UNIVERSITY OF
WATERLOO

Data-Intensive Distributed Computing
CS 431/631 451/651 (Fall 2021)

Part 9a: Mutable State (1/2)

Ali Abedi

# Mutable state

From sequential reads and append only writes
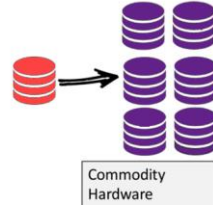
To random reads and writes

# Why not RDBMS?

- Does not scale out → expensive



- Does not support semi-structured data

# NoSQL
## (Not only SQL)

1. Horizontal scaling
2. Replicate/distribute data over many servers
3. Simple call interface
4. Weaker concurrency model than ACID
5. Flexible schemas

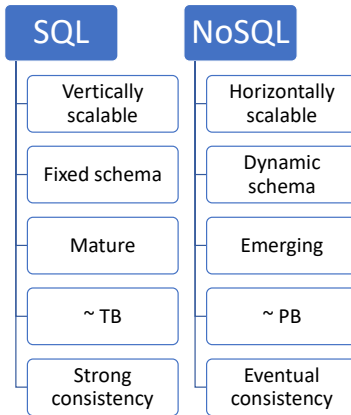But, don't blindly follow the hype…
Often, MySQL is what you really need!



HOW TO WRITE A CV

Leverage the NoSQL boom

Source: Cattell (2010). Scalable SQL and NoSQL Data Stores. *SIGMOD Record.*

# SQL vs. NoSQL

| SQL | NoSQL |
|---|---|
| Vertically scalable | Horizontally scalable |
| Fixed schema | Dynamic schema |
| Mature | Emerging |
| ~ TB | ~ PB |
| Strong consistency | Eventual consistency |

7

# (Major) Types of NoSQL databases



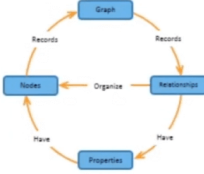| Key Value | Document-Based | Column-Based | Graph-Based |
|---|---|---|---|
| **Example:** Riak, Tokyo Cabinet, Redis server, Memcached, Scalaris | **Example:** MongoDB, CouchDB, OrientDB, RavenDB | **Example:** BigTable, Cassandra, Hbase, Hypertable | **Example:** Neo4J, InfoGrid, Infinite Graph, Flock DB |

# KEY-VALUE STORE

| Key | Value |
|-----|-------|
| K1 | AAA,BBB,CCC |
| K2 | AAA,BBB |
| K3 | AAA,DDD |
| K4 | AAA,2,01/01/2015 |
| K5 | 3,ZZZ,5623 |

10

# Three Core Ideas

Keeping track of the partitions?

Partitioning (sharding)
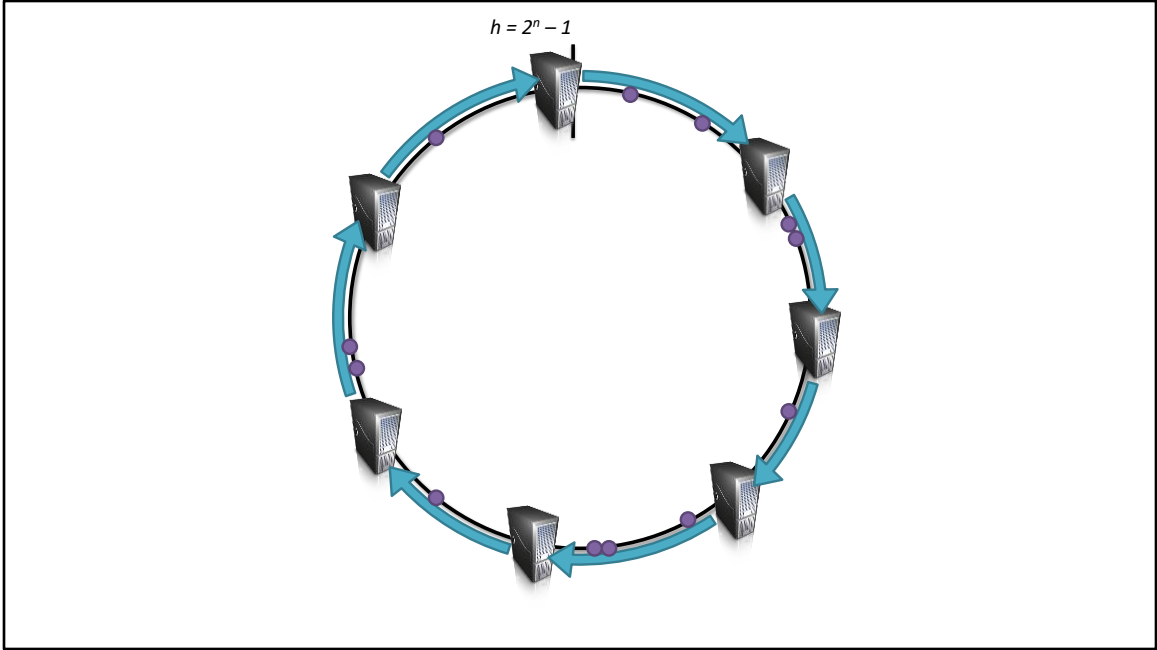To increase scalability and to decrease latency

Consistency?

Replication
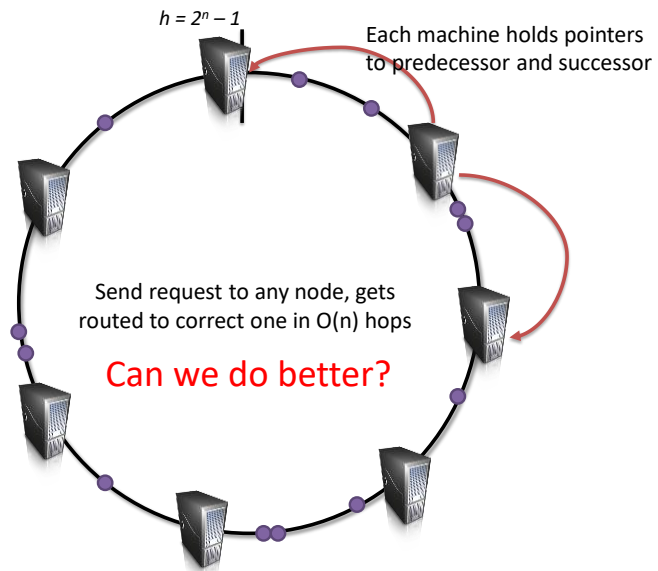To increase robustness (availability) and to increase throughput
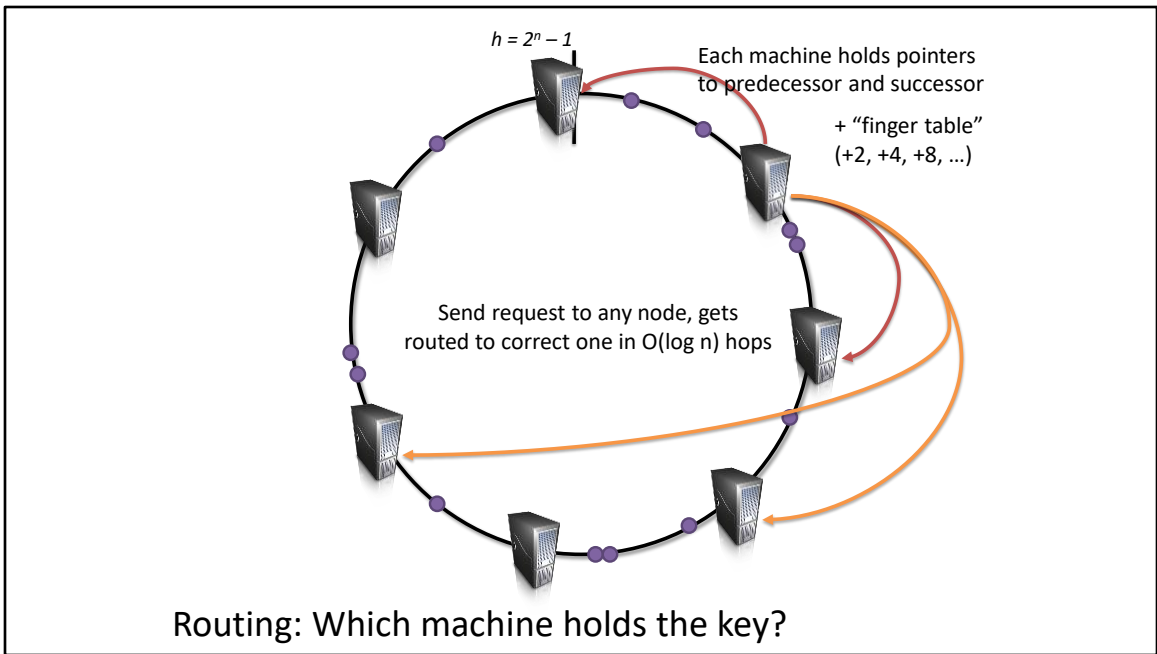
Caching
To reduce latency

# Chord distributed protocol
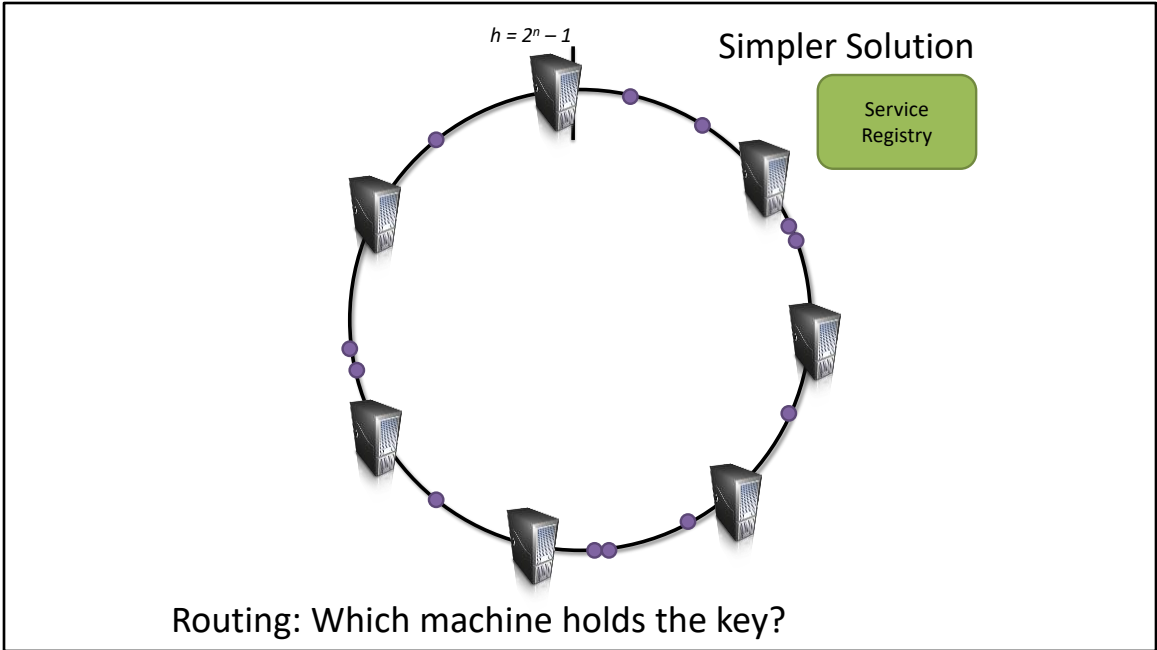
Stoica et al. (2001). Chord: A Scalable Peer-to-peer Lookup
Service for Internet Applications. *SIGCOMM.*
*And other resources …*

$h = 2^n - 1$

$h = 2^n - 1$

Each machine holds pointers to predecessor and successor

Send request to any node, gets routed to correct one in O(n) hops

Can we do better?

Routing: Which machine holds the key?

Routing: Which machine holds the key?

15

$h = 2^n - 1$

Simpler Solution

Service Registry

Routing: Which machine holds the key?

$h = 2^n - 1$

How do we rebuild the predecessor, successor, finger tables?

New machine joins: What happens?

$h = 2^n - 1$

Solution: Replication

Covered!

Covered!

Machine fails: What happens?

18

Bigtable

20

# Bigtable Applications

Gmail

Google's web crawl

Google Earth

Google Analytics

Data source and data sink for MapReduce

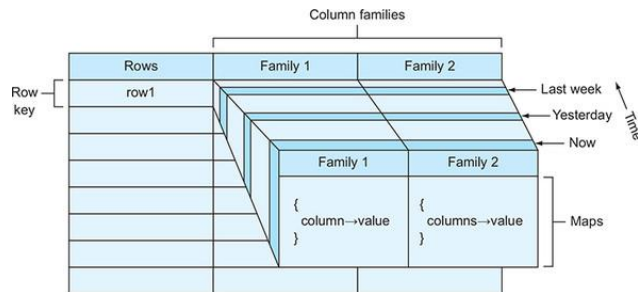HBase is the open-source implementation…

21

# Data Model

A table in Bigtable is a sparse, distributed, persistent  multidimensional sorted map

Map indexed by a row key, column key, and a timestamp
(row:string, column:string, time:int64) → uninterpreted byte array

Supports lookups, inserts, deletes
Single row transactions only

# Rows and Columns

Rows maintained in sorted lexicographic order
Applications can exploit this property for efficient row scans
Row ranges dynamically partitioned into tablets
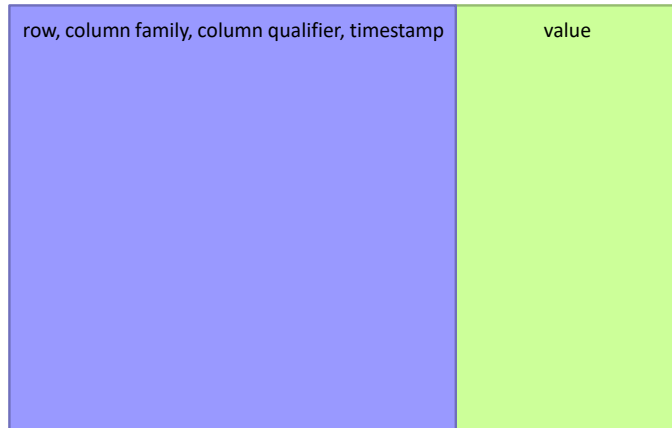
Columns grouped into column families
Column key = family:qualifier
Column families provide locality hints
Unbounded number of columns

At the end of the day, it's all key-value pairs!

23

# Key-Values

| row, column family, column qualifier, timestamp | value |
|---|---|

# Okay, so how do we build it?

| In Memory | On Disk |
|---|---|
| Mutability Easy | Mutability Hard |
| Small | Big |

25

# Log Structured Merge Trees

Writes → **MemStore** → Reads

What happens when we run out of memory?

# Log Structured Merge Trees

Writes ──────▶ **MemStore** ──────────────▶ Reads

Memory
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
Disk     *Flush to disk*
⬇

**Store**

Immutable, indexed, persistent, key-value pairs

What happens to the read path?

27

# Log Structured Merge Trees

Writes ——→ **MemStore** ——→ *Merge* ○ ——→ Reads

Memory

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Disk    *Flush to disk* ↓

**Store**

Immutable, indexed, persistent, key-value pairs
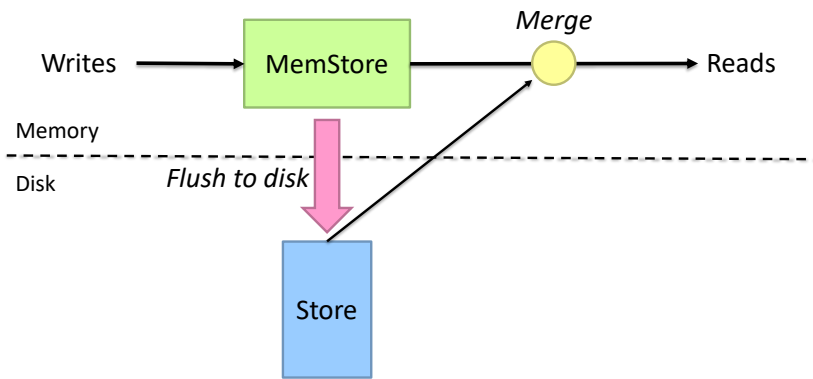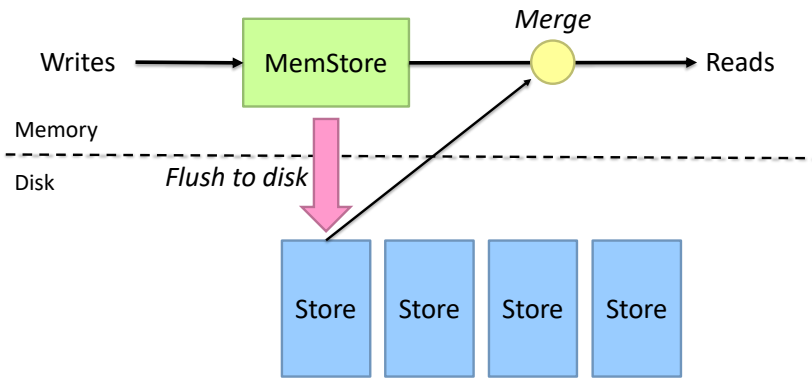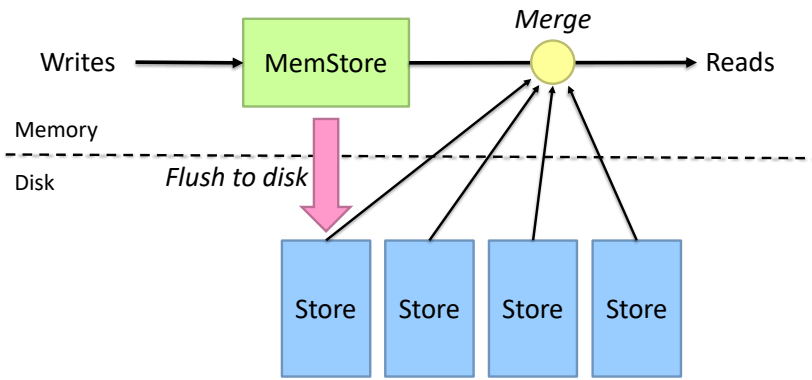
What happens as more writes happen?

**28**

# Log Structured Merge Trees



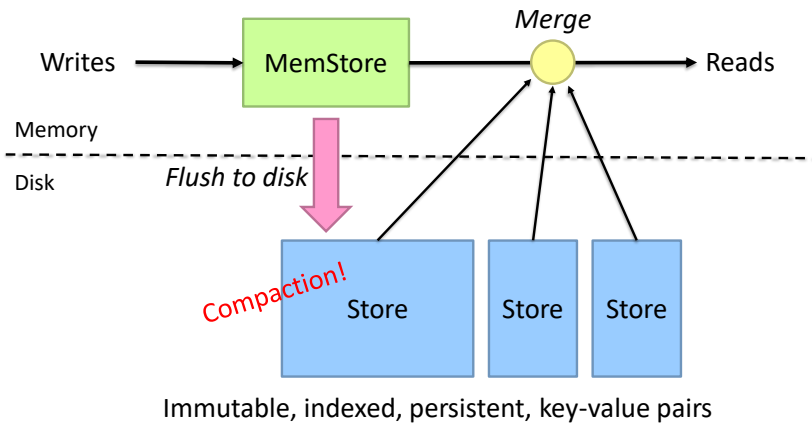Immutable, indexed, persistent, key-value pairs

What happens to the read path?

**29**

# Log Structured Merge Trees

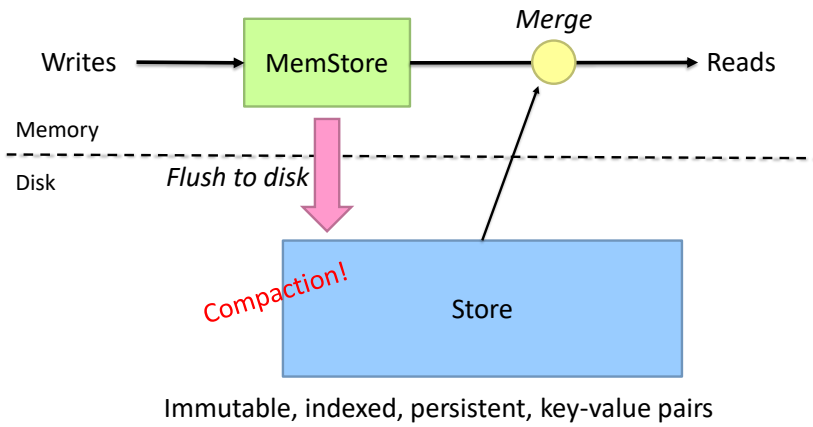Writes → **MemStore** → *Merge* (◯) → Reads

Memory

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Disk

*Flush to disk*

| Store | Store | Store | Store |

Immutable, indexed, persistent, key-value pairs

What's the next issue?

30

# Log Structured Merge Trees

Writes → MemStore → *Merge* → Reads

Memory

*Flush to disk*

Disk

Compaction!

Store   Store   Store

Immutable, indexed, persistent, key-value pairs

31

# Log Structured Merge Trees



32

# Log Structured Merge Trees



Immutable, indexed, persistent, key-value pairs

One final component…

33

# Log Structured Merge Trees

The complete picture…



*Merge*

Writes ⟶ ○ ⟶ MemStore ⟶ ○ ⟶ Reads

Memory
- - - - - - - - - - - - - - - - - - - - - - - - -
Disk

*Flush to disk*

*Logging for persistence* — WAL | Store | Store | Store

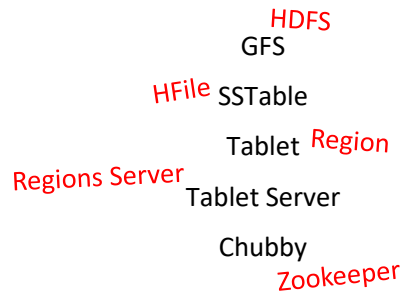Immutable, indexed, persistent, key-value pairs

*Compaction!*

**34**

# Log Structured Merge Trees
The complete picture…

Okay, now how do we build a distributed version?

# HBase Bigtable building blocks

HDFS
GFS

HFile SSTable

Tablet Region

Regions Server
Tablet Server

Chubby
Zookeeper

36

# SSTable *HFile*

Persistent, ordered immutable map from keys to values
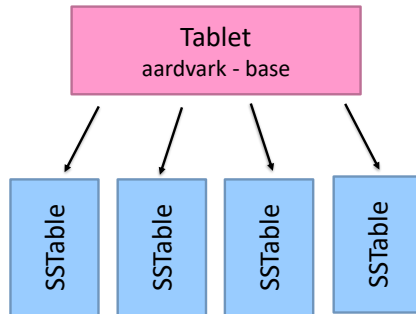Stored in GFS: replication "for free"

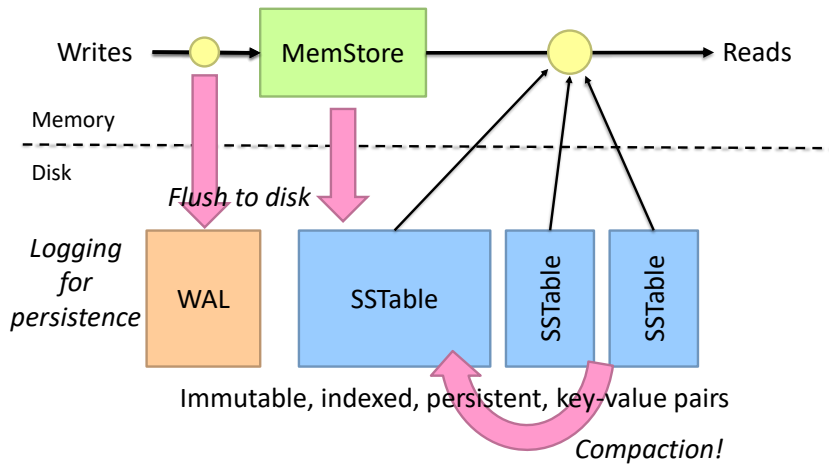Supported operations:
Look up value associated with key
Iterate key/value pairs within a key range

# Region Tablet

Dynamically partitioned range of rows
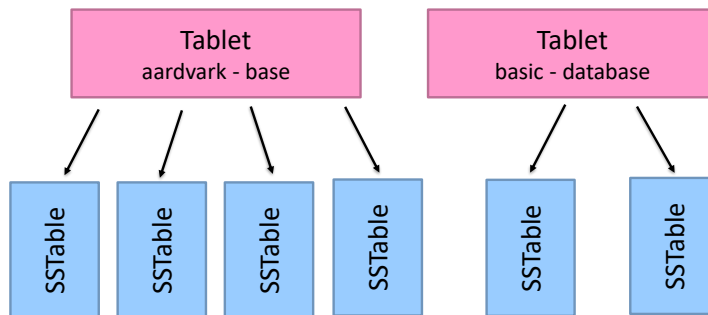Comprised of multiple SSTables

**Region Server**
# Tablet Server

Writes → MemStore → Reads

Memory
- - - - - - - - - - - - - - - - - - - - - - - - -
Disk

*Flush to disk*

*Logging for persistence*

WAL    SSTable    SSTable    SSTable

Immutable, indexed, persistent, key-value pairs

*Compaction!*

39

Region          Region Server
# Tablet to Tablet Server Assignment

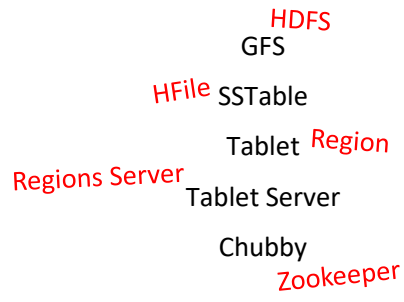Each tablet is assigned to one tablet server at a time
Exclusively handles read and write requests to that tablet

What happens when a tablet grow too big?
What happens when a tablet server fails?

We need a lock service!

41

# HBase Bigtable building blocks

HDFS
GFS

HFile SSTable

Tablet Region

Regions Server
Tablet Server

Chubby
Zookeeper

# Architecture

Client library

Bigtable master <span style="color:red">HMaster</span>

Tablet servers
<span style="color:red">Regions Servers</span>

**43**

# Bigtable Master

Roles and responsibilities:
Assigns tablets to tablet servers
Detects addition and removal of tablet servers
Balances tablet server load
Handles garbage collection
Handles schema changes

Tablet structure changes:
Table creation/deletion (master initiated)
Tablet merging (master initiated)
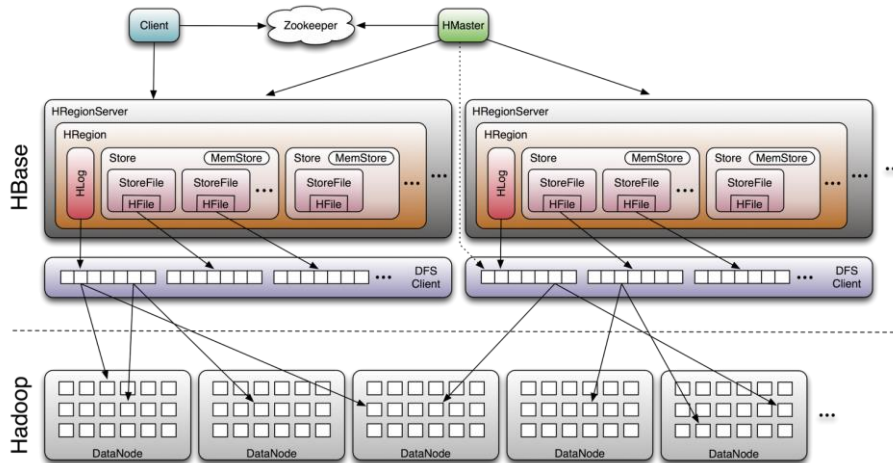Tablet splitting (tablet server initiated)

# HBase

**45**