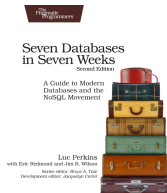


Lesson 5.3: Apache HBase

Instructor: Dr. GP Saggese - gsaggese@umd.edu

- Content in slides
- Web
 - 2006, BigTable paper
 - <https://hbase.apache.org/>
 - <https://github.com/apache/hbase>
- Good overview:
 - Seven Databases in Seven Weeks, 2e



(Apache) HBase



- HBase = **H**adoop Data**B**ase
 - Supports large tables on commodity hardware clusters
 - Column-oriented DB
 - Part of Apache Hadoop ecosystem
 - Uses Hadoop filesystem (HDFS)
 - HDFS modeled after Google File System (GFS)
 - HBase based on Google BigTable
 - Google BigTable runs on GFS, HBase runs on HDFS
 - Used by Google, Airbnb, eBay
- **When to use HBase**
 - For large DBs (e.g., many 100 GBs or TBs)
 - With at least 5 nodes in production
- **Applications**
 - Large-scale online analytics
 - Heavy-duty logging
 - Search systems (e.g., Internet search)
 - Facebook Messages (based on Cassandra)
 - Twitter metrics monitoring

HBase: Features

- Data versioning
 - Store versions of data
- Data compression
 - Compress and decompress on-the-fly
- Garbage collection for expired data
- In-memory tables
- Atomicity at row level
- Strong consistency guarantees
- Fault tolerant for machines and network
 - Write-ahead logging
 - Write data to in-memory log before disk
 - Distributed configuration
 - Nodes rely on each other, not centralized source

From HDFS to HBase

- **Different types of workloads for DB backends**
 - **OLTP** (On-Line Transactional Processing)
 - Read and write individual data items in a large table
 - E.g., update inventory and price as orders come in
 - **OLAP** (On-Line Analytical Processing)
 - Read large data amounts and process it
 - E.g., analyze item purchases over time
- **Hadoop FileSystem (HDFS) supports OLAP workloads**
 - Provide a filesystem with large files
 - Read data sequentially, end-to-end
 - Rarely updated
- **HBase supports OLTP interactions**
 - Built on HDFS
 - Use additional storage and memory to organize tables
 - Write tables back to HDFS as needed

HBase Data Model

- **Warning:** HBase uses names similar to relational DB concepts, but with different meanings
- A **database** consists of multiple tables
- Each **table** consists of multiple rows, sorted by row key
- Each row contains a *row key* and one or more column families
- Each **column family**
 - Contains multiple columns (family:column)
 - Defined when the table is created
- A **cell**
 - Uniquely identified by (table, row, family:column)
 - Contains metadata (e.g., timestamp) and an uninterpreted array of bytes (blob)
- **Versioning**
 - New values don't overwrite old ones
 - put() and get() allow specifying a timestamp (otherwise uses current time)

```
\# HBase Database: from table name to Table.  
Database = Dict[str, Table]
```

```
\# HBase Table.  
table: Table = {  
    # Row key  
    'row1': {  
        # (column family:column) → value  
        'cf1:col1': 'value1',  
        'cf1:col2': 'value2',  
        'cf2:col1': 'value3'  
    },  
    'row2': {  
        ... # More row data  
    }  
}  
database = {'table1': table}
```

```
\# Querying data.  
(value, metadata) = \  
table['row1']['cf1:col1']
```

Example 1: Colors and Shape

- Table with:
 - 2 column families: “color” and “shape”
 - 2 rows: “first” and “second”
- Row “first”:
 - 3 columns in “color”: “red”, “blue”, “yellow”
 - 1 column in “shape”: shape = 4
- Row “second”:
 - No columns in “color”
 - 2 columns in “shape”
- Access data using row key and

column (family:qualifier)



	row keys	column family "color"	column family "shape"
row	"first"	"red": "#F00" "blue": "#00F" "yellow": "#FF0"	"square": "4"
row	"second"		"triangle": "3" "square": "4"

```
table = {  
  'first': {  
    # (column family, column) → value  
    'color': {'red': '#F00',  
              'blue': '#00F',  
              'yellow': '#FF0'}  
    'shape': {'square': '4'}  
  },  
  'second': {  
    'shape': {'triangle': '3',  
              'square': '4'}  
  }  
}
```

Why All This Convoltuted Stuff?

- **A row in HBase is like a mini-database**

- A cell has many values
- Data stored sparsely

- **Rows in HBase are "deeper" than in relational DBs**

- Relational DBs: rows have many column values (fixed array with types)
- HBase: rows like a two-level nested dictionary with metadata (e.g., timestamp)

- **Applications**

- Store versioned website data
- Store a wiki

	row keys	column family "color"	column family "shape"
row	"first"	"red": "#F00" "blue": "#00F" "yellow": "#FF0"	"square": "4"
row	"second"		"triangle": "3" "square": "4"

Example 2: Storing a Wiki

- **Wiki (e.g., Wikipedia)**
 - Contains pages
 - Each page has a title, article text varies over time
- **HBase data model**
 - Table name → wikipedia
 - Row → entire wiki page
 - Row keys → wiki identifier (e.g., title or URL)
 - Column family → text
 - Column → " (empty)
 - Cell value → article text

	row keys (wiki page titles)	column family "text"
row (page)	"first page's title"	"": "Text of first page"
row (page)	"second page's title"	"": "Text of second page"

```
wikipedia_table = {  
    # wiki id.  
    'Home': {  
        # Column family:column $to$ value  
        ':text': 'Welcome to the wiki!',  
    },  
    'Welcome page': {  
        ... # More row data  
    }  
}  
  
Database = Dict[str, Table]  
database: Database = {'wikipedia':  
    wiki_table}  
(article, metadata) = \  
wiki_table['Home']['text']
```


Example 2: Storing a Wiki

- **Add data**

- Columns don't need to be predefined when creating a table
- The column is defined as `text`

```
> put 'wikipedia', 'Home', 'text', 'Welcome!'
```

- **Query data**

- Specify the table name, the row key, and optionally a list of columns

```
> get 'wikipedia', 'Home', 'text'  
text: timestamp=1295774833226,  
value=Welcome!
```

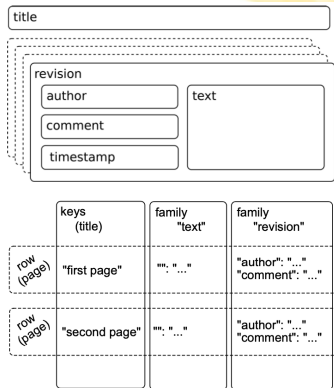
- HBase returns the timestamp (ms since the epoch 01-01-1970 UTC)

	row keys (wiki page titles)	column family "text"
row (page)	"first page's title"	"": "Text of first page"
row (page)	"second page's title"	"": "Text of second page"

```
wikipedia_table = {  
    # wiki id.  
    'Home': {  
        # Column family, column → value  
        'text': 'Welcome to the wiki!'  
    },  
    'Welcome page': {  
        ... # More row data  
    }  
}  
  
Database = Dict[str, Table]  
database: Database = {'wikipedia':  
    wiki_table}  
(queried_value, metadata) = \  
wiki_table['Home']['text']
```

Example 2: Improved Wiki

- **Improved wiki using versioning**
- A page
 - Uniquely identified by title
 - Can have multiple revisions
- A revision
 - Made by an author
 - Optionally contains a commit comment
 - Identified by timestamp
 - Contains text
- **HBase data model**
- Add family column "revision" with multiple columns
 - E.g., author, comment,
- Timestamp automatically binds article text and metadata
- Title not part of revision
 - Fixed and uniquely identifies page (like a primary key)
 - To change title, re-write entire row



Data in Tabular Form

Name			Home		Office	
Key	First	Last	Phone	Email	Phone	Email
101	Florian	Krebsbach	555-1212	florian@wobegon.org	666-1212	fk@phc.com
102	Marilyn	Tollerud	555-1213		666-1213	
103	Pastor	Inqvist			555-1214	inqvist@wel.org

- Fundamental operations
 - CREATE table, families
 - PUT table, rowid, family:column, value
 - PUT table, rowid, whole-row
 - GET table, rowid
 - SCAN table (*WITH filters*)
 - DROP table

Nested Data Representation

Name			Home		Office	
Key	First	Last	Phone	Email	Phone	Email
101	Florian	Krebsbach	555-1212	florian@wobegon.org	666-1212	fk@phc.com
102	Marilyn	Tollerud	555-1213		666-1213	
103	Pastor	Inqvist			555-1214	inqvist@wel.org

```
GET People:101
{
  Timestamp: T403;
  Name: {First="Florian", Last="Krebsbach"},
  Home: {Phone="555-1212", Email="florian@wobegon.org"},
  Office: {Phone="666-1212", Email="fk@phc.com"}
}
```

```
GET People:101:Name
{First="Florian", Last="Krebsbach"}
```

```
GET People:101:Name:First
"Florian"
```

Column Family vs Column

- **Adding a column**
 - Cheap
 - Done at run-time
- **Adding a column family**
 - Not at run-time
 - Requires table copy (expensive)
 - Indicates data storage method
 - Easy to add: map
 - Hard to add: static array
 - E.g., mongoDB document vs Relational DB column
- **Why differentiate column families vs columns?**
 - Why not store all row data in one column family?
 - Each column family configured independently, e.g.,
 - Compression
 - Performance tuning
 - Stored together in files
 - Designed for specific data types
 - E.g., timestamped web data for search engine

Consistency Model

- **Atomicity**
 - Update entire rows atomically or not at all
 - Independent of column count
- **Consistency**
 - GET returns a complete row from the table's history
 - Weak/eventual consistency
 - Check timestamp for certainty
 - SCAN
 - Includes all data written before scan
 - May include updates since start
- **Isolation**
 - Concurrent vs sequential semantics
 - Not guaranteed beyond a single row
 - Row is the atom of information
- **Durability**
 - Successful writes are durable on disk

Checking for Row or Column Existence

- HBase uses Bloom filters to check row or column existence
 - Acts like a cache for keys
 - Track presence without querying
- **Hashset complexity**
 - Unbounded space for data storage
 - No false positives
 - $O(1)$ average/amortized time
- **Bloom filter implementation**
 - Probabilistic hash set
 - Array of bits initially set to 0
 - Hash new data, set bits to 1
 - To test data presence, compute hash, check bits
 - All bits 0: data not seen
 - All bits 1: likely seen (false positive possible)
- **Bloom filter complexity**
 - Constant space usage
 - False positives possible (no false negatives)
 - $O(1)$ time complexity

Write-Ahead Log (WAL)

- Technique used by DBs
 - Provide atomicity and durability
 - Protect against node failures
 - Equivalent to journaling in file systems
- HBase and Postgres use WAL
- **WAL mechanics**
- Updated state of tables:
 - Not written to disk immediately
 - Buffered in memory
 - Written to disk as checkpoints periodically
- **Problem**
 - Server crash during this period loses state
- **Solution**
 - Use append-only disk-resident data structure
 - Log operations since last checkpoint in WAL
 - Clear WAL when tables stored to disk

Storing Variable-Length Data in Dbs

SQL Table

```
People(ID: Integer, FirstName: CHAR[20], LastName: CHAR[20], Phone: CHAR[8])  
UPDATE People SET Phone="555-3434" WHERE ID=403;
```

ID	FirstName	LastName	Phone
101	Florian	Krebsbach	555-3434
102	Marilyn	Tollerud	555-1213
103	Pastor	Ingvist	555-1214

- Each row: 52 bytes
- Move to next row:
fseek(file,+52)
- Get to Row 401: fseek(file,
401*52)
- Overwrite data in place

HBase Table

```
People(ID, Name, Home, Office)  
PUT People, 403, Home:Phone, 555-3434
```

```
{  
  101: {  
    Timestamp: T403;  
    Name: {First="Florian", Middle="Garfield", Last="Krebsbach"},  
    Home: {Phone="555-1212", Email="florian@wobegon.org"},  
    Office: {Phone="666-1212", Email="fk@phc.com"}  
  },  
  ...  
}
```

Need to use
pointers



HBase Implementation

- **How to store the web on disk?**
- **HBase is backed by HDFS**
 - Store each table (e.g., Wikipedia) in one file
 - “One file” means one gigantic file stored in HDFS
 - HDFS splits/replicates file into blocks on different servers
- Idea in several steps:
 - **Idea 1: Put entire table in one file**
 - Overwrite file with any cell change
 - Too slow
 - Idea 2: One file + WAL
 - Better, but doesn't scale to large data
 - Idea 3: One file per column family + WAL
 - Getting better!
 - Idea 4: Partition table into regions by key
 - Region = chunk of rows [a, b)
 - Regions never overlap

Idea 1: Put the Table in a Single File

- How do we do the following operations?
 - CREATE, DELETE (easy / fast)
 - SCAN (easy / fast)
 - GET, PUT (difficult / slow)

Table People(Name, Home, Office) { 101: { Timestamp: T403; Name: {First="Florian", Middle="Garfield", Last="Krebsbach"}, Home: {Phone="555-1212", Email="florian@wobegon.org"}, Office: {Phone="666-1212", Email="fk@phc.com"} }, 102: { Timestamp: T593; Name: {First="Marilyn", Last="Tollerud"}, Home: {Phone="555-1213"}, Office: {Phone="666-1213"} }, ... }

File "People"

Idea 2: One File + WAL

Table People(Name, Home, Office)

PUT 101:Office:Phone = "555-3434" PUT 102:Home:Email = mt@yahoo.com
....

WAL for Table People

- Changes are applied only to the log file
- The resulting record is cached in memory
- Reads must consult both memory and disk

Memory Cache for Table People

101

102

GET People:101

GET People:103

PUT People:101:Office:Phone = "555-3434"

Timestamp: T403;

Idea 2 Requires Periodic Table Update

101: {Timestamp: T403;Name: {First="Florian", Middle="Garfield", Last="Krebsbach"}},Home: {Phone="555-1212", Email="florian@wobegon.org"},Office: {Phone="666-1212", Email="fk@phc.com"}}, 102: {Timestamp: T593;Name: { First="Marilyn", Last="Tollerud"}},Home: { Phone="555-1213" },Office: { Phone="666-1213" }}, . . .

Table for People on Disk (Old)

PUT 101:Office:Phone = "555-3434" PUT 102:Home:Email = mt@yahoo.com
...

WAL for Table People:

101: {Timestamp: T403;Name: {First="Florian", Middle="Garfield", Last="Krebsbach"}},Home: {Phone="555-1212", Email="florian@wobegon.org"},Office: {Phone="555-3434", Email="fk@phc.com"}},102: {Timestamp: T593;Name: { First="Marilyn", Last="Tollerud"}},Home: { Phone="555-1213", Email="my@yahoo.com" },
...

Idea 3: Partition by Column Family

Data for Column Family Name

Tables for People on Disk (Old)

PUT 101:Office:Phone = "555-3434" PUT 102:Home:Email = mt@yahoo.com
...

WAL for Table People

Tables for People on Disk (New)

- Write out a new copy of the table, with all of the changes applied
- Delete the log and memory cache
- Start over

Data for Column Family Home

Data for Column Family Office

Data for Column Family Home (Changed)

Data for Column Family Office (Changed)

Data for Column Family Name



Final HBase Data Layout

