

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

# Accumulo – Extensions to Google's Bigtable Design

Adam Fuchs

National Security Agency  
Computer and Information Sciences Research Group

March 29, 2012

# Contents

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

- 1 Design Drivers
- 2 Apache Accumulo
  - Intro to Bigtable
  - Iterators
  - FATE
  - Major Compaction
- 3 Design Patterns
- 4 Fin

# Progress

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators  
FATE  
Major  
Compaction

Design  
Patterns

Fin

## 1 Design Drivers

## 2 Apache Accumulo

- Intro to Bigtable
- Iterators
- FATE
- Major Compaction

## 3 Design Patterns

## 4 Fin

# Design Drivers

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators  
FATE  
Major  
Compaction

Design  
Patterns

Fin

Analysis of big data is central to our customers' requirements, in which the strongest drivers are:

- *Scalability*: The ability to do twice the work at only (about) twice the cost.
- *Adaptability*: The ability to rapidly evolve the analytical tools available in an operational environment, building upon and enhancing existing capabilities.

From these directives we can derive the following requirements:

- Simplicity in the overall architecture to encourage collaboration and ameliorate learning curve.
- Generic design patterns to store and organize data whose format we don't control.
- Generic discovery analytics to retrieve and visualize generic data.
- Solutions for common sub-problems, such as multi-level security and enforcement of legal restrictions, built into the infrastructure.

# Optimization

Accumulo

Adam Fuchs

Design Drivers

Apache

Accumulo

Intro to Bigtable

Iterators

FATE

Major

Compaction

Design

Patterns

Fin

... is a secondary concern, given:

- hundreds of evolving applications,
- hundreds of changing data sources,
- non-trivial data volumes,
- many complicated **interactions**.

Instead, we need a generic platform that is *cheap, simple, scalable, secure, and adaptable*, with *pretty good* performance.



# Progress

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major

Compaction

Design  
Patterns

Fin

1 Design Drivers

2 Apache Accumulo

- Intro to Bigtable
- Iterators
- FATE
- Major Compaction

3 Design Patterns

4 Fin

# Apache Accumulo

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin



- First code written in Spring of 2008
- Open-sourced as an Apache Software Foundation incubator podling in September, 2011
- Graduated to Top-Level Project in March, 2012
- Mostly a clone of Bigtable, but includes several notable features:
  - Iterators: a framework for processing sorted streams of key/value entries
  - Cell-level Security: mandatory, attribute-based access control with key/value granularity
  - Fault-Tolerant Execution Framework (FATE)
  - A compaction scheduler with nice properties

# Progress

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

1 Design Drivers

2 Apache Accumulo

- Intro to Bigtable
- Iterators
- FATE
- Major Compaction

3 Design Patterns

4 Fin



# Basic Data Type

Accumulo

Adam Fuchs

Design Drivers

Apache

Accumulo

Intro to Bigtable

Iterators

FATE

Major

Compaction

Design

Patterns

Fin

An Accumulo Key is a 5-tuple, including:

- **Row:** controls *Atomicity*
- **Column Family:** controls *Locality*
- **Column Qualifier:** controls *Uniqueness*
- **Visibility:** controls *Access* (unique to Accumulo)
- **Timestamp:** controls *Versioning*

## Sample Entries

Row	: Col. Fam.	: Col. Qual.	: Visibility	: Timestamp	⇒ Value
Adam	: Favorites	: Food	: (Public)	: 20090801	⇒ Sushi
Adam	: Favorites	: Programming Language	: (Private)	: 20090830	⇒ Java
Adam	: Favorites	: Programming Language	: (Private)	: 20070725	⇒ C++
Adam	: Friends	: Bob	: (Public)	: 20110601	⇒
Adam	: Friends	: Joe	: (Private)	: 20110601	⇒

# Tablets

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

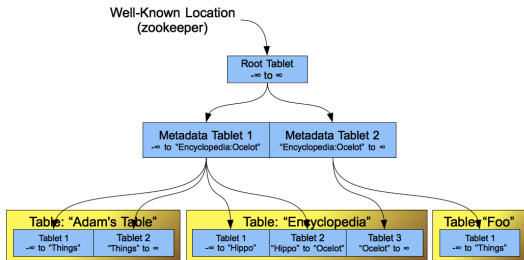
FATE

Major  
Compaction

Design  
Patterns

Fin

- Collections of key/value pairs form Tables
- Tables are partitioned into Tablets
- Metadata tablets hold info about other tablets, forming a three-level hierarchy
- A Tablet is a unit of work for a Tablet Server



# Distributed Processes

Accumulo

Adam Fuchs

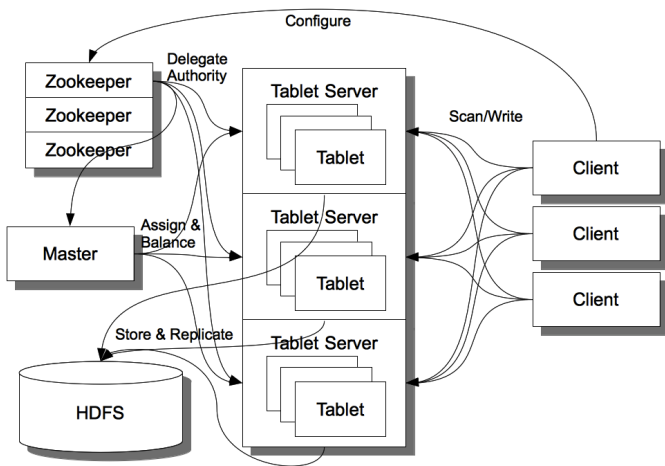
Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators  
FATE  
Major  
Compaction

Design  
Patterns

Fin



# Progress

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

**Iterators**

FATE

Major  
Compaction

Design  
Patterns

Fin

- 1 Design Drivers
- 2 **Apache Accumulo**
  - Intro to Bigtable
  - **Iterators**
  - FATE
  - Major Compaction
- 3 Design Patterns
- 4 Fin

# Tablet Server Composition

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

## Quick and loose definitions:

*Table*: A map of keys to values with one global sort order among keys.

*Tablet*: A row range within a Table.

*Tablet Server*: The mechanism that hosts Tablets, providing the primary functionality of Bigtable or Accumulo.

Tablet servers have several primary functions:

- 1 Hosting RPCs (read, write, etc.)
- 2 Managing resources (RAM, CPU, File I/O, etc.)
- 3 Scheduling background tasks (compactions, caching, etc.)
- 4 Handling key/value pairs

Category 4 is almost entirely accomplished through the *Iterator framework*.

# Tablet Server Data Flow

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

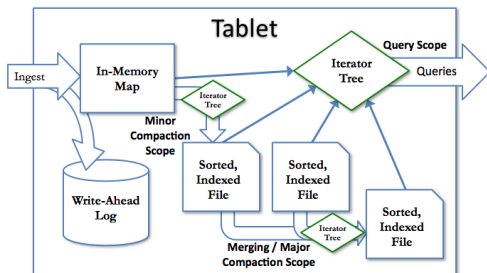
Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin



## Iterator Uses

- File Reads
- Block Caching
- Merging
- Deletion
- Isolation
- Locality Groups
- Range Selection
- Column Selection
- Cell-level Security
- Versioning
- Filtering
- Aggregation
- Partitioned Joins

# Iterators

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

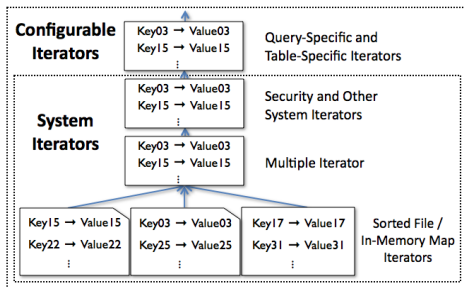
Intro to Bigtable  
Iterators

FATE  
Major  
Compaction

Design  
Patterns

Fin

- An *Iterator* is an object that provides an ordered stream of entries (key/value pairs), and supports basic *selection* and *filtering* methods.
- Core Iterators provide a basic view of a tablet's entries, implementing:
  - File Reads
  - Block Caching
  - Merging
  - Deletion
  - Isolation
  - Locality Groups
  - Range Selection
  - Column Selection
  - Cell-level Security
- Application-level Iterators modify table semantics to provide custom views, persisted or otherwise:
  - Versioning
  - Filtering
  - Aggregation
  - Partitioned Joins



# Modified Key/Value Pair Definition

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major

Compaction

Design  
Patterns

Fin

An Accumulo Key is a 5-tuple, including:

- **Row:** controls *Atomicity*
- **Column Family:** controls *Locality*
- **Column Qualifier:** controls *Uniqueness*
- **Visibility:** controls *Access* (unique to Accumulo)
- **Timestamp:** controls *Versioning*

## Sample Entries

Row	: Col. Fam.	: Col. Qual.	: Visibility	: Timestamp	⇒ Value
Adam	: Favorites	: Food	: (Public)	: 20090801	⇒ Sushi
Adam	: Favorites	: Programming Language	: (Private)	: 20090830	⇒ Java
Adam	: Favorites	: Programming Language	: (Private)	: 20070725	⇒ C++
Adam	: Friends	: Bob	: (Public)	: 20110601	⇒
Adam	: Friends	: Joe	: (Private)	: 20110601	⇒



# Visibility Label Syntax and Semantics

Accumulo

Adam Fuchs

Design Drivers

Apache

Accumulo

Intro to Bigtable

Iterators

FATE

Major

Compaction

Design

Patterns

Fin

## Document Labels

$Doc_1 : (\text{Federation})$   
 $Doc_2 : (\text{Klingon|Vulcan})$   
 $Doc_3 : (\text{Federation\&Human\&Vulcan})$   
 $Doc_4 : (\text{Federation\&}(\text{Human|Vulcan}))$

## User Authorization Sets

$CptKirk : \{\text{Federation, Human}\}$   
 $MrSpock : \{\text{Federation, Human, Vulcan}\}$

## Syntax

$WORD \Rightarrow [a-zA-Z0-9_]+$   
 $CLAUSE \Rightarrow AND$   
 $\Rightarrow OR$   
 $AND \Rightarrow AND \& \& AND$   
 $\Rightarrow (CLAUSE)$   
 $\Rightarrow WORD$   
 $OR \Rightarrow OR | OR$   
 $\Rightarrow (CLAUSE)$   
 $\Rightarrow WORD$

## Semantics

$$\frac{(T \Rightarrow \tau) \wedge (\tau \in A)}{(T, A) \models \text{true}} \text{ term}$$
$$\frac{(T \Rightarrow T_1 \& T_2) \wedge ((T_1, A) \models \text{true}) \wedge ((T_2, A) \models \text{true})}{(T, A) \models \text{true}} \text{ and}$$
$$\frac{(T \Rightarrow T_1 | T_2) \wedge (((T_1, A) \models \text{true}) \vee ((T_2, A) \models \text{true}))}{(T, A) \models \text{true}} \text{ or}$$
$$\frac{(T \Rightarrow (T_1)) \wedge (T_1 \models \text{true})}{(T, A) \models \text{true}} \text{ paren}$$

# Cell-Level Security Iterator

Accumulo

Adam Fuchs

Design Drivers

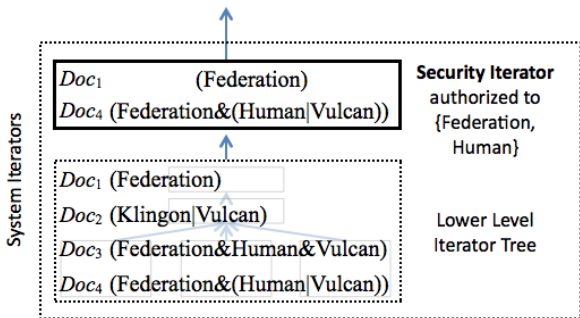
Apache  
Accumulo

Intro to Bigtable  
Iterators

FATE  
Major  
Compaction

Design  
Patterns

Fin



# Aggregation

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major

Compaction

Design  
Patterns

Fin

Goals: Count the number of times a word appears in a dynamic corpus, and count the number of documents that contain a given word.

## Sample Corpus

*Doc*<sub>1</sub> : "foo and bar are common variable names"

*Doc*<sub>2</sub> : "one cannot live on bar food alone"

*Doc*<sub>3</sub> : "Mr.T pities the fool at the bar"

*Doc*<sub>4</sub> : "someone should invent the kung foo bar"

## Input Key/Value Pairs:

Row	Column	Value
alone	<i>Doc</i> <sub>2</sub>	1
and	<i>Doc</i> <sub>1</sub>	1
are	<i>Doc</i> <sub>1</sub>	1
at	<i>Doc</i> <sub>3</sub>	1
bar	<i>Doc</i> <sub>1</sub>	1
bar	<i>Doc</i> <sub>2</sub>	1
bar	<i>Doc</i> <sub>3</sub>	1
bar	<i>Doc</i> <sub>4</sub>	1
cannot	<i>Doc</i> <sub>2</sub>	1
common	<i>Doc</i> <sub>1</sub>	1
foo	<i>Doc</i> <sub>1</sub>	1
foo	<i>Doc</i> <sub>4</sub>	1
food	<i>Doc</i> <sub>2</sub>	1
fool	<i>Doc</i> <sub>3</sub>	1
invent	<i>Doc</i> <sub>4</sub>	1
kung	<i>Doc</i> <sub>4</sub>	1
live	<i>Doc</i> <sub>2</sub>	1
Mr.T	<i>Doc</i> <sub>3</sub>	1
names	<i>Doc</i> <sub>1</sub>	1
on	<i>Doc</i> <sub>2</sub>	1
one	<i>Doc</i> <sub>2</sub>	1
should	<i>Doc</i> <sub>4</sub>	1
someone	<i>Doc</i> <sub>4</sub>	1
pities	<i>Doc</i> <sub>3</sub>	1
the	<i>Doc</i> <sub>3</sub>	1
the	<i>Doc</i> <sub>3</sub>	1
the	<i>Doc</i> <sub>4</sub>	1
variable	<i>Doc</i> <sub>1</sub>	1

# A Simple Aggregator

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

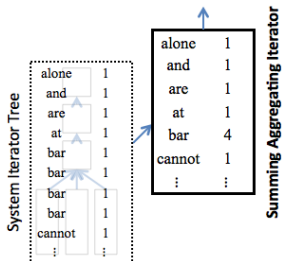
FATE

Major  
Compaction

Design  
Patterns

Fin

- Aggregators replace the "versioning" functionality of a table
- Any associative, commutative operations on the values for a given key can be encoded in an aggregator
- Aggregators can persist an aggregation of the entries written to the table
- Aggregators are significantly more efficient than a read-modify-write loop due to "lazy" aggregation



# Composing Multiple Iterators

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

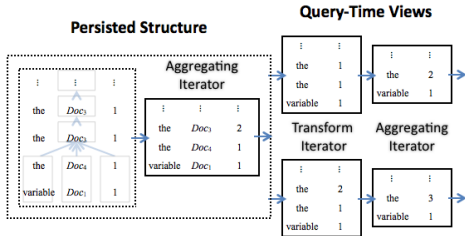
FATE

Major  
Compaction

Design  
Patterns

Fin

- We can compose multiple Iterators by streaming the results of one Iterator through another Iterator
- Partial aggregation for the persisted view keeps the table small
- Additional iterators and aggregators implement different discovery analytics at query time



# Accumulo vs. HBase Atomic Increment

## Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

- HBase performs a server-side *upsert* (read-modify-write), taking advantage of previous value being resident in write-cache
- Accumulo buffers inserts and aggregates lazily but consistently, taking advantage of merge-tree data streams
- Both methods implement the same atomic increment semantics
- Performance varies wildly...

# Increment Performance Comparison

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

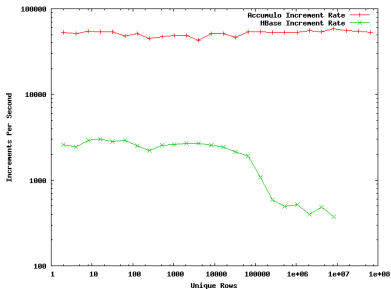
FATE

Major  
Compaction

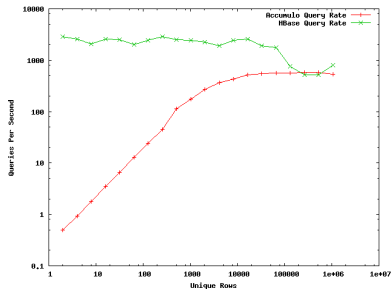
Design  
Patterns

Fin

## Write Performance



## Read Performance



- Aggregator wins for write performance with many different keys
- Upsert wins for read performance with a small number of keys
- Can we use both approaches?

# Multi-Term Query with Document Partitioning

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major

Compaction

Design  
Patterns

Fin

Goal: Find all of the documents that contain the words "foo" and "bar".

## Partitioned Corpus

<i>Doc</i> <sub>1</sub> :	"foo and bar are common variable names"	}	<i>Partition</i> <sub>1</sub>
<i>Doc</i> <sub>2</sub> :	"one cannot live on bar food alone"		
<i>Doc</i> <sub>3</sub> :	"Mr.T pities the fool at the bar"	}	<i>Partition</i> <sub>2</sub>
<i>Doc</i> <sub>4</sub> :	"someone should invent the kung foo bar"		



# Document Partitioning

Accumulo

Adam Fuchs

Design Drivers

Apache

Accumulo

Intro to Bigtable

Iterators

FATE

Major

Compaction

Design

Patterns

Fin

Divide and Conquer:

Row	ColFam	ColQual
<i>Part<sub>1</sub></i>	alone	<i>Doc<sub>2</sub></i>
<i>Part<sub>1</sub></i>	and	<i>Doc<sub>1</sub></i>
<i>Part<sub>1</sub></i>	are	<i>Doc<sub>1</sub></i>
<i>Part<sub>1</sub></i>	at	<i>Doc<sub>3</sub></i>
<i>Part<sub>1</sub></i>	bar	<i>Doc<sub>1</sub></i>
<i>Part<sub>1</sub></i>	bar	<i>Doc<sub>2</sub></i>
<i>Part<sub>1</sub></i>	bar	<i>Doc<sub>3</sub></i>
<i>Part<sub>1</sub></i>	cannot	<i>Doc<sub>2</sub></i>
<i>Part<sub>1</sub></i>	common	<i>Doc<sub>1</sub></i>
<i>Part<sub>1</sub></i>	foo	<i>Doc<sub>1</sub></i>
<i>Part<sub>1</sub></i>	food	<i>Doc<sub>2</sub></i>
<i>Part<sub>1</sub></i>	fool	<i>Doc<sub>3</sub></i>
<i>Part<sub>1</sub></i>	live	<i>Doc<sub>2</sub></i>
<i>Part<sub>1</sub></i>	Mr. T	<i>Doc<sub>3</sub></i>
<i>Part<sub>1</sub></i>	names	<i>Doc<sub>1</sub></i>
<i>Part<sub>1</sub></i>	on	<i>Doc<sub>2</sub></i>
<i>Part<sub>1</sub></i>	one	<i>Doc<sub>2</sub></i>
<i>Part<sub>1</sub></i>	pities	<i>Doc<sub>3</sub></i>
<i>Part<sub>1</sub></i>	the	<i>Doc<sub>3</sub></i>
<i>Part<sub>1</sub></i>	variable	<i>Doc<sub>1</sub></i>

Row	ColFam	ColQual
<i>Part<sub>2</sub></i>	bar	<i>Doc<sub>4</sub></i>
<i>Part<sub>2</sub></i>	foo	<i>Doc<sub>4</sub></i>
<i>Part<sub>2</sub></i>	invent	<i>Doc<sub>4</sub></i>
<i>Part<sub>2</sub></i>	kung	<i>Doc<sub>4</sub></i>
<i>Part<sub>2</sub></i>	should	<i>Doc<sub>4</sub></i>
<i>Part<sub>2</sub></i>	someone	<i>Doc<sub>4</sub></i>
<i>Part<sub>2</sub></i>	the	<i>Doc<sub>4</sub></i>

# Partitioned Join Iterator

Accumulo

Adam Fuchs

Design Drivers

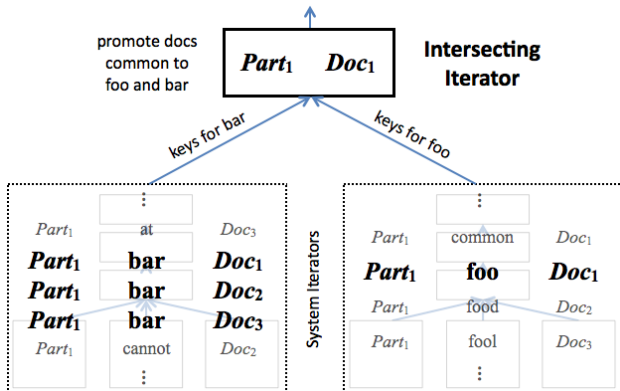
Apache  
Accumulo

Intro to Bigtable  
Iterators

FATE  
Major  
Compaction

Design  
Patterns

Fin



# Wikipedia Search Engine Experiment

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

## Goals:

- Create a generic text indexing platform
- Support a complex query language (i.e. mappable from Lucene)
- Scale to multiple nodes
- Support low-latency updates
- Support automatic balancing and fail-over

## Data

- Three languages of Wikipedia: EN, ES, DE
- 5.9 million articles
- 2.37 billion (word,document) tuples
- 11.8 GB (compressed)

## Cluster

- 10 Nodes
- 30 TB disk (60x500GB drives)
- 120 cores
- 320 GB RAM

# Wikipedia Search Results

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

- Tested on conjunctions of high-degree terms
- Retrieved entire contents of articles matching queries
- Paging possible for ultra-low latency response time

## Query Performance

Query	Samples (seconds)					Matches	Result Size
"old" and "man" and "sea"	4.07	3.79	3.65	3.85	3.67	22,956	3,830,102
"paris" and "in" and "the" and "spring"	3.06	3.06	2.78	3.02	2.92	10,755	1,757,293
"rubber" and "ducky" and "ernie"	0.08	0.08	0.10	0.11	0.10	6	808
"fast" and ("furious" or "furriest")	1.34	1.33	1.30	1.31	1.31	2,973	493,800
"slashdot" and "grok"	0.06	0.06	0.06	0.06	0.06	14	2,371
"three" and "little" and "pigs"	0.92	0.91	0.90	1.08	0.88	2,742	481,531

## Documents per Term

Term	Cardinality
ducky	795
ernie	13,433
fast	166,813
furious	10,535
furriest	45
grok	1,168

Term	Cardinality
in	1,884,638
little	320,748
man	548,238
old	720,795
paris	232,464
pigs	8,356

Term	Cardinality
rubber	17,235
sea	247,231
slashdot	2,343
spring	125,605
the	3,509,498
three	718,810

# Iterator Summary

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

**Iterators**

FATE

Major  
Compaction

Design  
Patterns

Fin

Iterators provide a modular implementation of Tablet Server functionality, resulting in:

- Reduced complexity of Tablet Server code
- Increased unit testability
- Simple extensibility for specialized applications

# Progress

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators

**FATE**

Major  
Compaction

Design  
Patterns

Fin

1 Design Drivers

2 **Apache Accumulo**

- Intro to Bigtable
- Iterators
- **FATE**
- Major Compaction

3 Design Patterns

4 Fin

# The Perils of Distributed Computing

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators

FATE  
Major  
Compaction

Design  
Patterns

Fin

## Dealing with failures is hard!

- Operations like table creation are logically atomic, but consist of multiple operations on distributed systems.
- Resource locking (via mutex, semaphores, etc.) provides some sanity.
- Distributed systems have many complicated failure modes: clients, master, tablet servers, and dependent systems can all go offline periodically.
- *Who is responsible for unlocking locks when any component can fail?*
- *How do we know it's safe to unlock a lock?*

# Accumulo Testing Procedures

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators

FATE  
Major  
Compaction

Design  
Patterns

Fin

## Testing Frameworks

- **Unit:** Verify correct functioning of each module separately
- **System:** Perform correctness and performance tests on a small running instance
- **Load/Scale:** Generate high loads at scale and measure performance and correctness
- **Random Walk:** Randomly, repeatedly, and concurrently execute a variety of test modules representative of user activity on an instance at scale
- **Simulation:** Evaluate the model to gauge expected performance

## Other Considerations

- Scoping tests to include server-side code, client-side code, dependent processes, etc.
- Code coverage vs. path coverage
- Static vs. dynamic analysis
- Simulating failures of distributed components
- Strange failure modes (often hardware/physics-related)



# Fault-Tolerant Executor

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators

**FATE**  
Major  
Compaction

Design  
Patterns

Fin

- If a process dies, previously submitted operations continue to execute on restart.
- FATE serializes every task in Zookeeper before execution.
- The Master process uses FATE to execute table operations and administrative actions.
- FATE eliminates the single point of failure.

# Adampotence

Accumulo

Adam Fuchs

Design Drivers

Apache

Accumulo

Intro to Bigtable

Iterators

FATE

Major

Compaction

Design

Patterns

Fin

- Idempotent:  $f(f(x)) = f(x)$
- Adampotent:  $f(f'(x)) = f(x)$ ,  
where  $f'(x)$  denotes partial execution of  $f(x)$

# REPO: Repeatable Persisted Operation

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

```
public interface Repo<T> extends Serializable {  
  
    long isReady(long tid, T environment) throws Exception;  
  
    Repo<T> call(long tid, T environment) throws Exception;  
  
    void undo(long tid, T environment) throws Exception;  
  
}
```

- call() returns next op, null if done
- call(), undo(), and isReady() must be idempotent
- undo() should clean up any possible partial execution of isReady() or call()

# FATE API

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators

FATE

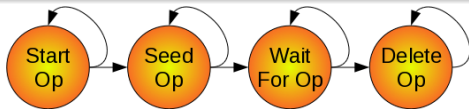
Major  
Compaction

Design  
Patterns

Fin

## Client API

```
long startTransaction();  
void seedTransaction(long tid, Repo op);  
TStatus waitForCompletion(long tid);  
Exception getException(long tid);  
void delete(long tid);
```



# FATE Execution State Model

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

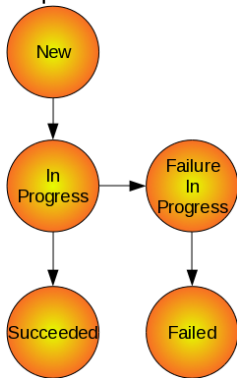
Intro to Bigtable  
Iterators

FATE  
Major  
Compaction

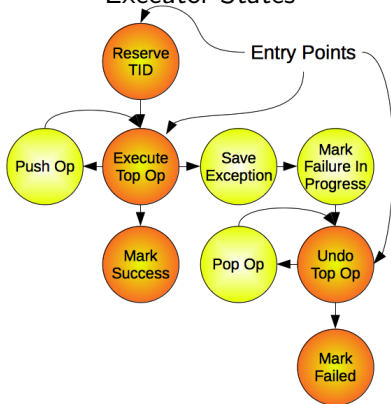
Design  
Patterns

Fin

## Operation States



## Executor States



# CreateTable FATE Op

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators

FATE  
Major  
Compaction

Design  
Patterns

Fin

## Steps for CreateTable Operation:

- 1 Reserve a Table ID
- 2 Set Table Permissions
- 3 Populate Configuration in Zookeeper
  - Reentrantly lock table
  - Relate table name to table ID
- 4 Create HDFS Directory
- 5 Populate Metadata Table Entries
- 6 Finish Create Table
  - Notify Master of new tablet(s)
  - Unlock table

# FATE Admin Tool

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major

Compaction

Design  
Patterns

Fin

```
$ ./bin/accumulo org.apache.accumulo.server.fate.Admin print
```

```
txid: 59c0403614dc0c39 status: IN_PROGRESS op: RenameTable locked: [] locking: [W:cz] top: RenameTable
txid: 37539f8d61548764 status: IN_PROGRESS op: ChangeTableState locked: [] locking: [W:cz] top: ChangeTableState
txid: 02f8323a3136e60d status: IN_PROGRESS op: TableRangeOp locked: [] locking: [W:cz] top: TableRangeOp
txid: 044015732e97eec1 status: IN_PROGRESS op: CompactRange locked: [] locking: [R:cz] top: CompactRange
txid: 6ce9dd63f9d51448 status: IN_PROGRESS op: CompactRange locked: [] locking: [R:cz] top: CompactRange
txid: 417cb9b60e44ecd9 status: IN_PROGRESS op: TableRangeOp locked: [] locking: [W:cz] top: TableRangeOp
txid: 5e7c5284a4677d6c status: IN_PROGRESS op: DeleteTable locked: [] locking: [W:cz] top: DeleteTable
txid: 6633d3d841d66995 status: IN_PROGRESS op: TableRangeOp locked: [W:cz] locking: [] top: TableRangeOpWait
```

- Monitoring tool for FATE operations
- Supports debugging, such as with deadlocks
- Helps recovery from failed clients

# FATE Summary

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators

FATE  
Major  
Compaction

Design  
Patterns

Fin

- FATE provides generic fault tolerance for administrative actions
- With FATE, we removed custom synchronization code for a dozen procedures
- Table-level locking is now low risk
- Improves testability
- Reduces complexity
- Increases modularity

## FATE Operations

- BulkImport
- ChangeTableState
- CloneTable
- CompactRange
- CreateTable
- DeleteTable
- RenameTable
- TableRangeOp
- DisconnectLogger
- FlushTablets
- ShutdownTServer
- StopLogger



# Progress

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

- 1 Design Drivers
- 2 **Apache Accumulo**
  - Intro to Bigtable
  - Iterators
  - FATE
  - **Major Compaction**
- 3 Design Patterns
- 4 Fin

# Major Compaction Efficiency

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

Major Compaction: Noun. The tablet operation that merges multiple files into one file.

- Overly aggressive major compaction results in  $N^2$  write complexity
- Overly lazy major compaction results in disk thrashing during queries (or unavailable tablets)
- Tuning major compaction operations is a trade-off between ingest and query performance

# Accumulo Major Compaction Algorithm

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

- 1 let  $r \geq 1.0$  be some ratio
- 2  $F \Leftarrow$  all files referenced by a tablet
- 3 if  $F$  is empty then exit
- 4  $f_0 \Leftarrow$  biggest file in  $F$
- 5  $a \Leftarrow$  aggregate size of files in  $F$
- 6 if  $a > r|f_0|$  then compact all files in  $F$  and exit
- 7 otherwise, remove  $f_0$  from  $F$  and go to step 3

# Major Compaction Performance

Accumulo

Adam Fuchs

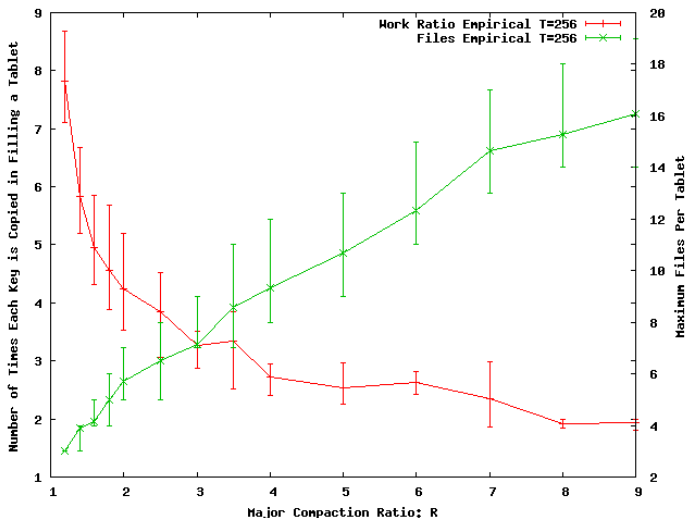
Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators  
FATE  
Major  
Compaction

Design  
Patterns

Fin



# Progress

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators  
FATE  
Major  
Compaction

Design  
Patterns

Fin

- 1 Design Drivers
- 2 Apache Accumulo
  - Intro to Bigtable
  - Iterators
  - FATE
  - Major Compaction
- 3 Design Patterns
- 4 Fin

# Design Patterns

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

Our use of Accumulo fundamentally differs from how we use RDBMS technology. In particular, Accumulo supports:

- Wide, sparse rows
- Indexes that span multiple columns

To adapt Accumulo for use in our applications, we have formalized several design patterns for Accumulo (or any Bigtable clone) including:

- Information Retrieval Patterns and Discovery Analytics
- Graph Analysis Patterns
- Machine Learning Patterns
- ...

# Event Table with Inverted Index

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators  
FATE  
Major  
Compaction

Design  
Patterns

Fin

**Table:**

Event Table

Inverted Index

**Row:**

<UUID>

<Term>

**Column Family:**

<Type>

<Type> + <Field>

**Column Qualifier:**

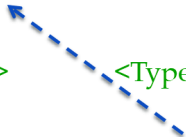
<Field>

<UUID>

**Value:**

<Term>

<Digest of Event>



# Inverted Index Flow

Accumulo

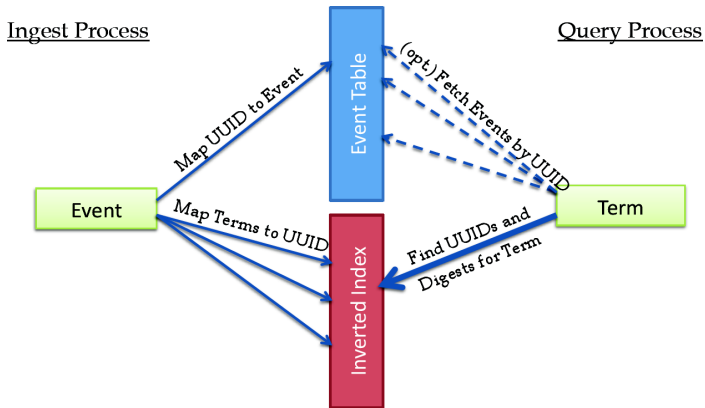
Adam Fuchs

Design Drivers

Apache  
Accumulo  
Intro to Bigtable  
Iterators  
FATE  
Major  
Compaction

Design  
Patterns

Fin





# Document Partitioned Index

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators  
FATE  
Major  
Compaction

Design  
Patterns

Fin

Table:

Indexed Event Table

Row:

<Partition ID>

Column Family:

"Event"

"Index"

"Geo"

Column Qualifier  
(2-Tuples):

<UUID>

<Term>

<Morton Coordinate>

<Field>

<UUID>

<UUID>

Value:

<Term>

Event and Index  
Records Co-Partitioned!

# Document Partitioned Flow

Accumulo

Adam Fuchs

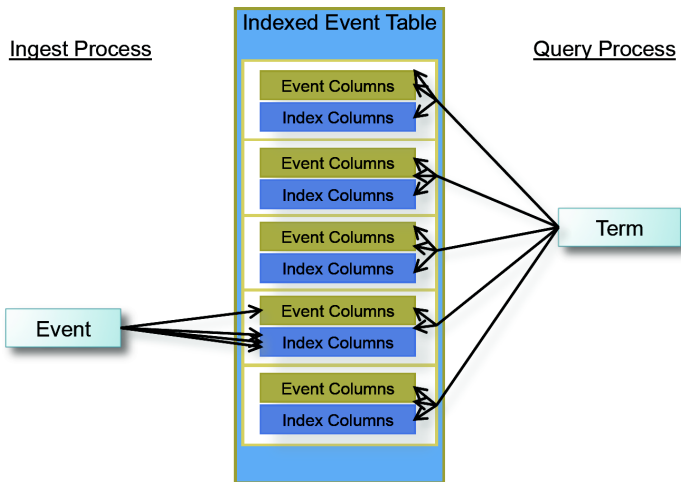
Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators  
FATE  
Major  
Compaction

Design  
Patterns

Fin



# Multidimensional Index

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

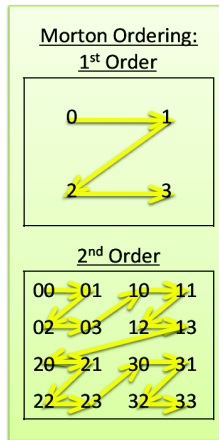
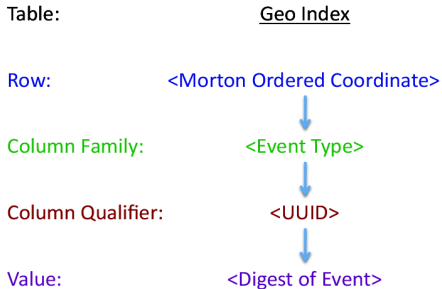
Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin



See also: <http://en.wikipedia.org/wiki/Geohash>

# Graph Table

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

Table:

Graph Table

Row:

<Node ID>

Column Family:

"Node Info"

"Out Edges"

"In Edges"

Column Qualifier  
(Tuples):

<Field>

<Node ID>

<Node ID>

<Edge ID>

<Edge ID>

Value:

<Value>

<Edge Info>

<Edge Info>

# Progress

Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable  
Iterators  
FATE  
Major  
Compaction

Design  
Patterns

Fin

- 1 Design Drivers
- 2 Apache Accumulo
  - Intro to Bigtable
  - Iterators
  - FATE
  - Major Compaction
- 3 Design Patterns
- 4 Fin

# Other Accumulo Features

## Accumulo

Adam Fuchs

Design Drivers

Apache  
Accumulo

Intro to Bigtable

Iterators

FATE

Major  
Compaction

Design  
Patterns

Fin

Check out Apache Accumulo (<http://accumulo.apache.org/>) for interesting implementations of:

- Merging Tablets
- Table Cloning: Hard link-style table copying
- Relative Key Encoded RFile file format
- Adaptive locality groups
- Isolation over scans of wide rows
- Bulk loading
- Logical time
- Client-side threading models for batch writes and scans
- Merging minor compactions
- Distributed write-ahead log