# SCYLLA:

# NoSQL at Ludicrous Speed

主讲人：ScyllaDB软件工程师 贺俊

# Today we will cover:

+ **Intro: Who we are, what we do, who uses it**

+ Why we started ScyllaDB

+ Why should you care

+ How we made design decisions to achieve no-compromise performance and availability

# Introduction

+ **Founded by KVM hypervisor creators**

+ **Q2 2014 - Pivot to the database world**

+ **Q3 2015 - Decloak during Cassandra Summit 2015, Beta**

+ **Q1 2016 - General Availability**

+ **Q3 2016 - First Scylla Summit: 100+ Attendees**

+ **Q1 2017 - Completed B round**

+ **$25MM in funding**

+ **HQs: Palo Alto, CA; Herzelia, Israel**

+ **42+ employees, hiring!**

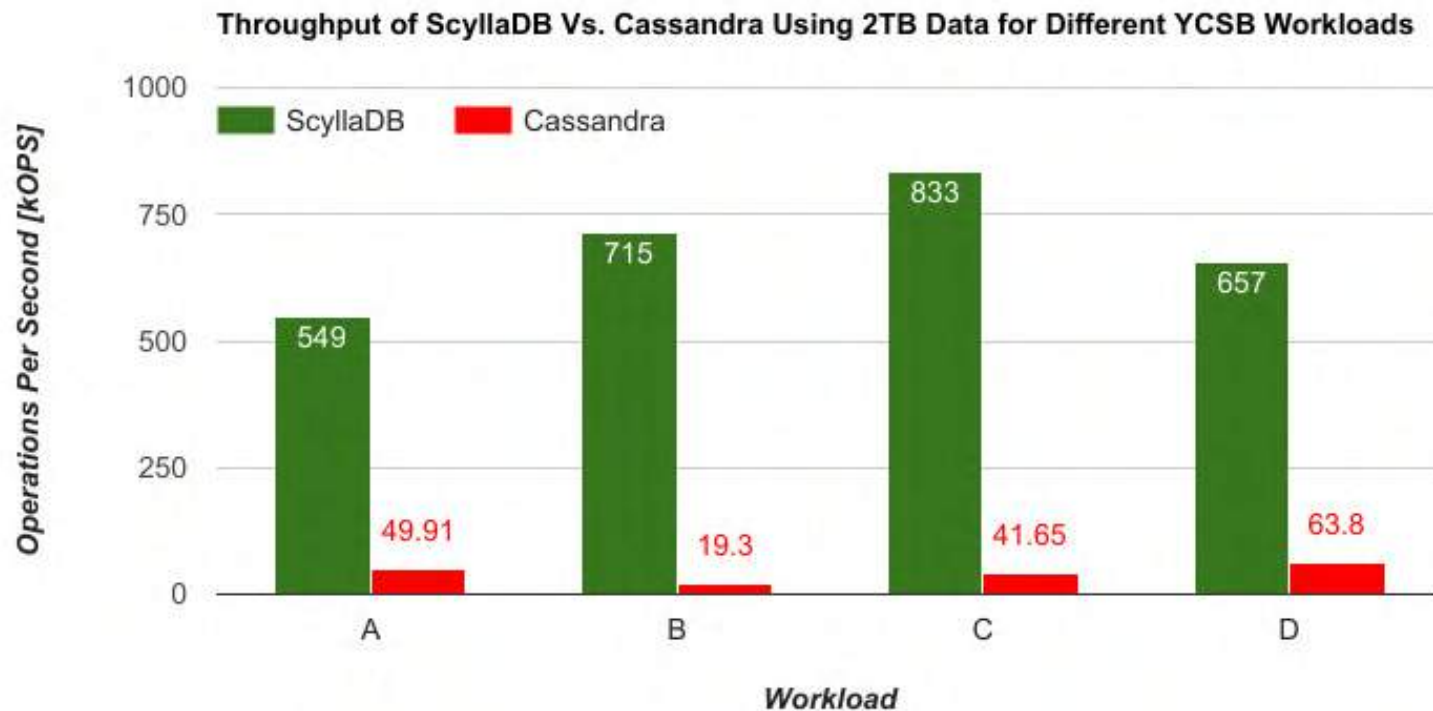Why ?@#$%$$%^?

Throughput of ScyllaDB Vs. Cassandra Using 2TB Data for Different YCSB Workloads
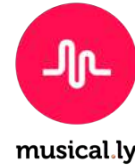
+   > 1 million OPS per node

+   < 1ms 99% latency

+   Auto tuned

+   Scale up and out

+   Open source

+   Large community (piggyback on Cassandra)

+   Blends in the ecosystem- Spark, Presto, time series, search, ..

## Cassandra shares #1 rank in HA

+ 1,000-node cluster
+ Flexible replication
+ Multi Datacenter
+ CQL language
+ Auto sharding
+ Wide rows
+ Lightweight Transactions

+ Homogeneous nodes
+ Spark integration, Presto
+ Vibrant Open Source community
+ More

# Where Scylla is deployed?

# Today we will cover:

+ **Intro: Who we are, what we do, who uses it**

+ Why we started ScyllaDB

+ Why should you care

+ How we made design decisions to achieve no-compromise performance and availability

# Why we started Scylla?

+ Originally it was about performance/efficiency only
+ Over time, we understood we can deliver more:
    + SLA between background and foreground tasks
    + Work well on any given hardware {back pressure}
    + Deliver consistent, low 99th percentile latency
    + Reduction in admin effort
    + Low latency under the face of failures (hot cache load balancing)
    + High observability

# Cassandra        Scylla

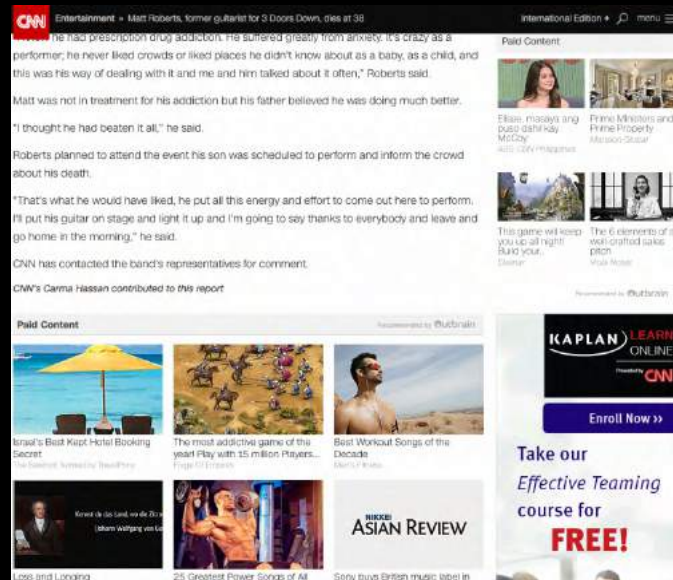| | Cassandra | Scylla |
|---|---|---|
| **Throughput:** | Cannot utilize multi-core efficiently | Scales linearly - shard-per-core |
| **Latency:** | High due to Java and JVM's GC | Low and consistent - own cache |
| **Complexity:** | Intricate tuning and configuration | Auto tuned, dynamic scheduling |
| **Admin:** | Maintenance impacts performance | SLA guarantee for admin vs serving |

# Today we will cover:

+ **Intro: Who we are, what we do, who uses it**

+ Why we started ScyllaDB

+ Why should you care

+ How we made design decisions to achieve no-compromise performance and availability

TalkingData

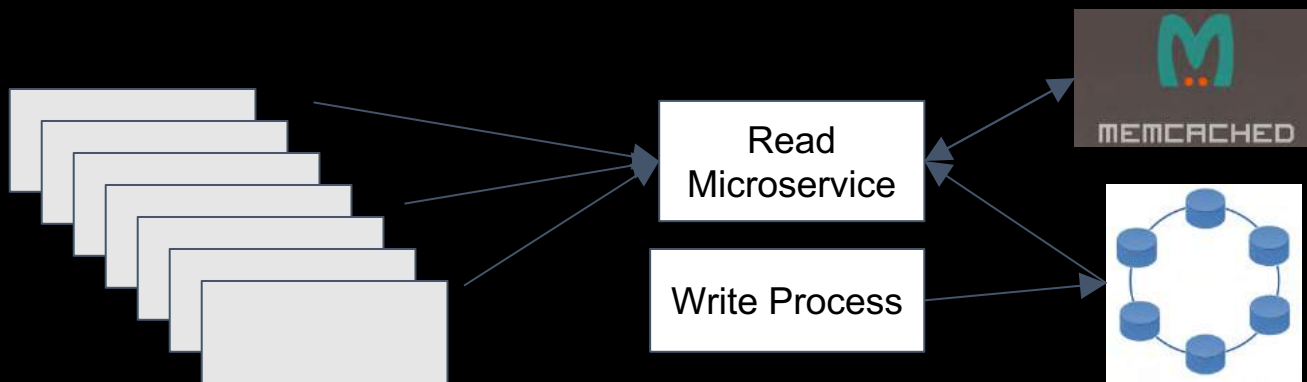# Case study: Document column family

- Outbrain is the world's largest content discovery platform.

- Over 557 million unique visitors from across the globe.

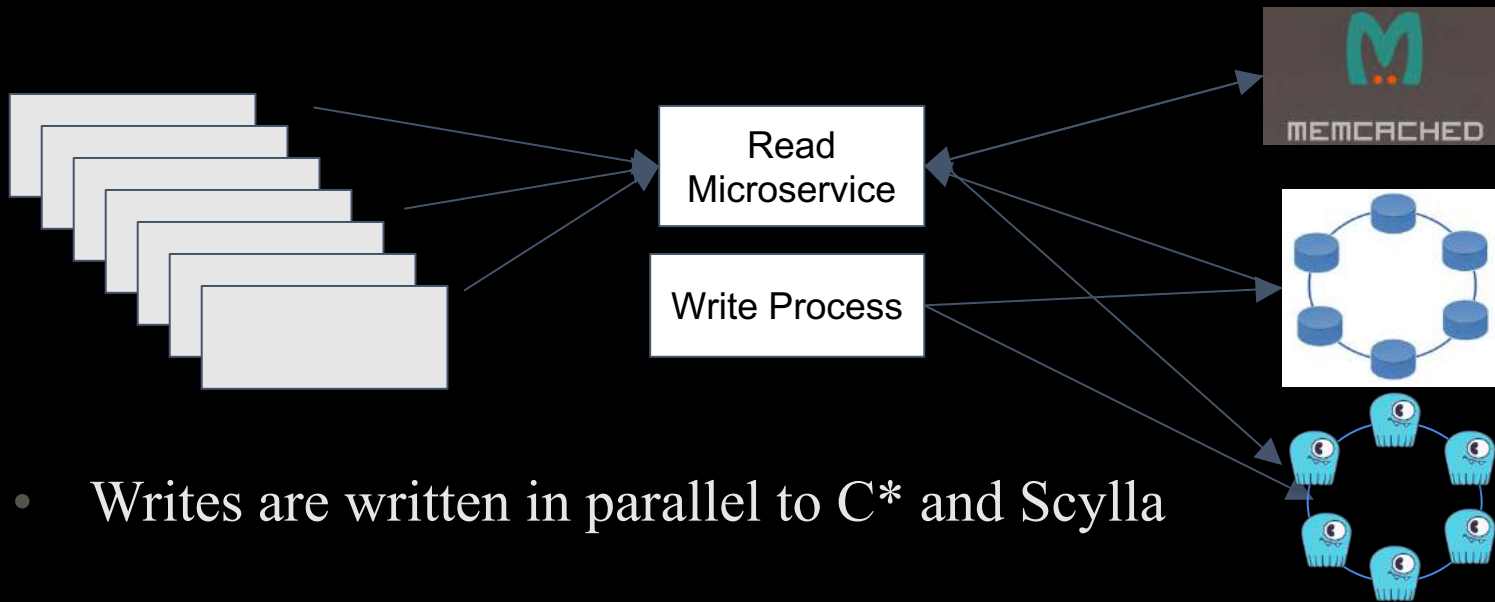- 250 billion personalized content recommendations every month.

# Outbrain: Cassandra plus Memcache



- First read from memcached, go to Cassandra on misses.

- Pain: 1) Stale data from cache 2) Complexity 3) Cold cache -> C* gets full volume
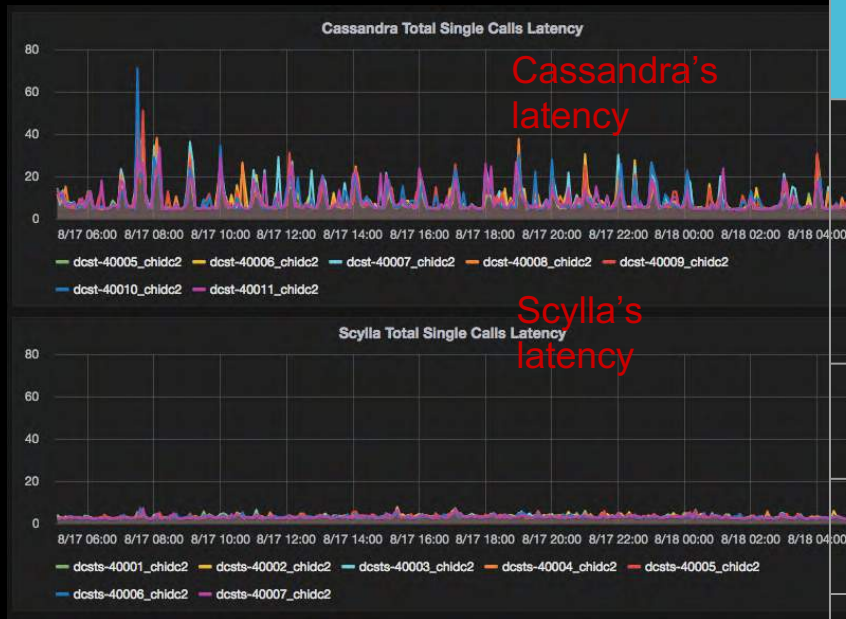
# Scylla/Cassandra side by side deployment



- Writes are written in parallel to C* and Scylla

- Reads are done in parallel:

1) Memcached + Cassandra 2) Scylla (no cache at all)

# Scylla (w/o cache) vs Cassandra + Memcached



Cassandra's latency

Scylla's latency

|  | Scylla | Cassandra | Diff% |
|---|---|---|---|
| **Requests/ Minute** | 12,000,000 | 500,000 (memcache handles 11,500,000) | 24X |
| **AVG Latency** | 4 ms | 8 ms | 2X |
| **Max Latency** | 8 ms | 35 ms | 3.5X |
| **Hardware** | 9 machines | 30+9 machines | 4.3X |

# What does it mean for a non Cassandra user?

+ Throughput, latency and scale benefits
+ Wide range of big data integration: {Kariosdb, Spark, JanusGraph, Presto, Kafka, Elastic}
+ Best HA/DR in the industry.
+ Stop using caches in front of the database
+ Consolidate HBase, Redis, MySQL, Mongo and others

# Assorted Quotes

## IBM

"ScyllaDB's NoSQL database offers a powerful combination of low latency and high availability, making it an attractive option for customers of our Watson Data Platform offering."

Derek Schoettle, General Manager, IBM Watson Data Platform

## Investing.com

"When we heard of a Cassandra drop-in-replacement we were skeptics. But very quickly we found it is all true—not only were the latency and GC issues completely gone, better hardware utilization allowed us to shrink the cluster size by half!"

Gabriel Mizrahi
CTO, Investing.com

**Read the Case Study**

## AppNexus

"We have a 47-node cluster across 5 data centers. With ScyllaDB we were able to reduce hardware cost and achieve great throughput and latency. Had we used Apache Cassandra for the same use case, we estimated that the cluster would have been at least twice as large."

Andrew Sweeney, VP of Engineering at AppNexus

## musical.ly

"Scylla reduced our latency to a level of single digit millisecond without changing a single line of code."
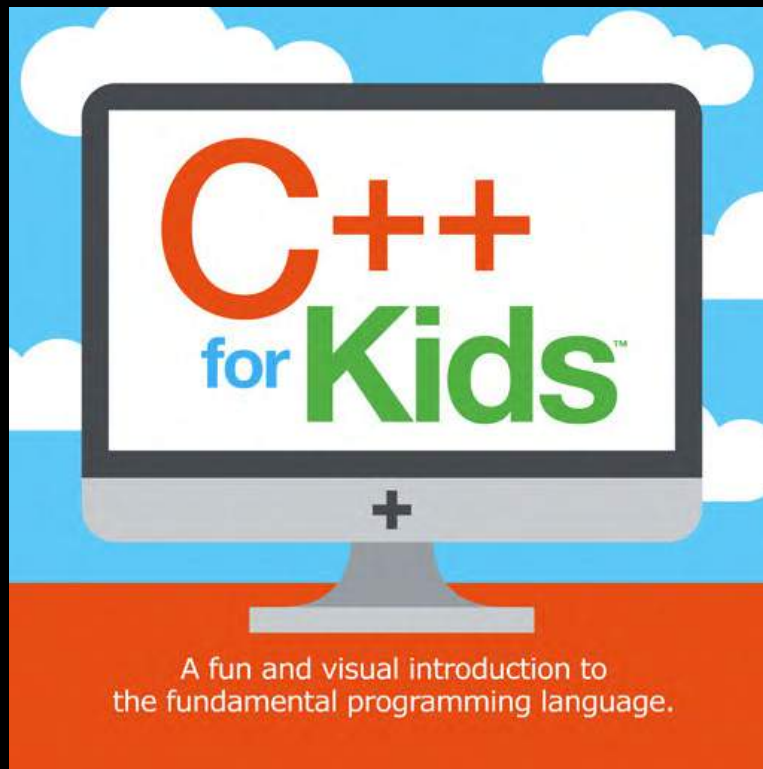
Terry Ma
Software Engineer, Musical.ly

# Today we will cover:

+ **Intro: Who we are, what we do, who uses it**

+ Why we started ScyllaDB

+ Why should you care

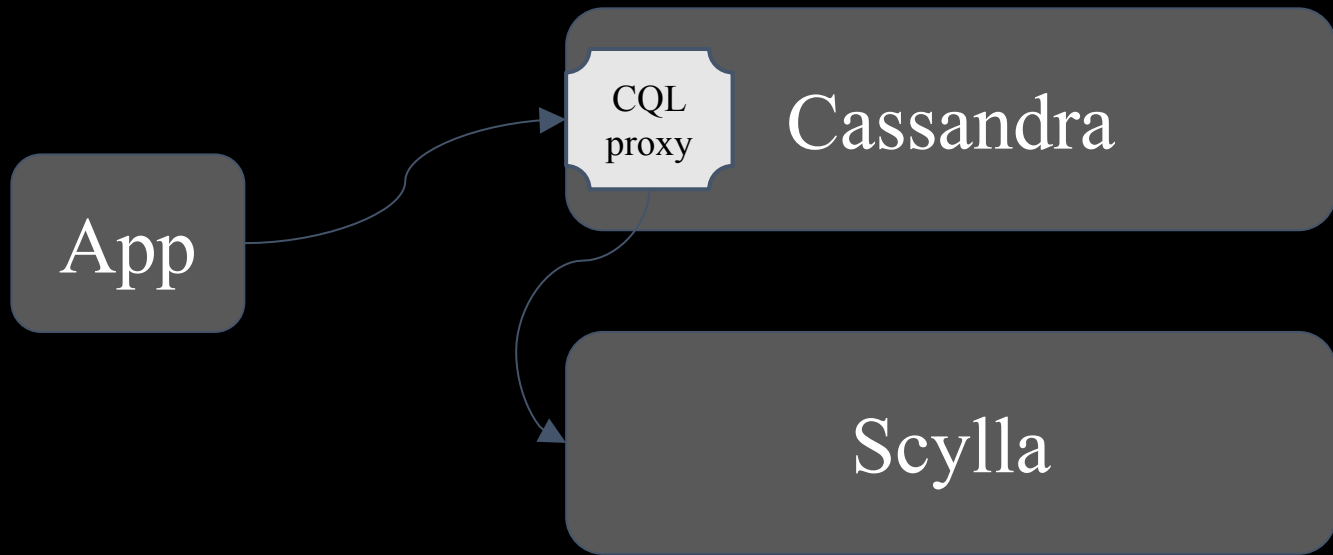+ How we made design decisions to achieve no-compromise performance and availability

# Design decisions: #2 Compatibility

- SSTable file format
- Configuration file format
- CQL language
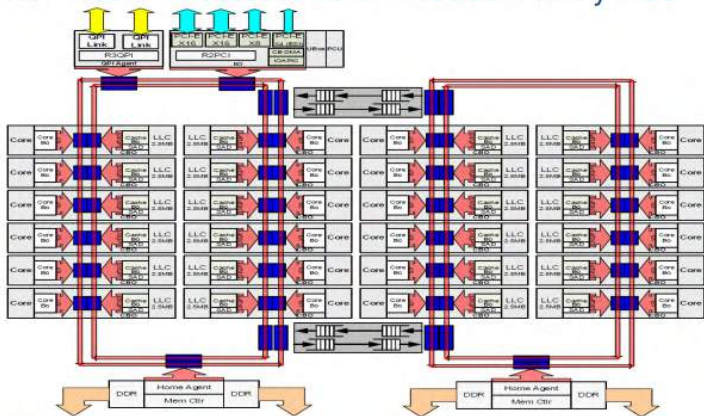- CQL native protocol
- JMX management protocol
- Management

# Double cluster - Migration w/o downtime



App

CQL proxy

Cassandra

Scylla

Intel® Xeon® Processor E5 v4 Product Family HCC

# Design decisions: #4 Shard per core

**Threads**

**Shards**

# Scylla has its own task scheduler

**Thread** is a function pointer

**Stack** is a byte array from 64k to megabytes

**Promise** is a pointer to eventually computed value

**Task** is a pointer to a lambda function

Thread

Stack

Scheduler

CPU

Promise

Task

Promise

Task

Promise

Task

Promise

Task

CPU

Context switch cost is high. Large stacks pollutes the caches

No sharing, millions of parallel events

'TalkingData

# SCYLLA IS DIFFERENT

- Thread per core
- Lock-free
- Task scheduler
- Reactor programing
- C++14

- Multi queue
- Poll mode
- Userspace TCP/IP

- NUMA friendly
- Log structured allocator
- Zero copy

- DMA
- Log structured merge tree
- DBaware cache
- Userspace I/O scheduler
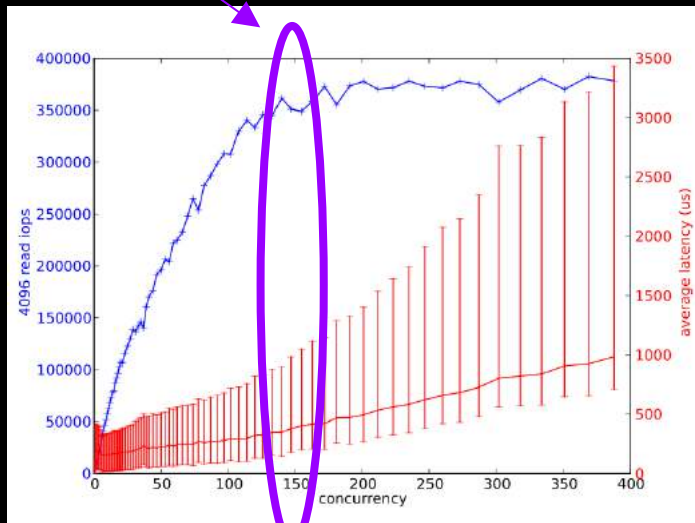
# Scylla vs C* latency by Kenshoo

Cassandra Streaming configuration

```
765    # Throttles all outbound streaming file transfers on this node to the
766    # given total throughput in Mbps. This is necessary because Cassandra does
767    # mostly sequential IO when streaming data during bootstrap or repair, which
768    # can lead to saturating the network connection and degrading rpc performance.
769    # When unset, the default is 200 Mbps or 25 MB/s.
770    # stream_throughput_outbound_megabits_per_sec: 200
```
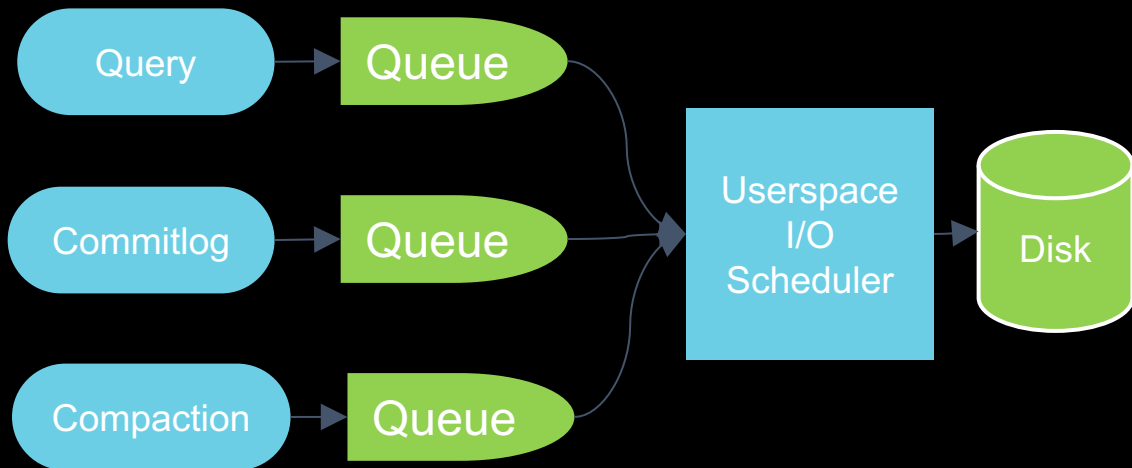
# Scylla I/O Scheduling
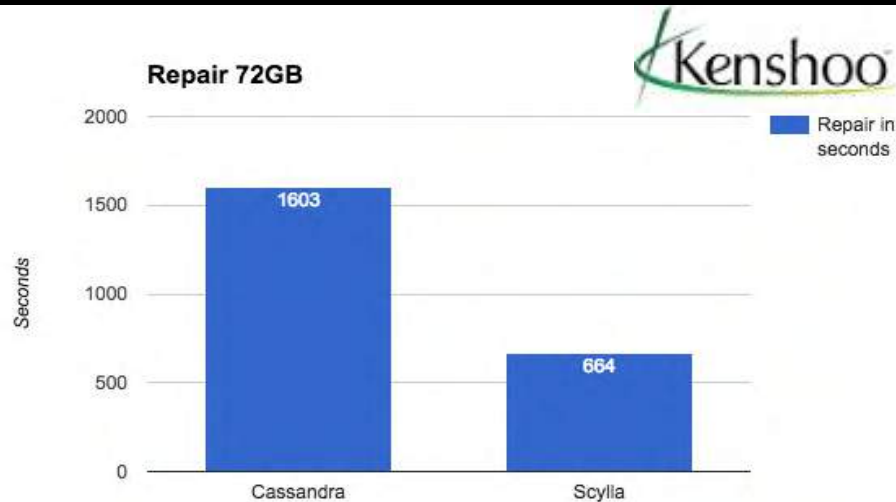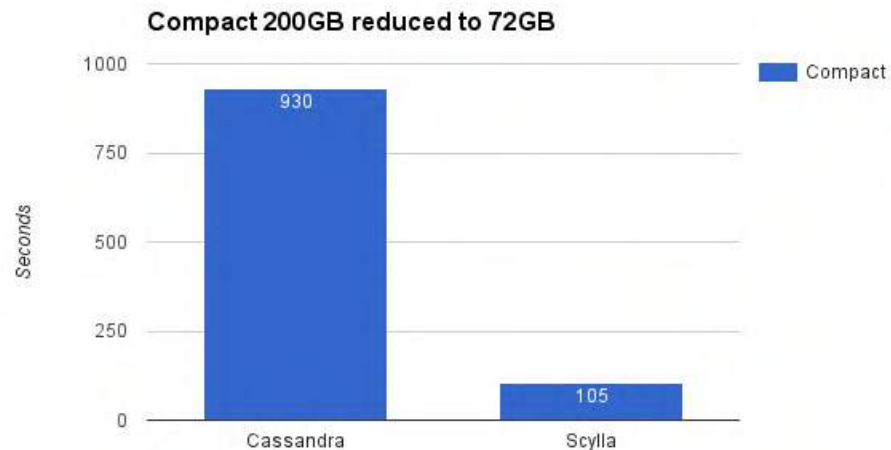
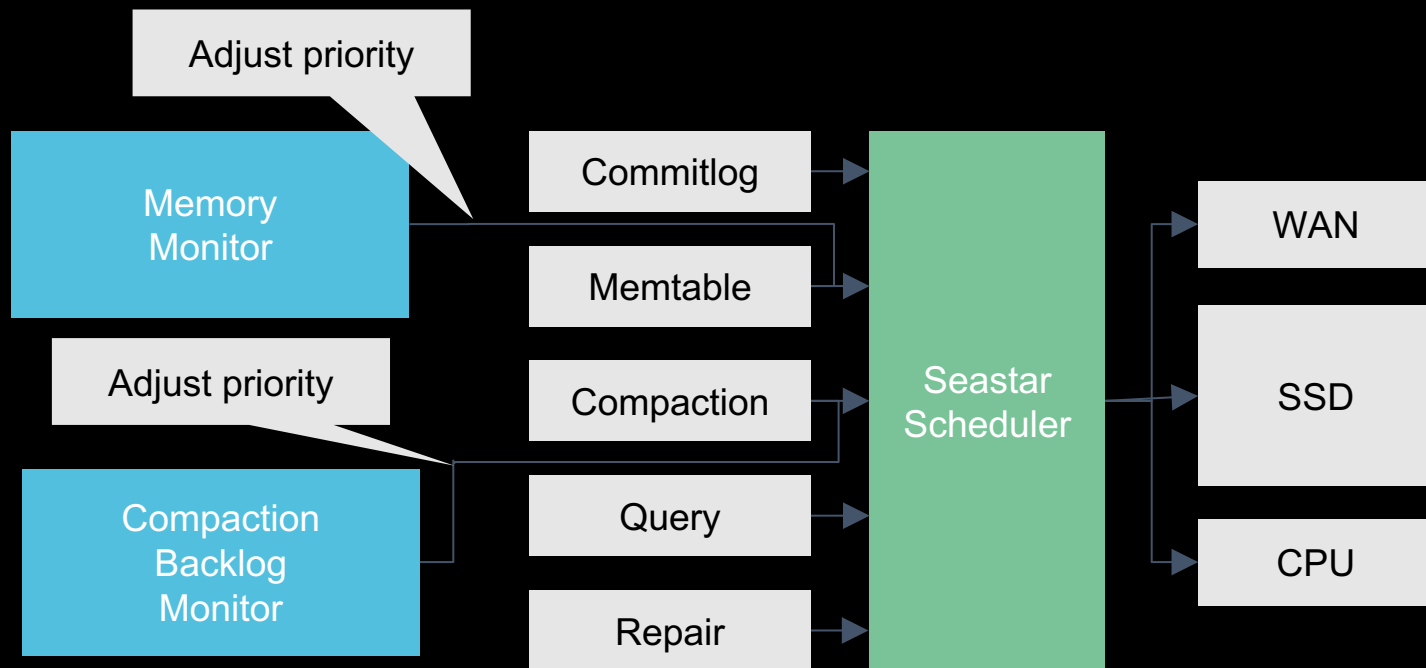Max useful disk concurrency



No queues     I/O queued in FS/device

Query → Queue

Commitlog → Queue

Compaction → Queue

Userspace I/O Scheduler → Disk

# I/O scheduler result by Kenshoo

# Design Decision: #7 Workload conditioning

# Workload Conditioning in practice



**Total Requests**

Disk can't keep up:
workload conditioning will figure out the right request rate

TalkingData

# Upcoming releases

+ Enterprise release, based on 1.6

+ 1.7 - May 2017

  - Counters

  - New intra-node sharding algorithm

  - SStableloader from 2.2/3.x

  - Debian

+ 2.0 – Sep 2017

  - Materialized views

  - Execution blocks (cpu cache optimization which boost performance)

  - Partial row cache (for wide row streaming)

  - Heat Weighted Load Balancing

# Scylla Beyond Cassandra



Core database



Vertical



Horizontal

TalkingData

# Q&A

## Resources

dor@scylladb.com (@DorLaor)

avi@scylladb.com (@AviKivity)

github.com/scylladb/scylla

scylladb.com/blog

@scylladb

http://bit.ly/2oHAfok

youtube.com/c/scylladb

slideshare.net/ScyllaDB

TalkingData

# THANKS

SCYLLA: NoSQL at Ludicrous Speed