, it splits into two child regions. Both child regions, representing one-half of the original region, are opened in parallel on the same Region server, and then the split is reported to the HMaster. For load balancing reasons, the HMaster may schedule for new regions to be moved off to other servers.



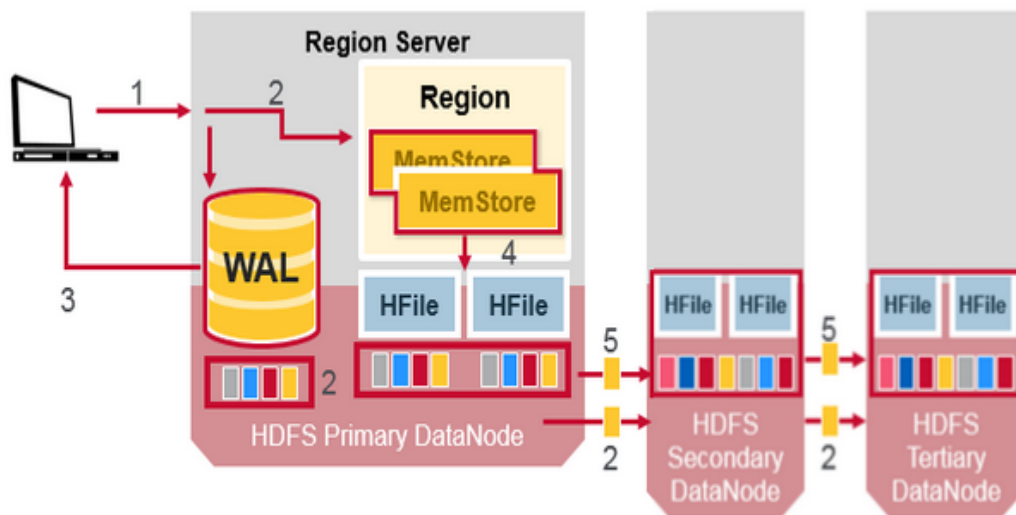
Read Load Balancing

Splitting happens initially on the same region server, but for load balancing reasons, the HMaster may schedule for new regions to be moved off to other servers. This results in the new Region server serving data from a remote HDFS node until a major compaction moves the data files to the Regions server's local node. HBase data is local when it is written, but when a region is moved (for load balancing or recovery), it is not local until major compaction.



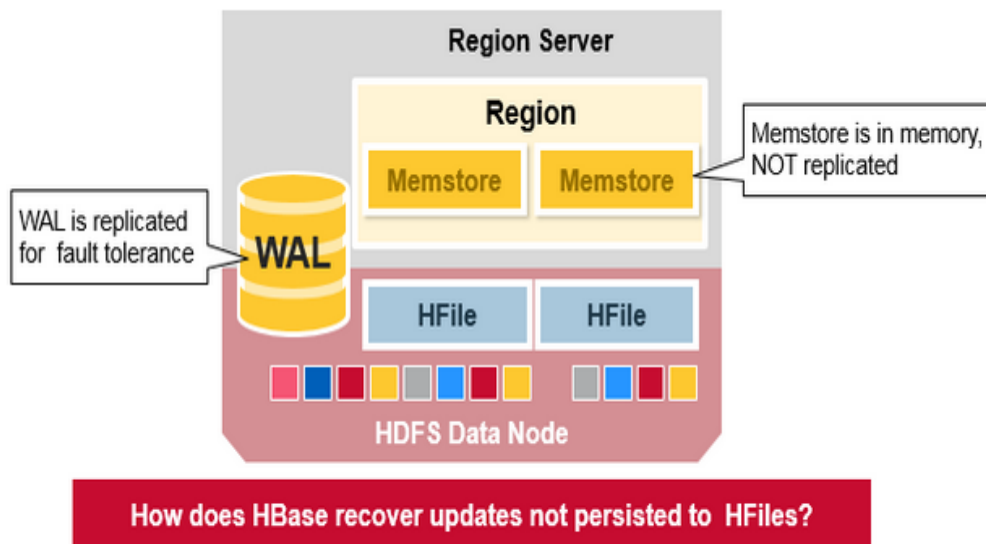
HDFS Data Replication

All writes and Reads are to/from the primary node. HDFS replicates the WAL and HFile blocks. HFile block replication happens automatically. HBase relies on HDFS to provide the data safety as it stores its files. When data is written in HDFS, one copy is written locally, and then it is replicated to a secondary node, and a third copy is written to a tertiary node.



HDFS Data Replication (2)

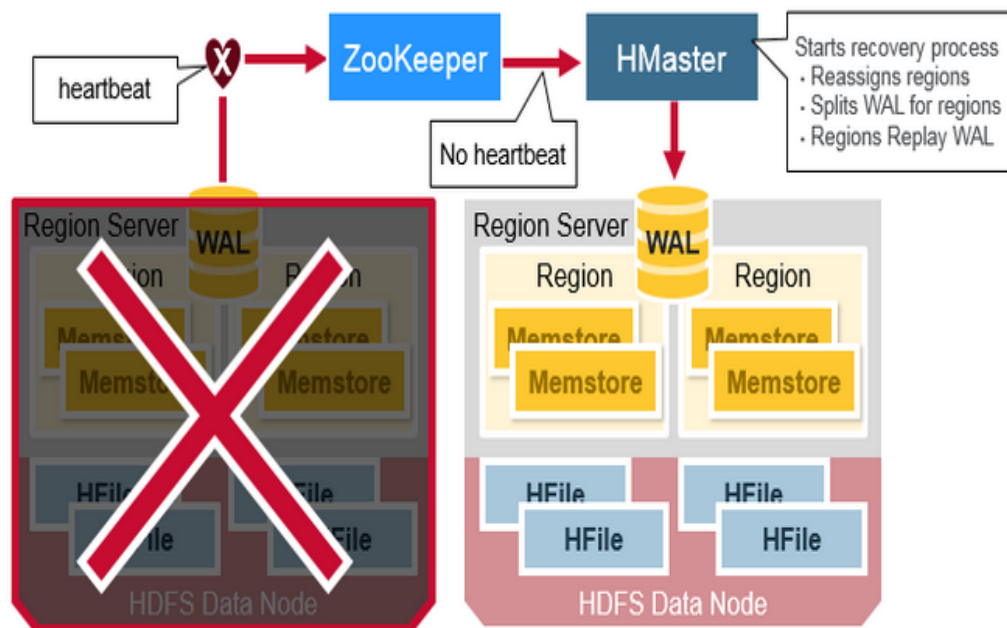
The WAL file and the Hfiles are persisted on disk and replicated, so how does HBase recover the MemStore updates not persisted to HFiles? See the next section for the answer.



HBase Crash Recovery

When a RegionServer fails, Crashed Regions are unavailable until detection and recovery steps have happened. Zookeeper will determine Node failure when it loses region server heart beats. The HMaster will then be notified that the Region Server has failed.

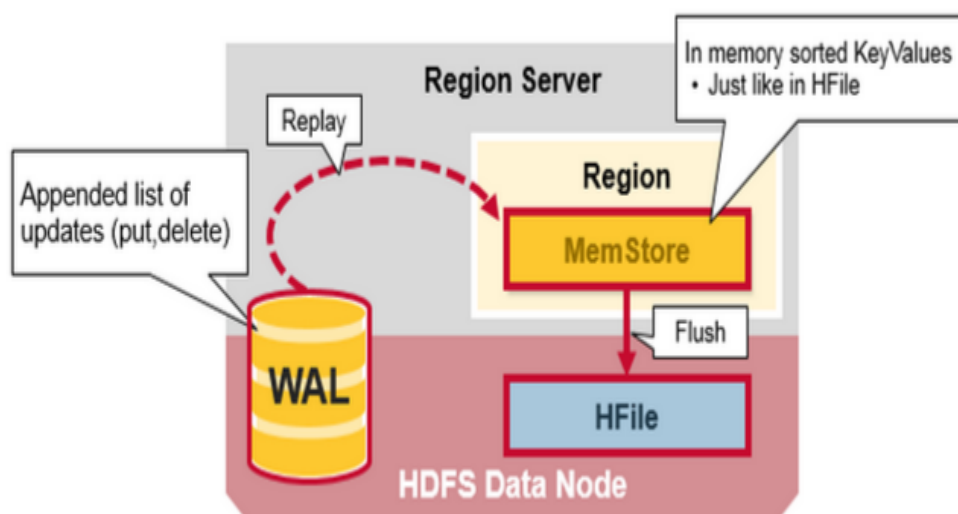
When the HMaster detects that a region server has crashed, the HMaster reassigns the regions from the crashed server to active Region servers. In order to recover the crashed region server's memstore edits that were not flushed to disk. The HMaster splits the WAL belonging to the crashed region server into separate files and stores these file in the new region servers' data nodes. Each Region Server then replays the WAL from the respective split WAL, to rebuild the memstore for that region.



Data Recovery

WAL files contain a list of edits, with one edit representing a single put or delete. Edits are written chronologically, so, for persistence, additions are appended to the end of the WAL file that is stored on disk.

What happens if there is a failure when the data is still in memory and not persisted to an HFile? The WAL is replayed. Replaying a WAL is done by reading the WAL, adding and sorting the contained edits to the current MemStore. At the end, the MemStore is flush to write changes to an HFile.



Apache HBase Architecture Benefits

HBase provides the following benefits:

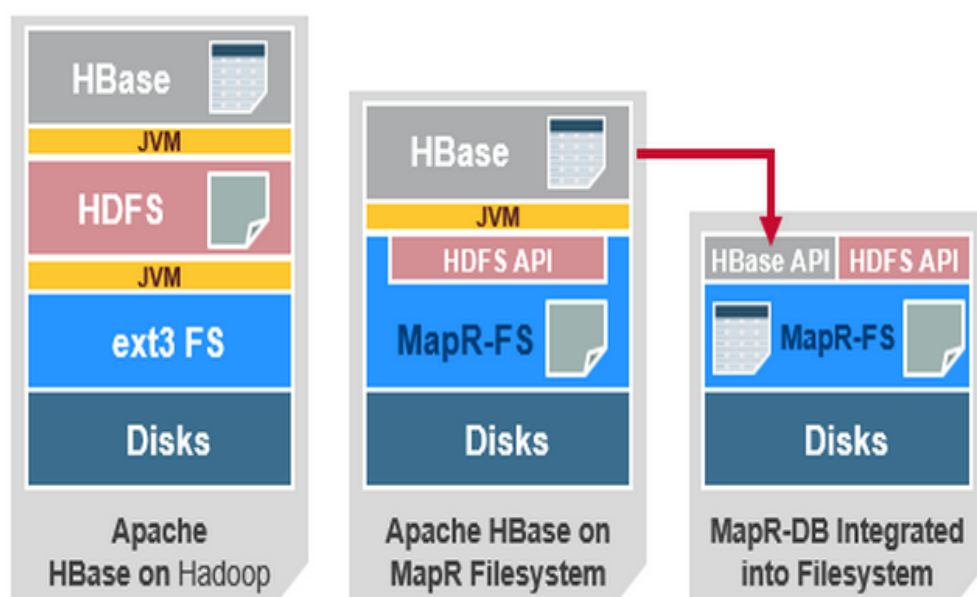
- **Strong consistency model**
 - When a write returns, all readers will see same value
- **Scales automatically**
 - Regions split when data grows too large
 - Uses HDFS to spread and replicate data
- **Built-in recovery**
 - Using Write Ahead Log (similar to journaling on file system)
- **Integrated with Hadoop**
 - MapReduce on HBase is straightforward

Apache HBase Has Problems Too...

- **Business continuity reliability:**
 - WAL replay slow
 - Slow complex crash recovery
 - Major Compaction I/O storms

MapR-DB with MapR-FS does not have these problems

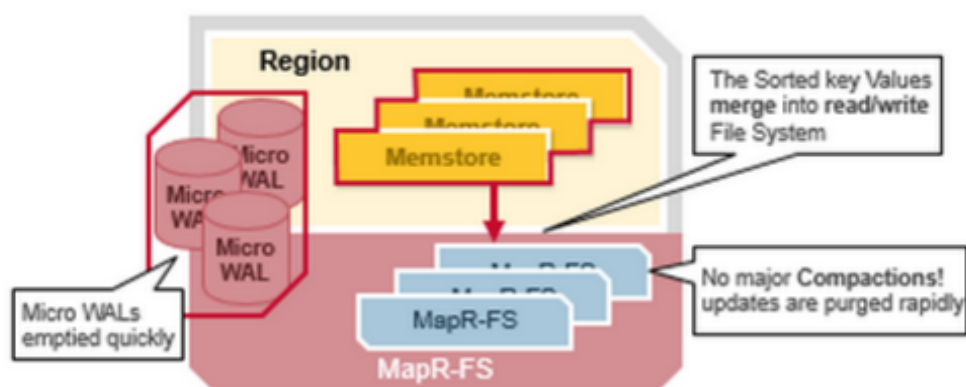
The diagram below compares the application stacks for Apache HBase on top of HDFS on the left, Apache HBase on top of MapR's read/write file system MapR-FS in the middle, and MapR-DB and MapR-FS in a Unified Storage Layer on the right.



MapR-DB exposes the same HBase API and the Data model for MapR-DB is the same as for Apache HBase. However the MapR-DB implementation integrates table storage into the MapR file system, eliminating all JVM layers and interacting directly with disks for both file and table storage.

MapR-DB offers many benefits over HBase, while maintaining the virtues of the HBase API and the idea of data being sorted according to primary key. MapR-DB provides operational benefits such as no compaction delays and automated region splits that do not impact the performance of the database. The tables in MapR-DB can also be isolated to certain machines in a cluster by utilizing the topology feature of MapR. The final differentiator is that MapR-DB is just plain fast, due primarily to the fact that it is tightly integrated into the MapR file system itself, rather than being layered on top of a distributed file system that is layered on top of a conventional file system.

Key differences between MapR-DB and Apache HBase



- Tables part of the MapR Read/Write File system
 - Guaranteed data locality
- Smarter load balancing
 - Uses container Replicas
- Smarter fail over
 - Uses container replicas
- Multiple small WALs
 - Faster recovery
- Memstore Flushes Merged into Read/Write File System
 - No compaction !

[You can take this free On Demand training to learn more about MapR-FS and MapR-DB.](#)

In this blog post, you learned more about the HBase architecture and its main benefits over NoSQL data store solutions. If you have any questions about HBase, please ask them in the comments section below.