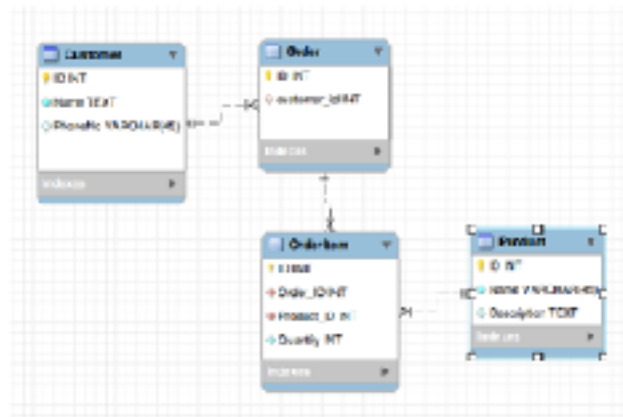# Why do we need NoSQL / HBase?
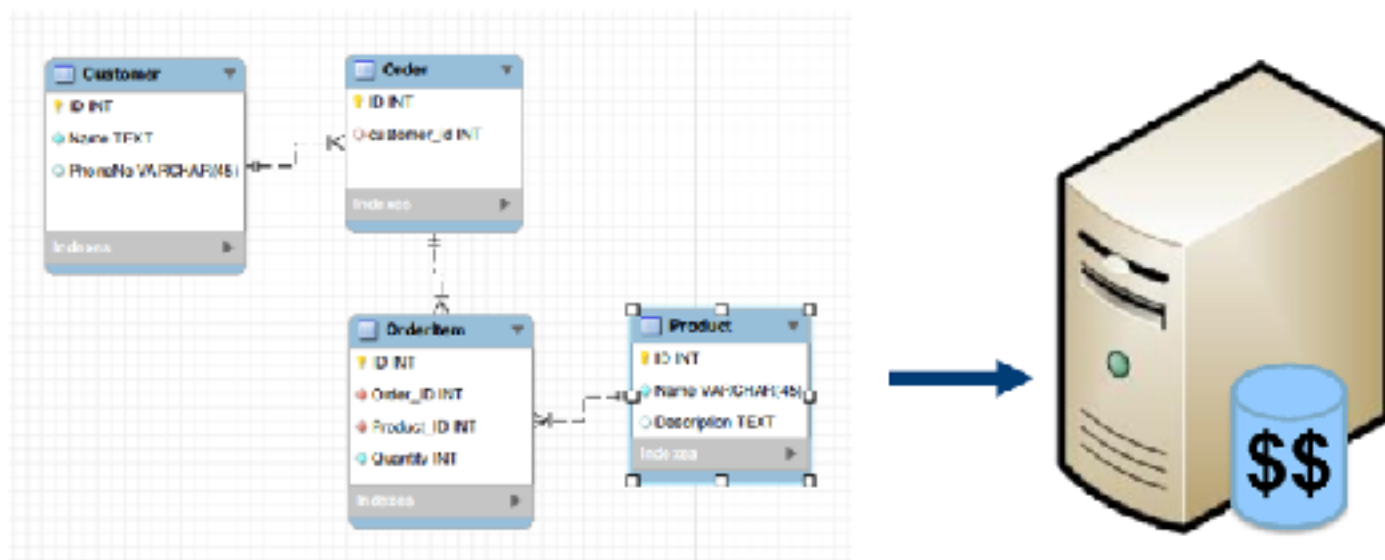## Relational Model

- Pros
  - **Standard** persistence model
  - **Transactions handle**
    - **concurrency**, **consistency**
  - efficient and robust structure for storing data
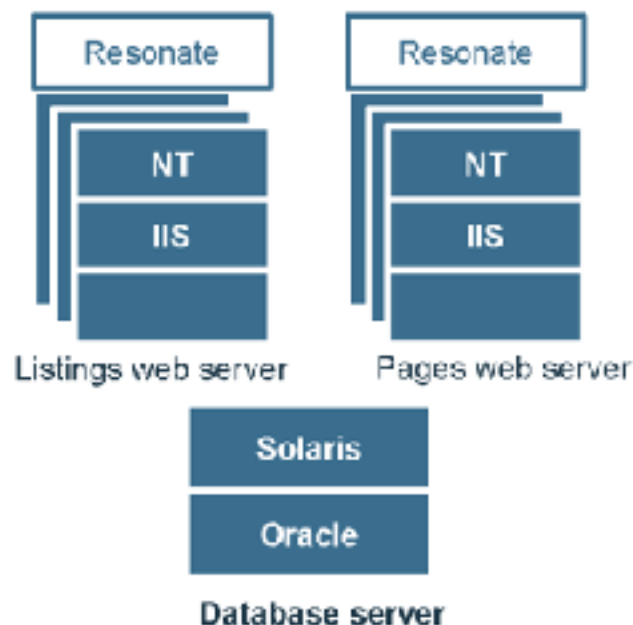
# RDBMS Scaling Up Example - eBay

- Back End Oracle Database server **scaled vertically to a larger machine** (Sun E10000)



http://www.addsimplicity.com/downloads/eBaySDForum2006-11-29.pdf

**MAPR** 9

Horizonal scale: split
table by rows into
partitions across a cluster

| Key | colB | colC |
|-----|------|------|
| val | val | val |
| xxx | val | val |

id 1-1000

| Key | colB | colC |
|-----|------|------|
| val | val | val |
| xxx | val | val |

id 1000-2000

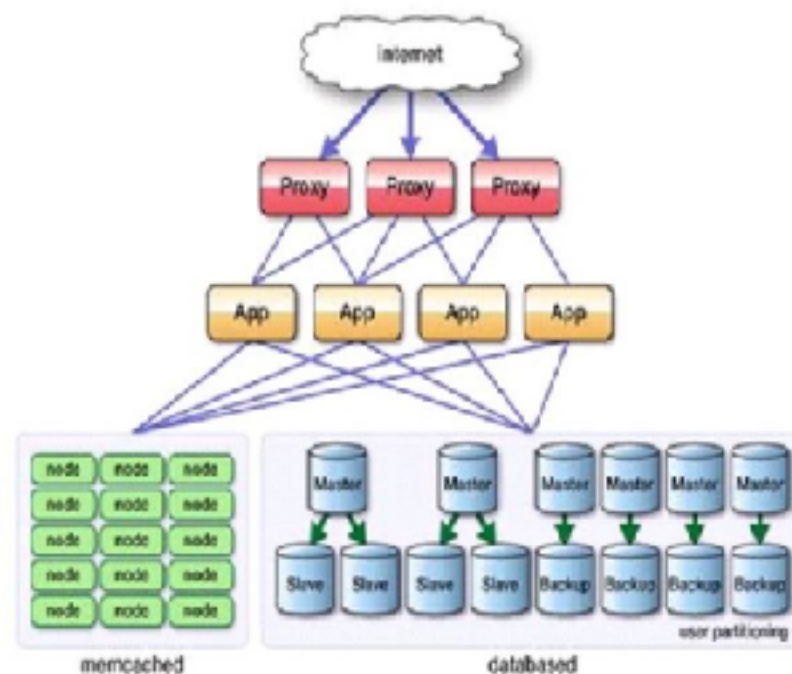| Key | colB | colC |
|-----|------|------|
| val | val | va |
| xxx | val | val |

id 2000=3000

- **Horizontal** scaling
  - Cheaper than vertical
  - parallel execution
- Relational databases were **not designed to do this** automatically

**MAPR** 10

# Facebook 2010

- 9000 memcache instances
- **4000 Shards mysql**



http://gigaom.com/2011/07/07/facebook-trapped-in-mysql-fate-worse-than-death/

# Hbase designed for Distribution, Scale, Speed

# Google

## Big Table

- Distributed Storage System
- Paper published in 2006.

## Google File System

## MapReduce

**Runs on commodity hardware**
**Designed to Scale**

**MAPR** 15

# HBase is a ColumnFamily oriented Database

| Customer id | Customer Address data | | | Customer order data | | |
|---|---|---|---|---|---|---|
| **RowKey** | **CF1** | | | **CF2** | | |
| | **colA** | **colB** | **colC** | **colA** | **colB** | **colC** |
| **axxx** | **Val** | | val | val | | val |
| **gxxx** | | val | | | val | |

Data is accessed and stored together:

- RowKey is the primary index
- Column Families group similar data by row key

# HBase is a Distributed Database

Put, Get by Key

Data is automatically distributed across the cluster

- **Key range** is used for horizontal partitioning

distributed data stored and accessed together:

- Pros
    - **scalable**
    - **Fast** Writes and Reads
- Cons
    - **No** joins
    - No dynamic queries
    - Need to **know** how data will be **queried in advance**

# Hbase Data Model

# HBase Data Model- Row Keys

| RowKey | Address | | | Order | | | | ... |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | street | city | state | Date | ItemNo | Ship Address | Cost | |
| smithj | val | | val | val | | | val | |
| spata | | | | | | | | |
| sxxxx | val | | | val | val | val | | |
| ... | | | | | | | | |
| turnerb | val | val | val | val | val | val | val | |
| ... | | | | | | | | |
| | val | | | | | | | |
| twlstr | val | | val | val | | | val | |
| ... | | | | | | | | |
| zaxx | val | val | val | val | val | val | | |
| zxxx | val | | | | | | val | |

Row Keys: identify the rows in an HBase table.

# Tables are split into Regions = contiguous keys



| RowKey | Address | | | Order | |
|--------|---------|--------|-------|-------|--------|
| | Name | Address | email | Date | ItemNo |
| **smithj** | val | | val | val | |
| **spata** | | | | | |
| | val | | | val | val |
| turnerb | | | | | |
| | val | val | val | val | val |
| tuxedoc | | | | | |
| twistr | val | | | | |
| ... | val | | val | val | |
| xaxx | | | | | |
| xxxx | val | val | val | val | val |

**Region1 Key Range**
- smithj
- sxxx

**Region2 Key Range**
- turnerb
- tuxedoc

**Region3 Key Range**
- twistr
- xaxx

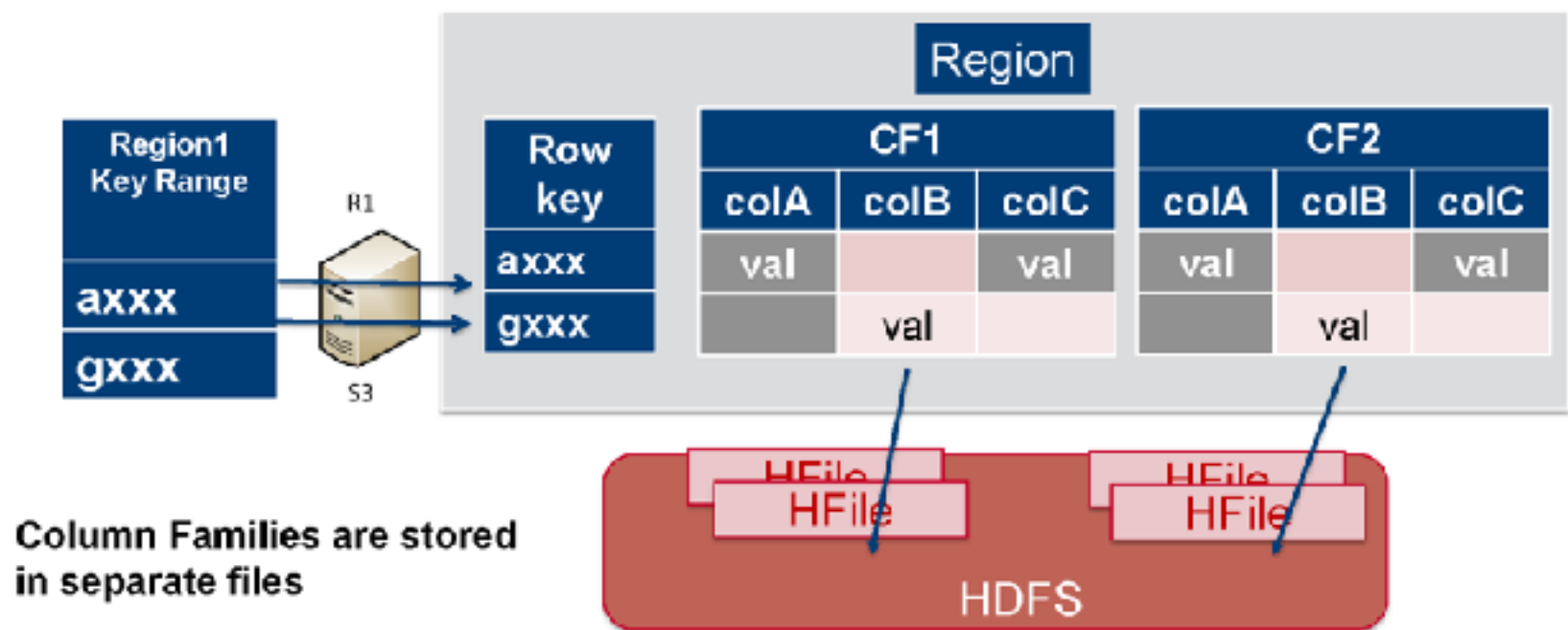R1 — S3

R2 R3 — S1

Tables are partitioned into **key ranges** (regions)
Region= served by nodes (Region Servers)
Regions are spread across cluster

# Column Families are stored Separately

- column families are stored and can be accessed separately



**Column Families are stored in separate files**

- Data is stored in **Key Value** format
- Value for each **cell** is specified by complete coordinates:
    - (Row key, ColumnFamily, Column Qualifier, timestamp ) => Value

Cell Coordinates= Key      Value

| Row key | Column Family | Column Qualifier | Timestamp | Value |
|---------|---------------|------------------|-----------|-------|
| Smithj | Address | city | 1391813876369 | Nashville |

# Logical Data Model vs Physical Data Storage



Logical Model

| RowKey | Address | | Order | |
|--------|---------|------|-------|--------|
| | Street | City | Date | ItemNo |
| smithj | Central Dr | Nashville | 2/2/16 | 10213A |

Key / Value

| Row Key | CF:Col | version | value |
|---------|--------|---------|-------|
| smithj | Address:street | 1 | Central Dr |

Physical Storage

| Row Key | CF:Col | version | value |
|---------|--------|---------|-------|
| smithj | OrerDate | 1 | 2/2/16 |

Physical Storage

© 2015 MapR Technologies   MAPR   26

# Sparse Data with Cell Versions

| | CF1:colA | CF1:colB | CF1:colC |
|---|---|---|---|
| **Row1** | @time7: value3 | | |
| **Row10** | @time2 value1 | @time2: value1 | |
| **Row11** | @time6: value2 | | |
| **Row2** | @time4: value1 | | @time4: value1 |

# Logical Data Model vs Physical Data Storage

| Row Key | CF1 | | CF2 | | Logical Model |
|---------|-----|-----|-----|-----|
| | ca | cb | ca | cd |
| ra | 1 | | 2 | |
| rb | | 3,4 | | |
| rc | 5 | | 6,7 | 8 |

**Column families are stored separately**
**Row keys, Qualifiers** are sorted **lexicographically**

### Physical Storage

**Key** **Value**

| Key | CF1:Col | version | value |
|-----|---------|---------|-------|
| ra | cf1:ca | 1 | 1 |
| rb | cf1:cb | 2 | 4 |
| rb | cf1:cb | 1 | 3 |
| rc | cf1:ca | 1 | 5 |

**Key** **Value**

| Key | CF2:Col | version | value |
|-----|---------|---------|-------|
| ra | cf2:ca | 1 | 2 |
| rc | cf2:ca | 2 | 7 |
| rc | cf2:ca | 1 | 6 |
| rc | cf2:cd | 1 | 8 |

# HBase Table is a Sorted Map of Maps

In Java this is the table:

Table Map of Rows

```
SortedMap<RowKey,
    SortedMap< ColumnFamily,
        SortedMap< ColumnName,
            SortedMap  < version, Value> >>>
```

Map of CF

Map of columns

Map of cells

| Key | CF:Col | version | value |
|---|---|---|---|
| smithj | Address:street | v2 | Main St |
| smithj | Address:street | v1 | Central Dr |
| spata | Address:street | v1 | High Ave |
| turnerb | Address:street | v1 | Cedar St |

| Key | CF:Col | version | value |
|---|---|---|---|
| smithj | Order:Date | v1 | 2/2/15 |
| spata | Order:Date | v1 | 1/31/14 |
| turnerb | OrderDate | v1 | 7/8/14 |

# HBase Table - SortedMap

```
<smithj,<Address, <street, <v1, Central Dr>>
                  <street, <v2, Main St>>
         <Order <Date, <v1, 2/2/15>>>
<spata,<Address, <street, <v1,High Ave>>
              <Order <Date, <v1, 1/31/14>>>
<turnerb,<Address, <street, <v1,Cedar St>>
              <Order <Date, <v1, 7/8/14>>>
```

| Key | CF:Col | version | value |
|---|---|---|---|
| smithj | Address:street | v2 | Main St |
| smithj | Address:street | v1 | Central Dr |
| spata | Address:street | v1 | High Ave |
| turnerb | Address:street | v1 | Cedar St |

| Key | CF:Col | version | value |
|---|---|---|---|
| smithj | Order:Date | v1 | 2/2/15 |
| shawa | Order:Date | v1 | 1/31/14 |
| turnerb | OrderDate | v1 | 7/8/14 |

- **Create Table, define Column Families before data is imported**

  - but not the rows keys or number/names of columns

- Low level API, technically more demanding

- **Basic data access operations (CRUD):**

  put       Inserts data into rows (both create and update)

  get       Accesses data from one row

  scan     Accesses data from a range of rows

  delete   Delete a row or a range of rows or columns

# Create HBase Table – Using HBase Shell

```
hbase> create '/user/user01/Customer', {NAME =>'Address'} , {NAME =>'Order'}

hbase> put '/user/user01/Customer', 'smithj', 'Address:street', 'Central Dr'

hbase> put '/user/user01/Customer', 'smithj', 'Order:Date', '2/2/15'

hbase> put '/user/user01/Customer', 'spata', 'Address:city', 'Columbus'

hbase> put '/user/user01/Customer', 'spata', 'Order:Date', '1/31/14'
```

| Row Key | Address | | | Order | |
|---------|---------|------|-------|------|--------|
| | street | city | state | Date | ItemNo |
| smithj | Central Dr | Nashville | TN | 2/2/15 | 10213A |
| spata | High Ave | Columbus | OH | 1/31/14 | 23401V |
| turnerb | Cedar St | Seattle | WA | 7/8/14 | 10938A |

# Create HBase Table – Using HBase Shell

```
hbase> get '/user/user01/Customer', 'smithj'

hbase> scan '/user/user01/Customer'

hbase> describe '/user/user01/Customer'
```
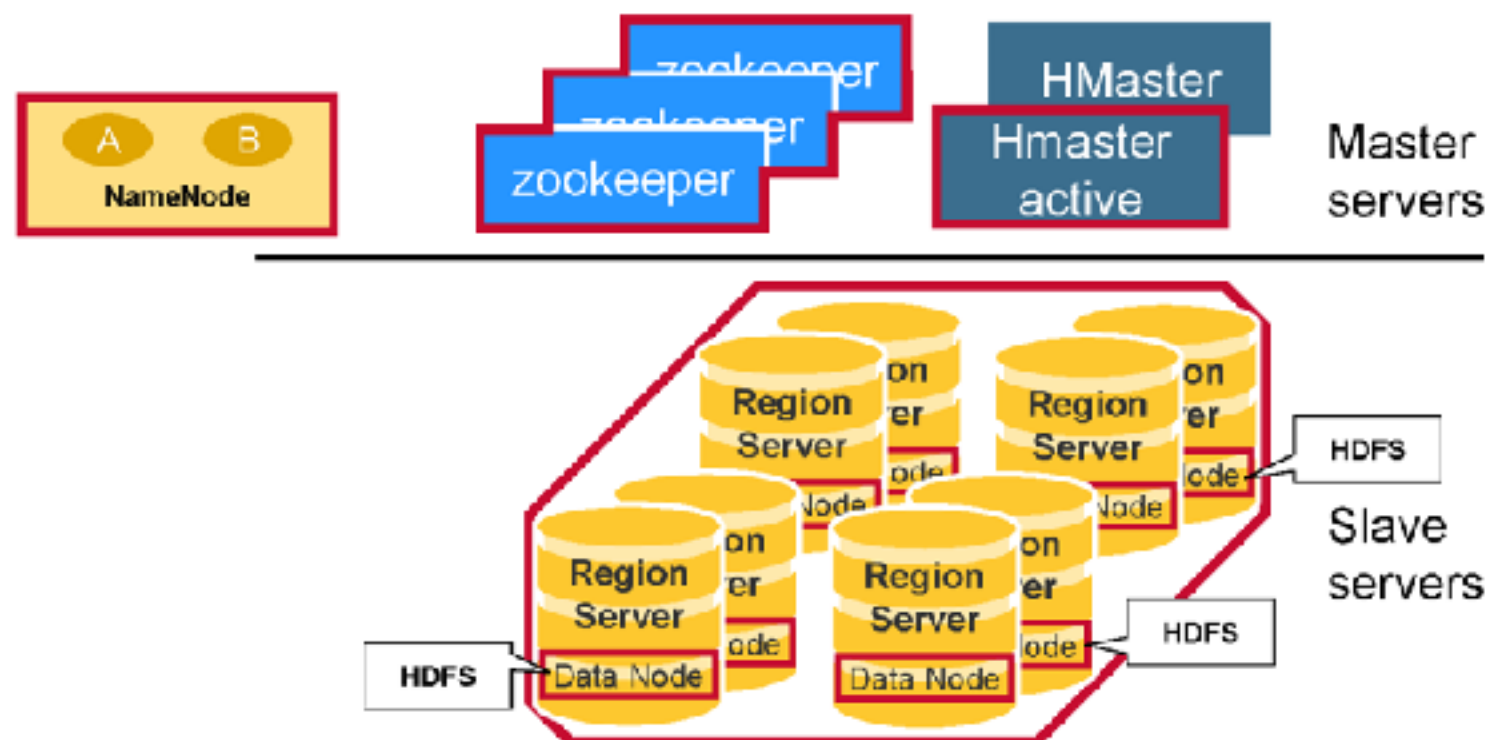
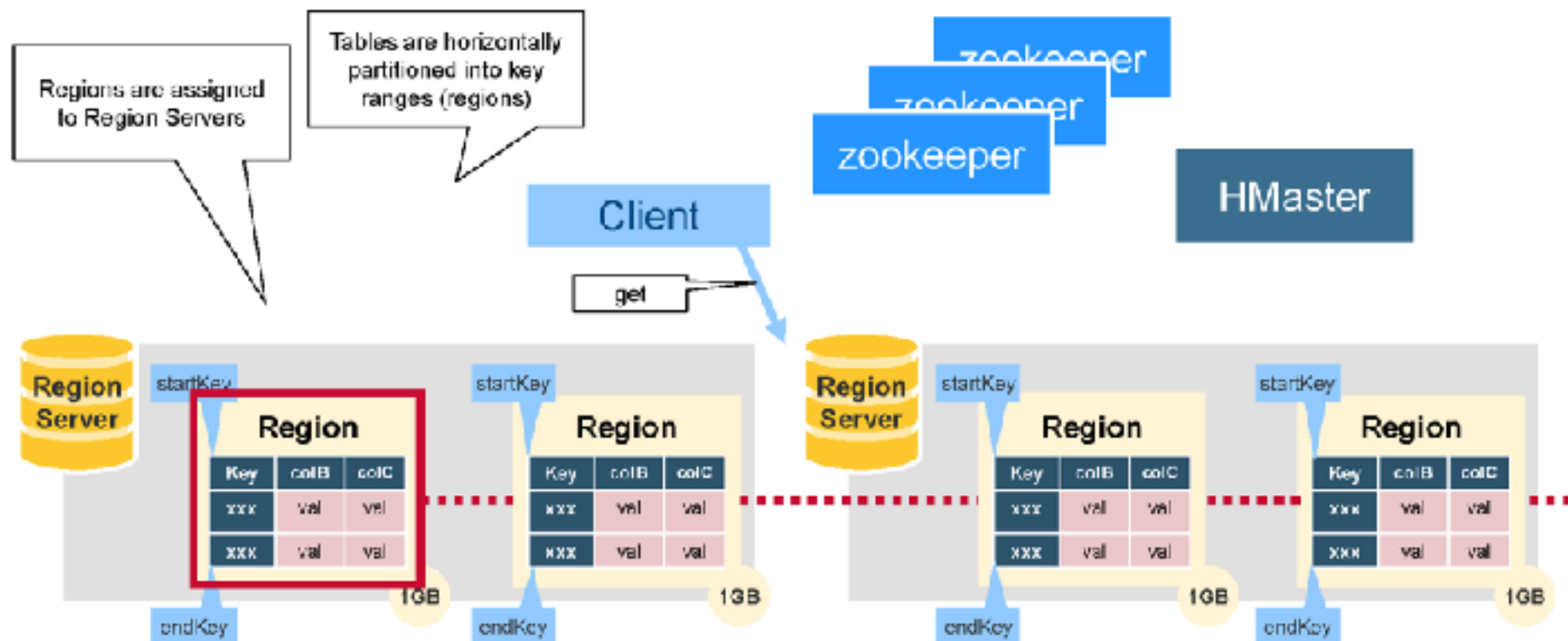# HBase Architecture
Data flow for Writes, Reads
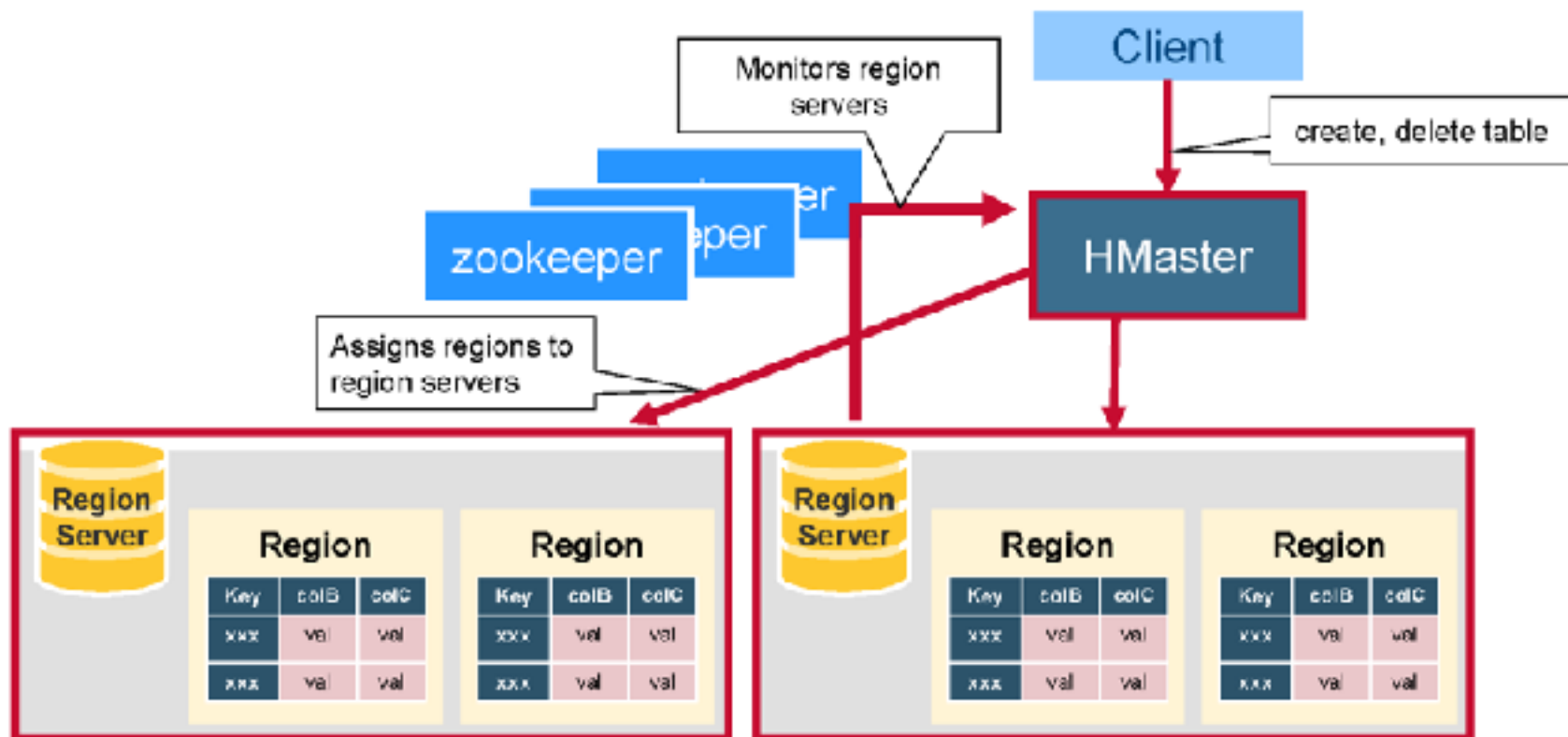Designed to Scale

# HBase Architectural Components



NameNode — A, B

zookeeper
zookeeper
zookeeper

HMaster
Hmaster active

Master servers

Region Server — Data Node
Region Server — Data Node
Region Server — Data Node
Region Server — Data Node

HDFS
HDFS
HDFS

Slave servers

# HBase HMaster

Client

Monitors region servers

create, delete table

zookeeper

HMaster

Assigns regions to region servers

| Region Server | | Region | | | | Region | | |
|---|---|---|---|---|---|---|---|---|
| | | Key | colB | colC | | Key | colB | colC |
| | | xxx | val | val | | xxx | val | val |
| | | xxx | val | val | | xxx | val | val |

| Region Server | | Region | | | | Region | | |
|---|---|---|---|---|---|---|---|---|
| | | Key | colB | colC | | Key | colB | colC |
| | | xxx | val | val | | xxx | val | val |
| | | xxx | val | val | | xxx | val | val |

# How the Components Work Together



- Active HMaster selection
- Region Server session

zookeeper

only 1 master is active

heartbeat

Ephemeral node

Hmaster active

heartbeat

Ephemeral node

**Region Server** — Data Node — Region

| Key | colB | colC |
|-----|------|------|
| xxx | val | val |
| xxx | val | val |

Region

| Key | colB | colC |
|-----|------|------|
| xxx | val | val |
| xxx | val | val |

**Region Server** — Data Node — Region

| Key | colB | colC |
|-----|------|------|
| xxx | val | val |
| xxx | val | val |

Region

| Key | colB | colC |
|-----|------|------|
| xxx | val | val |
| xxx | val | val |

# HBase Meta Table

**META** table

Meta table is used to find the Region for a given Table key

B tree

| Row key | Value |
|---|---|
| table,key,region | region server |

Region Server

Region Server

Region Server

| Region | | |
|---|---|---|
| Key | colB | colC |
| xxx | val | val |
| xxx | val | val |

| Region | | |
|---|---|---|
| Key | colB | colC |
| xxx | val | val |
| xxx | val | val |

| Region | | |
|---|---|---|
| Key | colB | colC |
| xxx | val | val |
| xxx | val | val |

| Region | | |
|---|---|---|
| Key | colB | colC |
| xxx | val | val |
| xxx | val | val |

| Region | | |
|---|---|---|
| Key | colB | colC |
| xxx | val | val |
| xxx | val | val |

| Region | | |
|---|---|---|
| Key | colB | colC |
| xxx | val | val |
| xxx | val | val |

Region Server Components

Region Server

Region

WAL

Memstore | Memstore

HDFS Data Node

Put

Ack

Updates are available to queries after put returns

Next updates are written to the Memstore

# HBase Memstore

Write cache in-memory

Sorted list of Key → Value

One per column family

**Region**

Memstore

Memstore

**Key** ← → **Value**

**Key** **Value**

| Key | CF1:Col | version | value |
|-----|---------|---------|-------|
| ra | cf1:ca | v1 | 1 |
| rb | cf1:cb | v2 | 4 |
| rb | cf1:cb | v1 | 3 |
| rc | cf1:ca | v1 | 5 |

| Key | CF2:Col | version | value |
|-----|---------|---------|-------|
| ra | cf2:ca | v1 | 2 |
| rc | cf2:ca | v2 | 7 |
| rc | cf2:ca | v1 | 6 |
| rc | cf2:cd | v1 | 8 |

# HBase Region Flush

**Region Server**

**Region**

Memstore          Memstore

**FLUSH**

HFile             HFile

**HDFS Data Node**

**WAL**

All memstores in region **flushed** to new HFiles on disk

HFile: sorted list of key → values on disk

Saves the last written sequence number so the system knows what was persisted so far.

# HBase HFile

Flushed quickly

**Sequential write**

Sorted list of Key → Value

HFile      HFile

HDFS Data Node

Moving Disk drive's head to a specific location, is slow.

**Key**    **Value**

| Key | CF1:Col | version | value |
|-----|---------|---------|-------|
| ra | cf1:ca | v1 | 1 |
| rb | cf1:cb | v2 | 4 |
| rb | cf1:cb | v1 | 3 |
| rc | cf1:ca | v1 | 5 |

**Key**    **Value**

| Key | CF2:Col | version | value |
|-----|---------|---------|-------|
| ra | cf2:ca | v1 | 2 |
| rc | cf2:ca | v2 | 7 |
| rc | cf2:ca | v1 | 6 |
| rc | cf2:cd | v1 | 8 |

# HBase HFile Structure

Lookups can be performed with a single disk seek

Hfile Index is kept in Block Cache (Memory)

Index is loaded when the HFile is opened

# HBase Read Merge

**1** First the scanner looks for the Row KeyValues in the Block cache

**2** Next the scanner looks in the Memstore

**3** If all row cells not in memstore or blockCache, look in HFiles

# HBase Read Merge



**Region Server**

**Region**

1 → **BlockCache**

2 → **Memstore**

3 → **HFile**, **HFile**, **HFile**

scanner

**WAL**

**HDFS Data Node**

MemStore creates multiple **small store files** over time when **flushing.**

**Read Amplification**

- multiple files have to be examined

# HBase Minor Compaction

Region Server

Region — Memstore

Region — Memstore

WAL

HDFS Data Node

HFile HFile HFile HFile HFile

HFile HFile HFile HFile HFile

Minor compaction

HFile HFile
HFile HFile
HFile HFile

# HBase Major Compaction

# Region Split



when region size >
hbase.hregion.max.
filesize → split

# Region Load Balancing

# HBase Crash Recovery

**Better than many NoSQL data store solutions, hence its popularity**

- **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