# Project Frankenstein

A multi-tenant, horizontally scalable Prometheus as a Service

Tom Wilkie (& Julius Volz)
Weaveworks, August 2016

**weave**works

# FRANKENSTEIN;

OR,

## THE MODERN PROMETHEUS.

---

By MARY W. SHELLY,

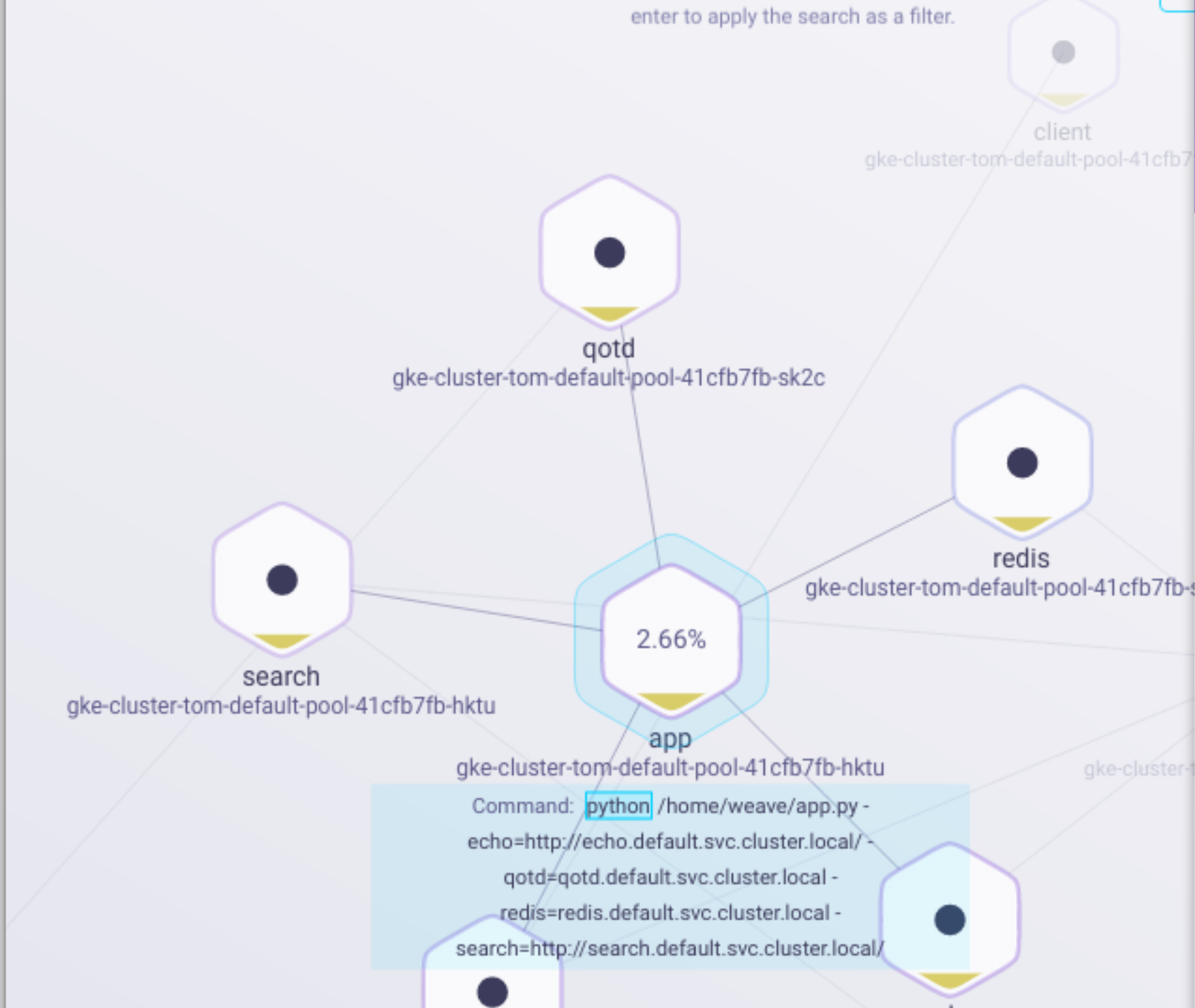AUTHOR OF 'THE LAST MAN,' 'PERKIN WARBECK,' &c.

*"the best way to visualise, manage & monitor your cloud native application"*

weavescope

Q python

Try "app", "id:b31e2a1a8c41", or "cpu > 2%". Hit
enter to apply the search as a filter.

client
gke-cluster-tom-default-pool-41cfb7...

qotd
gke-cluster-tom-default-pool-41cfb7fb-sk2c

redis
gke-cluster-tom-default-pool-41cfb7fb-...

search
gke-cluster-tom-default-pool-41cfb7fb-hktu

2.66%

app
gke-cluster-tom-default-pool-41cfb7fb-hktu

Command: python /home/weave/app.py -
echo=http://echo.default.svc.cluster.local/ -
qotd=qotd.default.svc.cluster.local -
redis=redis.default.svc.cluster.local -
search=http://search.default.svc.cluster.local/

echo
gke-cluster-tom-default-pool-41cfb7fb-hktu

frontend

Command: python /home/weave/echo.py

**localhost**:4040/#!/state/{"controlPipe":null,"nodeDetails":[{"id":"b31e2a1a8c417b34f0249b92a08b582d6d59ca52f15c36166...

Weave Scope - app

**app**

tomwilkie/app  gke-cluster-tom-default-p...
app-3461397960-iqtfe

STATUS

2.66 %                    32.6 MB

CPU                       MEMORY
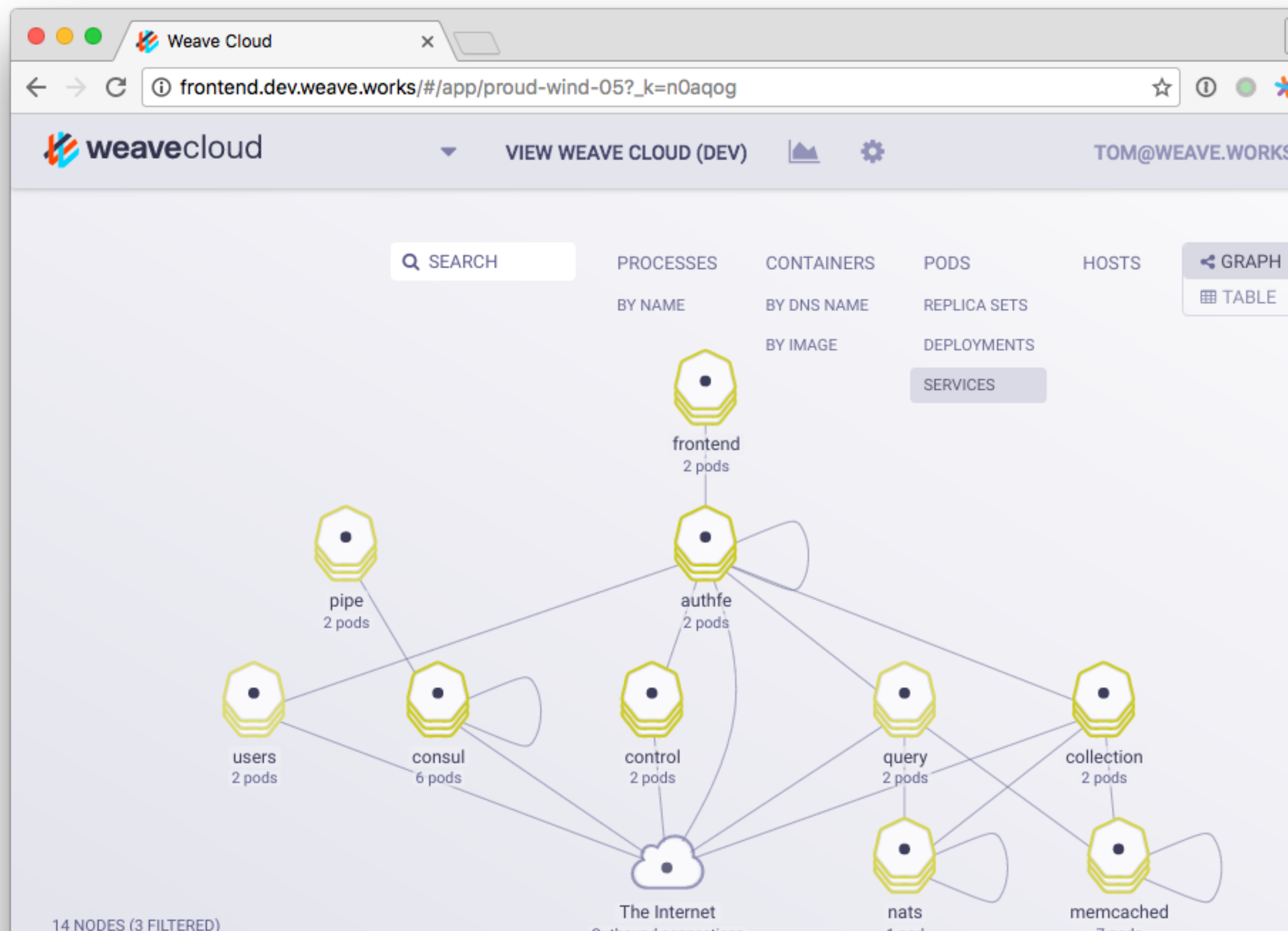
INFO

ID: b31e2a1a8c41
STATE: Up 10 minutes
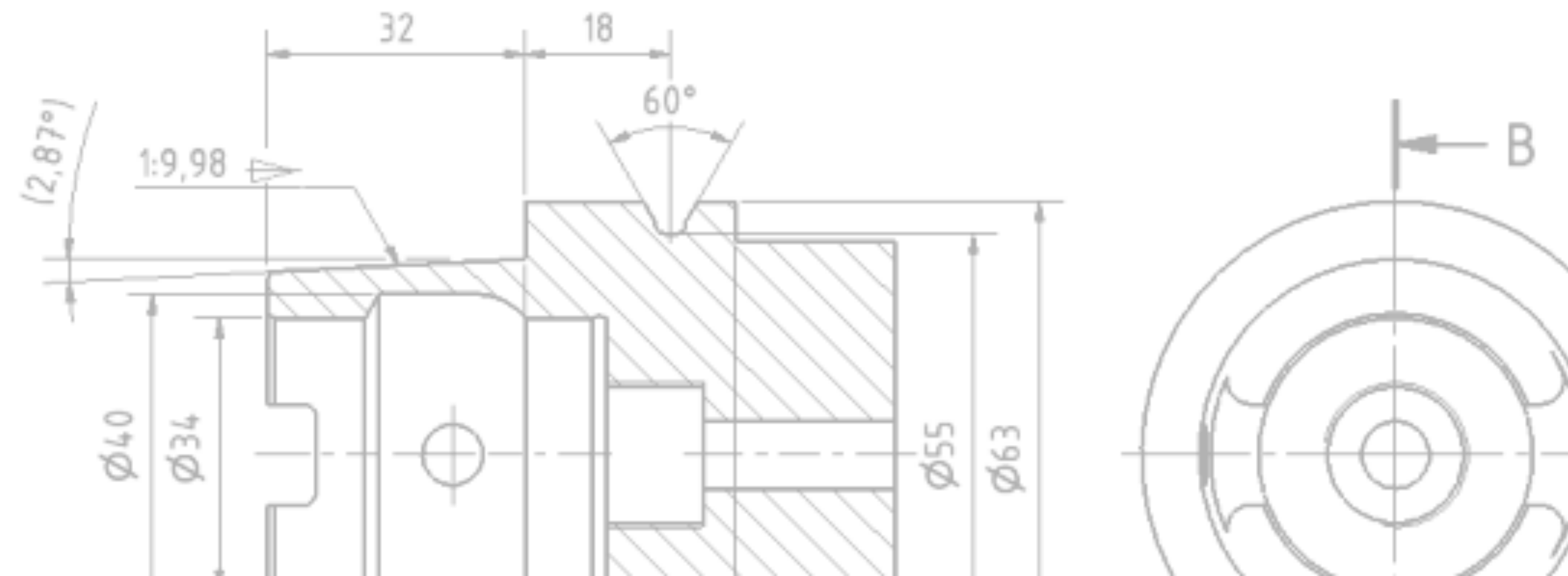COMMAND: python /home/weave/app.py -echo=ht...

+6

| INBOUND | PORT | ▼ COUNT |
|---|---|---|
| frontend | 80 | 37 |

| OUTBOUND | PORT | ▼ COUNT |
|---|---|---|
| qotd | 4446 | 53 |
| echo | 80 | 36 |
| search | 80 | 20 |
| redis | 6379 | 2 |

weaveworks

# Design

why not just run my own Prometheus?

- the as-a-service bit provides authentication and access control

- virtually infinite retention; all the state is managed for you, by us

- provide a different story around durability, HA and scalability

- (eventually) better query performance, especially for long queries
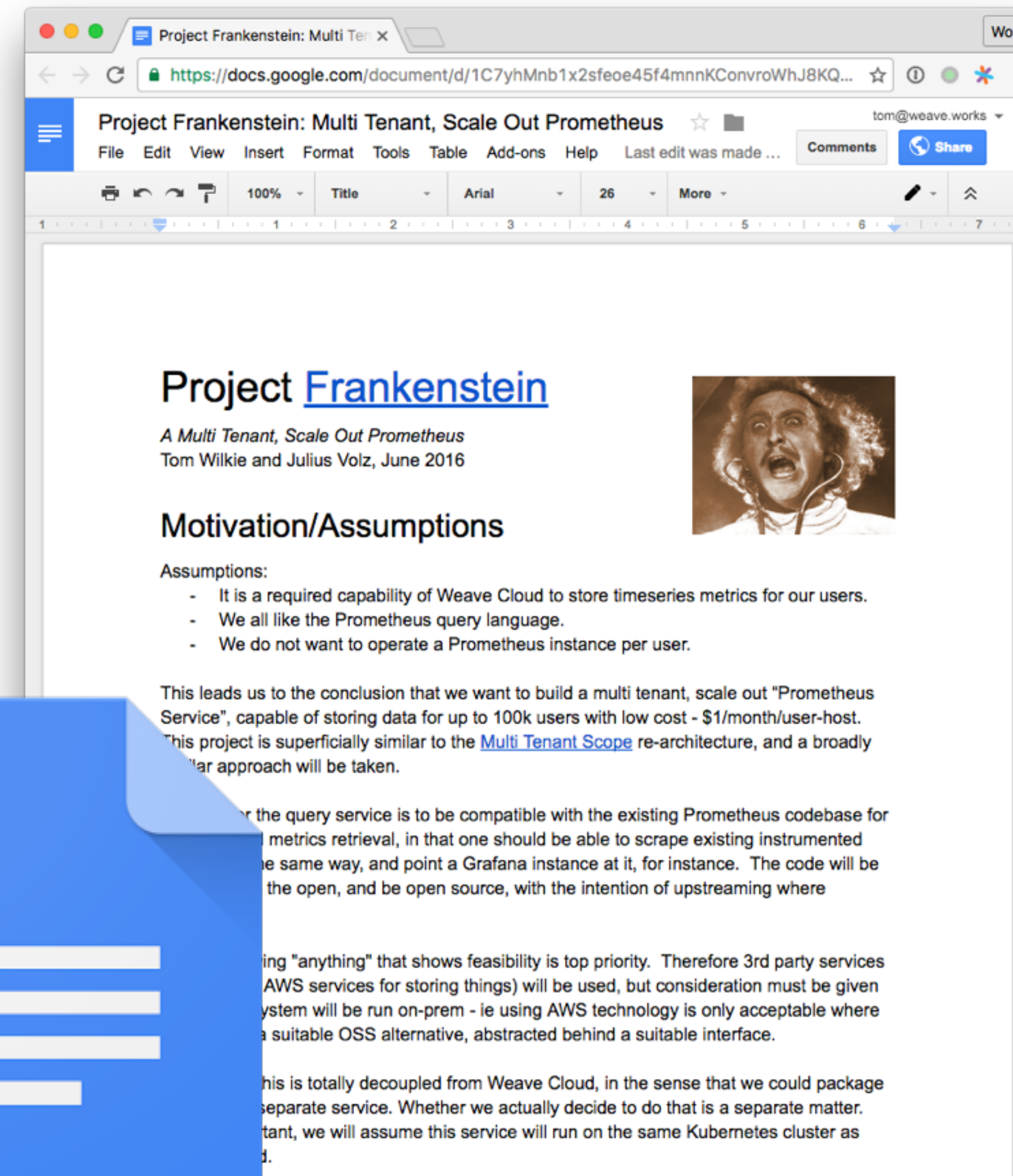
**weave**works

requirements:

1. API compatible with Prometheus

2. easy to operate and manage

3. tens of thousands of users, tens of millions samples/s

4. cost effective to run

5. reuse as much of Prometheus as possible

… so we can sell it

Aim: build proof of concept as quickly as possible

16/06    started design doc

22/06    circulated on list

22/06    initial commit

26/07    launch jobs

25/08    give talk!



weaveworks

**Your DC**

Retriever

scraping

your jobs

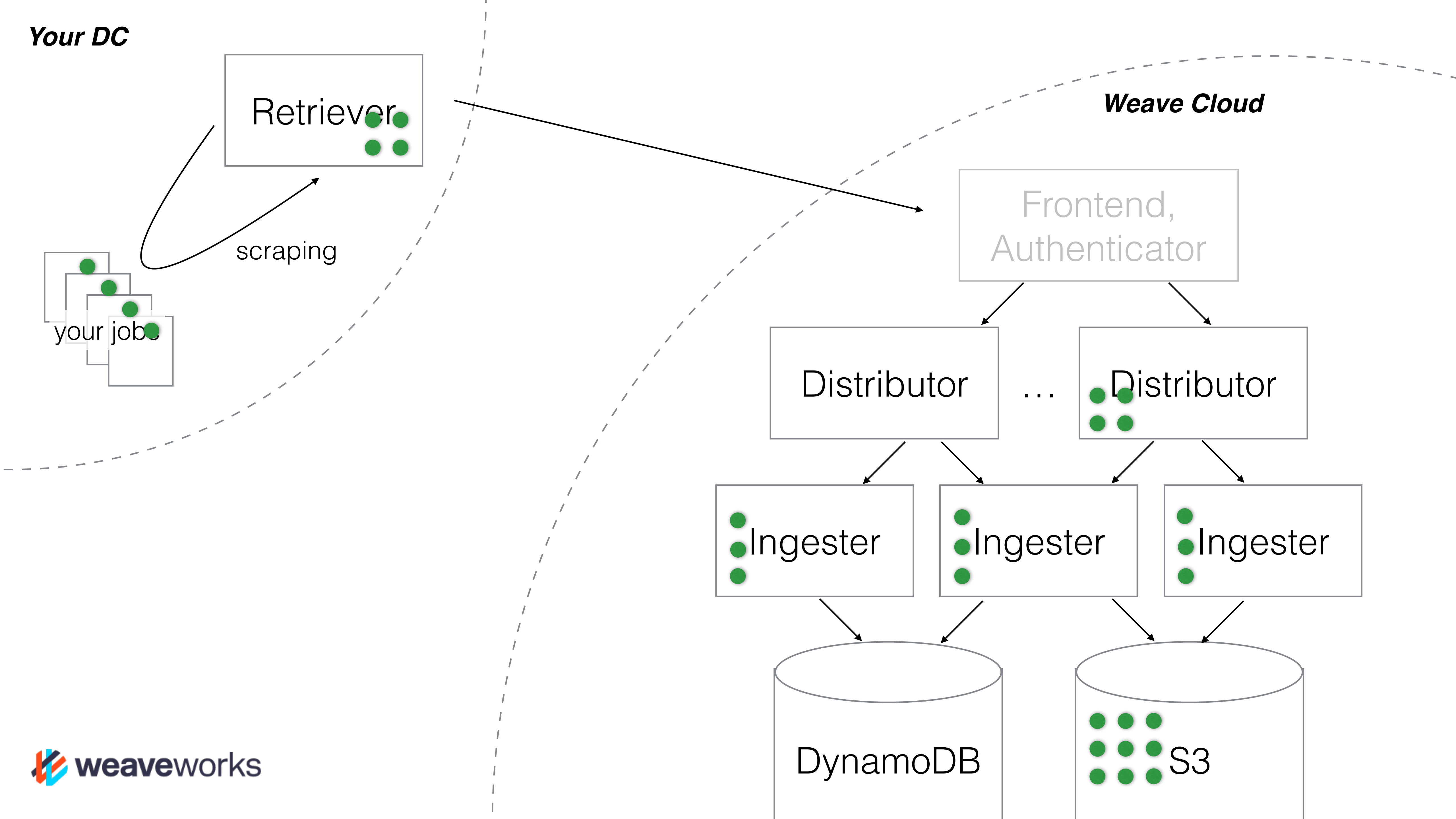**Weave Cloud**

Frontend, Authenticator

Distributor ... Distributor

Ingester    Ingester    Ingester

DynamoDB    S3

weaveworks

```
┌─────────────────────┐
│                     │
│     Retriever       │
│                     │
└─────────────────────┘
```

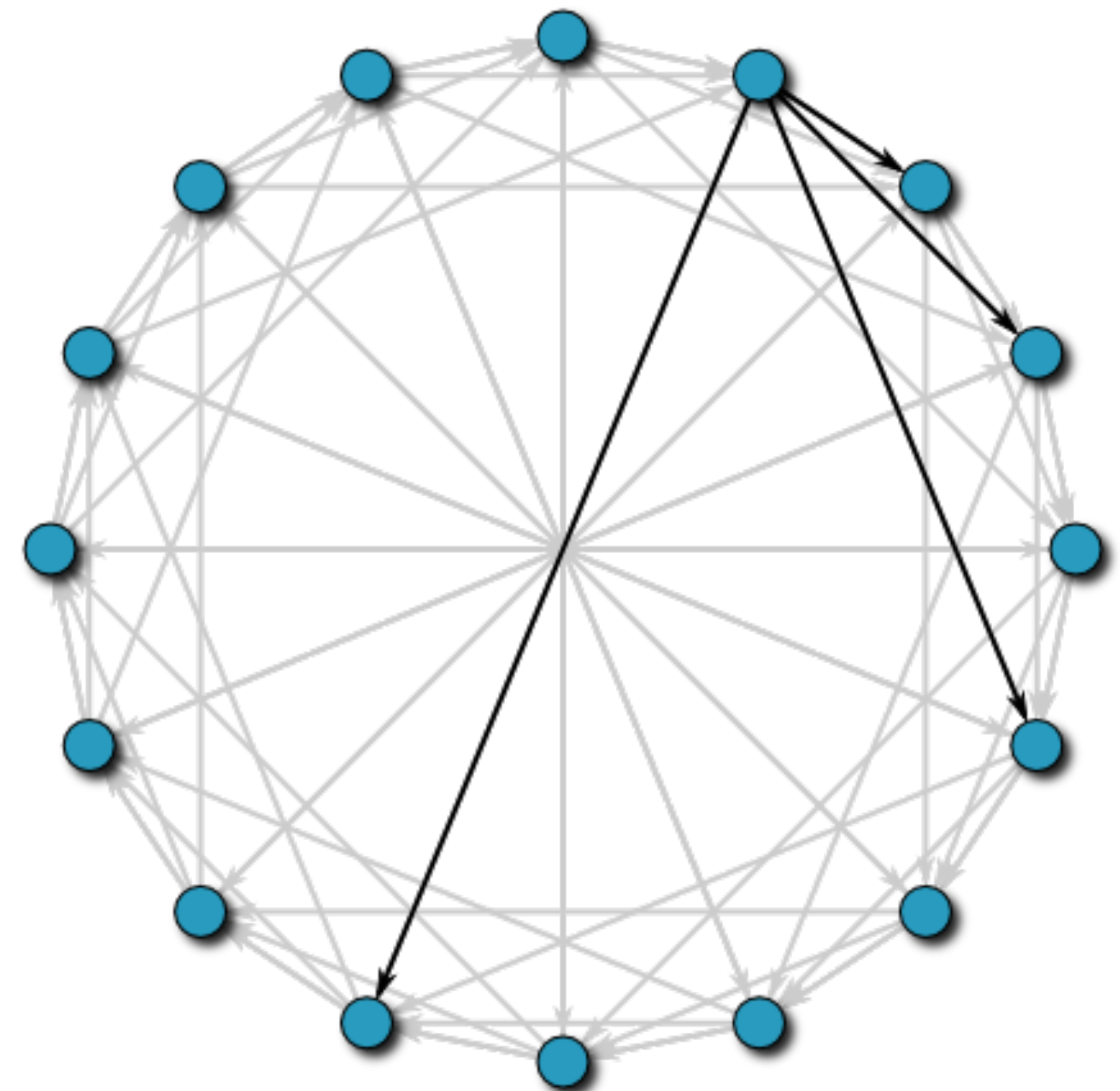Does scraping and relabelling.

Is a vanilla Prometheus plus:

- Brian Brazil's generic write PR ([#1487](#))

- Some modification to prevent local storage + indexing

```
/bin/prometheus -retrieval-only -storage.remote.generic-url=...
```

**weave**works

# Distributor

- Uses consistent hashing to assign timeseries to Ingesters

- Input to hash is (user ID, metric name)

- Tokens stored in Consul
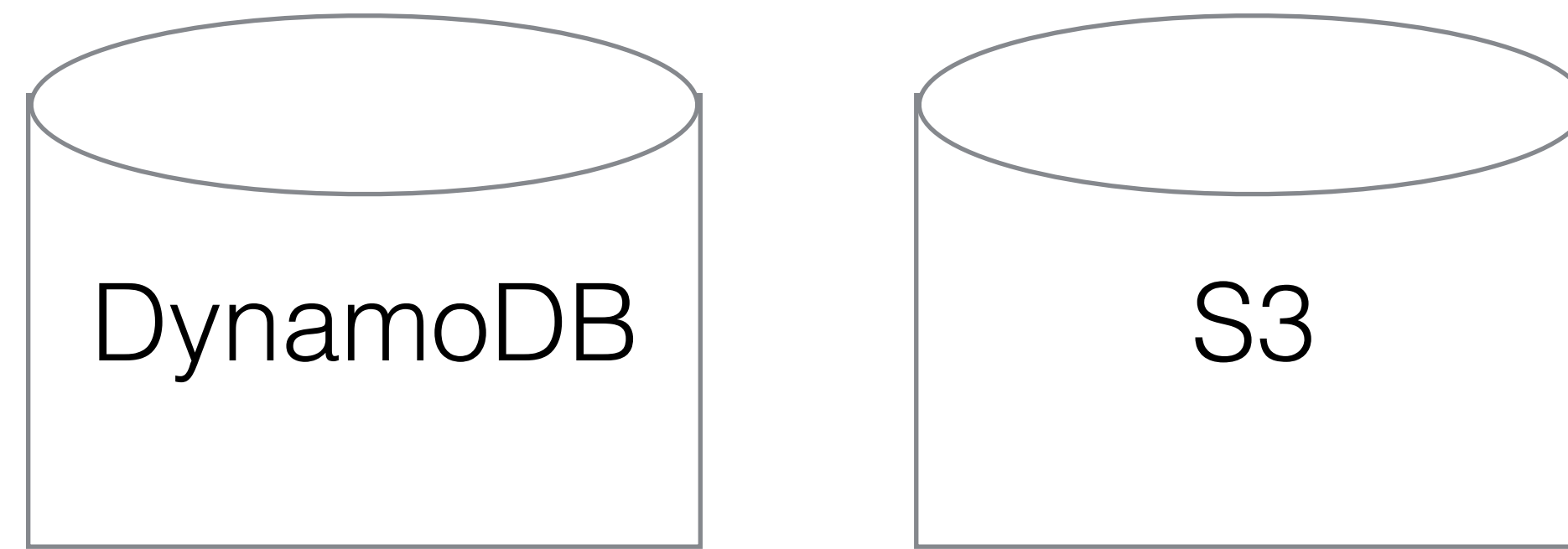
- Also currently handles queries



http://goo.gl/U9u1U2

weaveworks

DynamoDB          S3

External inverted index maintained in DynamoDB, chunks stored in S3

Item in DynamoDB looks like:

```
{
    hash key: "{user ID}:{metric name}:{hour}",
    range key: "{label name}:{label value}:{chunk ID}",
    metric: ...,
    from, through: ...,
    ID: ...,
}
```

weaveworks

# Evaluation

## The Good

- It works!  And in ~2 months.

- Seems pretty scalable, handling two clusters right now

- Query performance better than expected

## The Bad

- Hashing scheme means can't do queries that don't involve metric names.

- Possible to hotspot an ingester

**The Ugly**: the code…

weaveworks

# Demo

weaveworks

# Lots left to do…

**Features:**

- Recording rules

- Alerting & Alertmanager

**Reliability:**

- Replication between ingesters, commit log etc

- Ingestor lifecycle

- Separate query service?

**Performance:**

- Query parallelisation

- Background chunk coalescing

**Code:**

- Code cleanup

- Upstream appropriate changes

**weave**works

# Questions?

https://github.com/tomwilkie/prometheus

## Try it out!

Email help@weave.works for
instructions and to get on white list

**weave**works