

Monitoring the world: Scaling Thanos in Dynamic Prometheus Environments

Colin Douch, Observability Tech Lead @ Cloudflare



- Globally Distributed
- Highly Available
- ~50TB of metrics daily
- > 1000 sidecars
- 600 transient locations

- Globally Distributed
- Highly Available
- ~50TB of metrics daily
- > 1000 sidecars
- 600 transient locations



A Quick Introduction

Hi! I'm Colin

Observability Lead @ Cloudflare

Has bad opinions

Not British

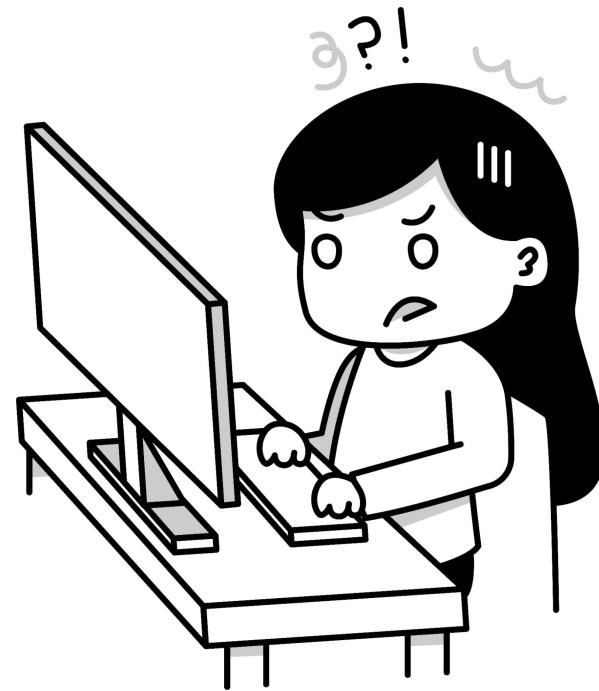
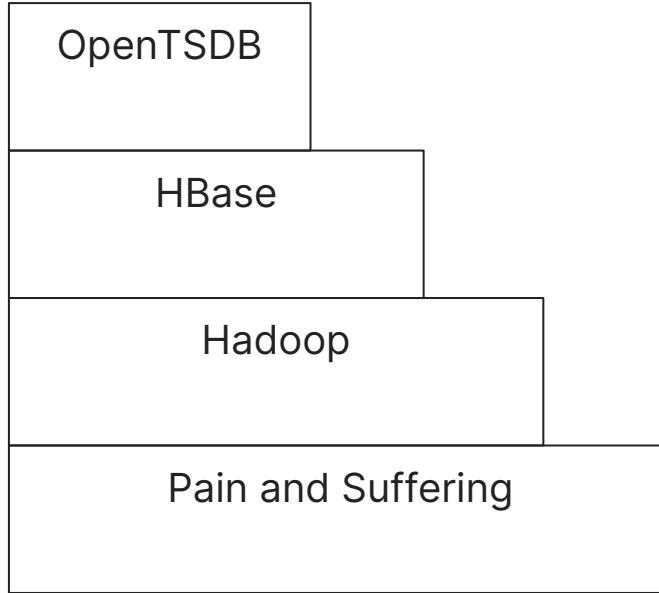


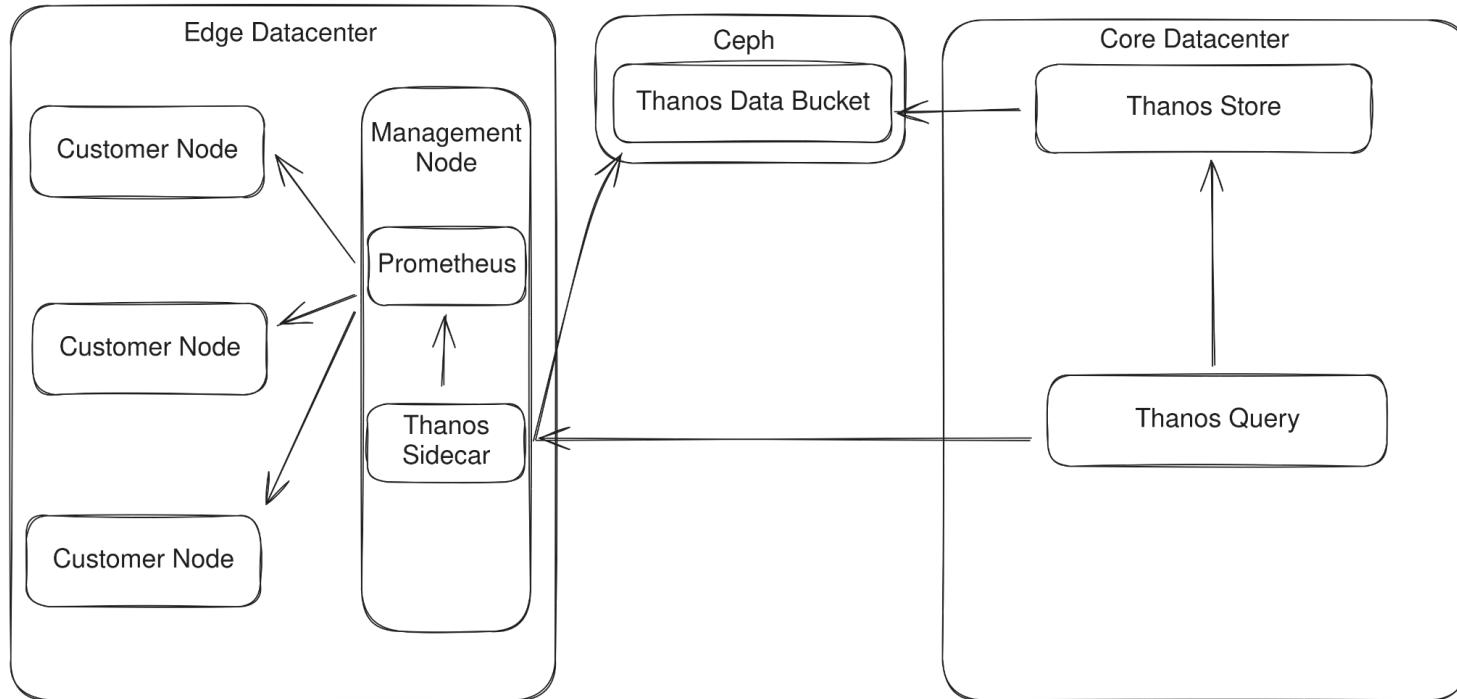


MAINTAINING OPENTSDB



ROLLING THE DICE ON A NEW UNTESTED TECHNOLOGY

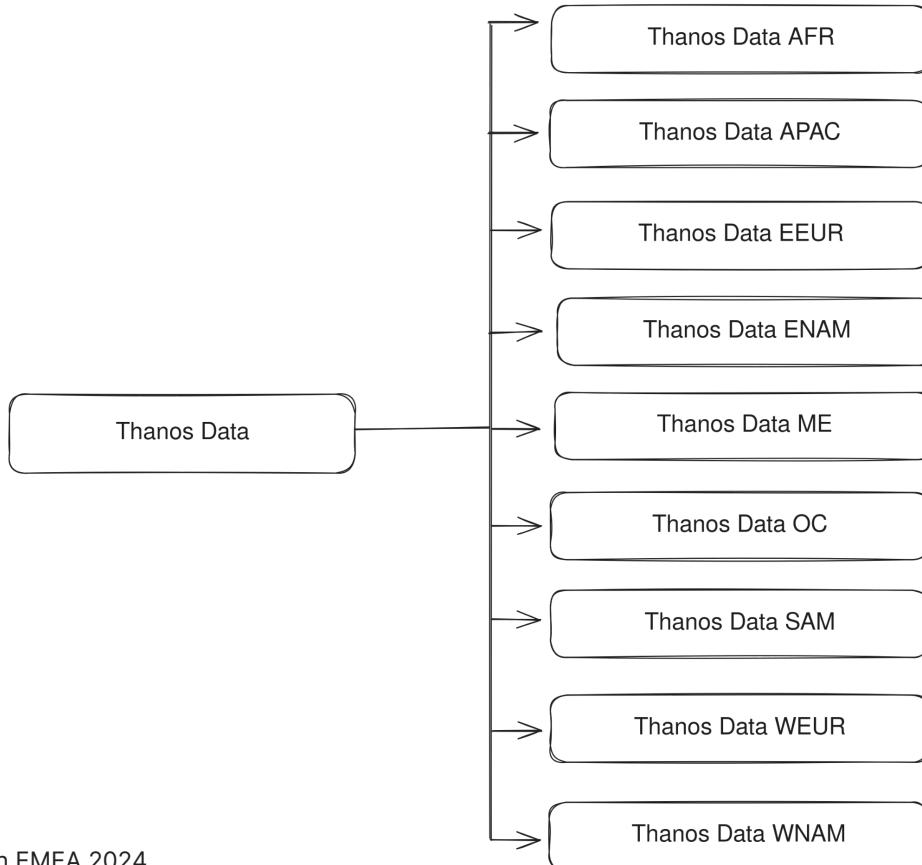


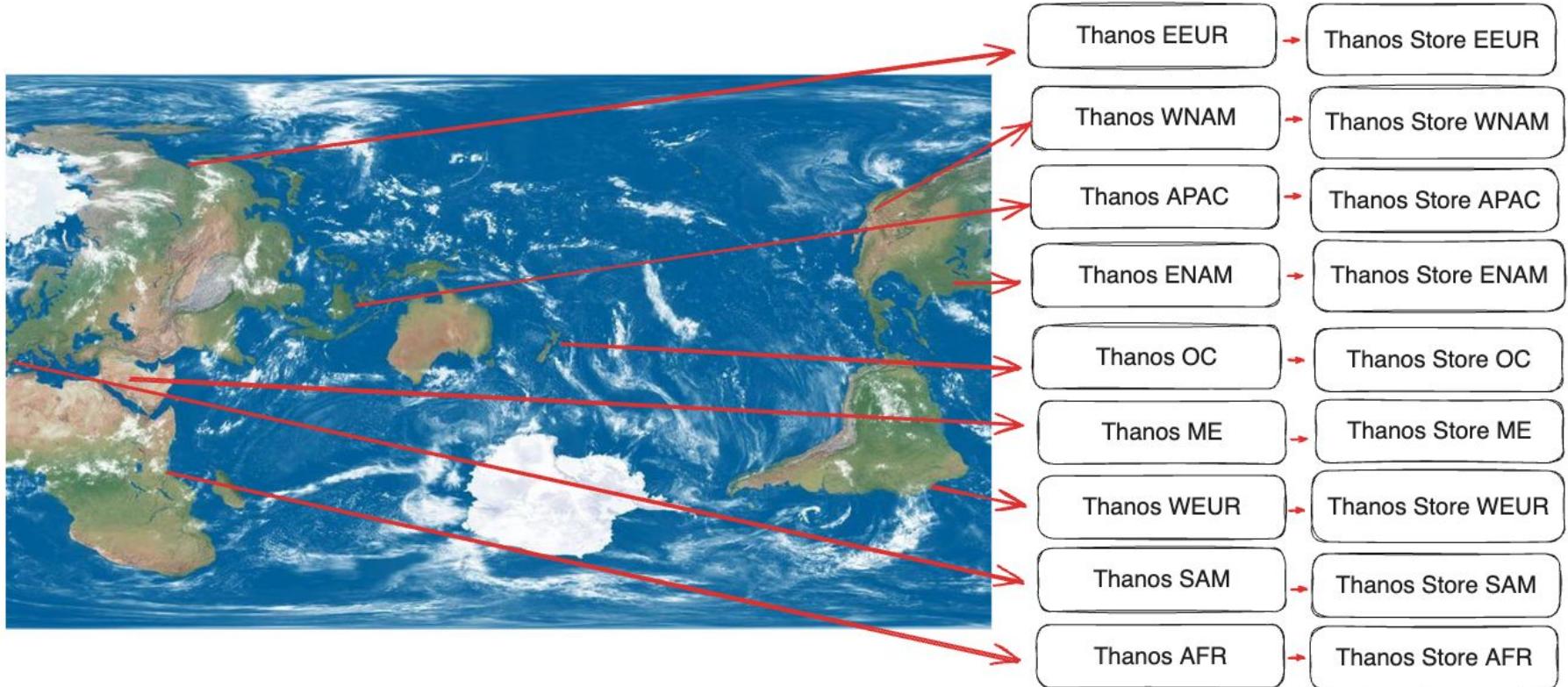


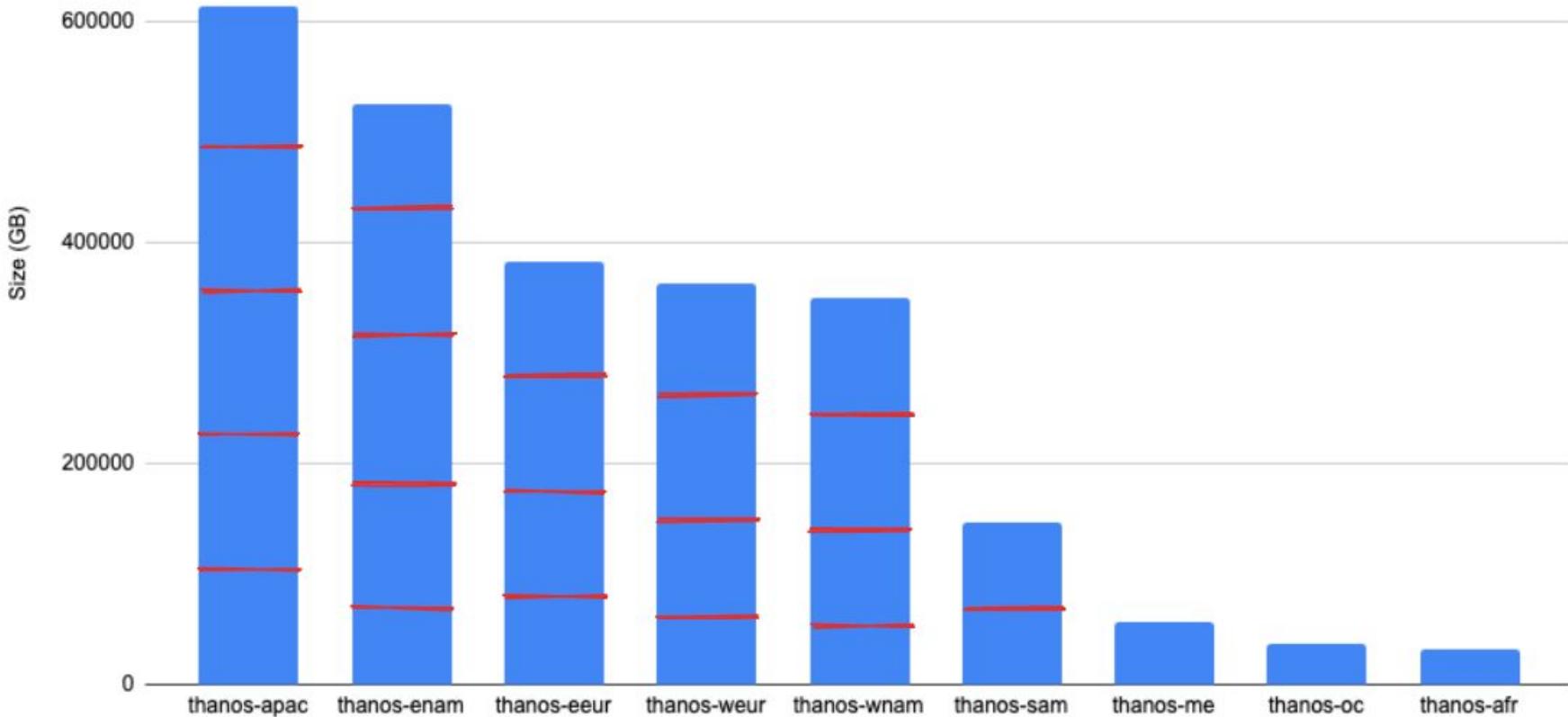
It's 2019 and we've got issues



- Stuck with Vertical Scaling, with a lack of time partitioning
- Compactors were running out of memory
- Everything was falling apart
- No one was having a good time





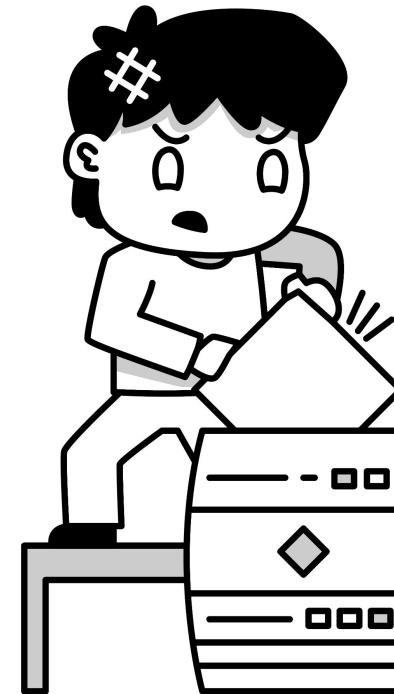


We were mostly doing **vibe
based scaling:**

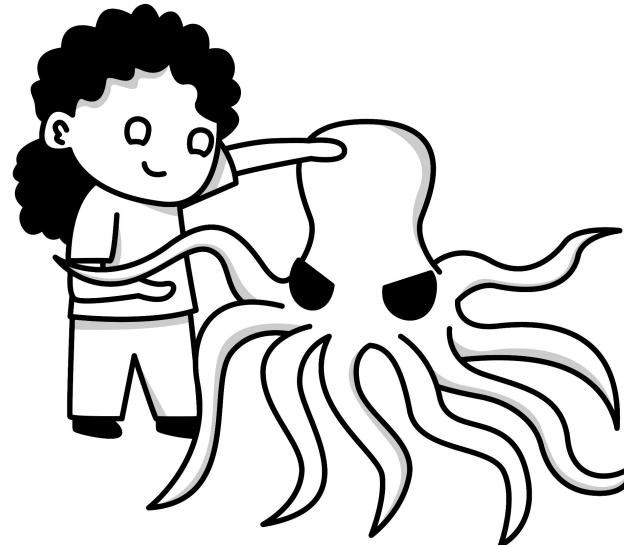
Create a new store when
things *felt* slow



Welcome to Scaling Thanos in Dynamic Prometheus Environments



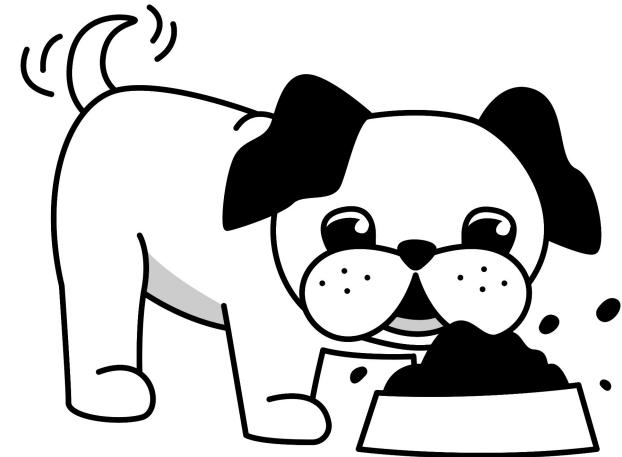
Ceph



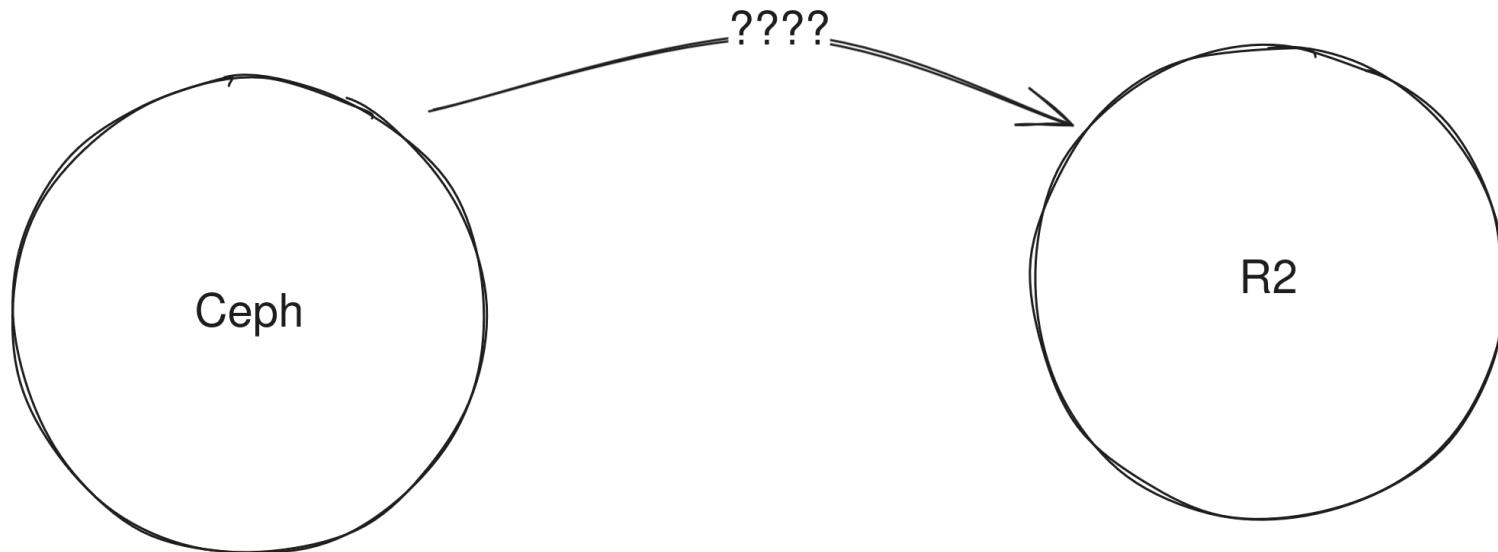
**The worst storage system,
except for all the others**

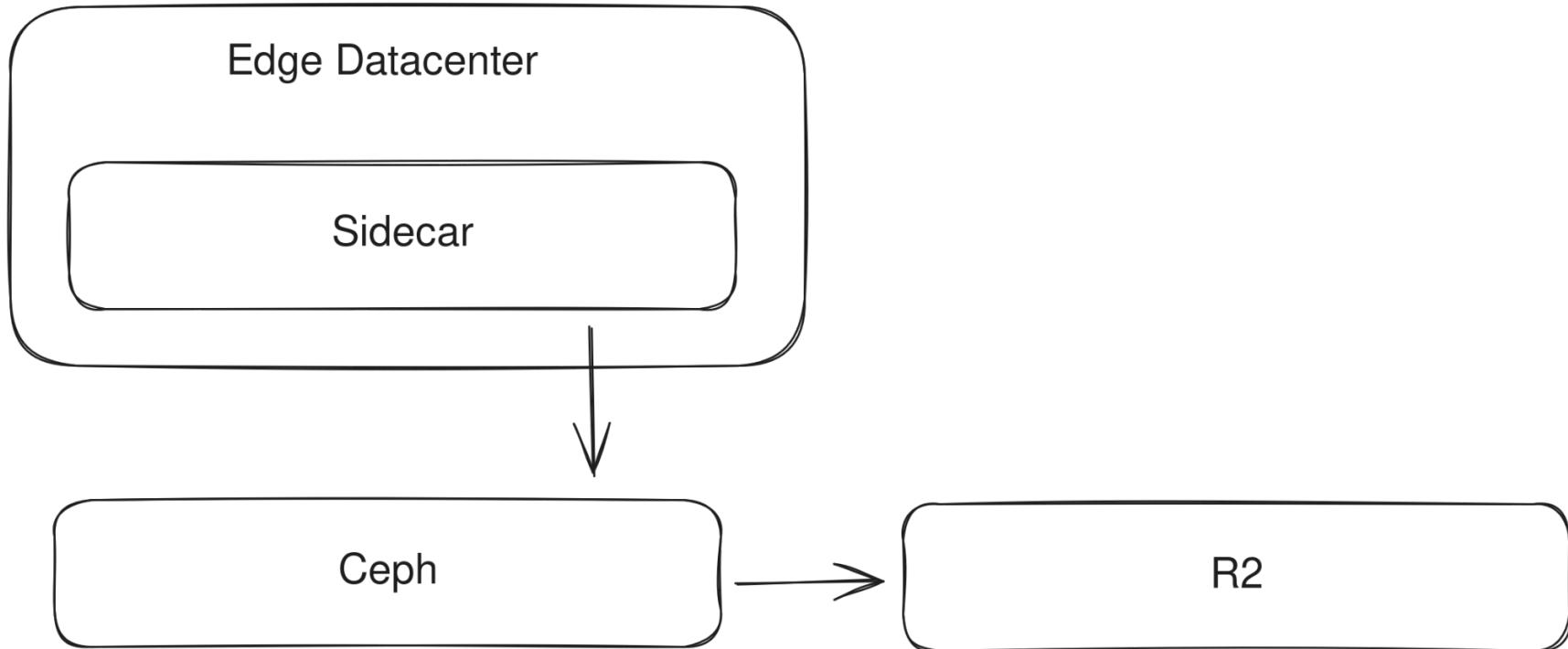
Cloudflare R2

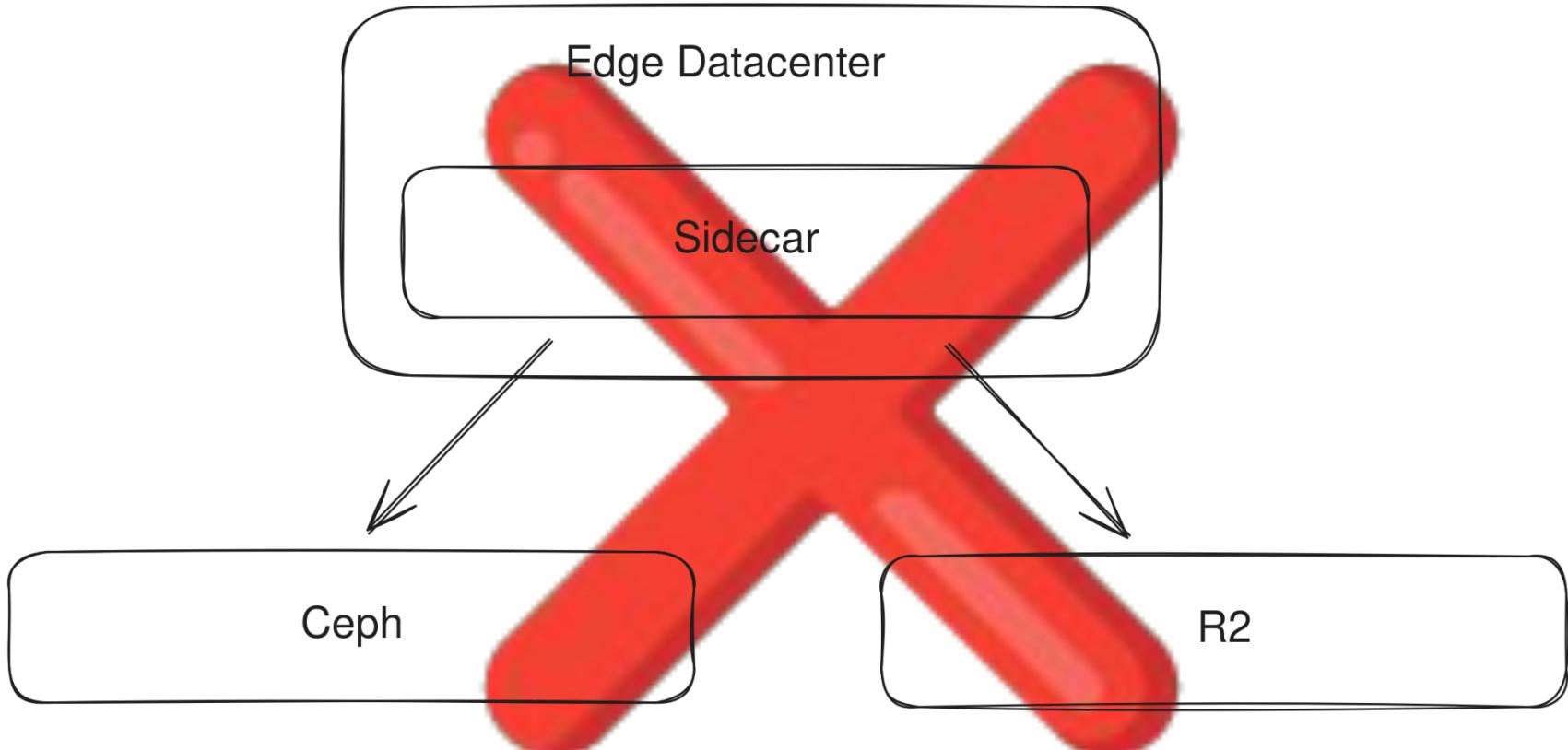
- We didn't have to manage it
- Dogfooding!
- Decentralising our storage



How do we migrate?







```
const (  
    // MetaFilename is the known JSON filename for meta information.  
    MetaFilename = "thanos.shipper.json"
```



Infinite Sadness

✖ Allow customizing the shipper metadata file name

Currently, the shipper metadata file is always called `thanos.shipper.json` in the Prometheus data directory. This precludes running multiple sidecars that upload to different object stores, as they will overwrite each other's metadata file.

This commit allows the metadata file name to be customized via a flag. The default is unchanged, but it can be overridden with the `--shipper.meta-file-name` flag.

As part of this, we update the signatures of `WriteMetaFile` and `ReadMetaFile` to take the full path of the metadata file, rather than just the directory, and updates the tests that go along with this.

Signed-off-by: sinkingpoint <colin@quirl.co.nz>

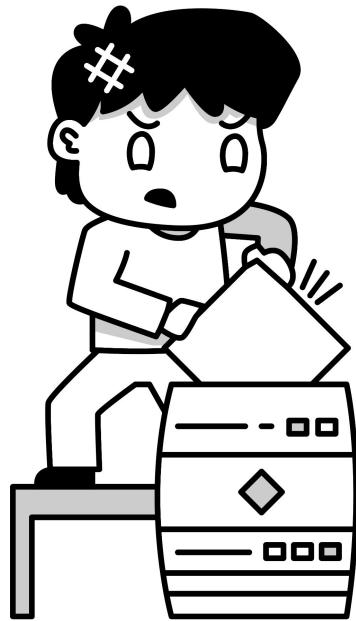
godep main (#6886)

godep v0.34.0 | ... v0.34.0-rc.0

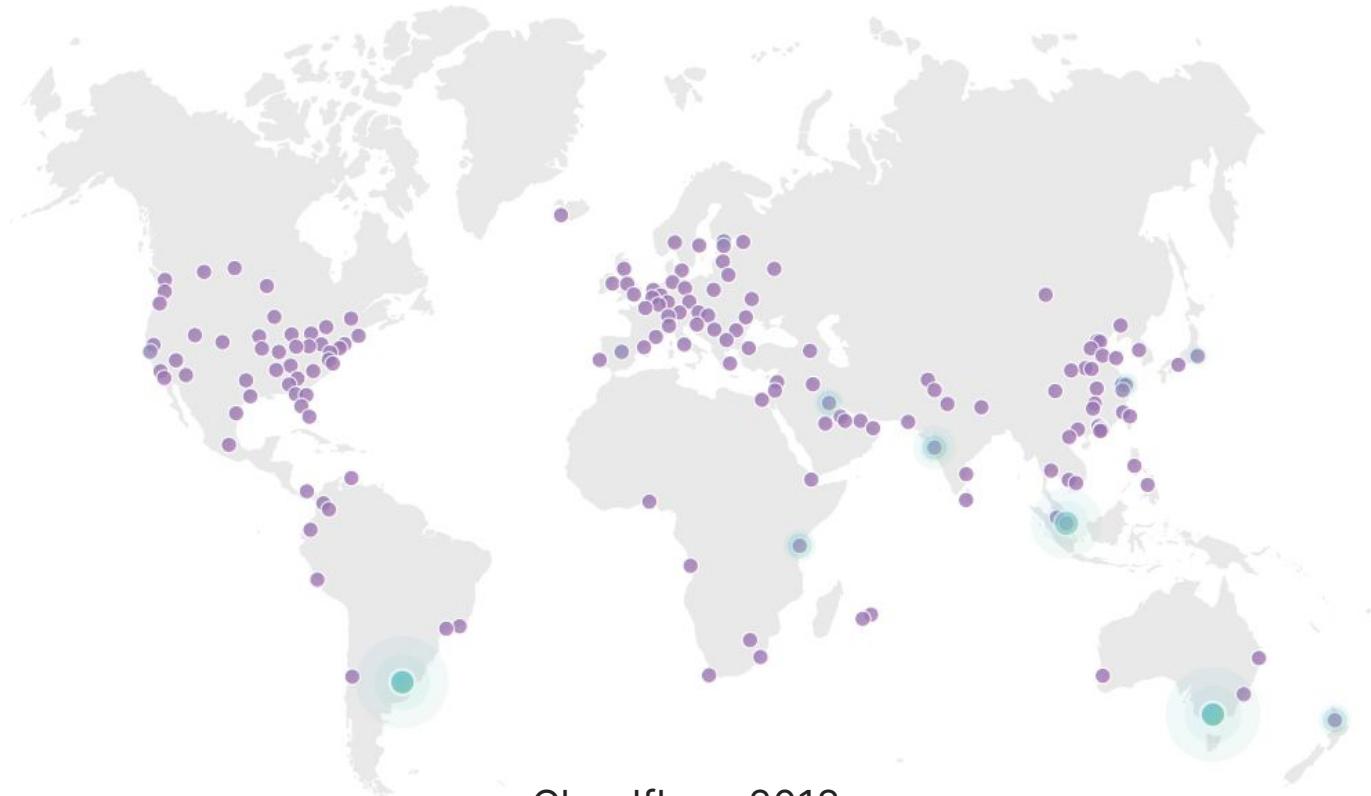
Buckets, buckets everywhere

And far too much data to drink





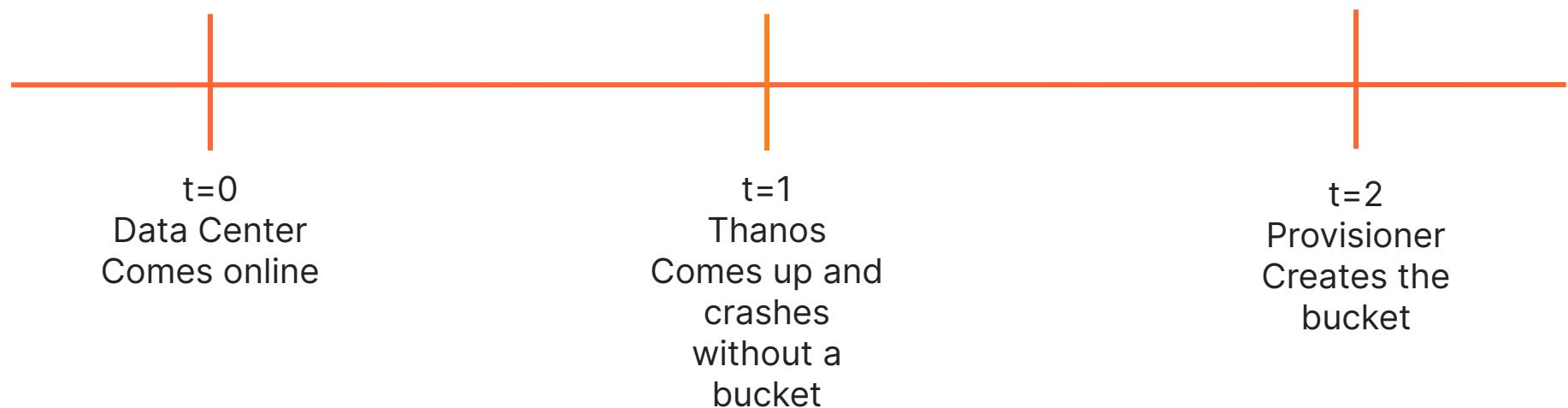
Bucket Granularity **puts a limit** on how much you can scale out because a store has to be able to handle a given time range



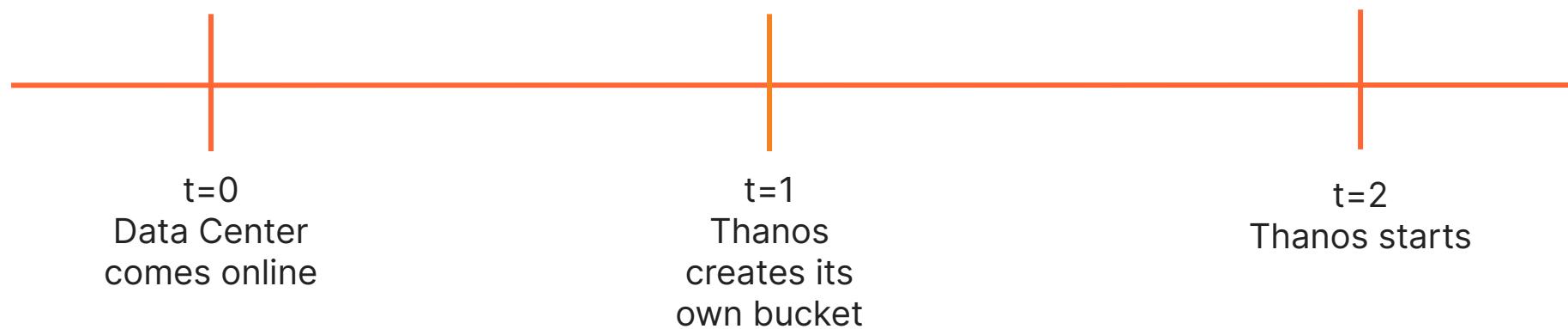
Cloudflare, 2018



Centralised provisioning leads to race conditions

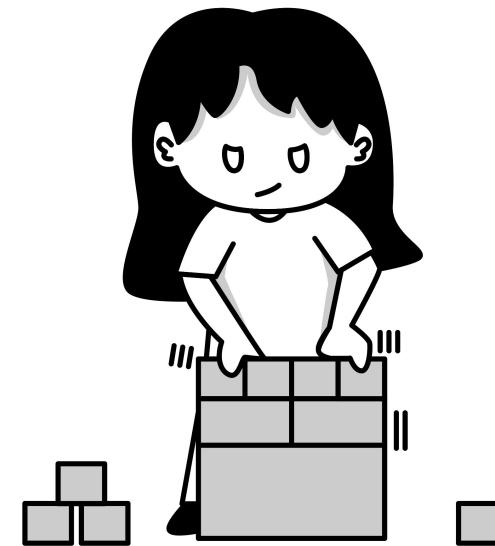


Distributed Provisioning works

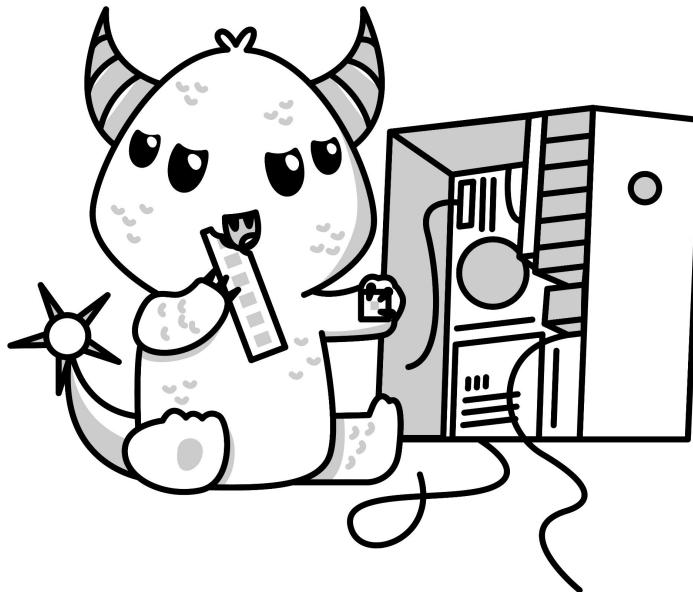


Compactors compacting

Oh no where did all my memory go



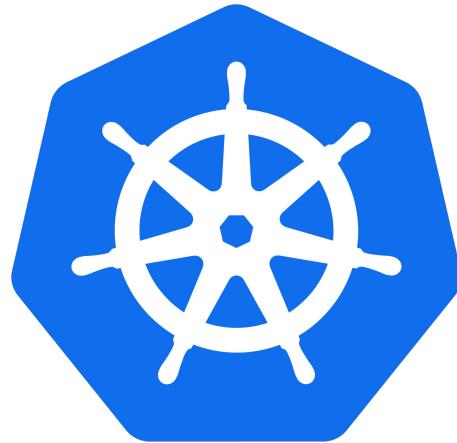
Compactors are:



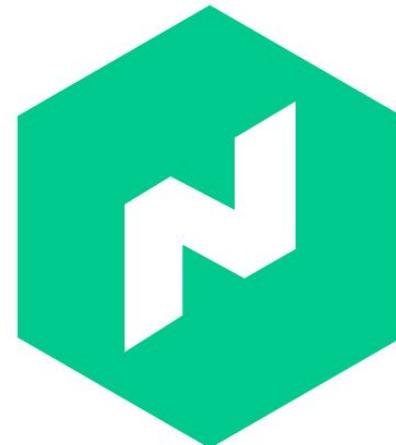
- *Bursty*
- *Heavy*
- *Kind of awkward*

But where to put them?

Centralised in core?



Distributed on the edge?

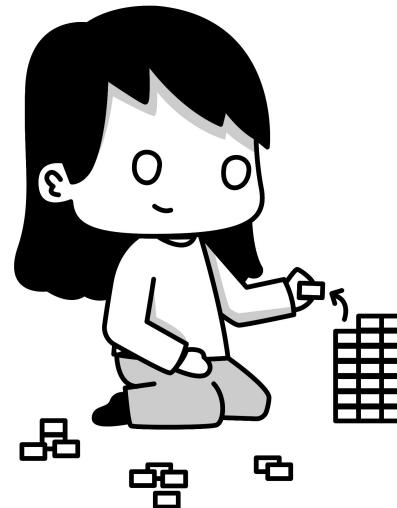


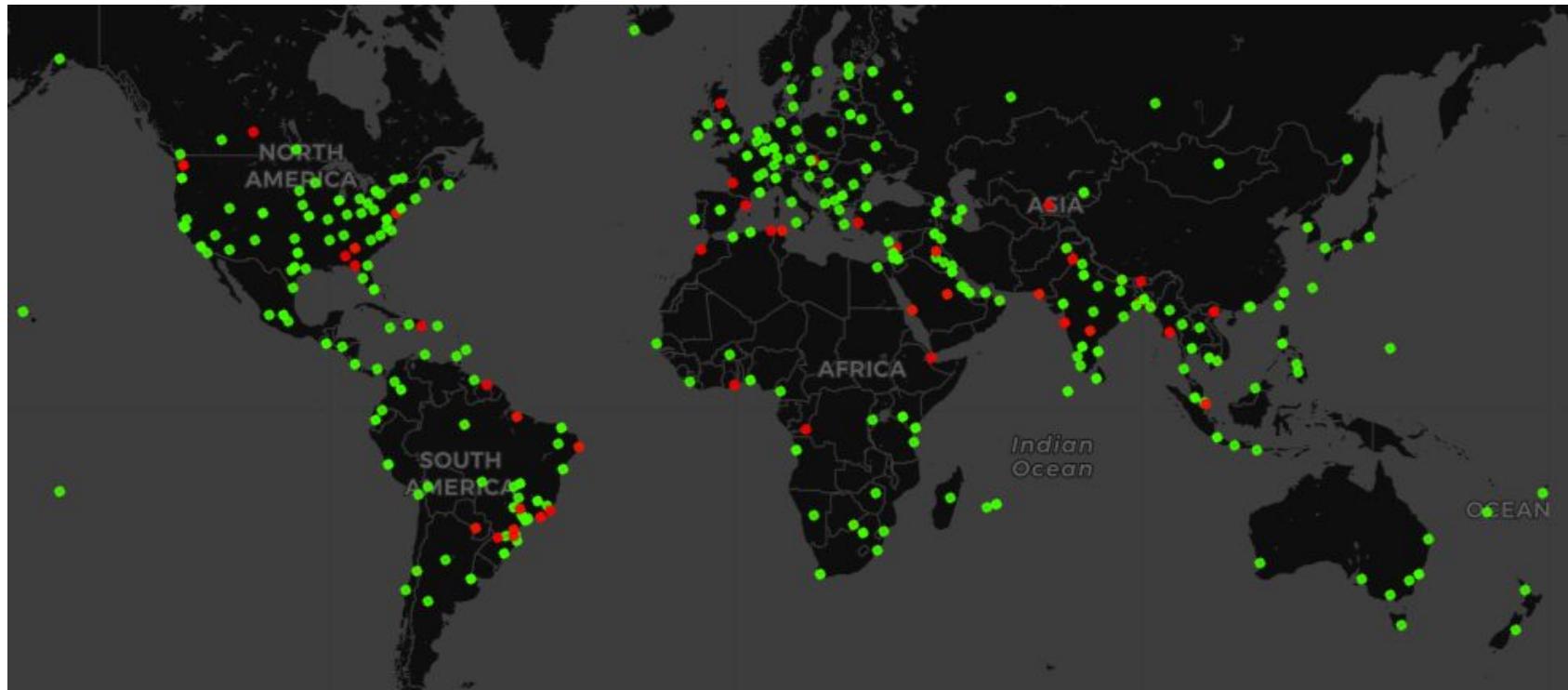


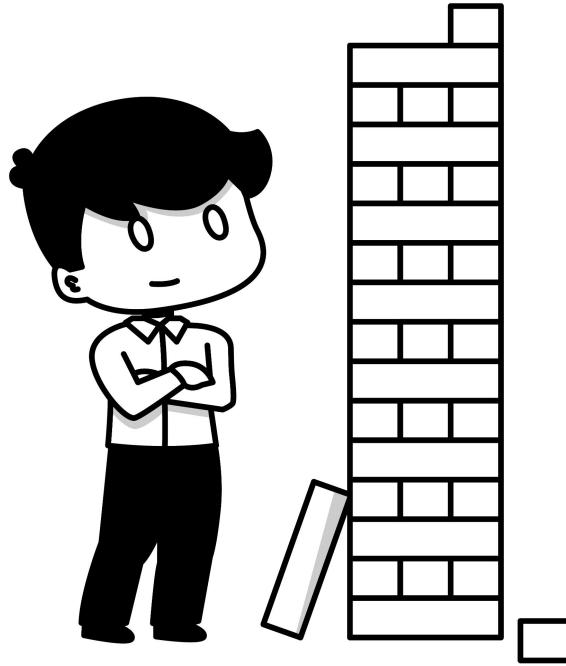
```
job "thanos-compactor" {
    periodic {
        cron = "{{ low_time_start_minute }} {{ low_time_start_hour }} * * * *"
    }
    group "compactor" {
        driver = "docker"
        config {
            image = "docker-registry/thanos:latest"
            command = "timeout"
            args = [
                "{{ low_time_length }}"
                "thanos"
                "compact"
            ] } }
    }
}
```

Stores Storing

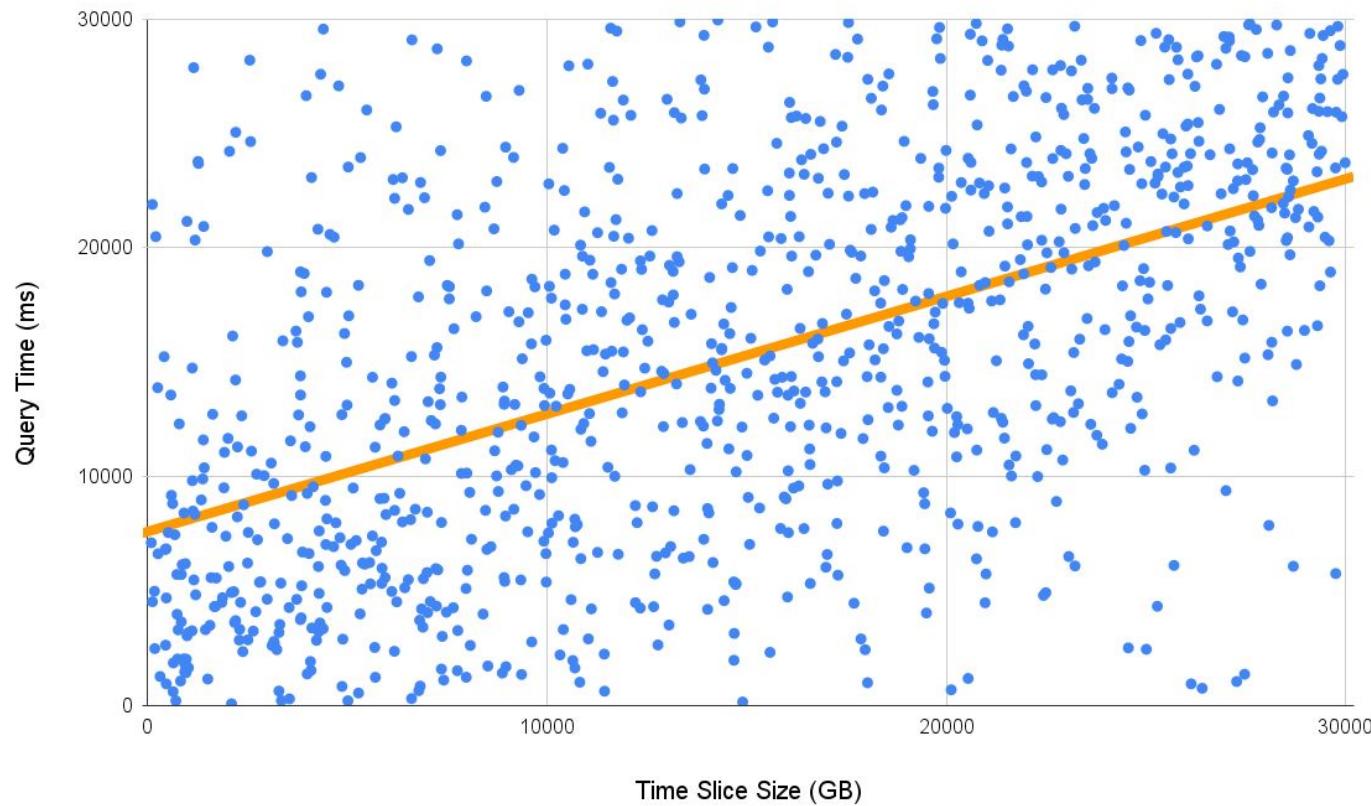
It turns out that querying halfway
around the world is slow

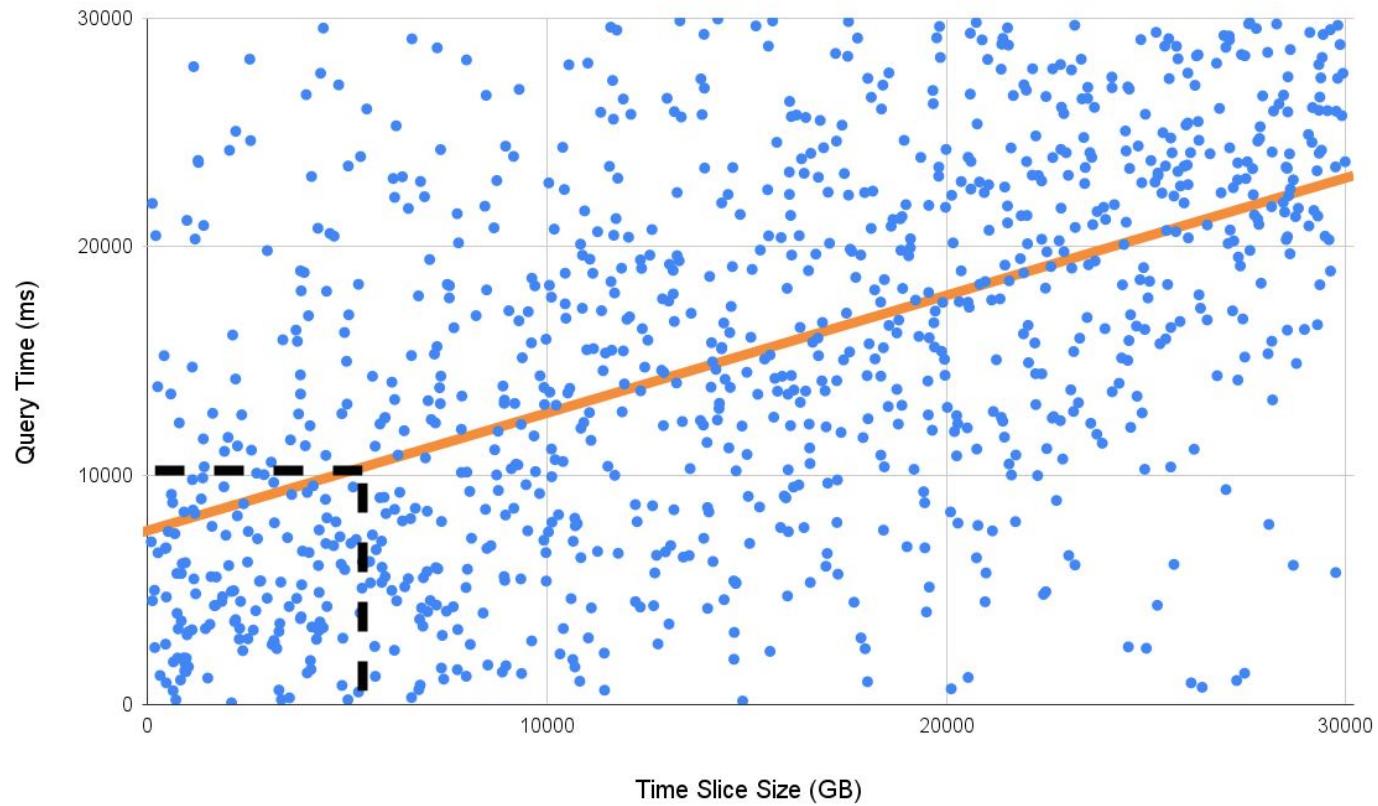


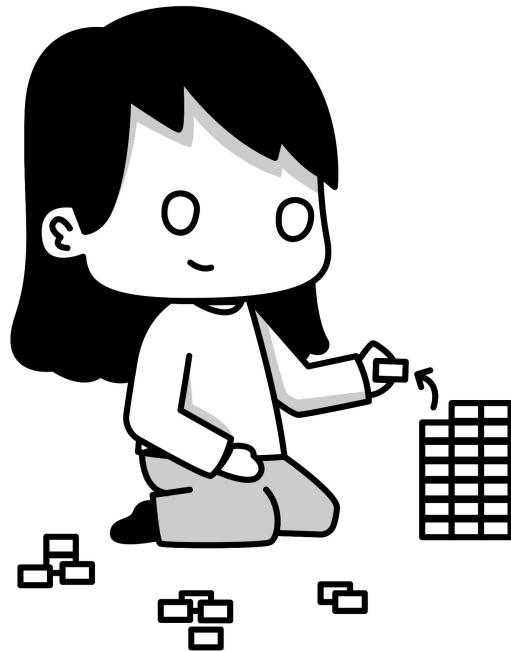
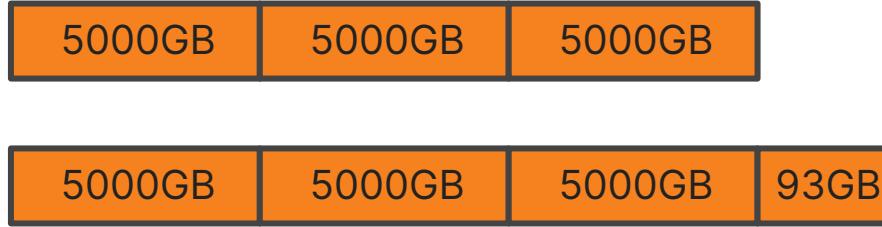


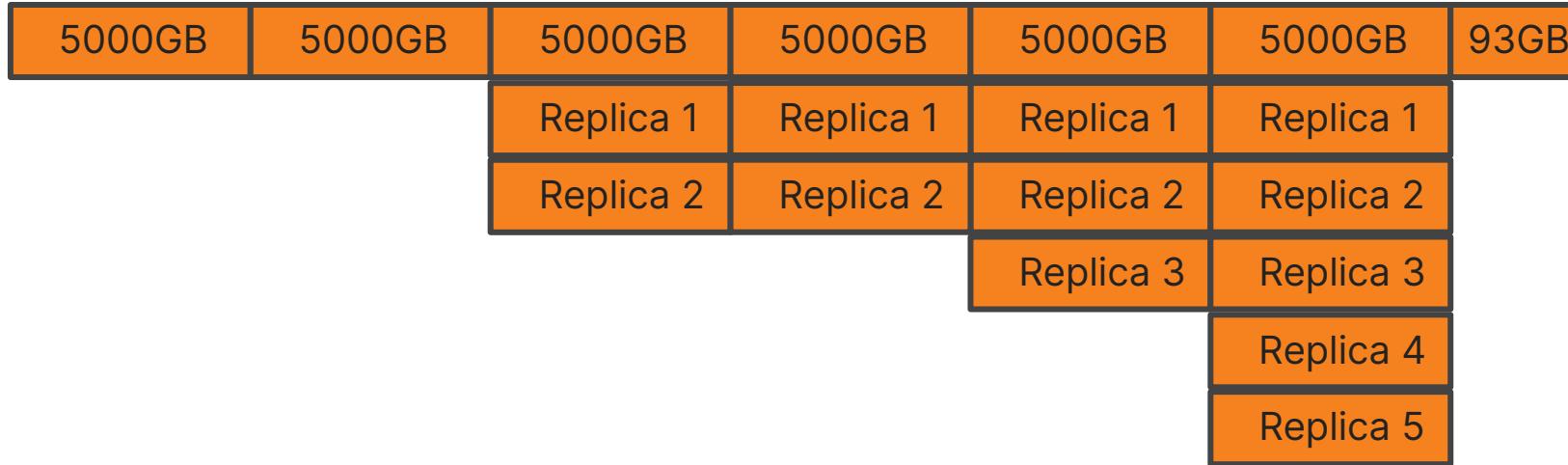


Stores have a **really frustrating** resource profile to scale



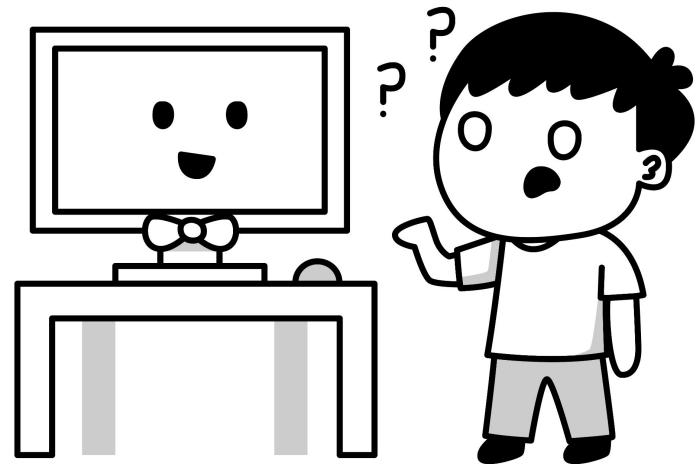


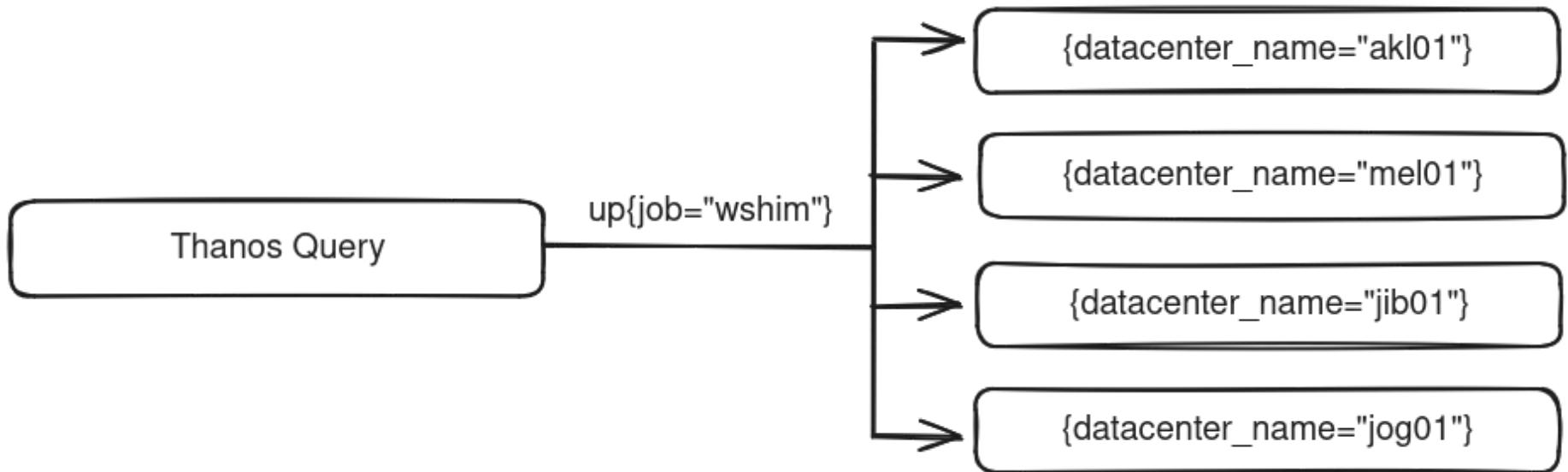


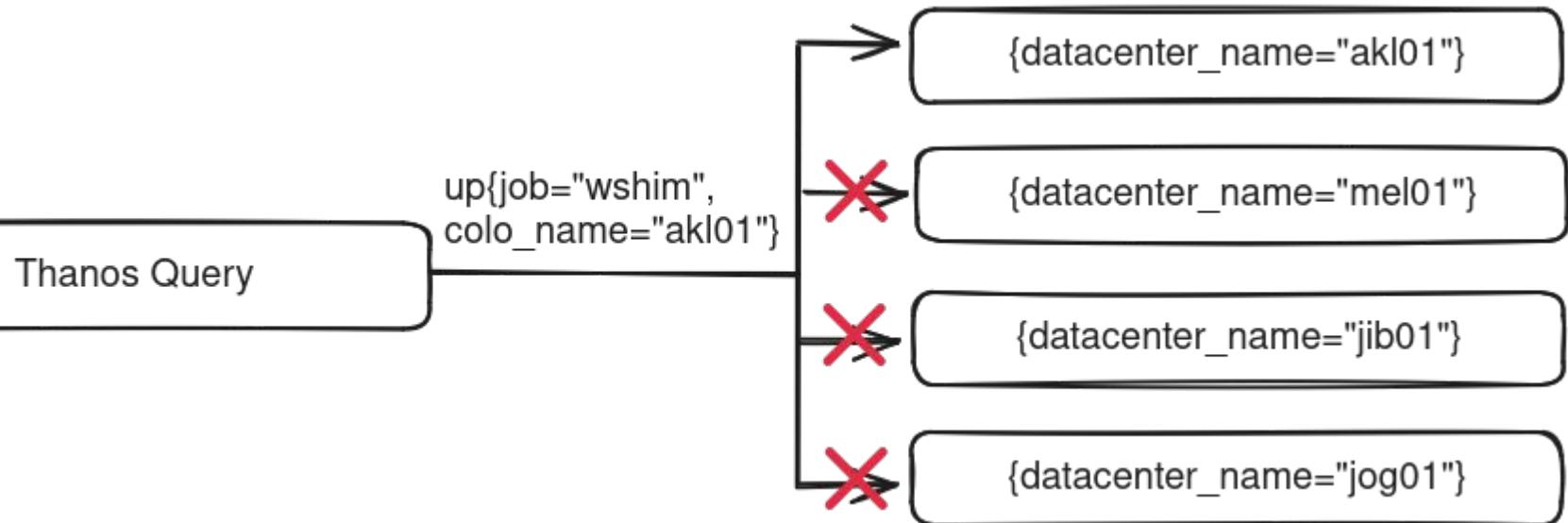


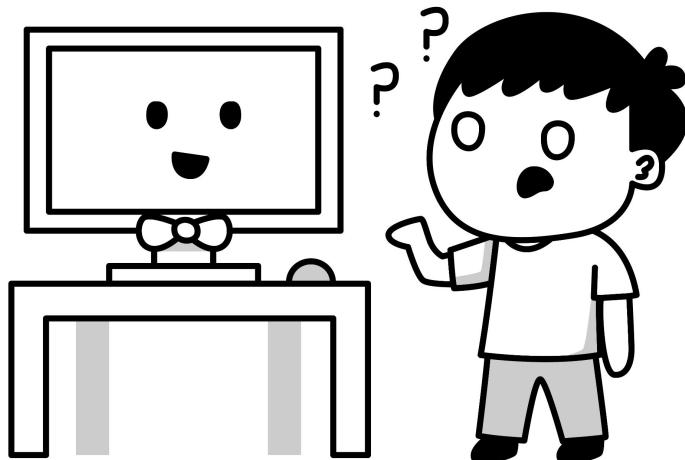
Querys Querying

Or: “Why is my Thanos query always timing out?!?!”









The **Label Enforcer** forces teams to use external label filters, and guides them towards fast queries

Bad: my_metric

Better: my_metric{datacenter_name="akl01"}

Bad: my_metric{datacenter_name="akl01"} + my_second_metric

Better: my_metric{datacenter_name="akl01"} +
my_second_metric{datacenter_name="akl01"}

Also OK: my_metric{datacenter_name=~"akl01|mel01"}

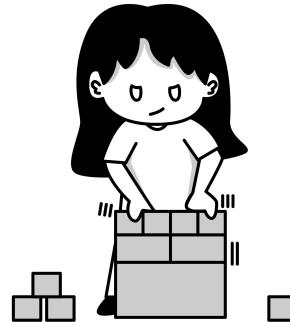
Also OK: my_metric{datacenter_name=~".+"}

That's about it

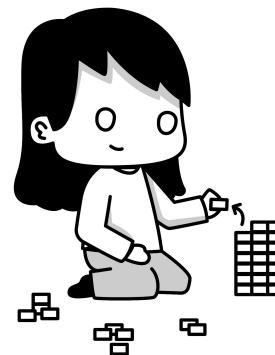
Automated bucket provisioning



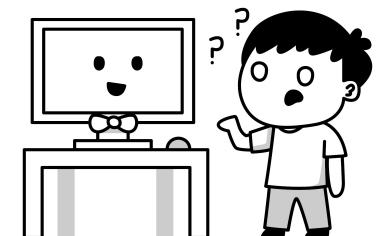
Compactors running with spare capacity



Stores that get broken up to meet our SLOs



Guiding engineers to fast Queries by default



Thanks!



Art by Alicia Edwards ^



Slides ^