# HBase Internals and Operations

**Biju Nair**
**Software Engineer**
**bnair10@bloomberg.net**

**Bloomberg**
Engineering

**TechAtBloomberg.com**

# Agenda

- Introduction to HBase

- Operating HBase

- Questions

**Bloomberg**

Engineering

# Bloomberg in a nutshell



The **Bloomberg Terminal** delivers a diverse array of information on a single platform to facilitate financial decision-making.

**Bloomberg**

Engineering

# Bloomberg technology by the numbers

- **5,000+** software engineers

- **150+** technologists and data scientists devoted to machine learning

- One of the largest private networks in the world

- **120 billion** pieces of data from the financial markets each day, with a peak of more than 10 million messages/second

- **2 million** news stories ingested / published each day (500+ news stories ingested/second)

- News content from **125K+ sources**

- Over **1 billion** messages and Instant Bloomberg (IB) chats handled daily

**TechAtBloomberg.com**

**Bloomberg**

Engineering

# HBase at Bloomberg

- Started with **v0.94.6**

- **>2 billion** reads per day

- **>1 billion** writes per day

- **51+ TB** of compressed data stored in HBase

**Bloomberg**

Engineering

# HBase Principles

- Ordered Key Value Store

- Distributed shared nothing

**Bloomberg**

Engineering

# Key Value

...

| | |
|---|---|
| Key-9999 | Value-a |
| Key-9998 | Value-b |
| Key-9997 | Value-c |
| Key-9996 | Value-d |
| Key-9995 | Value-e |
| Key-9994 | Value-f |

...

Bloomberg

Engineering

# Ordered Key Value



|  | ... | |
|---|---|
| Key-9999 | Value-a |
| Key-9998 | Value-b |
| Key-9997 | Value-c |
| Key-9996 | Value-d |
| Key-9995 | Value-e |
| Key-9994 | Value-a |
| Key-9993 | Value-g |
|  | ... | |

Lexicographic order

**Bloomberg**

Engineering

# Ordered Key Value

| | | |
|---|---|---|
| Key-9999 | Value-a | |
| Key-9998 | Value-b | |
| Key-9997 | Value-c | |
| Key-9996 | Value-d | Region |
| Key-9995 | Value-e | |
| Key-9994 | Value-a | |
| Key-9993 | Value-g | |

Lexicographic order

**Bloomberg**

Engineering

# Distributed Ordered Key Value

Bloomberg

Engineering

# Distributed Ordered Key Value

**Bloomberg**

Engineering

# Table Row View

# Table Row View

| | |
|---|---|
| R11\|col1\|1234567 | Value-A |
| R11\|col2\|1234567 | Value-B |
| R11\|col3\|1234567 | Value-C |
| R11\|col4\|1234567 | Value-D |

| R11 | Col1 | Col2 | Col3 | Col4 |
|---|---|---|---|---|
| | Value-A | Value-B | Value-C | Value-D |

**Bloomberg**

Engineering

# Versioning

| | |
|---|---|
| R11\|col1\|1234567 | Value-A1 |
| R11\|col1\|1234566 | Value-A |
| R11\|col2\|1234567 | Value-B |
| R11\|col3\|1234567 | Value-CC |
| R11\|col3\|1234563 | Value-C |
| R11\|col4\|1234567 | Value-DD |
| R11\|col4\|1234560 | Value-D1 |
| R11\|col4\|1234557 | Value-D |

Descending order

Bloomberg
Engineering

# Column Family

Bloomberg
Engineering

# ACIDity

- **A**tomic at row level

- **C**onsistent to a point in time before the request

- **I**solation through MVCC (reads) and row locks (mutations)

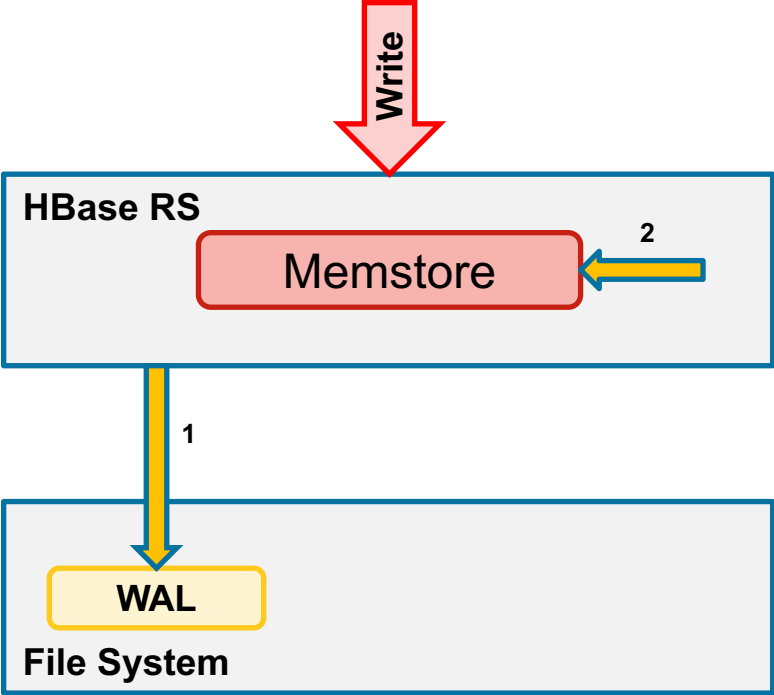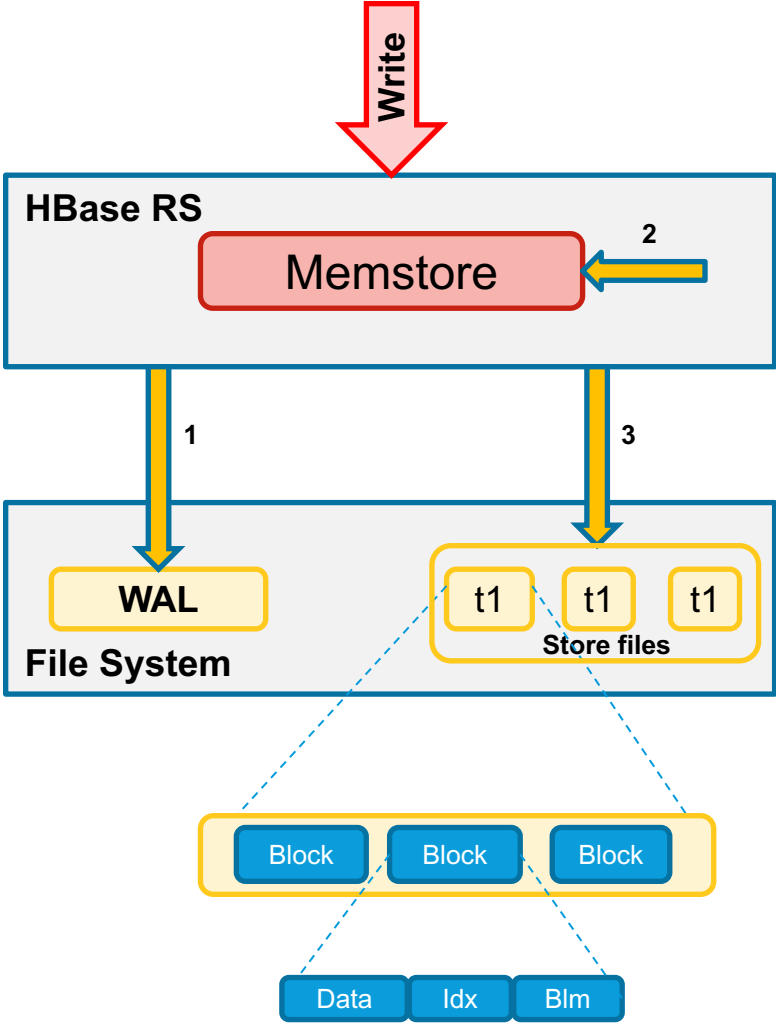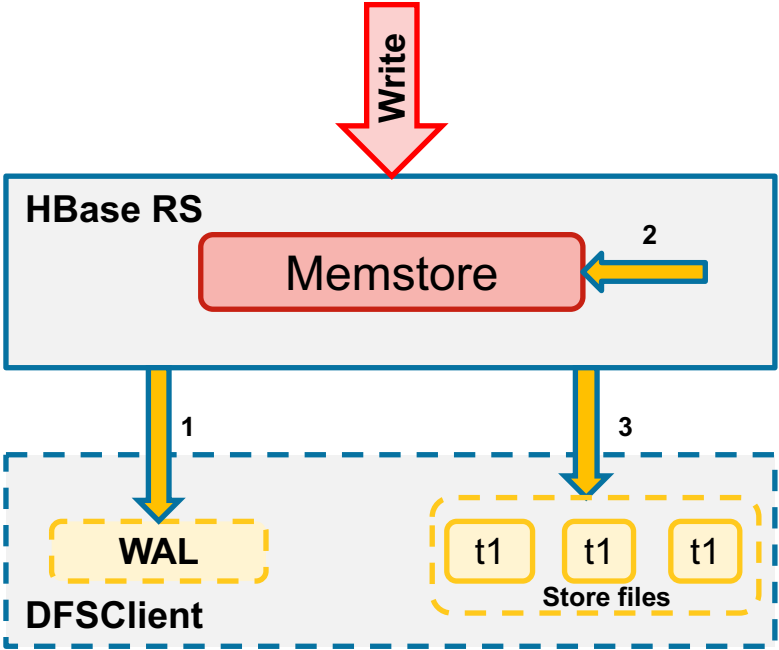- **D**urability is guaranteed for all successful mutations

**Bloomberg**

Engineering

# Namespace

**Bloomberg**

Engineering

# HBase Write

Write

**HBase Region Server**

1

**WAL**

**File System**

**Bloomberg**

Engineering

# HBase Write

HBase RS

Write

Memstore

2

1

WAL

File System

Bloomberg

Engineering

# HBase Write

Bloomberg
Engineering

# HBase Write

# HBase Write

Bloomberg
Engineering

# Compaction

# Compaction

**Bloomberg**

Engineering

# HBase Read

Read

**HBase RS**

Memstore

Block Cache | Block

1

**File System**

WAL

t1 | t1 | t1

Store files

# HBase Read

Bloomberg
Engineering

# HDFS Read

**Bloomberg**

Engineering

# HDFS Short-Circuit Read

**Bloomberg**

Engineering

# Large Read Cache

Bloomberg

Engineering

# Large Read Cache

Bloomberg

Engineering

# HBase Complete

Bloomberg

Engineering

# HBase Complete

**Bloomberg**

Engineering

# HBase Complete

**Bloomberg**

Engineering

# HBase Complete

# HBase Complete

**Bloomberg**
Engineering

# Region Server Failure

Bloomberg
Engineering

# Region Server Failure

**Bloomberg**

Engineering

# Region Server Failure

# Region Replication

# Region Replication

**Bloomberg**

Engineering

# Region Replication



https://www.youtube.com/watch?v=I6S-Vbs9WsU

Bloomberg

Engineering

# Load Balancing

**Bloomberg**

Engineering

# Load Balancing

# Balancer

- Region Count Cost
- Primary Region Count Cost
- Table Skew Cost
- Locality Cost
- Rack Locality Cost
- Region Replica Host Cost
- Region Replica Rack Cost
- Read Request Cost
- Write Request Cost
- Memstore Size Cost
- Storefile Size Cost
- Move Cost

**Bloomberg**

Engineering

# Other Features

- HBase Replication

- HBase multi-tenancy support
  - https://www.youtube.com/watch?v=bZjz2G38Ju0

- HBase Co-processors and Filters
  - https://www.slideshare.net/Hadoop_Summit/hbase-coprocessors-uses-abuses-solutions

**Bloomberg**

Engineering

# Operator View

# ZooKeeper Availability

- ZK Quorum

- One leader and remaining followers
  - — stat
  - — ruok
  - — mntr

- Test for availability
  - — e.g., List children of a znode

# HBase Availability

- Master Availability -  http://hmaster-node:16010/jmx
  — "name" : "Hadoop:service=HBase,name=Master,sub=Server","tag.isActiveMaster" : "true"

- Dead RegionServers
  —  "name" : "Hadoop:service=HBase,name=Master,sub=Server", "numDeadRegionServers" : 0

- Region In Transition
  — "name" : "Hadoop:service=HBase,name=Master,sub=AssignmentManger","ritCount" : 0

- Test for availability
  — e.g., Query system table by listing tables

**TechAtBloomberg.com**

**Bloomberg**

Engineering

# HDFS Availability

- Namenode Availability - [http://namenode-host:50070/jmx](http://namenode-host:50070/jmx)
  - "name" : "Hadoop:service=NameNode,name=FSNamesystem", "tag.HAState" : "active"
- Dead Datanodes
  - "name" : "Hadoop:service=NameNode,name=NameNodeInfo", "DeadNodes" : "{}"
- Missing Blocks
  - "name" : "Hadoop:service=NameNode,name=NameNodeInfo", "NumberOfMissingBlocks" : 0
- Percentage Used
  - "name" : "Hadoop:service=NameNode,name=NameNodeInfo", "PercentUsed" : 59
- Under replicated blocks
  - "name" : "Hadoop:service=NameNode,name=FSNamesystemState","UnderReplicatedBlocks":0
- Test for availability
  - e.g., Append data to a test file

**Bloomberg**

Engineering

# HBase Performance

- RegionServer JMX metrics - http://rs-node:60300/jmx
  - "name":"Hadoop:service=HBase,name=RegionServer,sub=Server"
    - Blockcache hit ratio
    - Request counts
    - Request response time
    - Compaction related metrics
    - Region count
    - Flush related metrics
    - Percentage of files local
    - Split related metrics
  - "name" : "Hadoop:service=HBase,name=RegionServer,sub=Tables",
    - Table level metrics

https://www.slideshare.net/MichaelStack4/hbaseconasia2018-track31-serving-billions-of-queries-in-millisecond-latencies

**TechAtBloomberg.com**

**Bloomberg**

Engineering

# JVM

- GC – JMX Metrics
  - "name" : "java.lang:type=GarbageCollector,name=ParNew",
  - "name" : "java.lang:type=GarbageCollector,name=ConcurrentMarkSweep",
- GC Logging
  - -verbose:gc
  - -XX:+PrintHeapAtGC
  - -XX:+PrintGCDetails
  - -XX:+PrintGCTimeStamps
  - -XX:+PrintGCDateStamps
  - -XX:+PrintGCApplicationStoppedTime
  - -XX:+PrintClassHistogram
  - -XX:+PrintGCApplicationConcurrentTime
  - -XX:+PrintTenuringDistribution
  - -Xloggc:

**TechAtBloomberg.com**

**Bloomberg**

Engineering

# OS/HW

- Memory

- CPU

- Disk

- Networking

**Bloomberg**

Engineering

# Logs

- ZooKeeper Log

- HDFS
  — Namenode log
  — Datanode log

- HBase
  — Master log
  — RegionServer log

- OS
  — Syslog

**Bloomberg**

Engineering

# Interacting with HBase

- HBase shell
  — DDL: create namespace/table, alter
  — Security: grant, revoke
  — DML: get, put, scan
  — Tools: assign, compact, balance
  — General: status

- HBase admin API

- HBase client API

**Bloomberg**

Engineering

# Data Backup / Restore

- Snapshot
  — hbase shell > snapshot 'table', 'table_mmddyy'

- Restore from snapshot
  — hbase shell > restore_snapshot 'table_mmddyy'

- Export Snapshot
  — $ hbase org.apache.hadoop.hbase.snapshot.ExportSnapshot

- CopyTable
  — hbase org.apache.hadoop.hbase.mapreduce.CopyTable

Bloomberg
Engineering

# Thank You!

**Acknowledgement: Apache HBase Community**
**Reference: http://hbase.apache.org**
**Connect with Hadoop Team: hadoop@bloomberg.net**

**Bloomberg**
**Engineering**

**TechAtBloomberg.com**

# We are hiring!

https://www.bloomberg.com/careers

## Questions?

Engineering Bloomberg

**TechAtBloomberg.com**