# HBase

**A Comprehensive Introduction**

James Chin, Zikai Wang
Monday, March 14, 2011
CS 227 (Topics in Database Management)
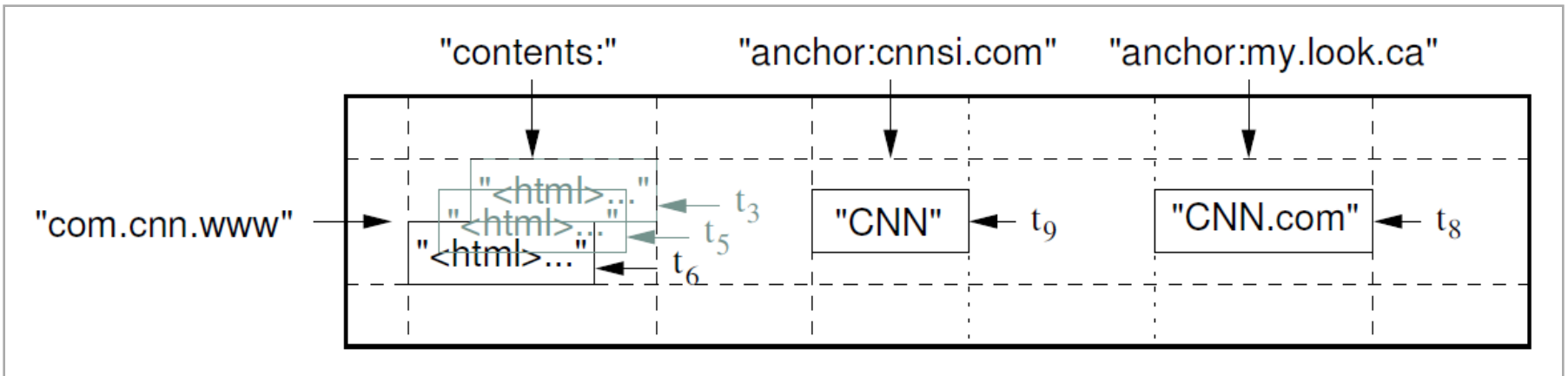CIT 367

H·BASE

# Overview

- Began as project by Powerset to **process massive amounts of data** for natural language search

- Open-source implementation of Google's **BigTable**
  - Lots of **semi-structured data**
  - Commodity Hardware
  - Horizontal Scalability
  - Tight integration with **MapReduce**

- Developed as part of Apache's **Hadoop** project and runs on top of **HDFS (Hadoop Distributed Filesystem)**
  - Provides **fault-tolerant** way of storing **large quantities of sparse data**.

- Non-relational, distributed database

- Column-Oriented

- Multi-Dimensional

- High Availability

- High Performance

# Data Model & Operators

# Data Model

- A **sparse**, **multi-dimensional**, **sorted** map
  - {row, column, timestamp} -> cell
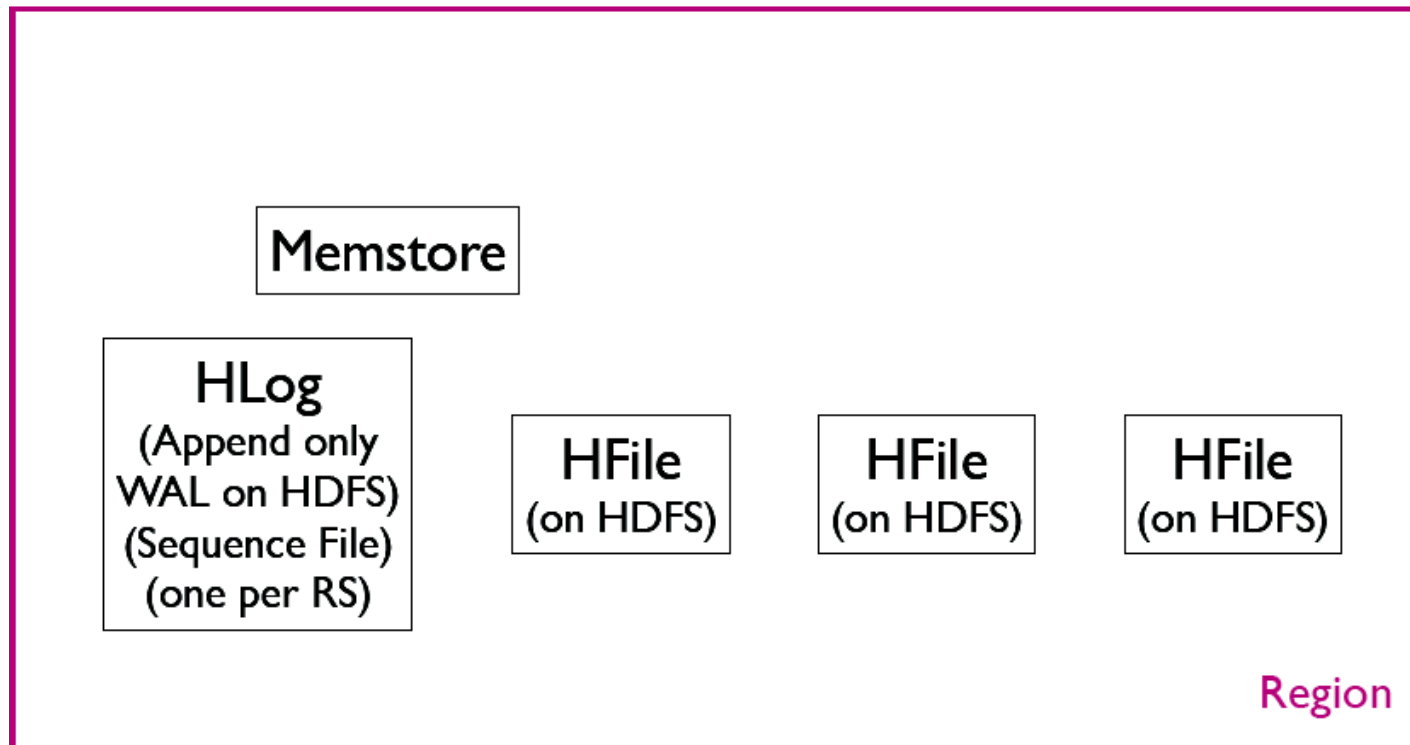- Column = **Column Family** : Column Qualifier



- Rows are **sorted lexicographically** based on row key
- **Region:** contiguous set of sorted rows
- HBase: a large number of columns, a low number of column families (2-3)

# Operators

- Operations are based on **row keys**

- **Single-row operations:**
  - Put
  - Get
  - Scan

- **Multi-row operations:**
  - Scan
  - MultiPut

- No built-in joins (use MapReduce)

# Physical Structures

- **Region:** unit of distribution and availability
- Regions are split when grown too large
- Max region size is a tuning parameter
  - Too low: prevents parallel scalability
  - Too high: makes things slow

- HBase has **no built-in support for secondary indexes**
- API only exposes operations by **row key**

| Row Key | Name | Position | Nationality |
|---------|------|----------|-------------|
| "1" | Nowitzki, Dirk | PF | Germany |
| "2" | Kaman, Chris | C | Germany |
| "3" | Gasol, Paul | PF | Spain |
| "4" | Fernandez, Rudy | SG | Spain |

- **Find all players from Spain?**
  - With built-in API, scan the entire table
  - Manually build a secondary index table
  - Exploit the fact that rows are sorted lexicographically by row key based on byte order

- ## Data Table:

| Row Key | Name | Position | Nationality |
|---|---|---|---|
| "1" | Nowitzki, Dirk | PF | Germany |
| "2" | Kaman, Chris | C | Germany |
| "3" | Gasol, Paul | PF | Spain |
| "4" | Fernandez, Rudy | SG | Spain |

- ## **Index table on nationality column**

  - a scan operation
  - start row = "Spain"
  - stop scanning: set a RowFilter with a BinaryPrefixComparator on the end value("Spain")
  - range queries are also supported

| Row Key | Dummy |
|---|---|
| "Germany 1" | Germany 1 |
| "Germany 2" | Germany 2 |
| "Spain 3" | Spain 3 |
| "Spain 4" | Spain 4 |

- **Find all power forwards from Spain?**
  - A composite index

- Row keys are **plain byte arrays**
  - Byte order = your desired order?
  - Convert strings, integers, floats, decimals carefully to bytes
  - Default sorting is ascending; if descending indexes are needed, reverse bit order

- **Lily's** HBase Indexing Library
    - Aids in building and querying indexes in HBase
    - Hides the details of playing with byte[] row keys

- **HBase + full text indexing and searching systems**
    - Apache Lucene (Apache Solr, elasticsearch)
    - Lily, HAvroBase (HBase + Solr), HBasene (HBase + Lucene)

# System Architecture

# APIs

- **Java**
  - Get, Put, Delete, Scan
  - IncrementColumnValue
  - TableInputFormat - MapReduce Source
  - TableOutputFormat - MapReduce Sink
- Rest
- Thrift
- Scala
- Jython
- Groovy DSL
- Ruby shell
- Java MR, Cascading, Pig, Hive

# ACID Properties

# ACID Properties

- HBase **not ACID-compliant**, but does guarantee certain specific properties

- **Atomicity**
  - All mutations are atomic within a row. Any put will either wholely succeed or wholely fail.
  - APIs that mutate several rows will *not* be atomic across the multiple rows.
  - The order of mutations is seen to happen in a well-defined order for each row, with no interleaving.

- **Consistency** and **Isolation**
  - All rows returned via any access API will consist of a complete row that existed at some point in the table's history.

## **Consistency** of Scans

- A scan is not a consistent view of a table. Scans do not exhibit snapshot isolation.
- Those familiar with relational databases will recognize this isolation level as "read committed".

## **Durability**

- All visible data is also durable data. That is to say, a read will never return data that has not been made durable on disk.
- Any operation that returns a "success" code (e.g. does not throw an exception) will be made durable.
- Any operation that returns a "failure" code will not be made durable (subject to the Atomicity guarantees above).
- All reasonable failure scenarios will not affect any of the listed ACID guarantees.

# Users

**Users:** Just to name a few…

- **Previous Solution:** Cassandra

- **Current Solution:** HBase

- **Why?** Cassandra's replication behavior

- Customer Indexing

- **Previous Solution:** offline process at a single node

- **Current Solution:**
  - Import user data into HBase
  - Periodically MapReduce job reading from HBase
  - Hits FlockDB and other internal services in mapper
  - Write data to sharded, replicated, horizontally scalable, in-memory, low-latency Scala service

- **Vs. Others:**
  - HDFS: Data is mutable
  - Cassandra: OLTP vs. OLAP?

# Users: Mozilla - Socorro

- **Socorro**, Mozilla's crash reporting system (https://crash-stats.mozilla.com/products)
  - Catches, processes, and presents crash data for Firefox, Thunderbird, Fennec, Camino, and Seamonkey.

- 2.5 million crash reports per week, 320GB per day

- **Previous Solution:** NFS (raw data), PostgreSQL (analyze results)
  - 15% of crash reports are processed

- **Current Solution:** Hadoop (processing) + HBase (storage)

# HBase vs. RDBMS

# HBase vs. RDBMS

| HBase | RDBMS |
|---|---|
| Column-oriented | Row oriented (mostly) |
| Flexible schema, add columns on the fly | Fixed schema |
| Good with sparse tables | Not optimized for sparse tables |
| No query language | SQL |
| Wide tables | Narrow tables |
| Joins using MR – not optimized | Optimized for joins (small, fast ones too!) |
| Tight integration with MR | Not really... |

# HBase vs. RDBMS (cont.)

| HBase | RDBMS |
|---|---|
| De-normalize your data | Normalize as you can |
| Horizontal scalability – just add hardware | Hard to shard and scale |
| Consistent | Consistent |
| No transactions | Transactional |
| Good for semi-structured data as well as structured data | Good for structured data |

# Questions?

# Thanks!