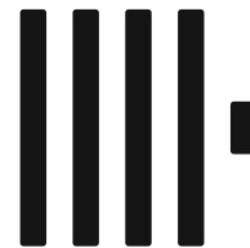
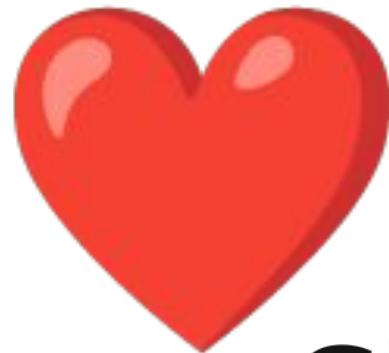


# Experiments in Backing Prometheus with Clickhouse

Colin Douch, Observability @ Cloudflare

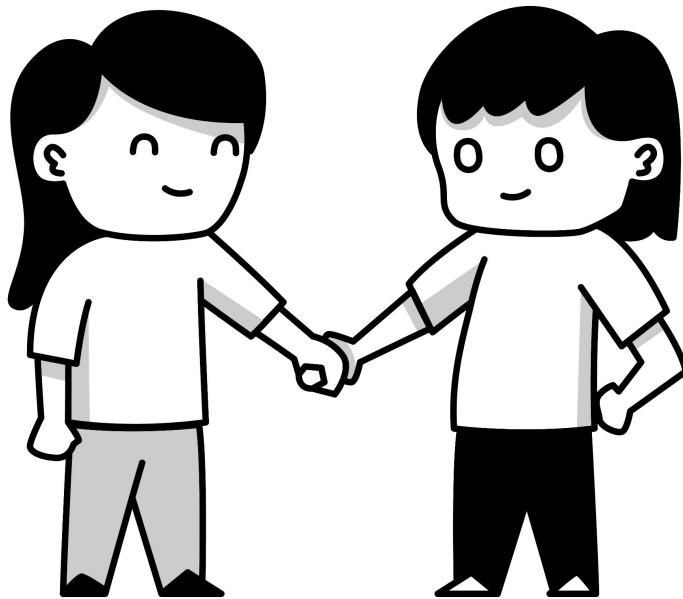




# ClickHouse

| Host     | Metric Name | node_status |
|----------|-------------|-------------|
| 46m34    | wshim_total | V           |
| 91ssds3  | errors      | V           |
| 36com360 | cats_count  | V           |
| 46dm3    | cpu_seconds | I           |
| 830m19   | total       | P           |





What if we  
could have  
***both***

# A Quick Introduction

Hi! I'm Colin

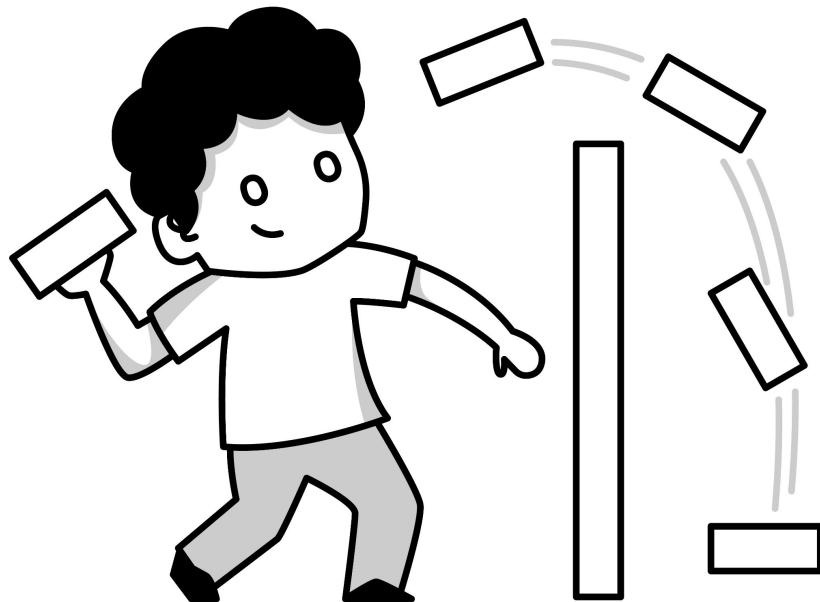
Observability Lead @ Cloudflare

Has bad opinions

Not British



But how?



**Remote write**  
Is where the  
magic happens

It turns out that  
implementing  
remote write it  
quite trivial



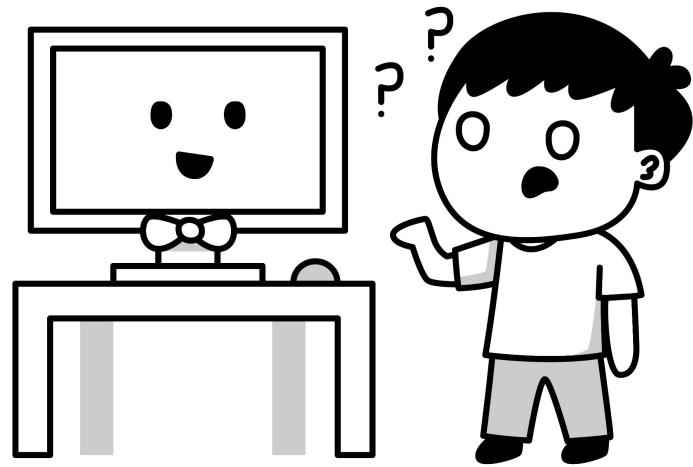
<https://github.com/sinkingpoint/prometheus-clickhouse-bridge>

For queries, we  
need to implement  
*remote read*



<https://github.com/sinkingpoint/prometheus-clickhouse-bridge>

# Welcome to Experiments in Backing Prometheus with Clickhouse

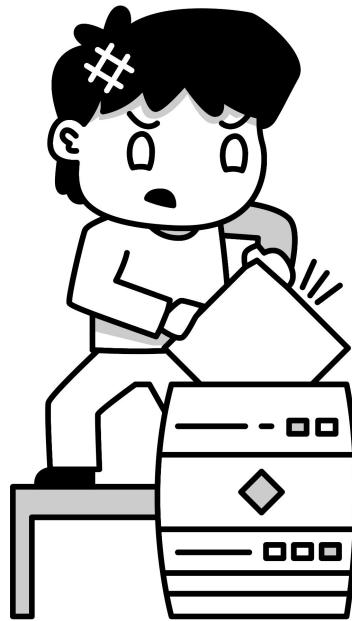


# Testing Environment

Let's be good scientists



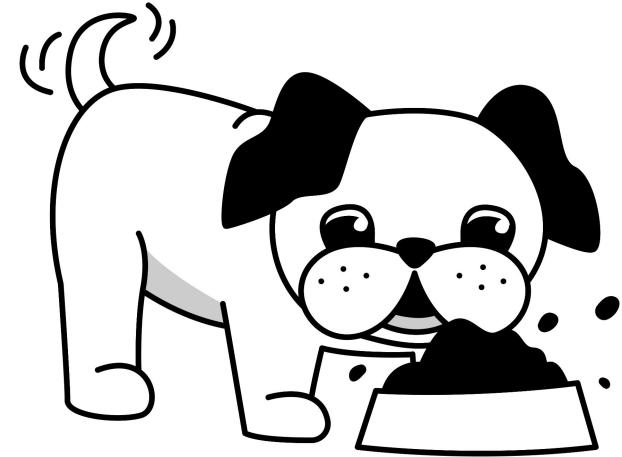
## Hardware



- 92 Cores x 1.9GHz
- 2.4TB NVME Drive
- 256 GB of RAM

## Timeseries Data

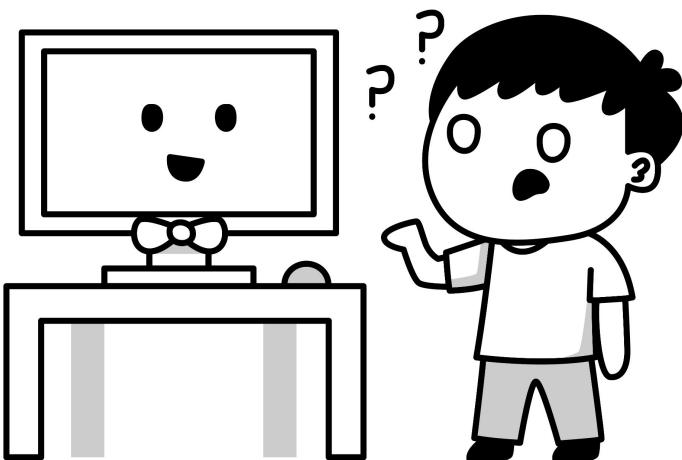
24 Hours of production metrics, across 672 individual services with a wide variety of usage patterns



But... These are  
already **completed**  
**blocks**...



<https://github.com/sinkingpoint/replay>



# 7 days worth of production queries

## Query Buckets

- Small - 1134 Queries
- Medium - 2475 Queries
- Large - 475 Queries



## Simple Queries:

```
rate(errors[2m])
```

## Medium Queries:

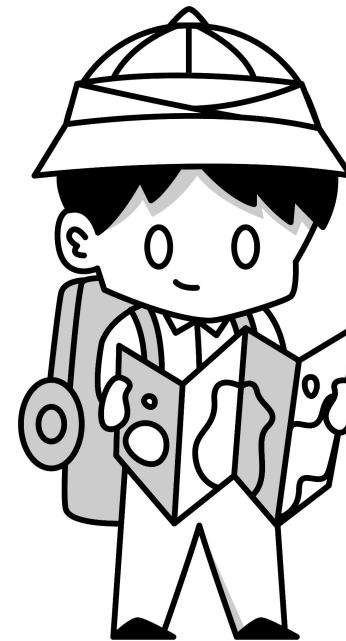
```
rate(errors[2m]) /  
rate(total[2m])
```

```
(-(((avg_over_time(snmp{direction="in",snmptype="transit"}[24h])
  or
  avg_over_time(snmp{direction="in",snmptype="transit"}[24h]))
  and on (name)
  metadata{environment="prod",is_mcp="1"})
 - on (description, name) (
  (avg_over_time(snmp{direction="in",snmptype="transit"}[24h])
  or
  avg_over_time(snmp{direction="in",snmptype="transit"}[24h]))
  and on (name)
  metadata{environment="prod",is_mcp="1"}))) *
((((
  avg_over_time(snmp{direction="in",snmptype="transit"}[24h])
  or
  avg_over_time(snmp{direction="in",snmptype="transit"}[24h]))
  and on (name) metadata{environment="prod",is_mcp="1"})
 - on (description, name) (
  avg_over_time(snmp{direction="in",snmptype="transit"}[24h])
  or
  avg_over_time(snmp{direction="in",snmptype="transit"}[24h]))
  and
  on (name) metadata{environment="prod"})))) ^ 0.5
```

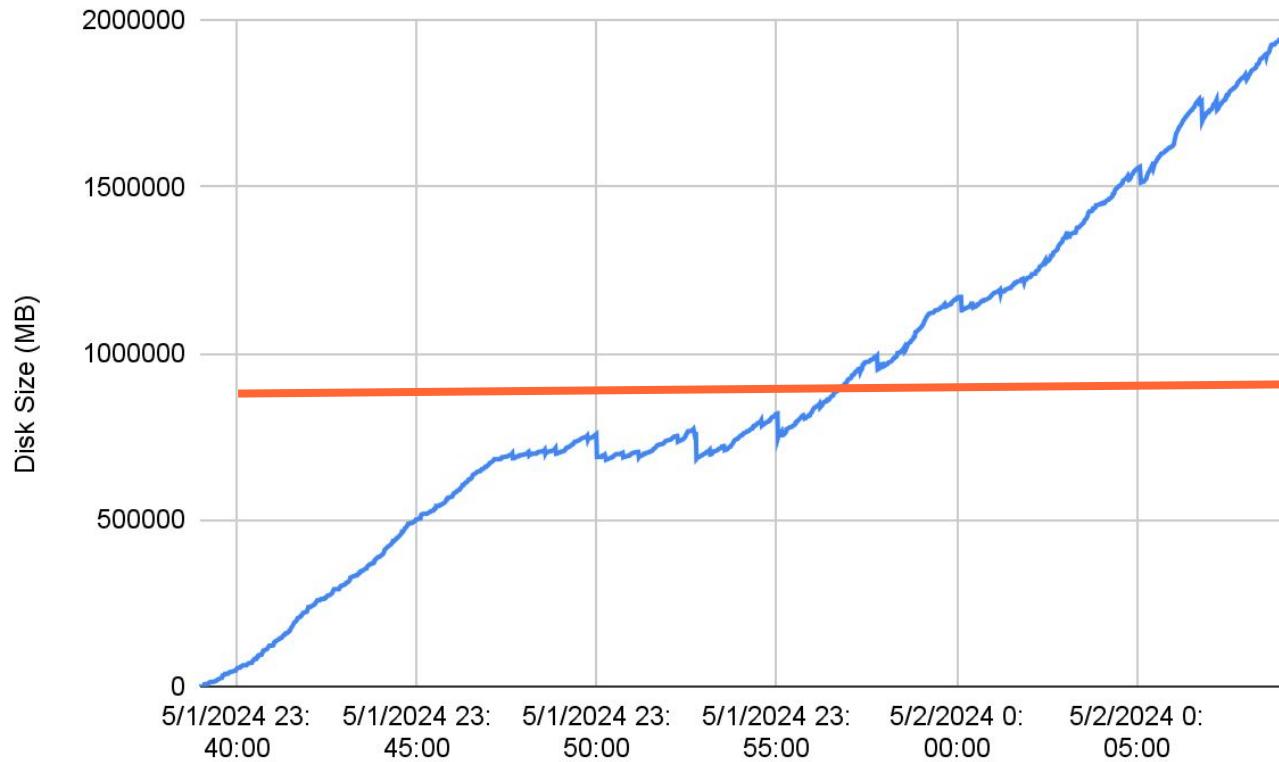


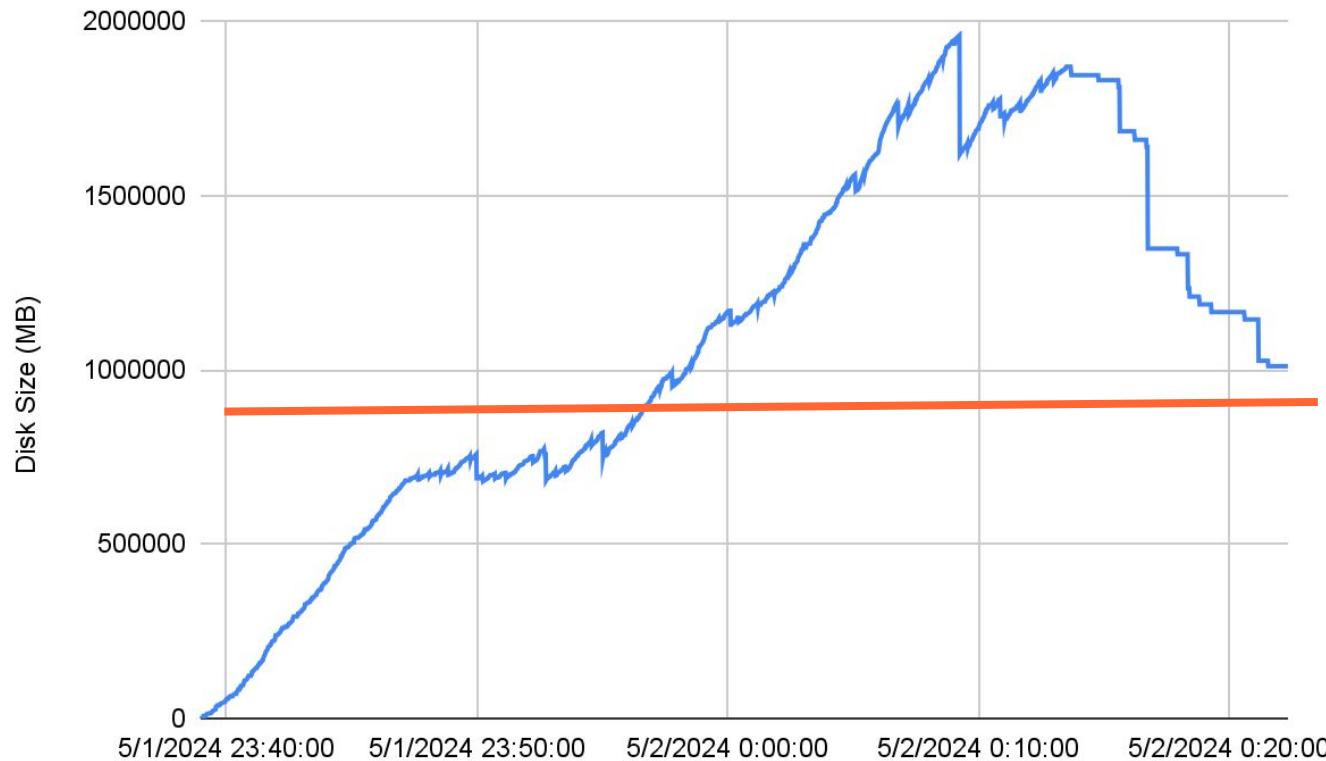
# Intrepid first steps

Let's see if we can pull some numbers

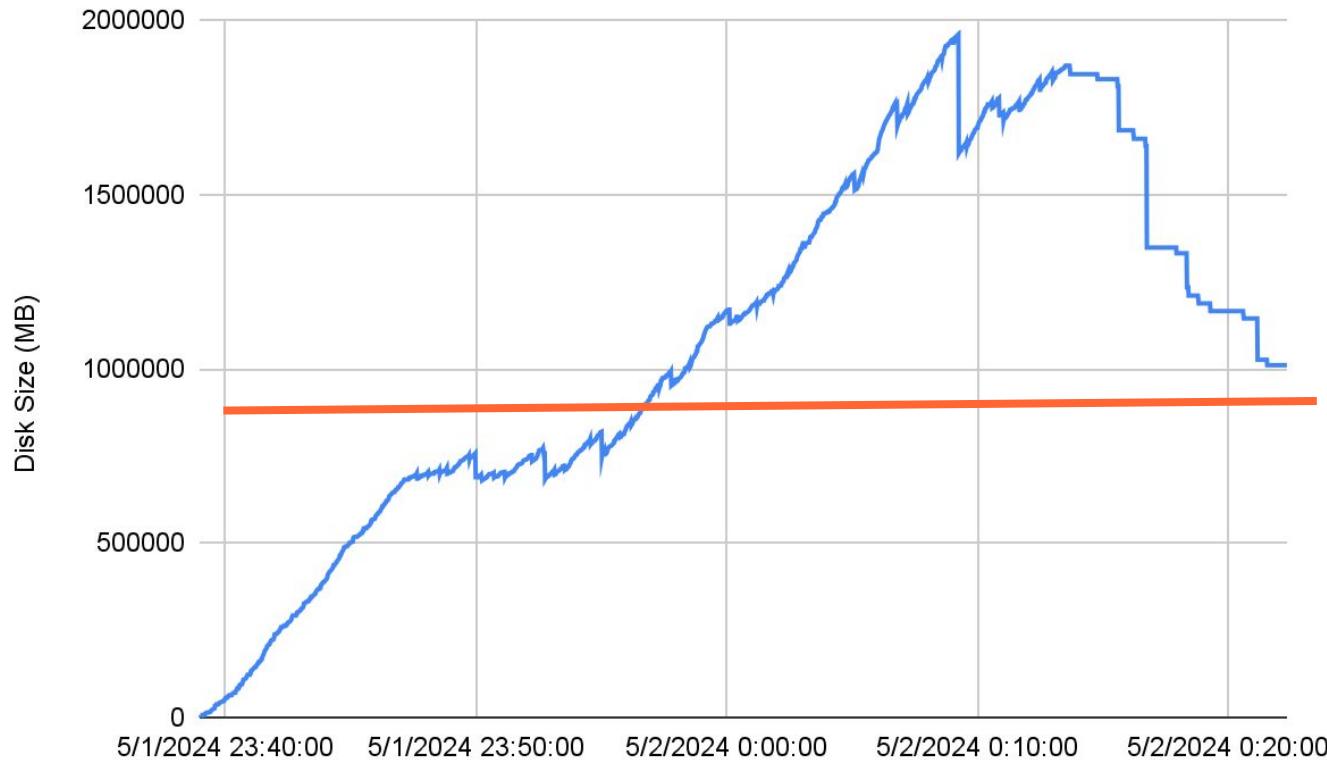


```
CREATE TABLE IF NOT EXISTS metrics (
    timestamp DateTime CODEC(DoubleDelta, ZSTD),
    name String CODEC(ZSTD),
    labels Map(String, String) CODEC(ZSTD),
    value Float64 CODEC (DoubleDelta, ZSTD),
) ENGINE = MergeTree() PRIMARY KEY (name, labels, timestamp);
```





```
CREATE TABLE IF NOT EXISTS metrics (
    timestamp DateTime CODEC(DoubleDelta, ZSTD),
    name String CODEC(ZSTD),
    labels Map(String, String) CODEC(ZSTD),
    value Float64 CODEC (DoubleDelta, ZSTD),
) ENGINE = MergeTree() PRIMARY KEY (name, labels, timestamp);
```



# Small Queries

# Medium Queries

# Large Queries

---

Median Query  
Time:

0.386s

---

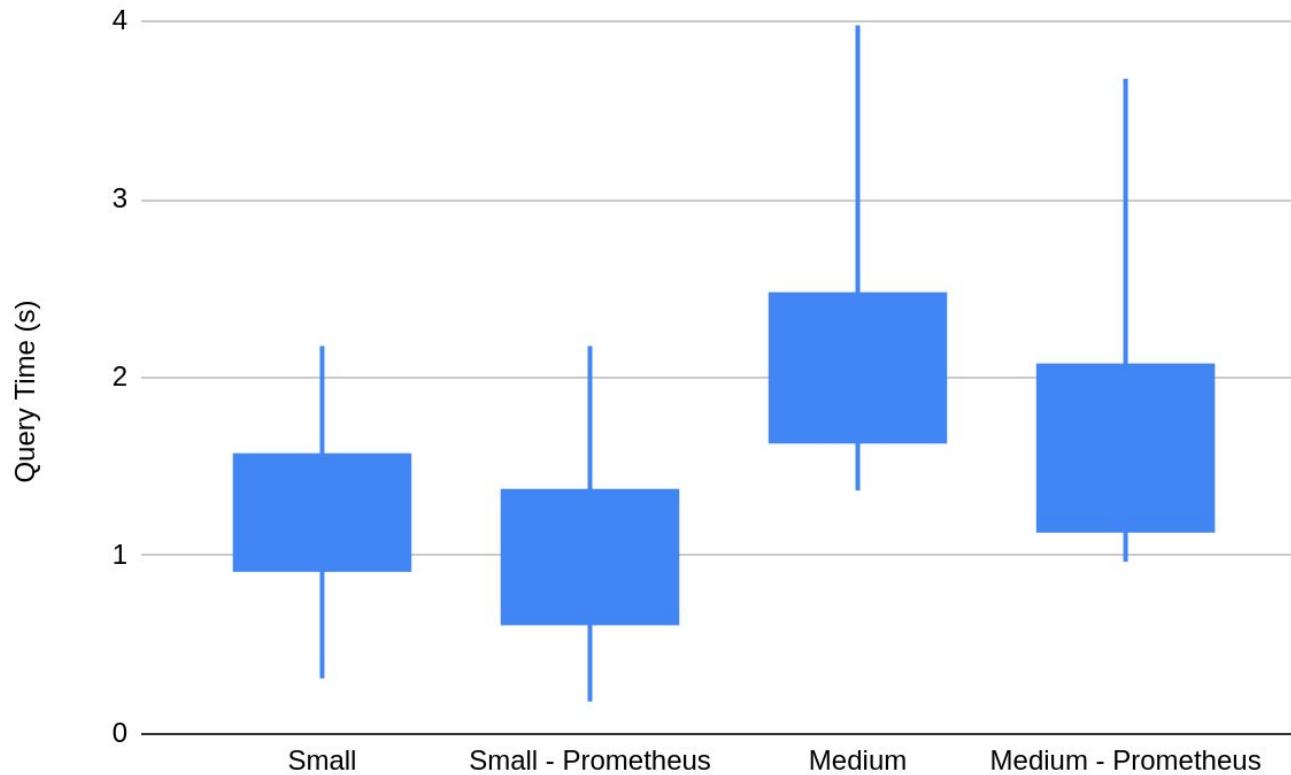
Median Query  
Time:

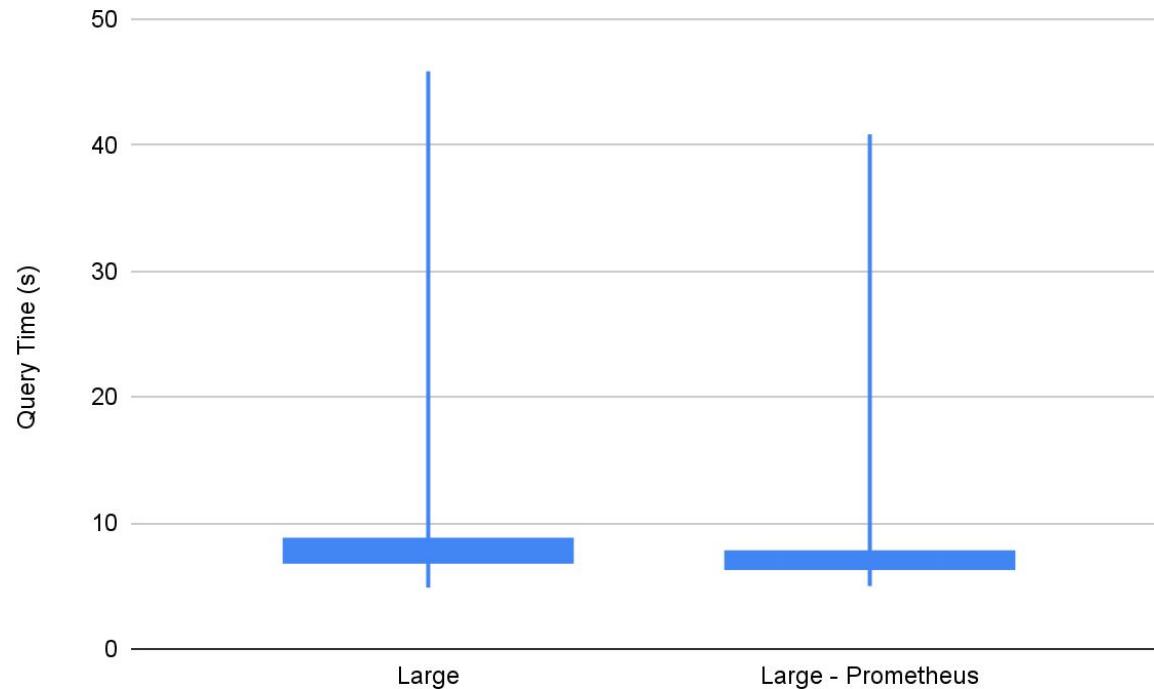
0.734s

---

Median Query  
Time:

1.623s





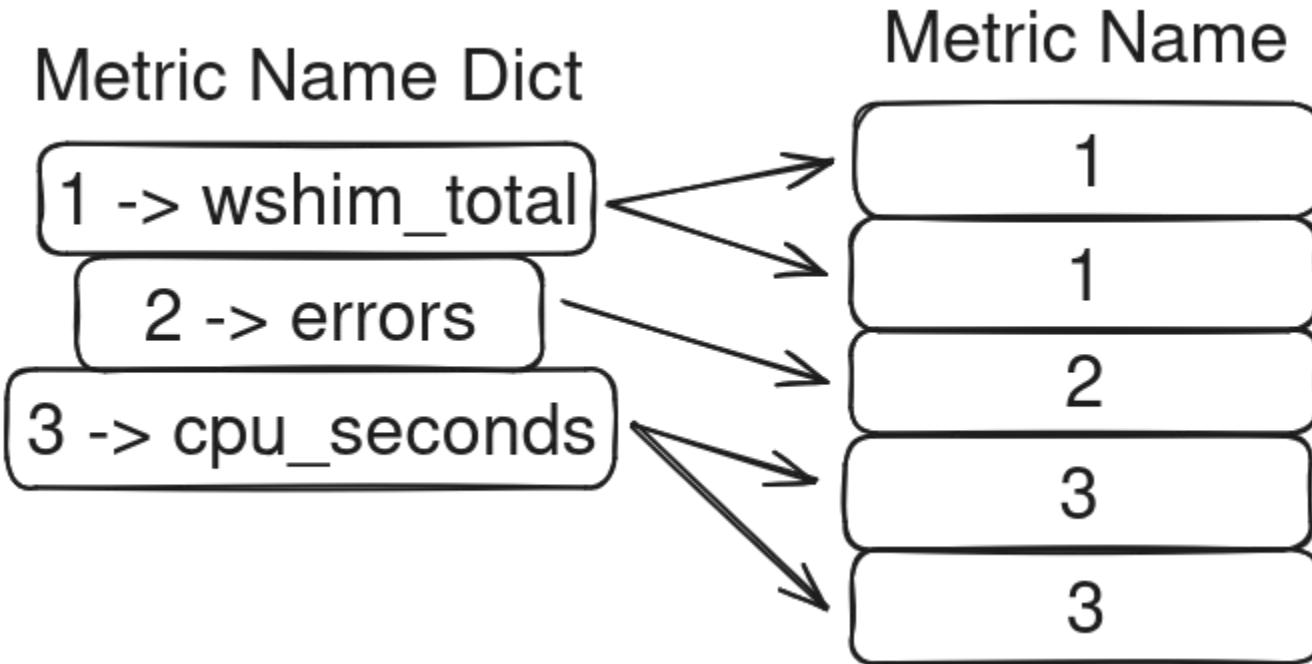
# Low Cardinality Fields

Wait... Don't we want *high* cardinality fields?

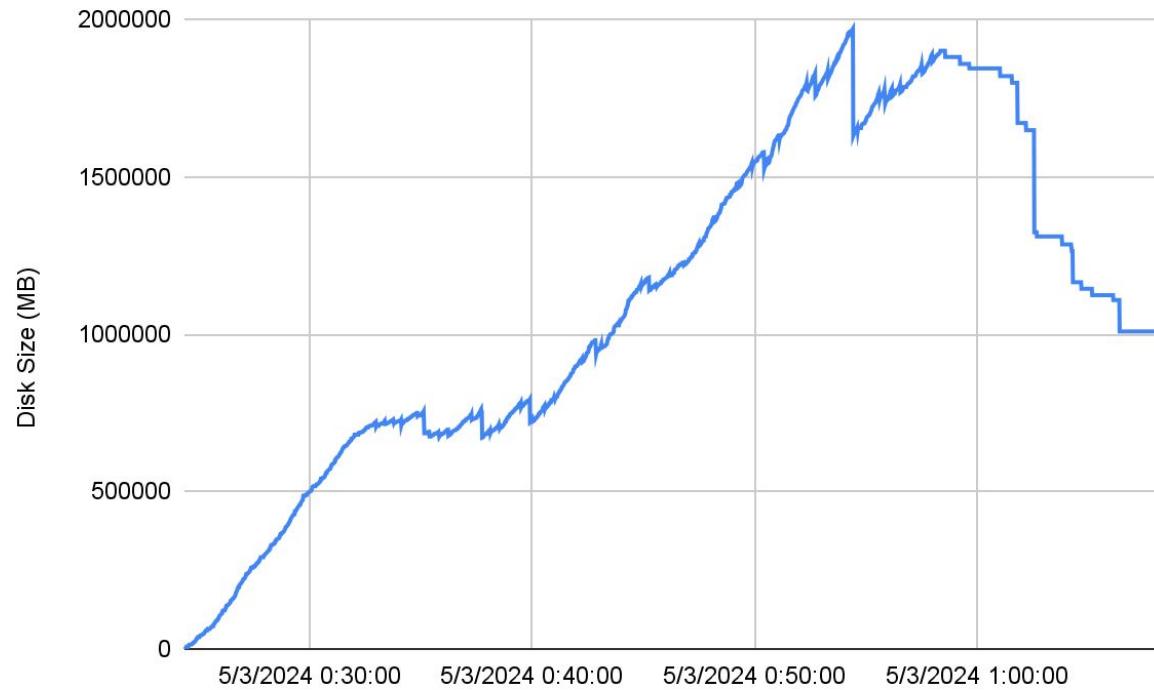


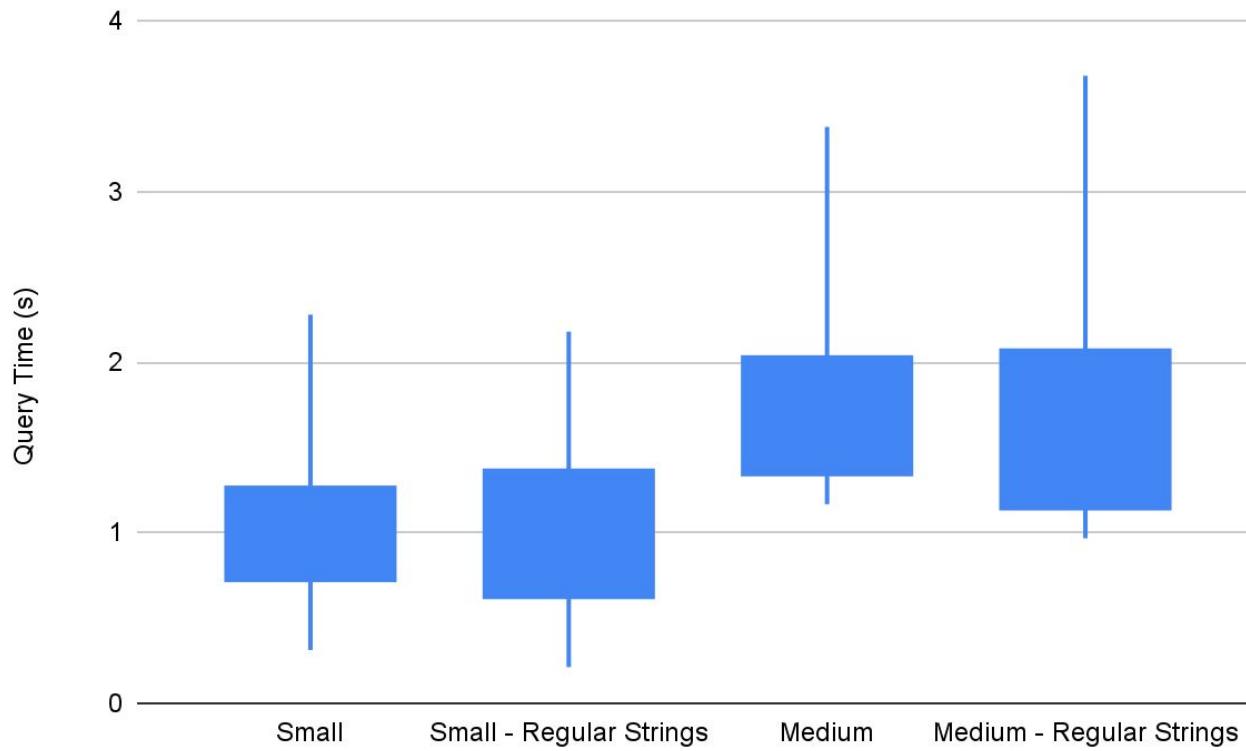
Low Cardinality to  
Prometheus is  
**vastly** different to  
Low Cardinality to  
Clickhouse





```
CREATE TABLE IF NOT EXISTS metrics (
    timestamp DateTime CODEC(DoubleDelta, ZSTD),
    name LowCardinality(String) CODEC(ZSTD),
    labels Map(LowCardinality(String), LowCardinality(String)
CODEC(ZSTD),
    value Float64 CODEC (DoubleDelta, ZSTD),
) ENGINE = MergeTree() PRIMARY KEY (name, labels, timestamp);
```





# Nested Fields

**Maps aren't the only way to store KVs**



```
CREATE TABLE IF NOT EXISTS metrics (
    timestamp DateTime CODEC(DoubleDelta, ZSTD),
    name LowCardinality(String) CODEC(ZSTD),
    labels Nested (
        key String,
        value String
    ) ,
    value Float64 CODEC (DoubleDelta, ZSTD),
) ENGINE = MergeTree() PRIMARY KEY (name, labels, timestamp);
```

## Map

name -> wshim\_total

instance=46m24

status=V

vs

## Nested fields

name

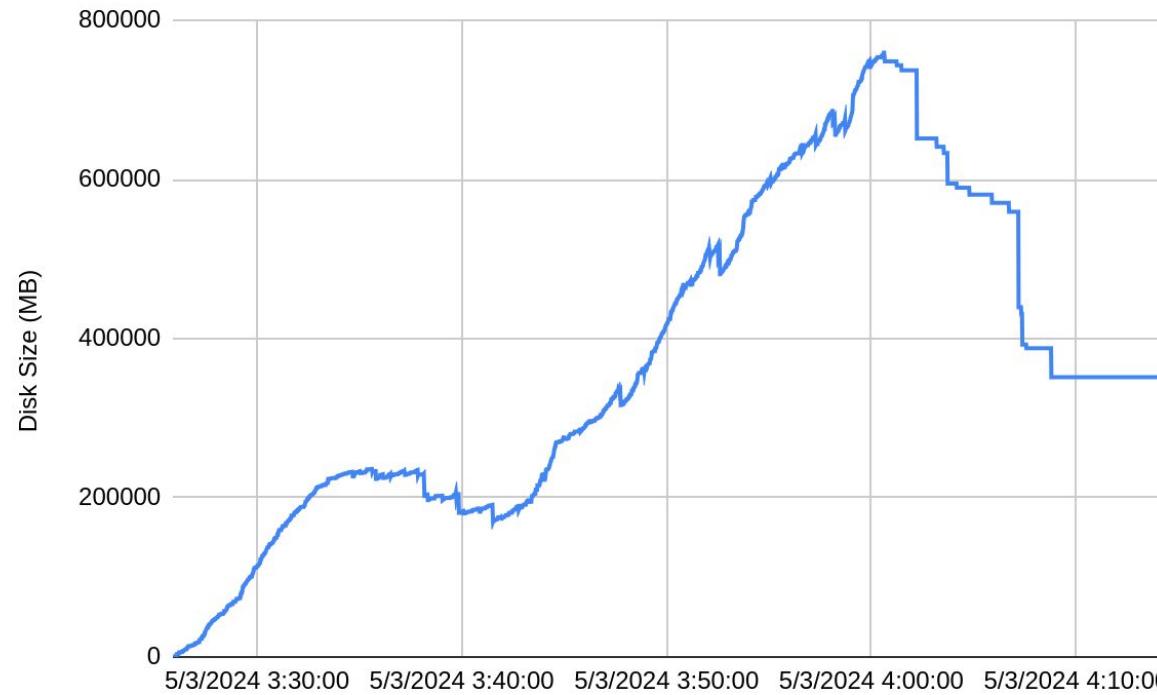
instance

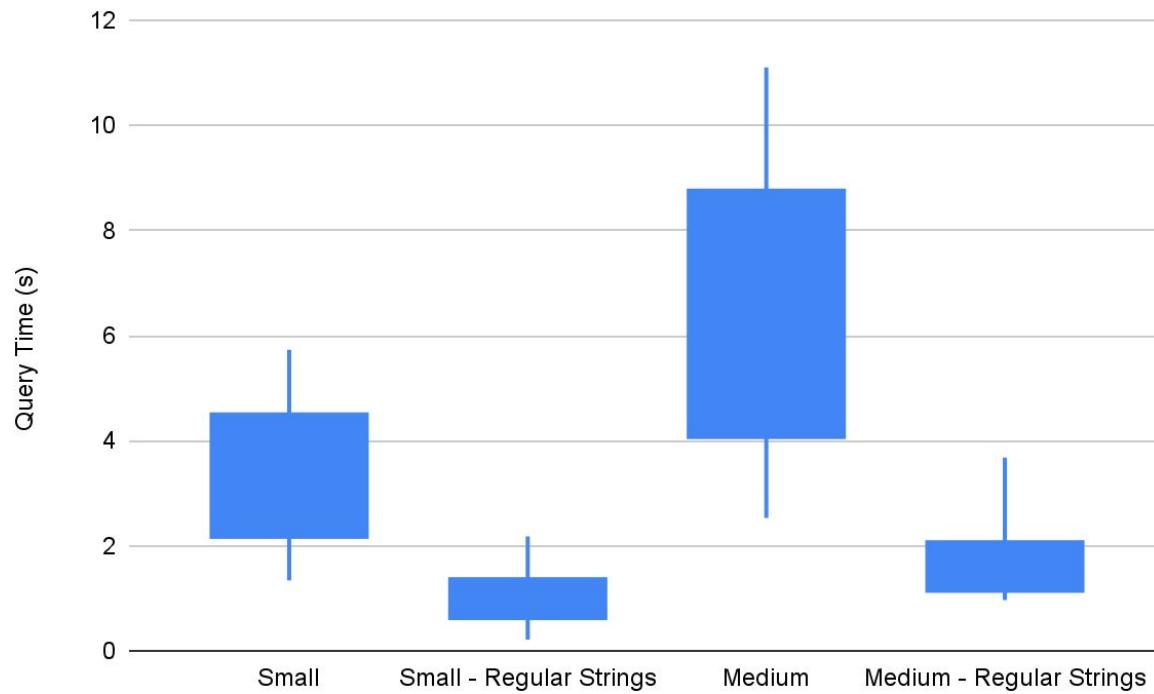
status

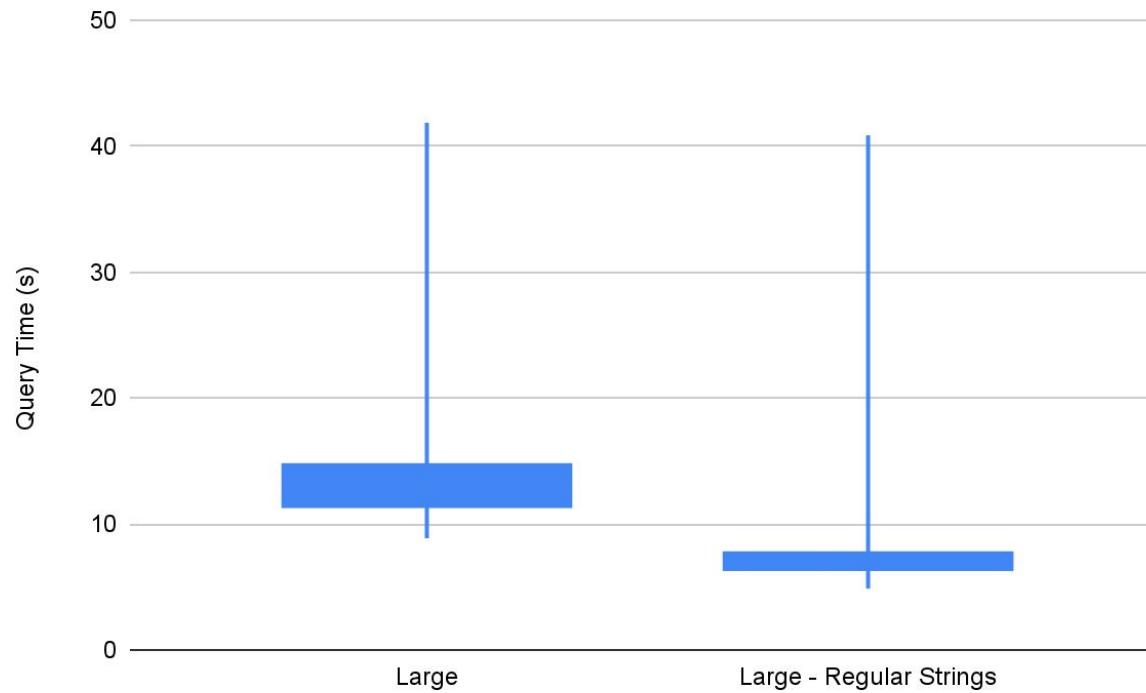
wshim\_total

46m24

V

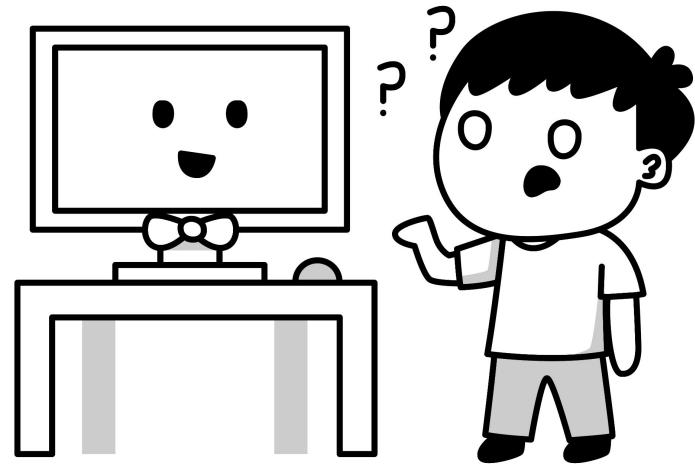






# Prometheus Encodings

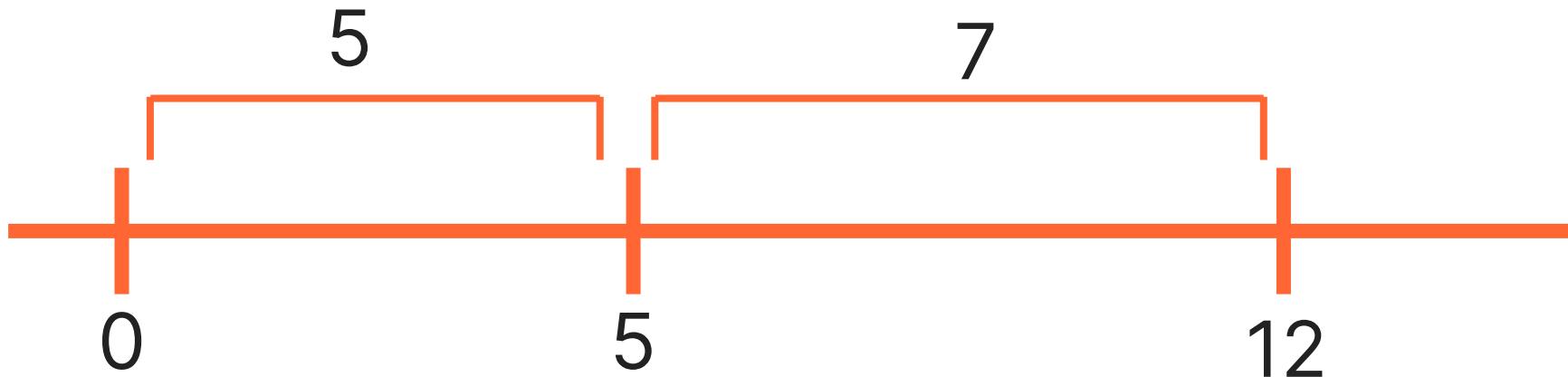
**How does this Prometheus thing work anyway?**



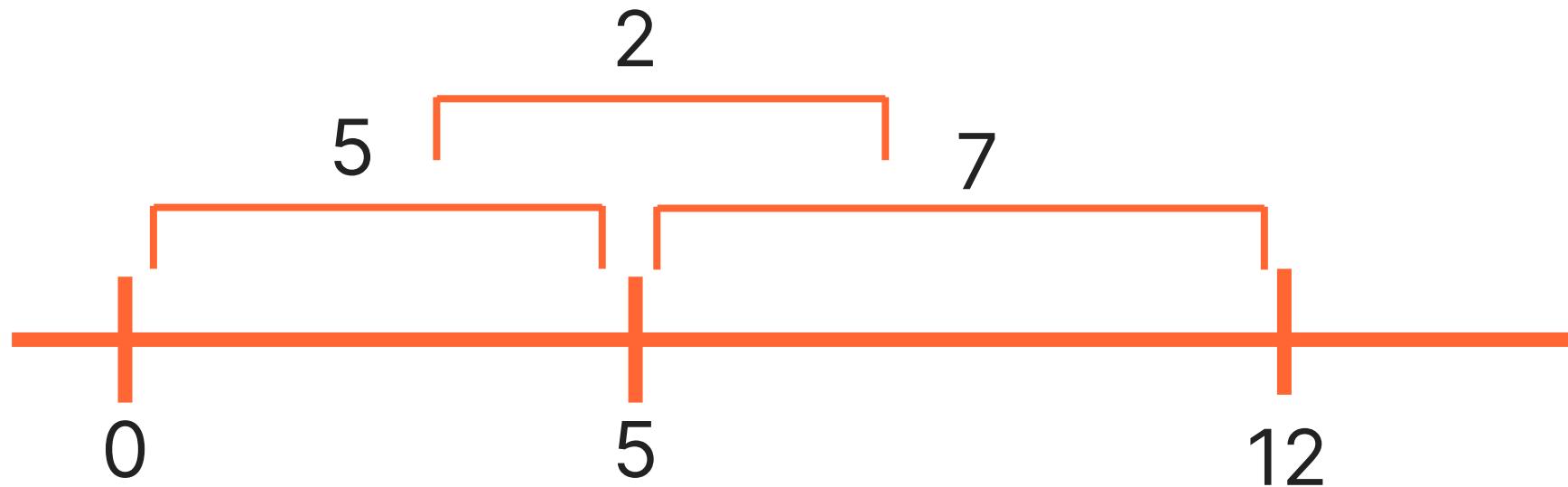


Who here knows  
how Prometheus  
encodes data?

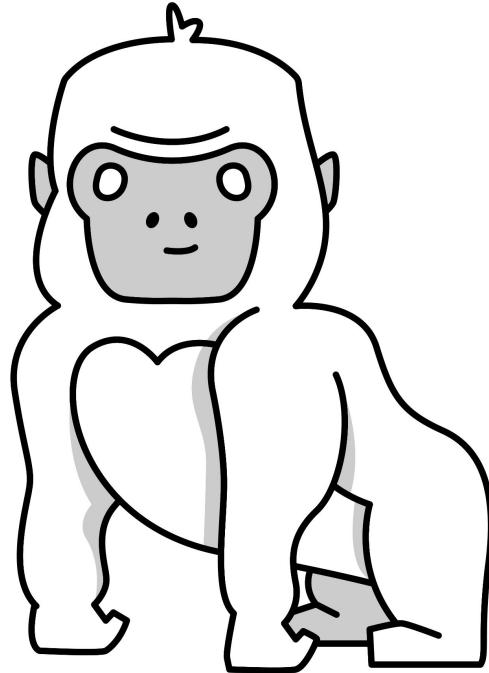
## Type 0: Delta Encoding



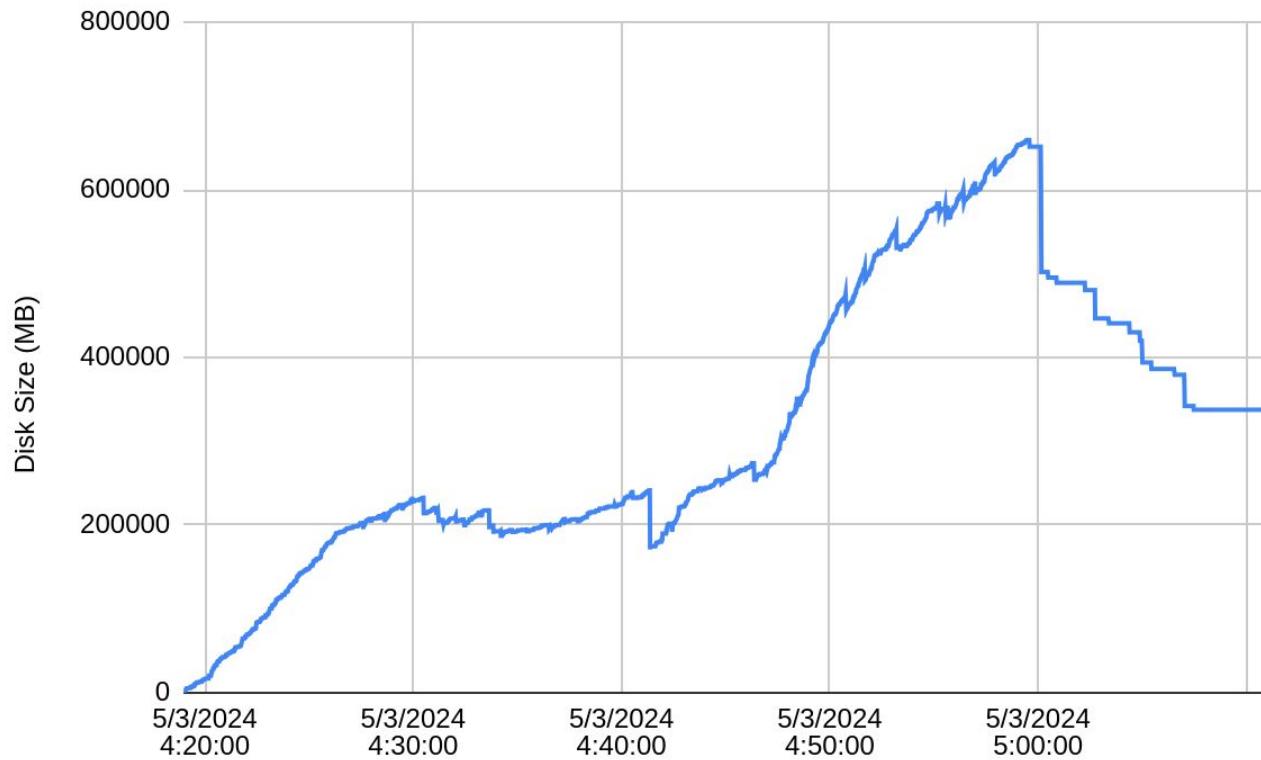
## Type 1: *Double Delta*

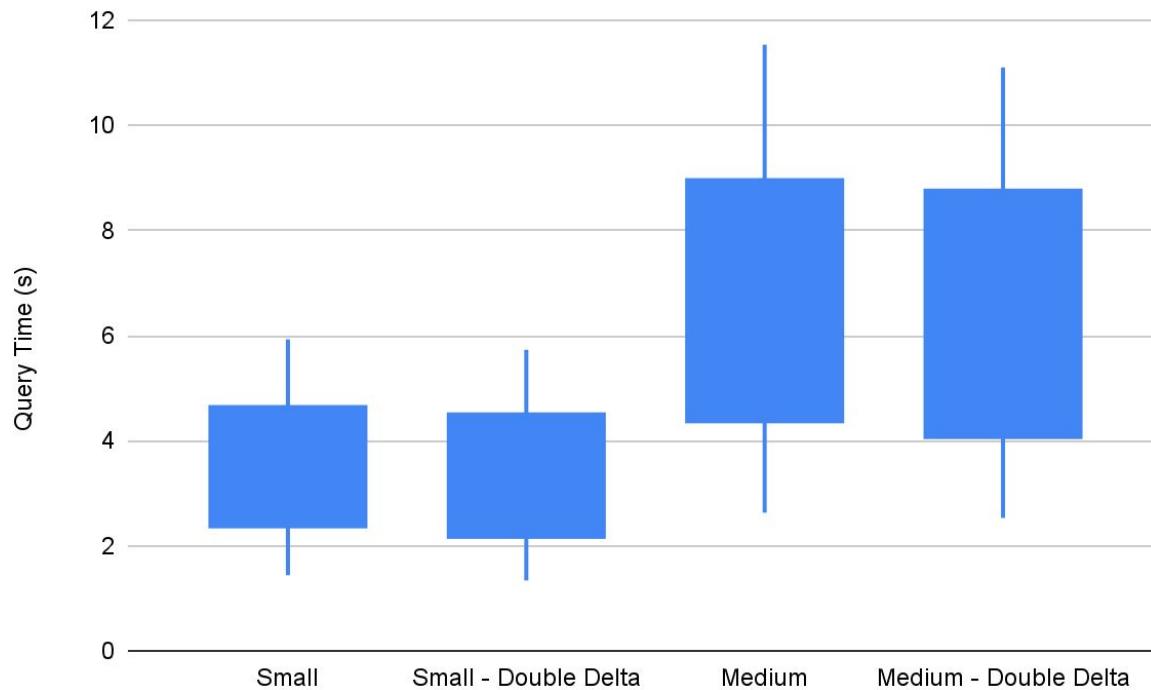


# Type 2: *Gorilla Encoding*



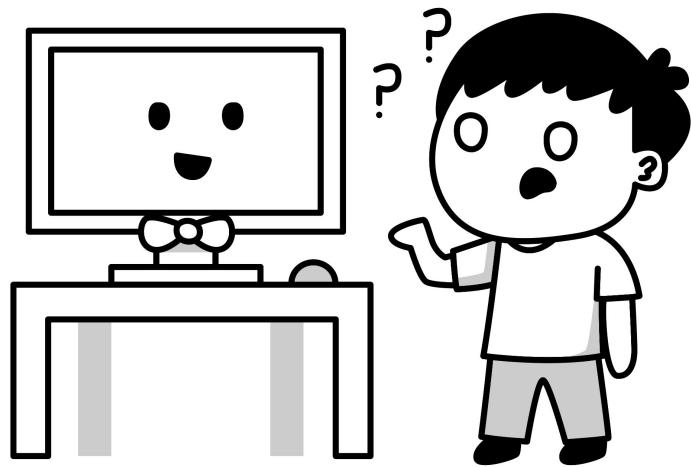
```
CREATE TABLE IF NOT EXISTS metrics (
    timestamp DateTime CODEC(DoubleDelta, ZSTD),
    name String CODEC(ZSTD),
    labels Nested (
        key String,
        value String
    ),
    value Float64 CODEC (Gorilla, ZSTD),
) ENGINE = MergeTree() PRIMARY KEY (name, labels, timestamp);
```





# Query Speeds

We've made it small, let's make it fast



## What's actually being slow?

Lots of **ors**:

```
errors_total{instance=~"41m3|41m4|42m5"}
```

Lots of **prefix matches**:

```
errors{instance=~"41.+"}  
errors{instance=~"41m3"}
```

Some **contains**:

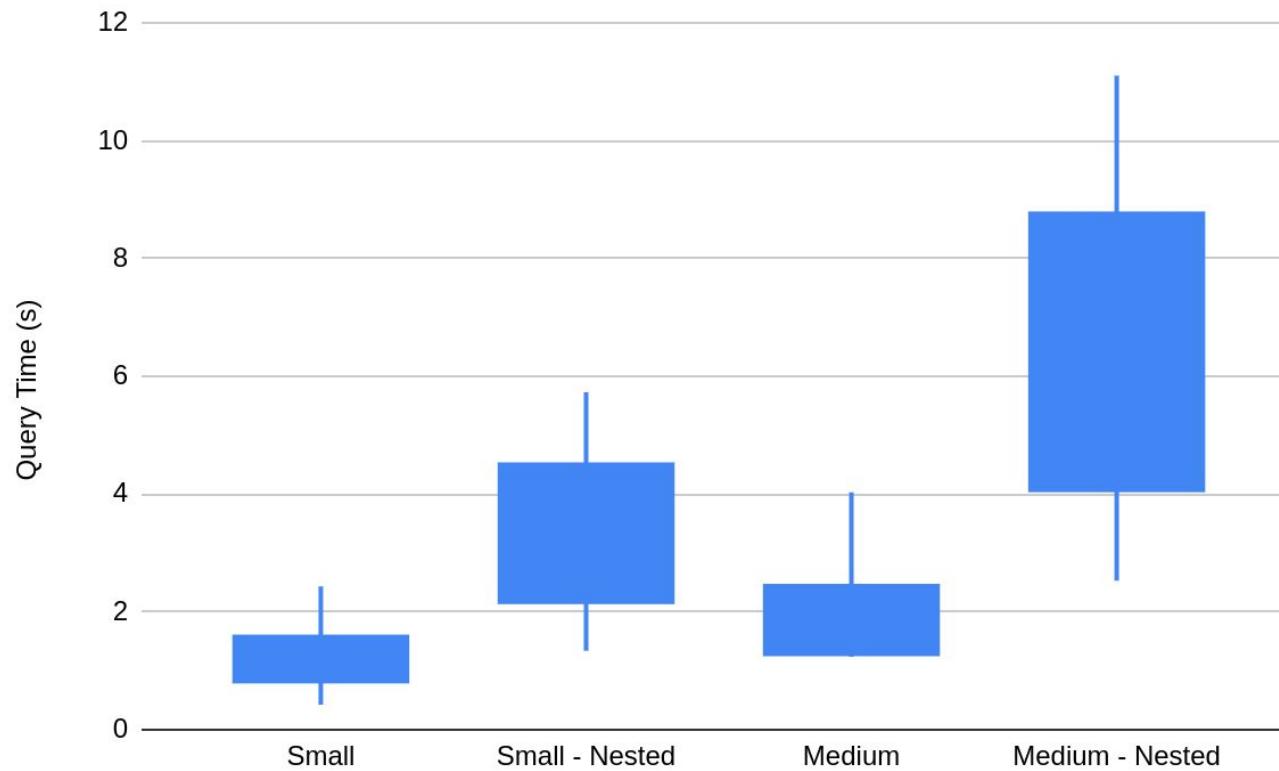
```
errors_total{instance=~".+m.+"}  
errors{instance=~".+m.+"}  
errors{instance=~".+m3"}
```

Very few actual regexes

instance=~"41m3|41m4" -> instance='41m3' or instance='41m4'

instance=~"41m.\*" -> startsWith(instance, '41m')

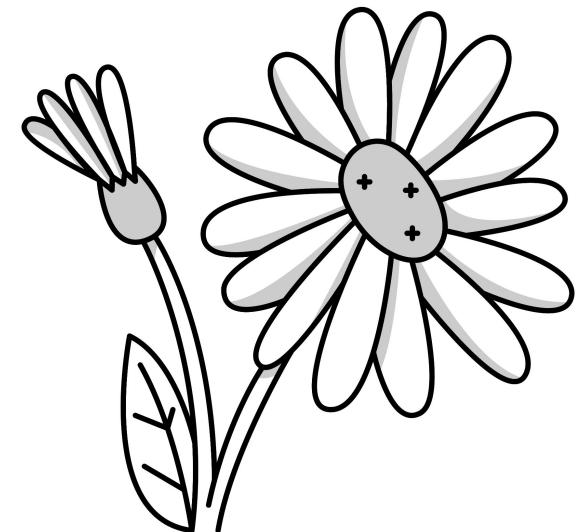
instance=~".+m.+" -> like(instance, '%m%')

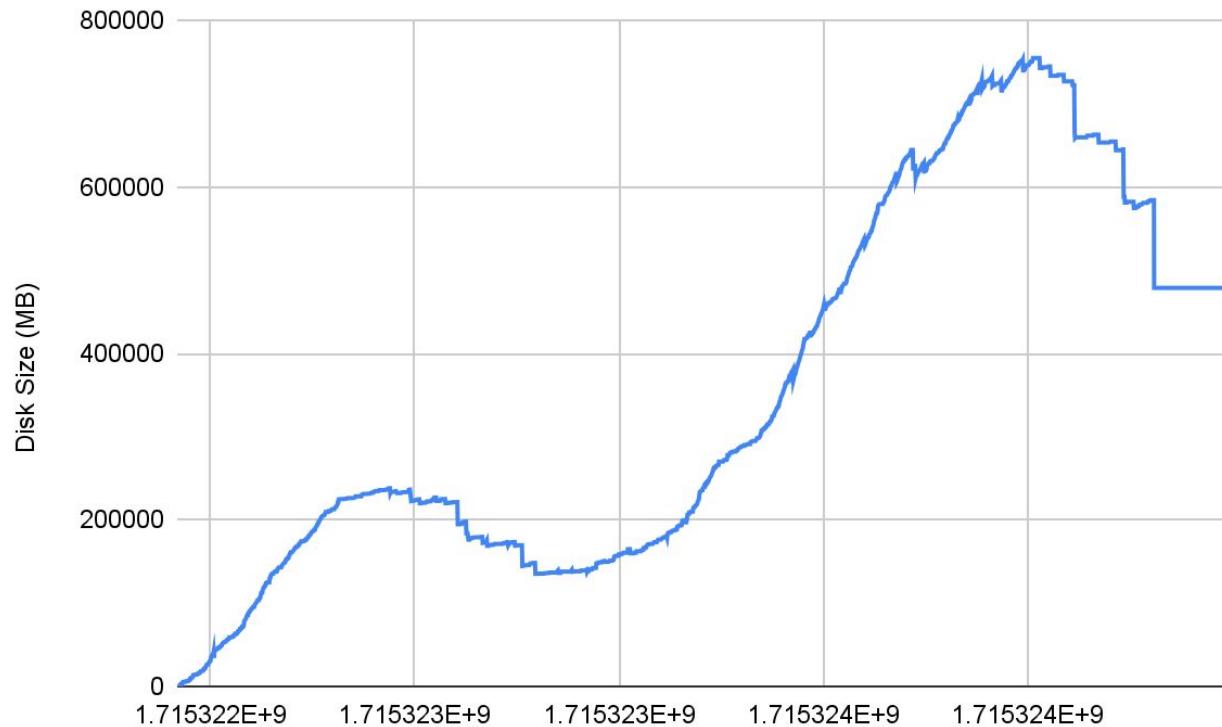


```
CREATE TABLE IF NOT EXISTS metrics (
    timestamp DateTime CODEC(DoubleDelta, ZSTD),
    name String CODEC(ZSTD),
    labels Nested (
        key String,
        value String
    ),
    value Float64 CODEC (DoubleDelta, ZSTD),
) ENGINE = MergeTree() PRIMARY KEY (name, labels, timestamp);
```

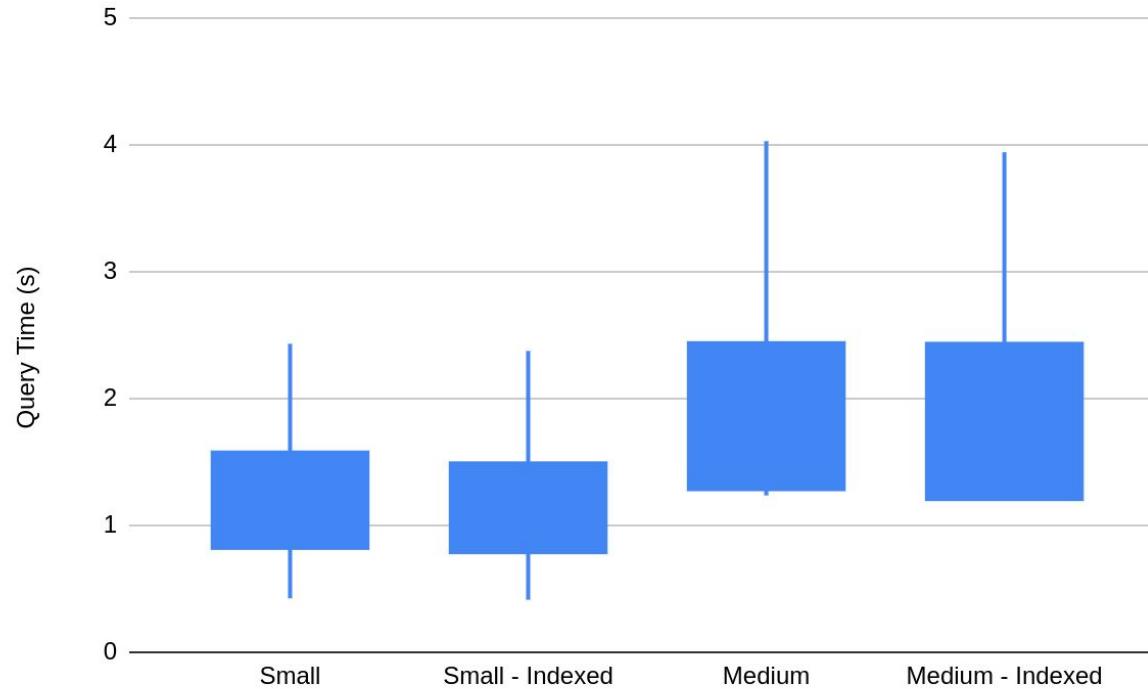
## Bloom Filter (ngrambf\_v1)

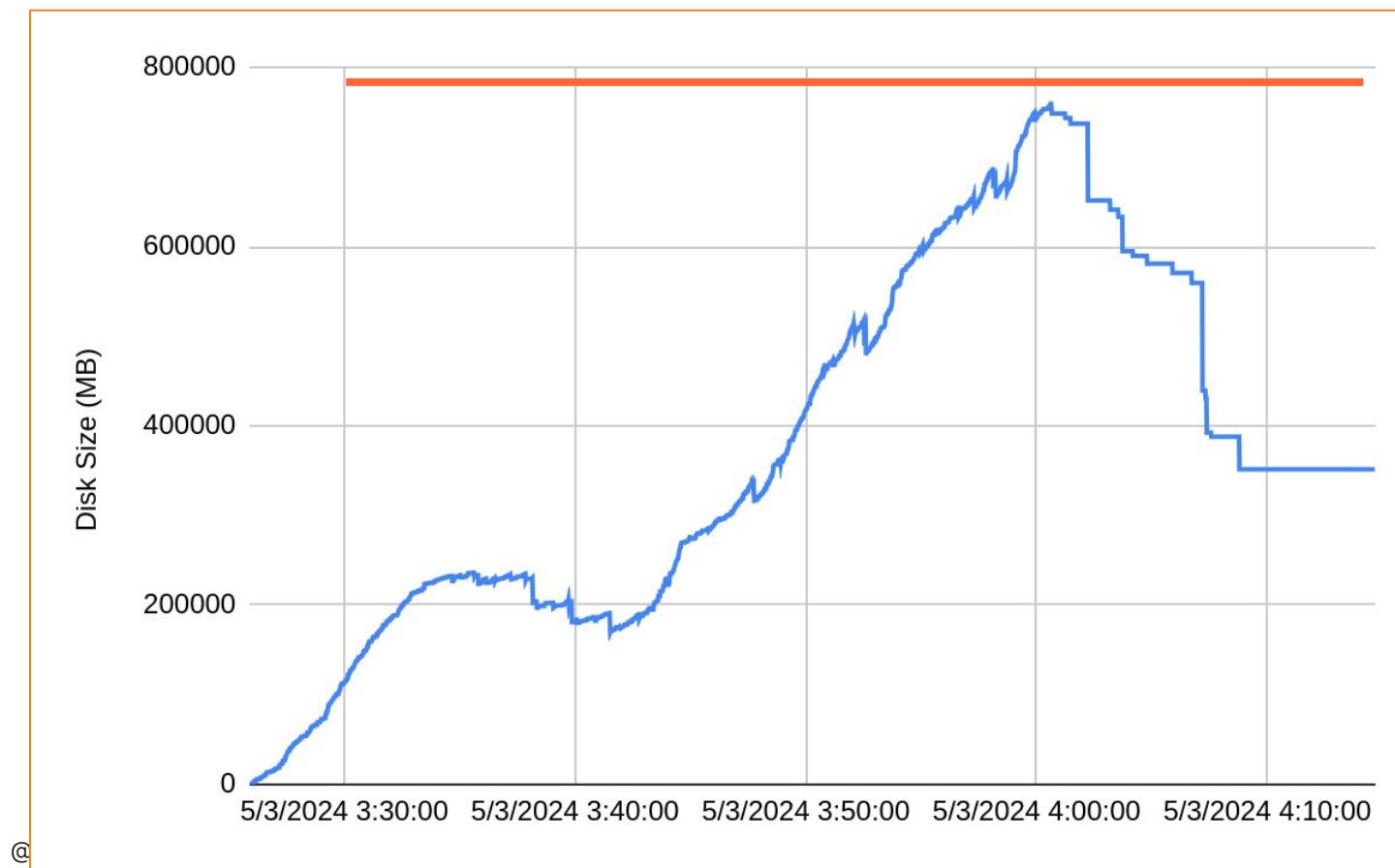
```
ALTER TABLE metrics
  ADD INDEX tags_index tags.value
  TYPE ngrambf_v1(5, 256, 3, 1)
  GRANULARITY 1;
```

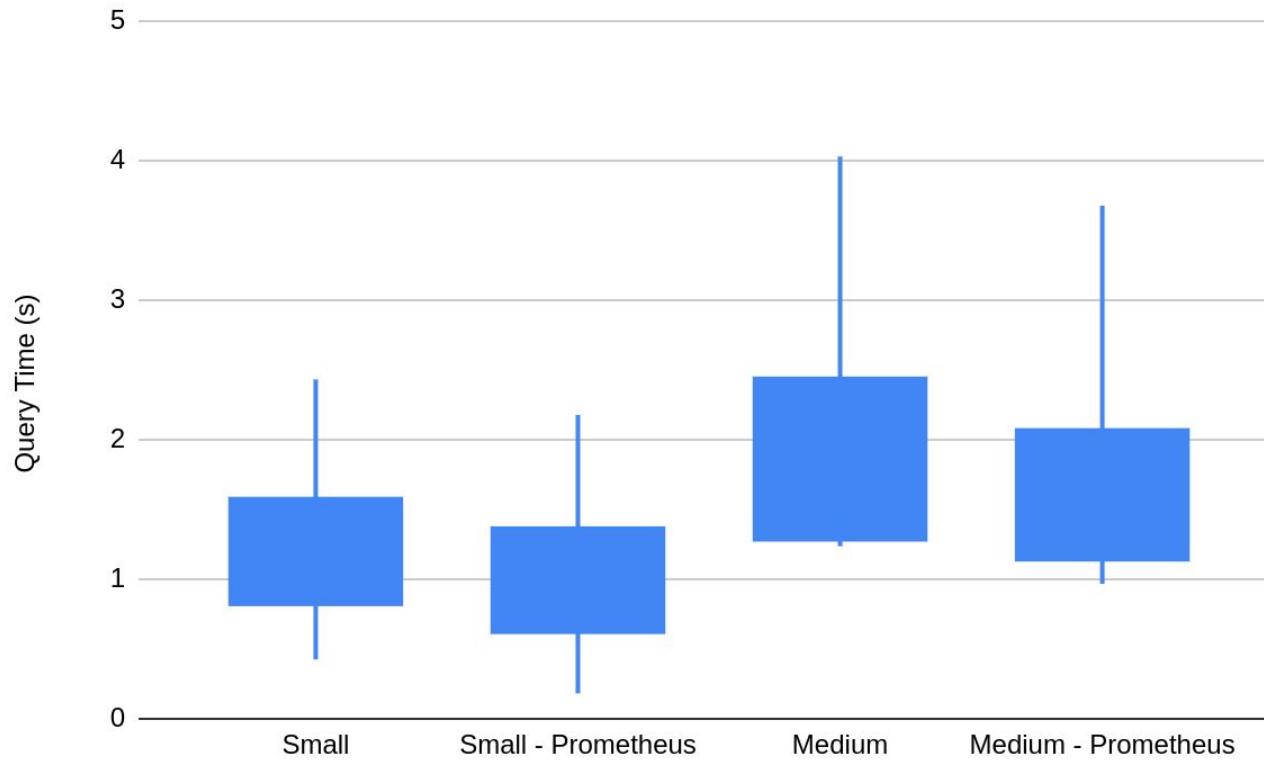




## Bloom Filter Query Times







## Future Work

- Testing in real world scraping scenarios
- Testing with Clickhouse's "ReplicatedMergeTree"
- Splitting Labels out into their own table with a "ReplacingMergeTree"
- Tuning our indexes further
- Investigating other compression algorithms

# Thanks!



Art by Alicia Edwards ^



Slides ^