

It's all about SCALE!!



Warning: fopen(/home/dodgers/public_html/./logs/oracle_error_log.txt) [function.fopen]: failed to open stream: Permission denied in /usr/local/apache/htdocs/include2007/oracle/db_oracle.inc.php on line 194

Cannot open Database Error Log, please check!! (/home/dodgers/public_html/./logs/oracle_error_log.txt)

Warning: fopen(/home/dodgers/public_html/./logs/oracle_error_log.txt) [function.fopen]: failed to open stream: Permission denied in /usr/local/apache/htdocs/include2007/oracle/db_oracle.inc.php on line 194

Cannot open Database Error Log, please check!! (/home/dodgers/public_html/./logs/oracle_error_log.txt)

Warning: fopen(/home/dodgers/public_html/./logs/oracle_error_log.txt) [function.fopen]: failed to open stream: Permission denied in /usr/local/apache/htdocs/include2007/oracle/db_oracle.inc.php on line 194

Cannot open Database Error Log, please check!! (/home/dodgers/public_html/./logs/oracle_error_log.txt)

Warning: fopen(/home/dodgers/public_html/./logs/oracle_error_log.txt) [function.fopen]: failed to open stream: Permission denied in /usr/local/apache/htdocs/include2007/oracle/db_oracle.inc.php on line 194

Cannot open Database Error Log, please check!! (/home/dodgers/public_html/./logs/oracle_error_log.txt)



訂購歷史紀錄

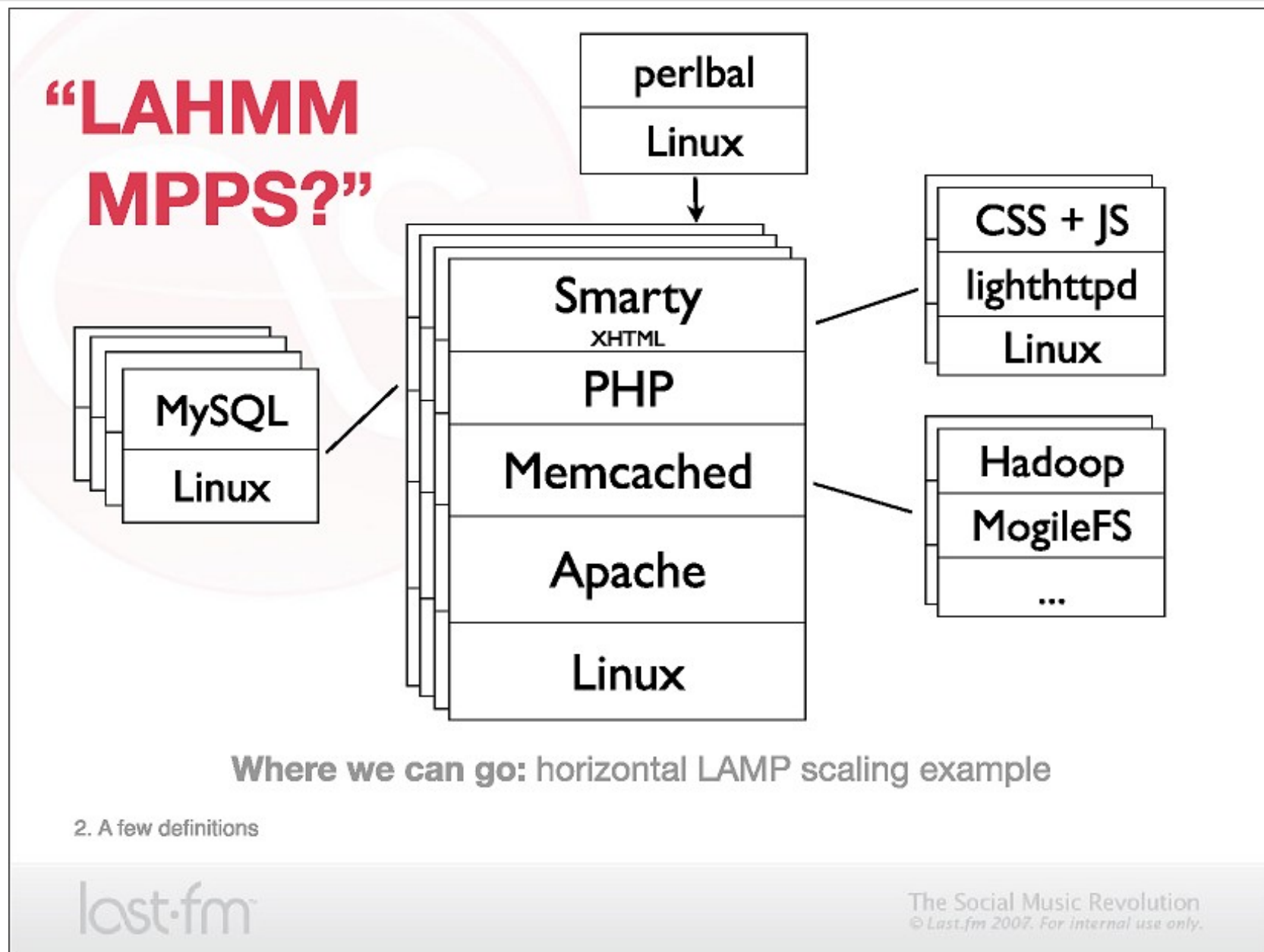


denied in /usr/local/apache/htdocs/include2007/oracle/db_oracle.inc.php on line 194

Cannot open Database Error Log, please check!! (/home/dodgers/public_html/./logs/oracle_error_log.txt)

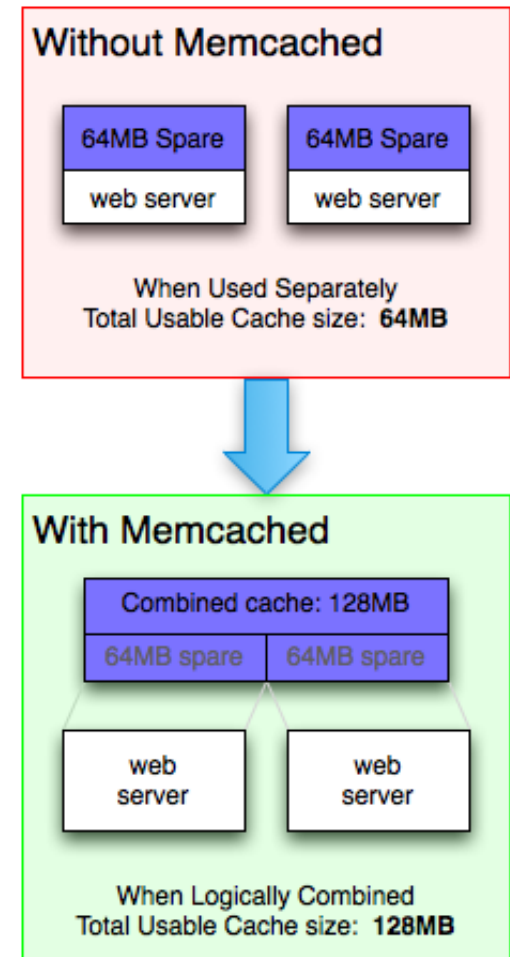
Warning: fopen(/home/dodgers/public_html/./logs/oracle_error_log.txt) [function.fopen]: failed to open stream: Permission

How to scale up web service in the past ?



Tools used by large scale websites

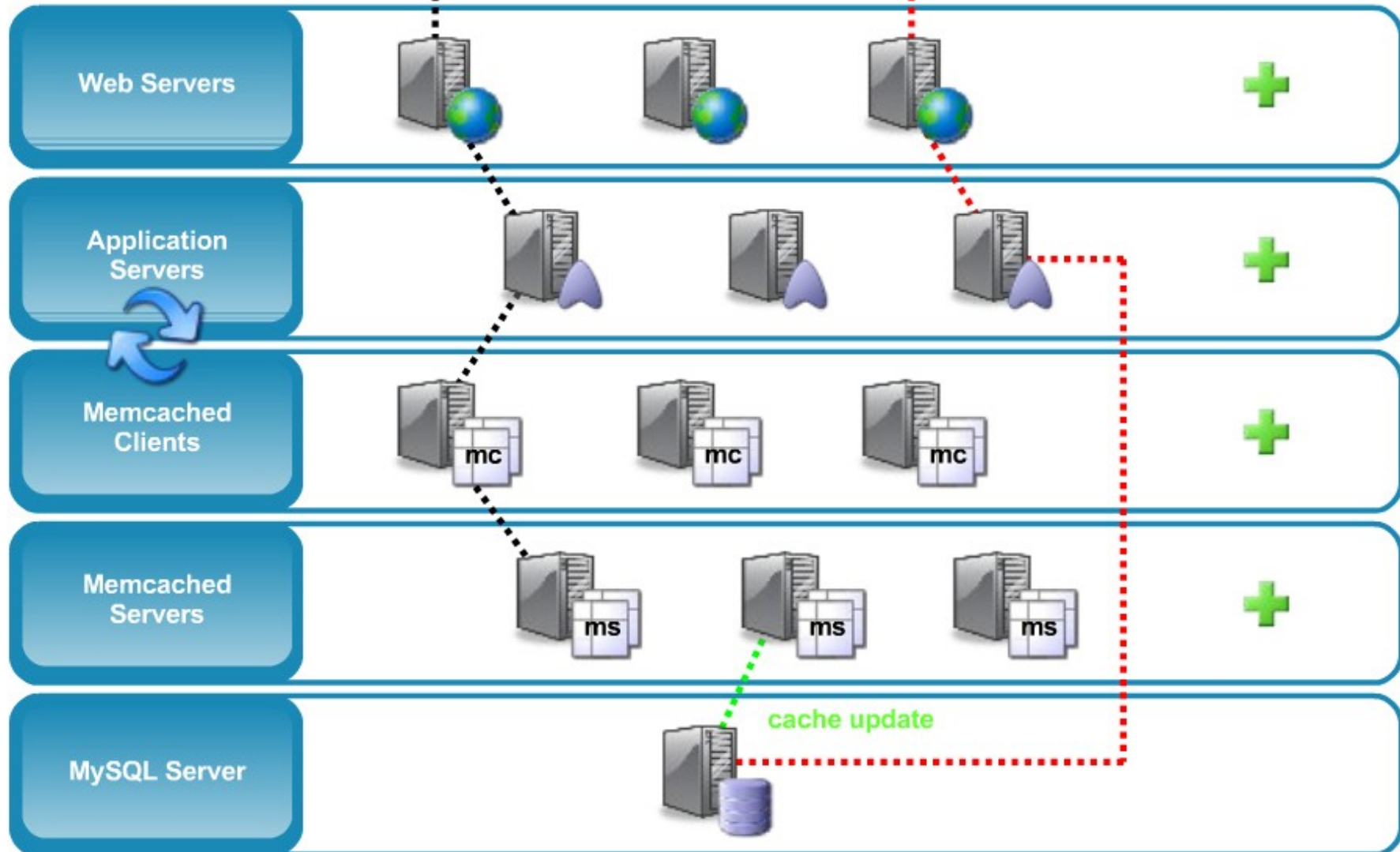
- Perlbal - <http://www.danga.com/perlbal/>
 - ◆ 多個網頁伺服器的負載平衡
 - ◆ Load balancer
- MogileFS - <http://www.danga.com/mogilefs/>
 - ◆ 分散式檔案系統
 - ◆ Distributed File System fo small files
 - ◆ 有公司認為 MogileFS 比起 Hadoop 適合拿來處理小檔案
- memcached - <http://memcached.org/>
 - ◆ 共享記憶體 ??
 - ◆ Share Memory
 - ◆ 把資料庫或經常讀取的部分，
用記憶體快取 (Cache) 方式存放
- Moxi - <http://code.google.com/p/moxi/>
 - ◆ Memcache 的 PROXY
- More Resource:
 - ◆ <http://code.google.com/p/memcached/wiki/HowToLearnMoreScalability>
 - ◆ <http://www.slideshare.net/techdude/scalable-web-architectures-common-patterns-and-approaches>



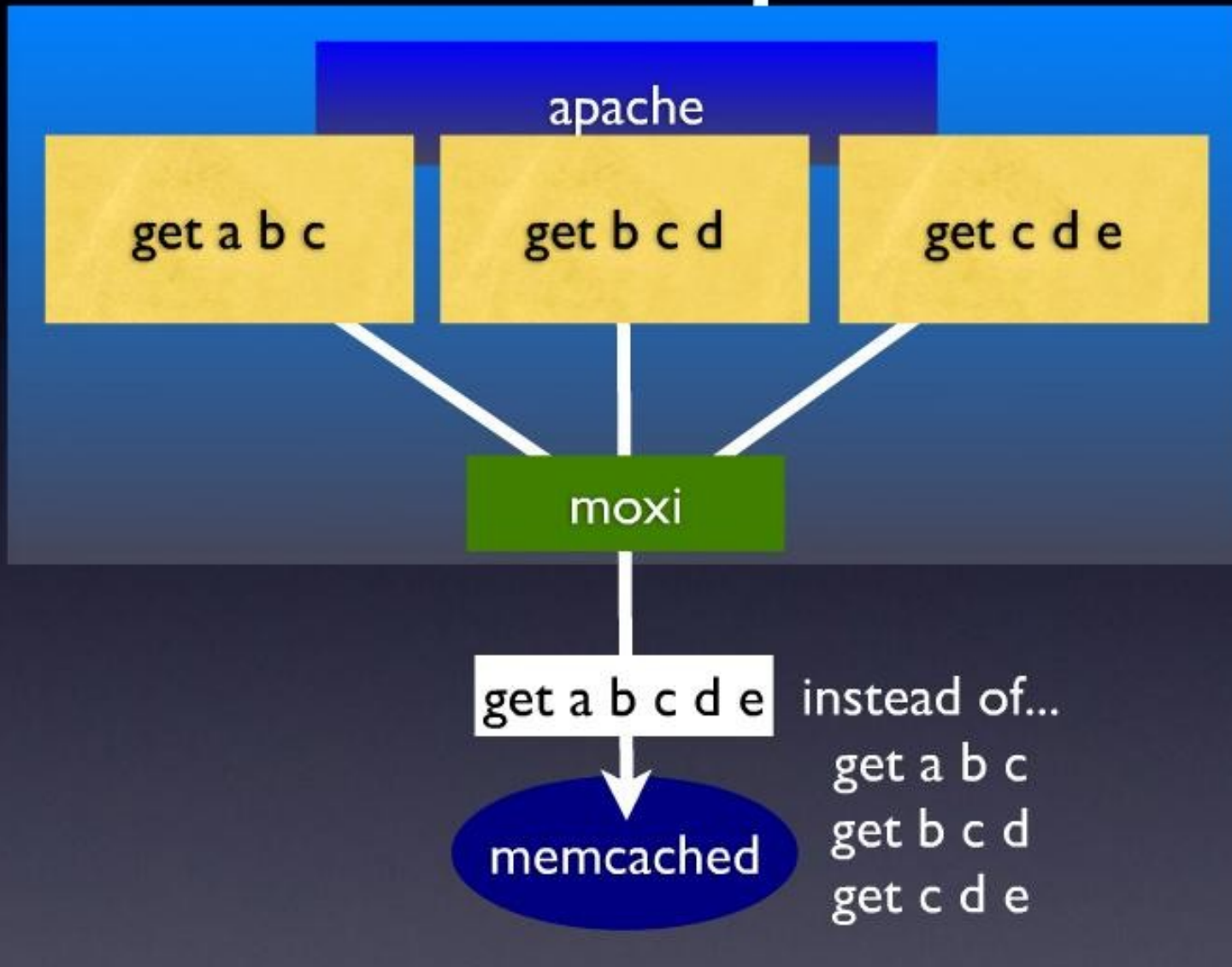
Memcached & MySQL

read

write

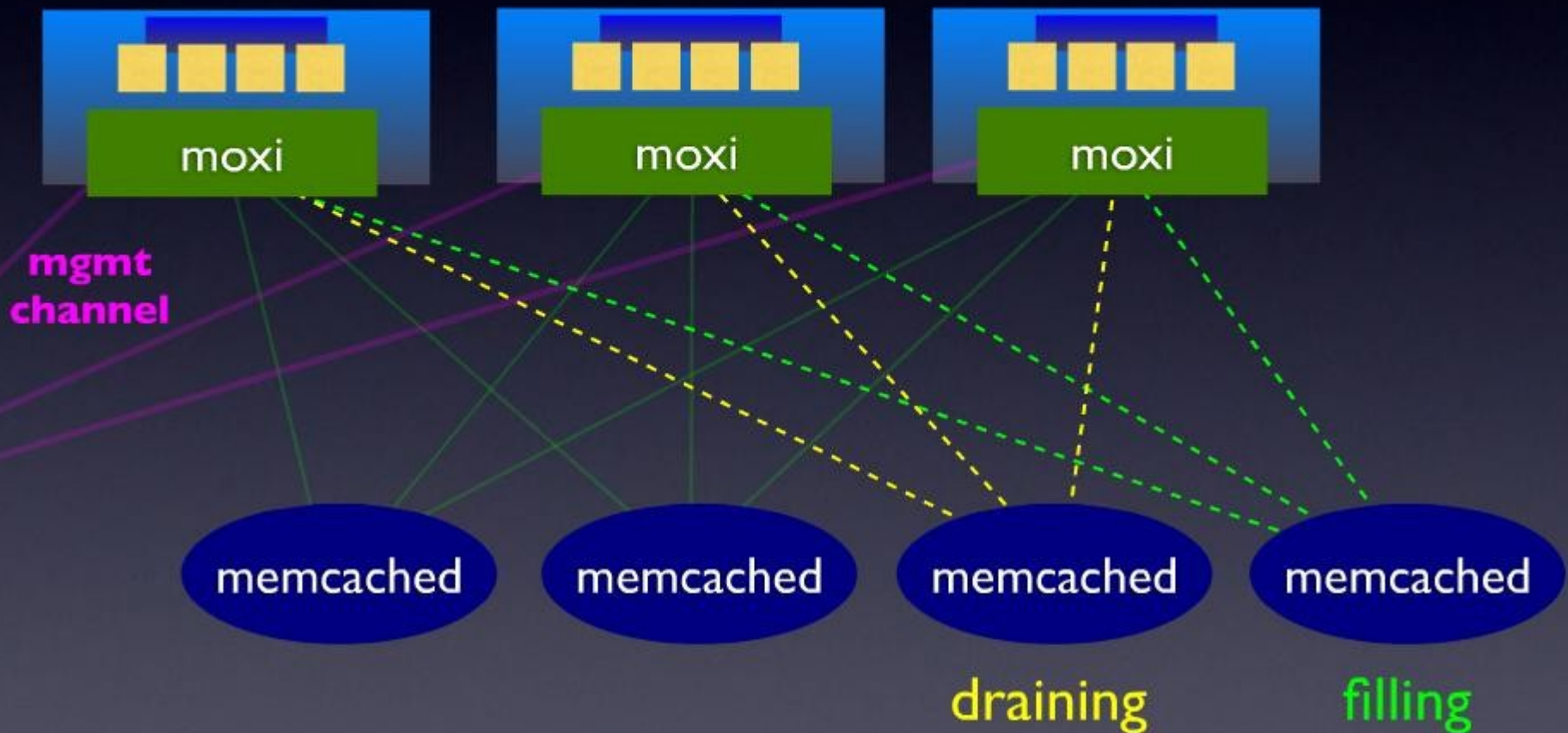


GET de-duplication



draining and filling

lazily migrate items from old server to new server



HBase is ..

- HBase is a distributed **column-oriented database** built on top of HDFS.
- A distributed data store that can scale horizontally to 1,000s of commodity servers and **petabytes** of indexed storage.
- Designed to operate on top of the Hadoop distributed file system (**HDFS**) or Kosmos File System (**KFS**, aka Cloudstore) for scalability, fault tolerance, and high availability.
- Integrated into the Hadoop **map-reduce** platform and paradigm.

Benefits

- Distributed storage
- Table-like in data structure
 - multi-dimensional map
- High scalability
- High availability
- High performance

Who use HBase

- Adobe
 - 內部使用 (Structure data)
- Kalooga
 - 圖片搜尋引擎 <http://www.kalooga.com/>
- Meetup
 - 社群聚會網站 <http://www.meetup.com/>
- Streamy
 - Migrate from MySQL to Hbase <http://www.streamy.com/>
- Trend Micro
 - 雲端掃毒架構 <http://trendmicro.com/>
- Yahoo!
 - 儲存文件 fingerprint 避免重複 <http://www.yahoo.com/>
- More - <http://wiki.apache.org/hadoop/Hbase/PoweredBy>

Backdrop

- Started toward by Chad Walters and Jim
- 2006.11
 - Google releases paper on **BigTable**
- 2007.2
 - Initial HBase prototype created as Hadoop contrib.
- 2007.10
 - First useable HBase
- 2008.1
 - Hadoop become Apache top-level project and HBase becomes subproject
- 2008.10~
 - HBase 0.18, 0.19 released

HBase Is Not ...

- Tables have **one primary index**, the *row key*.
- **No join operators.**
- Scans and queries can select a subset of available columns, perhaps by using a wildcard.
- There are three types of lookups:
 - Fast lookup using row key and optional timestamp.
 - Full table scan
 - Range scan from region start to end.

HBase Is Not ... (2)

- Limited atomicity and transaction support.
 - HBase supports **multiple batched mutations of single rows** only.
 - Data is unstructured and untyped.
- No accessed or manipulated via SQL.
 - Programmatic access via Java, REST, or **Thrift APIs**.
 - Scripting via JRuby.

Why Bigtable?

- Performance of RDBMS system is good for transaction processing but for very large scale analytic processing, the solutions are commercial, expensive, and specialized.
- Very large scale analytic processing
 - Big queries – typically range or table scans.
 - **Big databases (100s of TB)**

Why Bigtable? (2)

- Map reduce on Bigtable with optionally Cascading on top to support some relational algebras may be a cost effective solution.
- Sharding is not a solution to scale open source RDBMS platforms
 - Application specific
 - Labor intensive (re)partitionaing

Why HBase ?

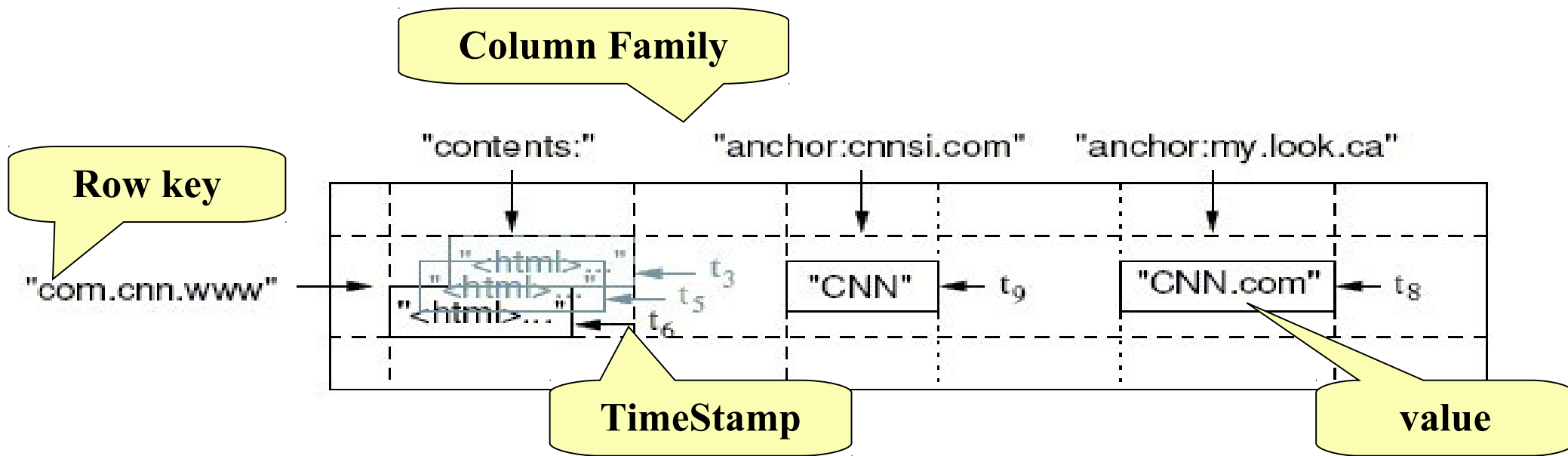
- HBase is a Bigtable clone.
- It is open source
- It has a good community and promise for the future
- It is developed on top of and has good integration for the Hadoop platform, if you are using Hadoop already.
- It has a Cascading connector.

HBase benefits than RDBMS

- *No real indexes*
- *Automatic partitioning*
- *Scale linearly and automatically* with new nodes
- *Commodity hardware*
- *Fault tolerance*
- *Batch processing*

Data Model

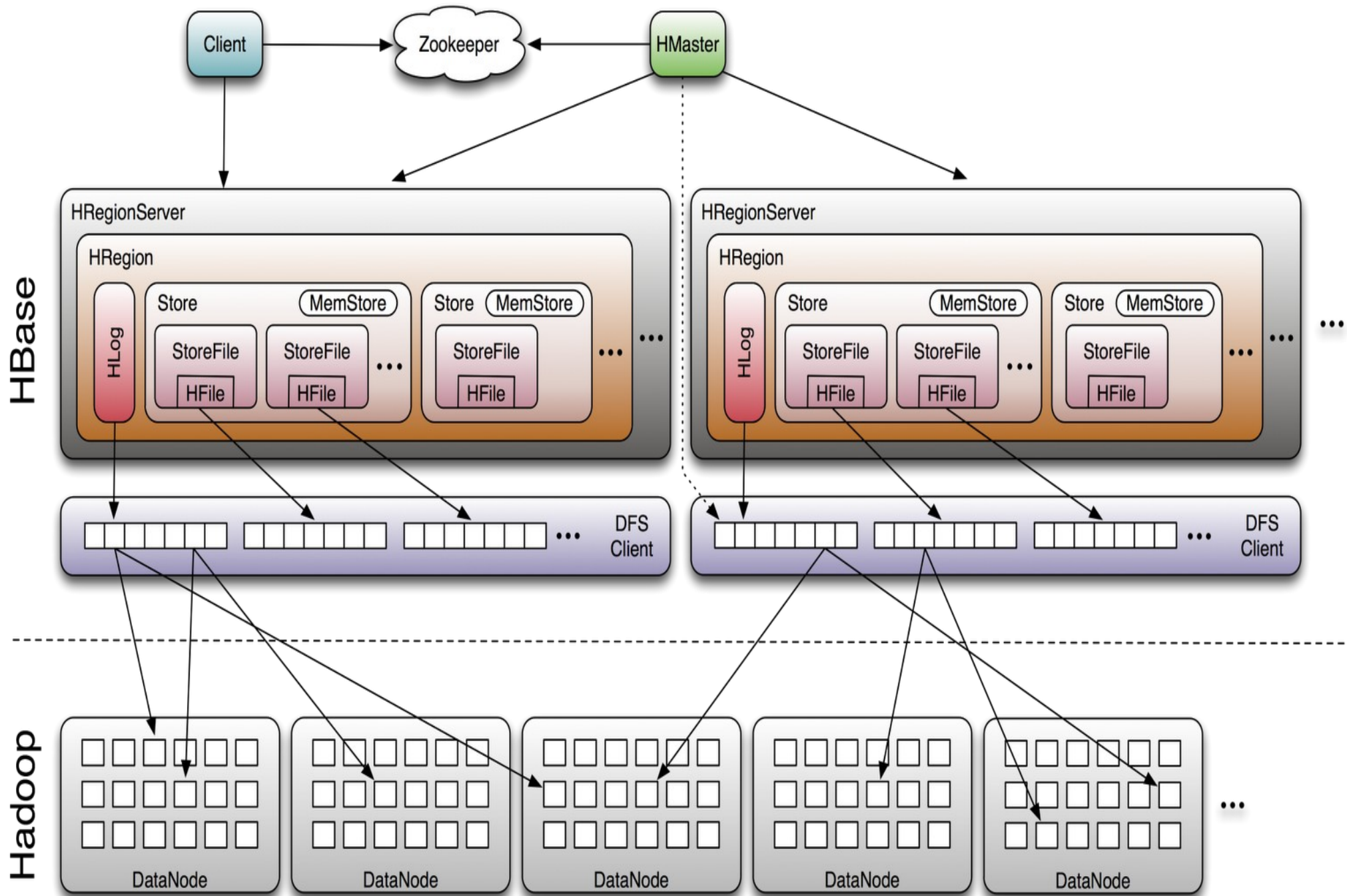
- Tables are sorted by **Row**
- Table schema only define it's *column families*.
 - Each family consists of any number of columns
 - Each column consists of any number of versions
 - Columns only exist when inserted, NULLs are free.
 - Columns within a family are sorted and stored together
- Everything except table names are byte[]
- **(Row, Family: Column, Timestamp) → Value**



Members

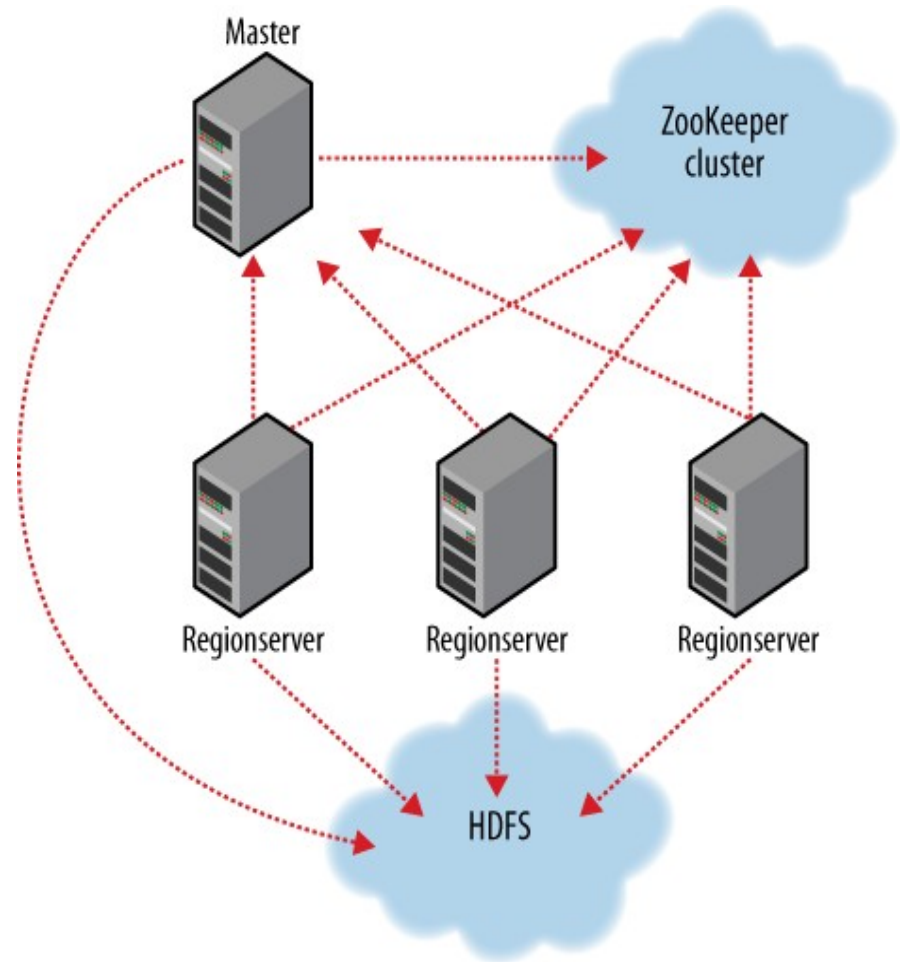
- *Master*
 - Responsible for monitoring region servers
 - Load balancing for regions
 - Redirect client to correct region servers
 - The current SPOF
- *regionserver slaves*
 - Serving requests(Write/Read/Scan) of Client
 - Send HeartBeat to Master
 - Throughput and Region numbers are scalable by region servers

Architecture



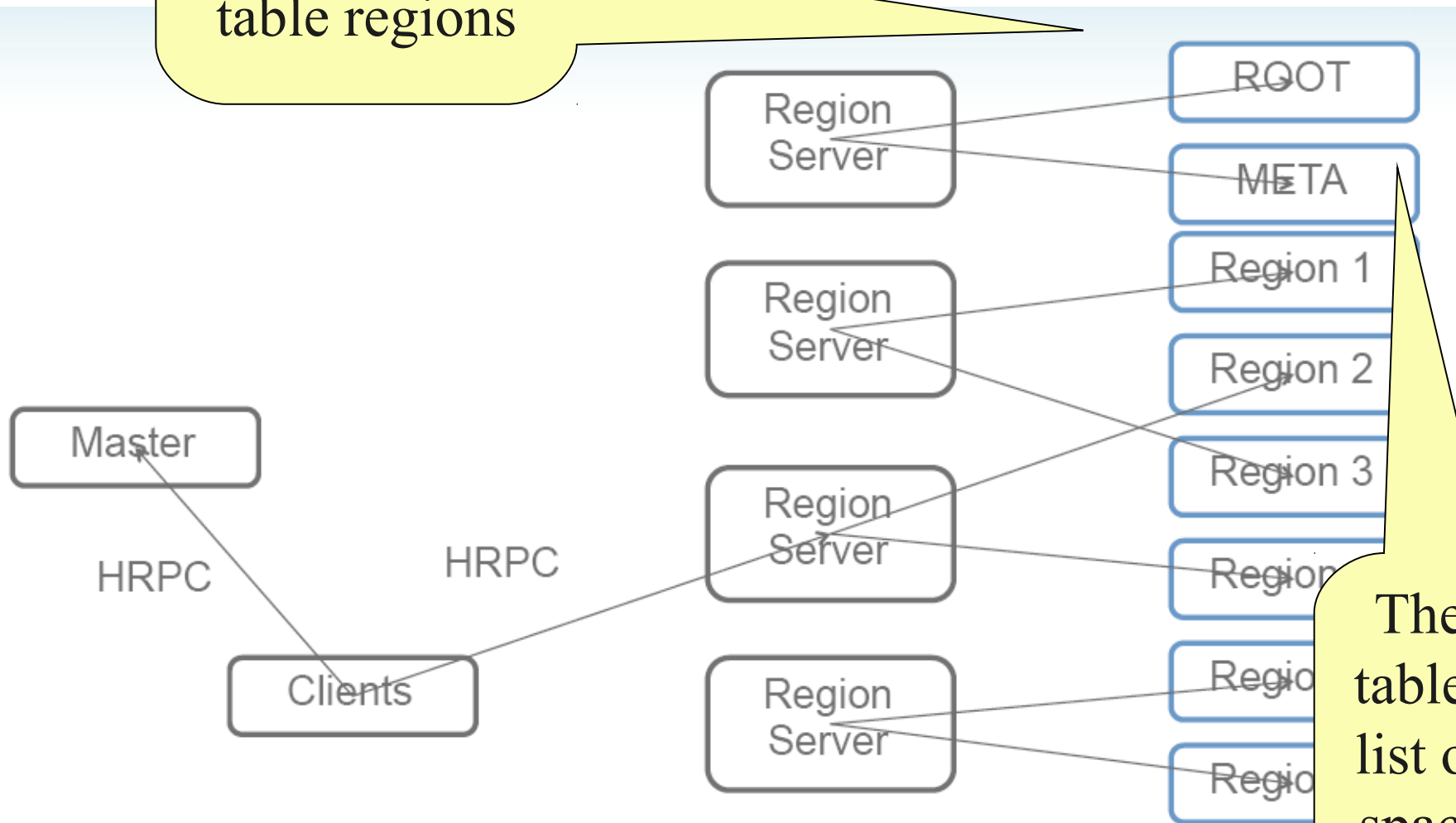
ZooKeeper

- HBase depends on ZooKeeper (Chapter 13) and by default it manages a ZooKeeper instance as the authority on cluster state



Operation

The `-ROOT-` table holds the list of `.META.` table regions



The `.META.` table holds the list of all user-space regions.



Questions?

Slides - <http://trac.nchc.org.tw/cloud>

Jazz Wang
Yao-Tsung Wang
jazz@nchc.org.tw



Powered by DRBL