data + ||| ClickHouse

# Singapore Meetup
July 11, 2024

# Thank you to our host!

# Tech Talks

- **The State of SQL-based Observability**
  Pradeep Chhetri, Site Reliability Engineer @ ClickHouse

- **ClickHouse: Powering Coinhall's
  Real-Time Blockchain Data Platform**
  Aaron Choo, Co-Founder & CTO @ Coinhall

- **Panel Q&A**

# The State of SQL-based Observability

ClickHouse

# Pradeep Chhetri
## Site Reliability Engineer @ ClickHouse

pradeep-chhetri-61a80935

p_chhetri

chhetripradeep

- I love playing with computers, trying out new softwares and databases.
- In my free time, i enjoy watching chess and football.

# Table of contents

**01** Introduction

**02** Challenges
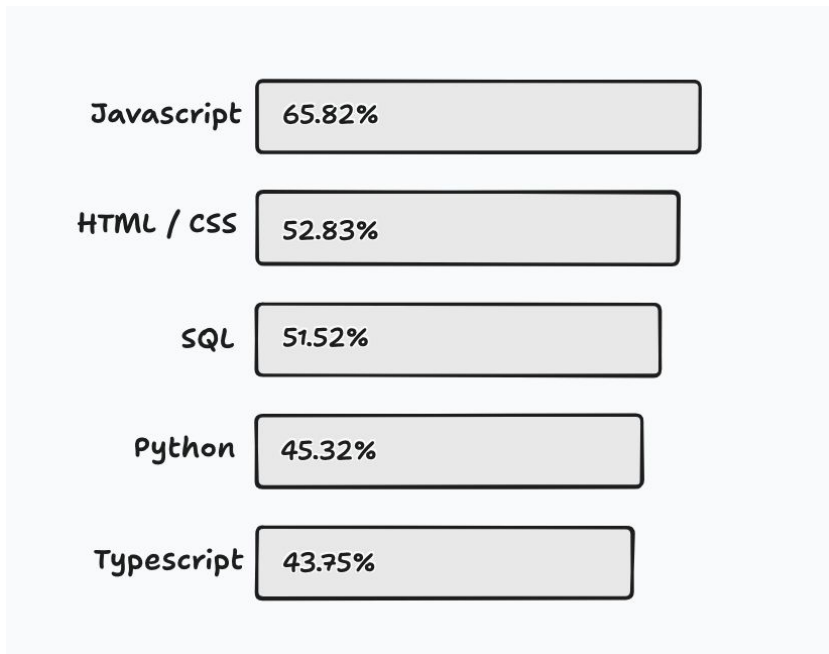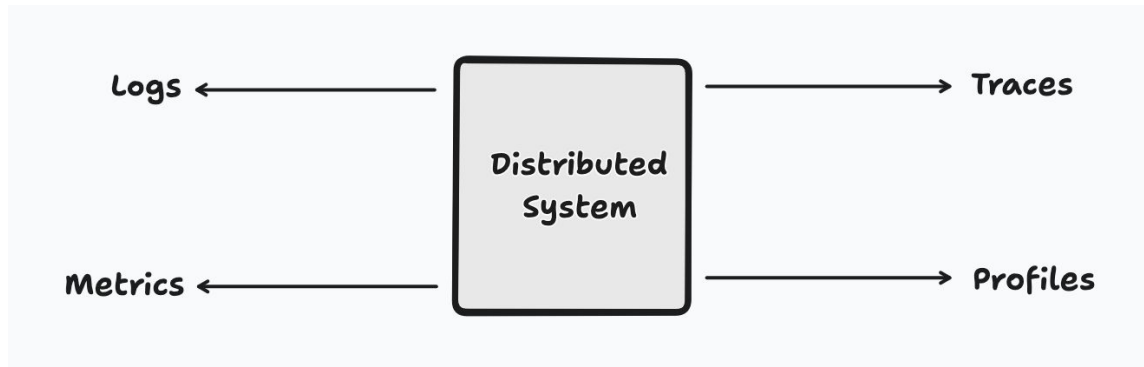
**03** Deployments

**04** Common Considerations

**05** Demo

**06** Questions

ClickHouse

# Introduction

# SQL

3rd most popular programming language

| Language | Percentage |
|----------|-----------|
| Javascript | 65.82% |
| HTML / CSS | 52.83% |
| SQL | 51.52% |
| Python | 45.32% |
| Typescript | 43.75% |

ClickHouse

# Observability

# Evolution of SQL Databases

| 50 years ago | 30 years ago | 20 years ago | 10 years ago | Present |
|---|---|---|---|---|

First SQL database invented

OLTP databases (Postgres, MySQL, SQLite)

OLAP databases Microsoft SSAS

Cloud DWHs (Redshift, Bigquery, Snowflake)

RealTime Databases (ClickHouse, Druid, Pinot)

ClickHouse

# Evolution of Observability



| 50 years ago | 30 years ago | 20 years ago | 10 years ago | Present |
|---|---|---|---|---|
| System Engineering | Centralised Logging (Syslog) | Enterprises (Splunk) | Open source Solutions (ELK, LOKI) | Open source standards (OpenTelemetry) |

ClickHouse

# Overlap of SQL Databases & Observability

| 50 years ago | 30 years ago | 20 years ago | 10 years ago | Present | |
|---|---|---|---|---|---|
| First SQL database invented | OLTP databases (Postgres, MySQL, SQLite) | OLAP databases Microsoft SSAS | Cloud DWHs (Redshift, Bigquery, Snowflake) | RealTime Databases (ClickHouse, Druid, Pinot) | SQL Based Observability |
| System Engineering | Centralised Logging (Syslog) | Enterprises (Splunk) | Open source Solutions (ELK, LOKI) | Open source standards (OpenTelemetry) | |

ClickHouse

# ClickHouse

# Challenges

# Observability Pipeline

# Challenges with Observability

High Cardinality

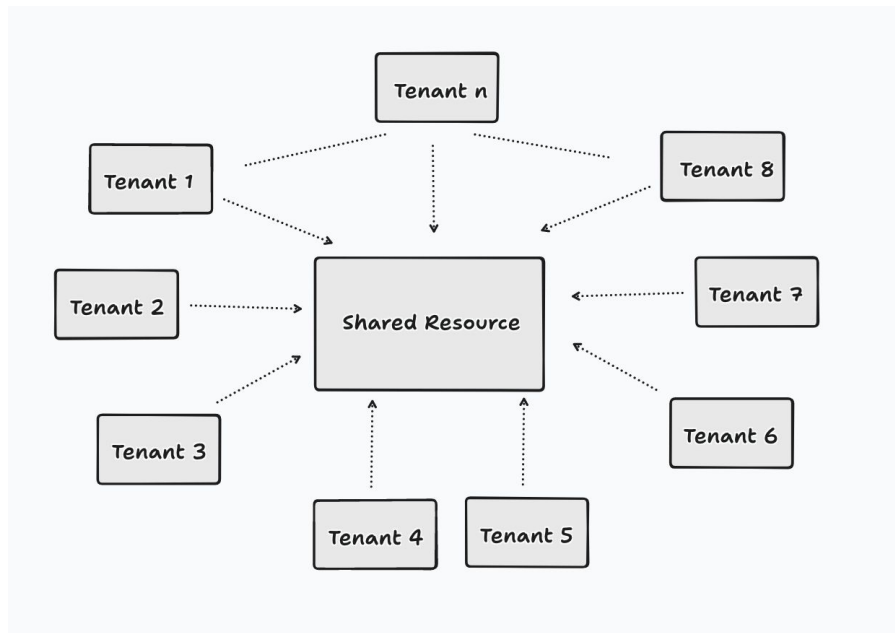| Customer ID | Customer Plan |
|---|---|
| 3A9D9780-0E89-4F3E-B299-459121D12ACC | startup |
| 8FDAFF2F-6EBF-4C07-B1A9-0D893D868B11 | enterprise |
| A61F1D6D-787C-434A-8B5C-69EF0D29A9FA | startup |
| 9884F532-8F49-432C-9D4C-A815ACB6A0F2 | enterprise |
| 1F9566B5-E2BA-467F-8F18-EC0EFCE6E39C | startup |

ClickHouse

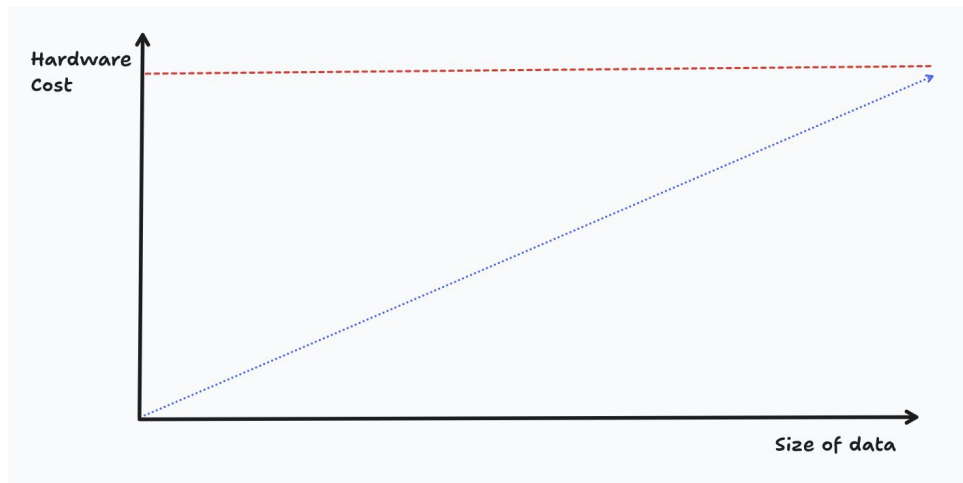# Challenges with Observability

Heavy Resource Utilization

# Challenges with Observability

Multi Tenancy Issues

# Challenges with Observability

Resources Cost

ClickHouse

# Solutions for Observability Challenges

Infinite Cardinality Support

Optimized Resource Utilization

Compute & Storage Separation

Support for Quota, Priority, Resource Management

Efficient Data Compression

Scale easily from 1 byte to 1000 petabytes

ClickHouse

"Observability is just another data problem."

# What is ClickHouse?

**open-source**

**column-oriented**

**distributed**

**OLAP database**

Developed since 2009

Open sourced in 2016

35k+ Github stars

1k+ contributors

300+ releases

Best for aggregations

Files per column

Sorting and indexing

Background merges

Replication

Sharding

Multi-master

Analytics use cases

Aggregations

Visualizations

Mostly immutable data

ClickHouse

# Real-world Deployments

# HTTP & DNS Analytics Platform

**Architecture**

Shippers → Kafka → ClickHouse

**Wins**

Efficient Ingestion & Compression

Improved Throughput & Latency

https://blog.cloudflare.com/http-analytics-for-6m-requests-per-second-using-clickhouse

CLOUDFLARE

ClickHouse

# Log Analytics Platform

**Architecture**

Log shippers → Kafka → ClickHouse → Kibana

QueryBridge to convert lucene to sql queries

**Wins**

Speed of ingestion, cost control

**Tradeoffs**

Stack administration, UI development

https://www.uber.com/blog/logging/

https://presentations.clickhouse.com/meetup40/uber.pdf

# Log Analytics Platform

```
CREATE TABLE <table_name>
(
        // Common metadata fields.
        _namespace              String,
        _timestamp              Int64,
        hostname                String,
        zone                    String,
        ...

        // Raw log event.
        _source                 String,

        // Type-specific field names and field values.
        string.names            Array(String),
        string.values           Array(String),
        number.names            Array(String),
        number.values           Array(Float64),
        bool.names              Array(String),
        bool.values             Array(UInt8),

        // Materialized fields
        bar.String              String,
        foo.Number              Float64,
        ...
)
...
```

Uber

ClickHouse

# Adopting OLAP store for tracing

**Architecture**

OpenTelemetry → ClickHouse → Kibana

**Wins**

Data compression, open source licensing

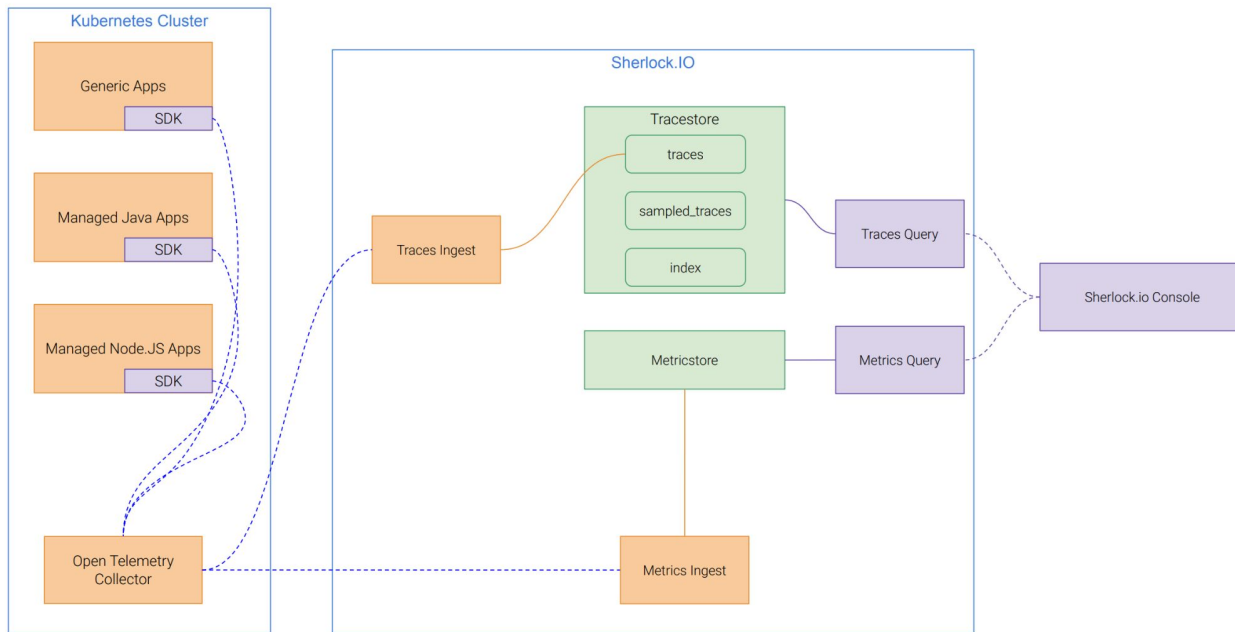**Tradeoffs**

Managing tiered OTel collectors

https://kccnceu2024.sched.com/event/1YeNu

# Adopting OLAP store for tracing

# Dogfooding ClickHouse across o11y

**Architecture**

OpenTelemetry → ClickHouse → Grafana
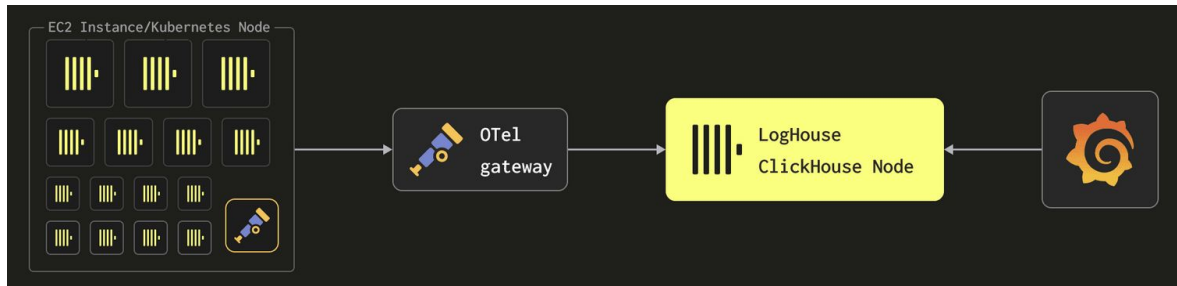
**Wins**

Granular log retention, Saved money on Datadog

**Tradeoffs**

1.5 FTEs to build stack

https://clickhouse.com/blog/building-a-logging-platform-with-clickhouse-and-saving-millions-over-datadog

**ClickHouse**

# Dogfooding ClickHouse across o11y

# Common considerations

ClickHouse

# Query language considerations

*"SQL is not compact enough compared to domain-specific query languages"*

ClickHouse

# A simple query ?

```
source=events level="warning"
| STATS avg(duration) BY level
| FIELDS level, avg(duration) AS avg_dur
| sort - avg_dur | head 10
```

```
GET events/_search
{
  "size": 0,
  "_source": false,
  "track_total_hits": -1,
  "aggregations": {
    "groupby": {
      "composite": {
        "size": 10,
        "sources": [
          {
            "4e8796da": {
              "terms": {
                "field": "level.keyword",
                "missing_bucket": true,
                "order": "asc"
              }
            }
          }
        ]
      },
      "aggregations": {
        "c3318afb": {
          "avg": {
            "field": "duration"
          }
        }
      }
    }
  }
}
```

ClickHouse

# In good old SQL

```sql
SELECT
    level,
    avg(duration) AS dur
FROM events
GROUP BY level
ORDER BY dur DESC
```

ClickHouse

# Schema considerations

It helps to stop thinking about *"metrics"*, *"logs" and "traces"* separately and just think them as *"wide events"*

## All you need is Wide Events, not "Metrics, Logs and Traces"

IVAN BURMISTROV
FEB 15, 2024

https://isburmistrov.substack.com/p/all-you-need-is-wide-events-not-metrics

https://news.ycombinator.com/item?id=39529775
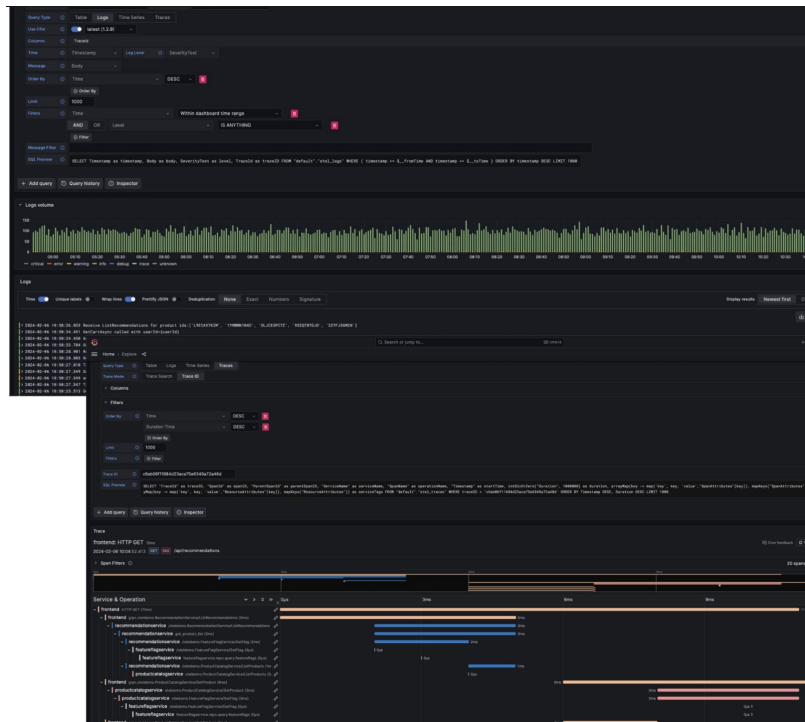
# Visualization considerations

Grafana

Apache Superset

Perses

Metabase

Build your own

ClickHouse

# Multi-tenancy considerations

*Example: Uber*

Consider datastore
ability to limit resources
by table, user, session
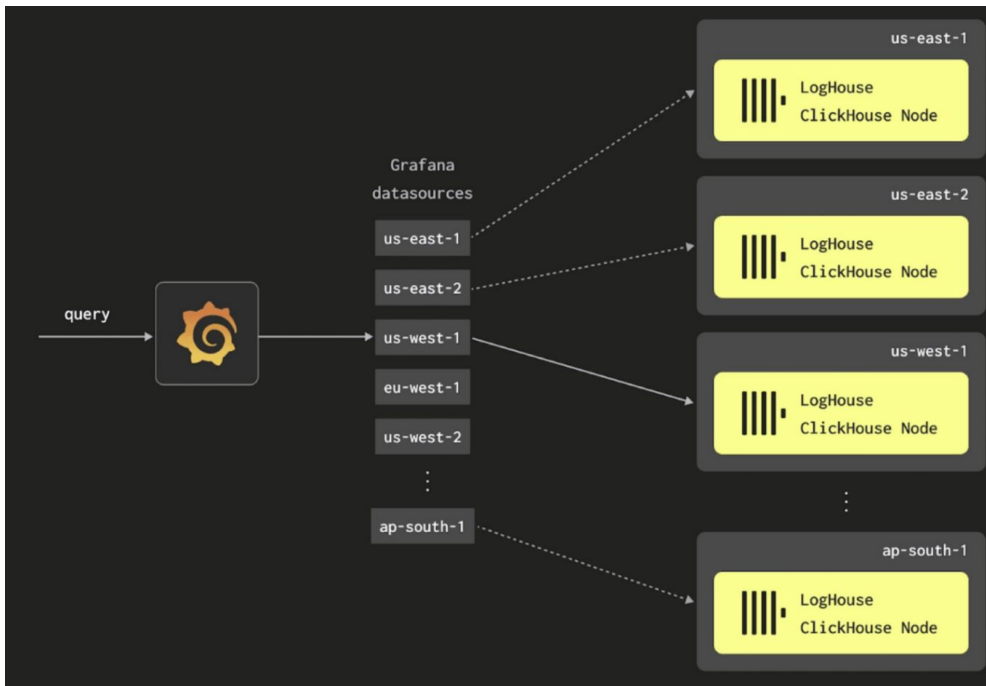
## Unified Multi-Tenant Storage Platform

- ClickHouse natively supports zero lock contention among concurrent reads and writes
- Service placement: single-tenant vs multi-tenant
  - Isolate heavy log producers, heavy log consumers
  - Co-locate everything else
  - Limit the impact of co-location, add service in order-by
- Workload isolation
  - Configure query parallelism per query
  - Eventually limit total query resource usage per node
  - Query cost accounting, defense against expensive queries

# Multi-region deployment considerations

*Example: ClickHouse Inc.*

Per region data collection, cross-region queries

Resilient to AZ outage but not region outage

# Choice of analytical datastore matters

| | ClickHouse | Druid | Pinot | BigQuery |
|---|---|---|---|---|
| **Real-time speed** | ✔ Best | ✔ Ok | ✔ Ok | ✗ Poor |
| **Compression** | ✔ Best | ✔ Ok | ✔ Ok | ✔ Better |
| **Sep storage & compute** | ✔ | ✗ | ✗ | ✔ |
| **Interoperability**<br>- OTel<br>- Grafana | ✔<br>✔ | ✗<br>✔ (no logging & tracing support) | ✗<br>✗ | ✗<br>✔ (no logging & tracing support) |
| **SQL compliance** | ✔ Good | ✗ Poor | ✗ Poor | ✔ Best |
| **TCO improvements** | ✔ 10x | 1–2x | 1–2x | depends on how often you query |

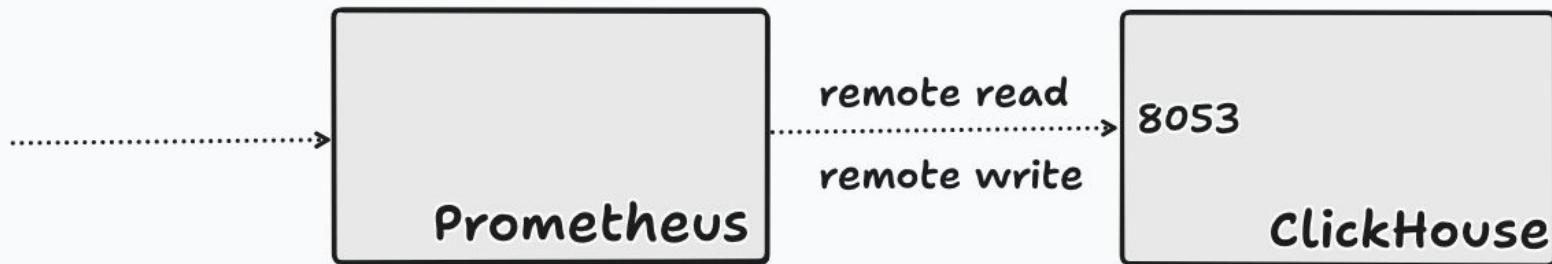https://benchmark.clickhouse.com

ClickHouse

# ClickHouse

# Demo

# ClickHouse Timeseries Engine

# ClickHouse Timeseries Engine

Pull Request: https://github.com/ClickHouse/ClickHouse/pull/64183

# Questions