# BlueRunner: Building an Email Service in the Cloud

Jun Rao

IBM Almaden Research Center

Apache Cassandra Committer

# The Team

## User Systems and Experience Research



Jerald Schoudt

Stefan Nusser

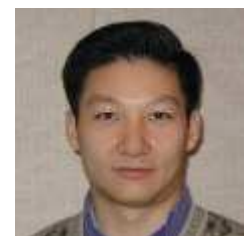Jan Pieper

Hernan Badenes

Julian Cerruti

## Information Management Research

Sandeep Tata

Bo Shekita

Jun Rao

2

# Outline

- BlueRunner overview
- Scalable Row Stores
- BlueRunner design in Cassandra
- Preliminary performance results
- Summary

# What's BlueRunner

- Research prototype for hosted emails at IBM ARC
- Browser-based email client + Cassandra backend
  - 3+ years on the client
  - ~1 year on the Cassandra backend
- Many advanced features in client
  - Scrolling
  - Foldering/Tagging
  - Sorting/Pivoting
  - Threading
  - Orienteering and Usability Improvements
- Backend
  - Thin client; most operations pushed to the backend
  - Designed for large-scale hosted environment
    - 100K mailboxes, each with 100K messages

# How does the client look?

# Why a new backend?

| Limitations in Traditional DBs |
| --- |
| Fault-tolerance relies on expensive reliable storage |
| Weak elasticity--- hard to grow a cluster incrementally |
| No automatic load-balancing |
| Rigid relational schema |
| No versioning support |

limiting scalability

mismatch for many apps

# The Cloud Landscape for Scalable Backend

**Linked in** ®
Voldemort

**Google** App Engine
**BigTable**

Yahoo!
PNUTS

apache **CouchDB** relax

Cassandra
**facebook**

**amazon** web services™
Amazon SimpleDB™ BETA

**Sun**

*Drizzle*

APACHE
HBase

**mongoDB**
{name: "mongo", type: "db"}

- **Flurry of activity in this space motivated by**
  - RDBMS too rigid/heavy for some apps
  - Existing RDBMS engines missing many key cloud requirements

# What is a Scalable Row Store?

- Middle ground btw a DBMS and a file system
  - Much simpler API then SQL
  - Designed to scale

| Limitations in Traditional DBs | Scalable Row Stores |
|---|---|
| Fault-tolerance relies on expensive reliable storage | Fault-tolerance done in software; replication on commodity disks |
| Weak elasticity--- hard to grow a cluster incrementally | Can grow a cluster incrementally and online |
| No automatic load-balancing | Built-in automatic load-balancing |
| Rigid relational schema | No strong schema required |
| No versioning support | Built-in versioning |

# Cassandra

- Google's Bigtable data model + Amazon's Dynamo scalable architecture

- Developed by Facebook in 2007

  - Used in production for a few apps (e.g., inbox search for 200M users)

- Became an Apache Incubator project early 2009

  - active community

  - additional committers from Rackspace and IBM

  - contributors from Digg, Twitter, etc
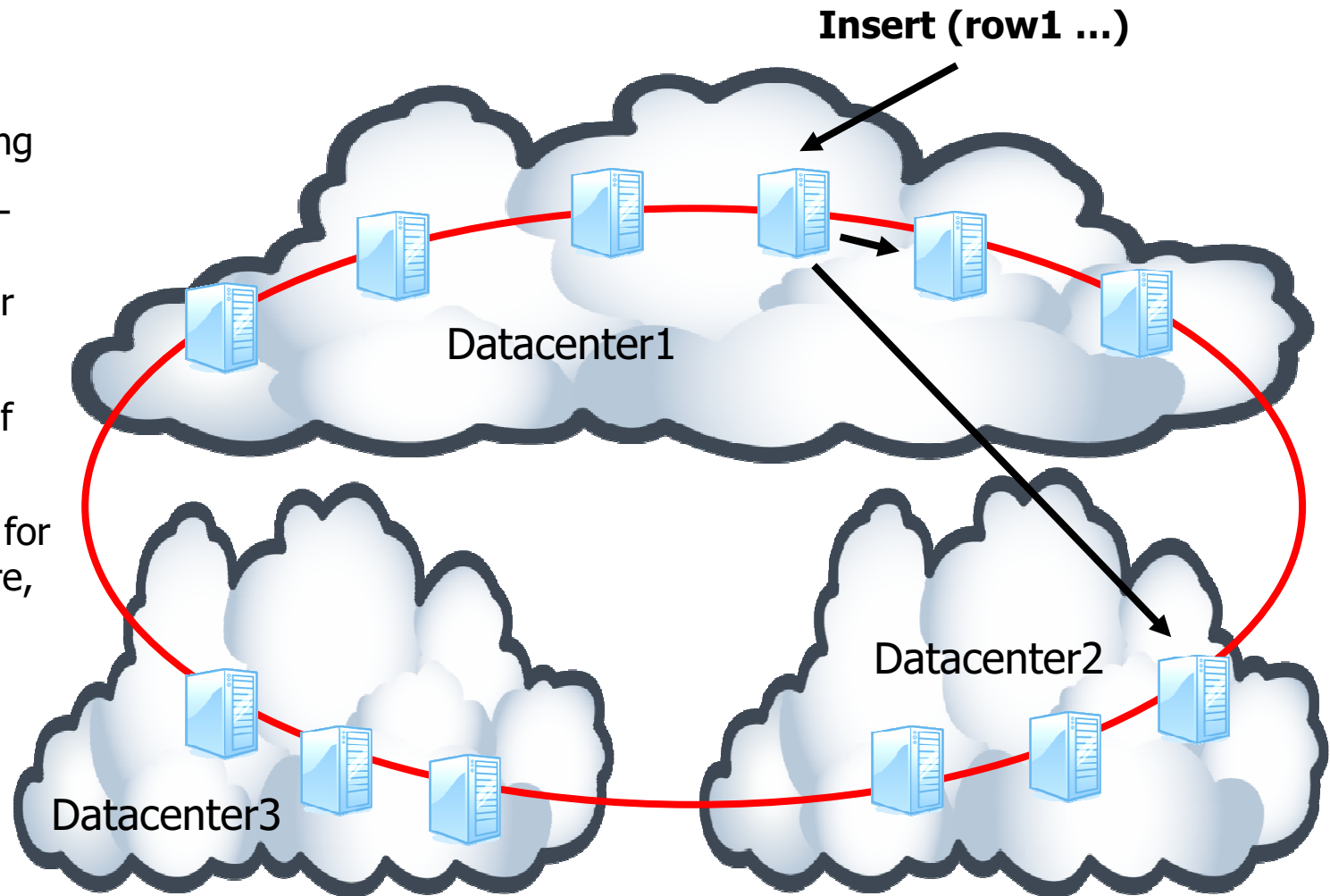
# Cassandra Data Model

- Familiar relational tables, rows, and columns, but more flexible
  - No upfront schema required
  - New columns can be added any time and columns can vary from row to row

|  | **row<br>key** | **col<br>name** | **col<br>value** | | |
|---|---|---|---|---|---|
| **row 1** | k127 | type: capacitor | farads: 12mf | cost: $1.05 | |
| **row 2** | k187 | type: resistor | ohms: 8k | label: banded | cost: $.25 |
| **row 3** | k217 | ... | ... | | |

- Columns grouped into Column Families
  - Column families are stored separately (like vertical partitioning)

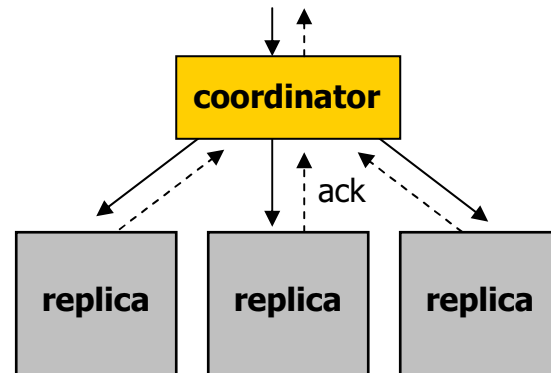# Cassandra Distributed Architecture

**Insert (row1 …)**

- Based on DHT ring

- Replication cross-rack and cross-datacenter (sync or async)

- No single point of failure

- Gossip protocols for membership, failure, DHT map, etc

Datacenter1

Datacenter2

Datacenter3

Scale, Fault-tolerance, Elasticity, Low-cost

# Eventual Consistency

- **CAP Theorem** [Brewer00]
  - Can only get 2 of **C**onsistency, **A**vailability, or **P**artition tolerance
- Cassandra relaxes **C** to eventual consistency
  - Emphasis is on performance and availability
  - Allow concurrent read/write to any replica – latest write wins on conflict
- **Knobs to tradeoff consistency and performance**
  - Writes are sent to all N replicas in parallel
  - Can choose to read from **R** replicas and wait for **W** acks for writes
  - Tune R and W to 1,2,3…,N for latency requirements

coordinator

ack

replica    replica    replica

coordinator acts as a
simple state machine
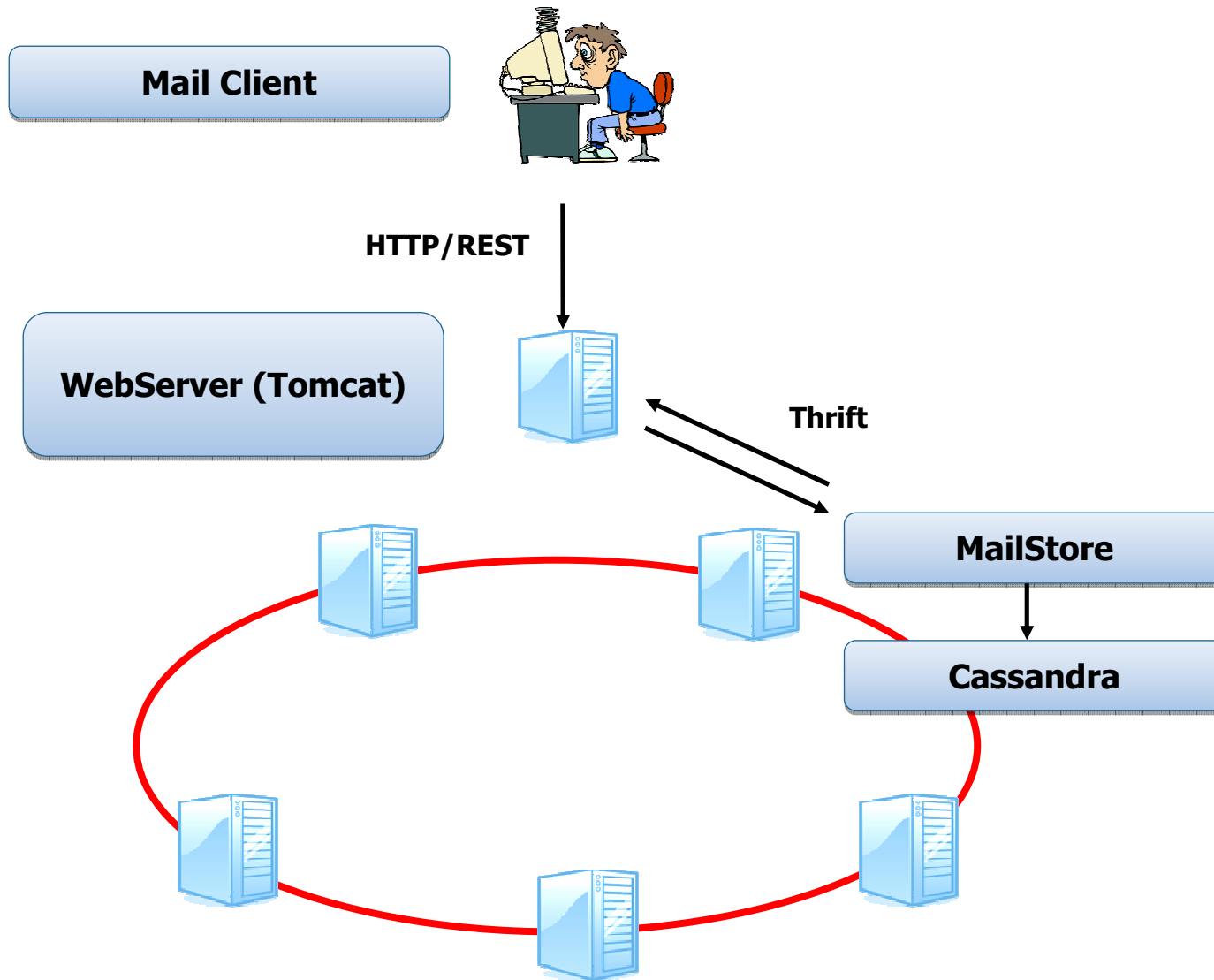
# Email Schema in Cassandra

- Row key - User id
  - All data for one user is on a single node
  - Currently no sharing of messages across users
- Column Families
  - MailList - Message id : full message
  - HeaderList - Message id : message headers + metadata
  - CollectionIndex - Collection id + sort key : message id
    e.g.,  Inbox/Date/2009-07-10-14:20:56 : message1000
           Inbox/Sender/Mike Brown : message1000
  - Others
    - CollectionMetadata, ThreadList, ThreadIndex
- Full message stored separately from index and metadata
- Data format - JSON

# Typical Operations

- Cassandra APIs used
  - get_column(row, CF, column)
  - get_columns(row, CF, columns[])
  - get_slice(row, CF, startColumn, asc/desc, count)
  - Efficient with row/column index support in Cassandra
- ListMessages
  - get_slice(Jun, CollectionIndex, Inbox/Date/current_date/, desc, 50) to obtain the first 50 messageIDs in Inbox
  - get_columns(Jun, HeaderList, messageID[])
- GetMessage
  - get_column(Jun, MailList, messageID)
- SortMessages by Sender
  - get_slice(Jun, CollectionIndex, Inbox/Sender//, asc, 50)
  - get_columns(Jun, HeaderList, messageID[])

# BlueRunner Deployment

**Mail Client**

**WebServer (Tomcat)**

HTTP/REST
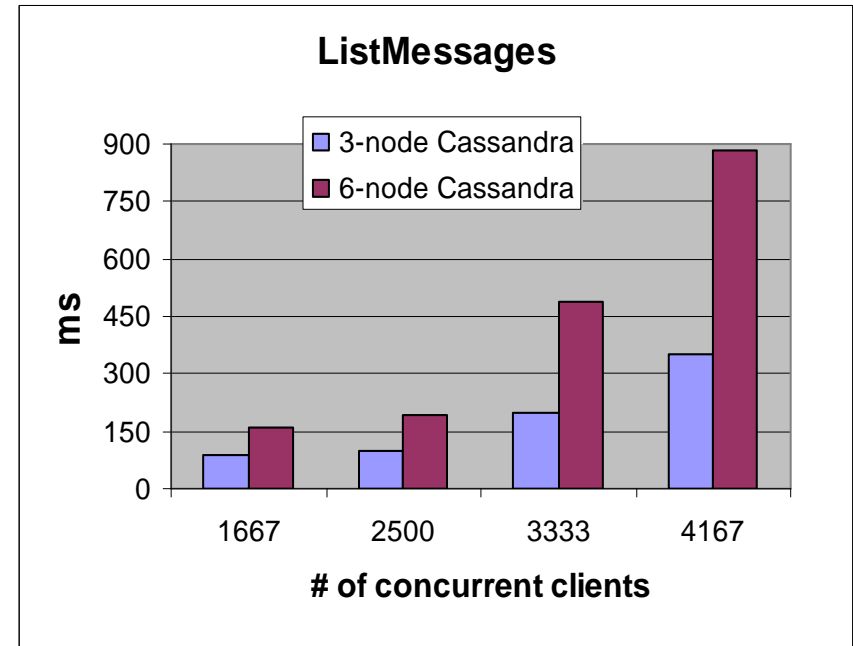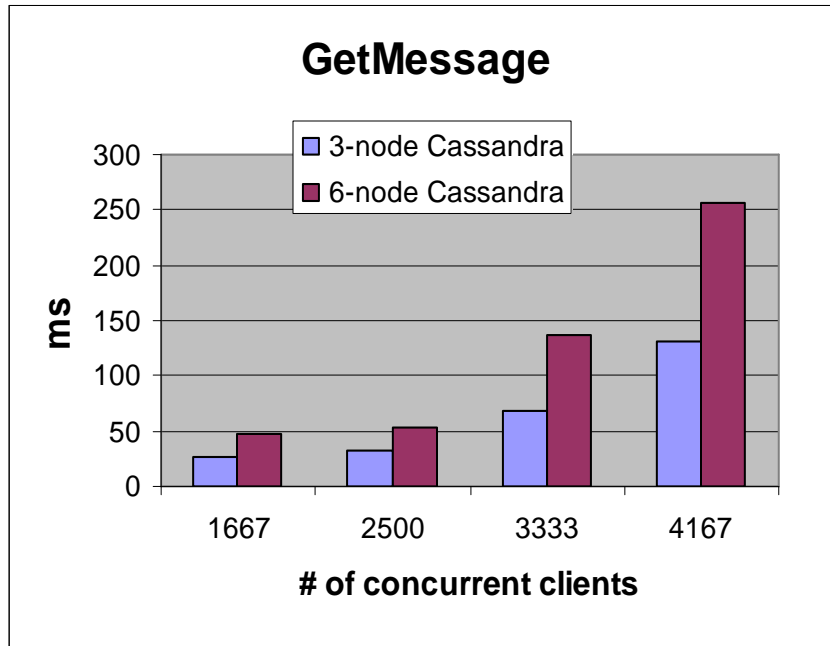
Thrift

**MailStore**

**Cassandra**

# Experimental Setup

- 6-node cluster and each node with
  - 2 quad-core CPUs
  - 16 GB memory
  - 5 SATA disks (1 for Cassandra commit log and 4 for data)
- Data
  - generated 1800 mailboxes per node
    - 250-16K messages per mailbox
    - ~50GB w/o replication
  - Cassandra replication set to 2
- Workload
  - Varying # of concurrent clients from 1600 to 4100 per node
  - Each client repeatedly
    - opens up an inbox
    - looks at a few message
    - go to sleep for a while

# Preliminary Result (median response time in ms)



**GetMessage** — ms vs # of concurrent clients (3-node Cassandra, 6-node Cassandra)

**ListMessages** — ms vs # of concurrent clients (3-node Cassandra, 6-node Cassandra)

- Able to sustain 2500 concurrent clients per node with reasonable response time
  - average: ~100 requests per sec per node

# Summary

- Cassandra-based backend very promising
  - Enabling scalability, availability, and elasticity
  - Flexible data model a good fit
- Future work
  - Many places can be improved
    - Alternative schema design
    - Secondary index/full-text index support
    - Enabling MapReduce-based analytics on the backend
  - Other potential collaborative apps on Cassandra
  - Research on better reasoning btw consistency and availability