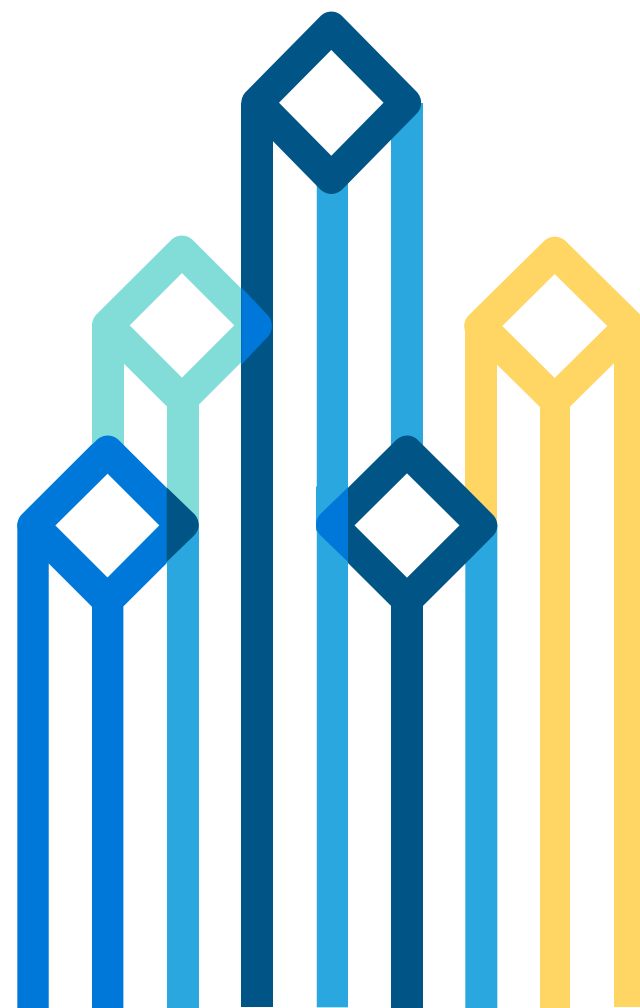




Apache HBase: Overview, Hands-On, and Use Cases

Apekshit Sharma
Esteban Gutierrez



Who we are

Apekshit Sharma

- Distributed Software Engineer, Cloudera
 - Software Engineer, Google
- Apache HBase contributor
- Performance improvements and configuration framework

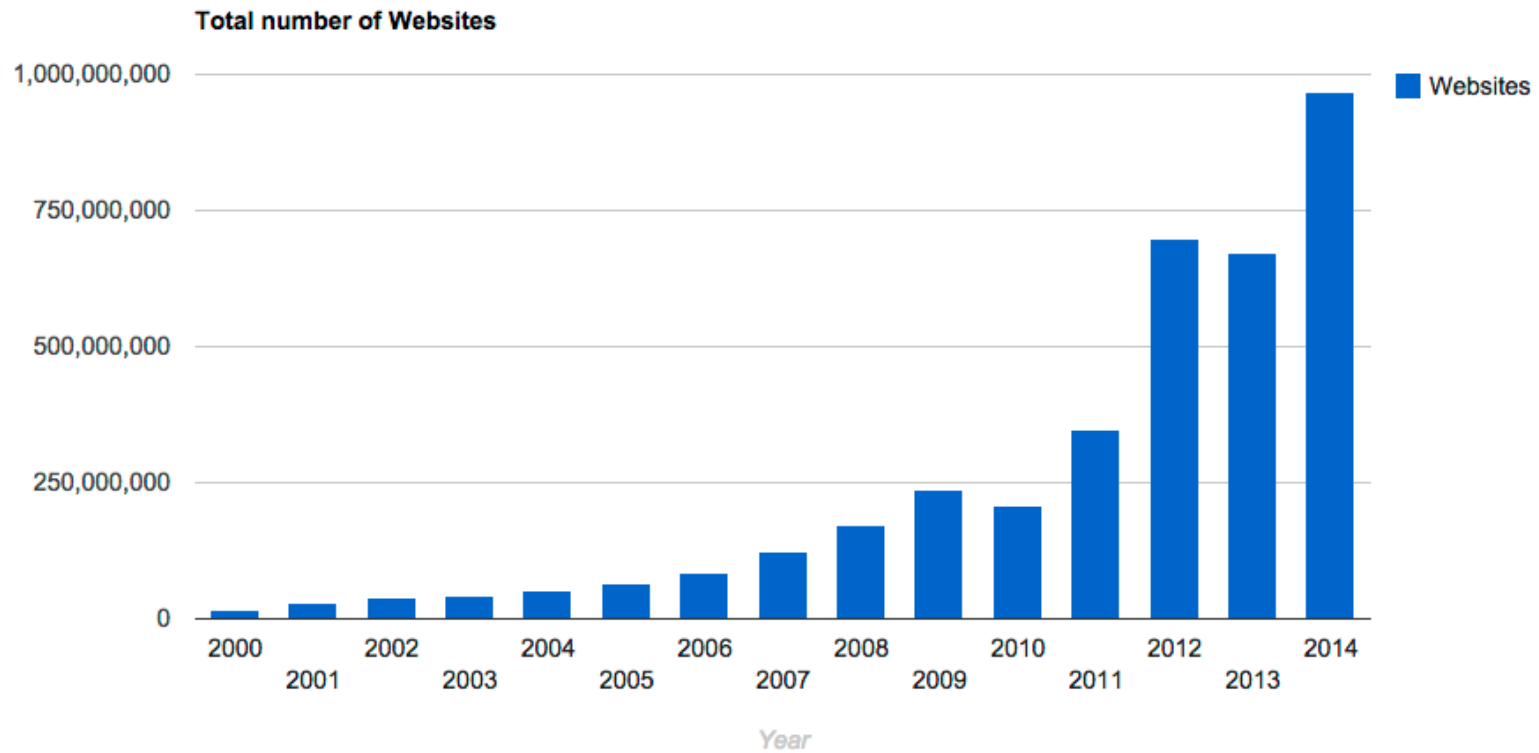
Esteban Gutierrez

- Software Engineer, Cloudera
- Apache HBase committer
- +15 years of experience on Internet scale platforms.

Contents

- Motivation
- Introduction to Apache HBase
 - Data model
- Hands-On: Installation, HBase shell
- A more in-depth introduction to Apache HBase
 - Apache Hadoop HDFS
 - HBase Architecture

Motivation



<http://www.internetlivestats.com/total-number-of-websites/>

Motivation

- Indexing the internet has challenges:
 - Scale
 - Volume
 - Rate
 - Diversity of content
 - URLs
 - High-resolution images
 - Video

Motivation

Bigtable: A Distributed Storage System for Structured Data

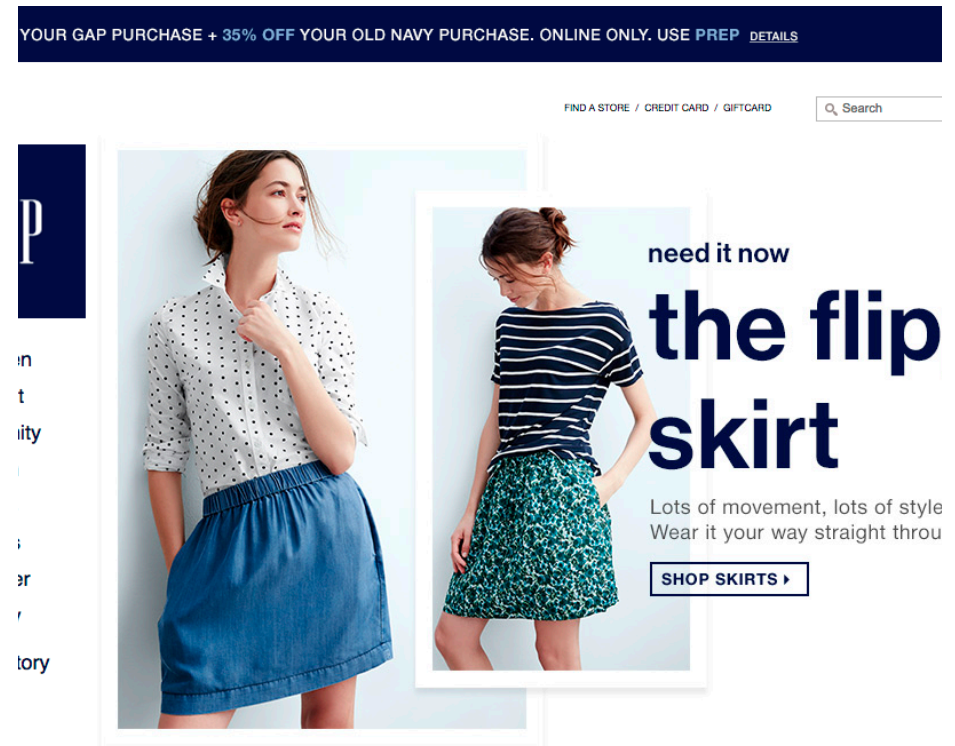
Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach
Mike Burrows, Tushar Chandra, Andrew Fikes, Robert E. Gruber

{fay,jeff,sanjay,wilsonh,kerr,m3b,tushar,fikes,gruber}@google.com

Google, Inc.

Motivation

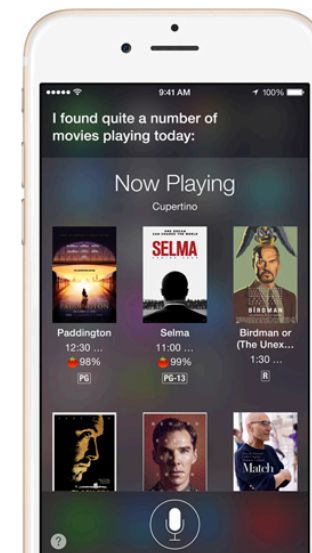
- What if you're not trying to index the internet?



Motivation

- Data for analytical processing
- User-facing real-time platforms

“What movies are playing today?”



Apache HBase



- **Open source**
- **Random access, Low latency and Consistent data store**
- **Horizontally scalable**
- **Built on top of Apache Hadoop**

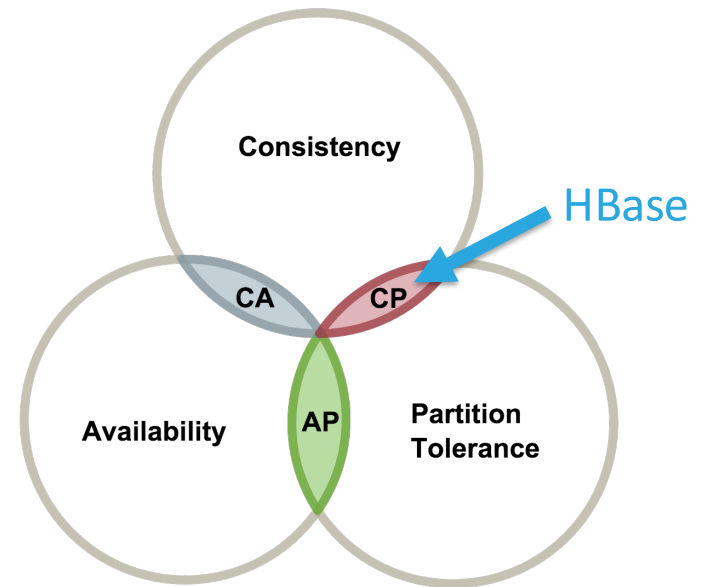
Apache HBase is open source

- [Apache 2.0 License](#)
- A community project with committers and contributors from diverse organizations
 - Cloudera, Facebook, Salesforce.com, Huawei, TrendMicro, eBay, HortonWorks, Intel, Twitter, ...
- Code license means anyone can modify and use the code.

Apache HBase is consistent

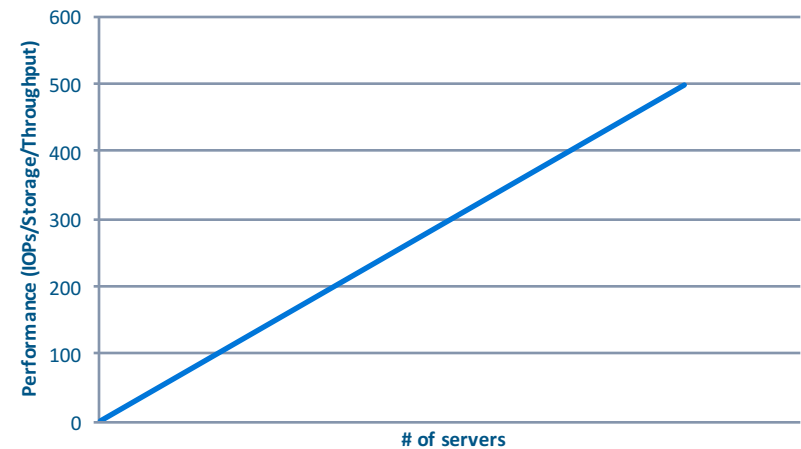
- CAP theorem

- Consistency
- Availability
- Partition tolerance



Apache HBase is horizontally scalable

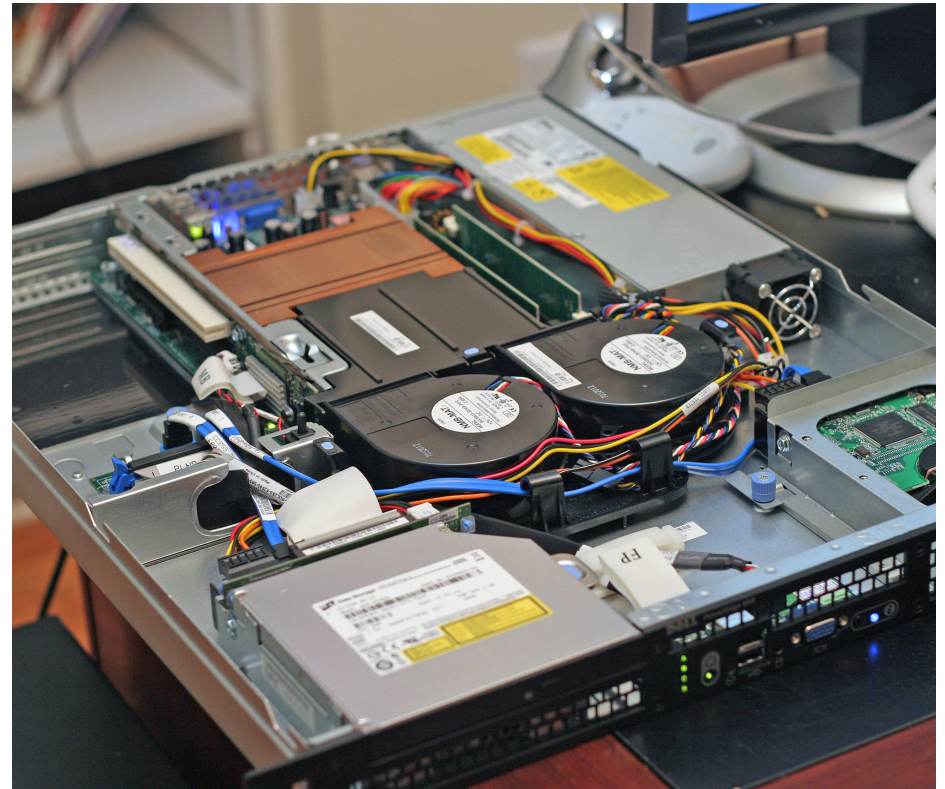
- Adding more servers **linearly** increases performance and capacity
 - Storage capacity
 - Input/output operations
- Store and access data on **commodity servers**



- **Largest cluster:** > 3000 nodes, > 100 PB
- **Average cluster:** 10-40 nodes, 100-400TB

Apache HBase is horizontally scalable

- Commodity servers (circa 2015)
 - 12-24 2-4TB hard disks
 - 2 (quad/hexa/octet)-core CPUs, 2-3 GHz
 - 64 - 512 GBs of RAM
 - 10 Gbps ethernet
 - \$5,000 - \$10,000 / machine



Data model

Columns

<Column family>:<column qualifier>

Row key	info:name	info:age	comp:base	comp:stocks
121	'tom'	'28'	'125k'	
145	'bob'	'32'	'110k'	'50' (ts=2012) '100' (ts=2014)

Data model

- Data is stored... in a big table
- Sorted **rows** with primary **row keys**. Supports billions of rows.
- **Column** : **<column family>:<column qualifier>**; Supports millions of columns
- **Cell** : intersection of row and column.
 - Can be empty. No storage/processing overheads
 - Can have different time-stamped values

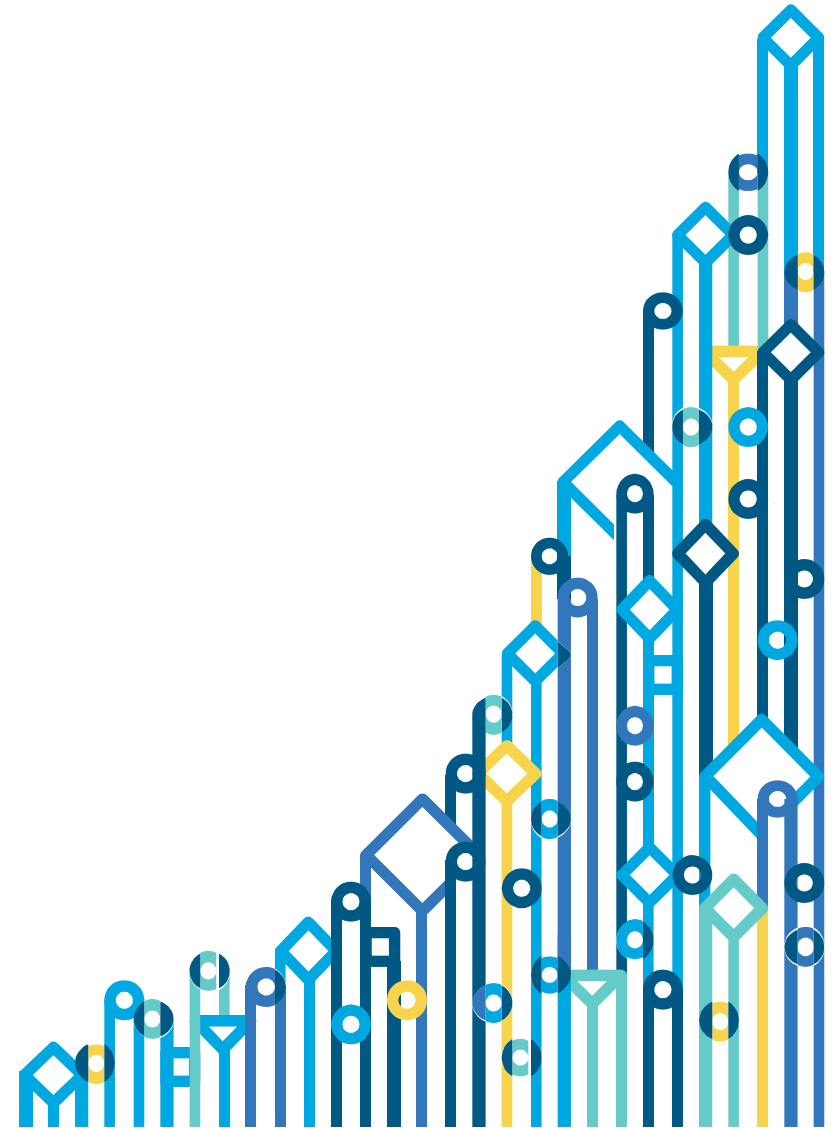


Hands On

Apache HBase installation

Play with HBase shell

<http://pastebin.com/60EUN75f>



Understanding basic architecture of



Hadoop (HDFS)

APACHE
HBASE

Apache Hadoop



Open source

Commodity servers

Horizontally scalable

Highly fault-tolerant

Massive processing power

2 Core Components



Apache Hadoop



History

The Google File System

Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung

Google*

2003

GFS

- **Distributed file system**
- **Commodity servers**
- **horizontally scalable**
- **Highly fault-tolerant**
- **Proprietary**

HDFS

- Distributed file system
- Commodity servers
- Horizontally scalable
- Highly fault-tolerant
- Open source

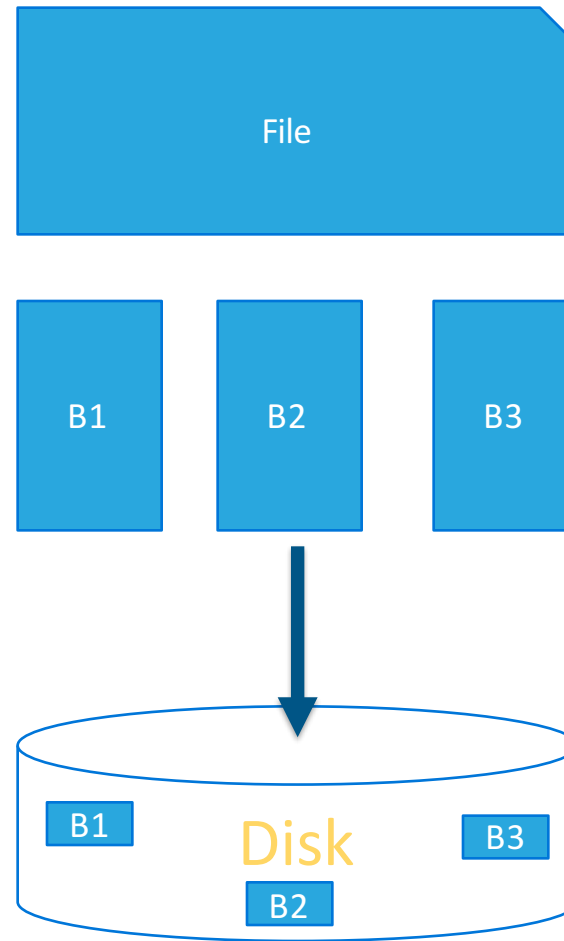
HDFS API

- File
 - Open, Close, Read, Write, Move, etc
- Directories
 - Create, Delete, etc
- Permissions
 - Owners, Groups, rwx permissions

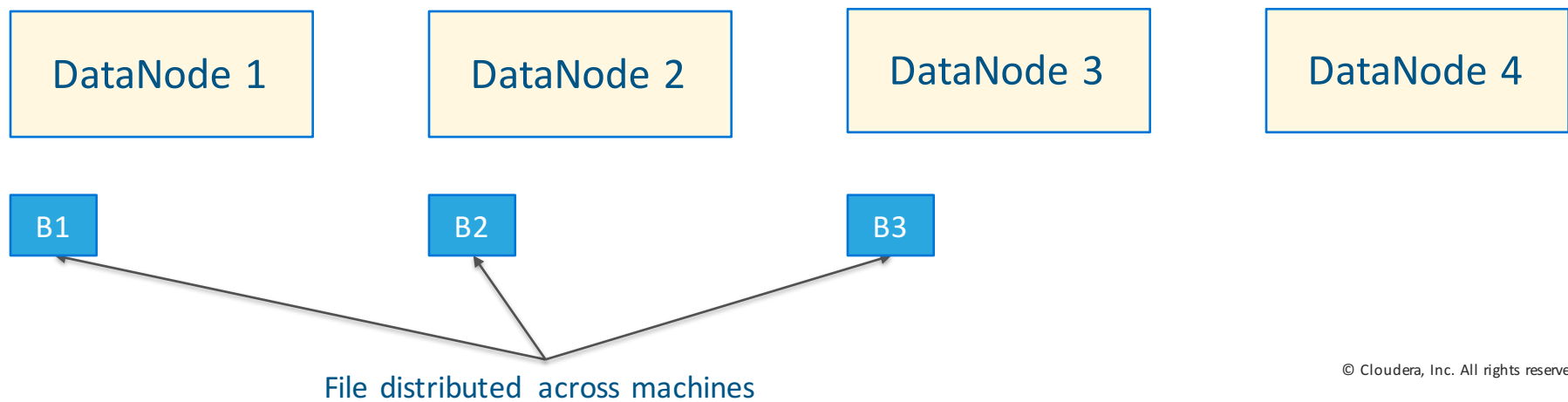
Basic Architecture of HDFS

Local file system

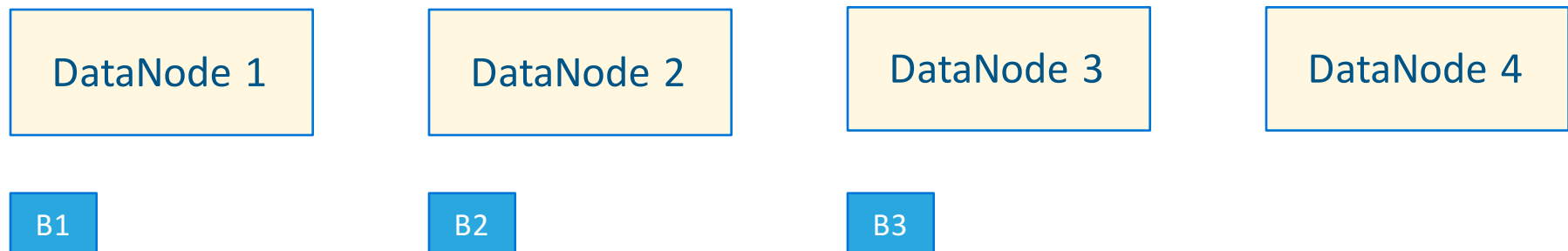
File system will split
the file into blocks



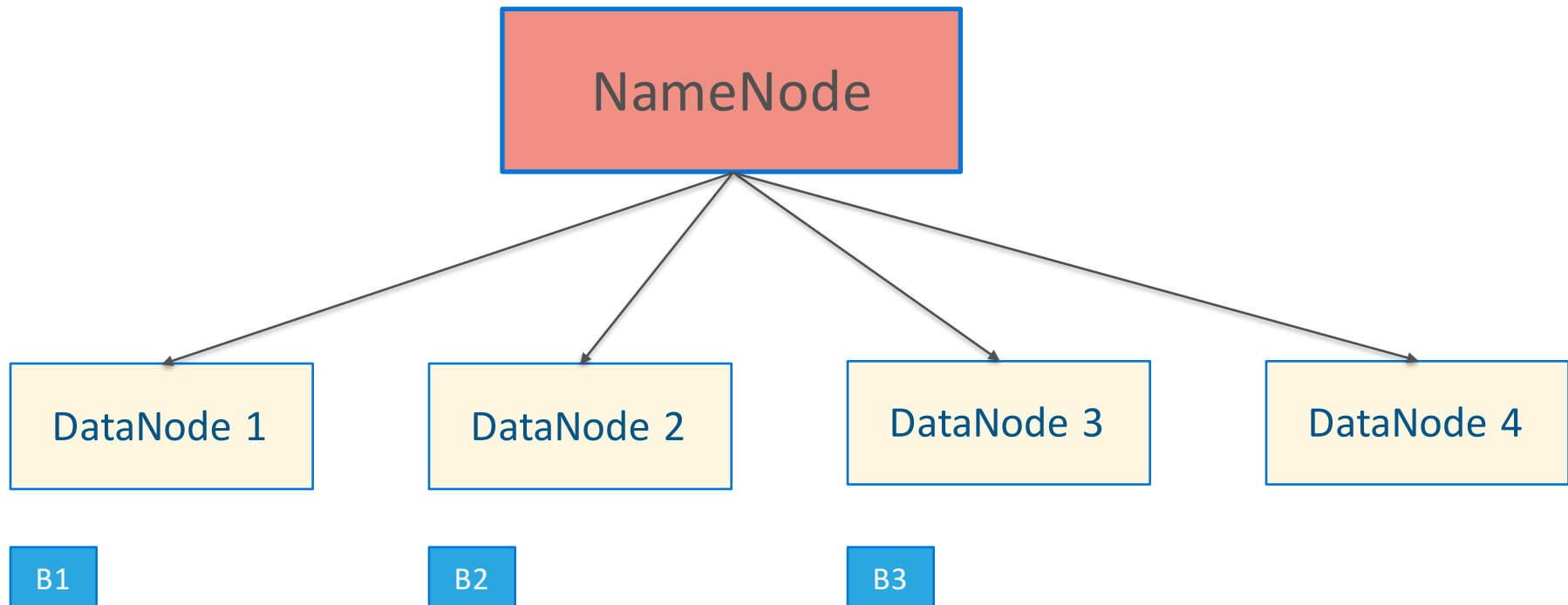
HDFS



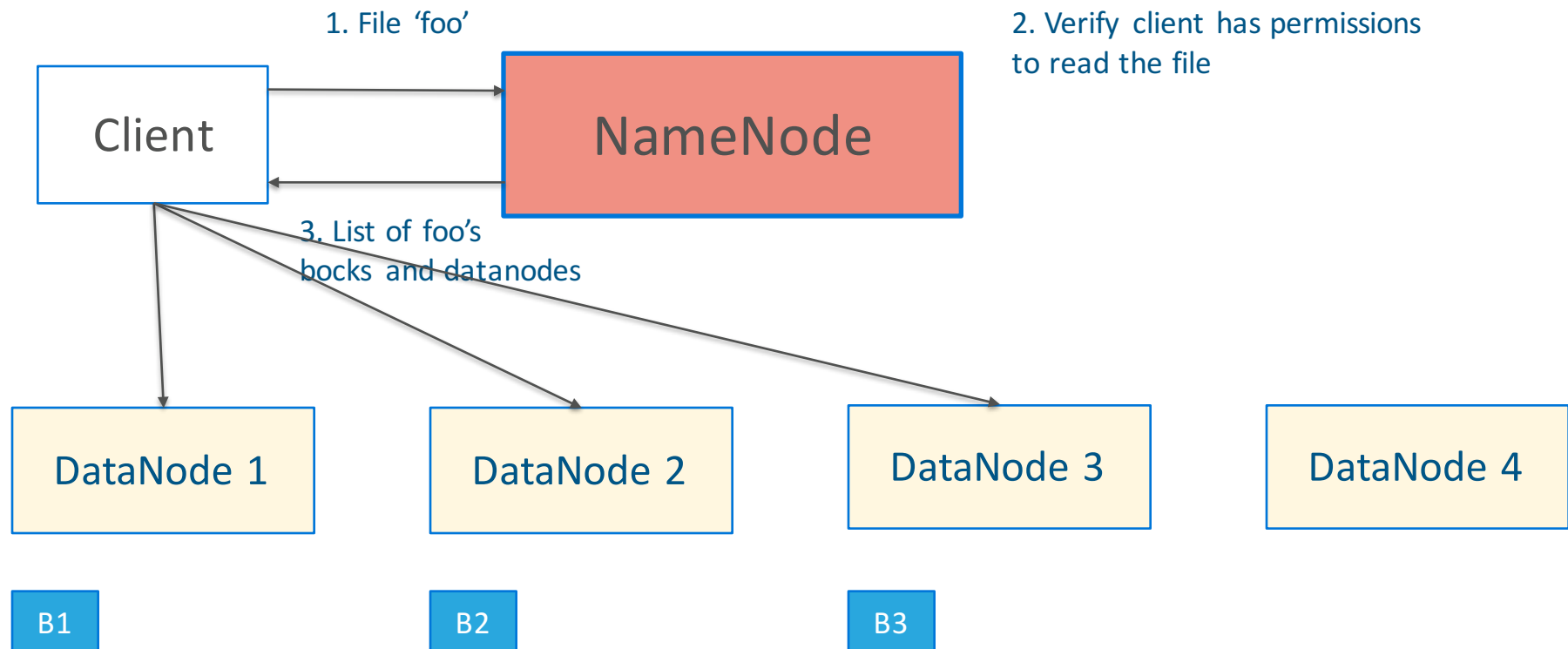
HDFS DataNode



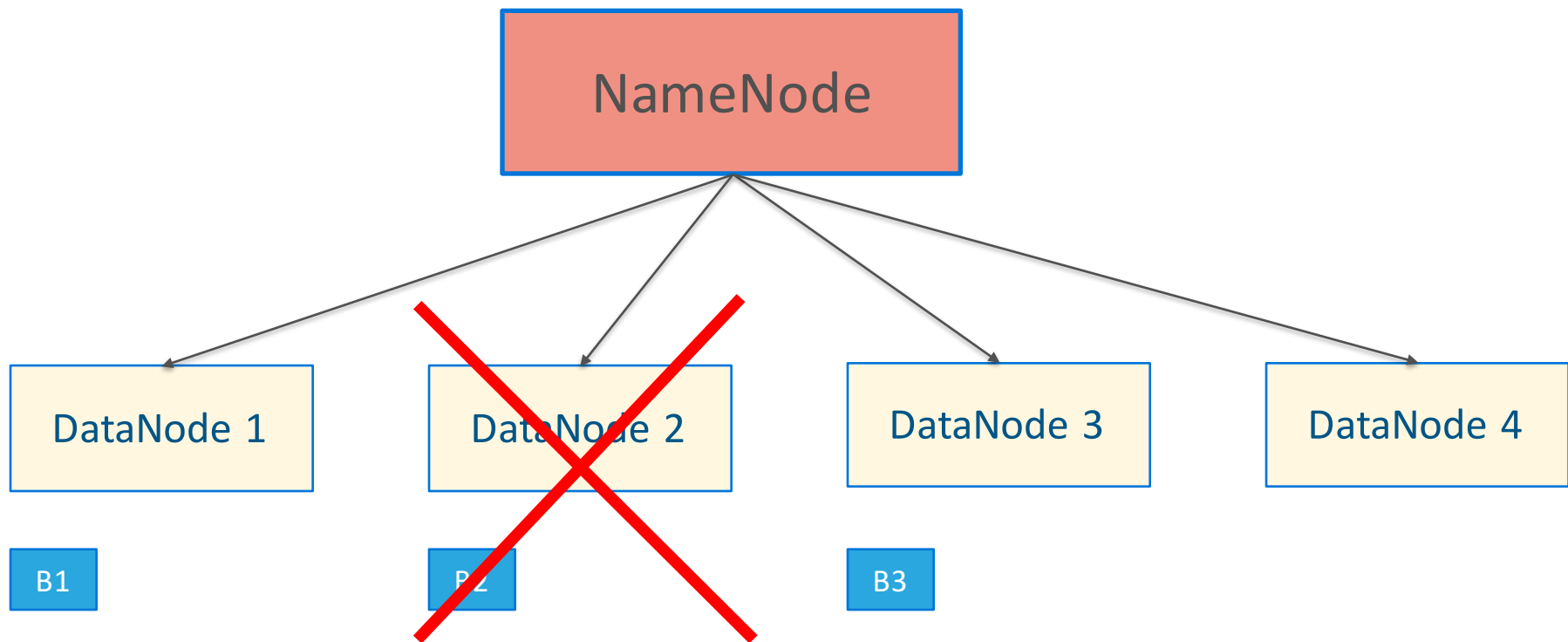
HDFS NameNode



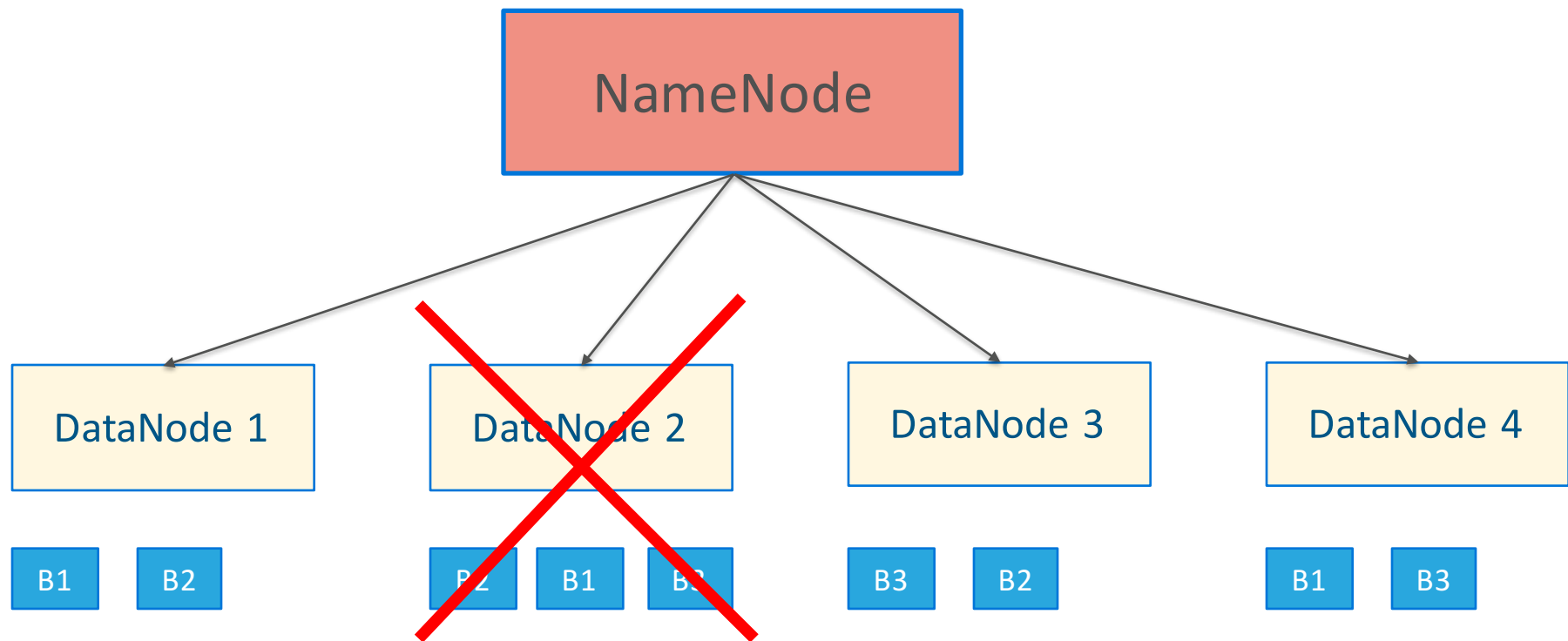
HDFS Reading a file



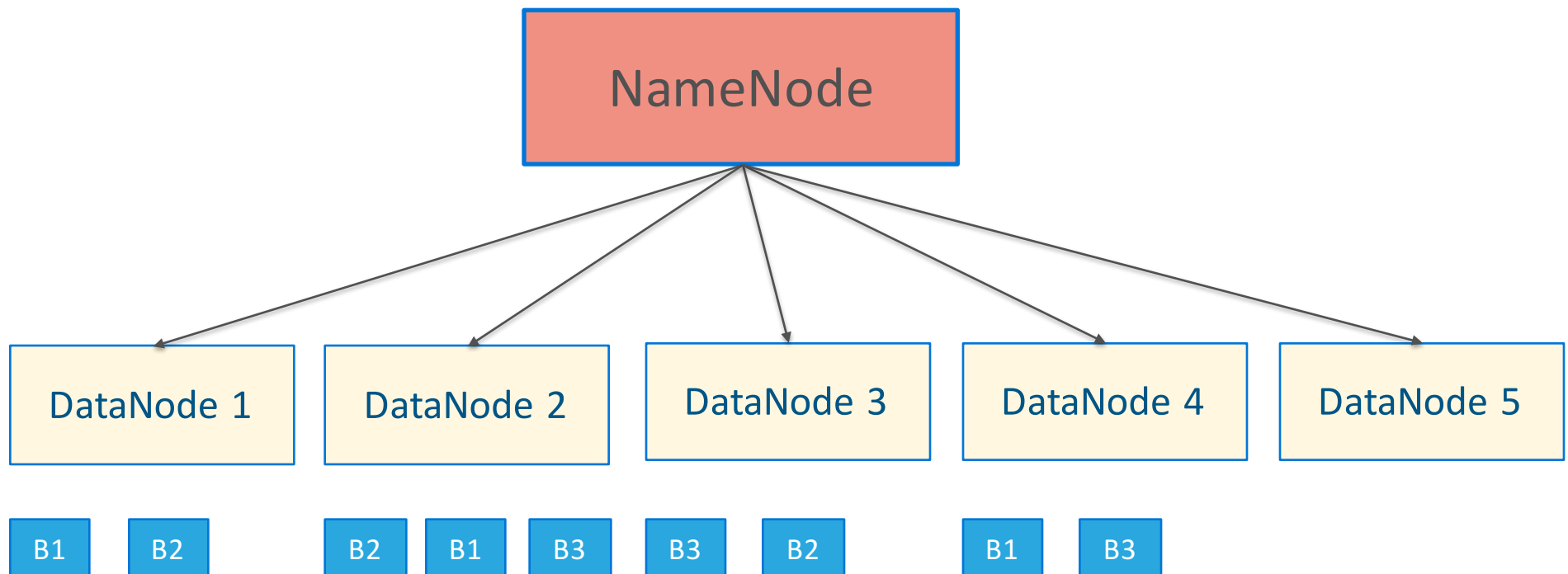
HDFS Fault tolerance



HDFS Redundancy



HDFS Horizontal Scalability



Let's look at some existing HDFS systems...

- **Yahoo! HDFS Clusters**

40k+ servers, 100k+ CPUs, 450PB data

- Spam filtering, web indexing, personalization, etc.

- **Facebook HDFS Cluster**

15TB new data per day

1200+ machines, 30PB in one cluster

- Datawarehouse, Messages, etc.

- **1-10PB HDFS Clusters**

Production critical workloads for websites, retail, finance, telcom, government, research ...

- **Sub PB HDFS Clusters**

PoCs, R&D, production in early stages ...

But.... there are restrictions!

It's not a magic wand!

Files are **append only**

- Access Model : Write-once-read-many
- Can not change existing contents

Not designed for small files

- HDFS block sizes are in MB (default 128MB)
 - Designed for typical GBs / TBs of file sizes
 - Normal files system have 4kb block size!
 - Name node becomes bottleneck because metadata grows too big

Summary

HDFS is a great distributed file system!

- Store massive data
- Scalable
- High throughput
- Fault tolerance

MapReduce

- Distributed processing framework
- Commodity machines
- Fault tolerance

MapReduce

- HBase uses and extends MR APIs for data access

Further reading: <http://hadoop.apache.org>



HBase Architecture

Unique id	Name	price	weight	store1	store2	store3
"1000000"	snickers	\$9.99	4 Oz	Yes	Yes	Yes
"3000000"	almonds	\$9.99	8 Oz	Yes	No	Yes
"8000000"	coke	\$9.99	16 Oz	Yes	Yes	Yes

HBase Architecture

Unique id	Name	price	weight	store1	store2	store3
"1000000"	snickers	\$9.99	4 Oz	Yes	Yes	Yes
"3000000"	almonds	\$9.99	8 Oz	Yes	No	Yes
"8000000"	coke	\$9.99	16 Oz	Yes	Yes	Yes
"4000000"	new	\$34.63	16 Oz	No	Yes	Yes

HBase Architecture

Unique id	Name	price	weight	store1	store2	store3
"1000000"	snickers	\$9.99	4 Oz	Yes	Yes	Yes
"3000000"	almonds	\$9.99	8 Oz	Yes	No	Yes
"8000000"	coke	\$9.99	16 Oz	Yes	Yes	Yes
"4000000"	foo	\$34.63	16 Oz	No	Yes	Yes
"5000000"	bar	\$22.54	16 Oz	Yes	Yes	Yes

HBase Architecture

Unique id	Name	price	weight	store1	store2	store3
"1000000"	snickers	\$9.99	4 Oz	Yes	Yes	Yes
"3000000"	almonds	\$9.99	8 Oz	Yes	No	Yes
"8000000"	coke	\$9.99	16 Oz	Yes	Yes	Yes
"4000000"	foo	\$34.63	16 Oz	No	Yes	Yes
"5000000"	bar	\$22.54	16 Oz	Yes	Yes	Yes
"9000000"	new1	\$2.5	16 Oz	Yes	Yes	Yes
"7000000"	new2	\$6.4	16 Oz	Yes	Yes	Yes
"2000000"	new3	\$6.4	16 Oz	Yes	Yes	Yes

HBase Architecture | Regions

["", "5000000")

Row Key	Name	brand	price	weight	store1	store2	store3
"1000000"	snickers	xxx	\$9.99	4 Oz	Yes	Yes	Yes
"2000000"	new3	xxx	\$6.4	16 Oz	Yes	Yes	Yes
"3000000"	almonds	xxx	\$9.99	8 Oz	Yes	No	Yes
"4000000"	foo	xxx	\$34.63	16 Oz	No	Yes	Yes

["5000000", "")

Row Key	Name	brand	price	weight	store1	store2	store3
"5000000"	bar	xxx	\$22.54	16 Oz	Yes	Yes	Yes
"7000000"	new2	xxx	\$6.4	16 Oz	Yes	Yes	Yes
"8000000"	coke	xxx	\$9.99	16 Oz	Yes	Yes	Yes
"9000000"	new1	xxx	\$2.5	16 Oz	Yes	Yes	Yes

HBase Architecture | RegionServer

Row Key	Name	price	weight
"1000000"	snickers	\$9.99	4 Oz
"2000000"	new3	\$6.4	16 Oz
"3000000"	almonds	\$9.99	8 Oz
"4000000"	foo	\$34.63	16 Oz



Server 12

Row Key	Name	price	weight
"5000000"	bar	\$22.54	16 Oz
"7000000"	new2	\$6.4	16 Oz
"8000000"	coke	\$9.99	16 Oz
"9000000"	new1	\$2.5	16 Oz



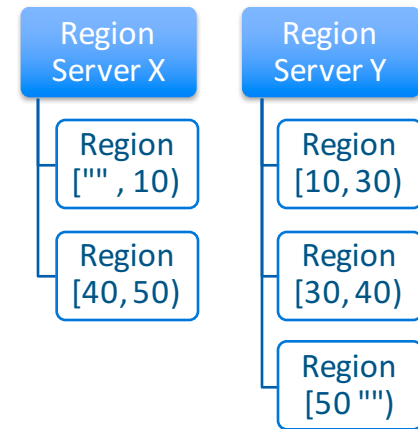
Server 7

HBase Architecture | Regions

- Horizontal split of tables
- Regions are served by RegionServers
- Automatically (based on configurations). Can be done manually too.

Relationship:

- A Region is only served by a single RegionServer at a time.
- A RegionServer can serve multiple Regions.



HBase Architecture

Row Key	Name	price	weight	store1	store2	store3
"1000000"	snickers	\$9.99	4 Oz	Yes	Yes	Yes
"2000000"	new3	\$6.4	16 Oz	Yes	Yes	Yes
"3000000"	almonds	\$9.99	8 Oz	Yes	No	Yes
"4000000"	foo	\$34.63	16 Oz	No	Yes	Yes
"5000000"	bar	\$22.54	16 Oz	Yes	Yes	Yes
"7000000"	new2	\$6.4	16 Oz	Yes	Yes	Yes
"8000000"	coke	\$9.99	16 Oz	Yes	Yes	Yes
"9000000"	new1	\$2.5	16 Oz	Yes	Yes	Yes

What if you have many more columns?

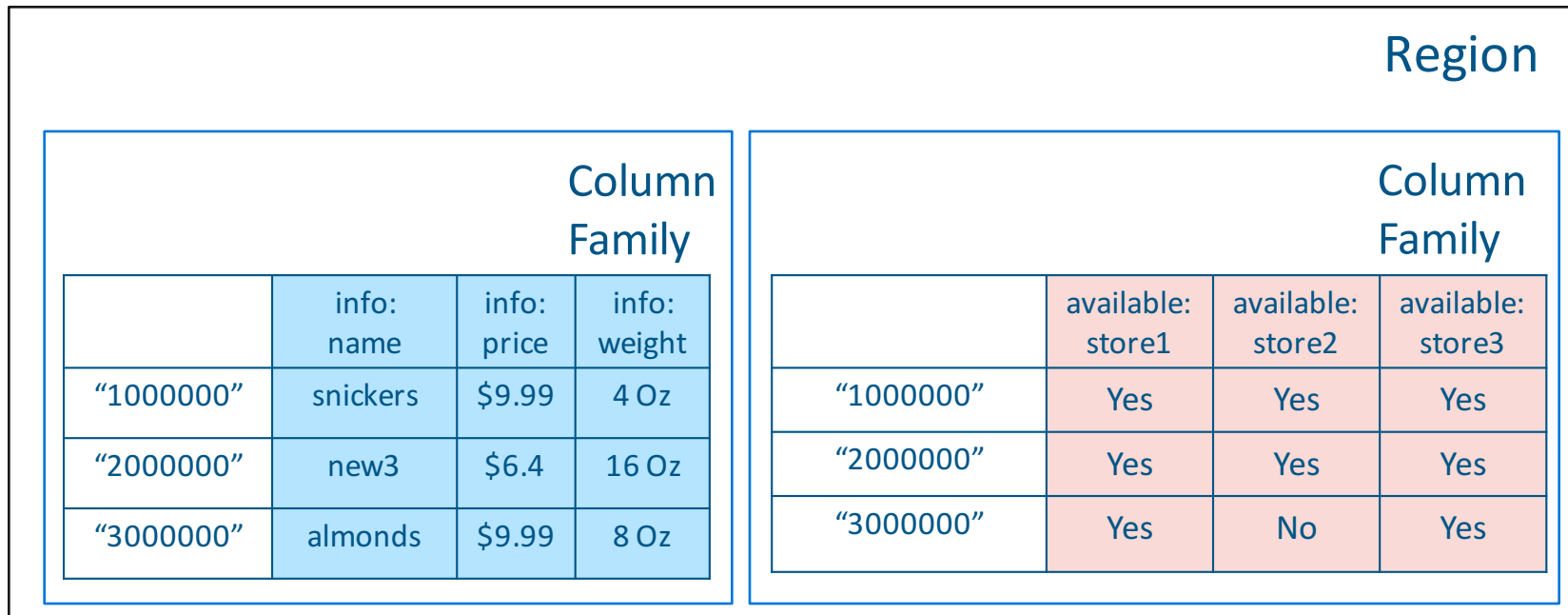
HBase Architecture | Column family

Row Key	info: name	info: price	info: weight	availability: store1	availability: store2	availability: store3
"1000000"	snickers	\$9.99	4 Oz	Yes	Yes	Yes
"2000000"	new3	\$6.4	16 Oz	Yes	Yes	Yes
"3000000"	almonds	\$9.99	8 Oz	Yes	No	Yes
"4000000"	foo	\$34.63	16 Oz	No	Yes	Yes
"5000000"	bar	\$22.54	16 Oz	Yes	Yes	Yes
"7000000"	new2	\$6.4	16 Oz	Yes	Yes	Yes
"8000000"	coke	\$9.99	16 Oz	Yes	Yes	Yes
"9000000"	new1	\$2.5	16 Oz	Yes	Yes	Yes

HBase Architecture | Column family

- Set of columns that have similar access patterns
- Data stored in separate files
- Tune performance
 - In-memory
 - Compression
 - Version retention policies
 - Cache priority
- Needs to be specified by the user

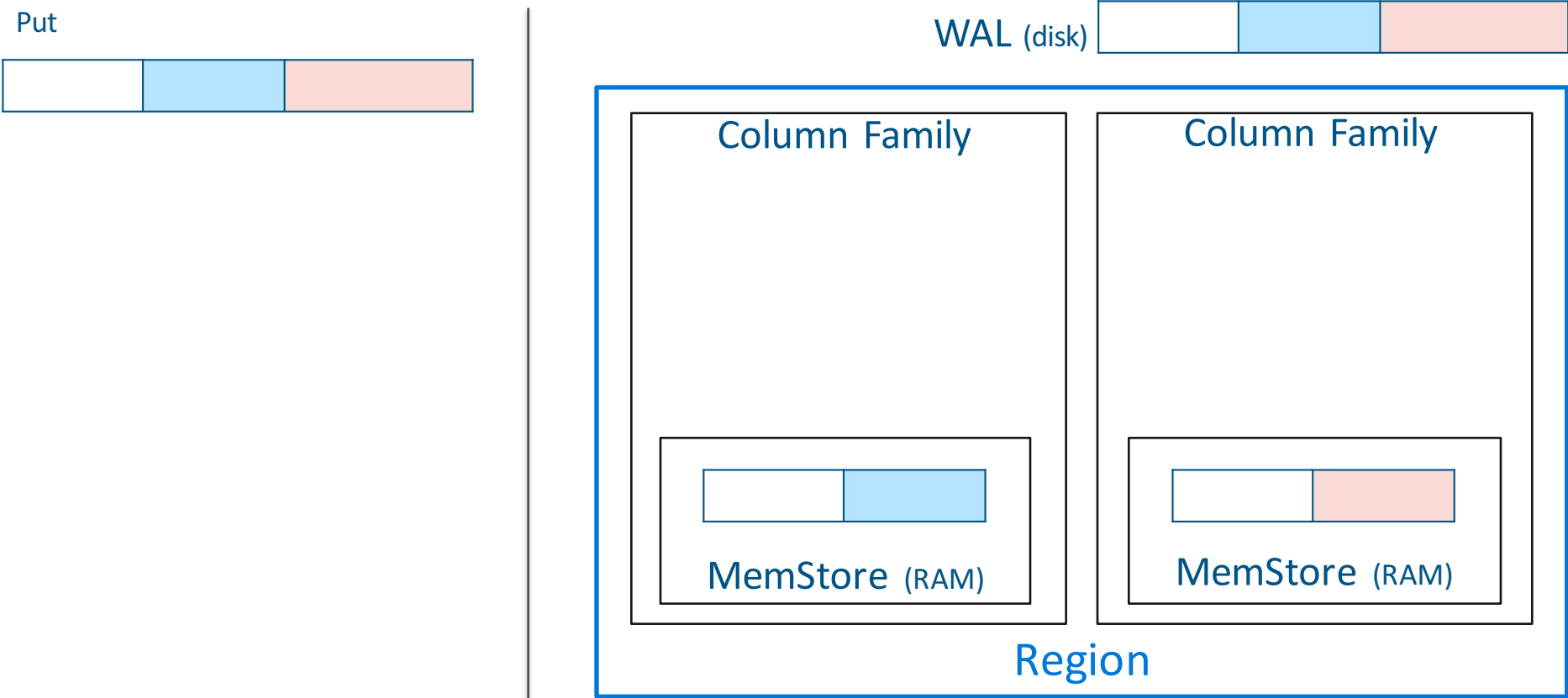
HBase Architecture



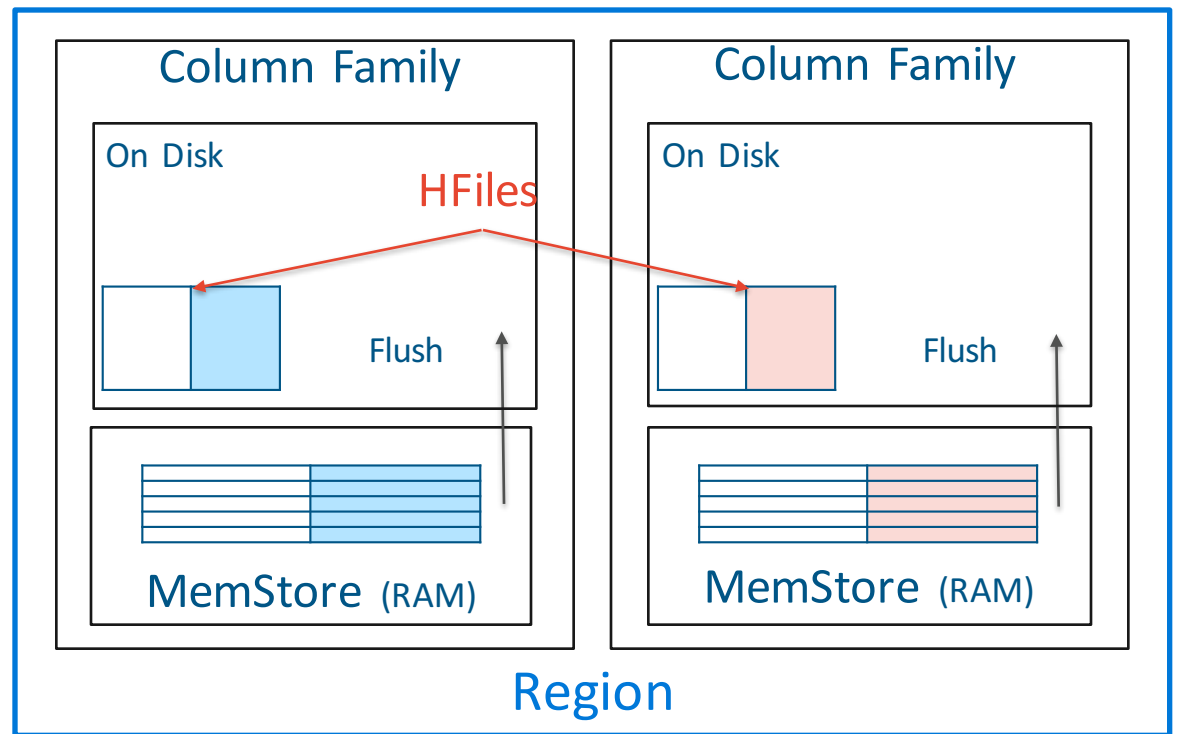
HBase Architecture | Put

1. Client creates a row to put.
2. Client checks which RegionServer hosts this row
3. Row is written into write-ahead log (WAL) on disk for persistence
4. Row is written to MemStore (in memory storage)

HBase Architecture | Put



HBase Architecture | HFile

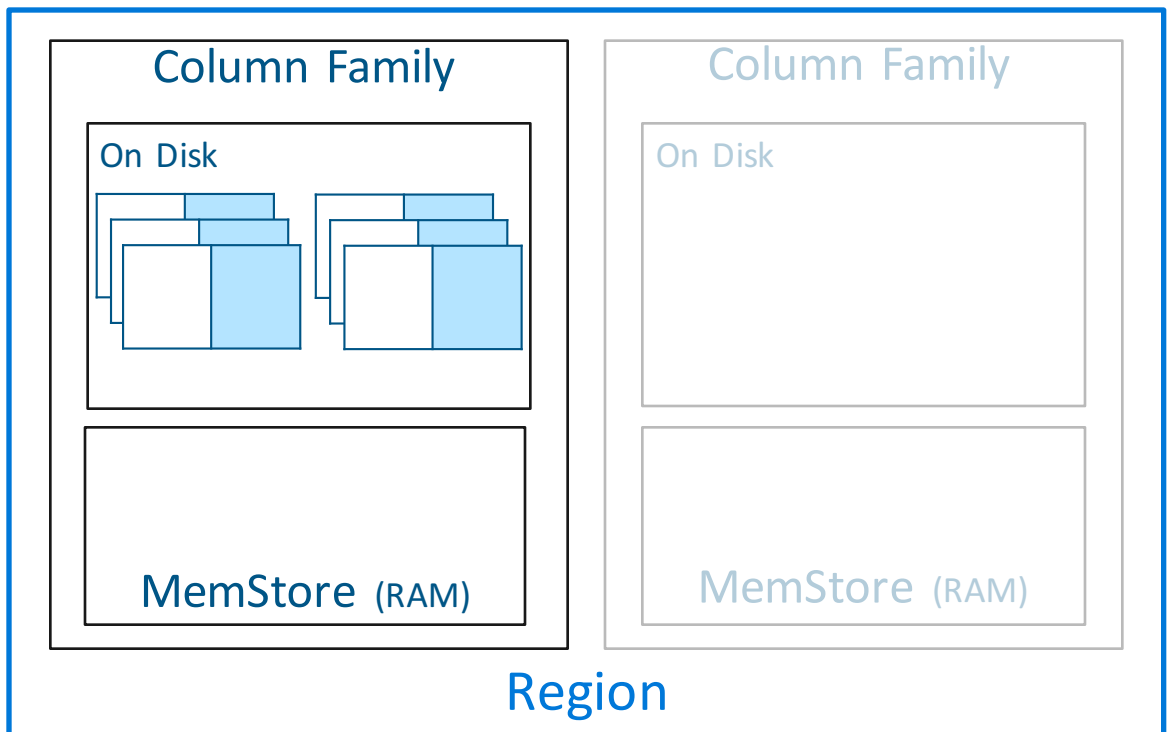


HBase Architecture | HFile

- As more rows are added to the table
 - Write to WAL
 - Write to MemStore
 - When MemStore gets full, its contents are flushed to HDFS as [Hfiles](#)

HBase Architecture | HFile

- Side effect of more HFiles in a Region
 - More disk seeks on read
 - More disk usage (overwritten value, deletes, etc)

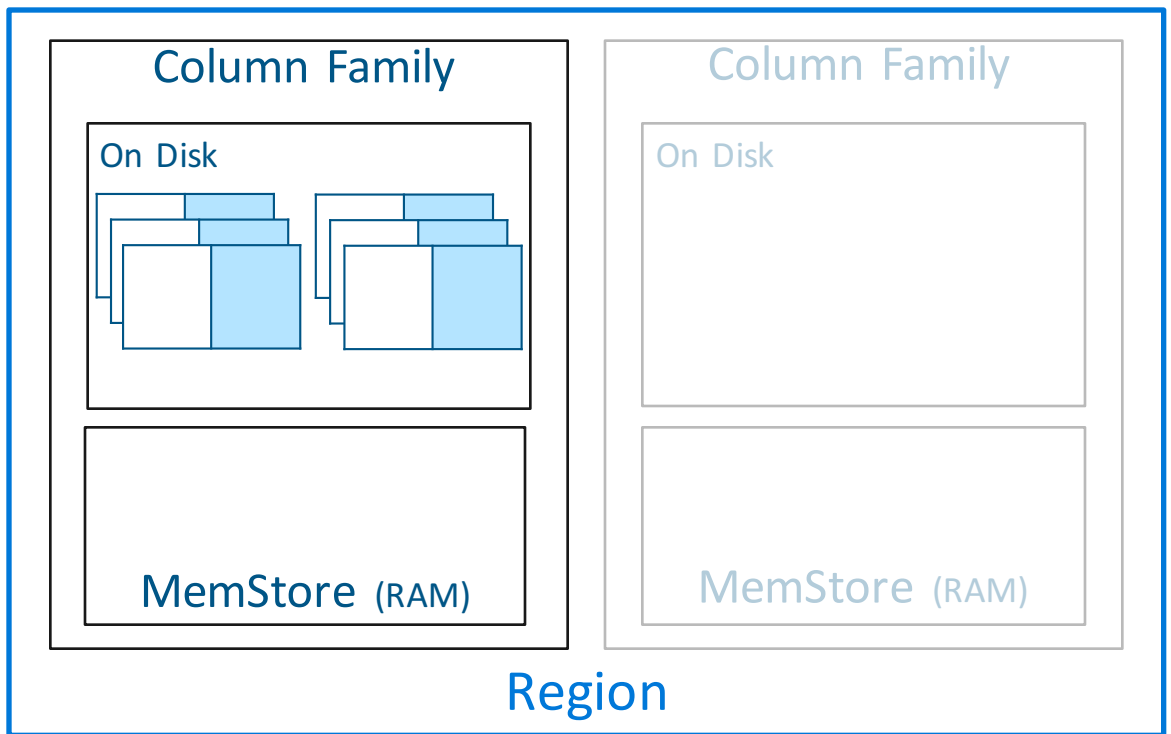


HBase Architecture | Compactions

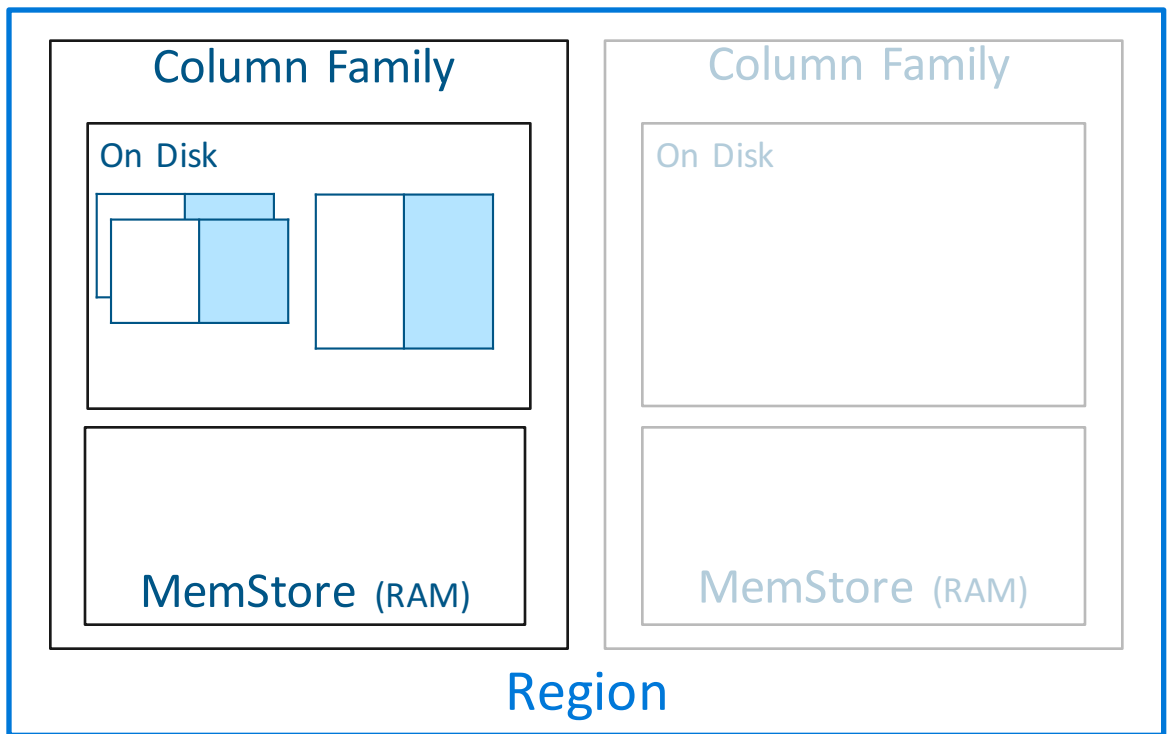
Merging multiple Hfiles into a single Hfile.

- Minor compactions
- Major compactions

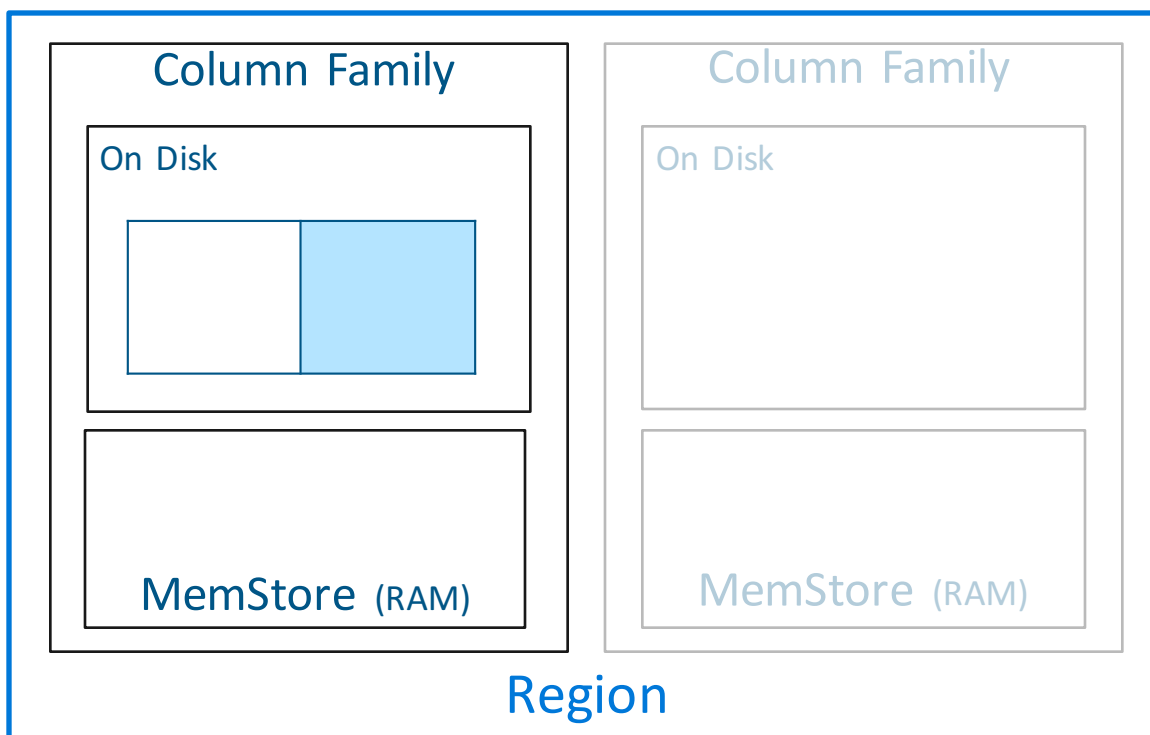
HBase Architecture



HBase Architecture



HBase Architecture



HBase Architecture | Compactions

- Minor compactions
 - Can be finely tuned using multiple configurations
 - Controlled by policy (pluggable)
- Major compactions
 - Periodic (set from configuration) or manually triggered.
 - Tend to be run during off-peak times.

HBase Architecture | Splits

- Eventually, Regions become imbalanced.
- HBase can split a Region into two.
- Daughter Regions are can be moved to a different RegionServer, if necessary.

Gives HBase horizontal scalability!

HBase Architecture | Splits

- Region has rows [0, 100)
- Splitting
 - Middle of current Region's row range
 - All column families are split

10		
20		
30		
40		
50		
60		
70		
80		
90		
100		

Region

RegionServer 2

HBase Architecture | Splits

- Splits
 - [0, 50)
 - [50, 100)

10		
20		
30		
40		
50		
60		
70		
80		
90		
100		

Region

RegionServer 2

HBase Architecture | Bulk Import

- This write path is durable, but if you're importing a lot of data
 - Every put goes into WAL, which means disk seeks. Lots of puts mean lots of disk seeks.
 - Lots of data into MemStores means lots of flushing to disk.
 - Lots of flushing to disk might mean lots of compactions.

So we have other better options like....

HBase Architecture | Bulk Loading

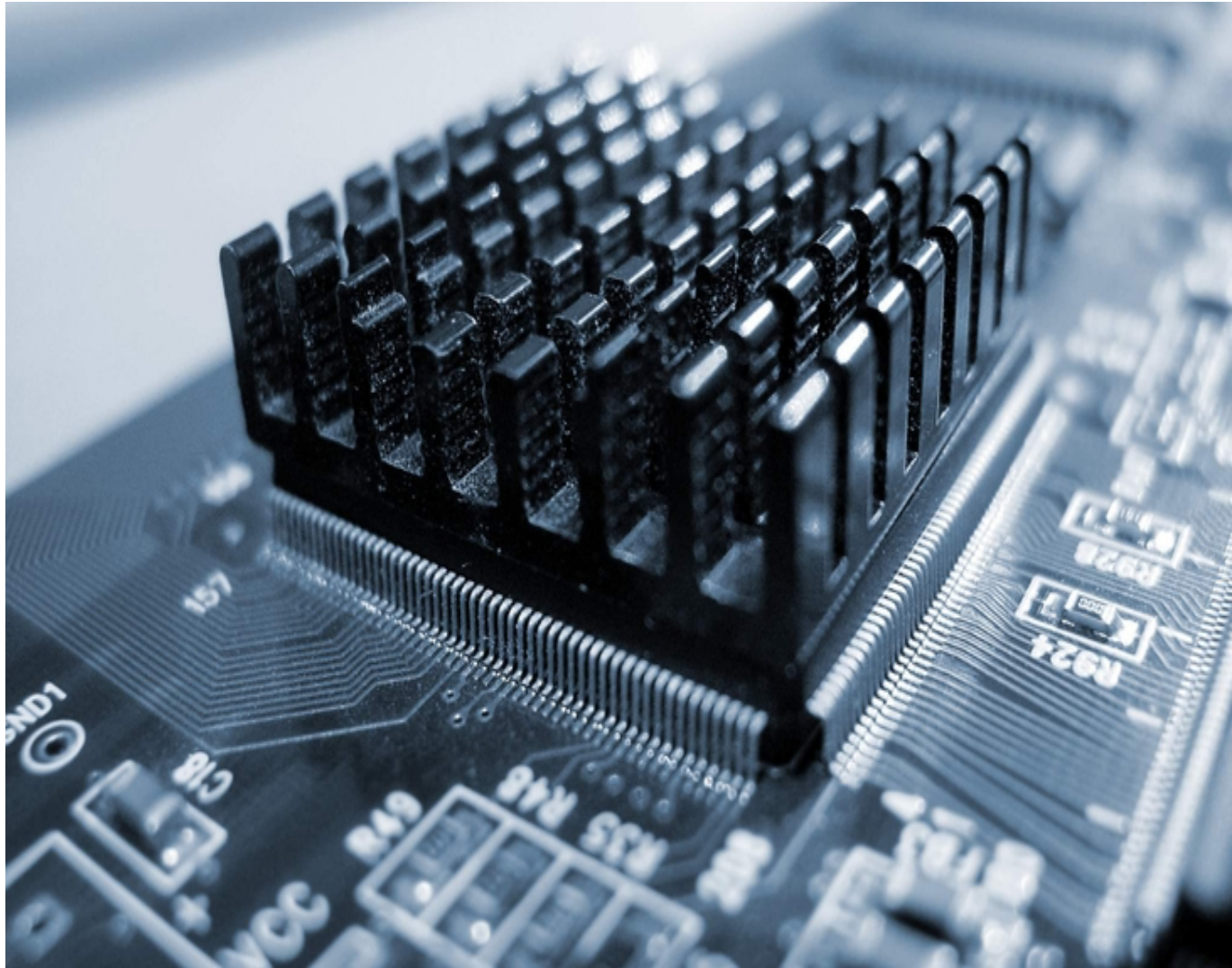
- Bypass conventional write path
 - Extract data from source
 - Transform data into HFiles (done with MapReduce job) directly
 - Tell RegionServers to serve these HFiles

HBase Architecture | APIs

- Java API
 - Most full-featured.
- REST API
 - Easily accessible.
- Thrift API
 - Support for many languages (e.g. C, C++, Perl, Ruby, Python).

Enough of Architecture

cloudera



That's all folks!



Try Hadoop Now
cloudera.com/live

cloudera **LIVE**

Hello, Cloudera Customers and Users!

These community forums are intended for developers and admins using Cloudera's Apache Hadoop-based platform to build and manage big data. We welcome your suggestions and feedback [here](#).

[Join this community](#) to get a 40% discount for [O'Reilly Media](#) print books, and 50% for e-books and videos (bundles not included) -- as well as

To participate in upstream open source projects, use their respective upstream mailing lists.

Ask a Question

Type your question here....

Continue

Community

News (2 Items)

Title	Posts
 Community Guidelines & News Latest Post - This community is now mobile-friendly	5
 Release Announcements Latest Post - Announcing: New Cloudera ODBC drivers for Impala a...	40

Get community
help or provide
feedback

cloudera.com/community



Sources

- A Survey of HBase Application Archetypes
 - Lars George, Jon Hsieh
 - <http://www.slideshare.net/HBaseCon/case-studies-session-7>
- Hadoop and HBase: Motivations, Use cases and Trade-offs
 - Jon Hsieh's slides



cloudera

Questions ?