: instaclustr

Processing 200K Transactions per Second with Apache Spark and Apache Cassandra

Ben Bromhead

Boston Apache Spark Meetup

3 May 2018

: instaclustr

Or... we built our own metrics/monitoring stack and it was worth it... but you probably shouldn't do it... probably

Ben Bromhead

Boston Apache Spark Meetup

3 May 2018

/usr/bin/whoami

⇔ instaclustr

- Ben Bromhead, CTO of Instaclustr
- We provide managed Cassandra, Spark and Kafka in the cloud (AWS, GCP, Azure & Softlayer).
- We provide support and services as well for those in private data centers.
- Manage and support 2k+ nodes.

Agenda



- Introduction to Cassandra
- Why Spark + Cassandra
- Problem background and overall architecture
- Implementation process & lessons learned
- What's next?







NoSQL database

- Highly available
- Master less
- Linear scalability
- Low latency

- No join
- Poor index
- Restricted filtering
- No ACID

- OLTP
- Data ingestion
- Design your requests first, your model second.

Introduction a Cassandra

⇔ instaclustr

Cassandra is a Distributed Hash Table



Introduction a Cassandra





Spark

⇔ instaclustr

Spark is a Distributed Big Data Processing Framework



Spark + Cassandra

⇔ instaclustr



Spark + Cassandra

⇔ instaclustr









Problem to solve....

Problem background

⇒ instaclustr

- How to efficiently monitor > 2000 servers all running Cassandra
 - Alerting
 - Metric history
 - Alert tuning
 - Graph / dashboard
 - Multi-tenant approach
- Off the shelf systems are available but:
 - Flexible enough?
 - Learn by using our technology
 - Optimizations opportunities.

Problem background



Friend: You know drinking is bad for you you should just use off the shelf

Me:



Implementation Approach

⇔ instaclustr

- 1. Collecting Metrics + Alert
- 2. Writing metrics
- 3. Rolling Up metrics
- 4. Presenting metrics
 - ~ 9(!) months (with quite a few detours and distractions)

I want it done by naptime



Solution Overview: instaclustr monitoring pipeline

Node

many

Node





Data model

⇔ instaclustr

- CREATE TABLE instametrics.events_raw_5m (
 - host text,
 - bucket_time timestamp,
 - service text,
 - time timestamp,
 - metric double,
 - state text,
 - PRIMARY KEY ((host, bucket_time, service), time)



Data Model



CREATE TABLE instametrics.host (<u>host</u> text PRIMARY KEY

CREATE TABLE instametrics.service_per_host (

host text,

service text,

PRIMARY KEY (host, service)

Writing metrics

Key lessons:

- Aligning Data Model with DTCS (now TWCS)
 - Initial design did not have time value in partition key
 - Settled on bucketing by 5 mins
 - Enables DTCS to work
 - Works really well for extracting data for roll-up
 - Adds complexity for retrieving data
- Batching of writes
 - Found batching of 200 rows per insert to provide optimal throughput and client load
- Controlling data volumes from column family metrics
 - Limited, rotating set of CFs per check-in
- Managing back pressure is important



⇔ instaclustr

Rolling Up metrics





 Developing functional solution was easy, getting to acceptable performance was hard (and time consuming) but seemed easy once we'd solved it



Data Model

⇔ instaclustr

CREATE TABLE instametrics_rollup.events_rollup_300 (

bucket_time timestamp,

host text,

service text,

time timestamp,

avg double,

max double,

min double,

state text,

PRIMARY KEY ((bucket_time, host, service), time)

Rolling Up metrics





- Developing functional solution was easy, getting to acceptable performance was hard (and time consuming) but seemed easy once we'd solved it
- Keys to performance?
 - Align raw data partition bucketing with roll-up timeframe (5 mins)
 - Use repartitionByCassandraReplica to align Spark partitions with Cassandra partitions
 - Use joinWithCassandra table to extract the required data 2-3x performance improvement over alternate approaches

⇔ instaclustr



Workers

Read Tuning:

spark.cassandra.input.fetch.size_in_rows
spark.cassandra.input.reads_per_sec

Write Tuning: spark.cassandra.output.throughput_mb_per_sec

5min – hourly – daily rollup

US East (Northern Virginia) · Amazon Web Services (VPC)					
ID	Address	State	CPU	Cores in Use	Memory in Use
worker-20170704000432-10.224.157.209-42317	10.224.157.209	ALIVE		2/4	1.8 GB / 20.0 GB
worker-20170710001210-10.224.16.0-43409	10.224.16.0	ALIVE	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	2/4	1.8 G8 / 20.0 GB
worker-20170710001333-10.224.173.178-40501	10.224.173.178	ALIVE	~~~~~~	2/4	1.8 GB / 20.0 GB

Presenting metrics



- Generally, just worked
- Main challenge was dealing with how to find latest data in rollup buckets when not all data is reported in each data set

Optimisation with Cassandra Aggregation



• Upgraded to Cassandra 3.7 and change code to use Cassandra aggregates: val RDDJoin = sc.cassandraTable[(String, String)]("instametrics", "service_per_host") .filter(a => broadcastListEventAll.value.map(r => a._2.matches(r)).foldLeft(false)(_ || _)) .map(a => (a._1, dateBucket, a._2)) .repartitionByCassandraReplica("instametrics", "events_raw_5m", 100) .joinWithCassandraTable("instametrics", "events_raw_5m", 100) .joinWithCassandraTable("instametrics", "events_raw_5m", SomeColumns("time", "state", FunctionCallRef("avg", Seq(Right("metric")), Some("avg")), FunctionCallRef("max", Seq(Right("metric")), Some("max")), FunctionCallRef("min", Seq(Right("metric")), Some("min")))).cache()

50% reduction in roll-up job runtime (from 5-6 mins to 2.5-3mins) with reduced CPU usage

Rolling Up metrics

⇔ instaclustr

"So I should just sit down here while you paint my por - oh you're done"



What's Next



- Riemann straight to Spark Streaming
 - Spark Streaming for 5 min roll-ups rather than save and extract

- Scale-out by adding nodes is working as expected
- Continue to add additional metrics to roll-ups as we add functionality
- Plan to introduce more complex analytics & feed historic values back to Reimann for use in alerting

Further info:

- () instaclustr
- Scaling Riemann: \checkmark https://www.instaclustr.com/blog/2016/05/03/post-500-nodes-high-availability-scalability-with-rie mann/
- **Riemann Intro:**

https://www.instaclustr.com/blog/2015/12/14/monitoring-cassandra-and-it-infrastructure-with-rie mann/

- Instametrics Case Study: https://www.instaclustr.com/project/instametrics/
- Multi-DC Spark Benchmarks: https://www.instaclustr.com/blog/2016/04/21/multi-data-center-sparkcassandra-benchmark-round <u>-2/</u>
- Top Spark Cassandra Connector Tips: https://www.instaclustr.com/blog/2016/03/31/cassandra-connector-for-spark-5-tips-for-success/
- Cassandra 3.x upgrade: https://www.instaclustr.com/blog/2016/11/22/upgrading-instametrics-to-cassandra-3/
- Cassandra Spark MLIB:

https://www.instaclustr.com/third-contact-monolith-part-c-pod/ Contact with a Monolith:

Part C - In the Pod

Friday 29th September 2017 by Paul Brebner

A simple classification problem: Will the Monolith react? Is it safe?! Maybe a cautious approach to a bigger version of the Monolith (2km long) in a POD that is only 2m in diameter...



Third Contact with a Monolith: Part C - In the Pod

instaclustr

⇒ instaclustr

Ben Bromhead CTO, Instaclustr ben@instaclustr.com

info@instaclustr.com

www.instaclustr.com

@instaclustr