# Two years in the making: What is new in Apache Cassandra 4.0?
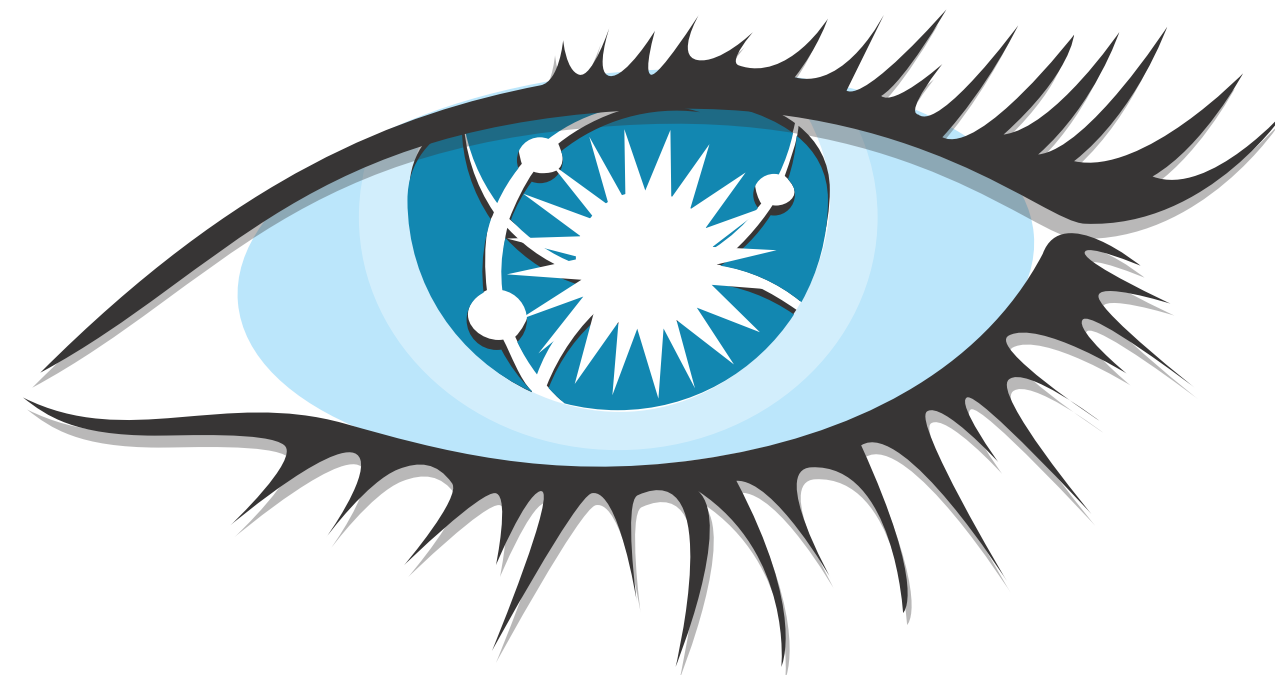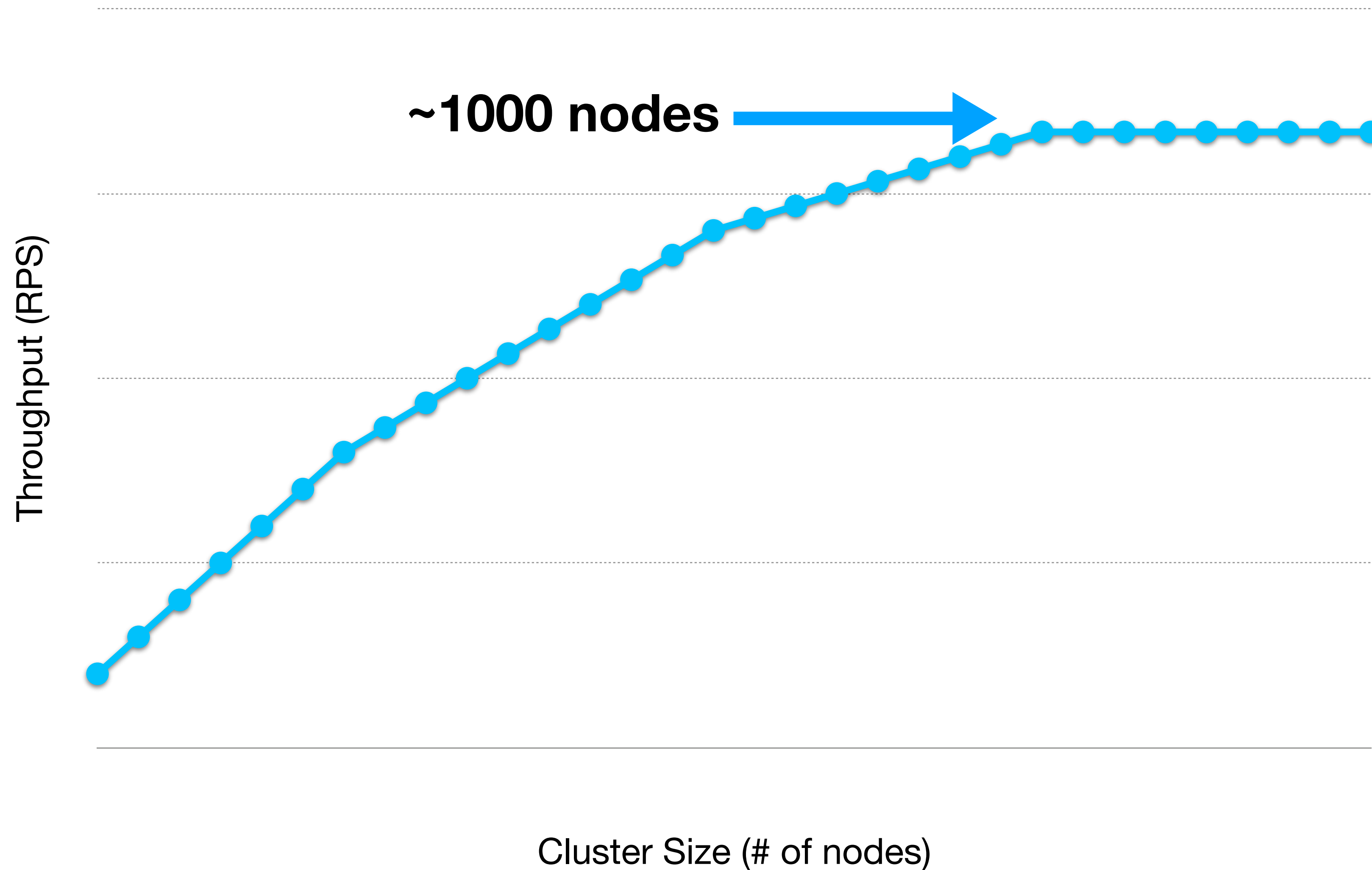
Dinesh A. Joshi

# Agenda

- History

- Performance & Efficiency

- Reliability & Stability

- Operability & Observability

- Summary

# Cassandra in Community

- Cassandra 2.x (stable, widely adopted)

- Cassandra 3.0.x (stable)

- Cassandra 3.11.x (ok, perf improvements over 3.0.x)

# Throughput vs Cluster Size

**~1000 nodes** →

Throughput (RPS)

Cluster Size (# of nodes)

# Cassandra 4.0 Changes

# Repair

# What is repair?

- Cassandra is an AP system (CAP theorem)

- Nodes can be out of sync

- Eventually consistent

- Repair ensures nodes are consistent

# Incremental Repair

- Pre 4.0 repair took long time

- Run repairs all the time!

- Shorter data reconciliation times

- Enables Transient Replication

# Impact

- Pre 4.0 repair took hours / days / weeks

- Run repair continuously

# Transient Replication

# Replication

**Replication Factor = 3**



**All nodes are <u>full</u> replicas**

# Transient Replication

- Voting with witnesses

- Cheap quorums

- Brings storage efficiency

# Impact

- Up to 33% reduction for RF=3

- Leverages Incremental Repair

# Async Internode Messaging

- Lower Latencies (**40%** lower avg **60%** lower p99)

- Memory Efficiency (~10x reduction)

- Scalable internode encryption (~4x throughput)

- Better throughput & response times (~2x vs 3.0)

# Zero Copy Streaming

- Speeds up all Streaming operations (~5x)

- IO Bound (Disk, NIC)

- Dramatically reduces MTTR

- Lowers operational cost

# Impact



Source: https://issues.apache.org/jira/browse/CASSANDRA-14765

# BTree Build Performance

- Affects hot path

- Bulk loads BTree

- ~2.6x speed up in throughput

# Reliability & Stability

- Checksummed Native Protocol

- Checksummed SSTable Metadata

# Audit Logging

- Logs everything

- Performant (Binary Logging)

- Helps in compliance for Enterprises

# Virtual Tables

# Virtual Tables

- Table backed by an API

- Queried through CQL

- NO JMX!

- Driver support

# Virtual Tables

```
cqlsh> select * from virtual.tables5 where keyspace_name = 'my_ks' and metric > 'memtable' and
metric < 'memtableZ' ALLOW FILTERING;

 keyspace_name | table_name                | metric              | value
---------------+---------------------------+---------------------+----------------
        my_ks |         monitoring_example |    memtableOnHeapSize |  {"value":95201}
        my_ks |         monitoring_example |   memtableOffHeapSize |  {"value":44811}
        my_ks |         monitoring_example | memtableLiveDataSize |  {"value":42128}
        my_ks |         monitoring_example | memtableColumnsCount |    {"value":248}
        my_ks |         monitoring_example |   memtableSwitchCount |      {"count":4}
...
```

# Virtual Tables

```
cqlsh> show VARIABLES ;

 variable                          | value
-----------------------------------+---------------------------------------------
                     authenticator |                         AllowAllAuthenticator
                        authorizer |                           AllowAllAuthorizer
                     auto_snapshot |                                         true
        batch_size_fail_threshold_in_kb |                                     50
        batch_size_warn_threshold_in_kb |                                      5
        batchlog_replay_throttle_in_kb |                                    1024
          cas_contention_timeout_in_ms |                                    1000
                      cluster_name |                                     snapshot
             column_index_size_in_kb |                                         64
              commit_failure_policy |                                         stop
               commitlog_directory |       /Users/jeff/.ccm/snapshot/node1/commitlogs
           commitlog_segment_size_in_mb |                                 33554432
                    commitlog_sync |                                     periodic
          commitlog_sync_period_in_ms |                                    10000
        compaction_throughput_mb_per_sec |                                     16
              concurrent_compactors |                                          2
```

# SSL Certificates Hot Reloading

- Certificates hot reload on update

- Optional Manual trigger via nodetool

- No disruption to live traffic

- Operators love it

# Full Query Logging (FQL)

- High performance query capture

- Replay tool (fqltool replay)

- Compare tool (fqltool compare)

- Useful for replaying traffic

# Summary (4.0 vs 3.0)

- Better operability

- Better scalability

- Better latencies

- Faster recovery