



# BigTable

## A Distributed Storage System for Structured Data

Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson  
C. Hsieh, Deborah A. Wallach, Mike Burrows,  
Tushar Chandra, Andrew Fikes, Robert E. Gruber

Presented by Michael Ford



Google<sup>TM</sup>  
code

You**Tube**



 **Blogger<sup>TM</sup>**

Google books 



# NoSQL

- Non SQL? Non relational? Not only SQL?
- DBMS without relational model
- Key Features:
  - Simplicity of data model
  - Control of availability
  - Horizontal scaling
  - Flexibility

# NoSQL – CAP Theorem

- Impossible for a distributed system to accomplish all three of:
  - Consistency
  - Availability
  - Partition Tolerance
- NoSQL systems often sacrifice **consistency** by choosing **eventual consistency** over ACID Transactions

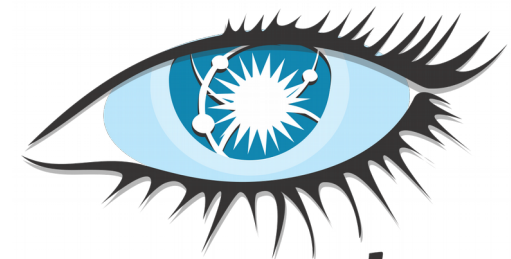




KEY-VALUE



Couchbase



cassandra



mongoDB®



Key	Value
"India"	{"B-25, Sector-58, Noida, India – 201301"}
"Romania"	{"IMPS Moara Business Center, Buftea No. 1, Cluj-Napoca, 400606", City Business Center, Coriolan Brediceanu No. 10, Building B, Timisoara, 300011"}
"US"	{"3975 Fair Ridge Drive. Suite 200 South, Fairfax, VA 22033"}



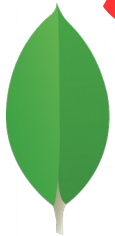
KEY-VALUE



DOCUMENT

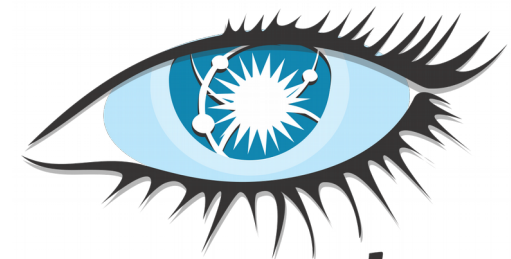


couchbase



STORE

mongoDB



cassandra



```
1 {officeName:"3Pillar Noida",
2 {Street: "B-25, City:"Noida", State:"UP", Pincode:"201301"}
3 }
4 {officeName:"3Pillar Timisoara",
5 {Boulevard:"Coriolan Brediceanu No. 10", Block:"B, Ist Floor", City: "Timisoara", Pinc
6 }
7 {officeName:"3Pillar Cluj",
8 {Latitude:"40.748328", Longitude:"-73.985560"}
9 }
```



ORACLE<sup>®</sup>  
Coherence



KEY-VALUE



INFINITEGRAPH<sup>®</sup>

GRAPH

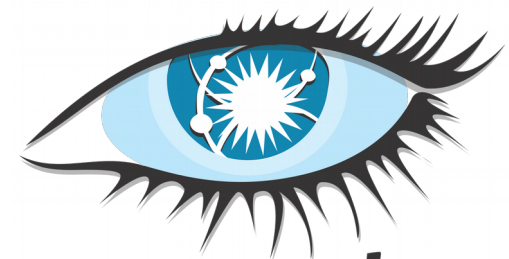


Couchbase

DOCUMENT  
STORE



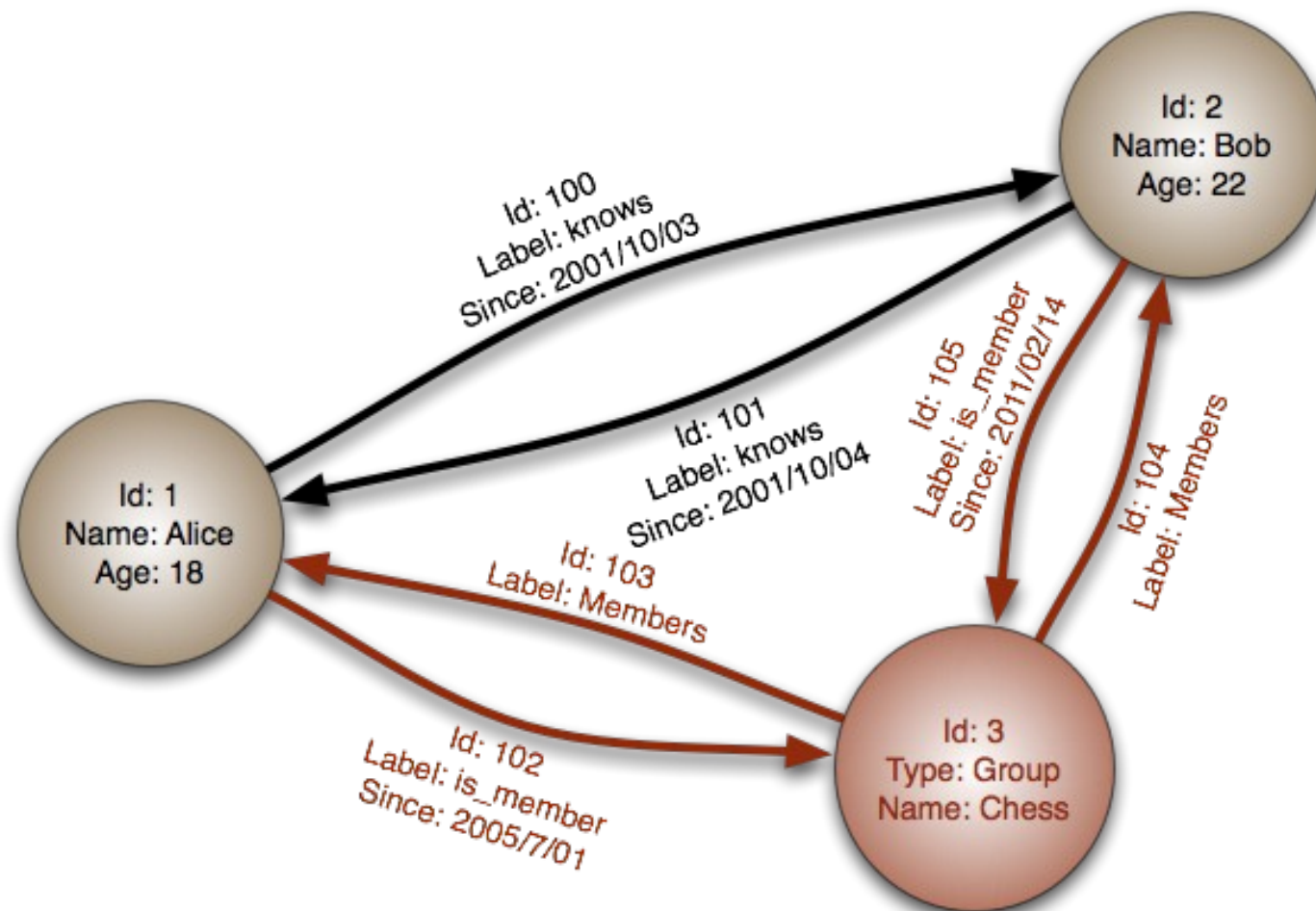
mongoDB<sup>®</sup>



cassandra

APACHE  
HBASE





ORACLE<sup>®</sup>  
Coherence



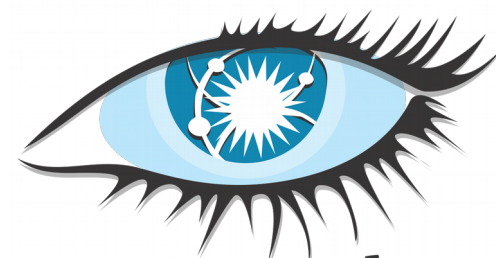
INFINITEGRAPH<sup>®</sup>



Couchbase



mongoDB<sup>®</sup>



cassandra

APACHE  
HBASE





# BigTable



# Overview

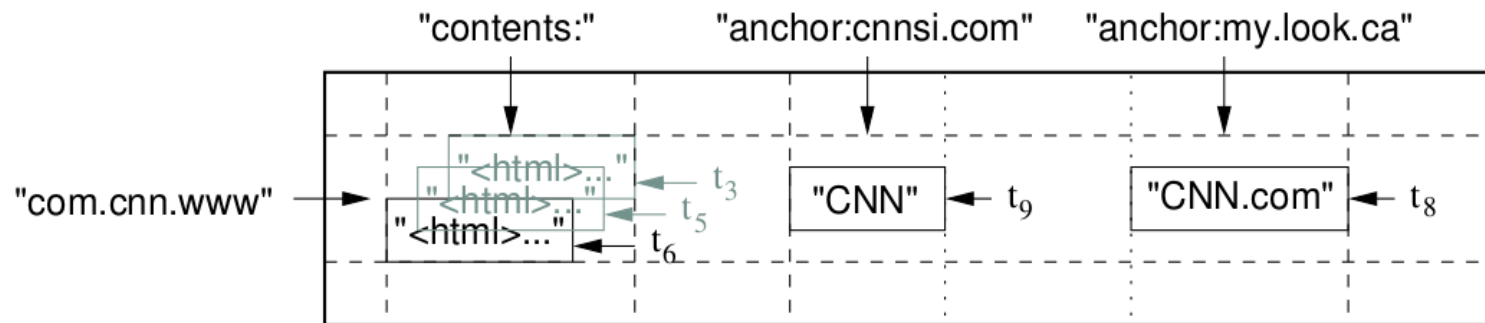
- Multi-dimensional, sparse, sorted map with tabular structure
- Dynamic control of schema:
  - Data layout
  - Data format
- Schema choice determines data locality
- Therefore user choices significantly affect speed



# Indexing

- Rows
  - Sorted alphabetically, keeps data together if using reverse naming
- Column-families
  - Need to provide column-family:column to retrieve data
  - Columns in family are compressed together so should be same type
  - Designed to have few families but many columns
- Timestamps
  - Default assignment real-time (microseconds)
  - Can be client assigned
  - Garbage collection can be set to automatically keep only n most recent





# Partitioning: Tablets

- Range of rows
- Unit of data distribution
- 100-200 MB each
- Automatic splitting

# Building Blocks

- Google File System
- SSTable
- Chubby



# ● Google File System

- Distributed data storage system
- Files divided into **fixed-size chunks**
- Chunks replicated 3 times across **chunk servers**
- Chunk location and metadata managed with **master server**

# Sorted String Table - SSTable

- Key-value map datatype
- Immutable
- Sorted sequentially for easy lookup



# Chubby

- Distributed lock service
- Determines lock status with **Paxos consensus algorithm**
- 5 chubby replicas
  - 1 of which is master
- Used for:
  - Ensure there is only 1 active master
  - Store location of root tablet
  - Discover and delete tablet servers
  - Store column family information
  - Store access control lists





# Implementation

# General Structure

## Tablet servers

- Manages set of tablets (10-1000 tablets)
- Manages read/write directly with client
- Splits tablets when too big

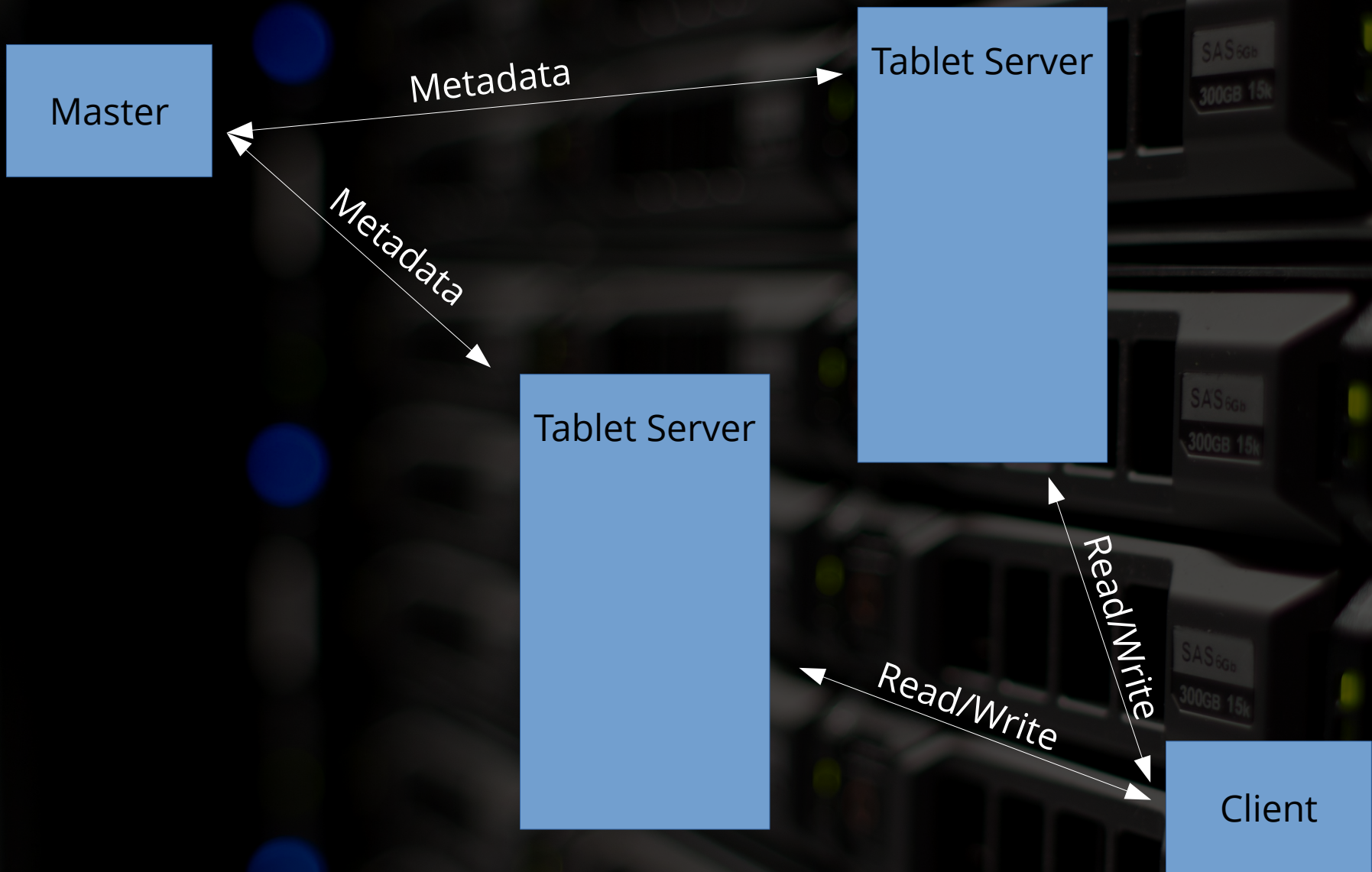
## Master server

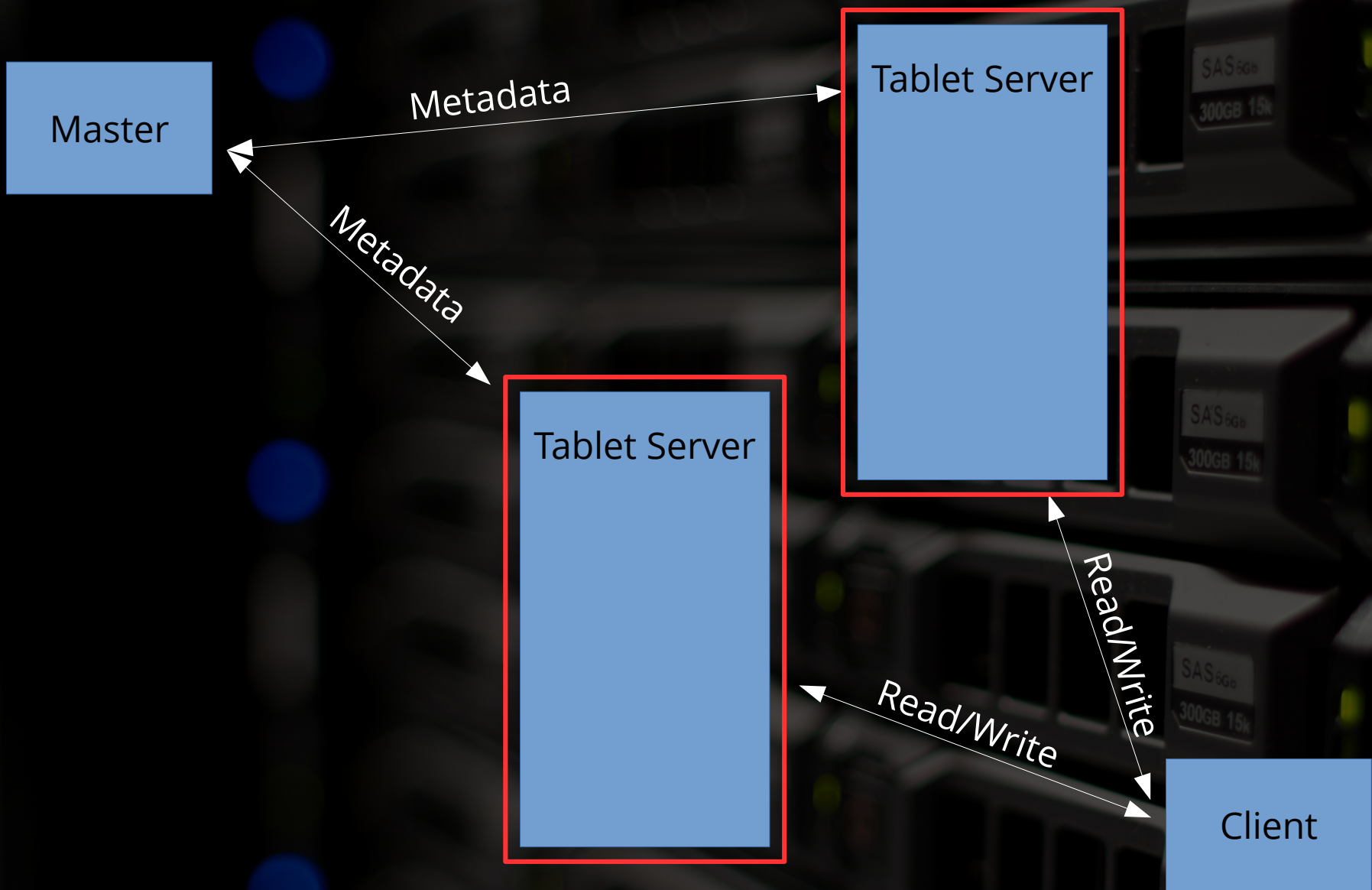
- Assigns tablets to tablet servers
- Managing tablet servers and load
- Garbage collection

## Client

- Can iterate over: column families, columns within a family, rows









# Tablet Servers

- Tablet assignment managed by master server
  - Acquires exclusive chubby lock in servers directory

Chubby file



# Tablet Servers

- Tablet assignment managed by master server
  - Acquires exclusive chubby lock in servers directory
- Tablet location storing
  - Chubby File
  - Root tablet
  - Metadata tablets
  - Tables containing data tablets



# Tablet Serving

## Tablet Storage

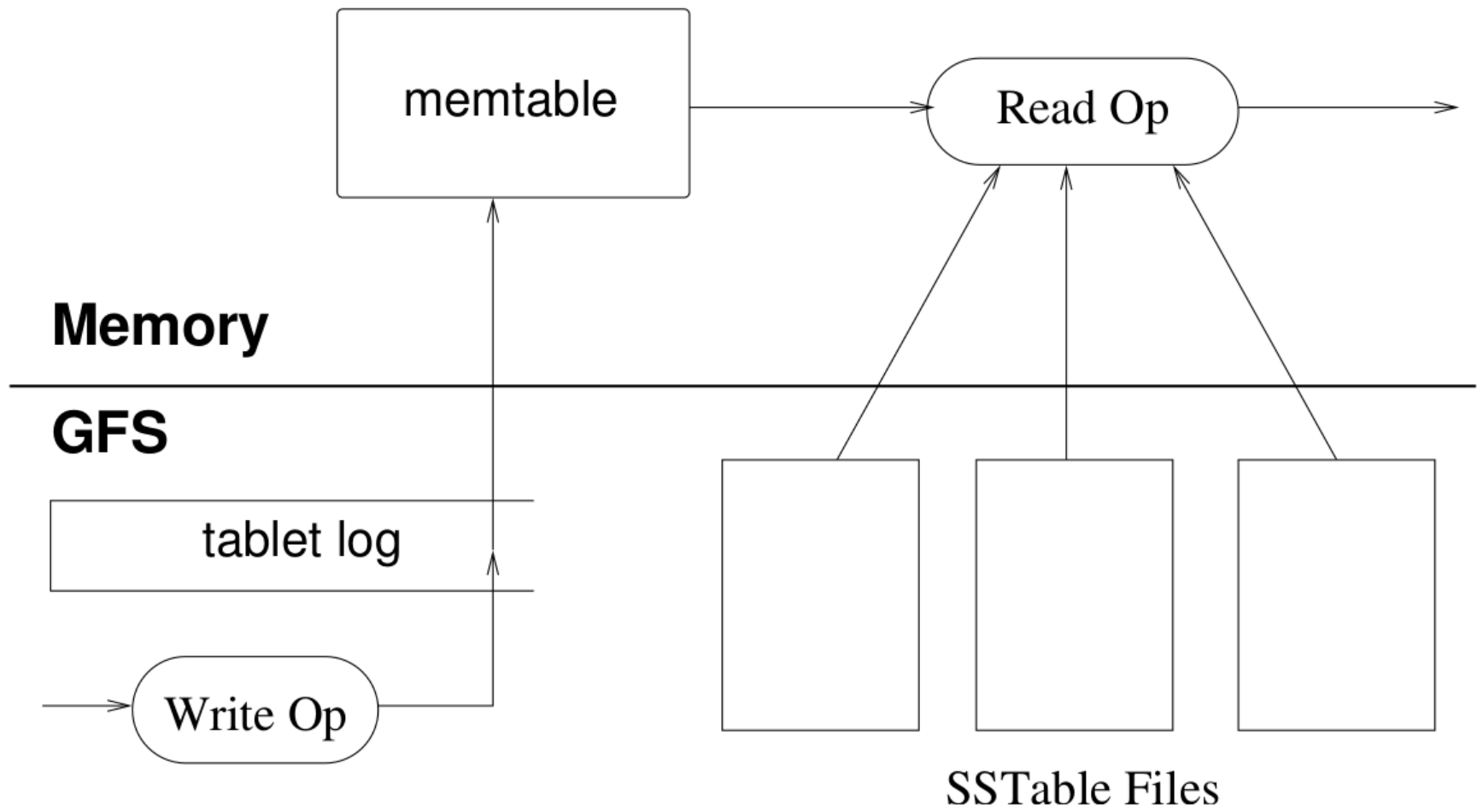
- Recent commits logged in *memtable*
- Old commits logged in SSTables

## Tablet Read

- SSTable locations from metadata
- Constructs tablet by merging SSTables and memtable using commits

## Tablet Write

- Committed to log
- Contents inserted into memtable





# Compaction

Problem: memtable fills up with commits

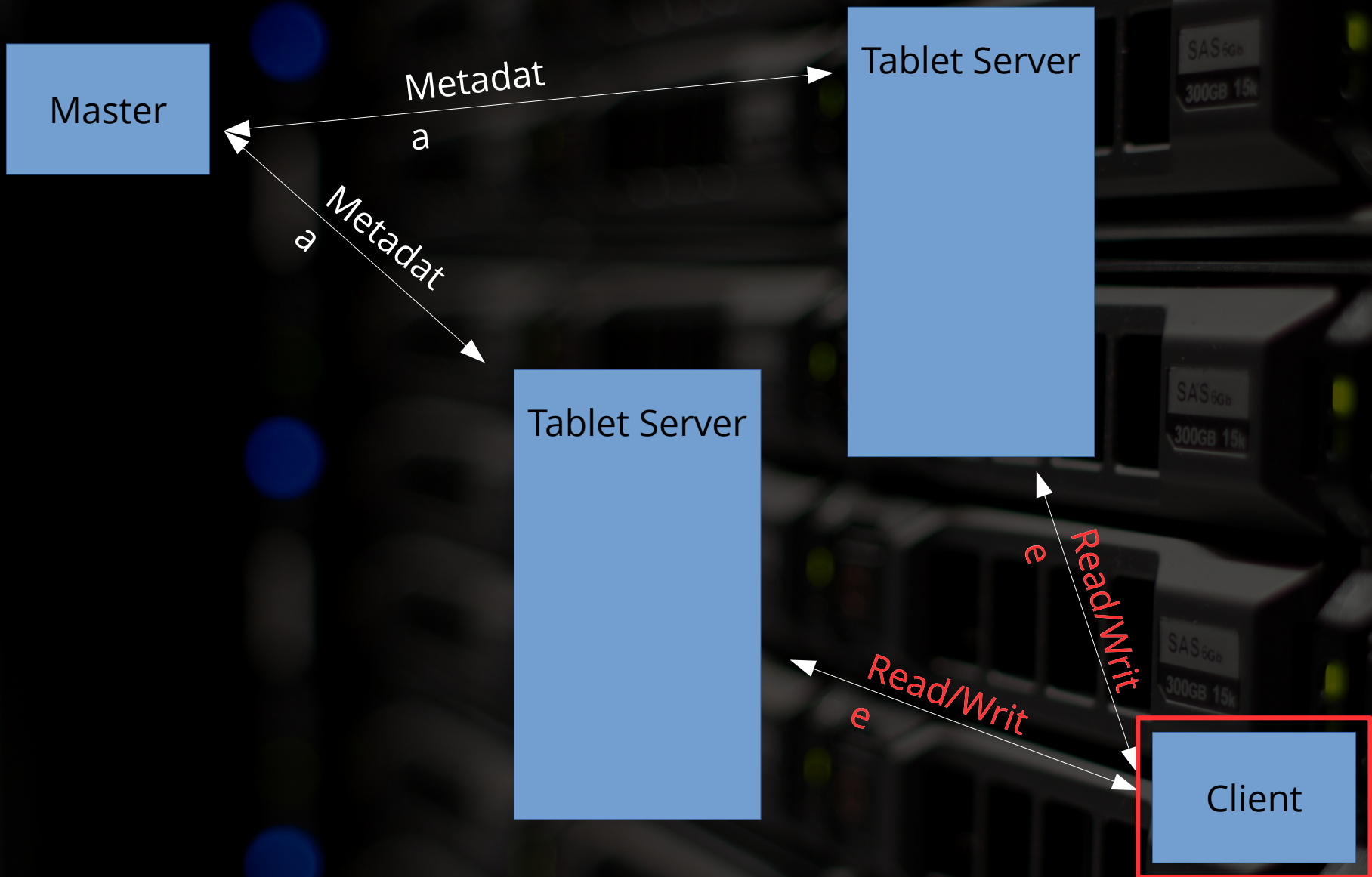
Minor compaction

- Memtable made into new SSTable
- New memtable created

Major compaction

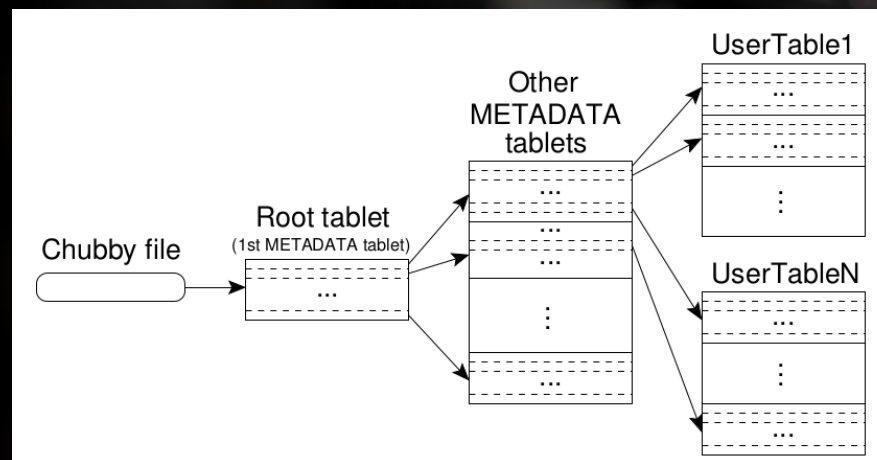
- Combines memtable and SSTables with no deletions



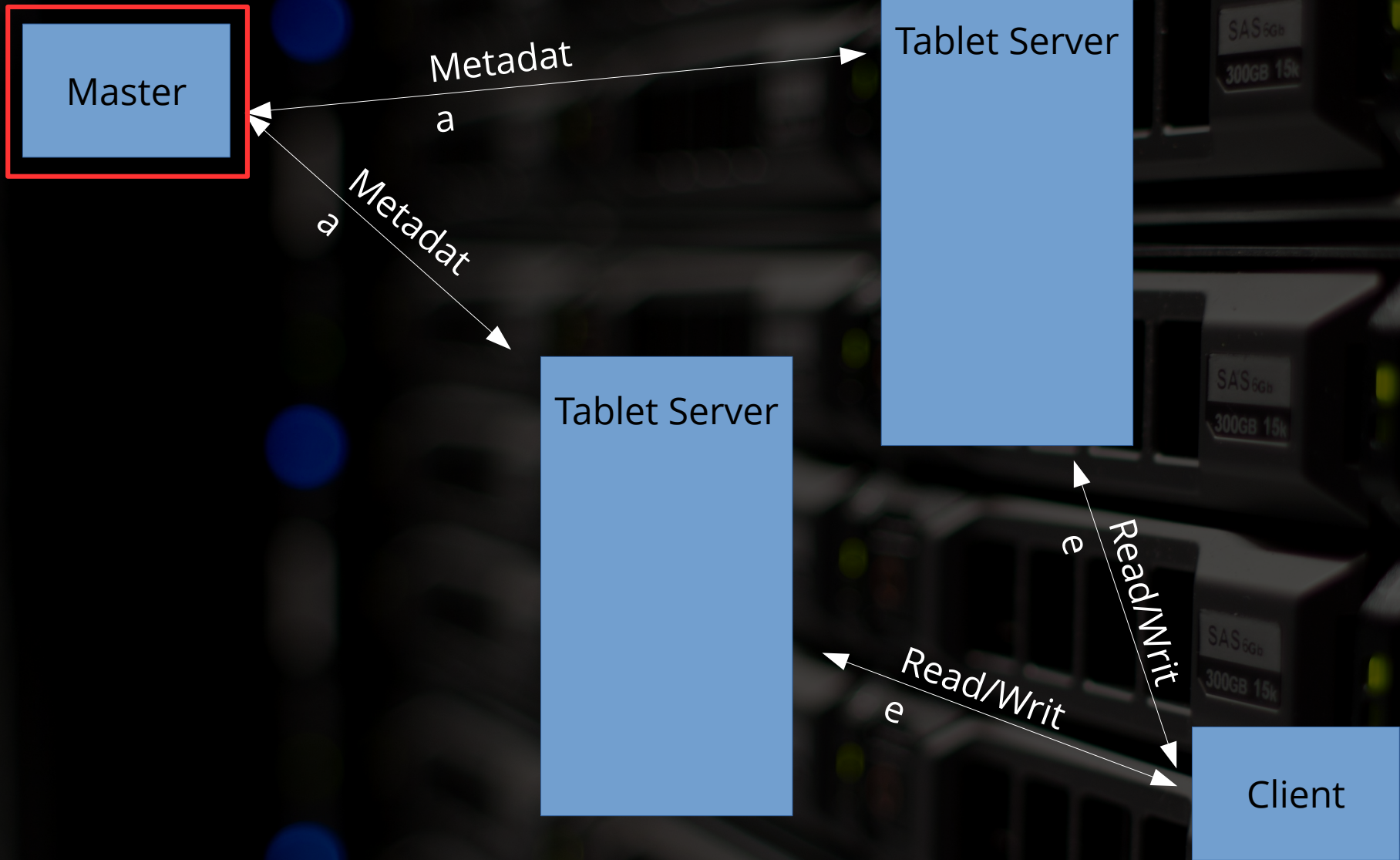


# Client Features

- Caches tablet locations
- If location inaccurate, moves up tablet server hierarchy
- Atomic single-row transactions











# Master Server

## Jobs

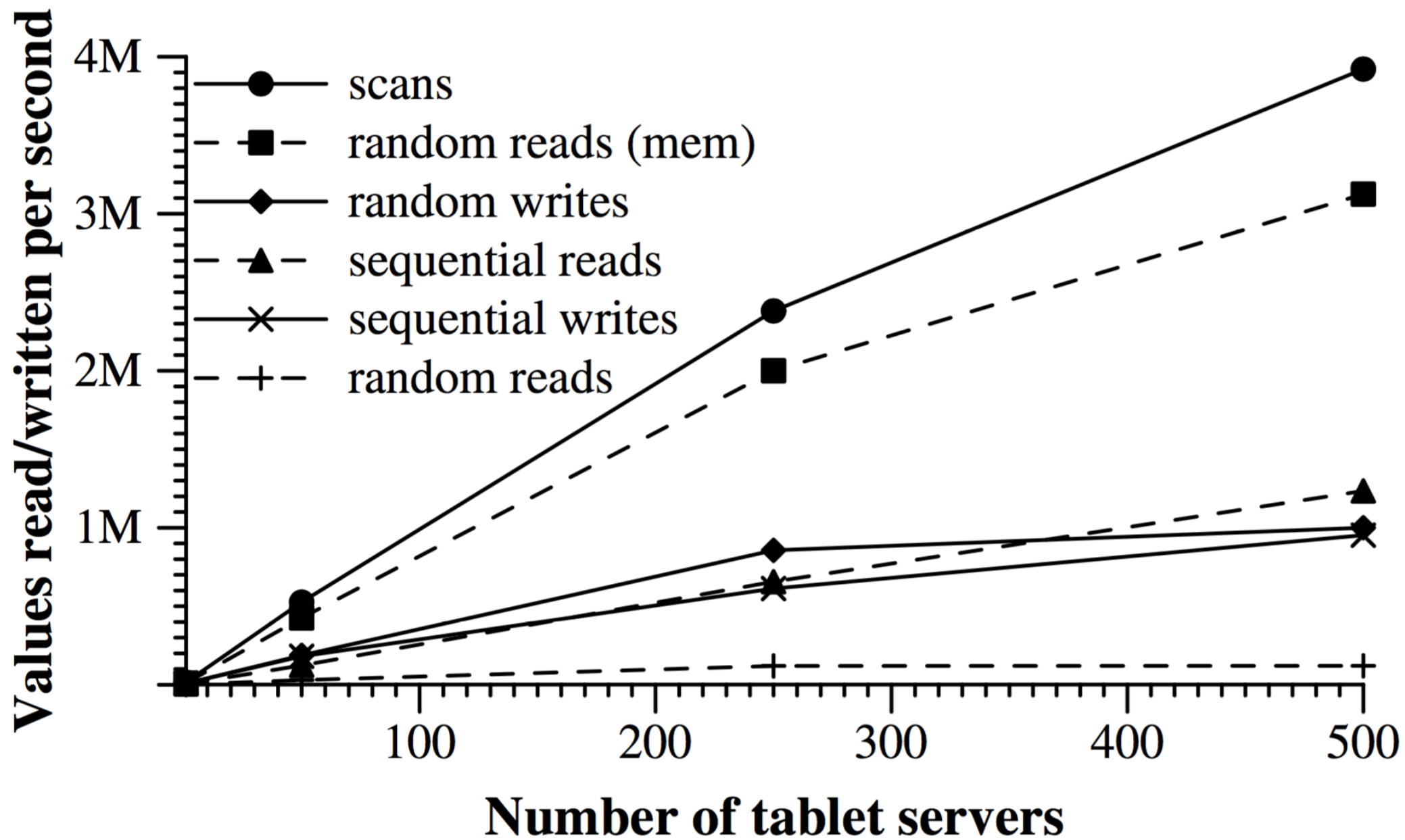
- 1) Managing tablet servers
- 2) Assigning tablets
- 3) Garbage collection

# Experiments and Examples

## Benchmarks

- Sequential write
- Sequential read
- Scan
- Random reads
- Random write

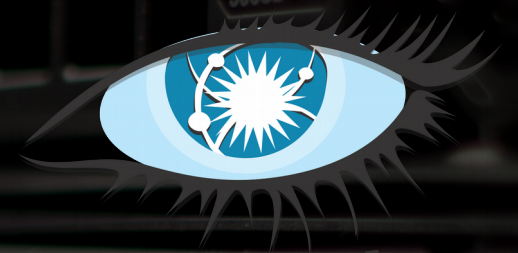






# Future directions

- Apache Cassandra, Hbase
- Google Cloud
- Google Spanner



*cassandra*

APACHE  
**HBASE**



Google  
Cloud Platform