# GrafanaCloud platform requirements

- Scalable for customers, but also scalable for our SRE!

- Fault tolerance and automated recovery

- Service discovery

- Horizontal Scaling

- Resource management

- Isolation between tenants

…. Kubernetes to the rescue; we're all in!

# Kubernetes: our not so secret weapon

- A consistent platform for on-prem and SaaS deployments

  - Shippable SaaS

- Fully managed options reduce SRE burden

  - GKE (Google Kubernetes Engine)

- Also run vanilla K8s on bare metal

  - Packet.net

- Or wherever our customers want us to be

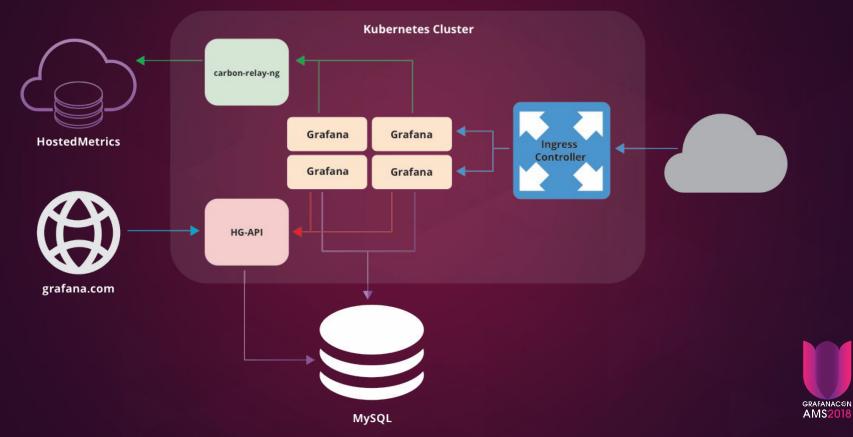  - Eg. Azure AKS, AWS EKS, colo, for GrafanaCloud Private Deployments

# Hosted Grafana

- A fully dedicated Grafana instance running the latest stable release

- One-Click installation of plugins from grafana.com

- Custom domain and authentication

- Anything config setting possible

- Who better to support it than the core Grafana team?

# Hosted Grafana
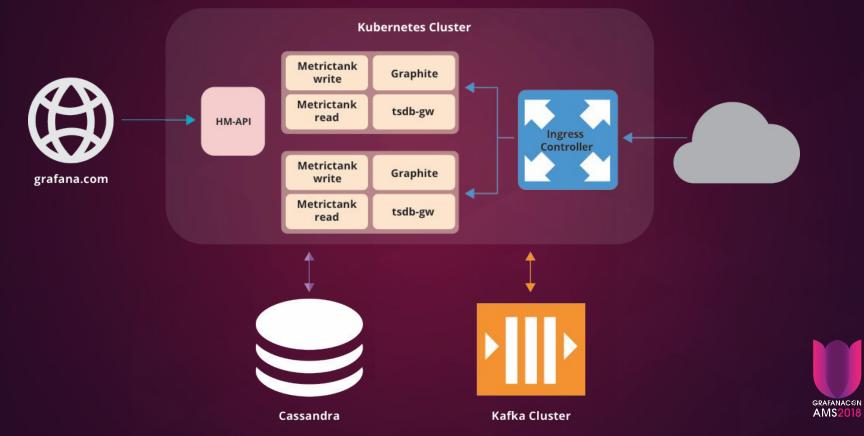
# Hosted Grafana Instance Dashboard

# Hosted Metrics

- Unlimited* Scale

- Support for large metric volume (hundreds of millions of DPM)

- Fast query response times to support alerting

- Tunable for different workloads (eg. retention, cache, redundancy)

- Fault tolerant

# Hosted Metrics

# Hosted Metrics - core components

GrafanaLabs metrictank: https://github.com/grafana/metrictank

- ○ Query engine compatible with Graphite and PromQL
  Keeps most data cached in memory for exceptionally fast query times
- ○ Compresses and aggregates data then saves it to the backend store
  Inspired by Facebook Gorilla (similar algo as Prometheus and InfluxDB) < 2 bytes per point

Apache Kafka: https://kafka.apache.org/
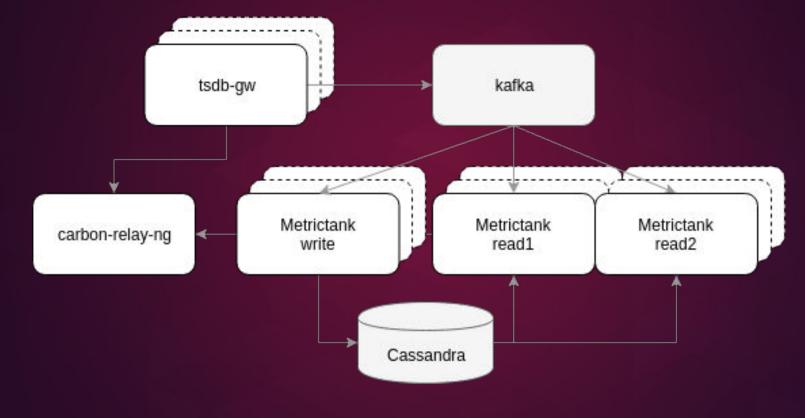
- ○ Distributed Queue
  Provides resilience; we always need to accept data

Apache Cassandra: http://cassandra.apache.org/ or Google Bigtable

- ○ Long term storage of metric data.
- ○ Horizontally scalable

GRAFANACON
AMS2018

# Hosted Metrics - Components

Hosted Metrics Customer Dashboard

# Cache Performance
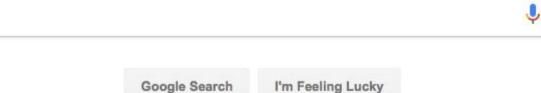
Kubernetes

Bigtable

+

# Google Cloud Bigtable

Misha Brukman
Product Manager

Google

# How do we …

… run **containerized workloads** at scale?

**Need:** Deploy, scale and upgrade microservices quickly and efficiently

**Solution:** Borg, Kubernetes (open source)

**Google Kubernetes Engine**

... build a petabyte-scale **analytics database**?

**Need:** Massive data index files took weeks to rebuild. We needed random read/write access

**Solution:** Bigtable

**Google Cloud Bigtable**

# Technologies to support Google products



2002 — GFS

2004 — MapReduce

2006 — Bigtable

2008 — Dremel

2009 — Colossus

2010 — Flume

2011 — Megastore

2012 — Spanner, MillWheel

2014 — F1

2015 — Borg

2016 — TensorFlow

… when scale is a solved problem

1 Billion users

# Technologies to support Google products

2002 — GFS

2004 — MapReduce

2006 — Bigtable

2008 — Dremel

2009 — Colossus

2010 — Flume

2011 — Megastore

2011 — Spanner

2013 — MillWheel

2014 — F1

2015 — Borg

2015 — TensorFlow

| 2002 | 2004 | 2006 | 2008 | 2010 | 2012 | 2014 | 2016 |

# Now available on Google Cloud Platform

## Compute

App Engine

Kubernetes Engine

Compute Engine

## Storage & Databases

Storage

Bigtable

Spanner

Cloud SQL

Datastore

## Big Data

BigQuery

Pub/Sub

Dataflow

Dataproc

Datalab

## Machine Learning

Vision API

ML Engine

Speech API

Translate API

Google Cloud Bigtable

# Google Cloud Bigtable

- Fully-managed NoSQL database

- **Built-in support for time series**

- Seamless scalability for throughput

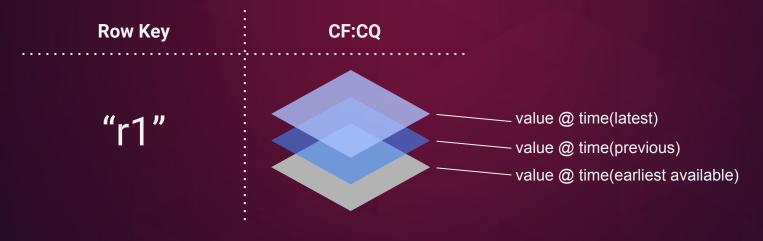- Learns and adjusts to access patterns

# Bigtable data model

- NoSQL (no-join) distributed key-value store, designed to scale-out
- has only one index (the row-key)
- supports atomic single-row transactions
- unwritten cells in do not take up any space

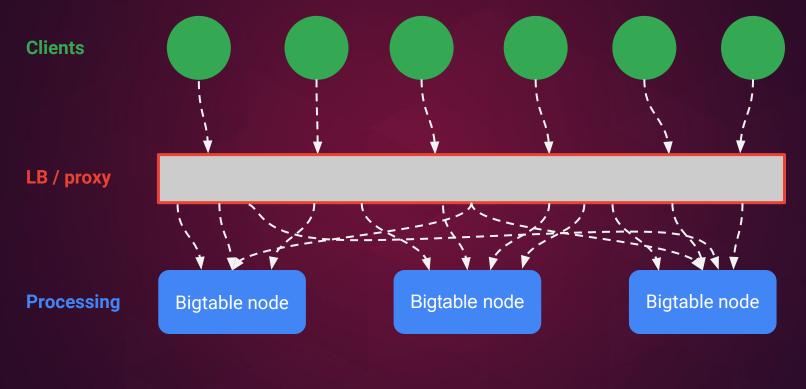| Row Key | Column-Family-1 | | Column-Family-2 | |
|---|---|---|---|---|
| | *Column-Qualifier-1* | *Column-Qualifier-2* | *Column-Qualifier-1* | *Column-Qualifier-2* |
| r1 | r1, cf1:cq1 | r1, cf1:cq2 | r1, cf2:cq1 | r1, cf2:cq2 |
| r2 | r2, cf1:cq1 | r2, cf1:cq2 | r2, cf2:cq1 | r2, cf2:cq2 |

# 3D database structure enables time series

- every cell is ***versioned*** (default is timestamp on server)
- garbage collection retains latest version (configurable)
- expiration (optional) can be set at column-family level
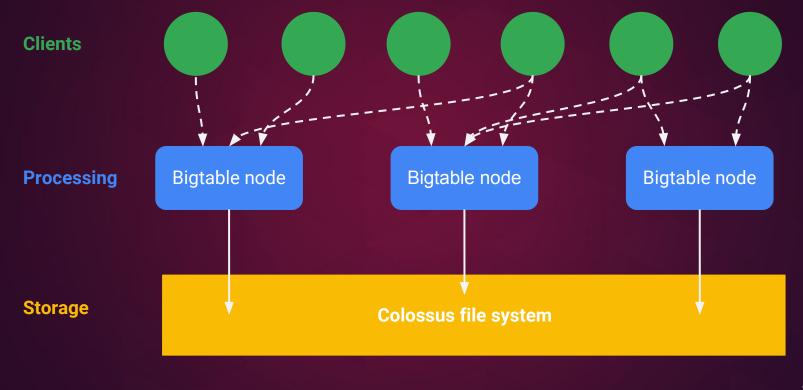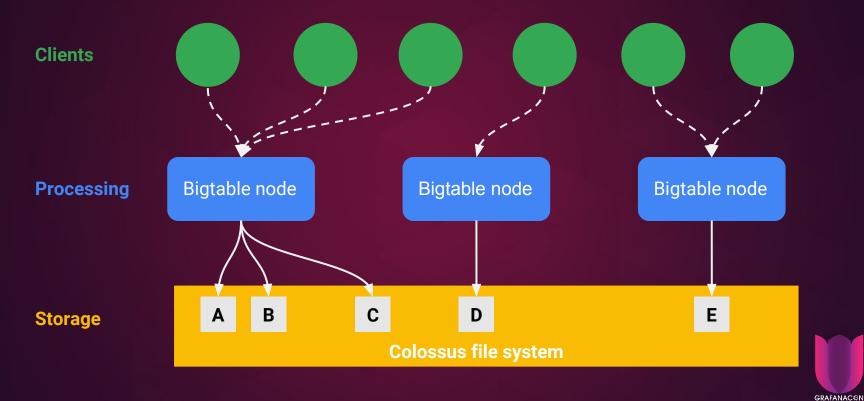- periodic compaction reclaims unused space from cells

**Row Key**                    **CF:CQ**

"r1"

value @ time(latest)

value @ time(previous)

value @ time(earliest available)

# Bigtable high-level architecture



**Clients**

**LB / proxy**

**Processing**

Bigtable node     Bigtable node     Bigtable node

GRAFANACON
AMS2018

Bigtable learns access patterns...

# Bigtable provides seamless resizing

# Bigtable provides linear scalability in performance

# Bigtable provides linear scalability in performance

# Great long tails

Single digit ms at the 99%

- Native scheduler protects serving
  path from compactions
- No garbage collection
- Very fast tablet reassignment

# Google Cloud Bigtable

- Fully-managed NoSQL database

- **Built-in support for time series**

- Seamless scalability for throughput

- Learns and adjusts to access patterns



GRAFANACON AMS2018

# Metrictank

Dieter Plaetinck
Principal Engineer

# Project
# Not product

# Data store
# Not database

# Genesis

(not the band)

# Requirements for Worldping TSDB

- Large scale (millions of points per second)
- Long term storage, rollups
- Resource efficient (cpu, memory, disk)
- Multi-tenant
- Open source
- Operationally friendly
- Proven technology
- Compatible with Graphite (or pluggable into Graphite)

# Didn't want to write yet another TSDB

# interesting bed time reading material

**Dieter Plaetinck** <dieter@raintank.io>                    9/14/15

to all-staff

looks like FB just released a paper describing their in house, in-mem, highly
compressed data store. they also compare it to whisper, influxdb and opentsdb.

http://www.vldb.org/pvldb/vol8/p1816-teller.pdf
https://twitter.com/armon/status/642803583050604549

**Torkel Ödegaard** <torkel@raintank.io>                    9/14/15

to Dieter

nice! bed time reading is always good to have :)

- github.com/dgryski/go-tsz
- NSQ (later Kafka)
- Cassandra
- (Elasticsearch for index)

# Timeline

- Sept 23, 2015 : First prototyping
- Dec 2015: Worldping production

Do we really want our own TSDB?

- 2016: Ad-hoc hosted metrics alpha's

Do we really want our own TSDB?

- Early 2017: Grafanacloud v1

Looks like it

- Early 2018: Grafanacloud v2

OK then. Can we add prometheus?

GRAFANACON
AMS2018

# metrictank

- service that reads from queue, compresses data to chunks. saves to DB
- Saves rollups
- Satisfies queries from memory and DB
- Input: Kafka (graphite, Prometheus, OpenTSDB, …)
- Input: direct Carbon, prometheus
- Whisper import
- Graphite function api (mix built-in and graphite-web)
- PromQL
- Can be deployed as eventually consistent cluster

# Integrating
# Not replacing

# Input options

- Kafka (carbon-relay-ng graphite, Prometheus, OpenTSDB, …)
- Plain carbon, prometheus (!!)
- Whisper importer

# Storage options

- Cassandra
- Bigtable
- (CosmosDB?)

# Output options

- Graphite api
- Prometheus api
- ...

# Data

- Chunk ringbuffer in memory
- LRU chunk cache in memory
- Storage plugin for persistence (Cassandra, …)
- Can reach ~100% memory hit rate

# Metadata (index)

- Plugin (Cassandra, …) for persistence
- Full in-memory copy
- Built-in expression handling, searching, tag index, autocomplete, etc

# Improve on Graphite

https://grafana.com/blog/2016/03/03/25-graphite-grafana-and-statsd-gotchas/

- Seamless changing of native data resolution
- Better support for churn (shortlived data)
- Multiple rollup functions, choice at query time (WIP)
- Automatic interval detection (WIP)

# Worse than Graphite

- Data must be mostly-ordered. No rewrite support
- No xFilesFactor yet

# Clustering

HA (replication)

&

horizontal scaling (partitioning/sharding)

# Clustering: HA (replication)

- Simply run # replicas desired (via orchestrator)
- Primary role (via config/orchestrator or API, not automatic)
- kafka/NSQ for tracking save state
- Kafka data backfill reduces time-to-ready

# Clustering: horizontal scaling (partitioning)

- Shard assignment tied to input (via config/orchestrator)
- Shard deterministically derived from metric name & metadata
- Index per node only for shards it "owns"
- Gossip for membership
- Queries can hit any instance, scatter+merge
- Kafka-lag based ready-state, priority, and min-available-shards

# Clustering limitation 1

primary status per instance, not shard

| node | A | B | C | D |
|------|------|------|------|------|
| shards | 0 1 | 0 1 | 2 3 | 2 3 |

# Clustering limitation 1

primary status per shard

| node | A | B | C | D |
|------|------|------|------|------|
| shards | **0** 1 | 0 **2** | **1** 3 | 2 **3** |

# Clustering limitation 2

- Rigid sharding scheme. Can't add/remove shards at will.
- => (live) cluster migration

# Clustering trade-offs

- [https://martin.kleppmann.com/2015/05/11/please-stop-calling-databases-cp-or-ap.html](https://martin.kleppmann.com/2015/05/11/please-stop-calling-databases-cp-or-ap.html)
- Kafka : very tuneable. Ours tuned for consistency -> buffering client side (rare)
- Cassandra : Eventually consistent. Tunable consistency latency trade-off
- eventually consistent. Everything streams in. Even when talking to MT directly
- Don't need transactions for monitoring data
- MT read instances depend on writers saving to Cassandra

# Use whatever makes sense for you

That's why Grafana supports graphite, influxDB, prometheus, cloudwatch, ….

That's why metrictank supports Cassandra, Bigtable, ….

# Tools

mt-aggs-explain

mt-explain

mt-index-cat

mt-index-migrate

mt-kafka-mdm-sniff

mt-kafka-mdm-sniff-out-of-order

mt-replicator-via-tsdb

mt-schemas-explain

mt-split-metrics-by-ttl

mt-store-cat

mt-update-ttl

mt-view-boundaries

mt-whisper-importer-reader

mt-whisper-importer-writer

# Tools

```
mt-index-cat -prefix statsd.prod -tags none -max-age 12h cass 'GET
http://metrictank/render?target=lowestCurrent(sumSeries({{.Name |
pattern}}),2)&from=-30min\nAuthorization: Bearer foo\n\n' \

| ./vegeta attack -rate 5 | ./vegeta report
```

# Fun under the hood stuff

- Golang issue #14812 GC bug
- Metrictank PR #136 Buffer reuse, custom json encoder, etc
- Golang contexts
- Jaeger tracing (opentracing)
- Automated chaos testing with docker-compose and pumba/tc
- profiletrigger

# Metrictank use cases

Large scale graphite installations

Long term storage prometheus

SaaS without vendor lock-in

Favor known database

# Conclusion

- Try it out, but beware
- Or try GrafanaCloud (SaaS or Private Deployment)
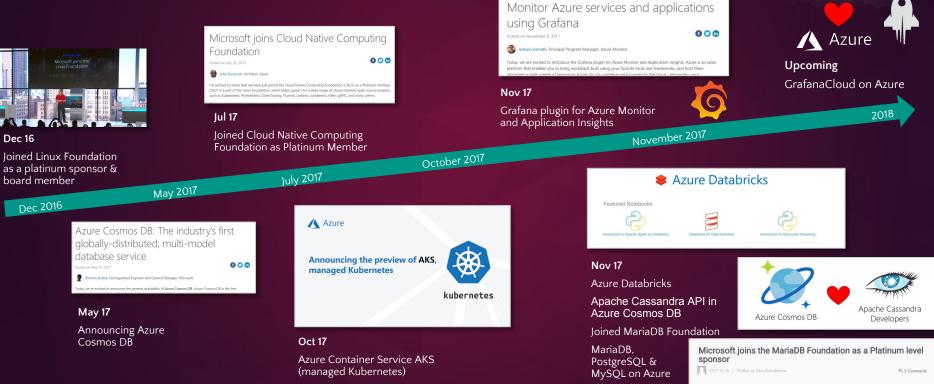
# Integrate with ecosystem
# Not divide and conquer

GRAFANACON
AMS2018

# Metrictank & Azure Cosmos DB

A globally distributed, massively scalable, multi-model database service

Table API

mongoDB®

Gremlin
G = (V, E)

cassandra

SQL

Column-family

Document

Key-value

Graph

Guaranteed low latency at the 99th percentile

Elastic scale out
of storage & throughput

Five well-defined consistency models

Turnkey global distribution

Comprehensive SLAs