

# Data-Intensive Distributed Computing

CS 431/631 451/651 (Winter 2019)

Part 7: Mutable State (1/2)

March 14, 2019

Adam Roegiest

Kira Systems

These slides are available at <http://roegiest.com/bigdata-2019w/>

This work is licensed under a Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States  
See <http://creativecommons.org/licenses/by-nc-sa/3.0/us/> for details



# Structure of the Course

Analyzing Text

Analyzing Graphs

Analyzing  
Relational Data

Data Mining

“Core” framework features  
and algorithm design

# The Fundamental Problem

We want to keep track of *mutable* state in a *scalable* manner

Assumptions:

State organized in terms of logical records

State unlikely to fit on single machine, must be distributed

MapReduce won't do!

Want more? Take a *real* distributed systems course!

# The Fundamental Problem

We want to keep track of *mutable* state in a *scalable* manner

Assumptions:

State organized in terms of logical records

State unlikely to fit on single machine, must be distributed

Uh... just use an RDBMS?

# What do RDBMSes provide?

Relational model with schemas

Powerful, flexible query language

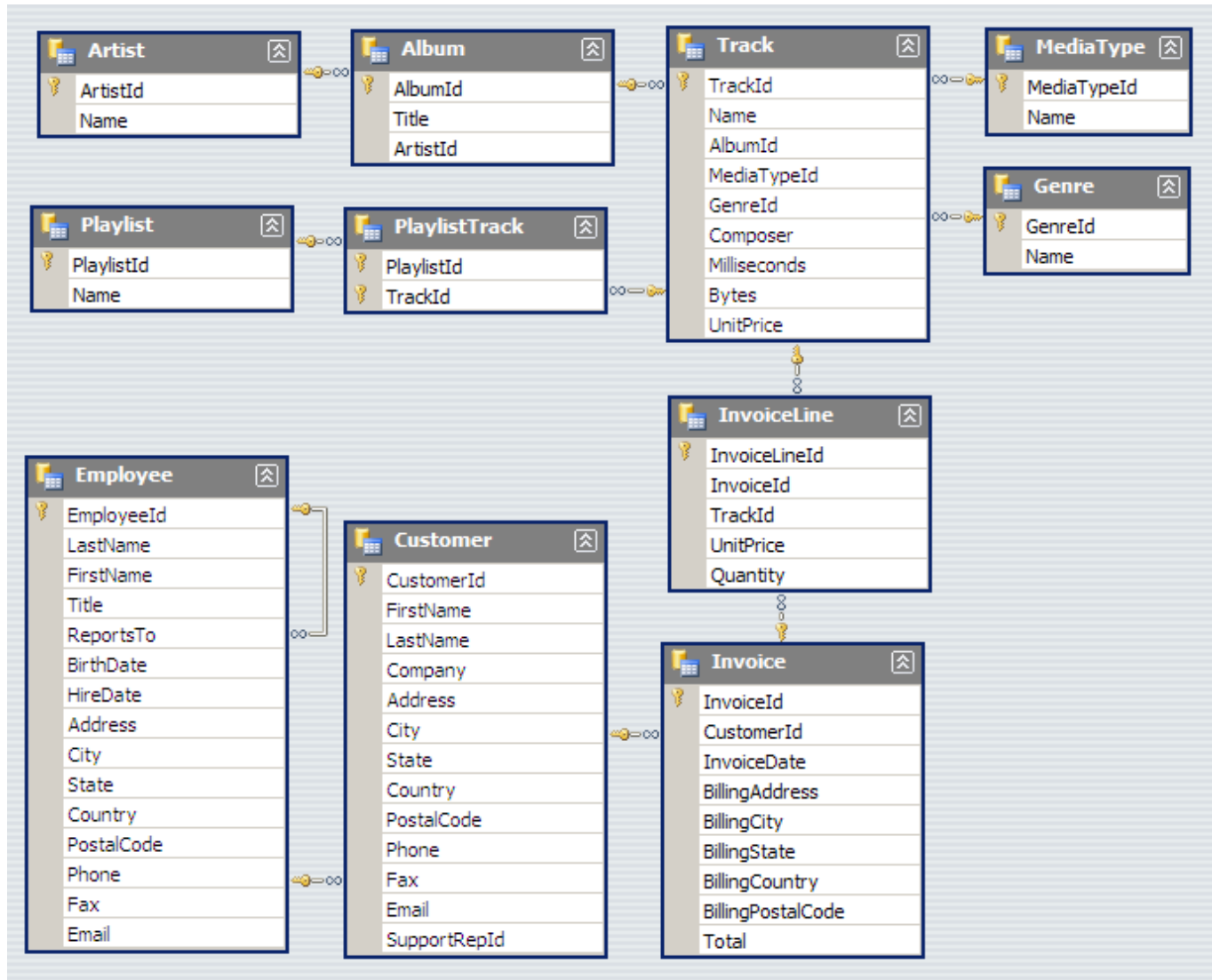
Transactional semantics: ACID

Rich ecosystem, lots of tool support

# RDBMSes: Pain Points



# #1: Must design up front, painful to evolve



Consistent field names?

This should really be a list...

```
{  
  "token": 945842,  
  "feature_enabled": "super_special",  
  "userid": 229922,  
  "page": "null",  
  "info": { "email": "my@place.com" }  
}
```

Is this really an integer?

Is this really null?

What keys? What values?

## JSON to the Rescue!

Flexible design doesn't mean *no* design!



#2: Pay for ACID!



# #3: Cost!



# What do RDBMSes provide?

Relational model with  
Powerful, flexible queries  
Transactional semantics  
Rich ecosystem, lots of



What if we want *a la carte*?

# Features *a la carte*?

What if I'm willing to give up consistency for scalability?

What if I'm willing to give up the relational model for flexibility?

What if I just want a cheaper solution?

Enter... NoSQL!

# HOW TO WRITE A CV



Leverage the NoSQL boom

# NoSQL

## (Not only SQL)

1. Horizontally scale “simple operations”
2. Replicate/distribute data over many servers
3. Simple call interface
4. Weaker concurrency model than ACID
5. Efficient use of distributed indexes and RAM
6. Flexible schemas

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

“web scale”



# (Major) Types of NoSQL databases

Key-value stores

Column-oriented databases

Document stores

Graph databases



# 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

*Consistency?*

Caching

To reduce latency

# Key-Value Stores



# Key-Value Stores: Data Model

Stores associations between keys and values

Keys are usually primitives

For example, ints, strings, raw bytes, etc.

Values can be primitive or complex: often opaque to store

Primitives: ints, strings, etc.

Complex: JSON, HTML fragments, etc.

# Key-Value Stores: Operations

Very simple API:

Get – fetch value associated with key

Put – set value associated with key

Optional operations:

Multi-get

Multi-put

Range queries

Secondary index lookups

Consistency model:

Atomic single-record operations (usually)

Cross-key operations: who knows?

# Key-Value Stores: Implementation

Non-persistent:

Just a big in-memory hash table

Examples: Redis, memcached

Persistent

Wrapper around a traditional RDBMS

Examples: Voldemort

What if data doesn't fit on a single machine?

# Simple Solution: Partition!

Partition the key space across multiple machines

Let's say, hash partitioning

For  $n$  machines, store key  $k$  at machine  $h(k) \bmod n$

Okay... But:

How do we know which physical machine to contact?

How do we add a new machine to the cluster?

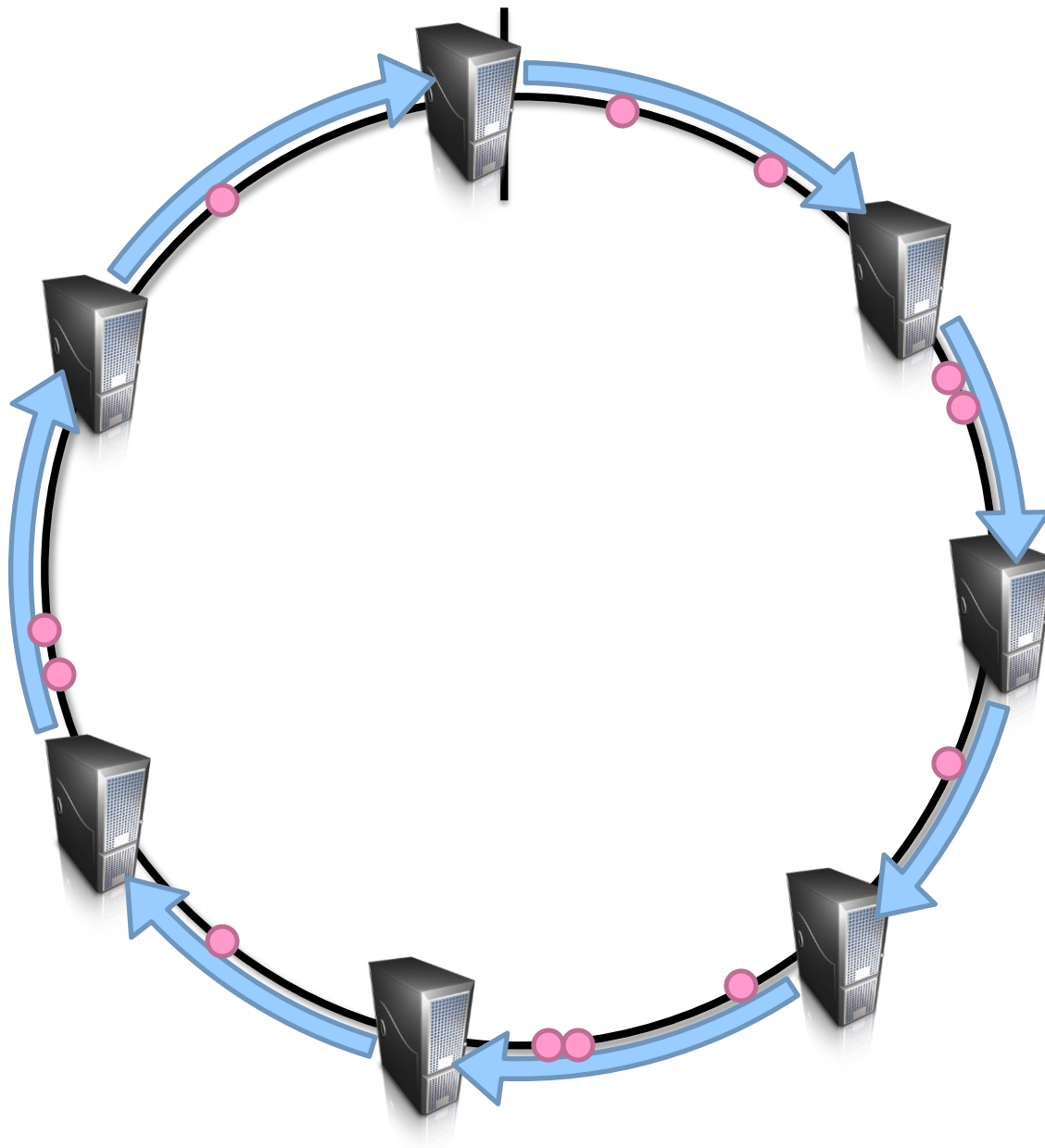
What happens if a machine fails?

# Clever Solution

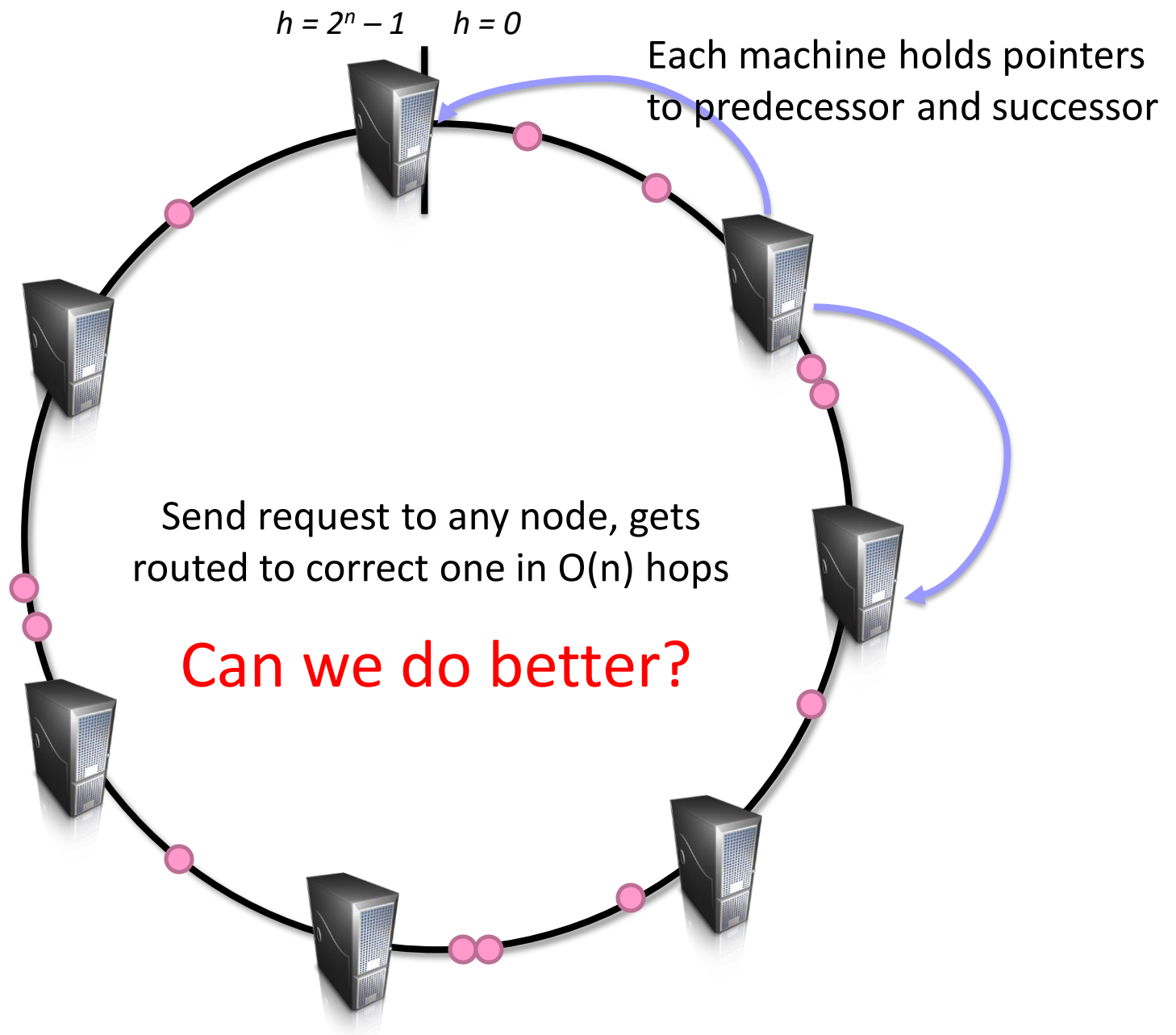
Hash the keys  
Hash the machines also!

**Distributed hash tables!**  
(following combines ideas from several sources...)

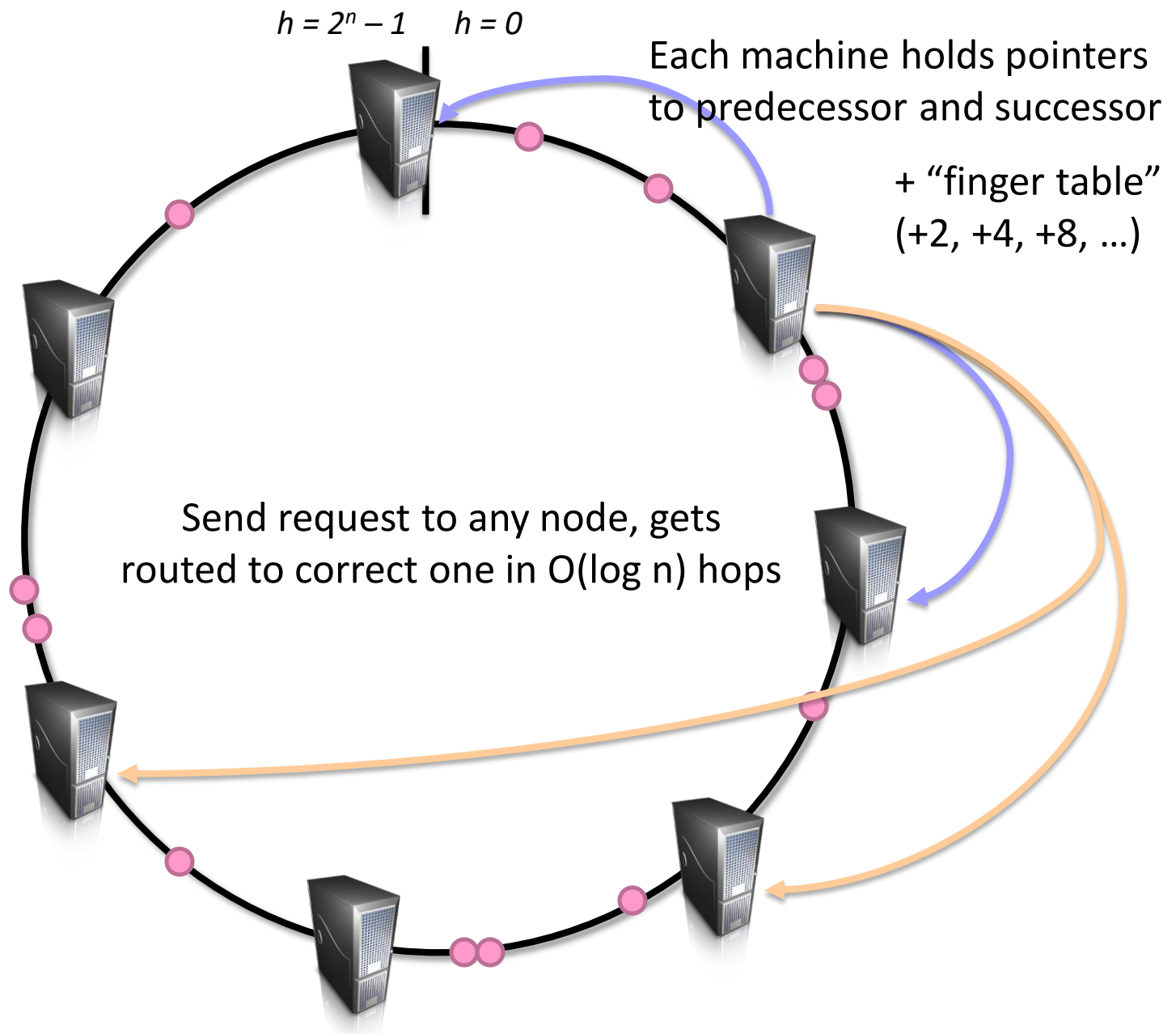
$h = 2^n - 1$     $h = 0$





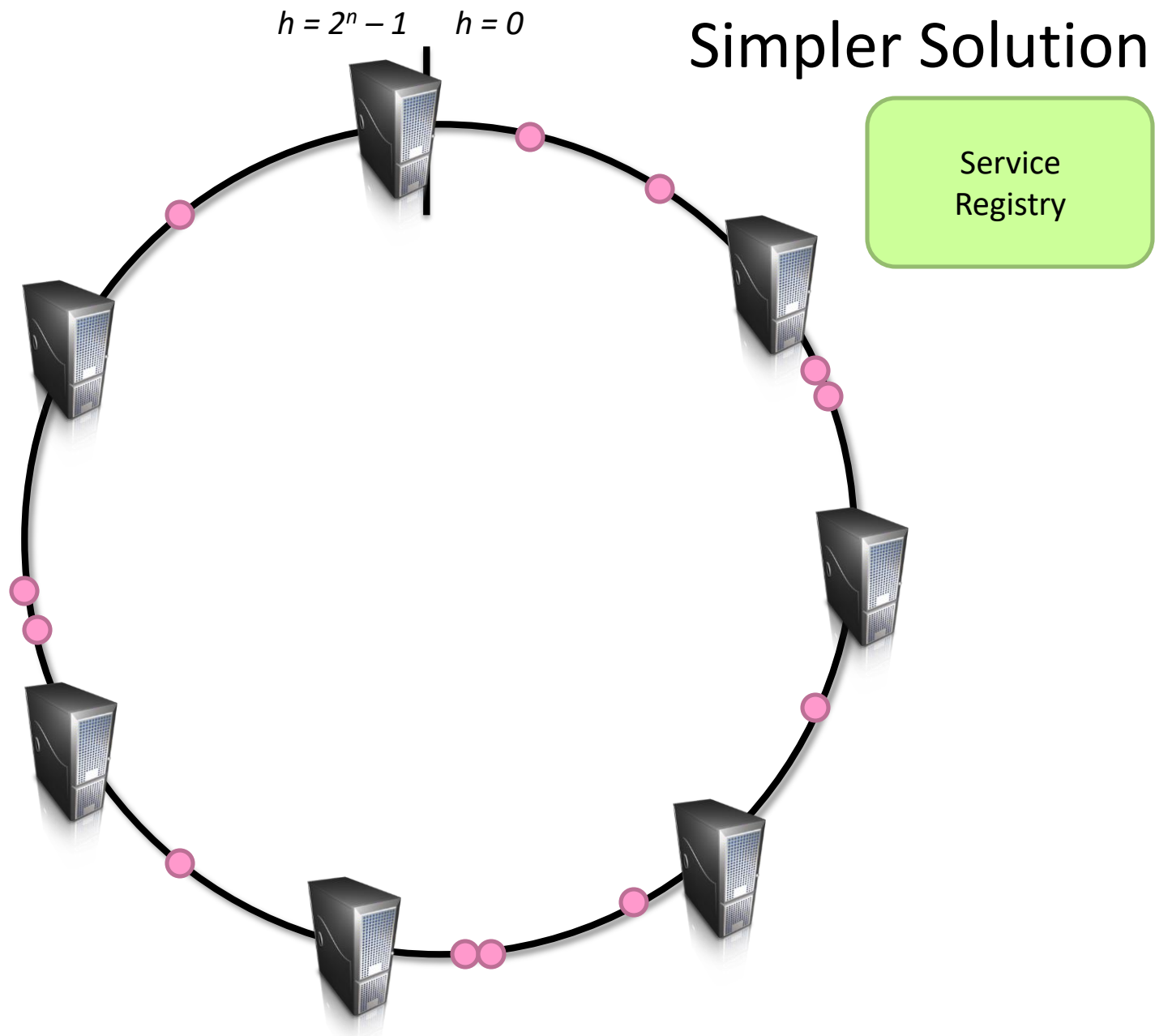


Routing: Which machine holds the key?



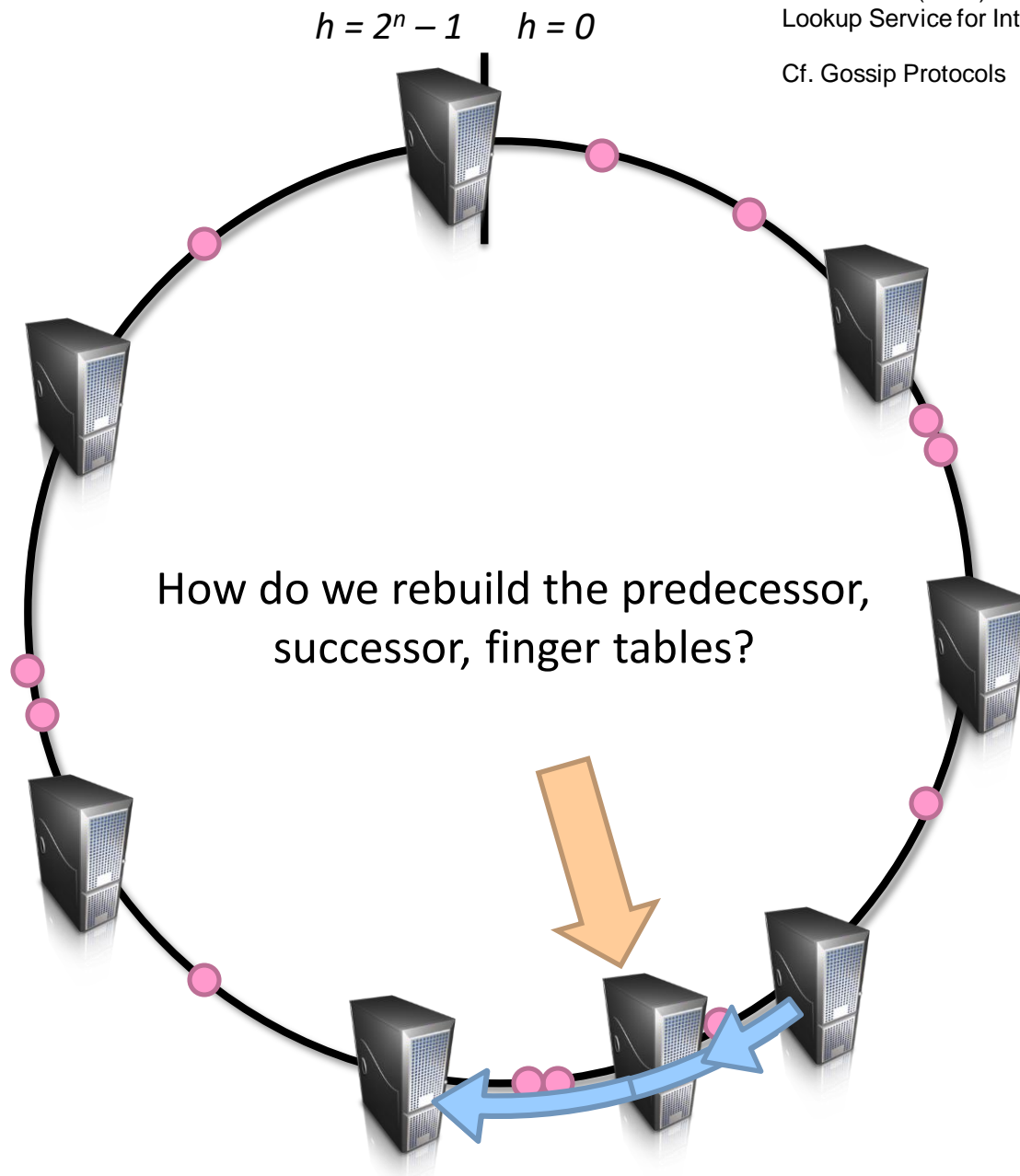
Routing: Which machine holds the key?

# Simpler Solution



Routing: Which machine holds the key?

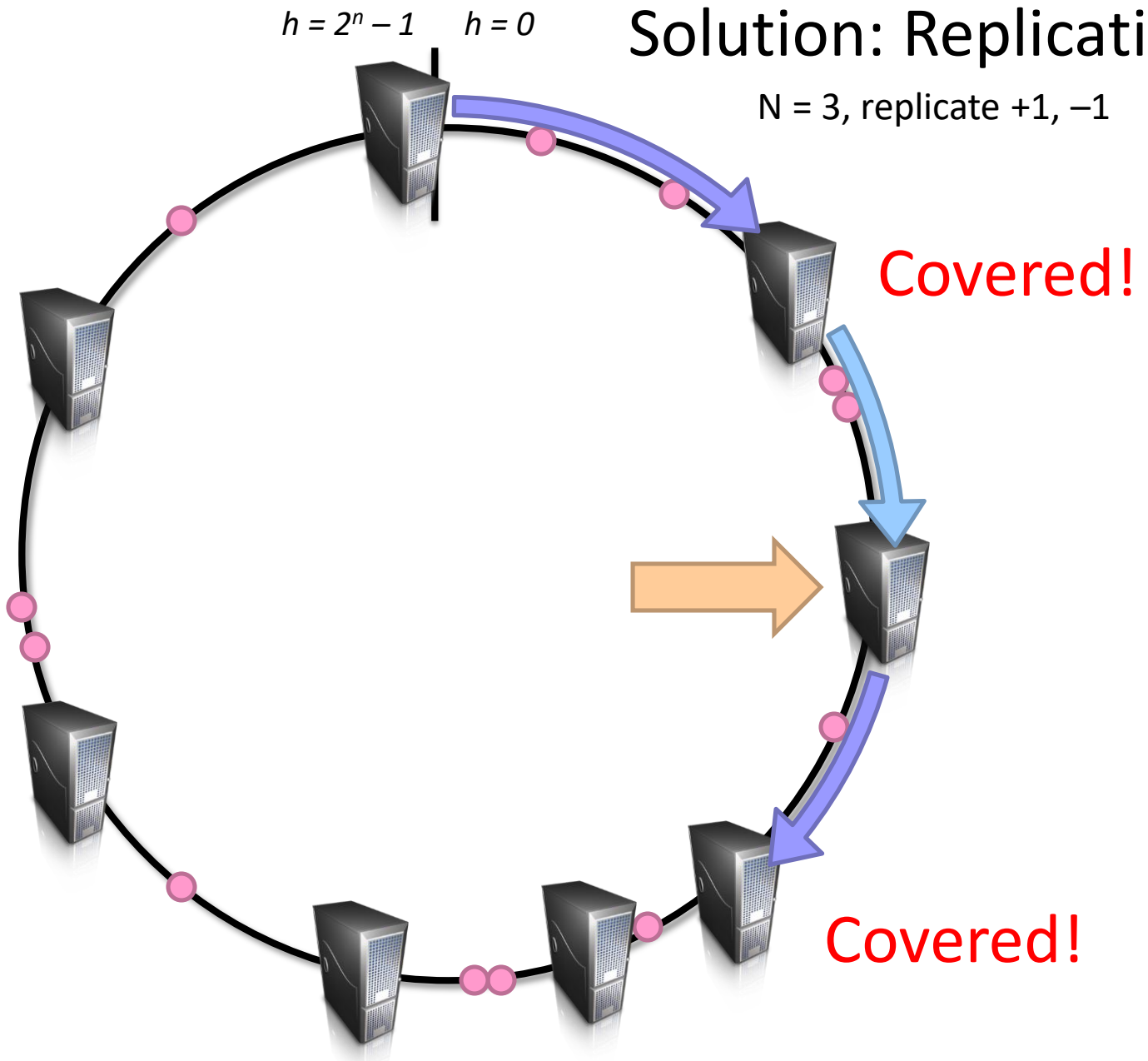
Cf. Gossip Protocols



New machine joins: What happens?

# Solution: Replication

$N = 3$ , replicate  $+1, -1$



Machine fails: What happens?

# 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

# Another Refinement: Virtual Nodes

Don't directly hash servers

Create a large number of virtual nodes, map to physical servers

Better load redistribution in event of machine failure

When new server joins, evenly shed load from other servers

# Bigtable





# Bigtable Applications

Gmail

Google's web crawl

Google Earth

Google Analytics

Data source and data sink for MapReduce

HBase is the open-source implementation...

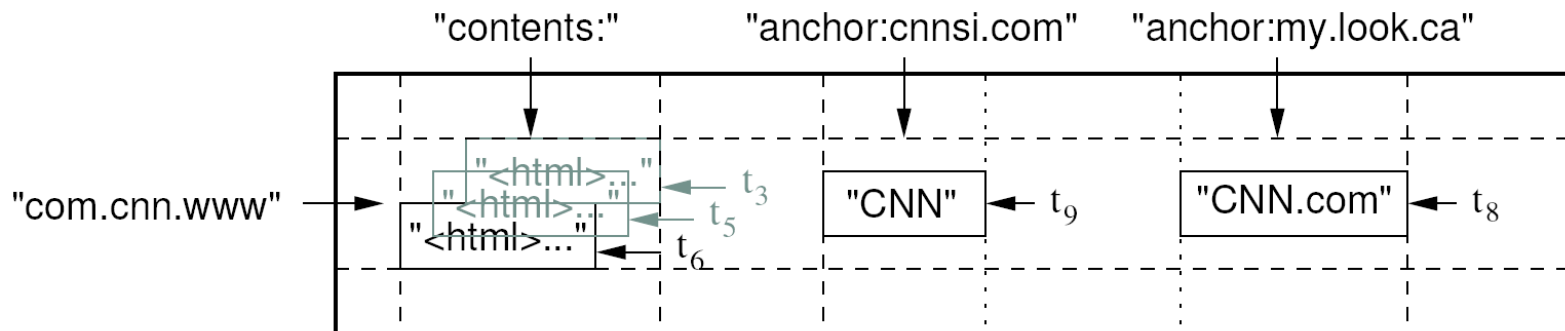
# Data Model

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



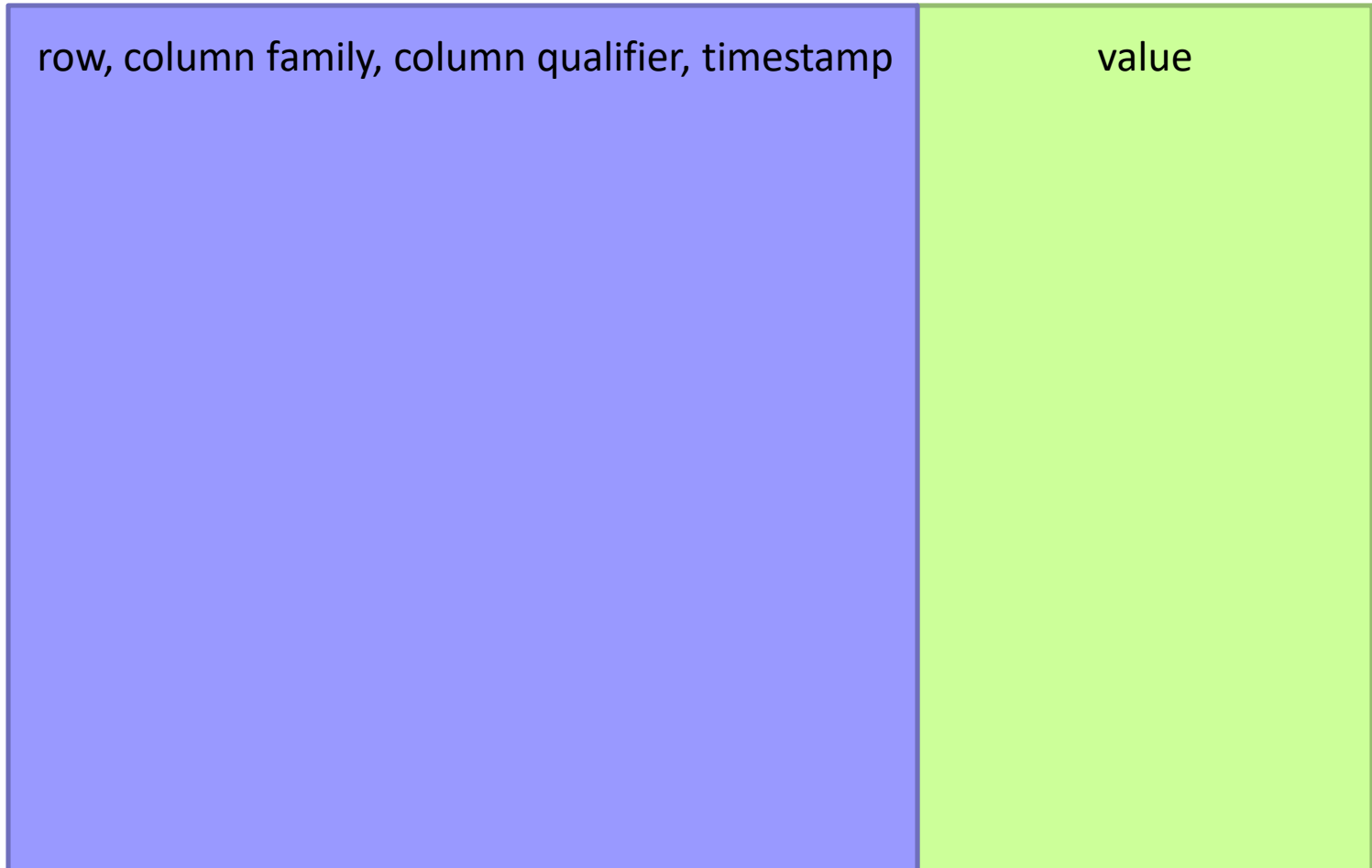
# Rows and Columns

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!

# Key-Values



Okay, so how do we build it?

In Memory

On Disk

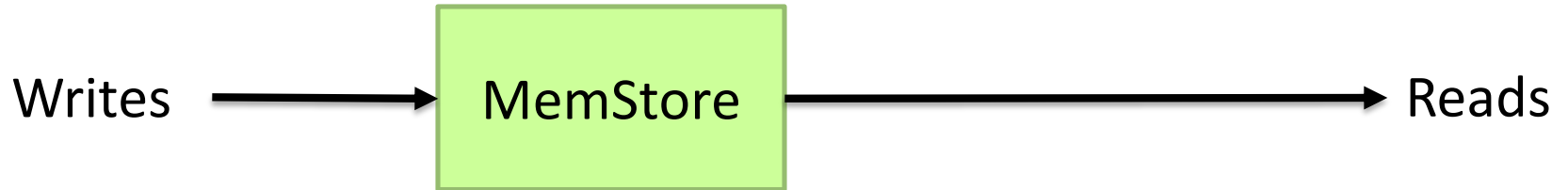
Mutability Easy

Mutability Hard

Small

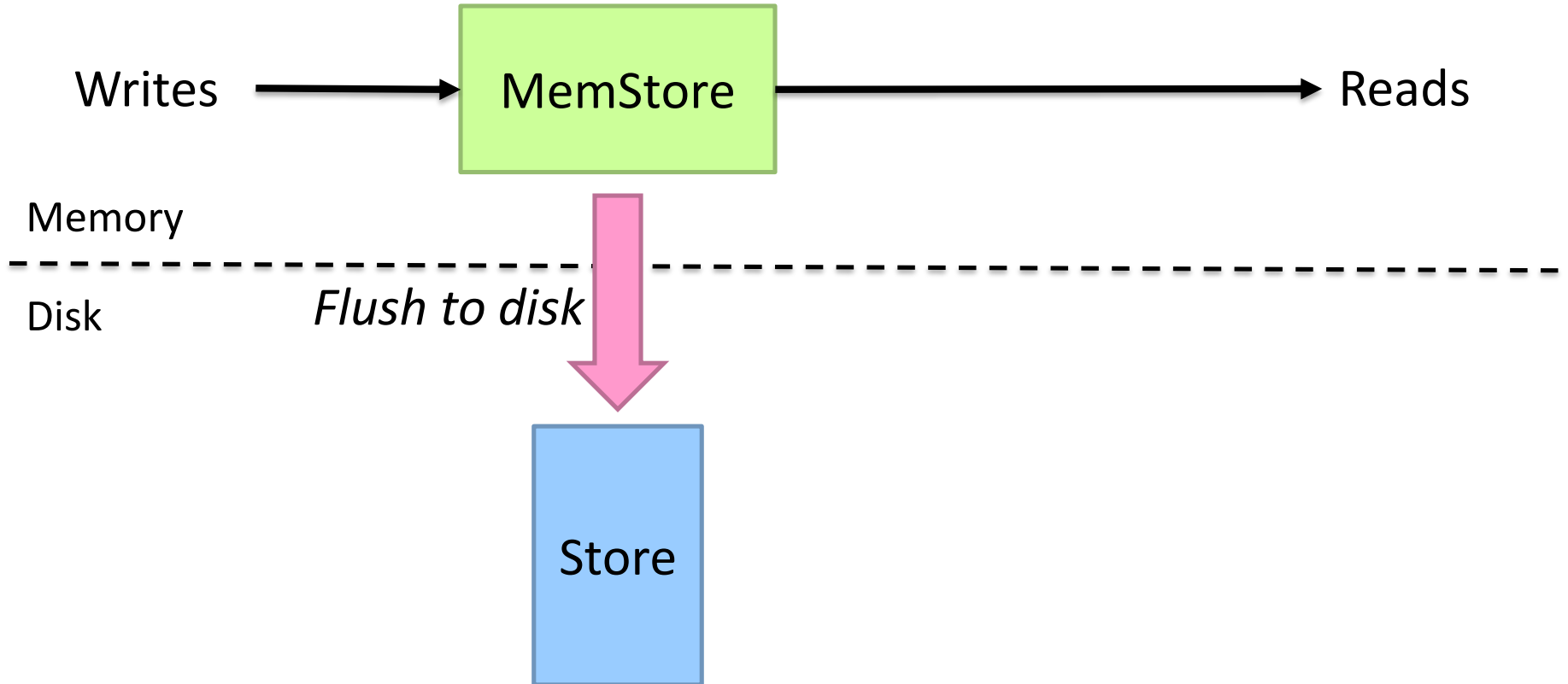
Big

# Log Structured Merge Trees



What happens when we run out of memory?

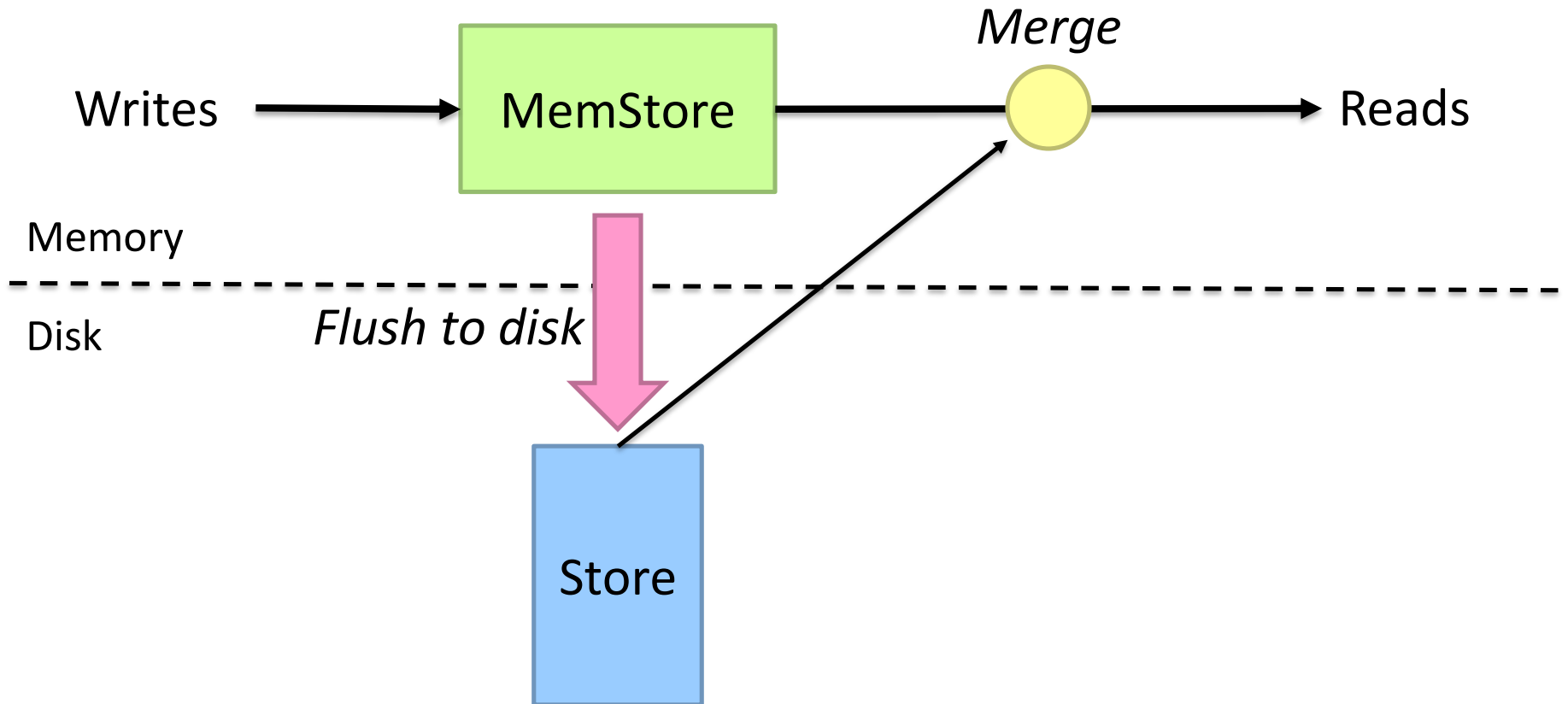
# Log Structured Merge Trees



Immutable, indexed, persistent, key-value pairs

What happens to the read path?

# Log Structured Merge Trees

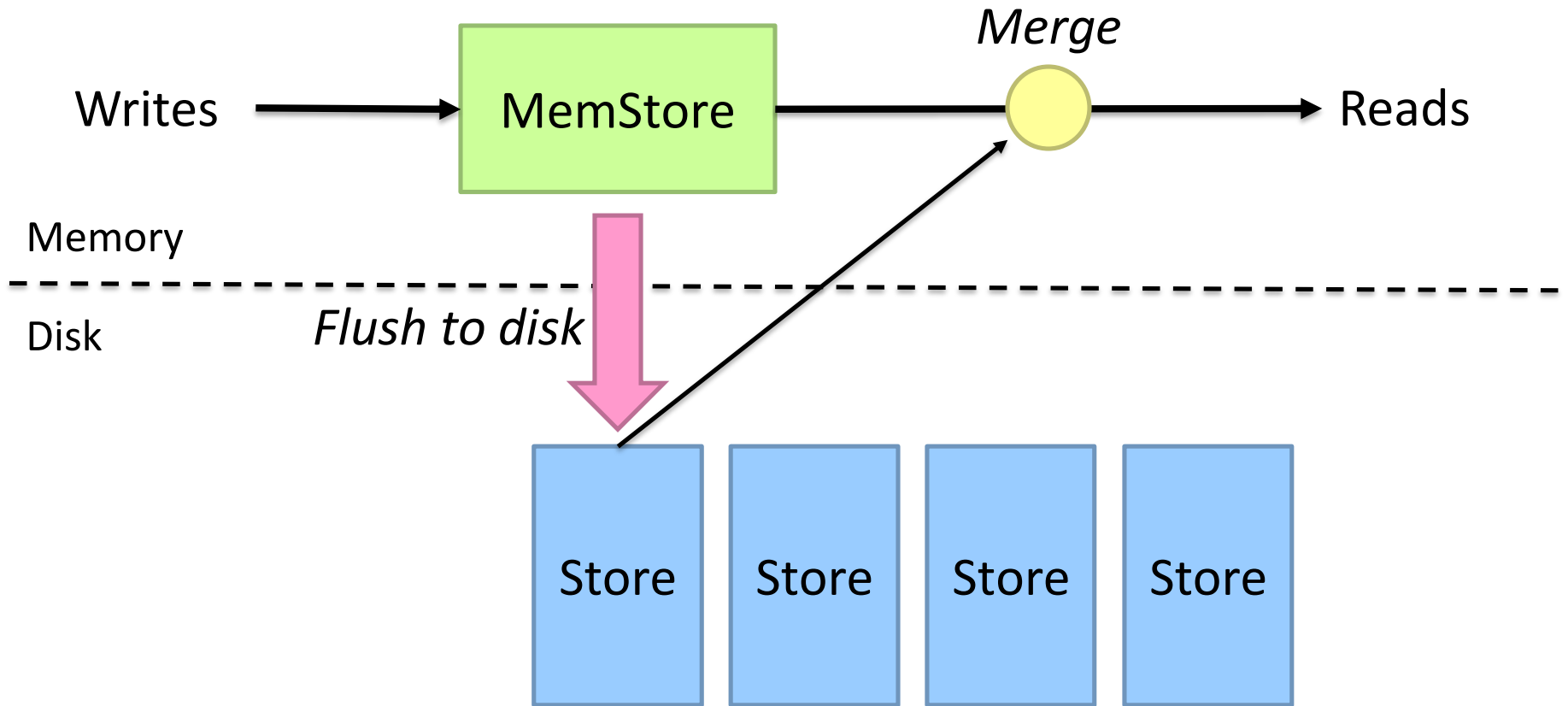


Immutable, indexed, persistent, key-value pairs

What happens as more writes happen?



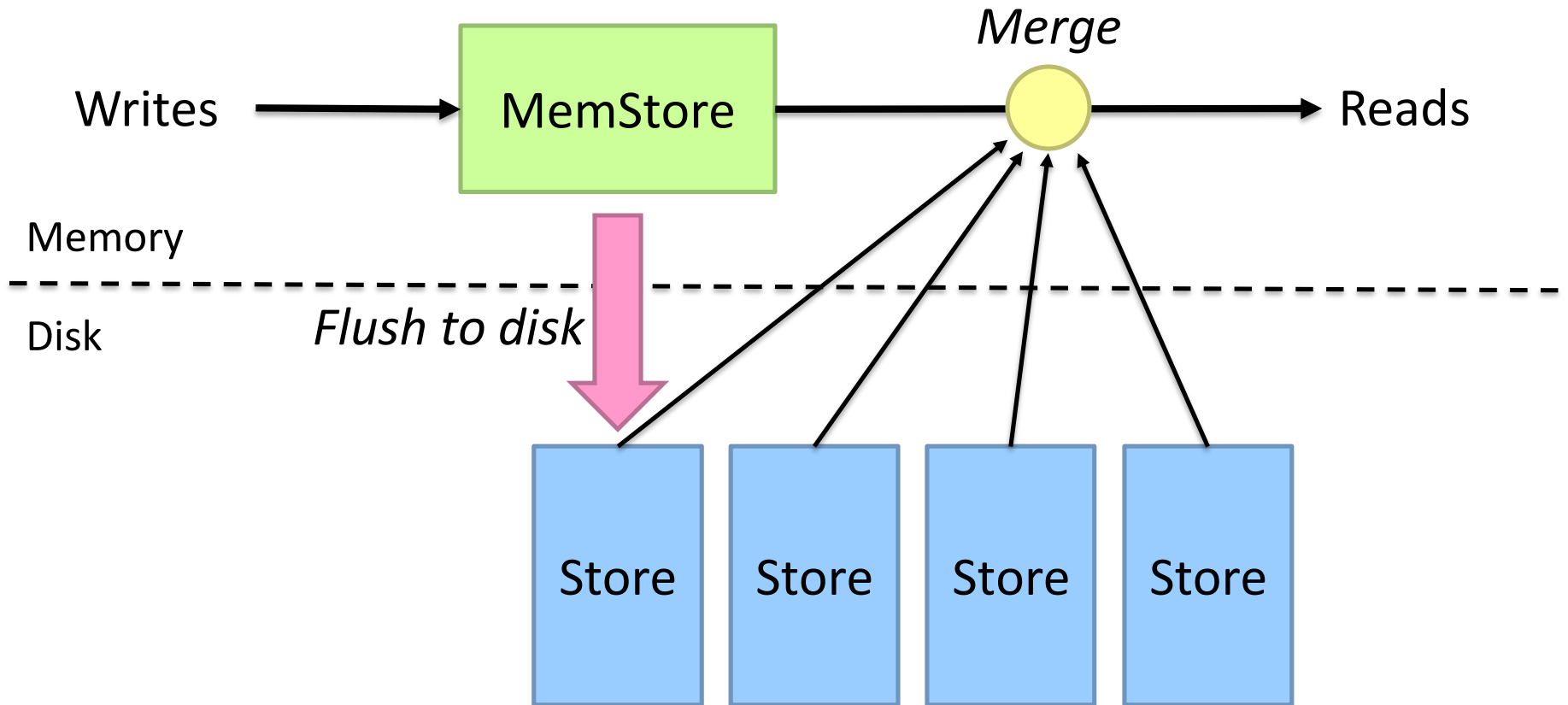
# Log Structured Merge Trees



Immutable, indexed, persistent, key-value pairs

What happens to the read path?

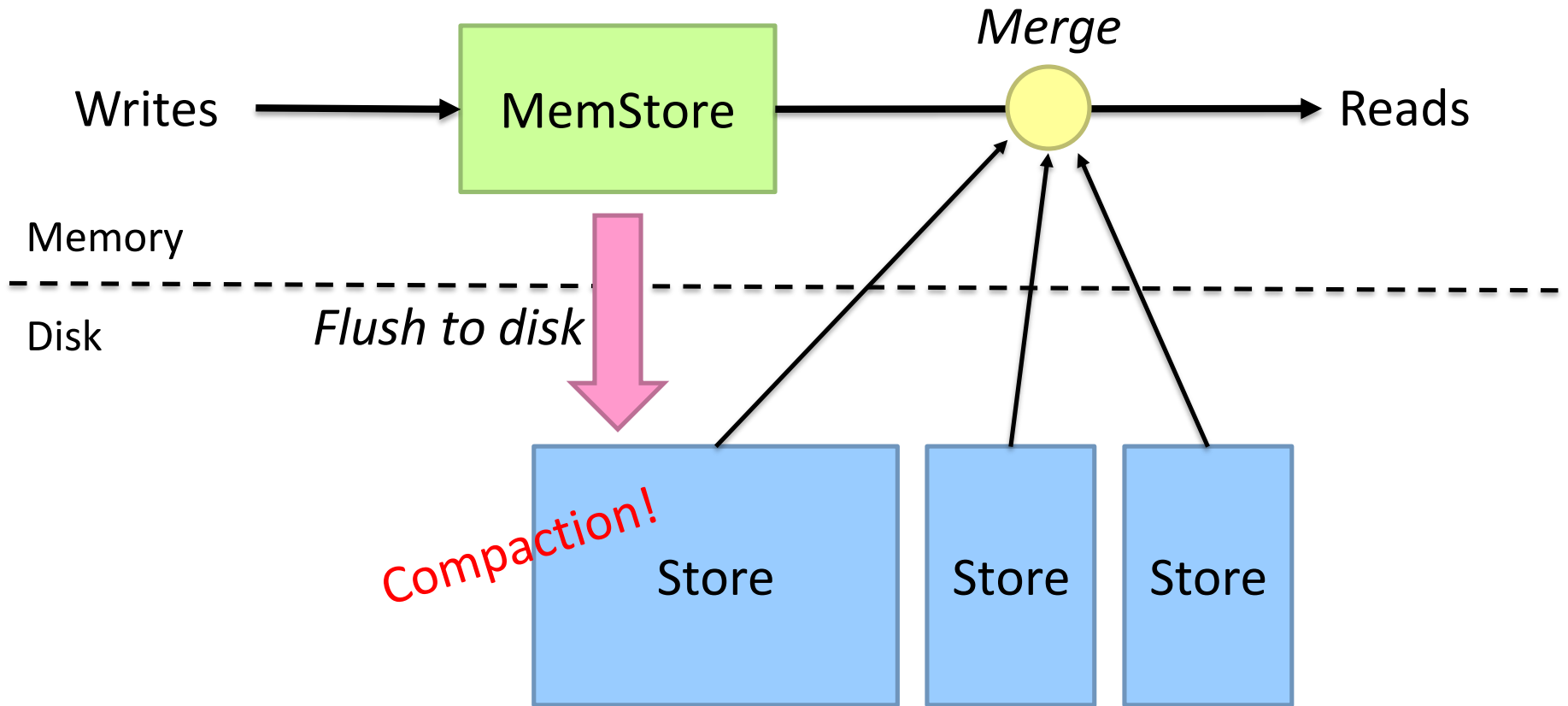
# Log Structured Merge Trees



Immutable, indexed, persistent, key-value pairs

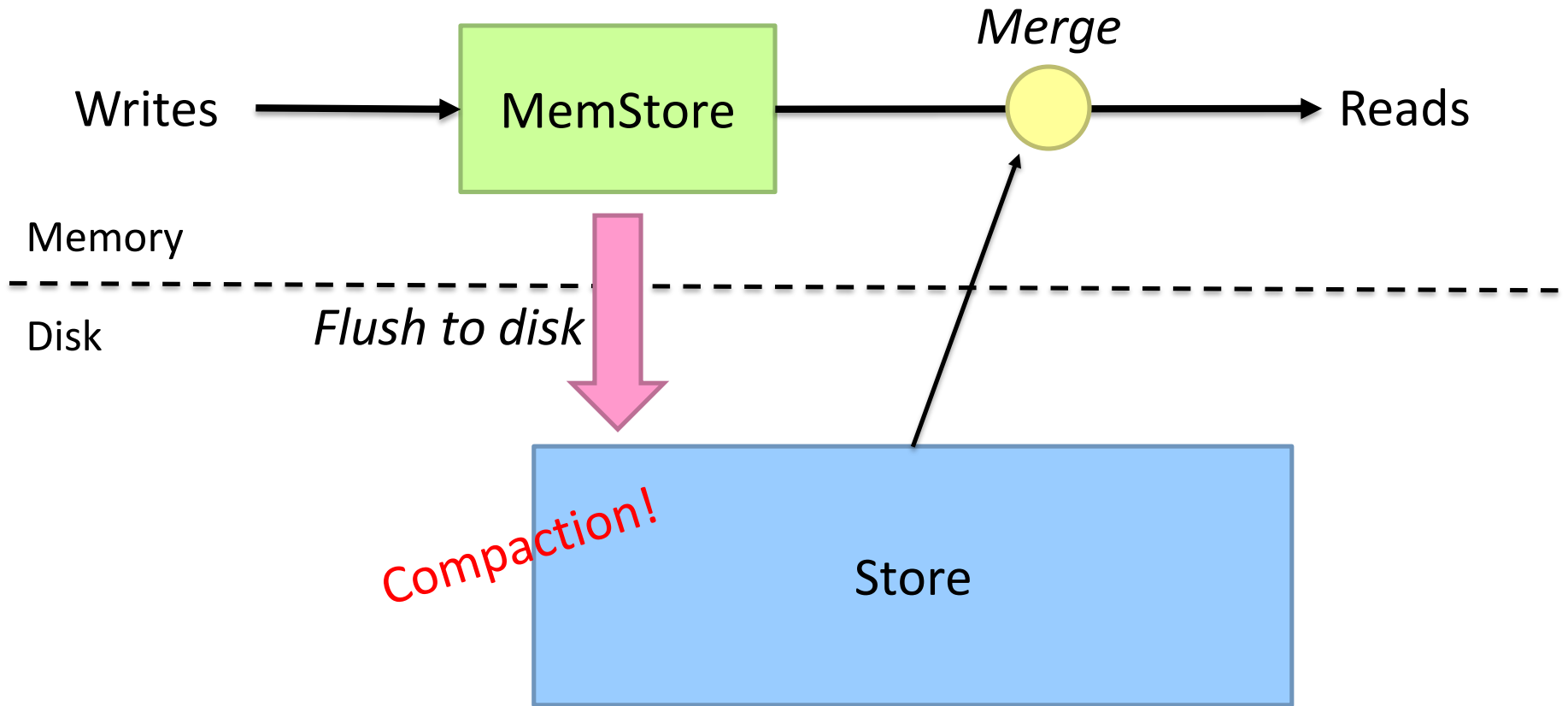
What's the next issue?

# Log Structured Merge Trees



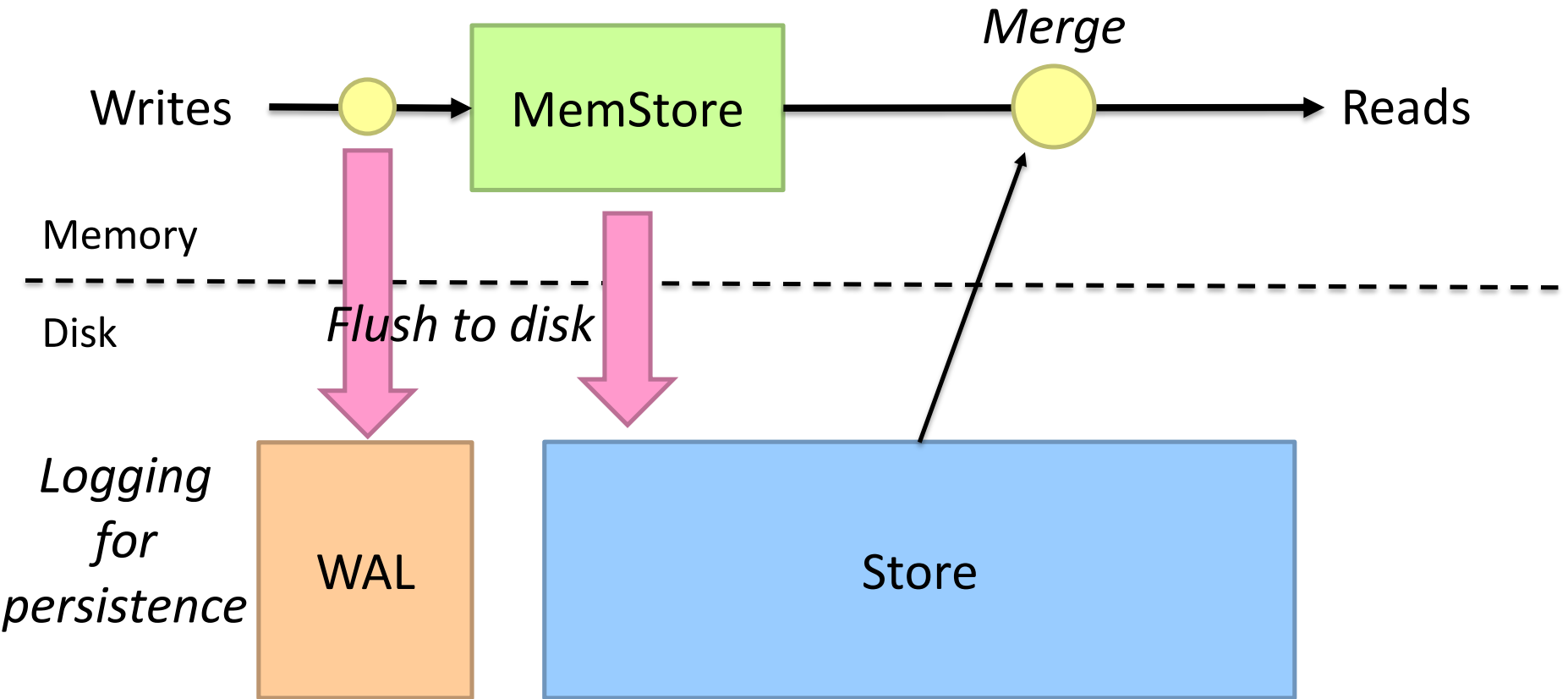
Immutable, indexed, persistent, key-value pairs

# Log Structured Merge Trees



Immutable, indexed, persistent, key-value pairs

# Log Structured Merge Trees

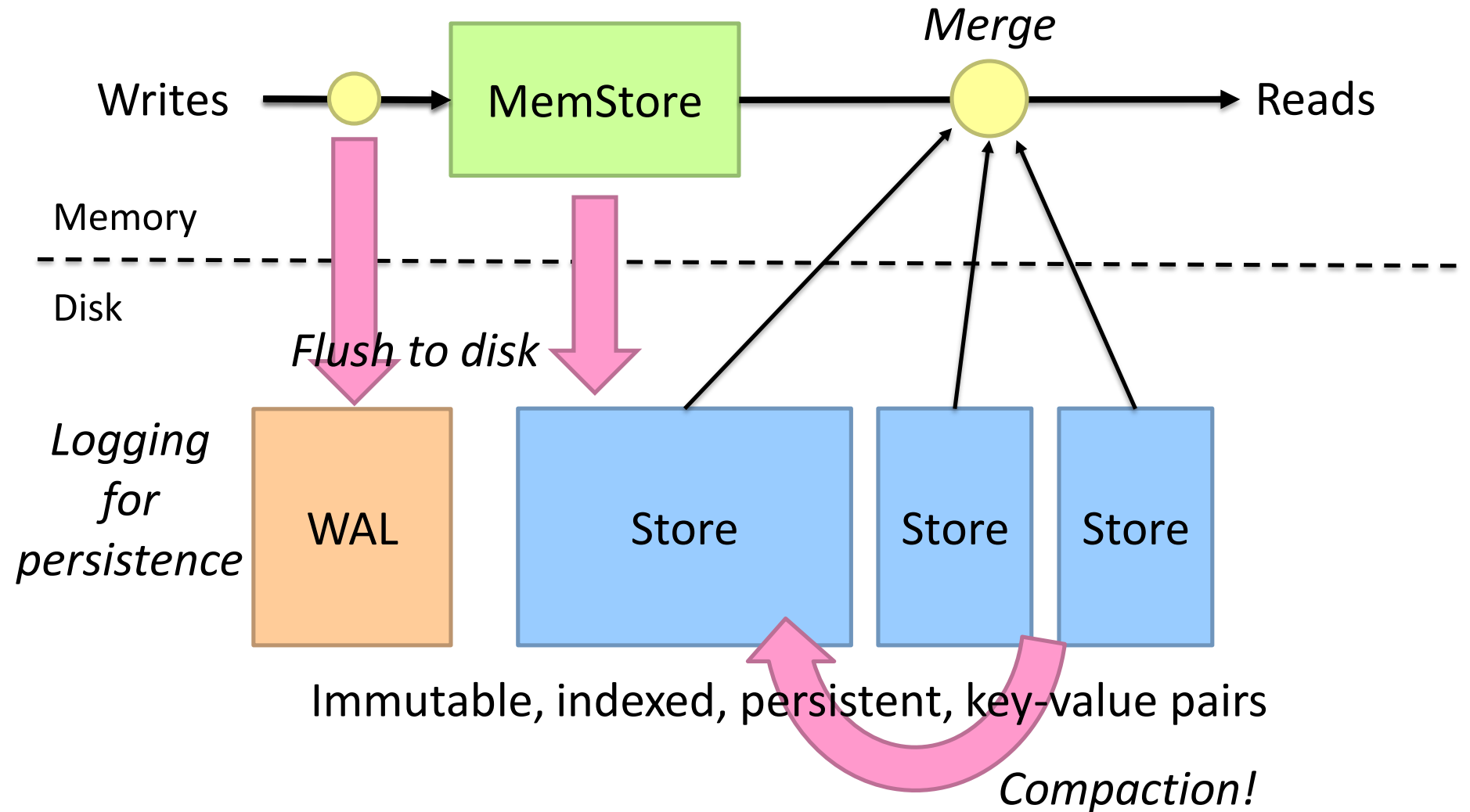


Immutable, indexed, persistent, key-value pairs

One final component...

# Log Structured Merge Trees

The complete picture...

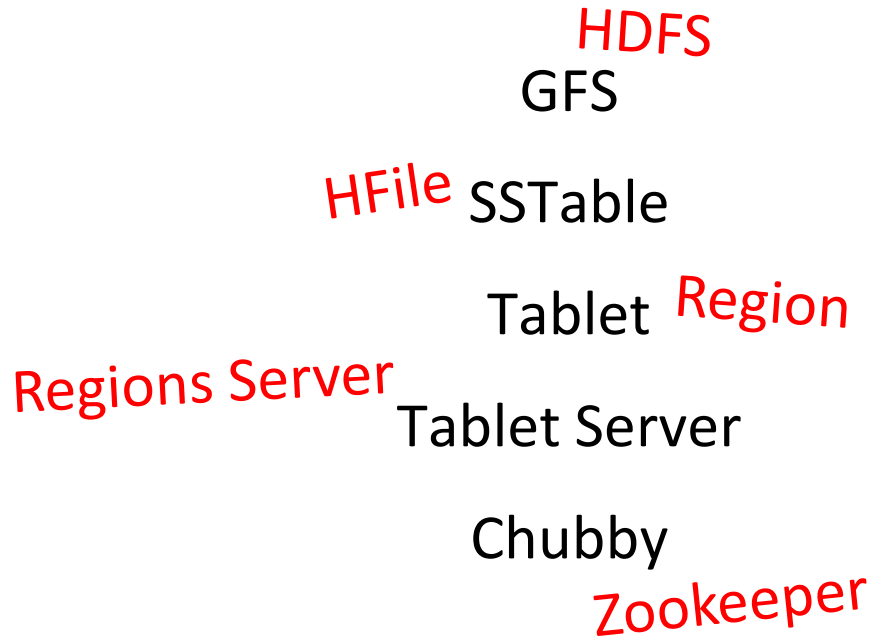


# Log Structured Merge Trees

The complete picture...

Okay, now how do we build a distributed version?

# HBase Bigtable building blocks





# SSTable <sup>HFile</sup>

Persistent, ordered immutable map from keys to values

Stored in GFS: replication “for free”

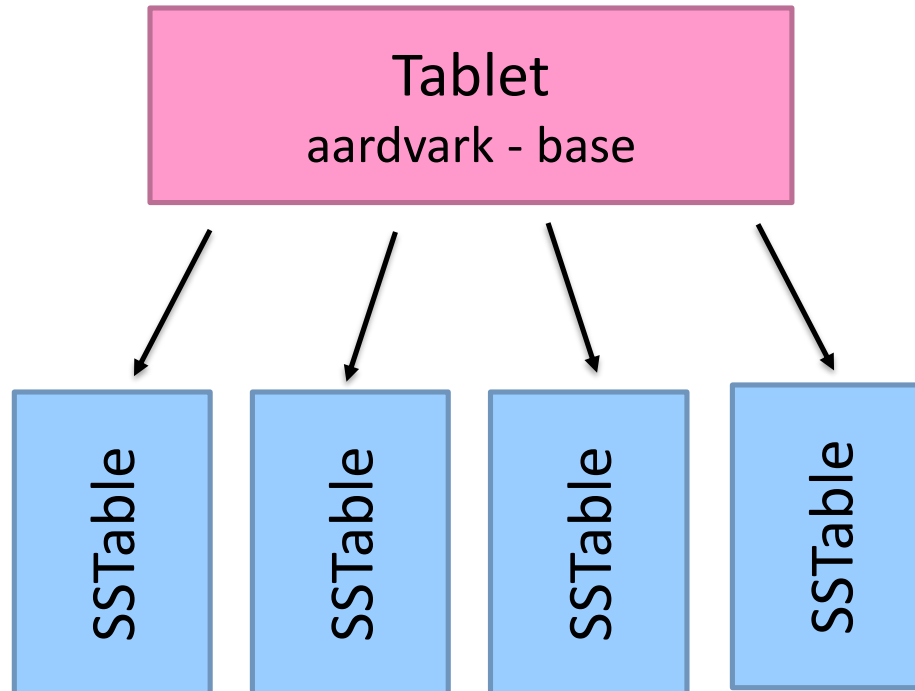
Supported operations:

Look up value associated with key

Iterate key/value pairs within a key range

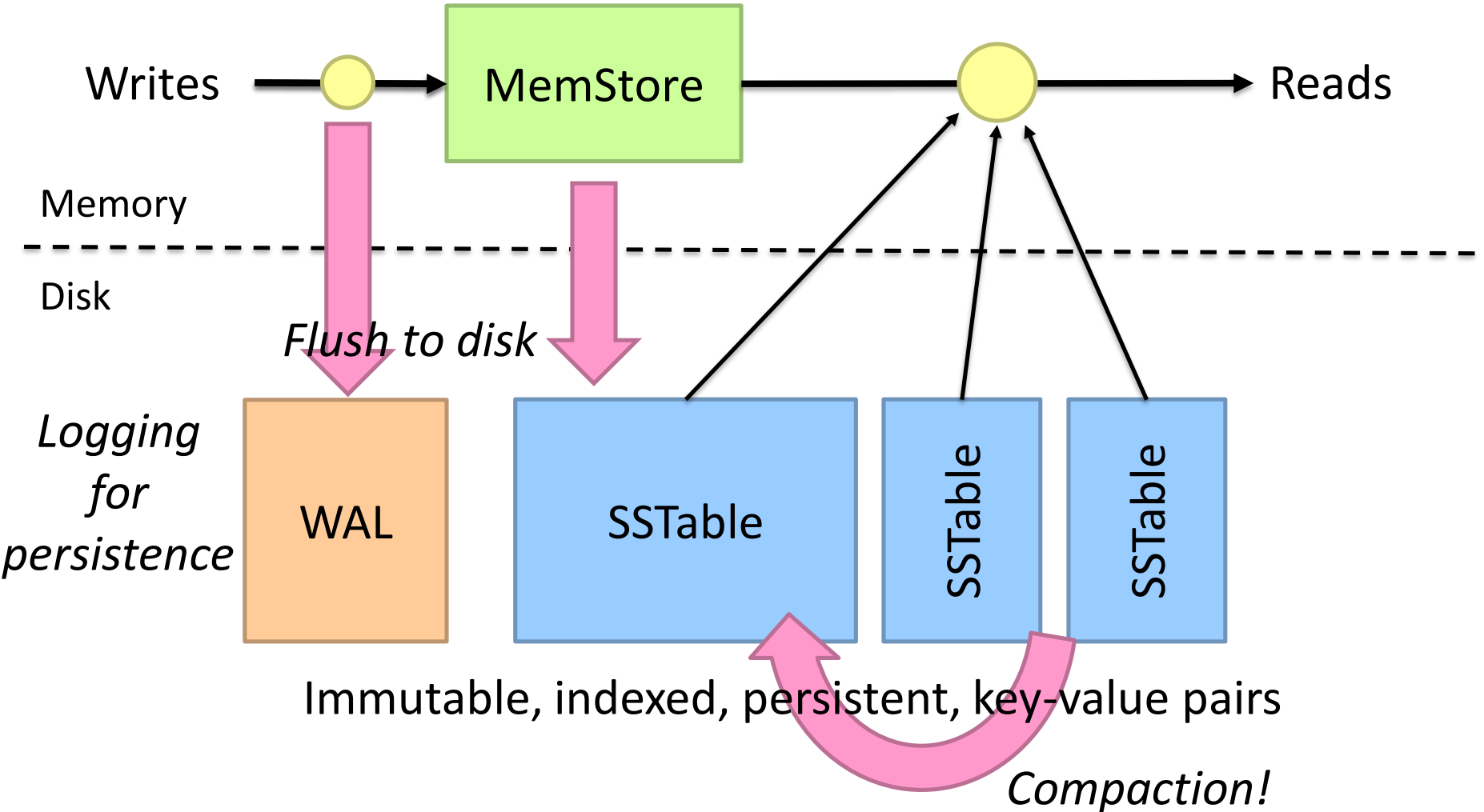
# Region Tablet

Dynamically partitioned range of rows  
Comprised of multiple SSTables



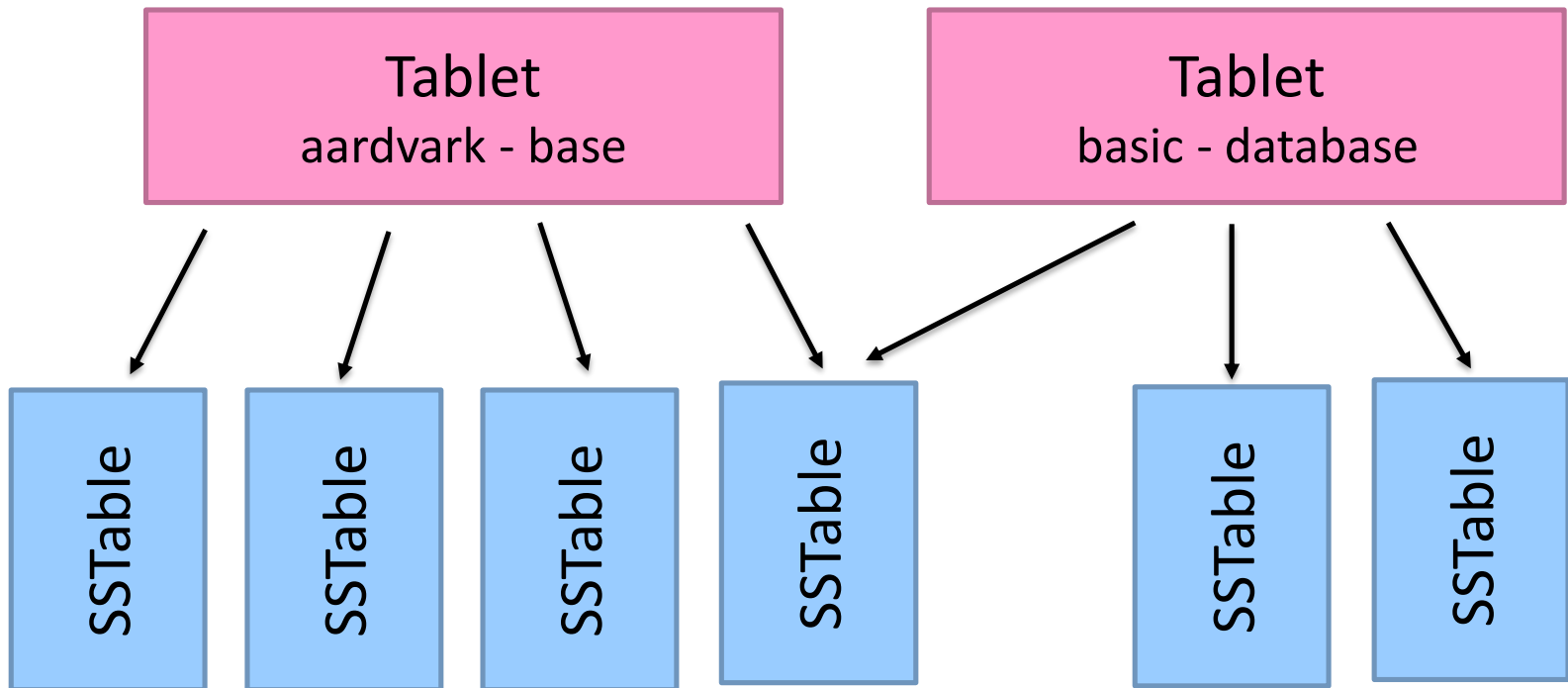
# Region Server

## Tablet Server



# Table

Comprised of multiple tablets  
SSTables can be shared between tablets



Region

Region Server

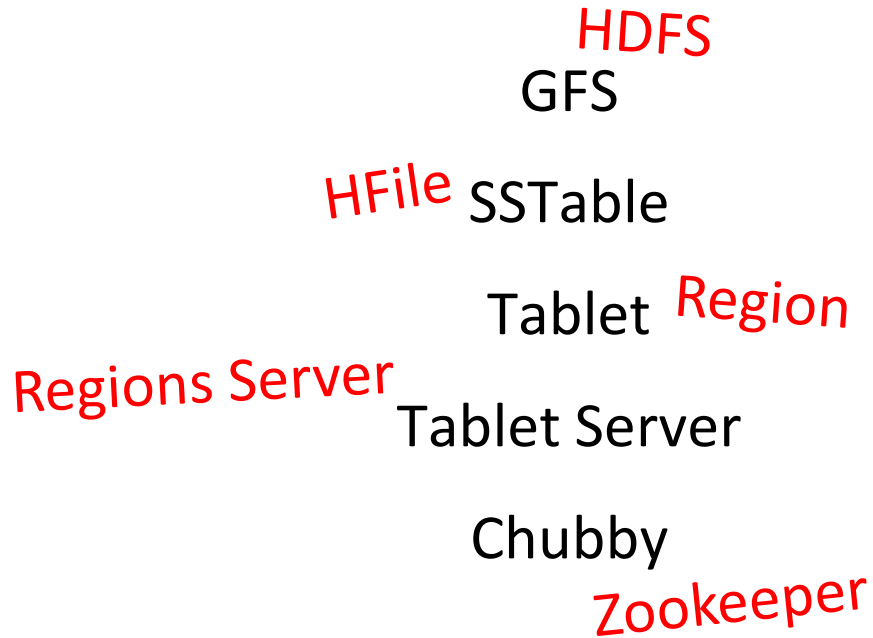
## Tablet to Tablet Server Assignment

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!

# HBase Bigtable building blocks



# Architecture

Client library

Bigtable master **HMaster**

Tablet servers

**Regions Servers**

# 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)



# Compactions

## Minor compaction

Converts the memtable into an SSTable  
Reduces memory usage and log traffic on restart

## Merging compaction

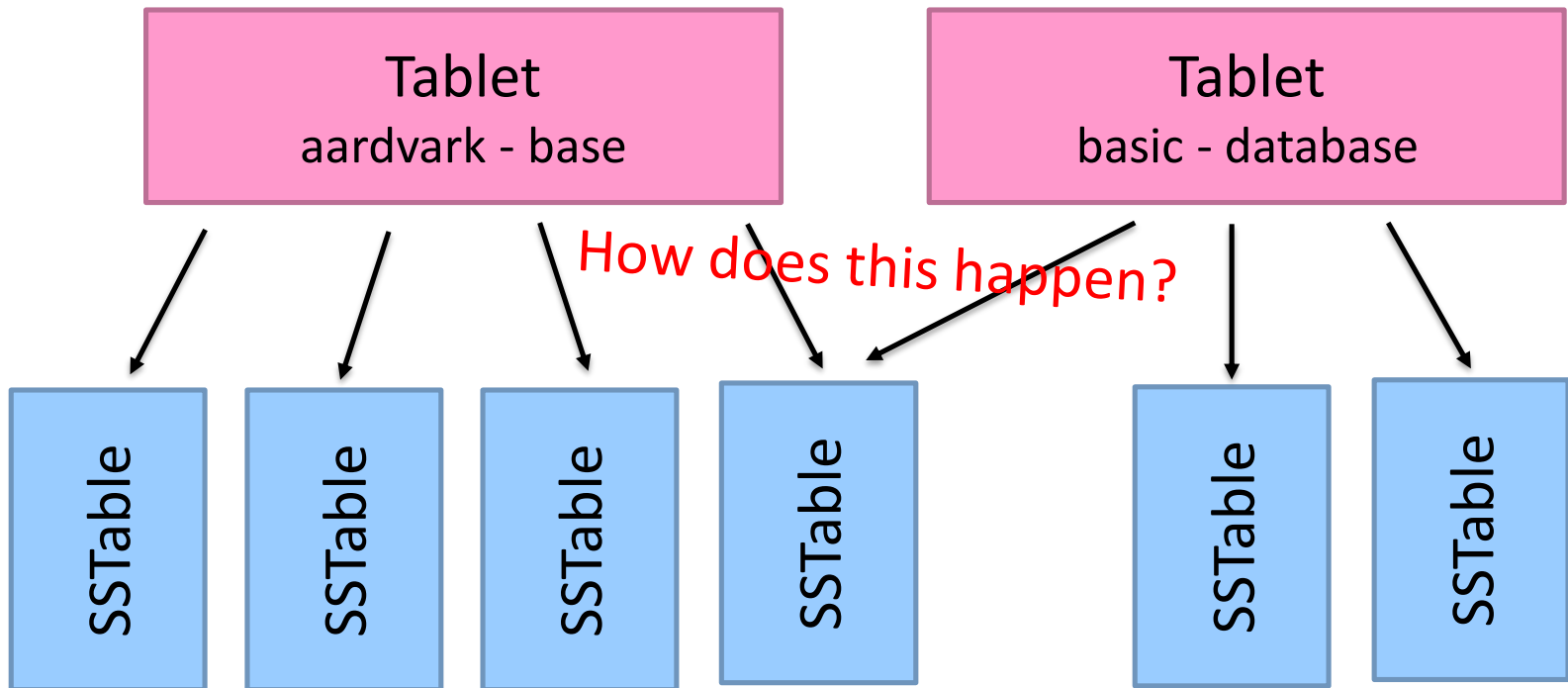
Reads a few SSTables and the memtable, and writes out a new SSTable  
Reduces number of SSTables

## Major compaction

Merging compaction that results in only one SSTable  
No deletion records, only live data

# Table

Comprised of multiple tables  
SSTables can be shared between tablets



# Three Core Ideas

*Keeping track of the partitions?*

Partitioning (sharding)

To increase scalability and to decrease latency

*Consistency?*

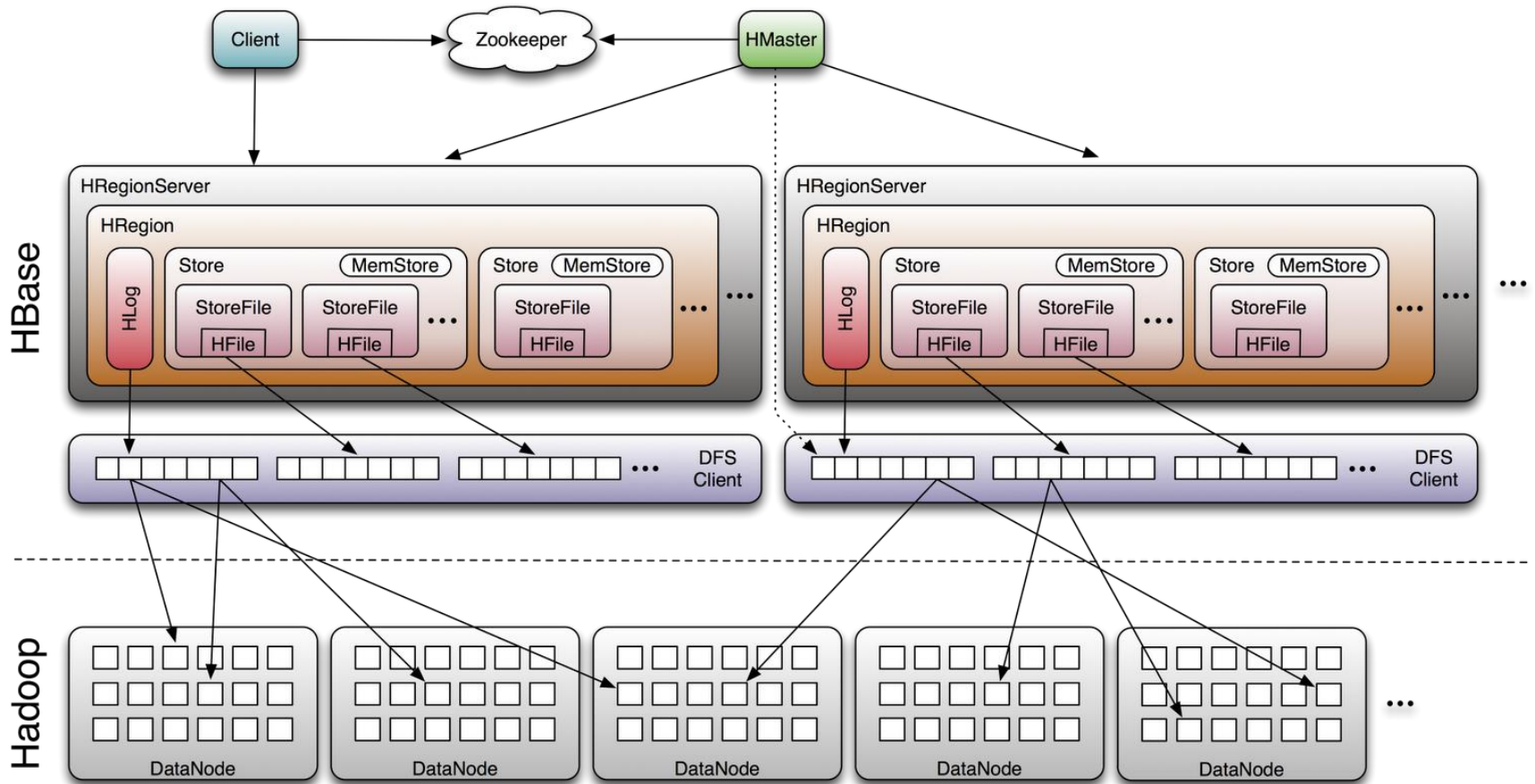
Replication

To increase robustness (availability) and to increase throughput

Caching

To reduce latency

# HBase





Source: Wikipedia (Japanese rock garden)