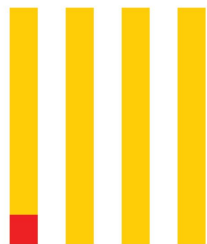# Experiments in Backing Prometheus with Clickhouse

**Colin Douch (@sinkingpoint), Observability @ Cloudflare**

# Clickhouse

- Columnar database out of Yandex

- Fast 🏎️

- Used for everything at Cloudflare



ClickHouse

# Prometheus



- Very efficient for timeseries data

- Has a few problems with cardinalities

- Lots of advice from various companies that boils down to "don't"

# Prometheus Clickhouse Bridge

- A quick remote-read/write bridge that takes remote read/write over HTTP and inserts into Clickhouse

- Doesn't support Exemplars or anything, just timeseries

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

# A bit of prior art: ❄️Epimetheus

- Made by Polar Signals

- Backed by FrostDB as an embedded DB

https://github.com/polarsignals/epimetheus

# A bit of prior art: Promhouse

- Made by Percona

- Backs Prometheus with Clickhouse, but doesn't use modern Clickhouse features

- Abandoned for the last 4 years

https://github.com/Percona-Lab/PromHouse

# A Schema

```
CREATE TABLE IF NOT EXISTS metrics (
    timestamp DateTime CODEC(Delta(4), ZSTD),
    name LowCardinality(String) CODEC(ZSTD),
    labels Map(String, String) CODEC(ZSTD),
    value Float64 CODEC (Gorilla, ZSTD),
) ENGINE = MergeTree() PRIMARY KEY (name, labels, timestamp);
```

# Results: Insertions

- Clickhouse's MergeTree engine really doesn't like small, frequent writes

- Dropping max_shards, increasing max_samples_per_send on the writing end really helped throughput

# Results: Storage

Clickhouse Disk: 6.85 bytes / sample

Prometheus Disk: 3.6 bytes / sample

# Results: Query Times

Median Query Time (Clickhouse): 1.782s
P99 (Clickhouse): 2.053s

Median Query Time (Prometheus): 1.875s
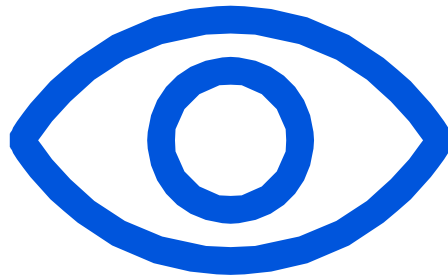P99 (Prometheus): 2.718s

# Caveats and Takeaways

Clickhouse's storage is less efficient than Prometheus for Timeseries data

Clickhouse being an external database can be scaled separately, and supports replication

This was specific to Cloudflare metrics, so might not be representative

# Thank You

CLOUDFLARE