

# Cassandra @ Signal

How to handle billion requests a day with Cassandra without losing sleep

Aki Colović  
Engineering Manager



Cut through the noise. [www.signal.co](http://www.signal.co)

# 7.2 devices/person

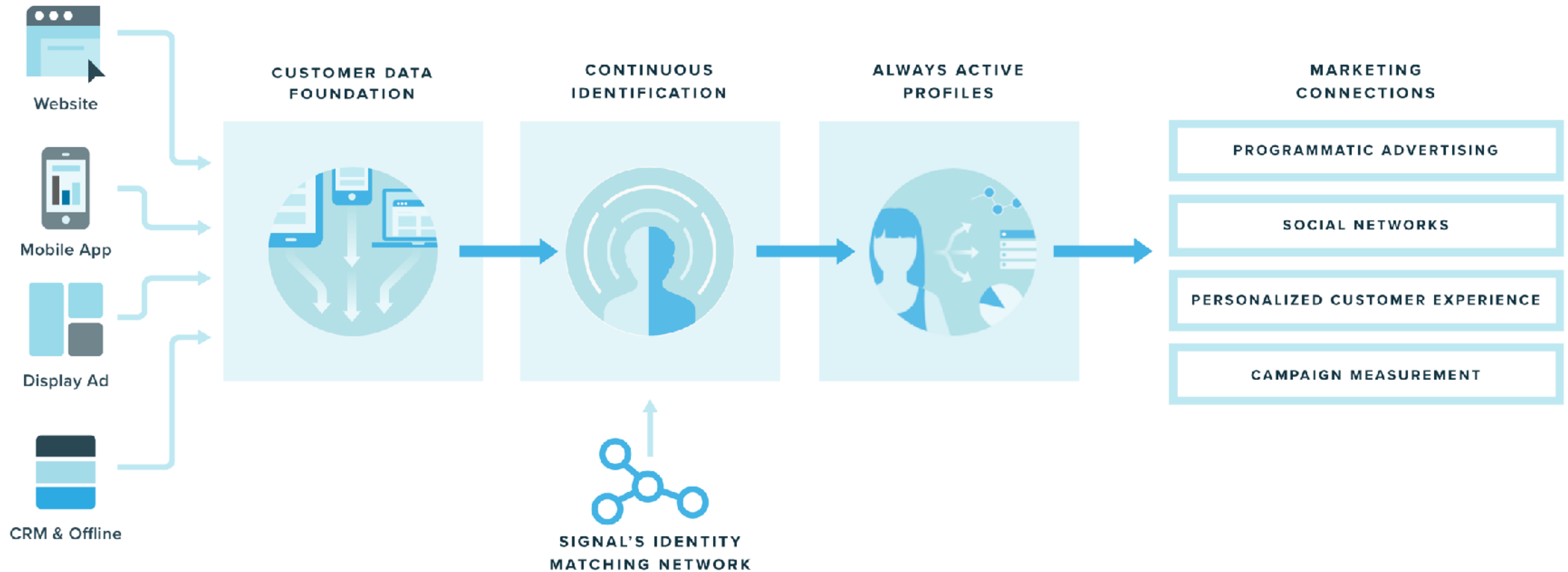




### **They do not have a single view of the customer:**

Just 6% of global marketers say they have attained the single view of the customer that is necessary to provide relevant, seamless experiences.

# Platform for Digital Marketing



# Our Partners



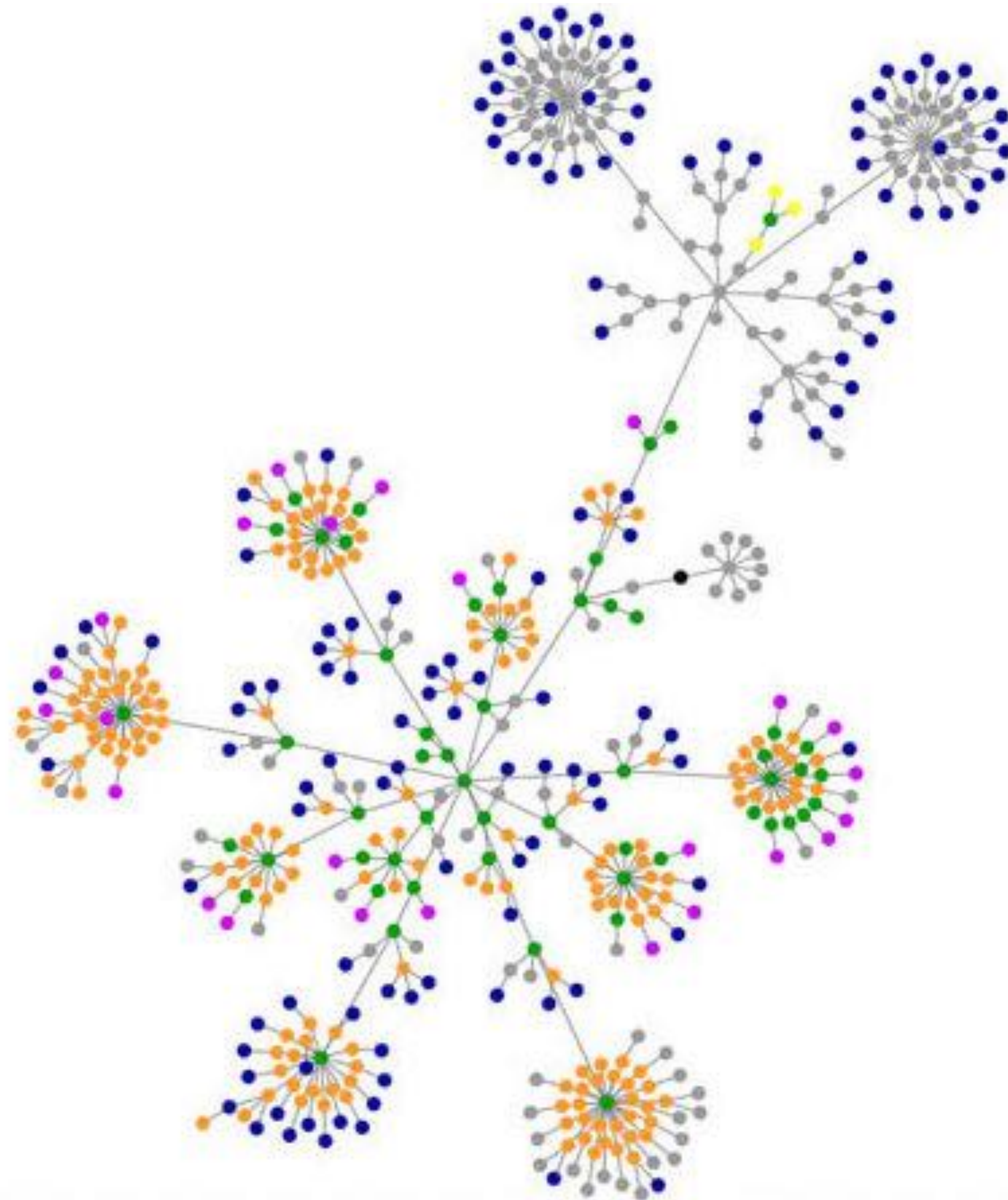
TURN

facebook®

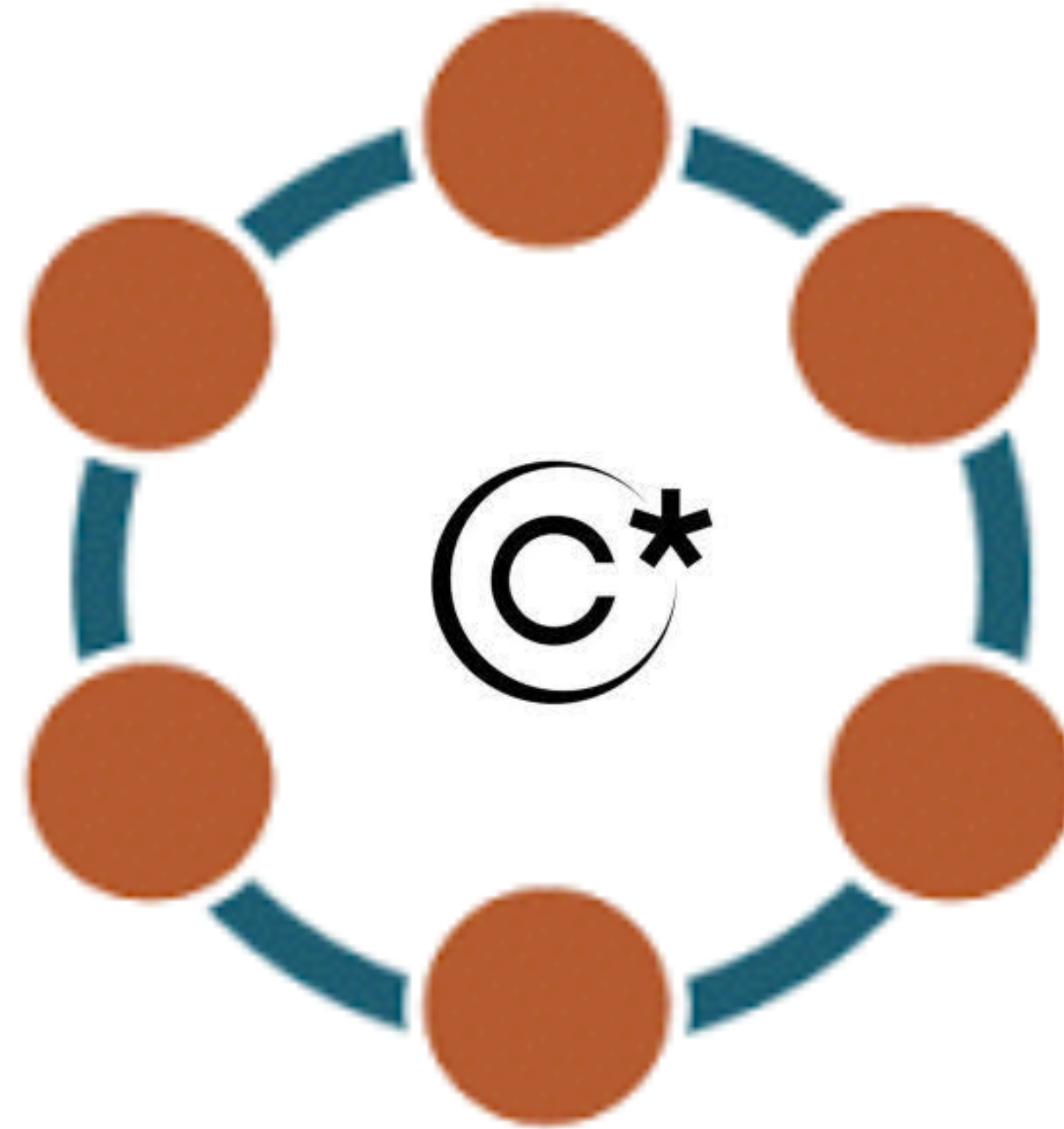




# Signal's Identity Graph



# Apache Cassandra





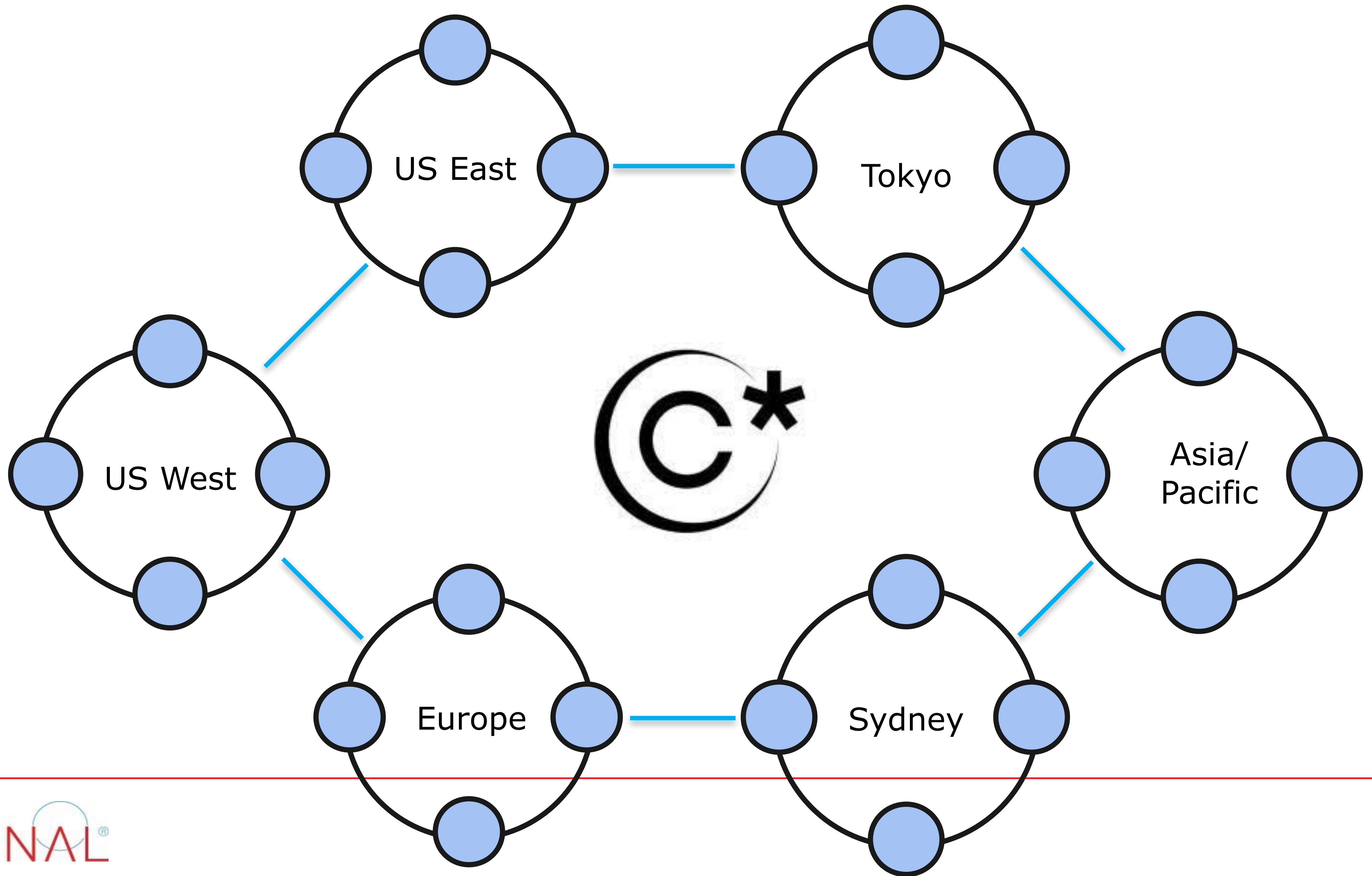
# Cassandra @ Signal

500+ Nodes

~1.5 billion queries per day

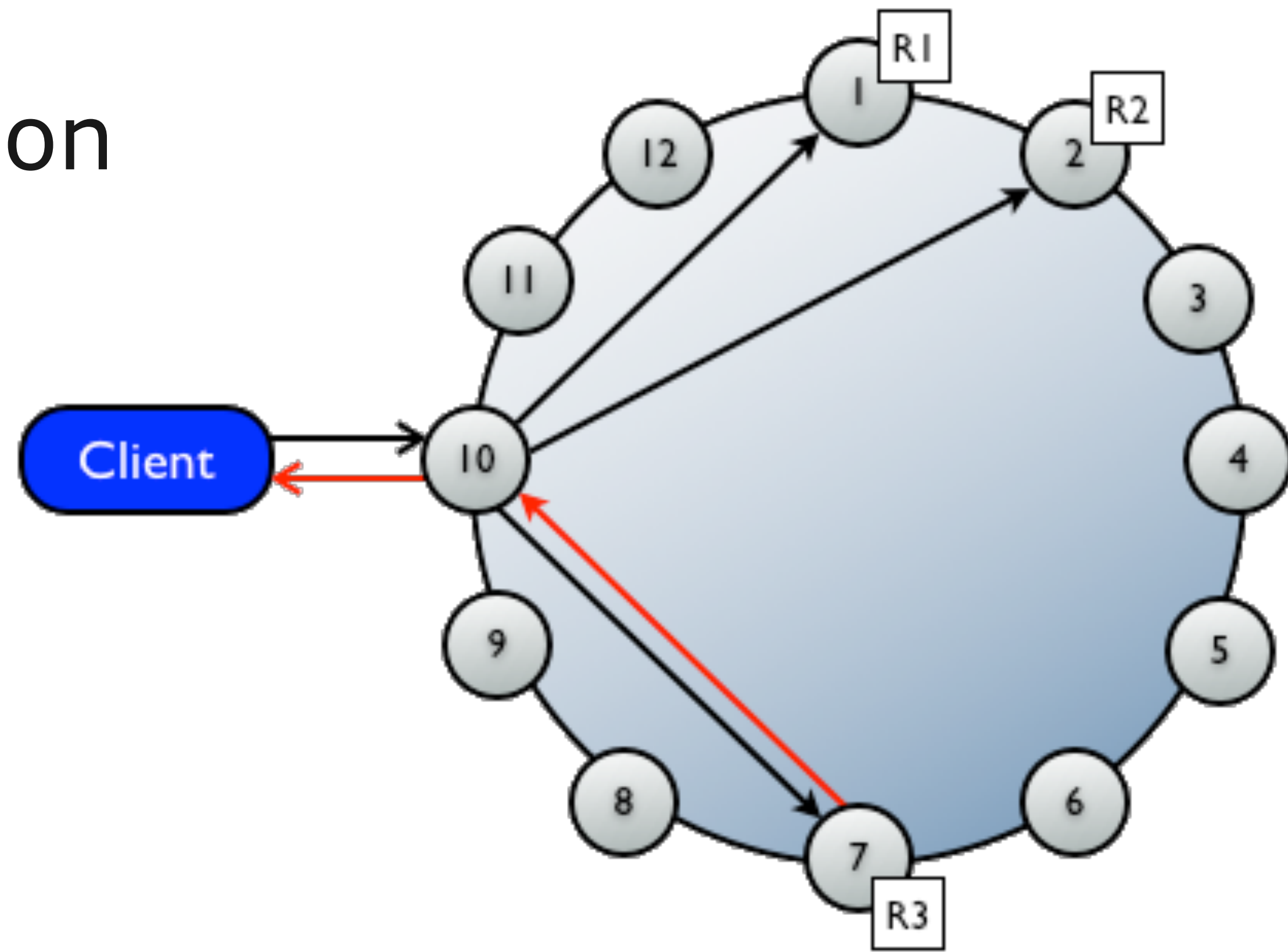
90.96+ TB of data





# What is Cassandra DB?

- distributed NoSQL database
- tunable consistency
- asynchronous replication
- decentralized
- high availability
- easy to scale



# Dealing With Failure

We've survived ...

Amazon swallowing nodes

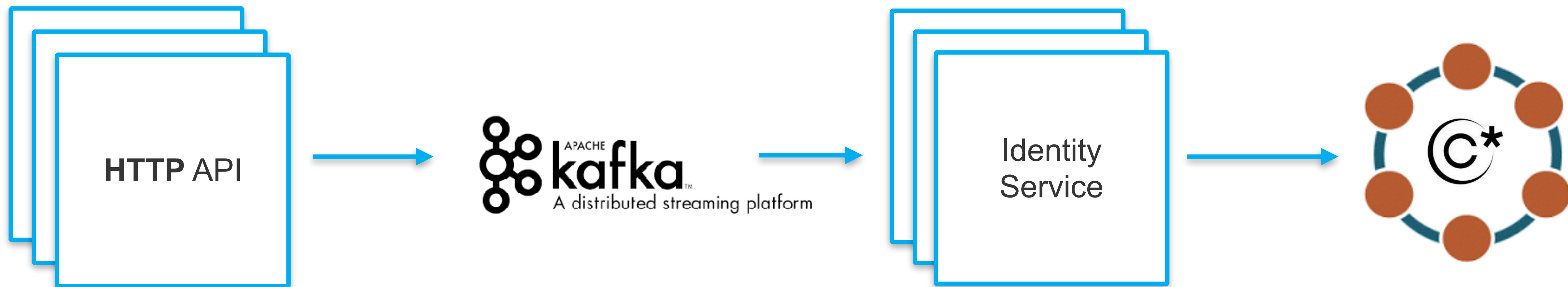
Amazon rebooting the internet

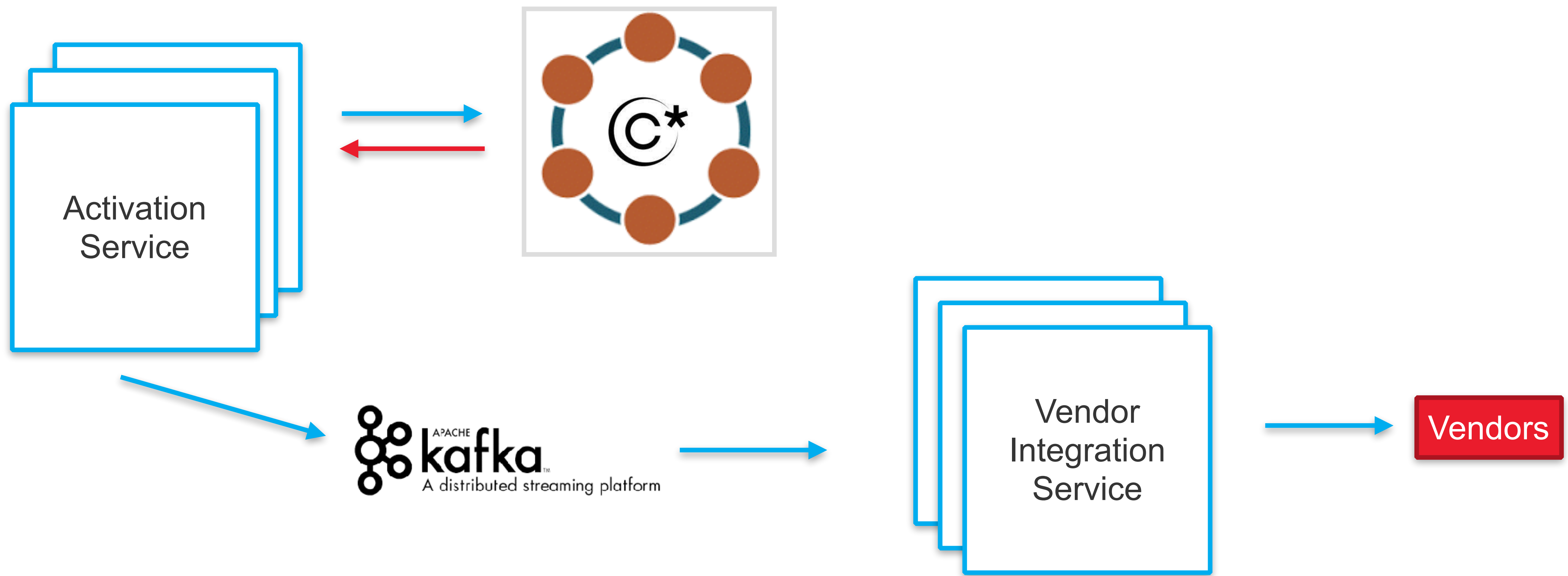
Disk Failures

Corrupt SSTables

Intra and inter region network issues

# Data Ingestion







# What About Downsides?



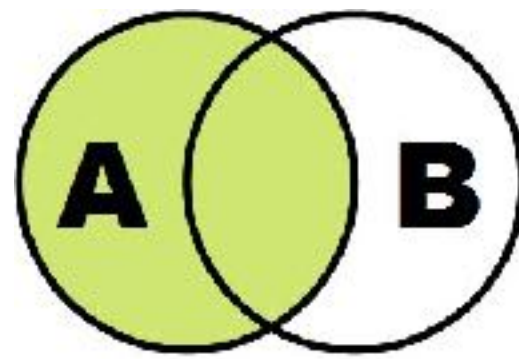
CQL != SQL

# Cassandra Query Language

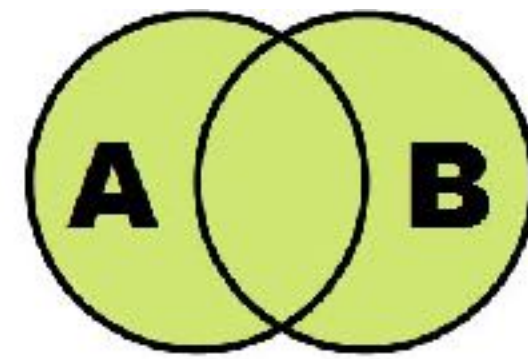
```
CREATE TABLE profiledata (  
  btid uuid,  
  identifier text,  
  data map<text, text>,  
  internal map<text, text>,  
  uids map<text, text>,  
  PRIMARY KEY (btid, identifier)  
) WITH  
  bloom_filter_fp_chance=0.010000 AND  
  caching='KEYS_ONLY' AND
```

```
SELECT FROM profiledata where btid=1a4b0130-7e0a-4b97-a368-91a5b439f6ae;
```

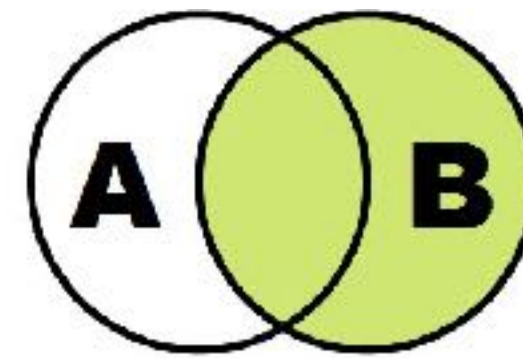
```
UPDATE profiledata_index SET btid=d4a44d4-53f8-495e-8db2-464889375ae3 WHERE entity_prefix ='s'
```



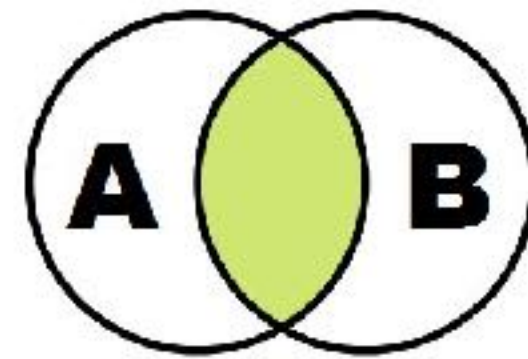
```
SELECT *
FROM A
LEFT JOIN B
ON A.id = B.id
```



```
SELECT *
FROM A
FULL OUTER JOIN B
ON A.id = B.id
```



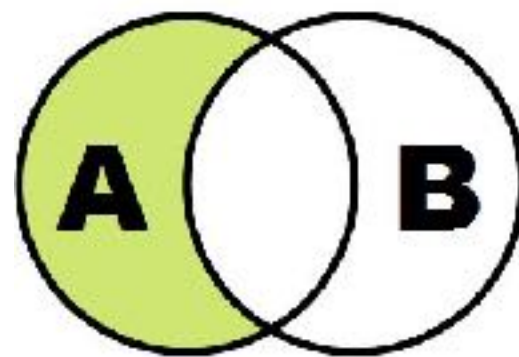
```
SELECT *
FROM A
RIGHT JOIN B
ON A.id = B.id
```



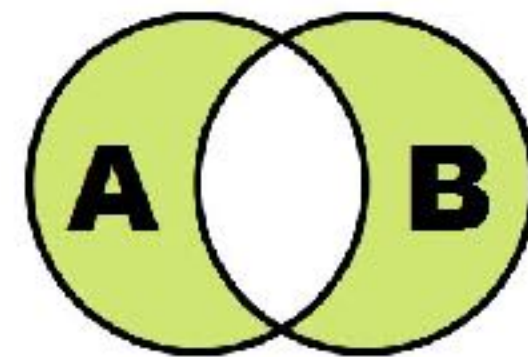
```
SELECT *
FROM A
INNER JOIN B
ON A.id = B.id
```

# No Joins, wait what?

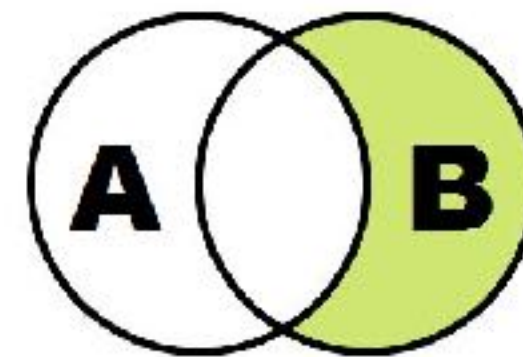
Copyright © 2012 www.matimattila.fi



```
SELECT *
FROM A
LEFT JOIN B
ON A.id = B.id
WHERE B.id IS NULL
```

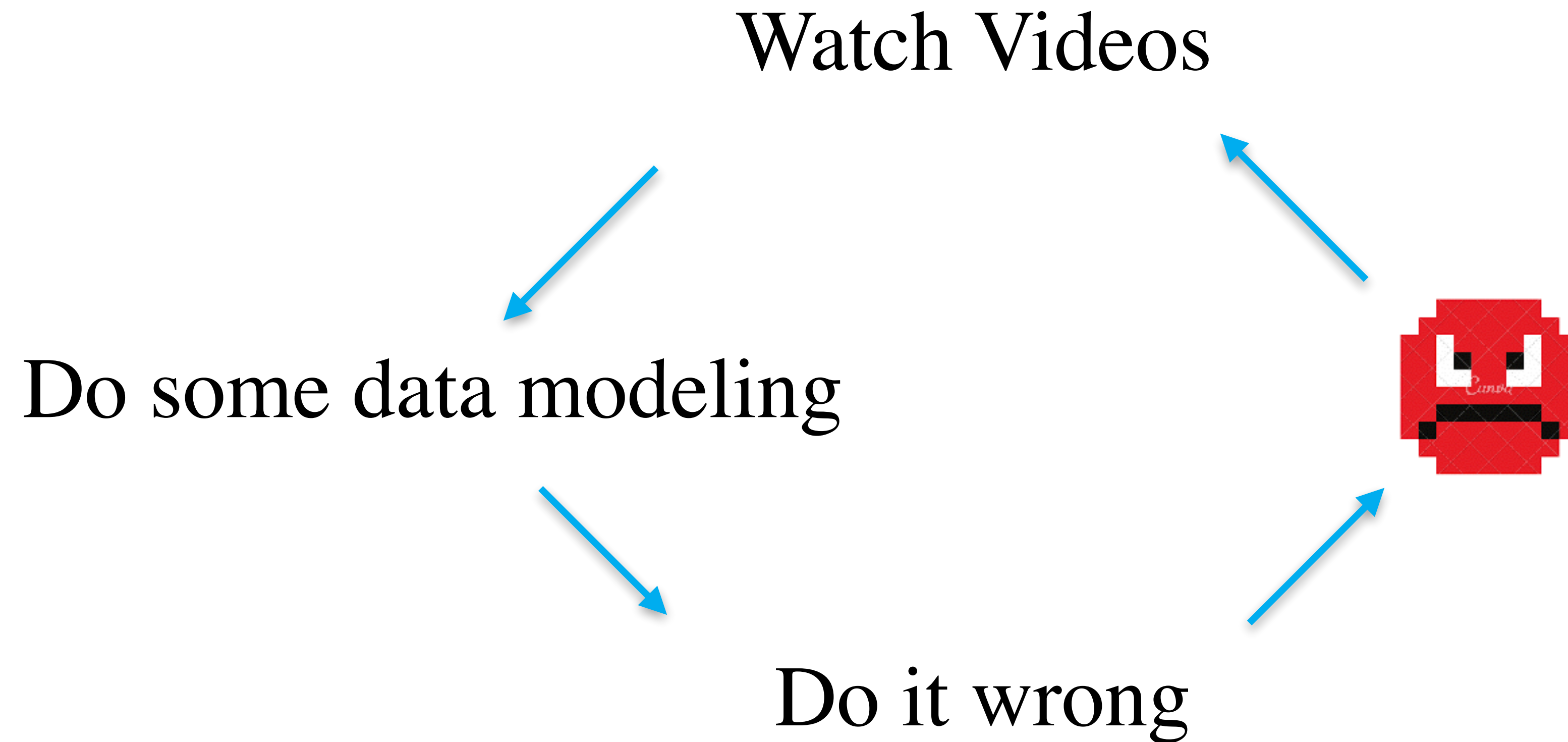


```
SELECT *
FROM A
FULL OUTER JOIN B
ON A.id = B.id
WHERE A.id IS NULL
OR B.id IS NULL
```



```
SELECT *
FROM A
RIGHT JOIN B
ON A.id = B.id
WHERE A.id IS NULL
```

# Data Modelling Pains





# Final thoughts on modeling...

"The best way to approach data modeling for Cassandra is to start with your queries and work backwards from there. Think about the actions your application needs to perform, how you want to access the data, and then design column families to support those access patterns."

# Secondary Index

- works well for low cardinality data
- does not scale
- material views( c\* 3.0+)

# Deleting is painful at scale

- tombstones
- disk not reclaimed immediately
- increased I/O pressure
- sometimes big expensive scans needed

# Testing at Scale

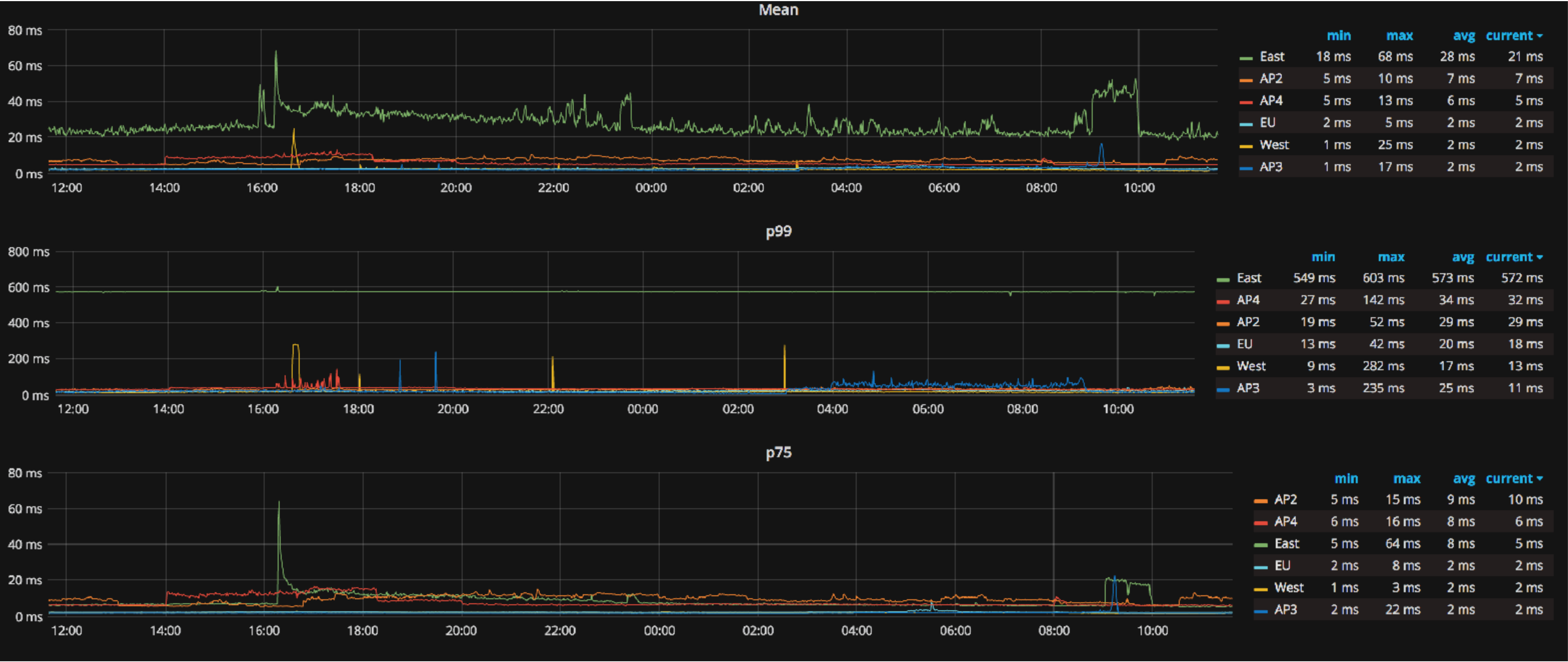
- test, test, test
- <https://github.com/jsevellec/cassandra-unit>
- <https://github.com/Netflix/chaosmonkey>

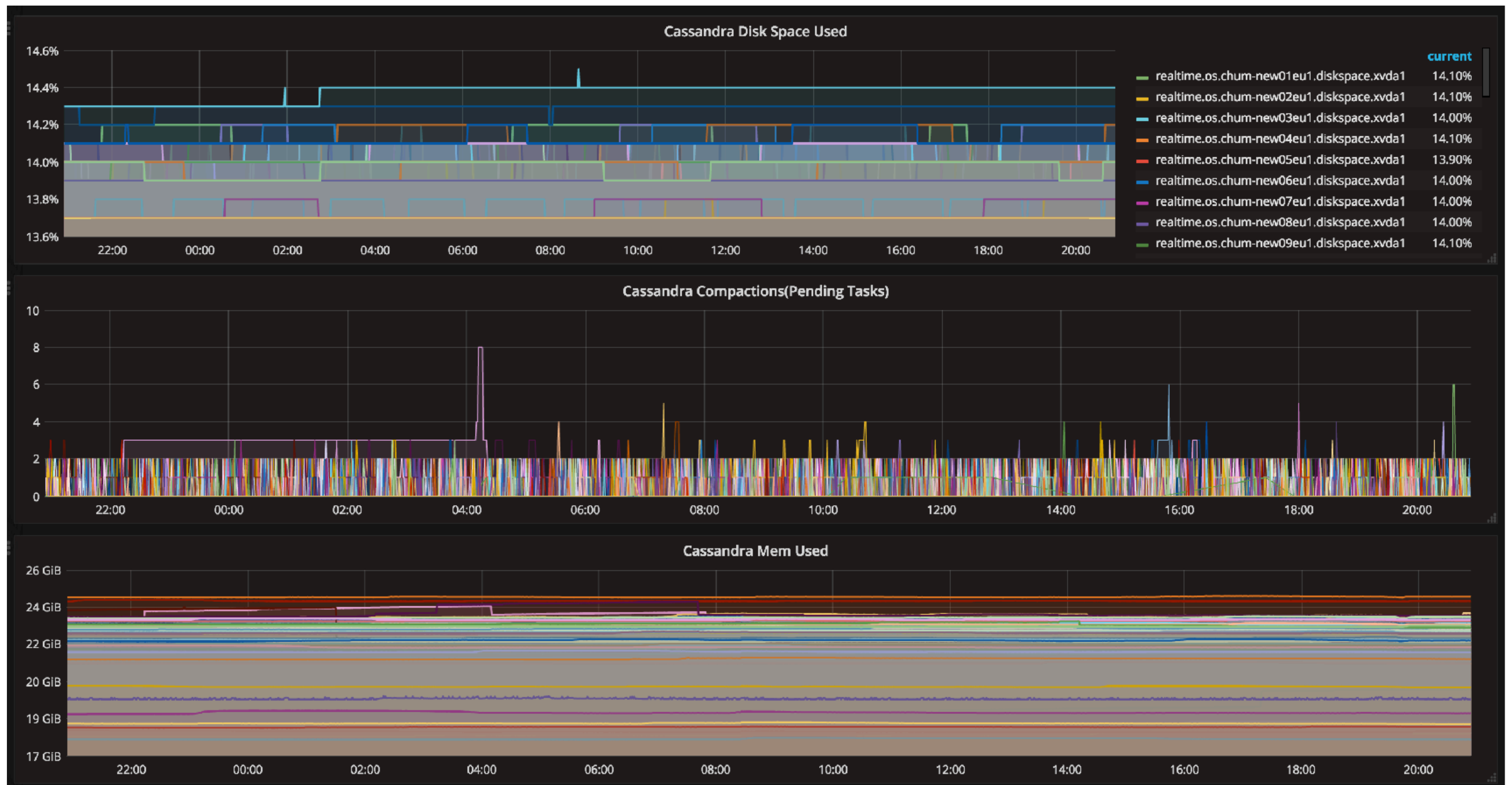


Monitoring will set you free!



# Monitor Read Latency

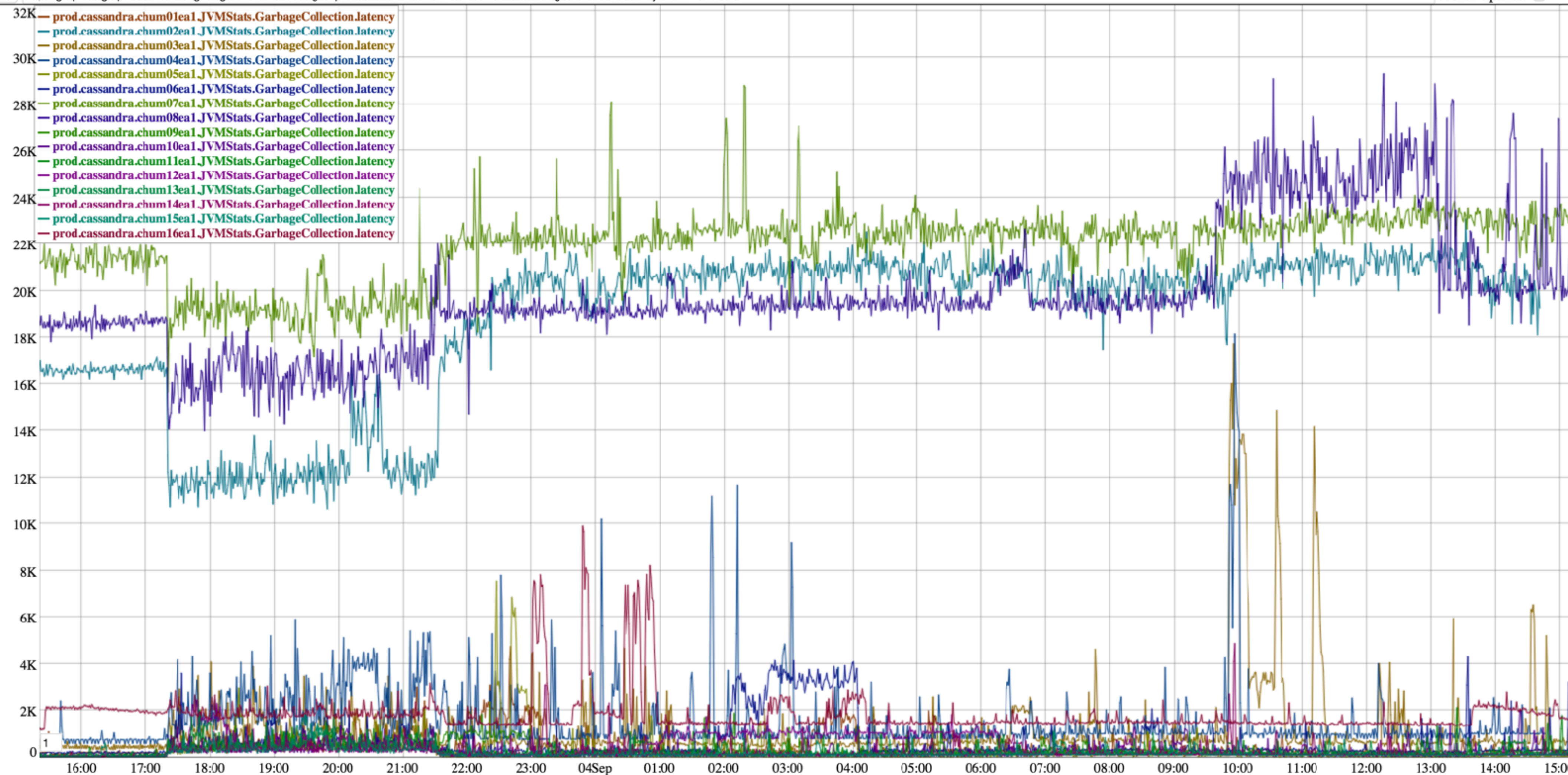






Go! [http://graphite.graph50ea1.thebrighttag.com/render/?target=prod.cassandra.chum\\*ea1.JVMStats.GarbageCollection.latency&from=-24hours](http://graphite.graph50ea1.thebrighttag.com/render/?target=prod.cassandra.chum*ea1.JVMStats.GarbageCollection.latency&from=-24hours)

Auto-update: ☐







**KEEP  
CALM  
'CAUSE YOU'RE  
NOT  
ALONE**





**signal.co/careers**

Hiring in **Chicago** and **NYC** office



**SIGNAL**®

**SIGNAL**®

Cut through the noise. [www.signal.co](http://www.signal.co)



# Thank You!



acolovic@signal.co



@akicolovic



Cut through the noise. [www.signal.co](http://www.signal.co)