

# Hadoop and Cassandra sitting in a tree...

Jake Luciani (@tjake)  
Strange Loop 2011



# K.I.S.S.I.N.G!?

Keeping

Integration

Simple

So

I

Never

Get-Woken-At-4am



# What makes Cassandra great?

Dynamo based

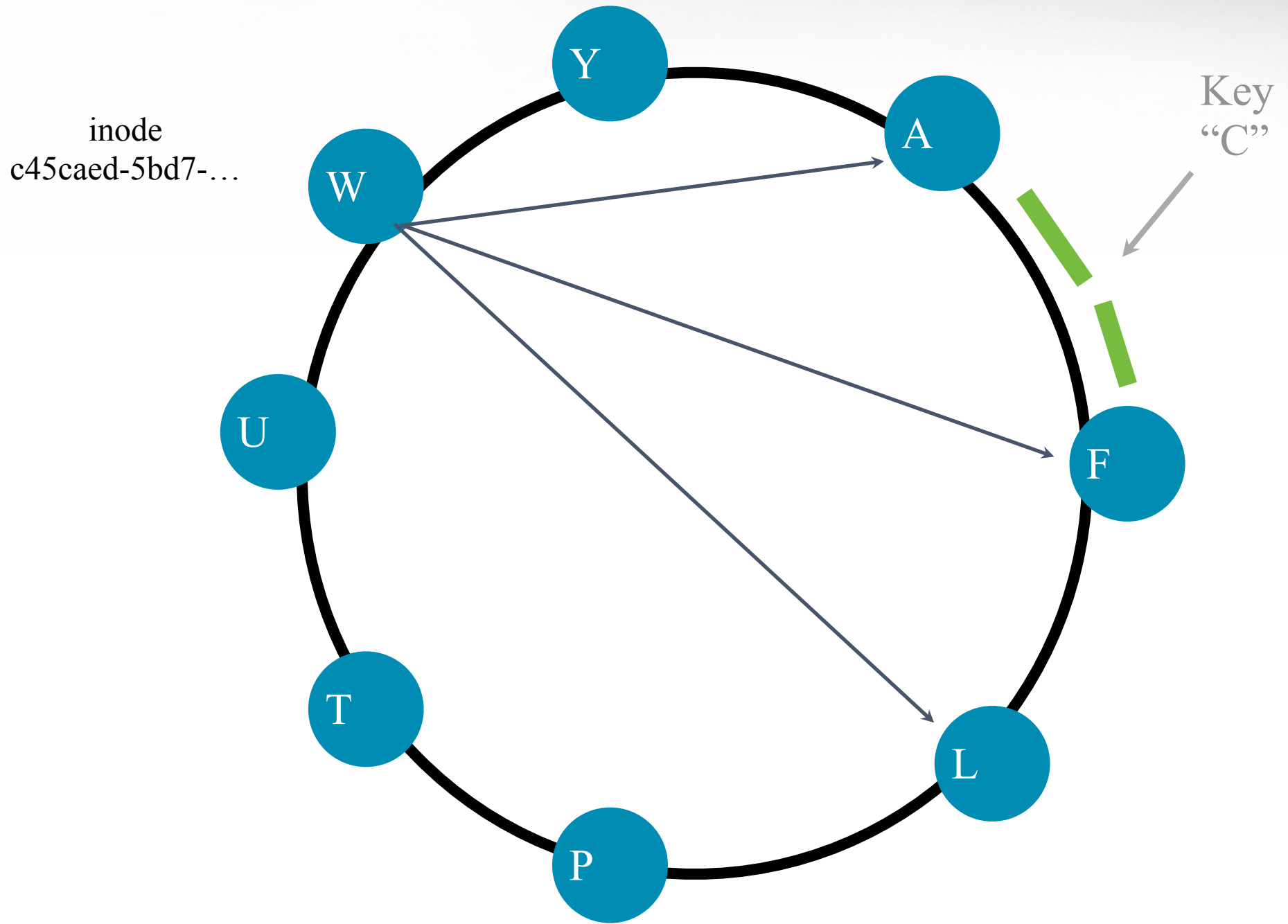
- Masterless
- Homogeneous
- Failure and Recovery at it's core
- Tunable consistency guarantees per read/write

BigTable based

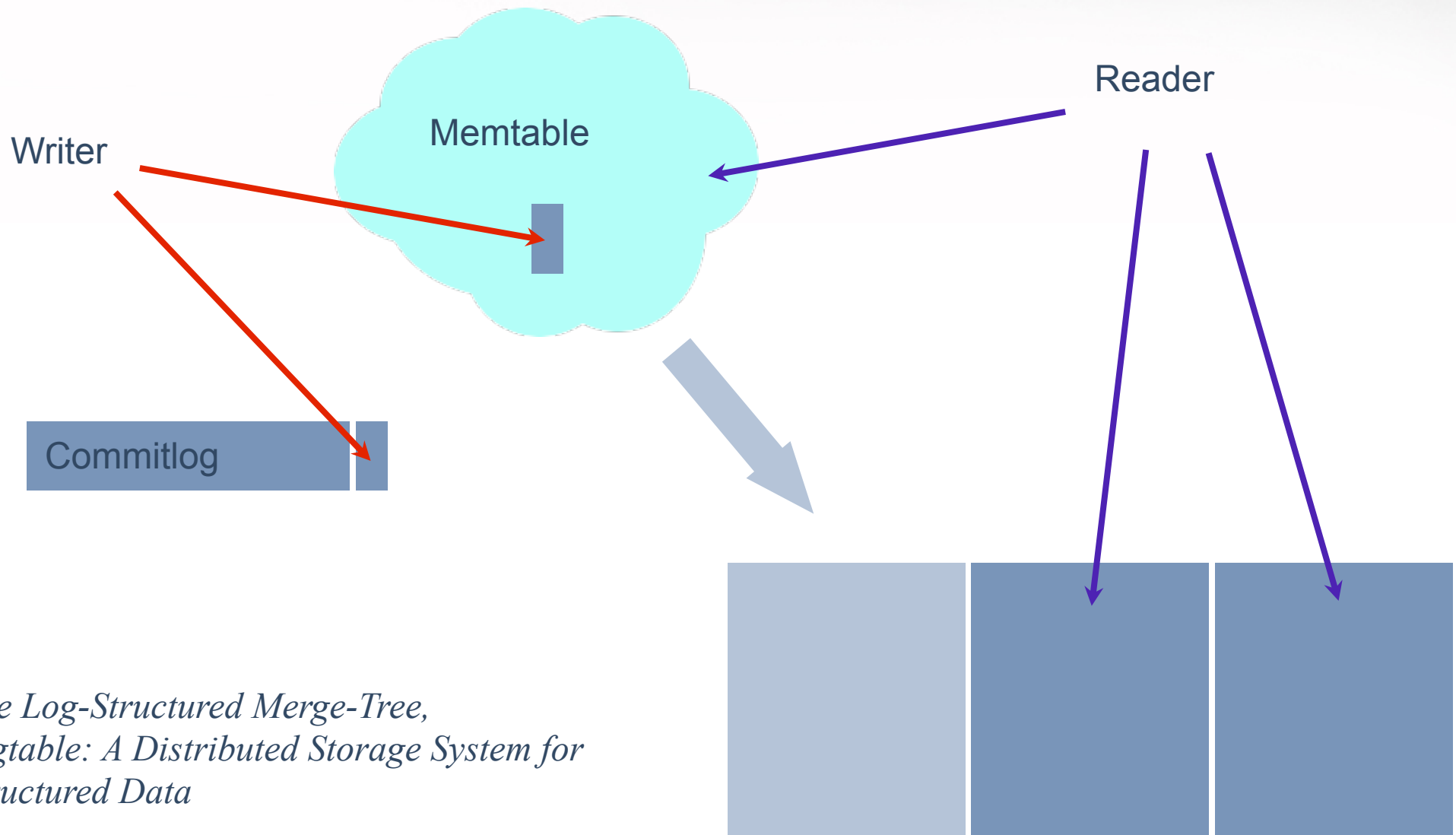
- Powerful data model
- Low Latency
- OLTP



# P2P replication (Dynamo)



# Locally-managed data (BigTable)



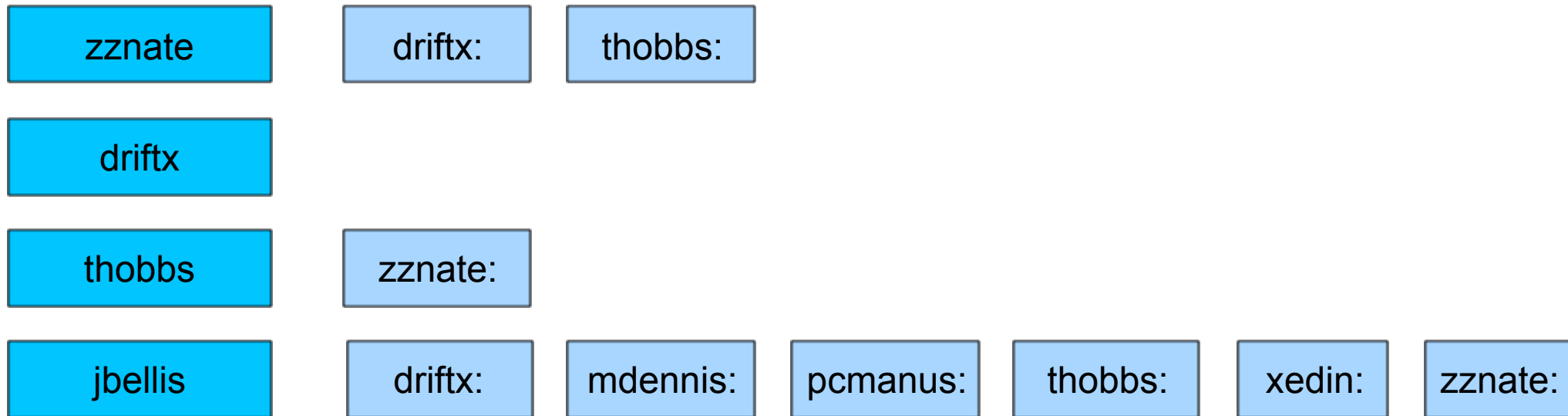
## Data Model (BigTable)

- ColumnFamilies contain rows + columns
- Sparse model: no rewriting for ALTER
- (Not really schemaless for a while now)

zznate	Password: *	Name: Nate	
driftx	Password: *	Name: Brandon	
thobbs	Password: *	Name: Tyler	
jbellis	Password: *	Name: Jonathan	Site: datastax.com

## Wide rows (BigTable)

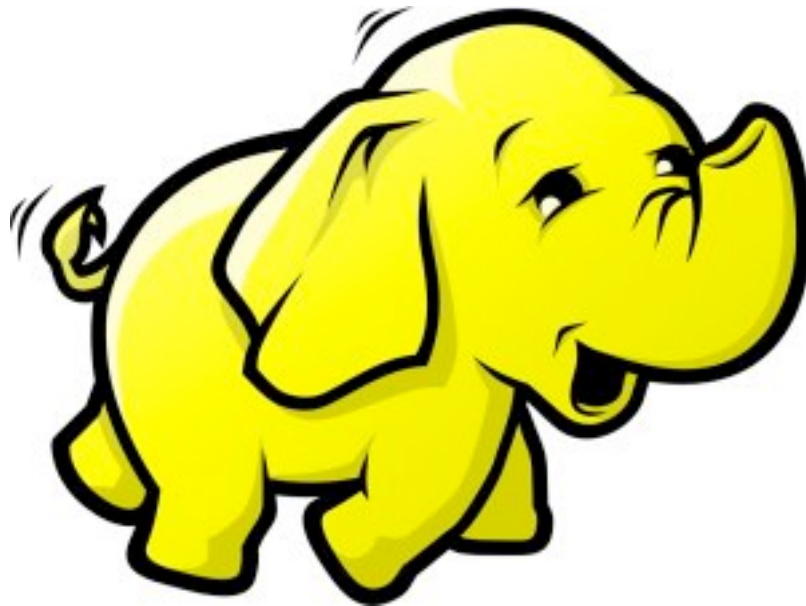
- Sparse model + sstable storage means we can have arbitrarily large rows
- Columns are sorted by comparator





# What makes Hadoop great?

- MapReduce Framework
- Pig
- Hive
- OLAP





# How can I use Cassandra with Hadoop?

- ColumnFamilyInputFormat
- ColumnFamilyOutputFormat
- Pig Driver
- Hive Driver
- Run TaskTrackers on Cassandra Nodes
  - Data Locality

# How do I run Hadoop jobs on data created in my application?

- *Run Hadoop on all nodes and access the data live?*

# How do I run Hadoop jobs on data created in my application?

- *Run Hadoop on all nodes and access the data live? **No. This will kill your OLTP performance.***
- *Use Scribe, Flume or Sqoop to put the data into HDFS, run my Hadoop jobs then load the data back into yet another system to consume the data?!*

# How do I run Hadoop jobs on data created in my application?

- *Run Hadoop on all nodes and access the data live?* **No. This will kill your OLTP performance.**
- *Use Scribe, Flume or Sqoop to put the data into HDFS, run my Hadoop jobs then load the data back into yet another system to consume the data?!* **Maybe. But there must be a better way.**
- Replicate the data from my live cluster to a secondary cluster, run Hadoop then replicate results back to the live cluster? **Yes!**

# The Traditional Hadoop Stack

## Master Nodes

Name Node

Secondary Name Node

Job Tracker

Hbase Master

ZooKeeper

MetaStore

## Slave Nodes

Data Node

Task Tracker

Region Server

## Client Nodes

Pig

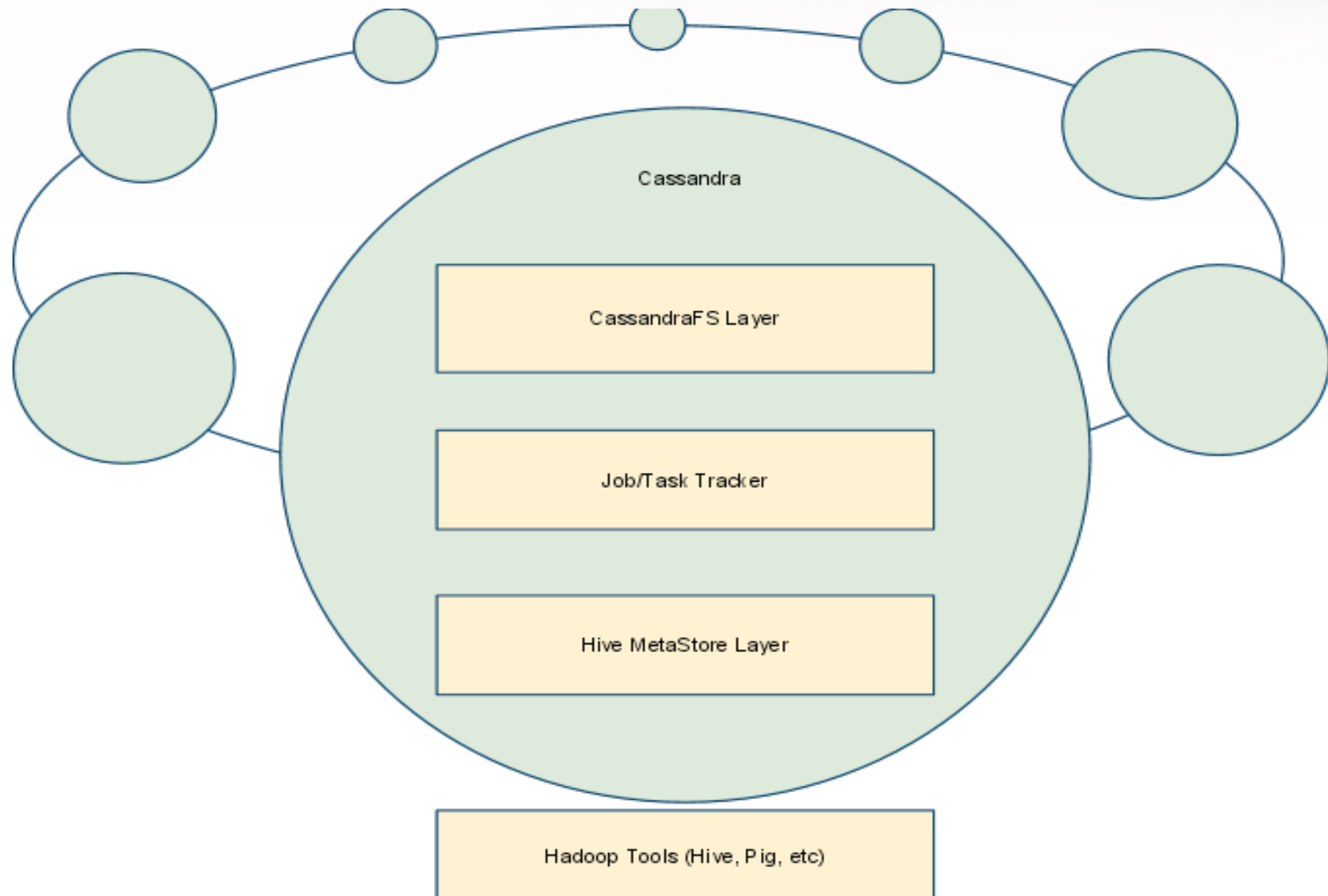
Hive

Region Server





# The Brisk Hadoop Ring





## Brisk

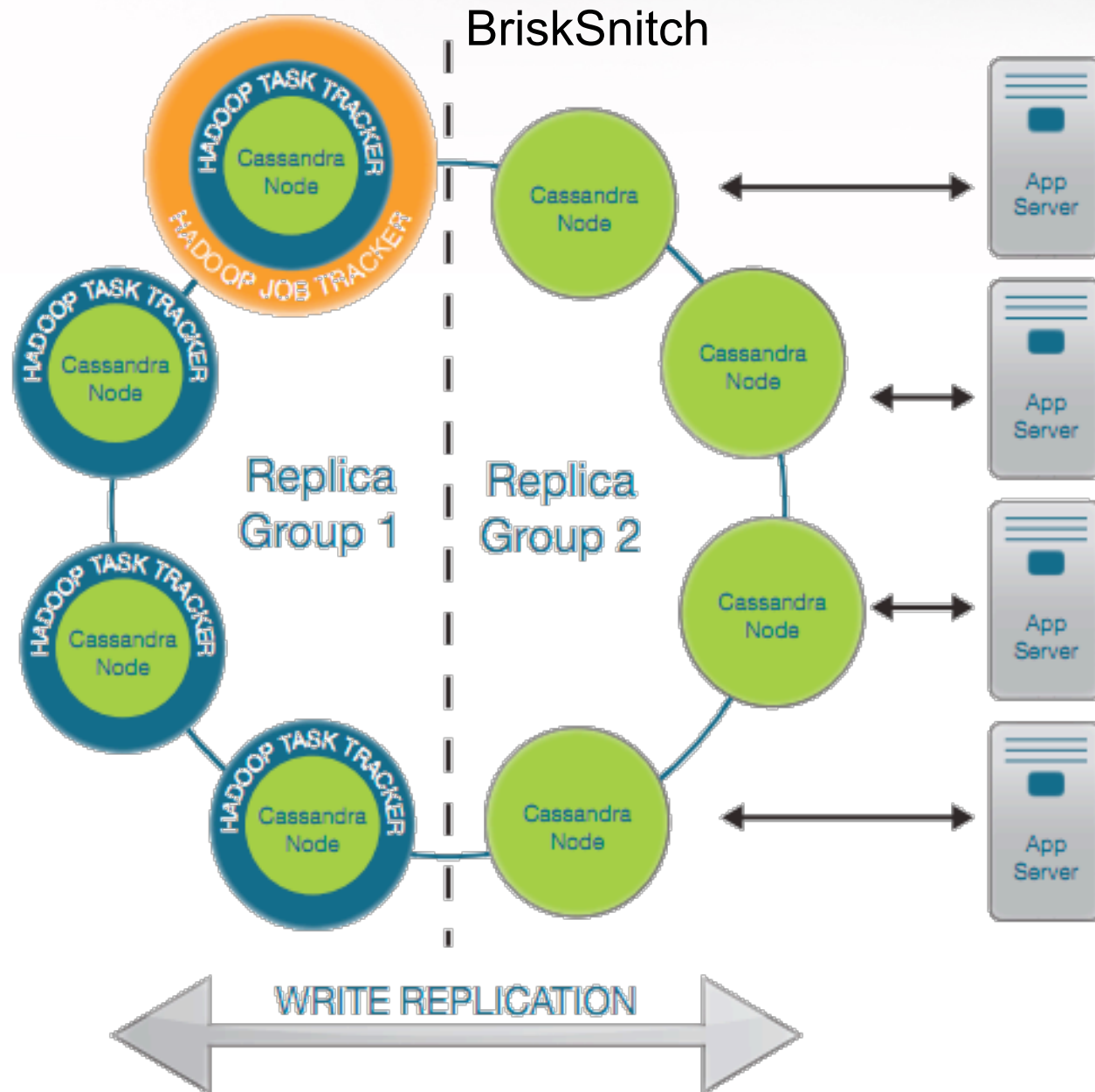
- Easy to deploy and operate
- No single points of failure
- Scale and change nodes with no downtime
- Cross-DC, multi-master clusters
- Allocate resources for OLAP vs OLTP
  - With no manual ETL
- Based on Cassandra 0.8 and Hadoop 0.20.203

# CassandraFS (HDFS Replacement)

- Built on ColumnFamilies
  - inode and sblocks
- Data stored as ByteBuffer internally — excellent fit for HDFS blocks
- Local reads mmap data directly (**no rpc**)
- Blocks are compressed with google snappy
- `hadoop distcp hdfs:///mydata cfs:///mydata`



# Hybrid Workloads



# Hive: CFS and ColumnFamilies

## CFS:

```
CREATE TABLE users (name STRING, zip INT);  
LOAD DATA LOCAL INPATH 'kv2.txt' OVERWRITE INTO TABLE users;  
select count(*), zip from invites group by bar;
```

## Column Family (fixed):

```
CREATE EXTERNAL TABLE Keyspace1.Users(name STRING, zip INT)  
STORED BY 'org.apache.hadoop.hive.cassandra.CassandraStorageHandler';
```

## Column Family (dynamic):

```
CREATE EXTERNAL TABLE Keyspace1.Users  
(row_key STRING, column_name STRING, value string)  
STORED BY 'org.apache.hadoop.hive.cassandra.CassandraStorageHandler';
```

## **Pig: CFS and ColumnFamilies**

### **CFS:**

```
grunt> data = LOAD 'cfs:///example.txt' using PigStorage() as  
(name:chararray, value:long);
```

### **ColumnFamily:**

```
data = LOAD 'cassandra://Demo1/Scores' using  
CassandraStorage() AS (key, columns: {T: tuple(name, value)});
```

### **ColumnFamily Slices:**

```
data = LOAD 'cassandra://Demo1/  
Scores&slice_start=M&slice_end=S' using CassandraStorage() AS  
(key, columns: {T: tuple(name, value)});
```



# Portfolio Manager Demo

- **Low-Latency Side (OLTP)**
  - Live tick prices for NASDAQ stocks
  - Thousands of Portfolios
- **Batch Analytics Side (OLAP)**
  - Historical Stock prices
  - Historical VAR (Value At Risk) calculation for each portfolio using Hive
  - VAR number is written back to Cassandra from Hive



# Portfolio Demo ColumnFamilies

## Portfolios

Portfolio1	GOOG	LNKD	P	AMZN	AAPL
	5	7	50	100	4

## LiveStocks

AAPL	LAST
	\$950.00

## StockHist

GOOG	2011-01-01	2011-01-02	2011-01-03	2011-01-04
	\$2.99	\$1.00	\$9.29	\$29.10

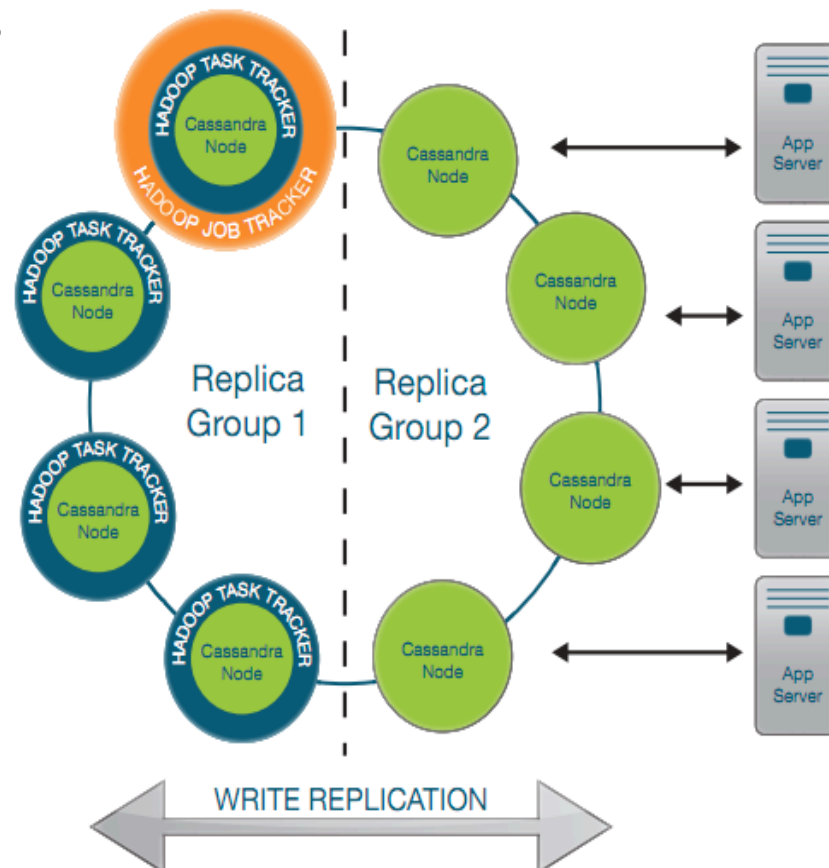
# Portfolio Demo Workloads

## OLAP:

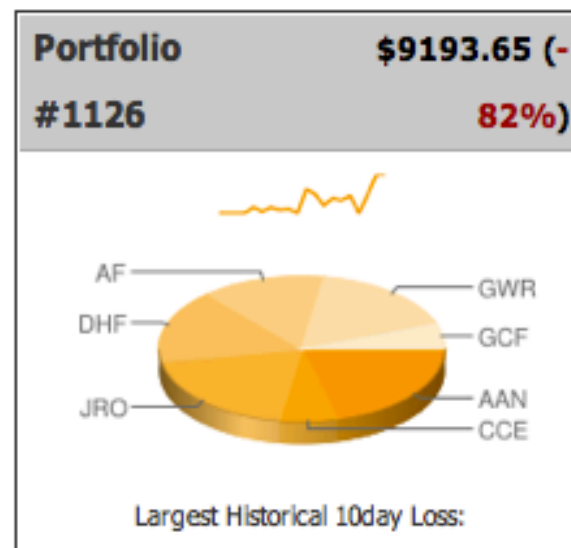
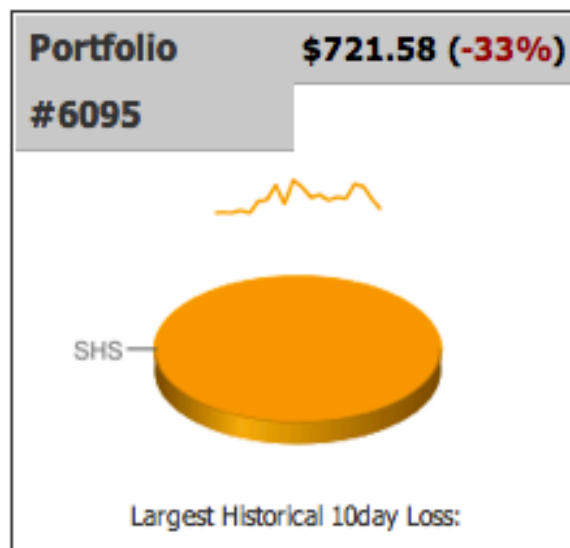
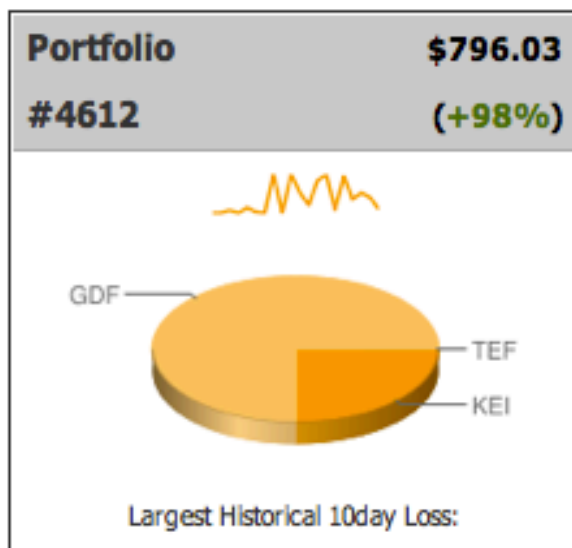
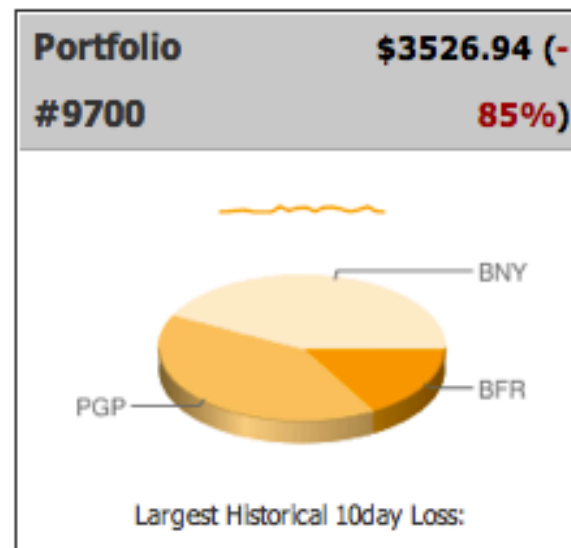
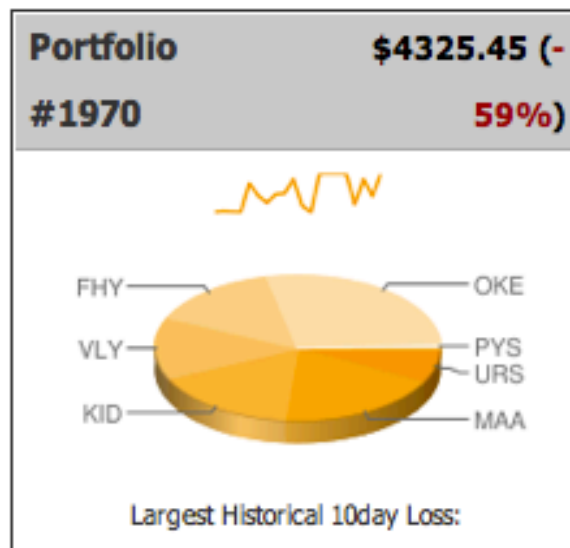
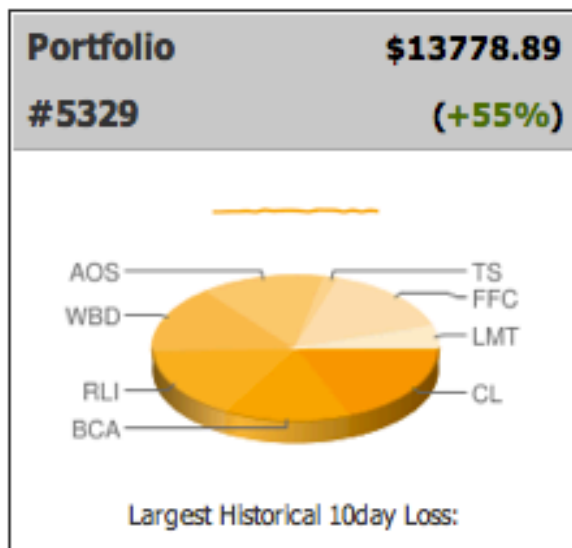
- Hive <- Cassandra
  - Portfolios
  - Historical Prices
- Intermediate Results
- Hive -> Cassandra
- Every N minutes

## OLTP:

- WebBased Portfolios
- Live Prices for today
- Historical Prices



# Portfolio Manager Demo



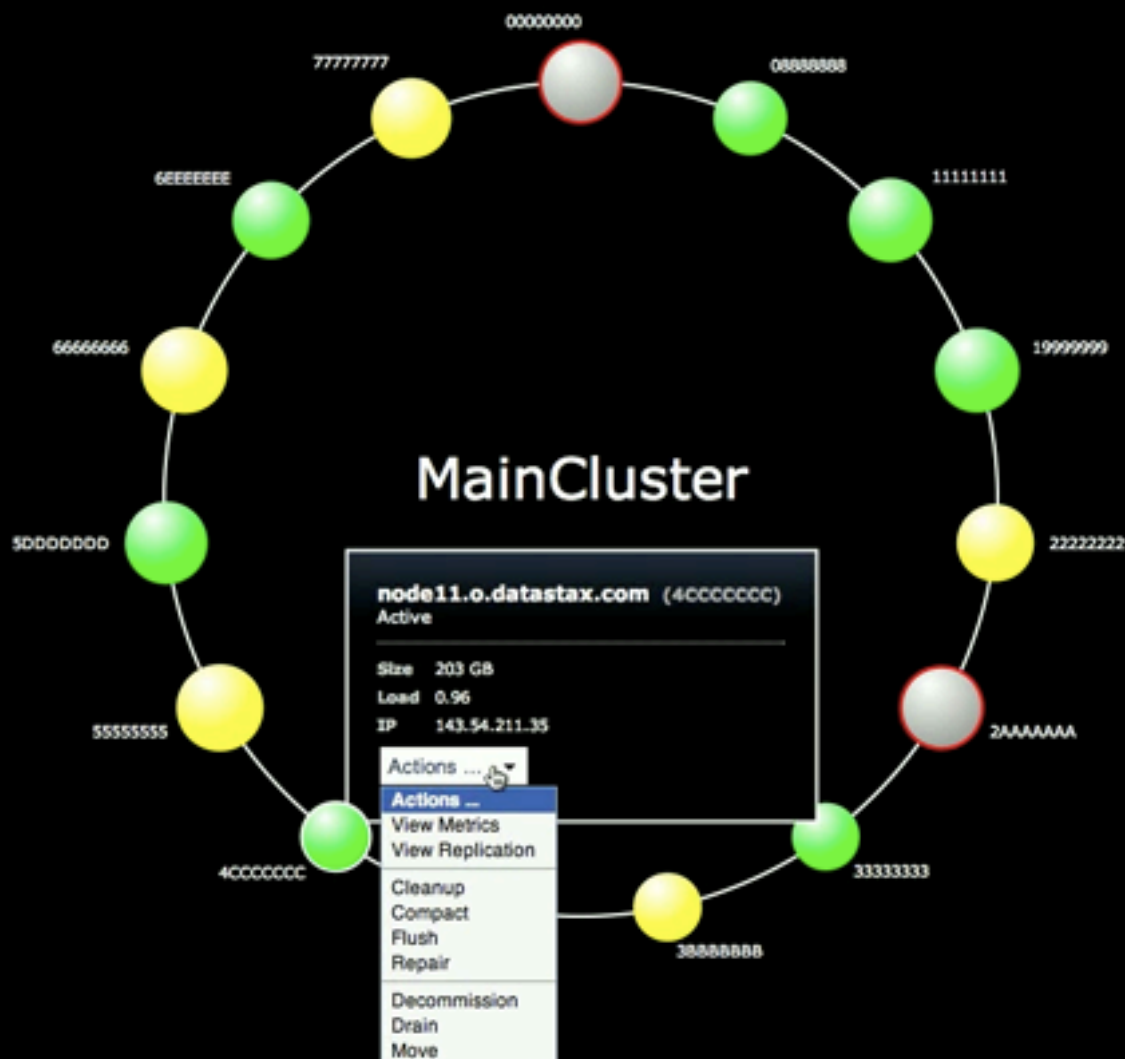
Ring

Physical

List

Add ▼

Balance Cluster



[Dashboard](#)
[Cluster Management](#)
[Events & Alerts](#)
[Performance](#)
[Data Modeling](#)
[Data Explorer](#)
[Hadoop Jobs](#)
[Refresh](#)

Auto-refreshes every minute.

[View Full Details](#)

	Job	Progress	Started	Duration	User
	TeraSort	<div> <div>33%</div> <div>Maps: 1/2 Reduces: 0/1</div> </div>	5/12/11 12:27 PM	40s	root
	TeraSort	<div> <div>100%</div> <div>Maps: 2/2 Reduces: 1/1</div> </div>	5/12/11 12:26 PM	34s	root
	TeraGen	<div> <div>100%</div> <div>Maps: 2/2 Reduces: 0/0</div> </div>	5/12/11 12:25 PM	44s	root
	PiEstimator	<div> <div>100%</div> <div>Maps: 10/10 Reduces: 1/1</div> </div>	5/12/11 12:19 PM	49s	riptano
	select count (*) from server_request(Stage-	<div> <div>100%</div> <div>Maps: 15/15 Reduces: 1/1</div> </div>	5/12/11 10:56 AM	1m 31s	riptano
	select count(*) size, method from myT...size(Stage-2)	<div> <div>100%</div> <div>Maps: 10/10 Reduces: 10/10</div> </div>	5/11/11 8:30 PM	57s	riptano
	select count(*) size, method from myT...size(Stage-1)	<div> <div>100%</div> <div>Maps: 2/2 Reduces: 10/10</div> </div>	5/11/11 8:29 PM	50s	riptano
	select count(*), method from myTest...method(Stagi	<div> <div>100%</div> <div>Maps: 2/2 Reduces: 10/10</div> </div>	5/11/11 8:24 PM	46s	riptano
	select count(*), method from myTest...method(Stagi	<div> <div>100%</div> <div>Maps: 2/2 Reduces: 10/10</div> </div>	5/11/11 8:12 PM	49s	riptano
	select count(*), method from myTest...method(Stagi	<div> <div>100%</div> <div>Maps: 2/2 Reduces: 1/1</div> </div>	5/11/11 8:04 PM	21s	riptano
	select min(duration) from myTest(Stage-1)	<div> <div>100%</div> <div>Maps: 2/2 Reduces: 1/1</div> </div>	5/11/11 8:02 PM	24s	riptano
	select max(duration) from myTest(Stage-1)	<div> <div>100%</div> <div>Maps: 2/2 Reduces: 1/1</div> </div>	5/11/11 8:02 PM	24s	riptano
	select count(*) from myTest(Stage-1)	<div> <div>100%</div> <div>Maps: 2/2 Reduces: 1/1</div> </div>	5/11/11 8:00 PM	25s	riptano

# Open Source

- <https://github.com/riptano/brisk>
- Apache License

**Thank You**

Jake Luciani

[jake@datastax.com](mailto:jake@datastax.com)

@tjake

