# Cache & Concurrency considerations for a high performance Cassandra

SriSatish Ambati

Performance, Riptano, Cassandra

Azul Systems & OpenJDK

@srisatish

# Trail ahead

Elements of Cache Performance

Metrics, Monitors

JVM goes to BigData Land!

Examples

Lucandra, Twissandra

Cassandra Performance with JVM Commentary

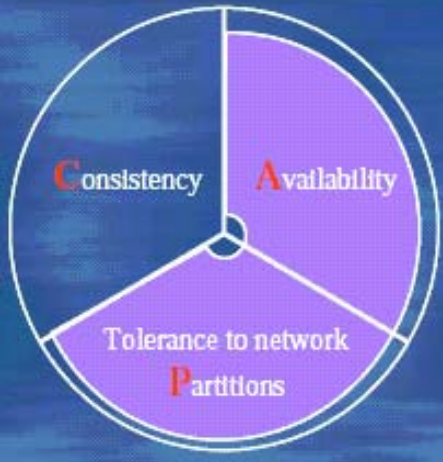Runtime Views

Non Blocking HashMap

Locking: concurrency

Garbage Collection

# A feather in the CAP

- Eventual Consistency
  - Levels
  - Doesn't mean data loss (journaled)
- SEDA
  - Partitioning, Cluster & Failure detection, Storage engine mod
  - Event driven & non-blocking io
  - Pure Java



**Forfeit Consistency**

Consistency  Availability

Tolerance to network Partitions

**Examples**
- Coda
- Web cachinge
- DNS

**Traits**
- expirations/leases
- conflict resolution
- optimistic

PODC Keynote, July 19, 2000

Count what is countable, measure what is measurable, and what is not measurable, make measurable
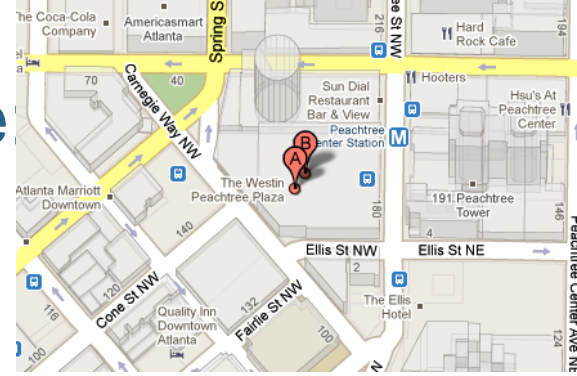
-Galileo

# Elements of Cache Performance Metrics

- Operations:
  - Ops/s: Puts/sec, Gets/sec, updates/sec
  - Latencies, percentiles
  - Indexing
- # of nodes – scale, elasticity
- Replication
  - Synchronous, Asynchronous (fast writes)
- Tuneable Consistency
- Durability/Persistence
- Size & Number of Objects, Size of Cache
- # of user clients

# Elements of Cache Performance "Think Locality"

- Hot or Not: The 80/20 rule.
  - A small set of objects are very popular!
  - What is the most RT tweet?
- Hit or Miss: Hit Ratio
  - How effective is your cache?
  - LRU, LFU, FIFO.. Expiration
- Long-lived objects lead to better locality.
- Spikes happen
  - Cascading events
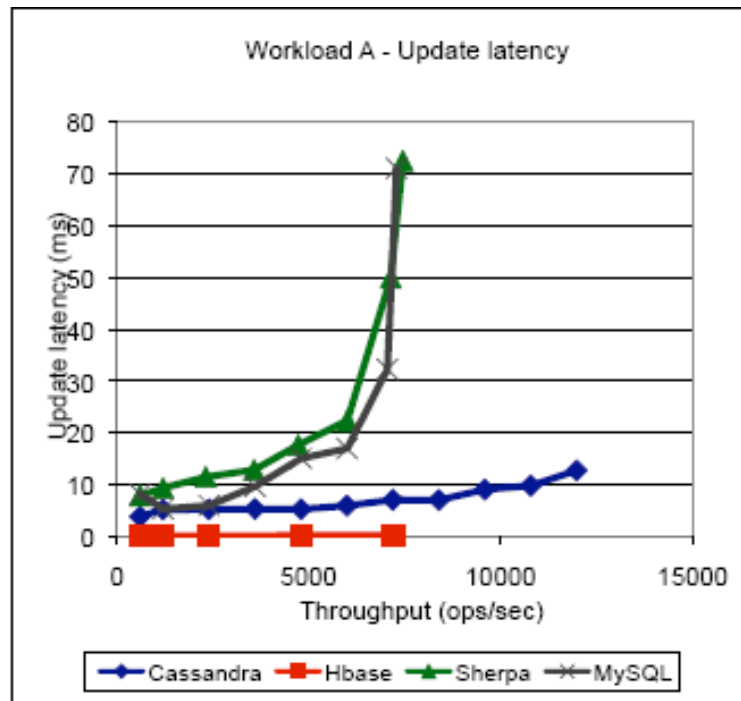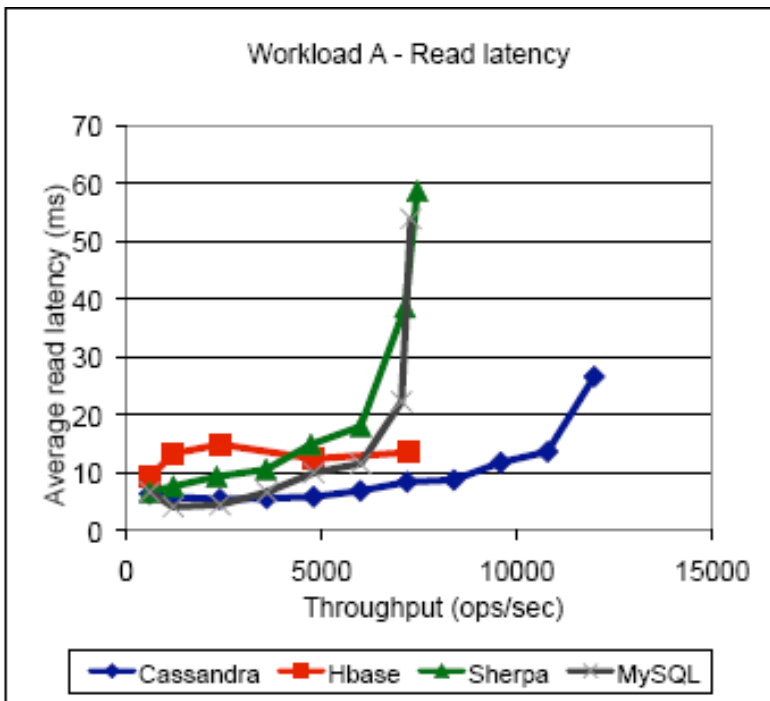  - Cache Thrash: full table scans

# Real World Performance

- Facebook Inbox
  - Writes:0.12ms, Reads:15ms @ 50GB data
- Twitter performance
  - Twissandra (simulation)
- Cassandra for Search & Portals
  - Lucandra, solandra (simulation)
- ycbs/PNUTS benchmarks
  - 5ms read/writes @ 5k ops/s (50/50 Update heavy)
  - 8ms reads/5ms writes @ 5k ops/s (95/5 read heavy)
- Lab environment
  - ~5k writes per sec per node, <5ms latencies
  - ~10k reads per sec per node, <5ms latencies
- Performance has improved in newer versions

# yahoo cloud store benchmark
## 50/50 – Update Heavy

# yahoo cloud store benchmark
## 95/5 – read heavy

# in BigData Land!

Limits for scale

- Locks : **synchronized**
  - Can't use all my multi-cores!
  - **java.util.collections** also hold locks
  - **Use non-blocking collections!**
- (de)Serialization is expensive
  - Hampers object portability
  - **Use avro, thrift!**
- Object overhead
  - average enterprise collection has 3 elements!
  - **Use byte[ ], primitives where possible!**
- Garbage Collection
  - Can't throw memory at the problem!
  - **Mitigate, Monitor, Measure foot print**

# Tools

- What is the JVM doing:
  - dtrace, hprof, introscope, jconsole, visualvm, yourkit, azul zvision
- Invasive JVM observation tools
  - bci, jvmti, jvmdi/pi agents, jmx, logging
- What is the OS doing:
  - dtrace, oprofile, vtune
- What is the network disk doing:
  - Ganglia, iostat, lsof, netstat, nagios

# furiously *fast* writes

client
issues
write

partitioner

find node

n2

n1

apply to memory

commit log

- Append only writes
  - Sequential disk access
- No locks in critical path
- Key based atomicity

# furiously *fast* writes

- Use separate disks for commitlog
  - Don't forget to size them well
  - Isolation difficult in the cloud..
- Memtable/SSTable sizes
  - Delicately balanced with GC
- `memtable_throughput_in_mb`

# Cassandra on EC2 cloud

# Cassandra on EC2 cloud

# Compactions

K1 < Serialized data >
K2 < Serialized data >
K3 < Serialized data >
--
--
--

Sorted

K2 < Serialized data >
K10 < Serialized data >
K30 < Serialized data >

--
--

Sorted

DELETED

K4 < Serialized data >
K5 < Serialized data >
K10 < Serialized data >

--
--

Sorted

## MERGE  SORT

Index File

Loaded in memory

K1   Offset
K5   Offset
K30  Offset
Bloom Filter

Sorted

K1 < Serialized data >
K2 < Serialized data >
K3 < Serialized data >
K4 < Serialized data >
K5 < Serialized data >
K10 < Serialized data >
K30 < Serialized data >

Data File

# Compactions

- Intense disk io & mem churn
- Triggers GC for tombstones
- Minor/Major Compactions
- Reduce priority for better reads
- Other Parameters -
    - `CompactionManager.`
    `minimumCompactionThreshold=xxxx`

# Example: compaction in realworld, cloudkick

# reads design

# reads performance

- BloomFilter used to identify the right file
- Maintain column indices to look up columns
  - Which can span different SSTables
- Less io than typical b-tree
- Cold read: Two seeks
  - One for Key lookup, another row lookup
- Key Cache
  - Optimized in latest cassandra
- Row Cache
  - Improves read performance
  - GC sensitive for large rows.
- Most (google) applications require single row transactions*

*Sanjay G, BigTable Design, Google.

# Client Performance Marshal Arts: Ser/Deserialization

- Clients dominated by Thrift, Avro
  - Hector, Pelops
- Thrift: upgrade to latest: 0.5, 0.4
- No news: java.io.Serializable is S.L..O….W
- Use "transient"
- avro, thrift, proto-buf
- Common Patterns of Doom:
  - Death by a million gets

# Serialization + Deserialization uBench

- http://code.google.com/p/thrift-protobuf-compare/wiki/BenchmarkingV2

total

| | |
|---|---|
| java-manual | 3234 |
| kryo-opt | 4375 |
| kryo | 4583 |
| protobuf/protostuff-core | 5094 |
| protobuf/protostuff-runtime | 5712 |
| protobuf | 6949 |
| protobuf/activemq+alt | 7199 |
| scala/sbinary | 7817 |
| thrift | 9228 |
| thrift-compact | 9532 |
| xml/manual-aalto | 11160 |
| json/jackson-manual | 11265 |
| avro | 12575 |
| json/protostuff-core | 12698 |
| avro-generic | 13233 |
| json/protostuff-runtime | 14117 |
| xml/javolution | 14950 |
| xml/manual-woodstox | 16244 |
| json/jackson-databind | 17187 |
| hessian | 23282 |
| aalto-xstream+c | 28258 |
| xml/manual-fastinfo | 31018 |
| woodstox-xstream+c | 37012 |
| fastinfo-xstream+c | 43989 |
| xml/xstream+c | 61663 |
| java-built-in | 74099 |
| scala/java-built-in | 113707 |
| json/protobuf | 141689 |
| json/google-gson | 218143 |
| aalto-xstream | 225905 |
| xml/xstream | 227399 |
| fastinfo-xstream | 227943 |
| woodstox-xstream | 228237 |

**Leading the Wave of Open Source**

# Adding Nodes

- New nodes
  - Add themselves to busiest node
  - And then Split its Range
- Busy Node starts transmit to new node
- Bootstrap logic initiated from any node, cli, web
- Each node capable of ~40MB/s
  - Multiple replicas to parallelize bootstrap
- UDP for control messages
- TCP for request routing

Leading the Wave
of Open Source

# inter-node comm

- Gossip Protocol
  - It's exponential
  - (epidemic algorithm)
- Failure Detector
  - Accrual rate phi
- Anti-Entropy
  - Bringing replicas to uptodate

# Bloom Filter: in full bloom

- "constant" time
- size:compact
- false positives
- Single lookup
  for key in file
- Deletion
- Improve
  - Counting BF
  - Bloomier filters

# Birthdays, Collisions & Hashing functions

- Birthday Paradox
   - For the N=21 people in this room
   - Probability that at least 2 of them share same birthday is ~0.47
- Collisions are real!
- An unbalanced HashMap behaves like a list O(n) retrieval
- Chaining & Linear probing
- Performance Degrades
-  with 80% table density

# the devil's in the details

# CFS

- All in the family!
- denormalize

ColumnFamilyStore
- logger : Logger
- cacheSavingExecutor : ScheduledThreadPoolExecutor
- flushSorter : ExecutorService
- flushWriter : ExecutorService
- postFlushExecutor : ExecutorService
- memtablesPendingFlush : Set<Memtable>
- table : Table
- columnFamily : String
- partitioner : IPartitioner
- mbeanName : String
- memtableSwitchCount : int
- fileIndexGenerator : AtomicInteger
- memtable : Memtable
- indexedColumns : SortedMap<byte[], ColumnFamilyStore>
- binaryMemtable : AtomicReference<BinaryMemtable>

# Memtable

- In-memory
- ColumnFamily specific
- **`throughput`** determines size before flush
- Larger memtables can improve reads



```
▽ © Memtable
   ⬛F logger : Logger
   ▫  isFrozen : boolean
   ▫ F currentThroughput : AtomicInteger
   ▫ F currentOperations : AtomicInteger
   ▫ F creationTime : long
   ▫ F columnFamilies : ConcurrentNavigabl
   ○ F cfs : ColumnFamilyStore
   ▫ F THRESHOLD : int
   ▫ F THRESHOLD_COUNT : int
   ● Memtable(ColumnFamilyStore)
   ●▵ compareTo(Memtable) : int
   ● getCurrentThroughput() : int
   ● getCurrentOperations() : int
   ▲ isThresholdViolated() : boolean
```

# SSTable

- MemTable "flushes" to a SSTable
- Immutable after
- Read: Multiple SSTable lookups possible
- Chief Execs:
  - SSTableWriter
  - SSTableReader



```
▽ ⓖ ᴬ SSTable
  ⚠ˢᶠ logger : Logger
  ⚲ˢᶠ COMPONENT_DATA : String
  ⚲ˢᶠ COMPONENT_INDEX : String
  ⚲ˢᶠ COMPONENT_FILTER : String
  ⚲ˢᶠ COMPONENT_STATS : String
  ⚲ˢᶠ COMPONENT_COMPACTED : String
  ⚲ˢᶠ TEMPFILE_MARKER : String
  ○ ᶠ descriptor : Descriptor
  ○ ᶠ components : Set<Component>
  ○ ᶠ metadata : CFMetaData
  ○ ᶠ partitioner : IPartitioner
  ◇ estimatedRowSize : EstimatedHistogram
  ◇ estimatedColumnCount : EstimatedHistogram
  ◇ ᶜ SSTable(Descriptor, CFMetaData, IPartitioner)
  ◇ ᶜ SSTable(Descriptor, Set<Component>, CFMetaData, IPartitioner)
  ● getEstimatedRowSize() : EstimatedHistogram
  ● getEstimatedColumnCount() : EstimatedHistogram
  ● ˢ conditionalDelete(Descriptor, Set<Component>) : boolean
  ● getFilename() : String
  ● getColumnFamilyName() : String
  ● getTableName() : String
```

# Write: Runtime threads

# Writes: runtime mem

| Name | Retained Size | |
|------|---------------|---|
| ⊟ ◻ java.lang.Thread [Stack Local, Thread] "CompactionExecutor: 1" native ID: 0xC99 | 90,548,336 | 8% |
| ⊟ ◻ org.apache.cassandra.io.sstable.SSTableWriter [Stack Local] | 90,538,944 | 8% |
| ⊟ ◻ org.apache.cassandra.io.util.BufferedRandomAccessFile | 67,109,216 | 6% |
| ◻ byte[67108864] = {5, 0, 2, 67, 48, 0, 0, 4, -109, 123, -34, -85, 19, -73, 0, 0, 0, :67,108,888 | 67,108,888 | 6% |
| ⊞ ◻ org.apache.cassandra.io.sstable.SSTableWriter$IndexWriter | 23,424,296 | 2% |
| ⊟ ◻ java.lang.Thread [Stack Local, Thread] "FlushWriter: 1" native ID: 0xC94 | 75,833,680 | 7% |
| ⊞ ◻ org.apache.cassandra.io.util.BufferedRandomAccessFile [Stack Local] | 67,109,216 | 6% |
| ⊟ ◻ org.apache.cassandra.io.sstable.SSTableWriter [Stack Local] | 8,719,960 | 1% |
| ⊟ ◻ org.apache.cassandra.io.sstable.SSTableWriter$IndexWriter | 8,714,376 | 1% |
| ⊟ ◻ org.apache.cassandra.io.util.BufferedRandomAccessFile | 8,388,960 | 1% |
| ◻ byte[8388608] = {0, 19, 50, 0, 0, 10, 48, 48, 56, 49, 49, 53, 49, 50, 53, 53, 0, | 8,388,632 | 1% |
| ⊟ ◻ java.util.concurrent.ConcurrentSkipListMap | 56,998,752 | 5% |
| ◻ java.util.concurrent.ConcurrentSkipListMap$HeadIndex | 56,998,640 | 5% |
| ⊞ ◻ org.apache.cassandra.io.CompactionIterator [Stack Local] | 32,538,968 | 3% |

| Paths from GC Roots: Alt+1 | Allocations: Alt+2 | Class Hierarchy: Alt+3 | Incoming References: Alt+4 | Quick Info: Alt+5 |

Class hierarchy for class selected in the upper table

| ▲ Name |
|--------|
| ⊟ C java.lang.Object |
| ⊟ C org.apache.cassandra.io.sstable.SSTable |
| ▶ C org.apache.cassandra.io.sstable.SSTableWriter |

# Example: Java Overheads

| Name | Objects | | Shallow Size | | Retained Size | |
|---|---|---|---|---|---|---|
| C byte[] | 4,655,099 | 24 % | 429,267,408 | 37 % | 429,267,408 | 37 % |
| C java.util.concurrent.ConcurrentSkipListMap$Node | 2,429,035 | 13 % | 97,161,400 | 8 % | 326,202,360 | 28 % |
| C org.apache.cassandra.db.Column | 1,734,935 | 9 % | 69,397,400 | 6 % | 128,156,760 | 11 % |
| C int[] | 1,086,052 | 6 % | 50,821,304 | 4 % | ≈ 50,821,304 | 4 % |
| C java.util.concurrent.ConcurrentSkipListMap$Index | 850,132 | 4 % | 34,005,280 | 3 % | ≈ 40,240,600 | 4 % |
| C java.math.BigInteger | 846,013 | 4 % | 40,608,624 | 4 % | ≈ 66,725,448 | 6 % |
| C org.apache.cassandra.dht.BigIntegerToken | 832,511 | 4 % | 19,980,264 | 2 % | ≈ 77,433,312 | 7 % |
| C org.apache.cassandra.db.DecoratedKey | 814,979 | 4 % | 26,079,328 | 2 % | ≈ 125,155,920 | 11 % |
| C java.util.concurrent.ConcurrentSkipListMap$HeadIndex | 539,996 | 3 % | 25,919,808 | 2 % | 237,476,888 | 21 % |
| C org.apache.cassandra.io.sstable.IndexSummary$KeyPosit | 469,441 | 2 % | 15,022,112 | 1 % | ≈ 101,348,704 | 9 % |
| C org.apache.thrift.protocol.TField | 376,899 | 2 % | 12,060,768 | 1 % | ≈ 12,060,768 | 1 % |
| C java.util.concurrent.atomic.AtomicInteger | 356,169 | 2 % | 8,548,056 | 1 % | ≈ 8,548,056 | 1 % |

Object Explorer: Alt+ | Generations: Alt+3 | Reachability: Alt+4 | Class Loaders: Alt+5 | Allocations: Alt+6 | Class Statics: Alt+

Objects selected in the upper table

Class name, string value, thread name or ID (Press "Enter" to apply / hint on syntax):

| Name | Retained Size | Shallow Size |
|---|---|---|
| ⊟ O java.util.concurrent.ConcurrentSkipListMap$Node | 63,392 | 40 |
| ⊞ **next** ⇨ O java.util.concurrent.ConcurrentSkipListMap$Node | 61,880 | 40 |
| ⊟ **value** ⇨ O org.apache.cassandra.db.ColumnFamily | 1,288 | 64 |
| ⊞ **<class>** ⇨ C org.apache.cassandra.db.ColumnFamily | 1,400 | 1,024 |
| ⊞ **columns** ⇨ O java.util.concurrent.ConcurrentSkipListMap | 1,160 | 88 |
| ⊞ **type** ⇨ O org.apache.cassandra.db.ColumnFamilyType | 32 | 32 |
| ⊞ **cfid** ⇨ O java.lang.Integer | 24 | 24 |

# writes: monitors

Group by [Waiting/blocked thread ▼]  then group by [Monitor class ▼]  ☐ Show blocked threads only

| Name | Time (ms) | | Count | |
|---|---|---|---|---|
| ⊟ Waiting thread 🔵 EXPIRING-MAP-TIMER-1 native ID: 0x4AA group: 'main' | 275,325 | 97% | 50 | 0% |
| └ on monitor of class 🆑 java.util.TaskQueue | 275,325 | 97% | 50 | 0% |
| ⊟ Waiting thread 🔵 EXPIRING-MAP-TIMER-2 native ID: 0x4AB group: 'main' | 275,268 | 97% | 50 | 0% |
| └ on monitor of class 🆑 java.util.TaskQueue | 275,268 | 97% | 50 | 0% |
| ⊟ Blocked thread 🔴 CompactionExecutor:1 native ID: 0x4B3 group: 'main' | 261,074 | 92% | 1,964 | 1% |
| ⊟ on monitor of class 🆑 sun.security.provider.Sun | 261,074 | 92% | 1,964 | 1% |
| ── held by thread 🟢 pool-1-thread-22 native ID: 0x1758 group: 'main' | 5,093 | 2% | 22 | 0% |
| ── held by thread 🟢 pool-1-thread-3 native ID: 0x1732 group: 'main' | 5,009 | 2% | 35 | 0% |
| ── held by thread 🟢 MutationStage:30 native ID: 0x4D3 group: 'main' | 4,670 | 2% | 39 | 0% |

| Name | Time (ms) | Count |
|---|---|---|
| ⊟ 🔍 java.security.**Provider.getService**(String, String) | 261,074 100% | 1,964 100 |
| ⊟ 🔍 org.apache.cassandra.utils.**ReducingIterator.computeNext**() | | |
| ⊟ 🔍 com.google.common.collect.**AbstractIterator.tryToComputeNext**() | | |
| ⊟ 🔍 com.google.common.collect.**AbstractIterator.hasNext**() | | |
| └ 🔍 java.lang.**Thread.run**() | | |

Group by [Monitor class ▼]  then group by [Waiting/blocked thread ▼]  ☐ Show blocked threads only

| Name | Time (ms) | | Count | |
|---|---|---|---|---|
| ⊟ Monitor of class 🆑 sun.security.provider.Sun | 114,760 | 100% | 39,477 | 57% |
| ⊟ was waited by thread 🔴 CompactionExecutor:1 native ID: 0x4B3 group: 'main | 105,233 | 92% | 765 | 1% |
| ── that was blocked by thread 🟢 pool-1-thread-22 native ID: 0x1758 group | 3,649 | 3% | 11 | 0% |
| ── that was blocked by thread 🟢 pool-1-thread-3 native ID: 0x1732 group: | 2,722 | 2% | 17 | 0% |
| ── that was blocked by thread 🟢 MutationStage:6 native ID: 0x4BB group: 'm | 2,570 | 2% | 14 | 0% |
| ── that was blocked by thread 🟢 MutationStage:18 native ID: 0x4C7 group: ' | 2,450 | 2% | 16 | 0% |
| ── that was blocked by thread 🟢 MutationStage:31 native ID: 0x4D4 group: ' | 2,404 | 2% | 12 | 0% |
| ── that was blocked by thread 🟢 pool-1-thread-10 native ID: 0x1740 group | 2,260 | 2% | 16 | 0% |
| ── that was blocked by thread 🟢 MutationStage:29 native ID: 0x4D2 group: ' | 2,219 | 2% | 15 | 0% |

| Name | Time (ms) | Count |
|---|---|---|
| ⊟ 🔍 java.security.**Provider.getService**(String, String) | 105,233 100% | 765 100% |
| ⊟ 🔍 org.apache.cassandra.utils.**ReducingIterator.computeNext**() | | |
| ⊟ 🔍 com.google.common.collect.**AbstractIterator.tryToComputeNext**() | | |
| ⊟ 🔍 com.google.common.collect.**AbstractIterator.hasNext**() | | |
| └ 🔍 java.lang.**Thread.run**() | | |

**Leading the Wave**
**of Open Source**

# U U I D

- java.util.UUID is slow
  - static use leads to contention

SecureRandom

- Uses /dev/urandom for seed initialization

  `-Djava.security.egd=file:/dev/urandom`

- PRNG without file is atleast 20%-40% better.
- Use TimeUUIDs where possible – much faster
- JUG – java.uuid.generator

- *http://github.com/cowtowncoder/java-uuid-generator*
- *http://jug.safehaus.org/*
- *http://johannburkard.de/blog/programming/java/Java-UUID-generators-compared.html*

Leading the Wave
of Open Source

# synchronized

- Coarse grained locks
- io under lock
- Stop signal on a highway
- java.util.concurrent does not mean no locks
- Non Blocking, Lock free, Wait free collections

# Scalable Lock-Free Coding Style

- Big Array to hold Data
- Concurrent writes via: CAS & Finite State Machine
  - No **locks**, no **volatile**
  - *Much* faster than locking under heavy load
  - Directly reach main data array in 1 step
- Resize as needed
  - Copy Array to a larger Array on demand
  - Use State Machine to help copy
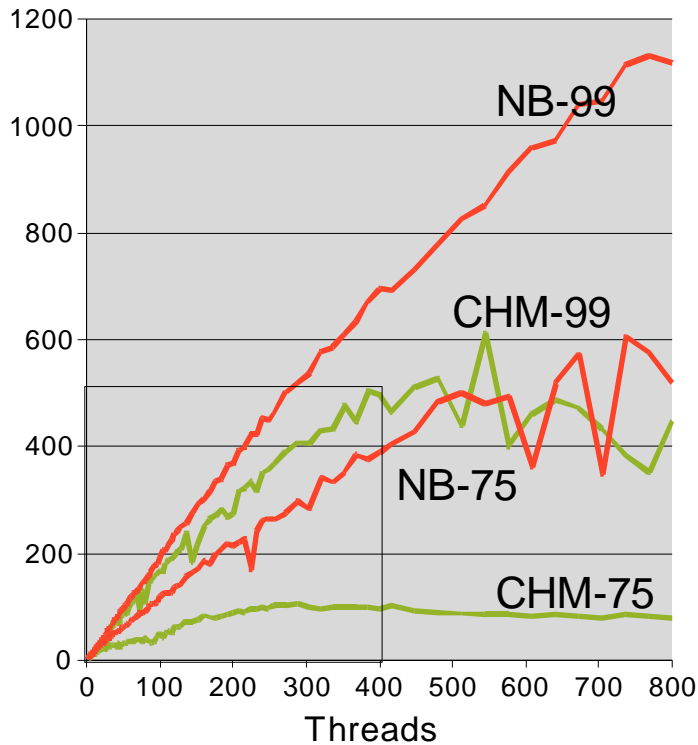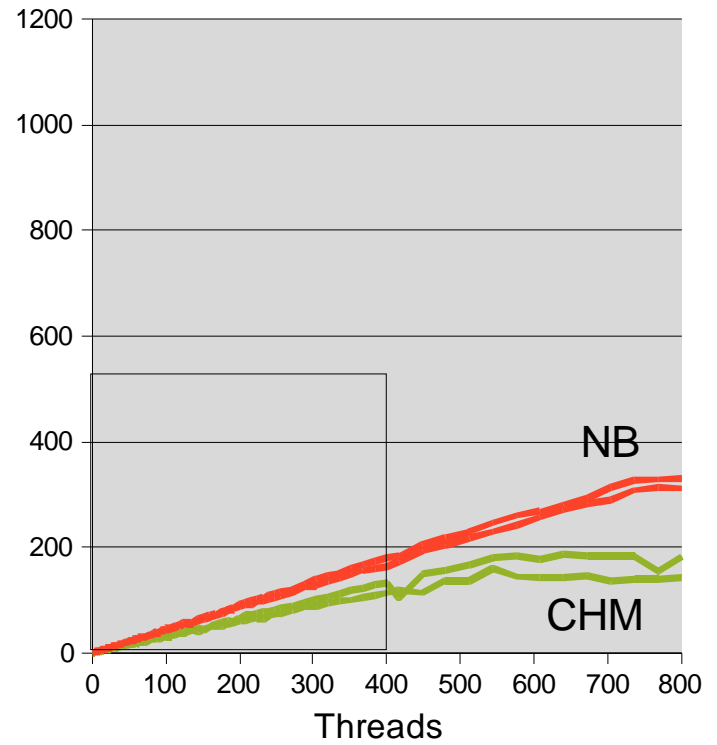  - "Mark" old Array words to avoid missing late updates

# Non-Blocking HashMap

## Azul Vega2 – 768 cpus

1K Table



1M Table

# Cassandra uses High Scale Non-Blocking Hashmap

```java
public class BinaryMemtable implements IFlushable
{
…
   private final Map<DecoratedKey,byte[]> columnFamilies =
    new NonBlockingHashMap<DecoratedKey, byte[]>();
/* Lock and Condition for notifying new clients about Memtable
    switches */
   private final Lock lock = new ReentrantLock(); Condition condition;
…
}
public class Table
{
…
   private static final Map<String, Table> instances =  new
    NonBlockingHashMap<String, Table>();
…
}
```
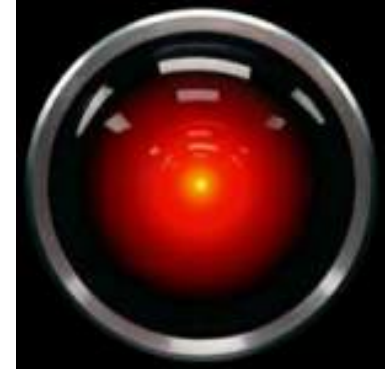
# GC-sensitive elements within Cassandra

- Compaction triggers System.gc()
  - Tombstones from files
- "GCInspector"
- Memtable Threshold, sizes
- SSTable sizes
- Low overhead collection choices

# Garbage Collection

- Pause Times
  if stop_the_word_FullGC > ttl_of_node
  => failed requests; failure accrual & node repair.
- Allocation Rate
  – New object creation, insertion rate
- Live Objects (residency)
  – if residency in heap > 50%
  – GC overheads dominate.
- Overhead
  – space, cpu cycles spent GC
- 64-bit not addressing pause times
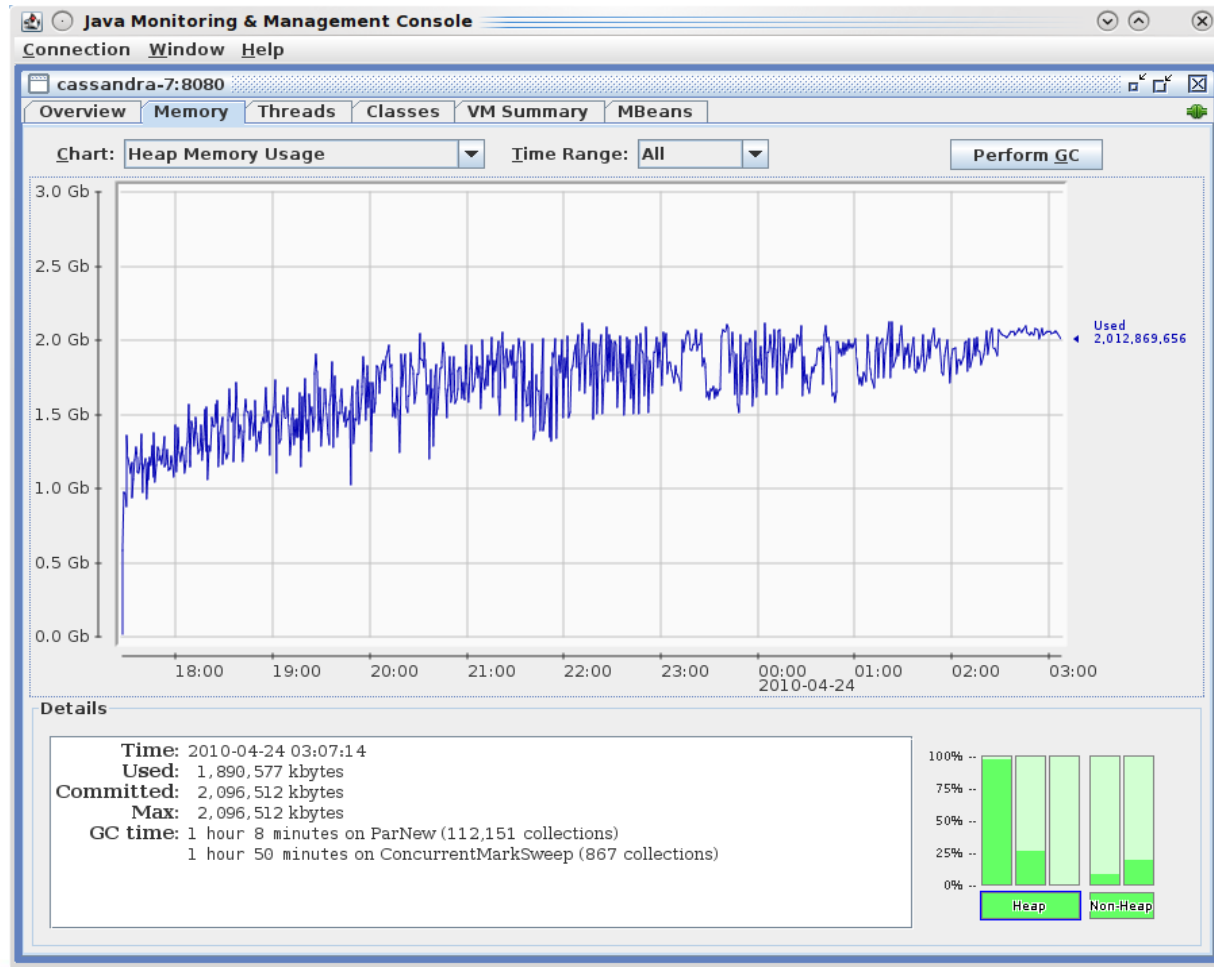  – Bigger is not better!

# Memory Fragmentation

- Fragmentation
  - Performance degrades over time
  - Inducing "Full GC" makes problem go away
  - Free memory that cannot be used
- Reduce occurrence
  - Use a compacting collector
  - Promote less often
  - Use uniform sized objects
- Solution – *unsolved*
  - Use latest CMS with CR:6631166
  - Azul's Zing JVM & Pauseless GC

# CASSANDRA-1014
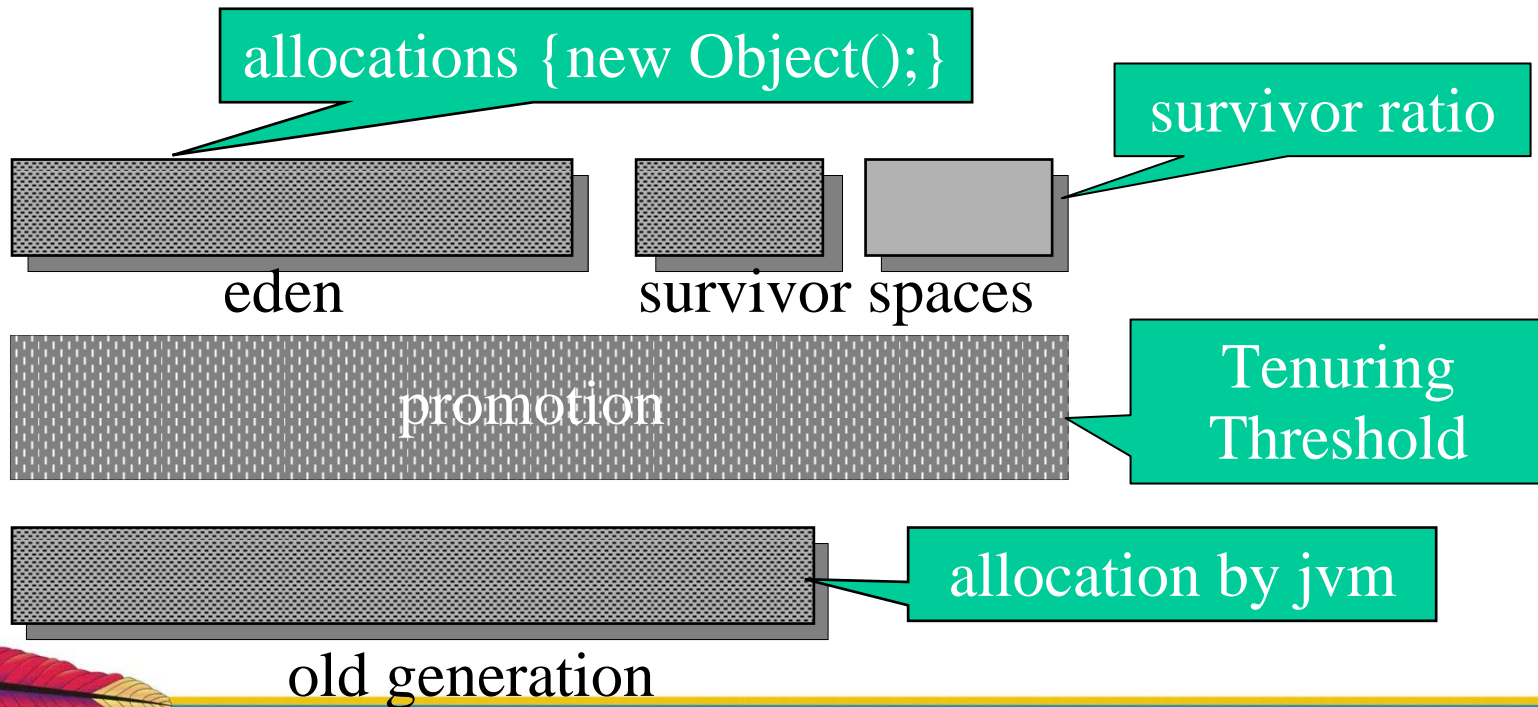
# Best Practices:
# Garbage Collection

- ## GC Logs are cheap even in production

  - -Xloggc:/var/log/cassandra/gc.log
  - -XX:+PrintGCDetails
  - -XX:+PrintGCTimeStamps -XX:+PrintTenuringDistribution
  - -XX:+PrintHeapAtGC

- ## Slightly expensive ones:

  - -XX:PrintFLSStatistics=2 -XX:CMSStatistics=1
  - -XX:CMSInitiationStatistics

# Sizing: Young Generation

- Should we set –Xms == -Xmx ?
- Use –Xmn (fixed eden)

allocations {new Object();}

survivor ratio

eden    survivor spaces

promotion

Tenuring
Threshold

allocation by jvm

old generation

# Tuning CMS

- Don't promote too often!
  - Frequent promotion causes fragmentation
- Size the generations
  - Min GC times are a function of Live Set
  - Old Gen should host steady state comfortably
- Parallelize on multicores:
  - -XX:ParallelCMSThreads=4
  - -XX:ParallelGCThreads=4
- Avoid CMS Initiating heuristic
  - -XX:+UseCMSInitiationOccupanyOnly
- Use Concurrent for System.gc()
  - -XX:+ExplicitGCInvokesConcurrent

# Summary

**Design & Implementation of Cassandra takes advantage of strengths while avoiding common JVM issues.**

- **Locks:**
  - **Avoids locks in critical path**
  - **Uses non-blocking collections, TimeUUIDs!**
  - **Still Can't use all my multi-cores..?**
    **>> Other bottlenecks to find!**
- **De/Serialization:**
  - **Uses avro, thrift!**
- **Object overhead**
  - **Uses mostly byte[ ], primitives where possible!**
- **Garbage Collection**
  - **Mitigate: Monitor, Measure foot print.**
  - **Work in progress by all jvm vendors!**

**Cassandra starts from a great footing from a JVM standpoint and will reap the benefits of the platform!**

# Q&A

## References

- *Verner Wogels, Eventually Consistent* http://www.allthingsdistributed.com/2008/12/eventually_consistent.html

- *Bloom, Burton H. (1970), "Space/time trade-offs in hash coding with allowable errors"*

- *Avinash Lakshman, http://static.last.fm/johan/nosql-20090611/cassandra_nosql.pdf*

- *Eric Brewer, CAP* http://www.cs.berkeley.edu/~brewer/cs262b-2004/PODC-keynote.pdf

- *Tony Printzeis, Charlie Hunt, Javaone Talk http://www.scribd.com/doc/36090475/GC-Tuning-in-the-Java*

- http://github.com/digitalreasoning/PyStratus/wiki/Documentation

- **http://www.cs.cornell.edu/home/rvr/papers/flowgossip.pdf**

- Cassandra on Cloud, http://www.coreyhulen.org/?p=326

Leading the Wave
of Open Source

Count what is countable, measure what is measurable, and what is not measurable, make measurable

-Galileo