



Supercharging Cassandra...

Tom Wilkie
Founder & VP Engineering
[@tom_wilkie](https://twitter.com/tom_wilkie)



Acunu

Before the Flood

1990

Small databases

BTree indexes

BTree File systems

RAID

Old hardware



Acunu

Two Revolutions

2010

Distributed, shared-nothing databases

Write-optimised indexes

Write-optimised indexes

BTree file systems

BTree file systems

RAID

RAID

...

New hardware

New hardware



Acunu

Bridging the Gap

2011

Distributed, shared-nothing databases

Castle

Castle

...

New hardware

New hardware



Acunu



Acunu

- Open API
- Management
- Deployment
- Monitoring

Cassandra



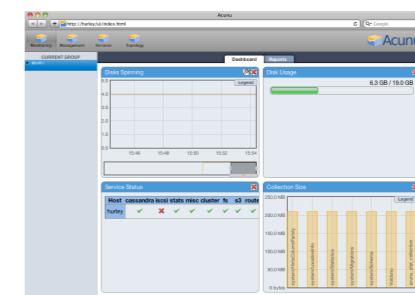
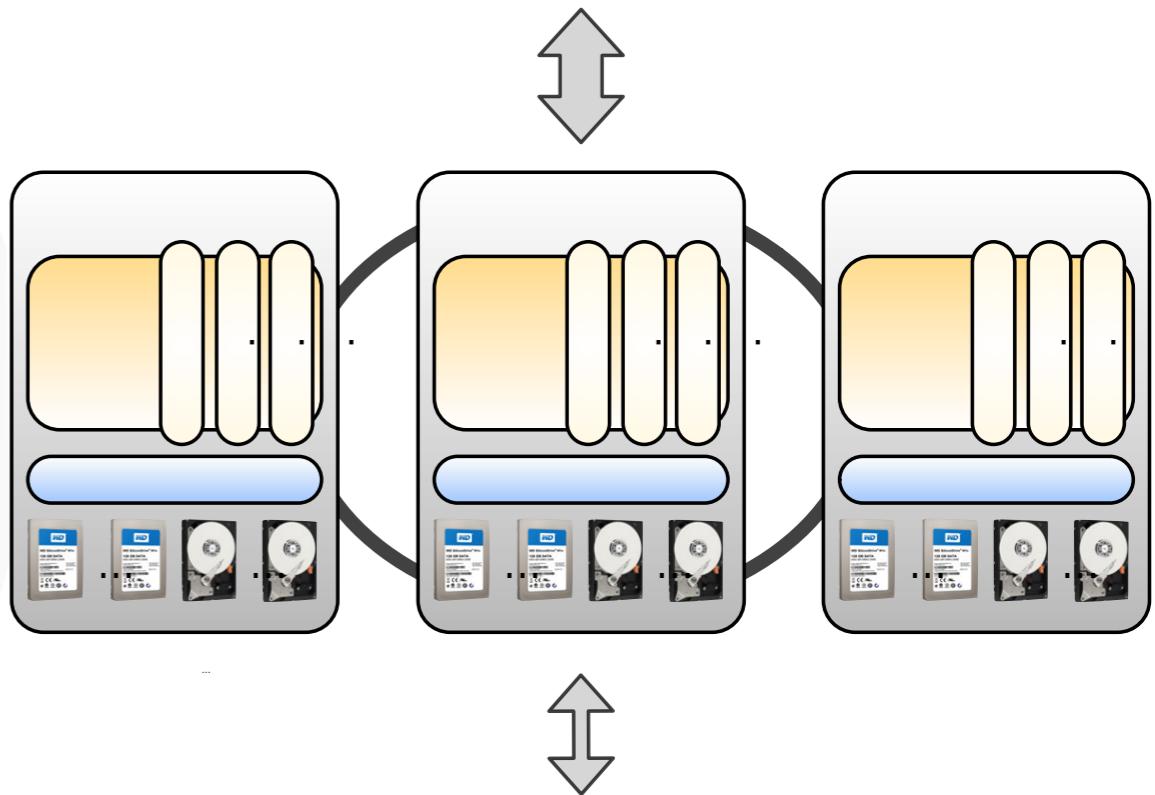
Memcached



Acunu Storage Core



Big Data
Applications



Cross-Cluster Management UI



Acunu

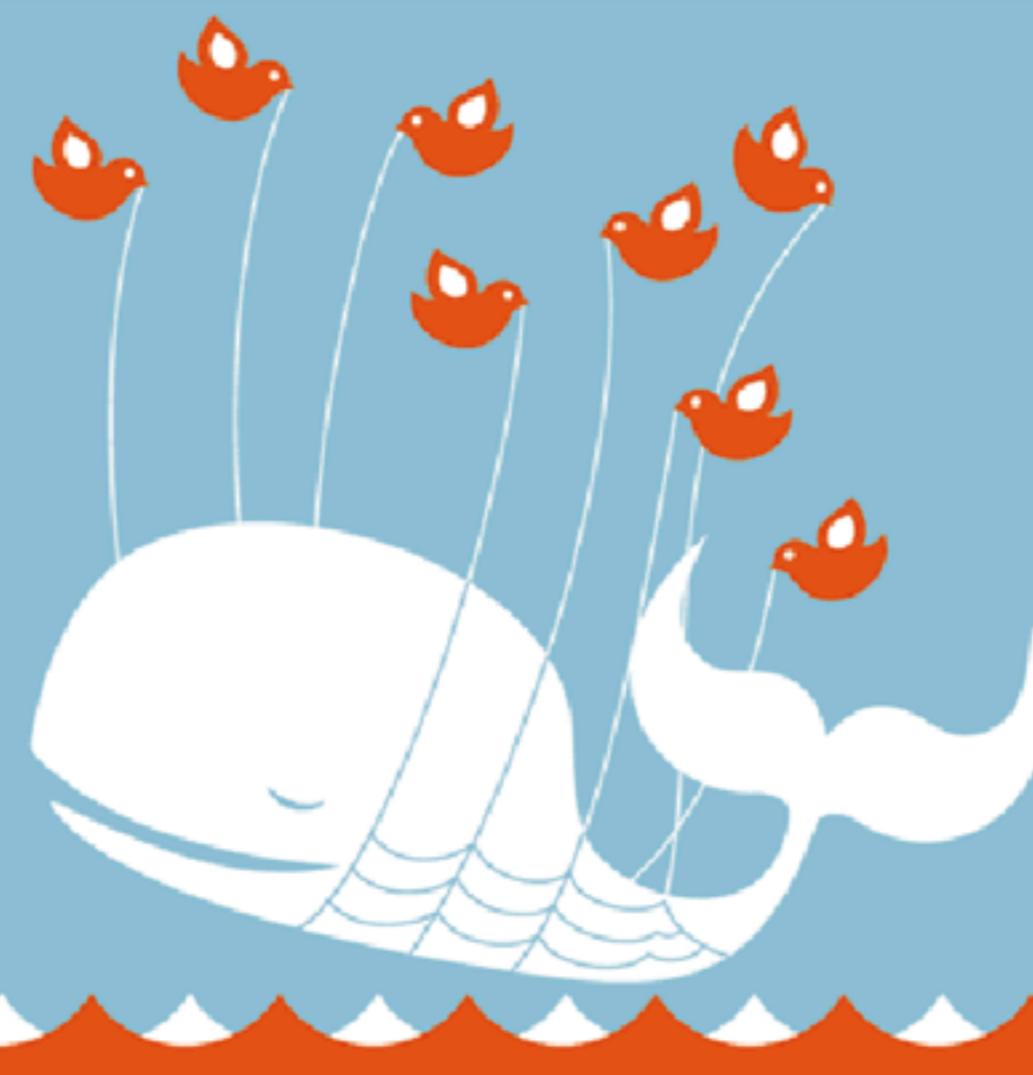
A vertical bar on the left side of the slide is divided into two horizontal sections: a yellow section at the top and a blue section at the bottom. The boundary between them is a smooth, curved line that starts from the left edge, dips slightly into the blue area, and then rises back towards the yellow area.

I. Predictability

[Home >](#)

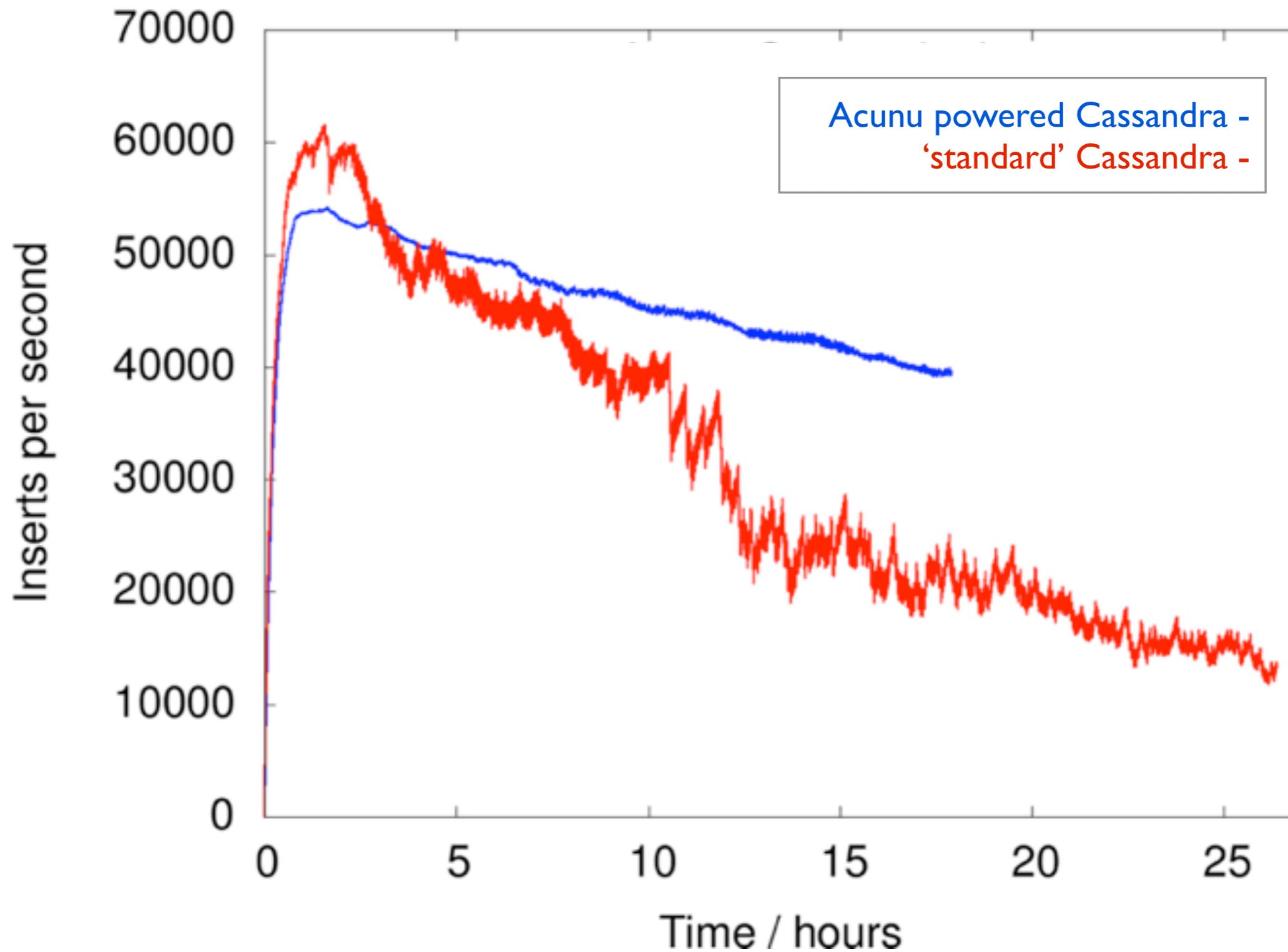
Twitter is over capacity.

Please wait a moment and try again. For more information, check out [Twitter Status »](#)

[English](#)[Deutsch](#)[Español](#)[Français](#)[Italiano](#)[日本語](#)

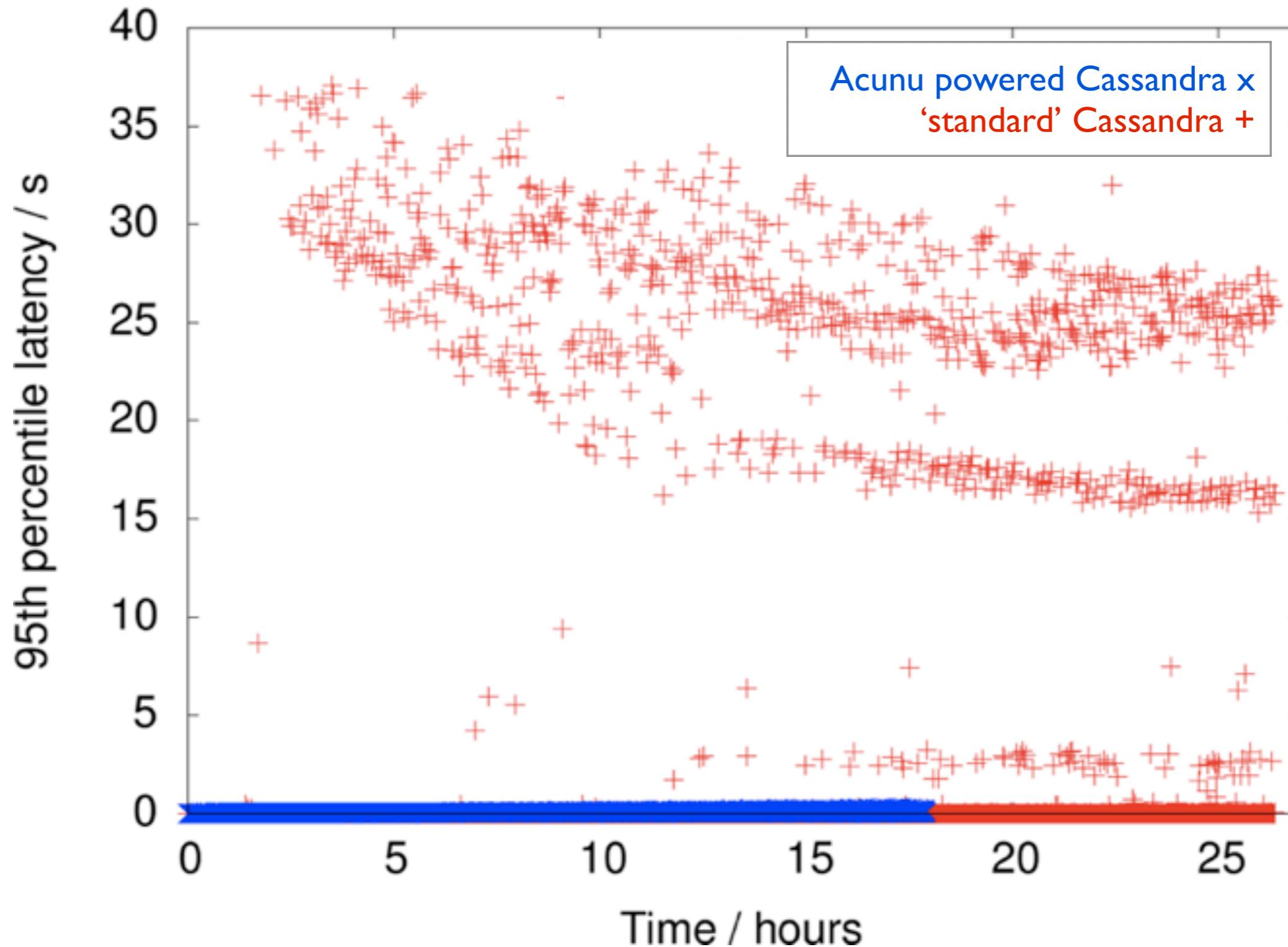
Small random inserts

Inserting 3 billion rows



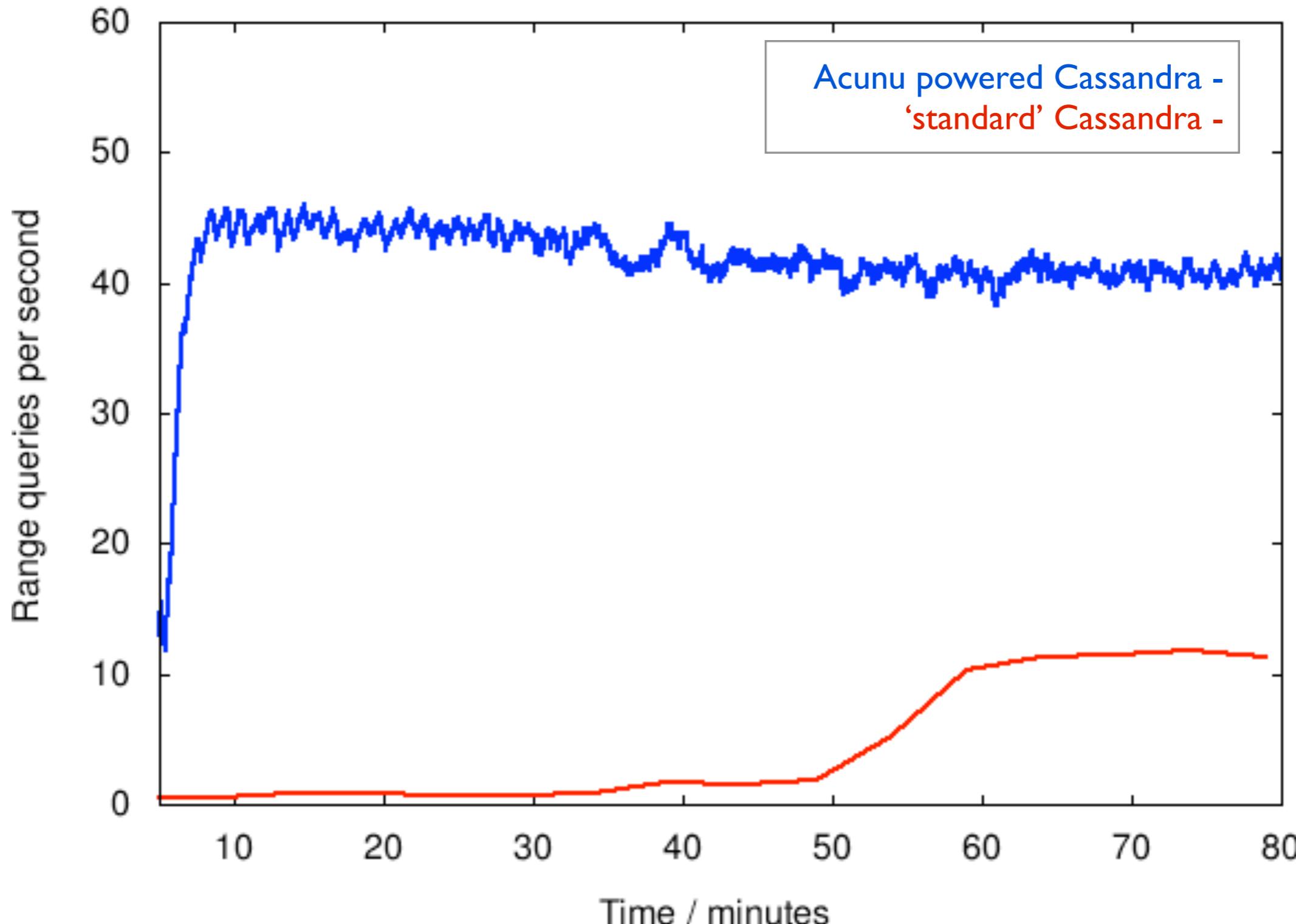
Insert latency

While inserting 3 billion rows



Small random range queries

Performed immediately after inserts

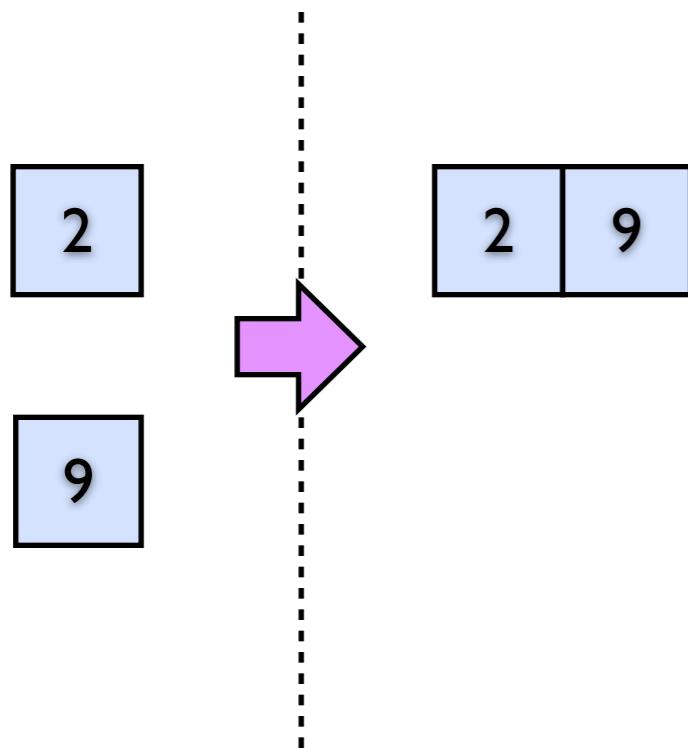


Performance summary

	Standard	Acunu	Benefits
inserts rate 95% latency	~32k/s ~32s	~45k/s ~0.3s	>1.4x >100x
gets rate 95% latency	~100/s ~2s	~350/s ~0.5s	>3.5x >4x
range queries 95% latency	~0.4/s ~15s	~40/s ~2s	>100x >7.5x

Doubling Array

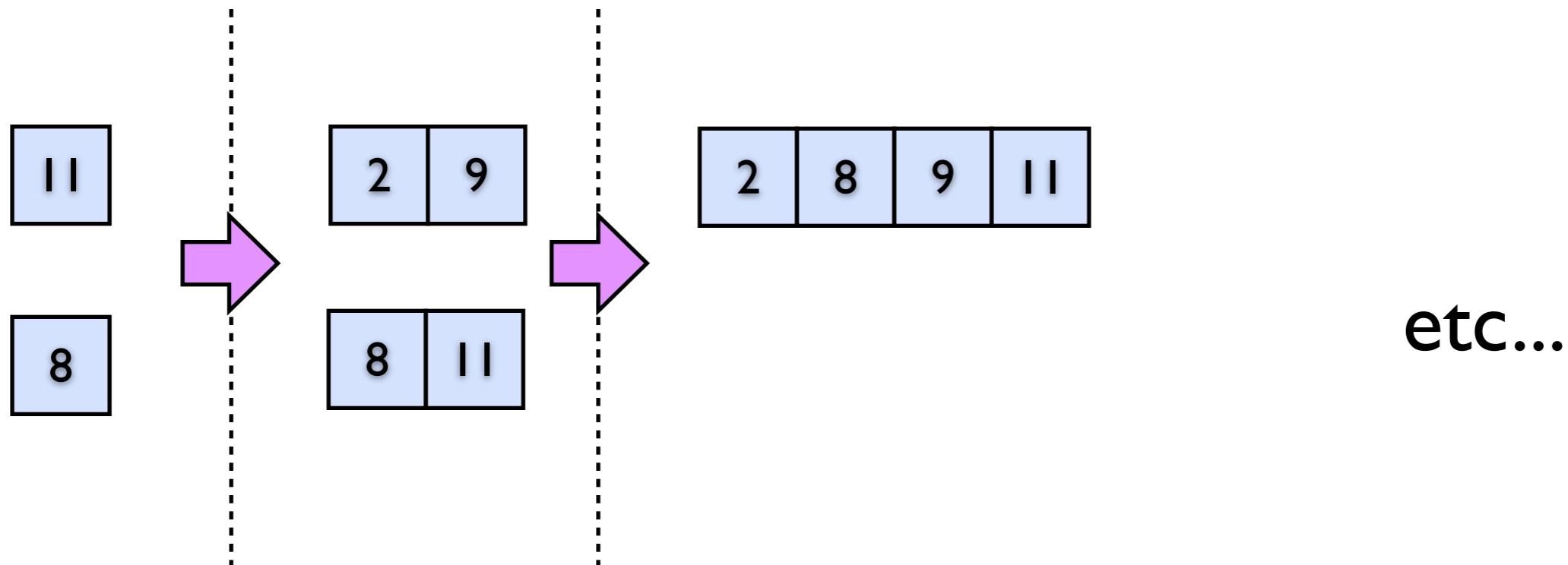
Inserts



Buffer arrays in memory
until we have $> B$ of them

Doubling Array

Inserts

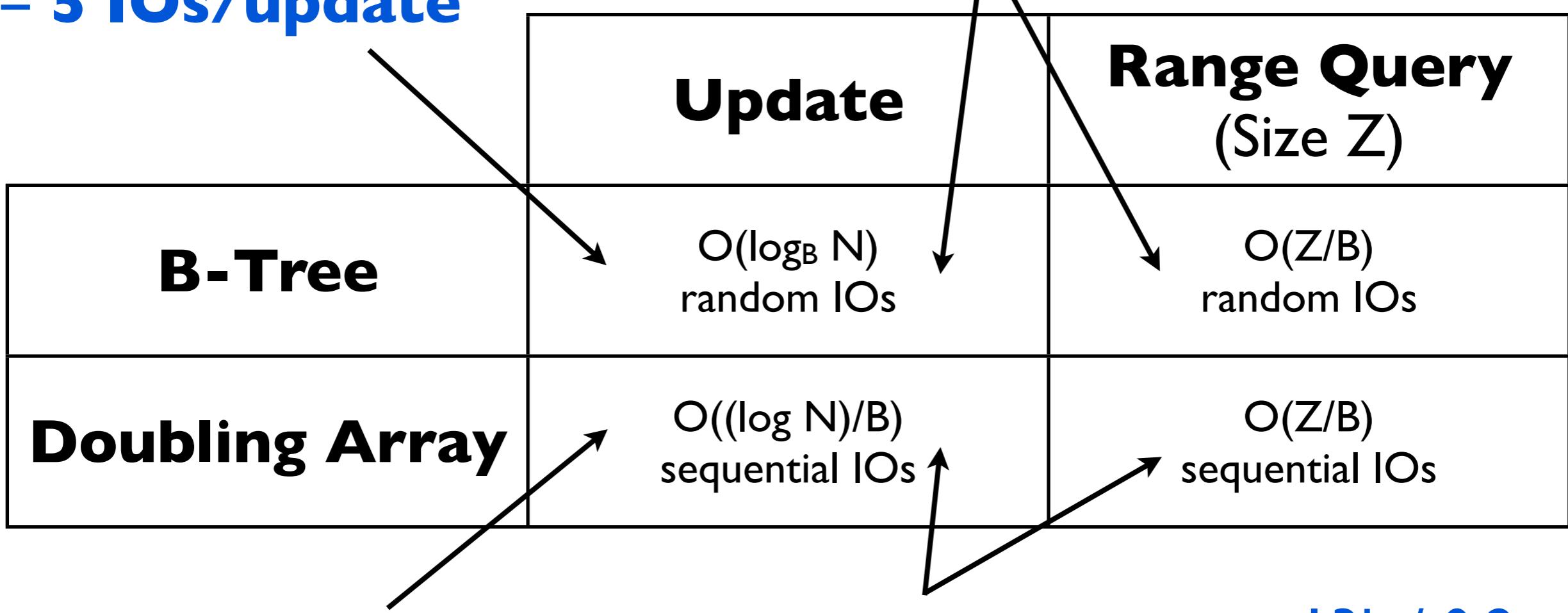


Similar to log-structured merge trees (LSM), cache-oblivious lookahead array (COLA), ...

Demo

<https://acunu-videos.s3.amazonaws.com/dajs.html>

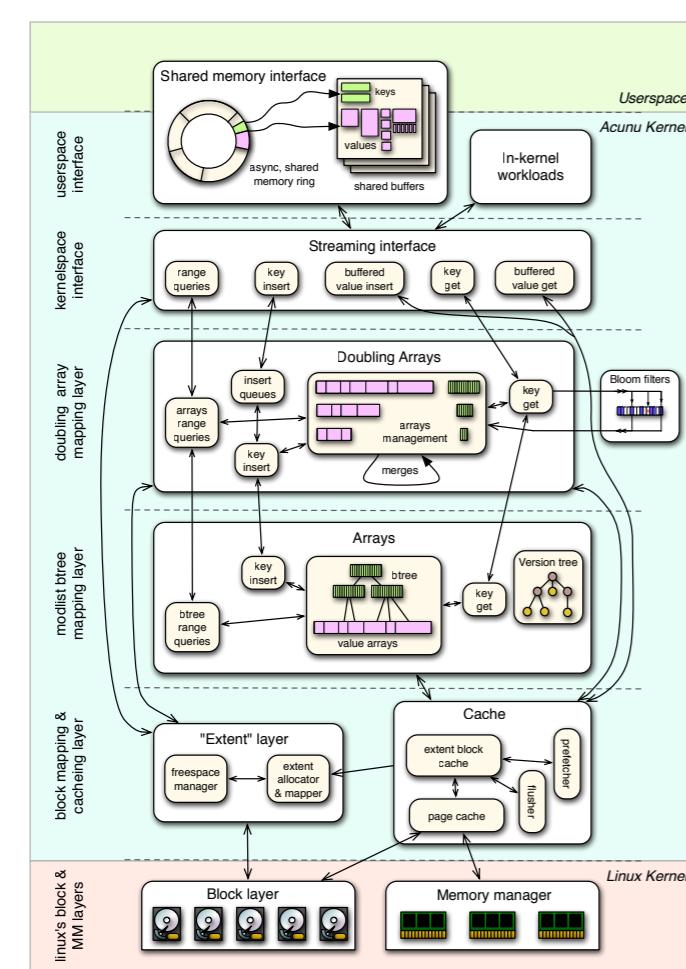
$8\text{KB} @ 100\text{MB/s, w/ } 8\text{ms seek}$
 $= 100 \text{ IOs/s}$
 $100 / 5$
 $\sim \log(2^{30})/\log 100$
 $= 5 \text{ IOs/update}$
 $= 20 \text{ updates/s}$



$\sim \log(2^{30})/100$
 $= 0.2 \text{ IOs/update}$
 $8\text{KB} @ 100\text{MB/s}$
 $= 13k \text{ IOs/s}$
 $13k / 0.2$
 $= 65k \text{ updates/s}$

$B = \text{"block size", say 8KB at 100 bytes/entry } \approx 100 \text{ entries}$

More



Castle

- Opensource (GPLv2, MIT for user libraries)
- <http://bitbucket.org/acunu>
- Loadable Kernel Module, targeting CentOS's 2.6.18
- <http://www.acunu.com/blogs/andy-twigg/why-acunu-kernel/>

SHOOTING
WALLS

<http://goo.gl/wXNDQ>

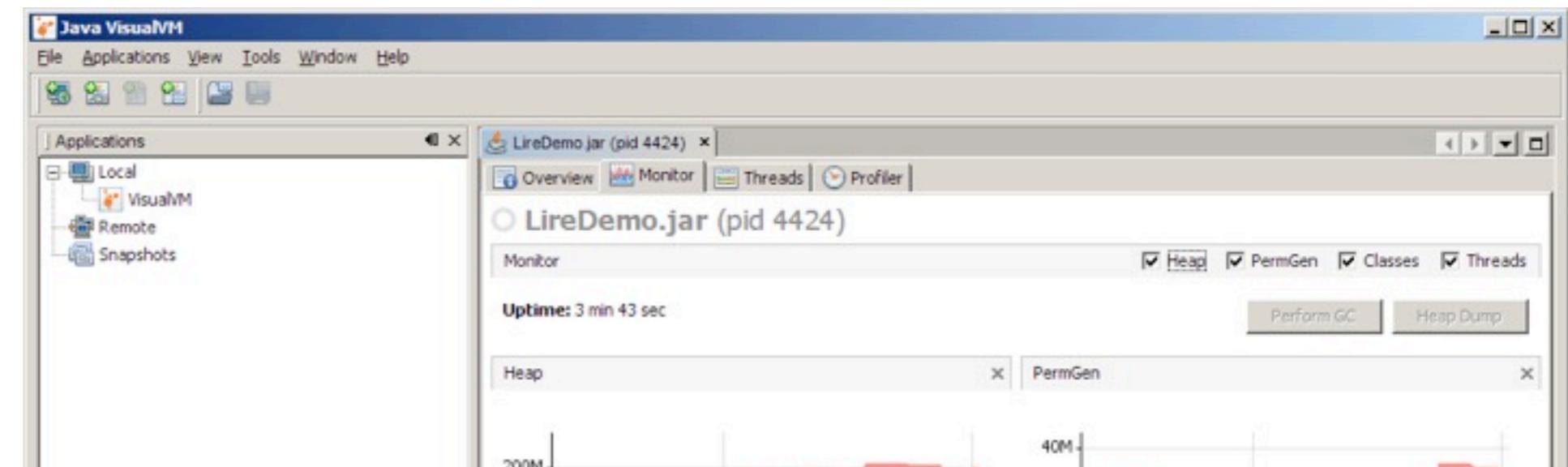
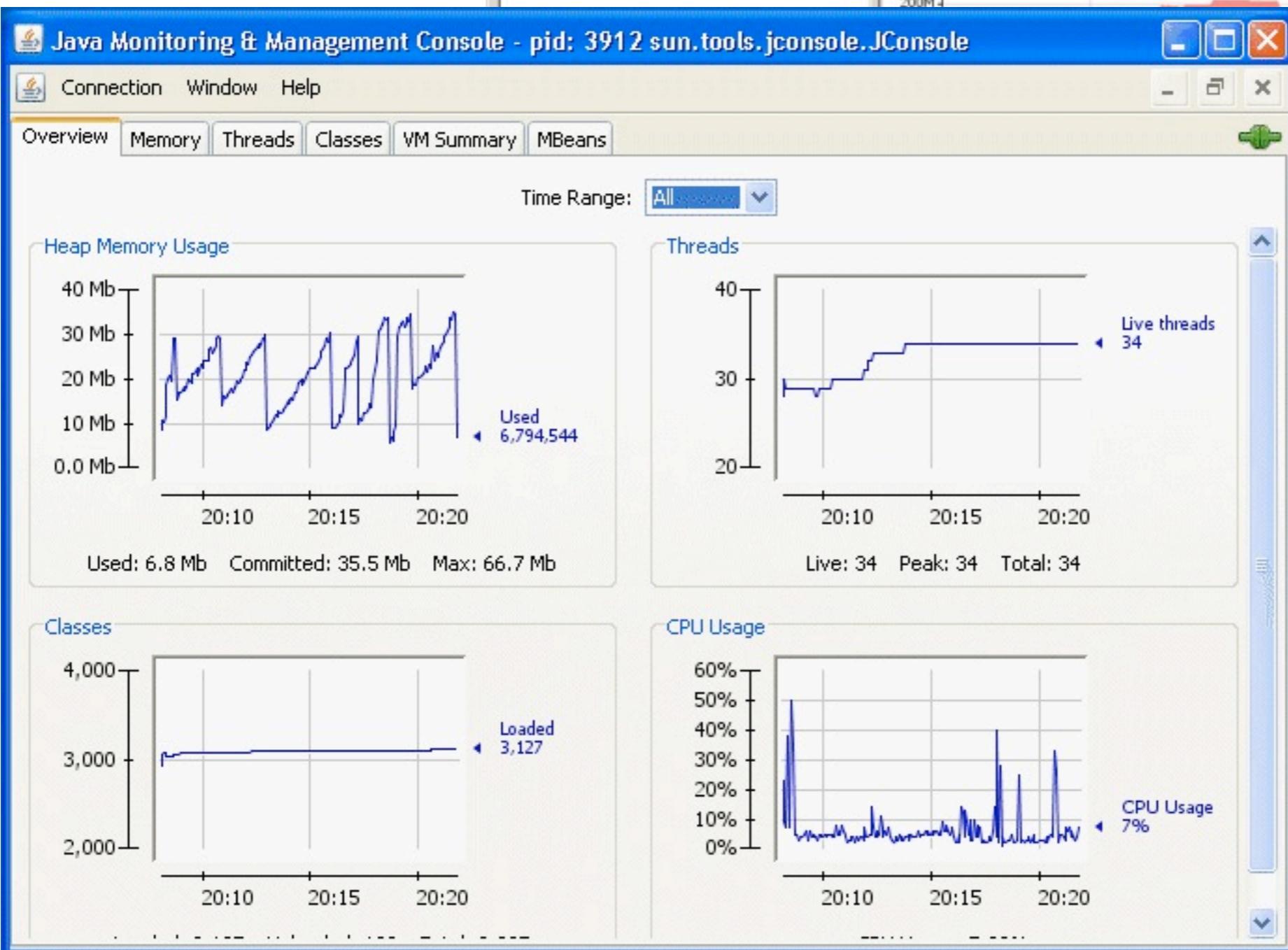


<http://goo.gl/gzihe>

A vertical decorative bar on the left side of the slide, divided into two horizontal sections: yellow at the top and blue at the bottom.

2. Monitoring

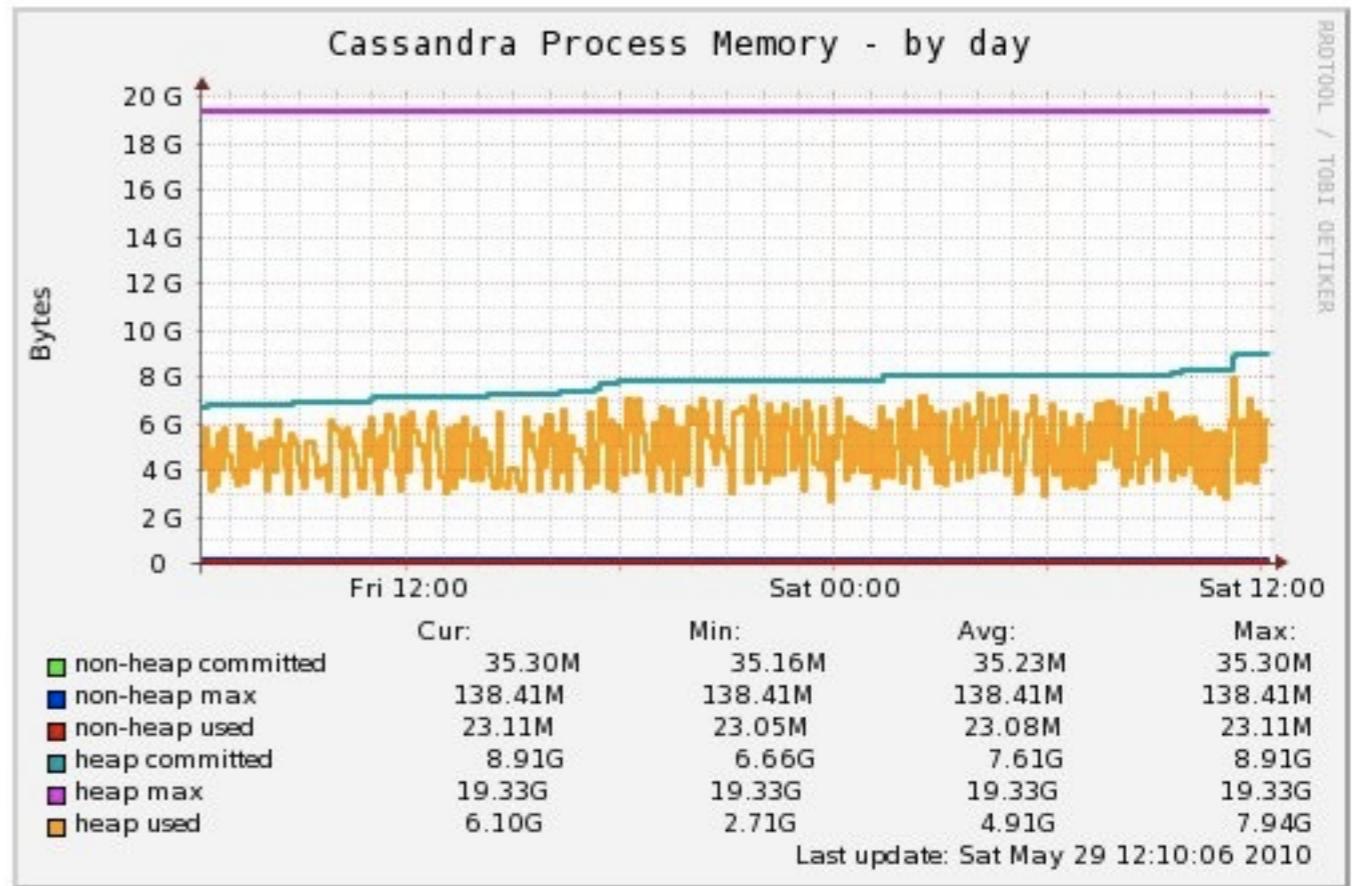
jQuery



VisualVM

Acunu

mx4j: Rest-JMX adapter



MX4J - MBean View - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

http://admc.com:8112/mbean?objectname=dom%3Aname

MX4J/Http Adaptor JMX Management Console

MX4J

Server view MBean view Timers Monitors Relations MLet About

MBean Description Attributes

Name	Description	Type	Value	New Value
AL	AL att descript.	java.util.ArrayList	View collection	Unknown type
D	D att descript.	double	7.8	<input type="text"/> set
DA	DA att descript.	ID	ID@20249	Unknown type
H	H att descript.	int	write-only attribute	<input type="text"/> set
HM	HM att descript.	java.util.HashMap	View map	Unknown type
I	I att descript.	java.lang.Integer	7	View array Read-only attribute Unknown type
SA	SA att descript.	Array of java.lang.String	View array	Set all

Operations

Munin, Nagios etc

hurley > Acunu Data Platform

http://hurley/ui/index.html#topology/data_centre/DC1

Google

Usual ▾ Recipes ▾ F1 ▾ Programming ▾ TeeVee ▾ Banks ▾ Acunu ▾ Quick ▾

Management Versions Topology Notifications

Acunu

RING

DATA CENTRES (1)

DC1

r1

hurley

« Hide Stats

Disk space (35.4 GB free)

CPU usage (last hour)

Disk graphs

Throughput (MB/s)

Disk iops

Disk Utilisation

/dev/xvda3 (2... 17.7 GB fr

Disk iops

Last Hour Last Day Last Week

Reads Writes

75 50 25 0

11:30 11:32 11:34 11:36 11:38

/dev/xvdb3 (... 17.7 GB fr

Disk iops

Last Hour Last Day Last Week

Reads Writes

75 50 25 0

10:50 11:00 11:10 11:20 11:30 11:40

jacob

« Hide Stats

Disk space (36.4 GB free)

CPU usage (last hour)

/dev/xvda3 (2... 18.2 GB fr

Throughput (MB/s)

Last Hour Last Day Last Week

Reads Writes

2.0

/dev/xvdb3 (... 18.2 GB fr

Throughput (MB/s)

Last Hour Last Day Last Week

Reads Writes

2.0

Acunu Data Platform default-7048

View Licence | Log Out

Loading "http://hurley/ui/index.html#topology/data_centre/DC1", completed 436 of 438 items

ryan > Acunu Data Platform

ryan/ui/index.html#management/hosts

Management Versions Topology

View Licence | Log Out Acunu

HOSTS

Add

Name	IP Address	Cassandra Partitioner	Cassandra Token	CPU usage (last hour)	Cassandra Status	S3 Status
ryan	10.2.129.17	RandomPartitioner	60610040111976956729740652633953450669			

Change Partitioner

ryan > Acunu Data Platform

ryan/ui/index.html#management/column_family/200f9319-fc23-2a92-54b2-75f0dcaad797

Management Versions Topology

View Licence | Log Out Acunu

HOSTS
ryan

COLLECTIONS

KEYSPACES
hello
test

DISKS

BUCKETS

REPORTS

General

Name	test
Comment	
Type	Standard
Column Comparator	BytesType
Subcolumn Comparator	
Collection UUID	44d07349-f498-5f21-6faa-afe7a2f9c9d4

Ops

Last Hour Last Day Last Week

Gets Puts Big Gets Big Puts Range Queries Keys

Throughput (MB/s)

Last Hour Last Day Last Week

Reads Writes

Latency

Last Hour Last Day Last Week

95th Percentile Read Latency 95th Percentile Write Latency Read Latency Write Latency

A vertical decorative bar on the left side of the slide, divided into two horizontal sections: yellow at the top and blue at the bottom.

3. Operations

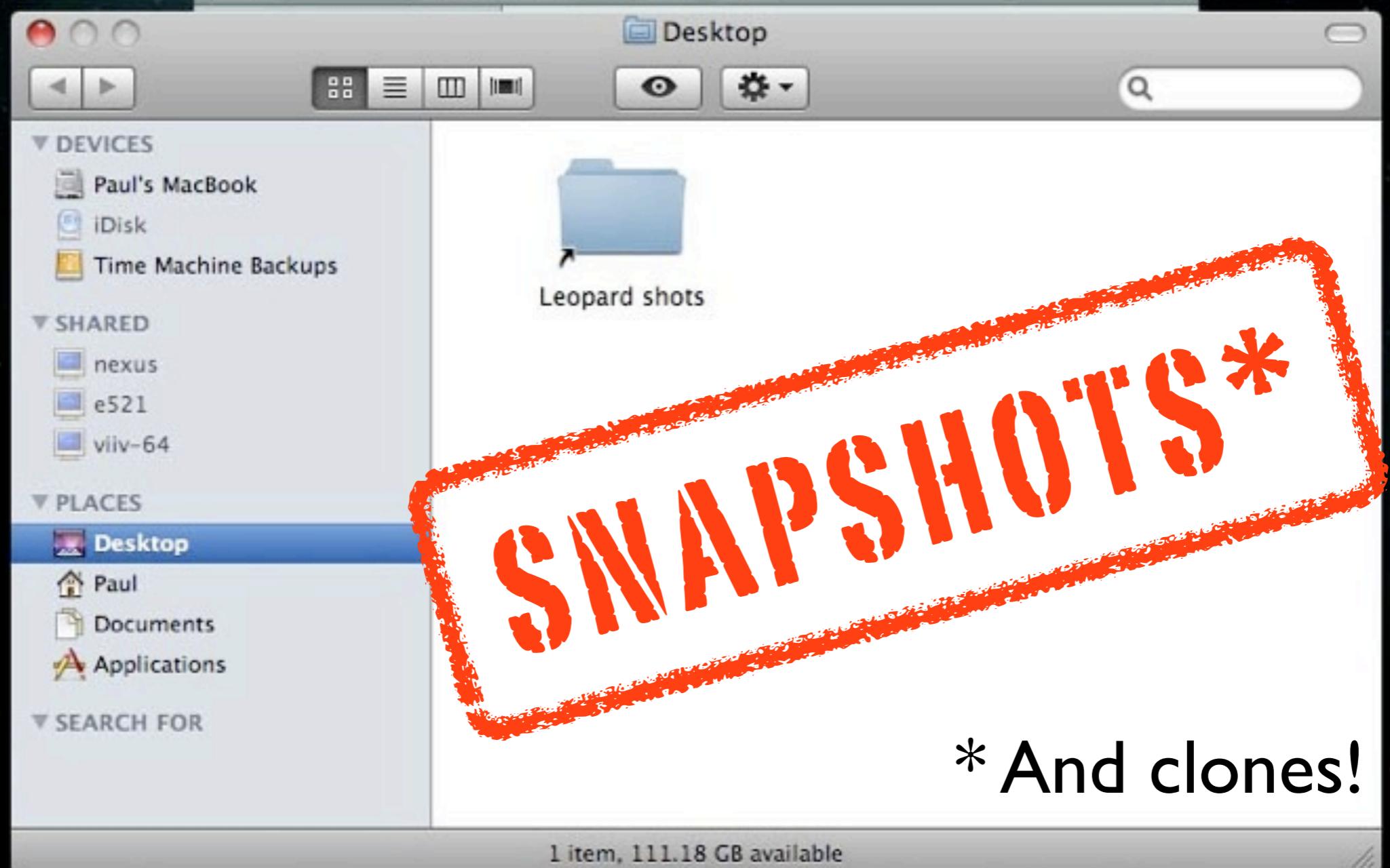


```
-bash-3.2$ nodetool
```

```
...
```

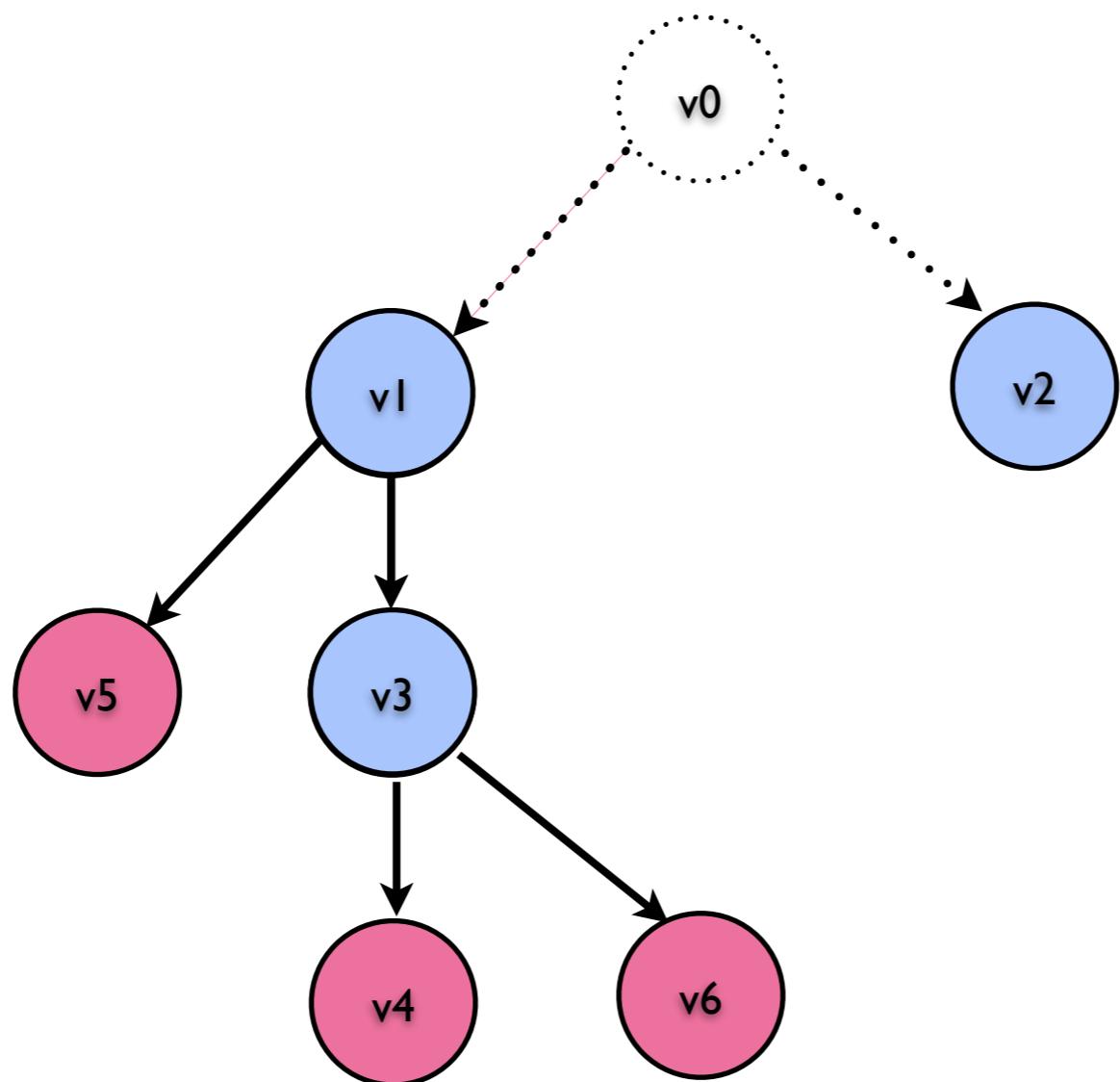
```
Available commands:
```

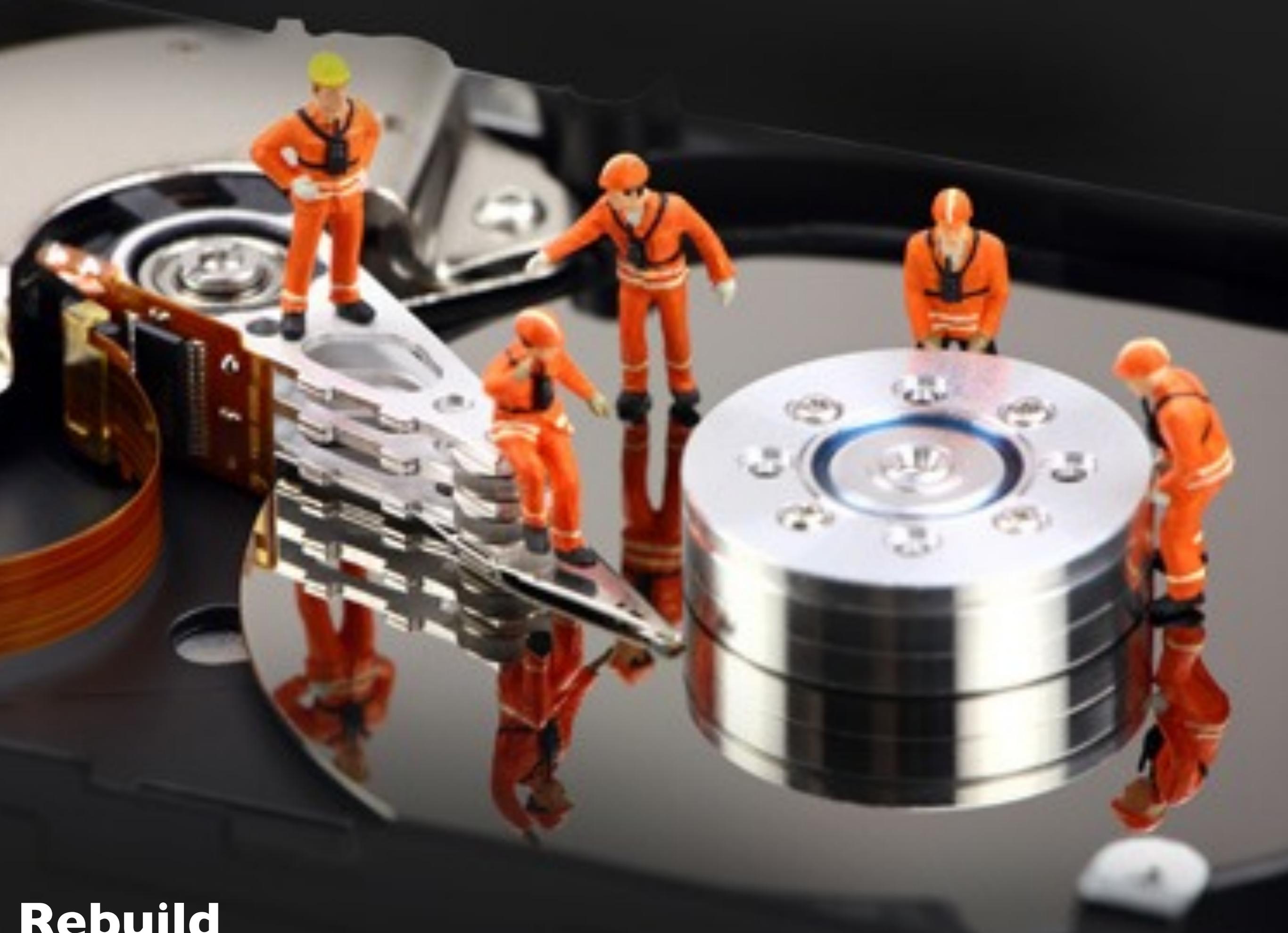
- ring
 - Print informations on the token ring
 - Join the ring
- join
 - Print node informations (uptime, load, ...)
 - Print statistics on column families
 - Print cassandra version
 - Print usage statistics of thread pools
 - Drain the node (stop accepting writes and flush all column families)
 - Decommission the node
- info
 - Print statistics on compactions
 - Disable gossip (effectively marking the node dead)
 - Reenable gossip
 - Disable thrift server
 - Reenable thrift server
- cfstats
 - Print network information on provided host (connecting node by default)
 - Move node on the token ring to a new token
- version
 - Show status of current token removal, force completion of
- tpstats
 - Set the MB/s throughput cap for compaction in the system
- drain
 - Take a snapshot of the specified keyspaces using
- decommission
 - Remove snapshots for the specified keyspaces
- compactionstats
 - Flush one or more column family
- netstats [host]
 - Repair one or more column family
- move <new token>
 - Run cleanup on one or more column family
- removetoken status|force|<token>
 - Force a (major) compaction on one or more column family
- setcompactionthroughput <value_in_mb>
 - Scrub (rebuild sstables for) one or more column family
- snapshot [keyspaces...] -t [snapshotName]
 - Invalidate the key cache of one or more column family
- clearsnapshot [keyspaces...] -t [snapshotName]
 - Invalidate the key cache of one or more column family
- flush [keyspace] [cfnames]
 - Get compactation thresholds for a given column family
- repair [keyspace] [cfnames]
 - Print statistic histograms for a given column family
- cleanup [keyspace] [cfnames]
 - Set the key and row cache capacities
- compact [keyspace] [cfnames]
 - Set the min and max compactation thresholds
- scrub [keyspace] [cfnames]
 - Invalidate the key cache of one or more column family
- invalidatekeycache [keyspace] [cfnames]
 - Invalidate the key cache of one or more column family
- invalidaterowcache [keyspace] [cfnames]
 - Invalidate the key cache of one or more column family
- getcompactionthreshold <keyspace> <cfname>
 - Print min and max compactation thresholds for a given column family
- cfdistograms <keyspace> <cfname>
 - Set the key and row cache capacities
- setcachecapacity <keyspace> <cfname> <keycachecapacity> <rowcachecapacity>
 - Set the min and max compactation thresholds
- setcompactionthreshold <keyspace> <cfname> <minthreshold> <maxthreshold>
 - Set the min and max compactation thresholds



SNAPSHOT'S*

* And clones!

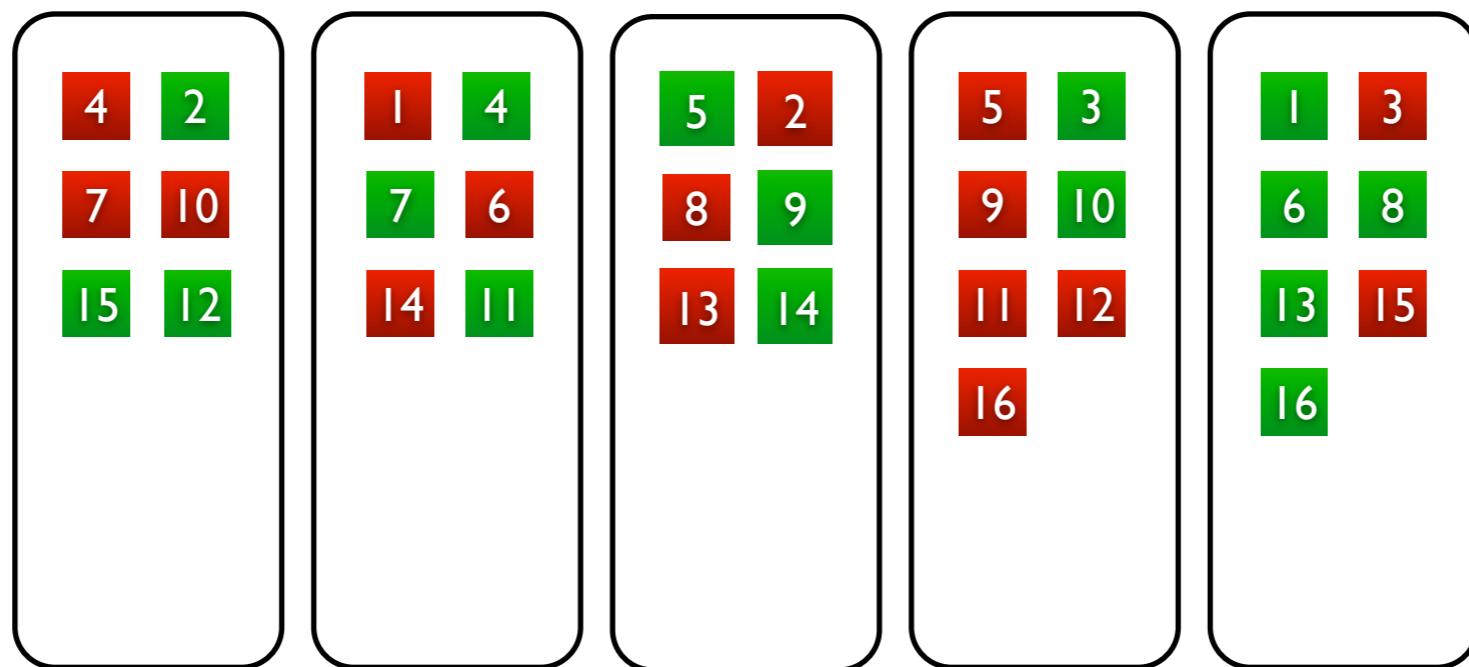




Rebuild

Disk Layout: RDA

random duplicate allocation



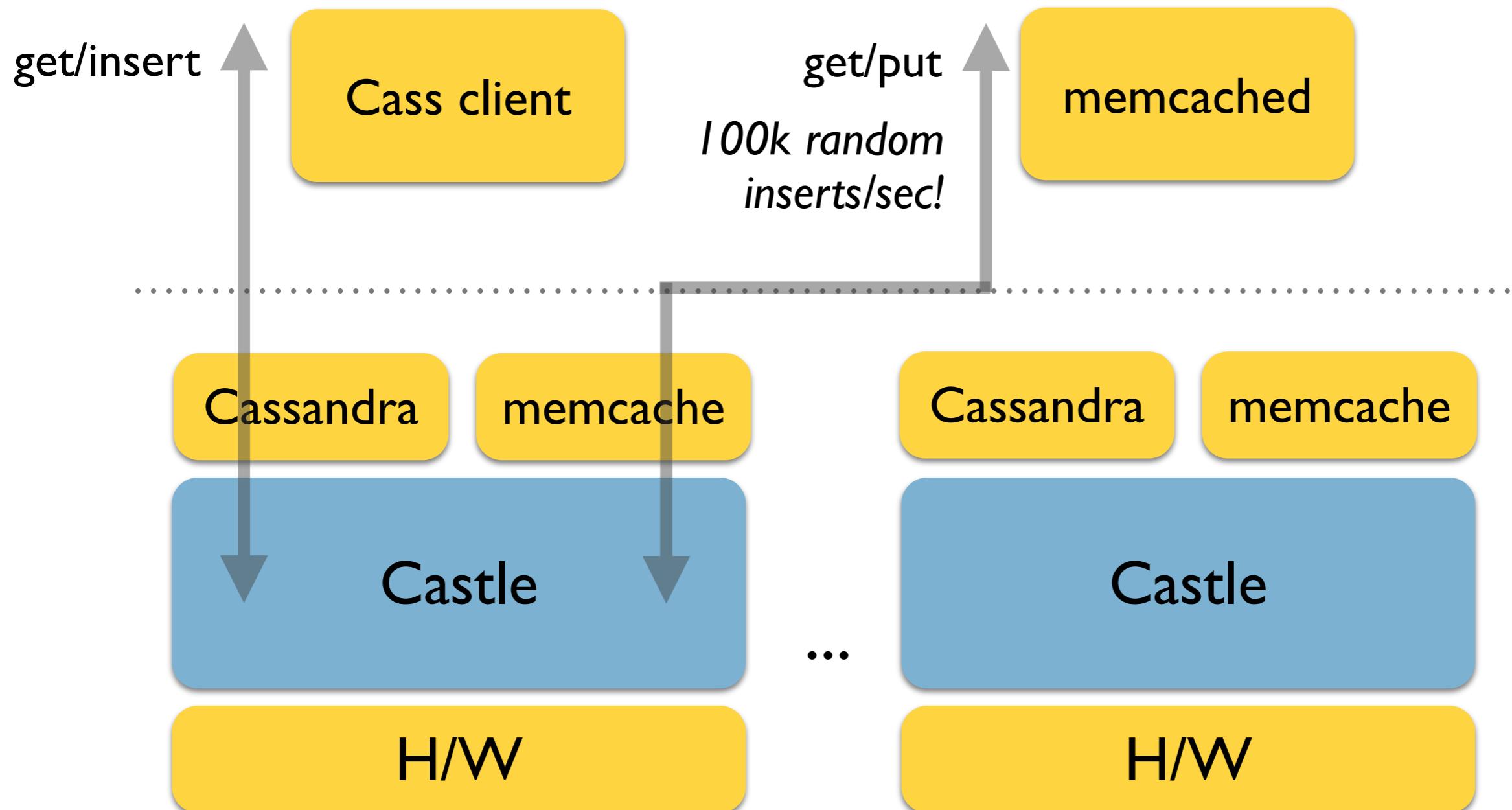
A large, abstract graphic on the left side of the slide features a yellow triangle at the top that tapers down to a wavy blue base, resembling a stylized wave or a rising sun.

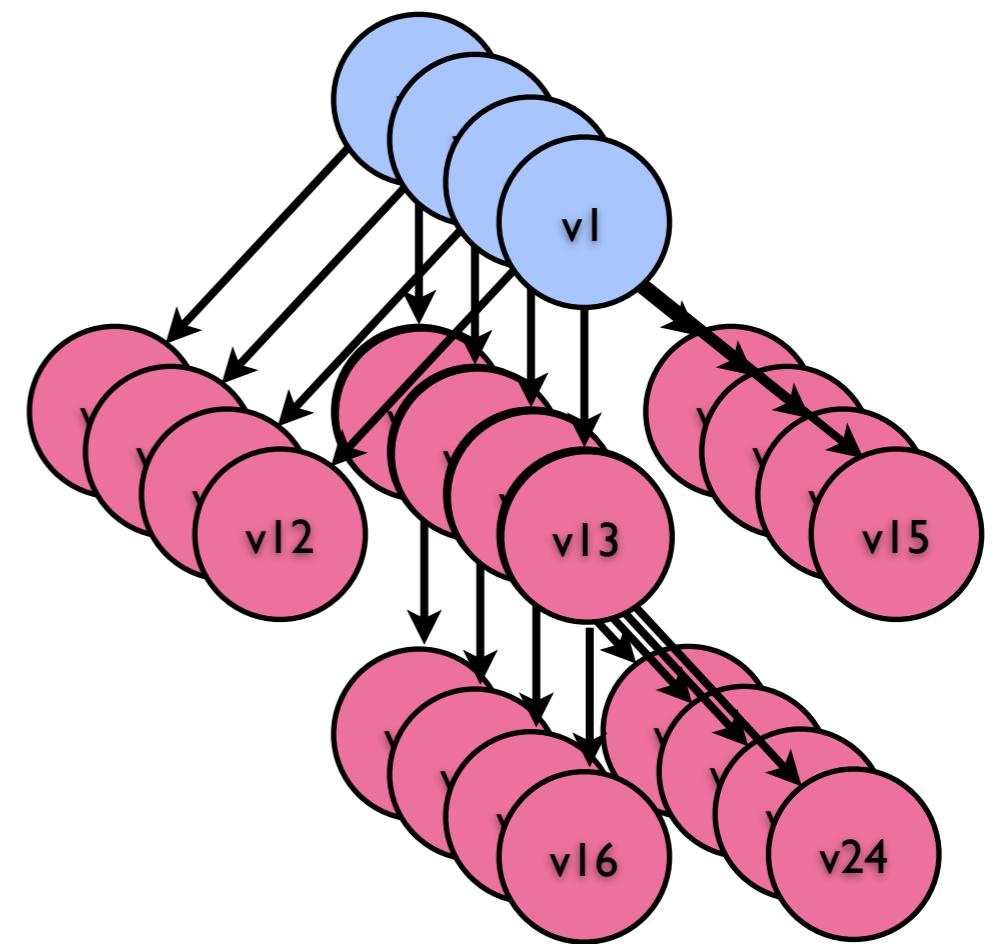
Future

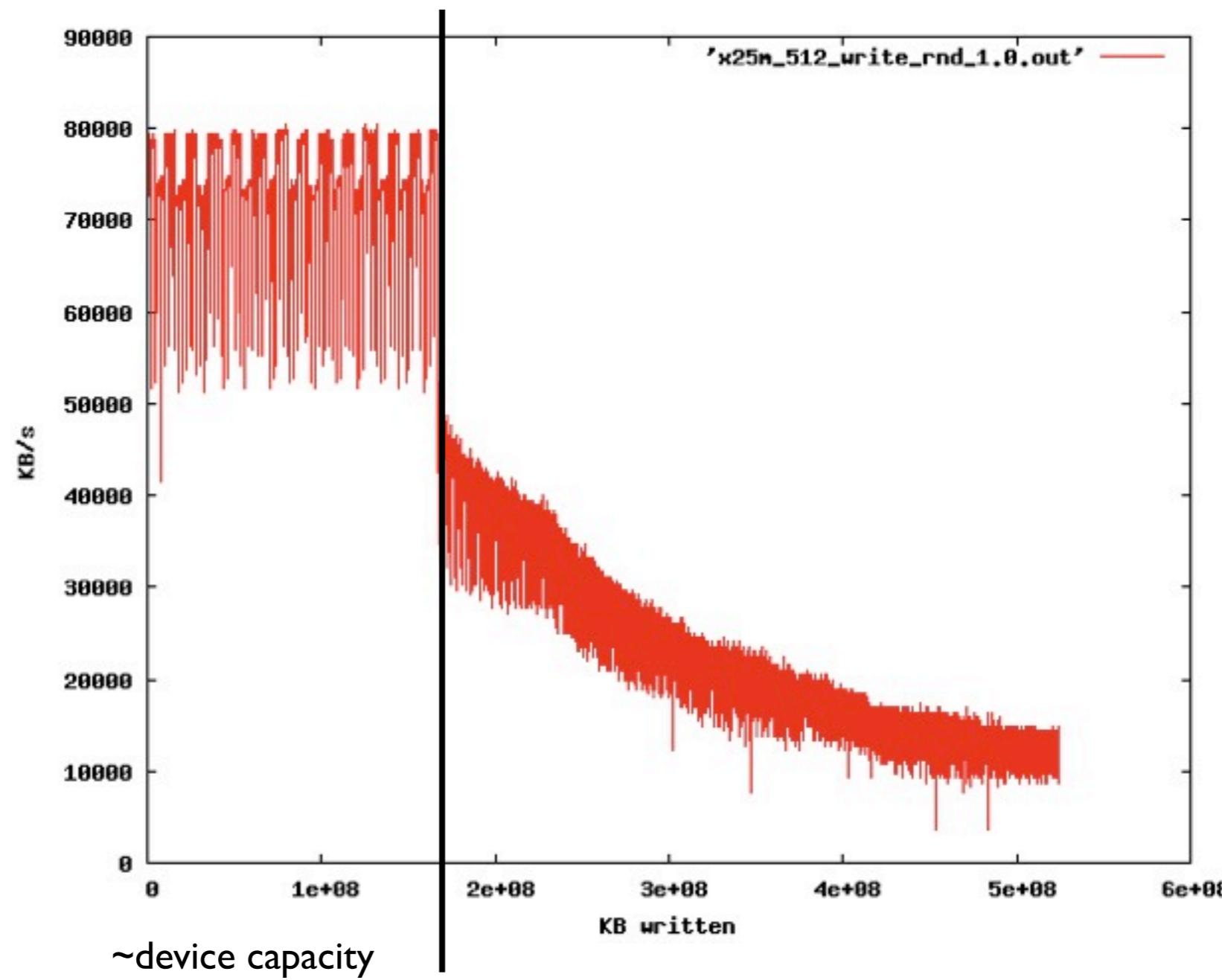


Acunu

Memcache + Cassandra







Beware the “write cliff”...

- Castle: Predictable Performance for Big Data
- Monitoring: distributed, multi-master tools, give you aggregated and summarised view of your cluster
- Snapshots & Clones: addressing real problems with new workloads
- RDA: lightening fast rebuilds for massive disks

Questions?

Tom Wilkie

@tom_wilkie

tom@acunu.com

<http://bitbucket.org/acunu>

<http://github.com/acunu>

<http://www.acunu.com/download>

<http://www.acunu.com/insights>



BIG DATA?
No Problem.



Acunu