

Distributed Tracing on ClickHouse

eBay Inc.

OSACON-2021

Sudeep Kumar, Amber Vaidya

November 2nd, 2021



eBay Scale



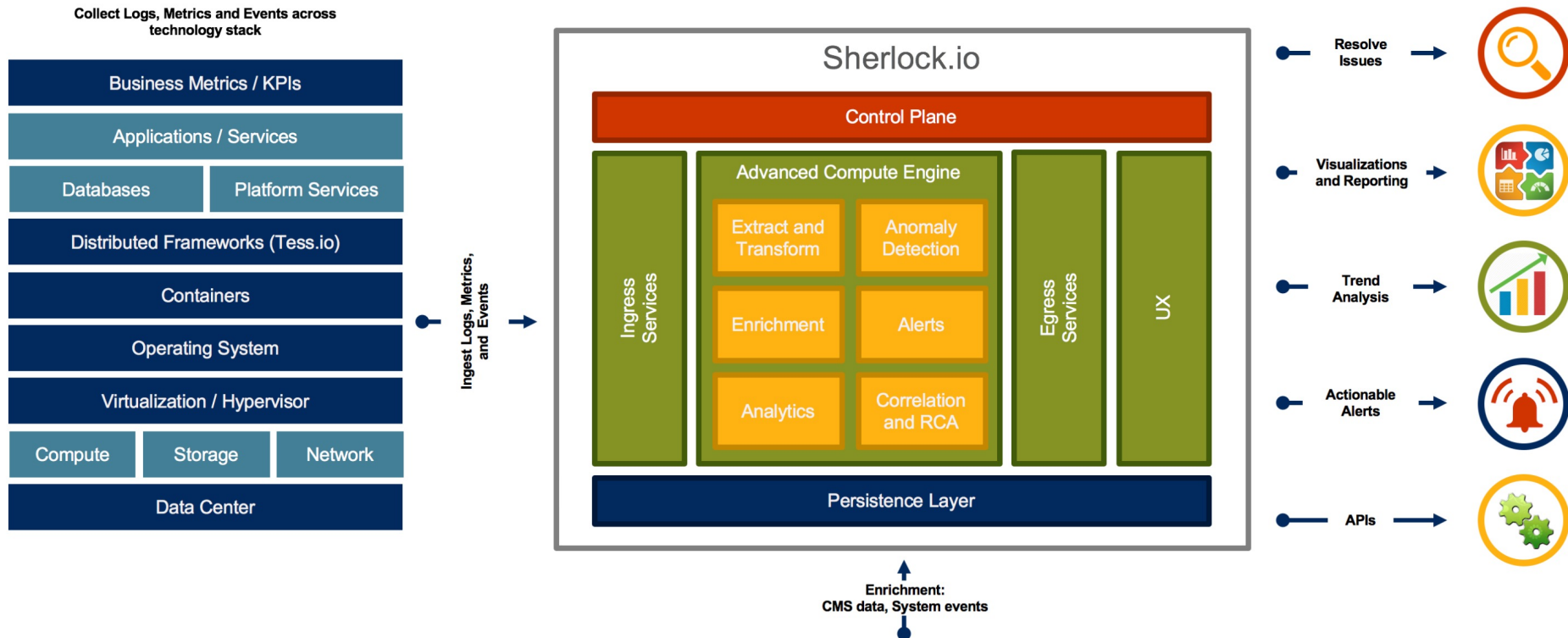
Source: Euromonitor International, 2017; Note: Regions are North America, Latin America, Europe, Asia-Pacific and Australasia/New Zealand. Latin America includes Brazil and Mexico.

- ❖ Global Commerce Leader with presence in more than 190 markets
- ❖ 154 Million Active Buyers
- ❖ \$19.5 Billion GMV in Q3 2021
- ❖ \$2.5 Billion in Revenue in Q3 2021

Sherlock.io

- ❖ Observability platform for EBAY's monitoring needs
- ❖ Supports different telemetry signals
 - *Metrics*
 - *Logs*
 - *Events*
 - *Traces*
- ❖ Scale of Operation (millions per second)

Platform Overview



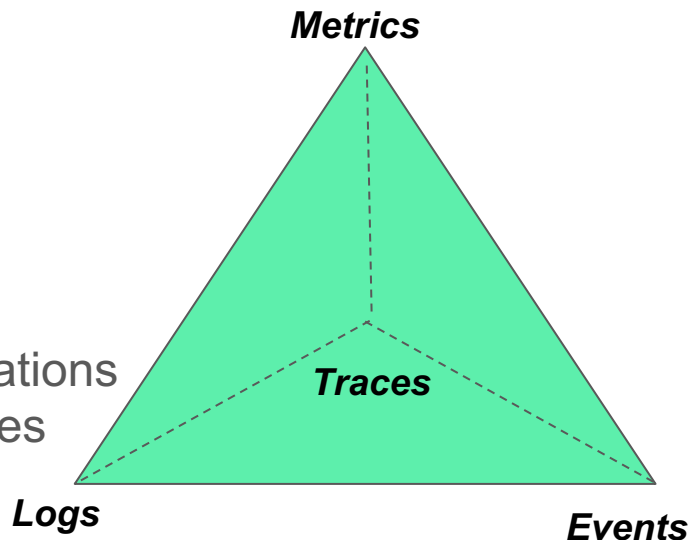
Tracing

- Monitoring signals

- Metrics
- Logs/Events
- Traces

- Distributed Tracing

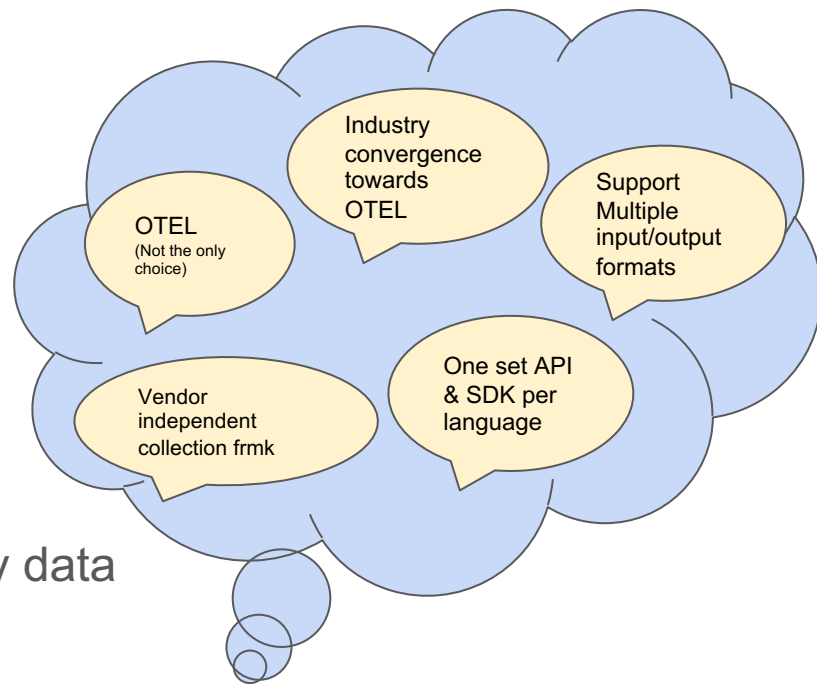
- Tracks request flows through multiple applications
- Helps identify errors/bugs, performance issues
- Service dependencies



Use-case: I know my 99% percentile latency of my service is 10s, but in what context was high latency observed. Also, was there an event or log associated with that context ?

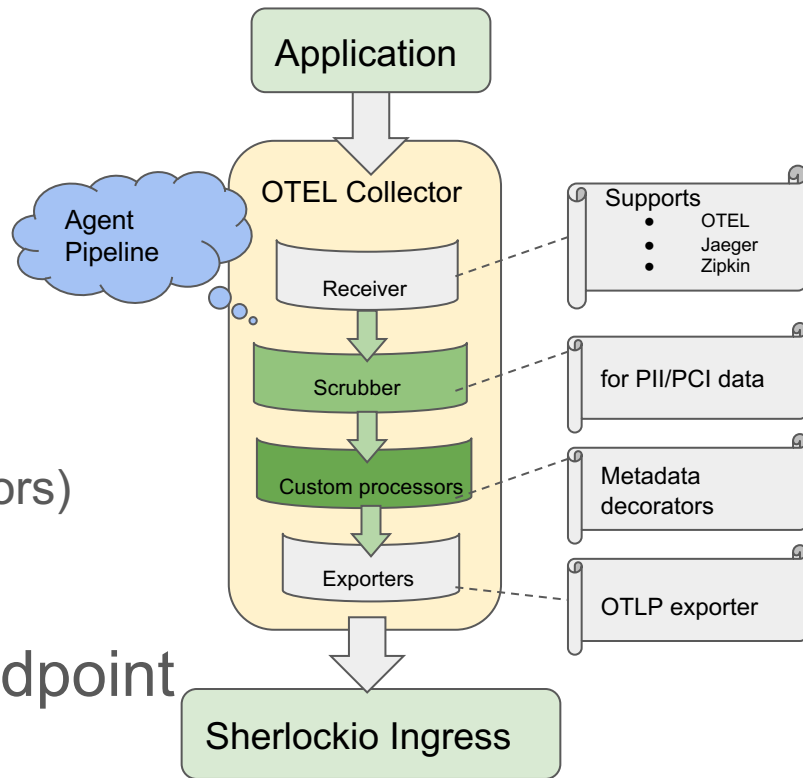
Open Telemetry (OTEL)

- Observability framework
 - CNCF project
 - Provides agents/libraries
- OTEL Collector
 - vendor-agnostic implementation
 - receive, process & export telemetry data



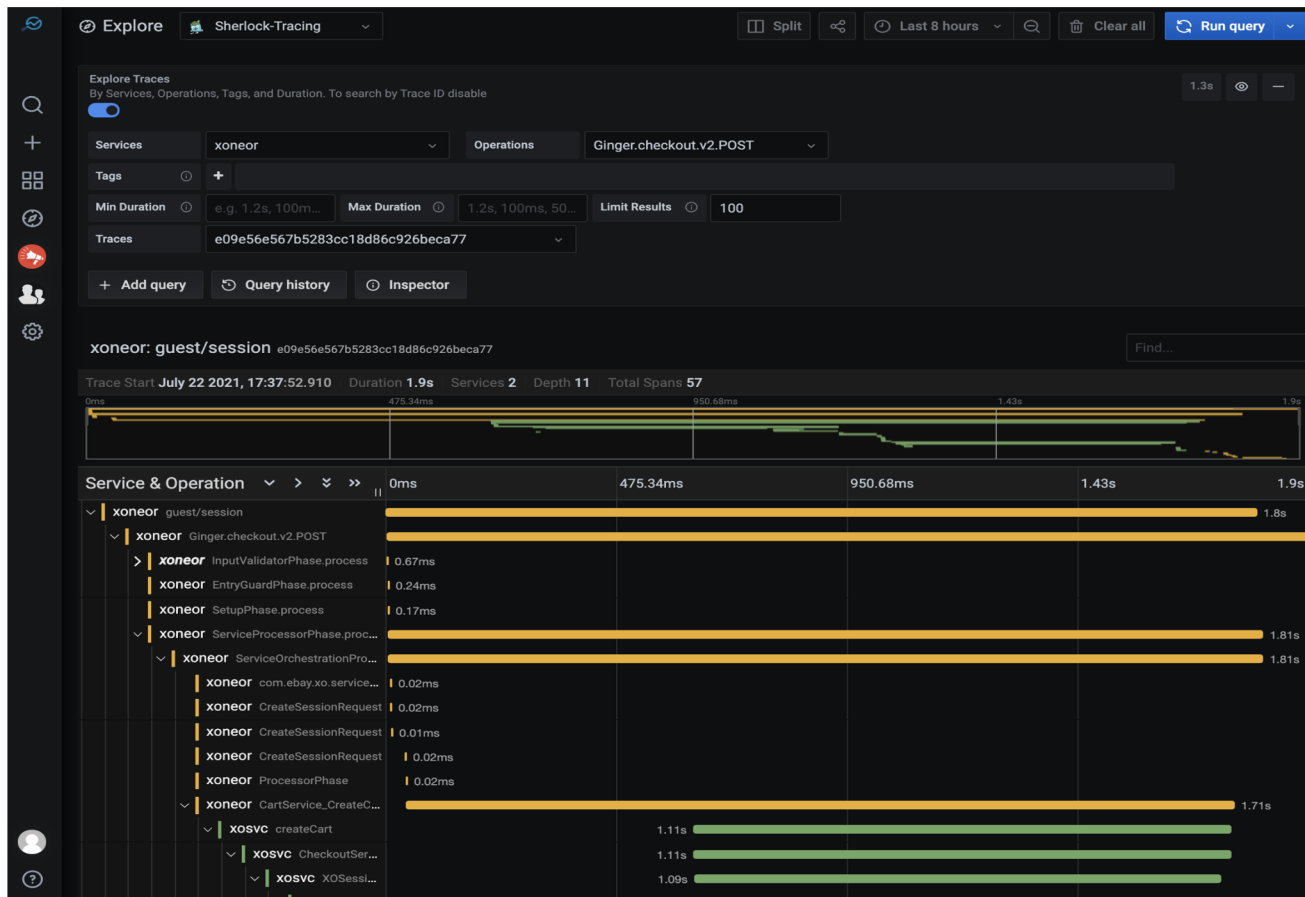
OTEL Collector

- Consists of receivers/exporters
 - Supports various i/o formats
- Intermediate Processors
 - Scrubbing
 - Metadata decorators (k8s processors)
 - RL/Sampling
- Exports to a custom ingestion endpoint
 - Otlp/Otlphttp format



Visualization

- Egress compliance with Jaeger Query APIs
- UX experience
 - Grafana jaeger plugin
 - Optional Jaeger UI



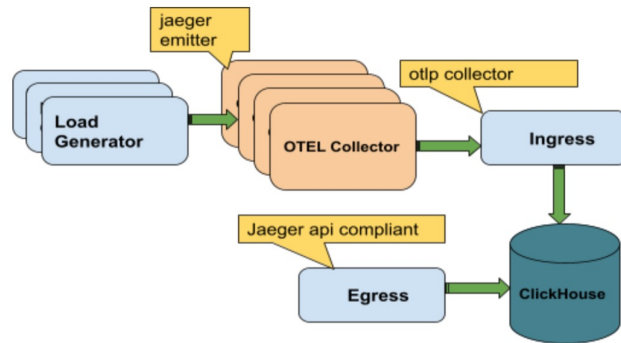
Storage evaluations

- Storage Options



- Criteria Used

- High ingestion volumes
- Can support Jaeger Query API Patterns
- Point Trace-Id lookups for a time interval



Some Observations

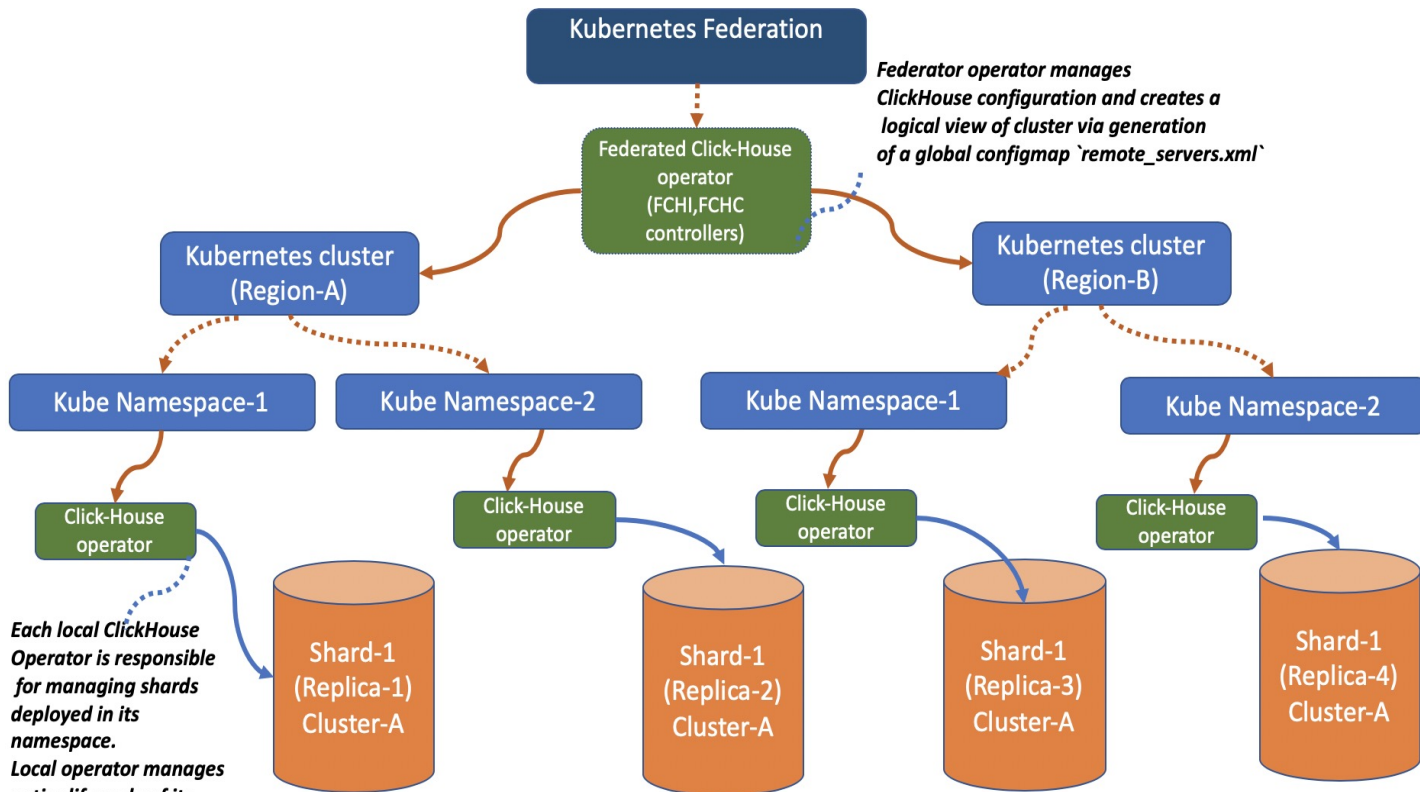
- Ingestion
 - 350K/sec - 4 VPU, 16GB
 - 3KB Span size – 20 Billion rows – 300GB
 - Avg Tag attribute length - 20
- Query patterns evaluated
 - Get Services/Operations
 - Get spans for given a service & Traceld
 - Get by Traceld
 - Search by Attributes (slowest) < 10 QPS

Why ClickHouse ?

- Open source & works well with structured data
- Works well for high ingestion use-cases
- Storage footprint is significantly lower
- Kubernetes operator support
- Extensive infrastructure monitoring in place
- Good inhouse ClickHouse expertise

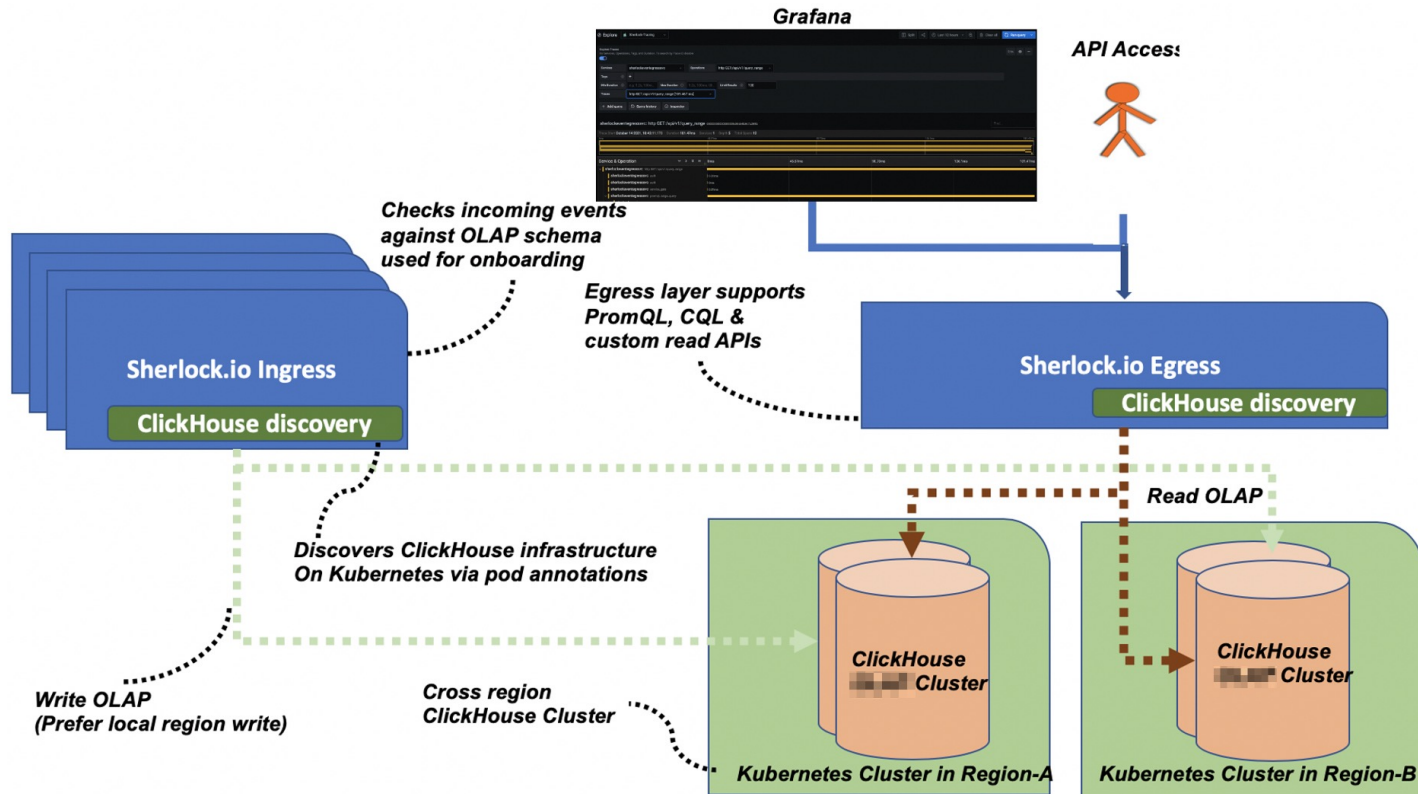


Deployment

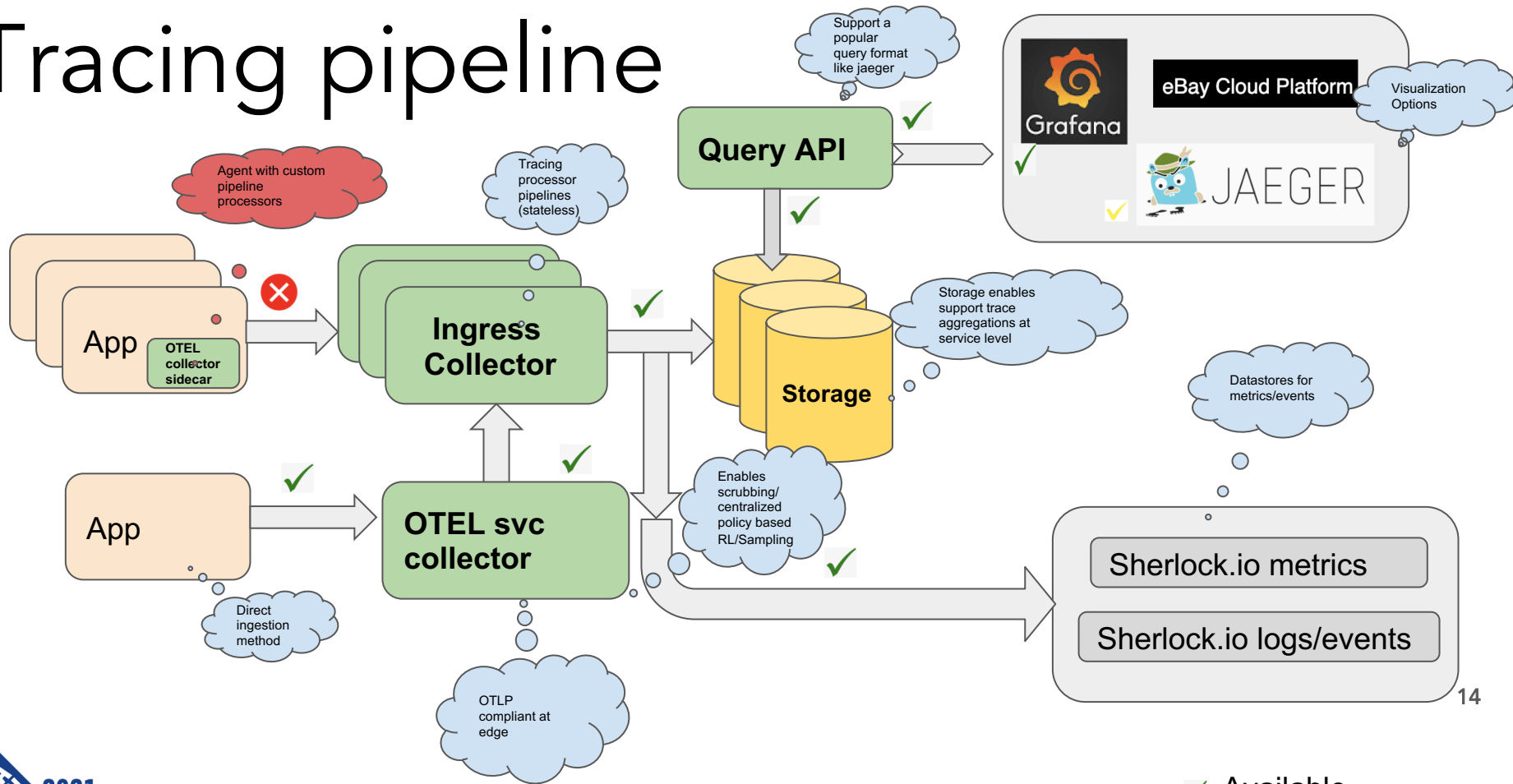


Each local ClickHouse Operator is responsible for managing shards deployed in its namespace. Local operator manages entire life cycle of its responsible shards

Deployment Contd..



Tracing pipeline



14

✓ Available
✓ Upcoming

ClickHouse Schema

Tables	Purpose	Engine
traces	<i>Fetch traces for a given service & op</i>	ReplicatedMergeTree
tracessvc	<i>Given a Trace, find all associated services</i>	ReplicatedReplacingMergeTree
svcop	<i>services to op lookup</i>	ReplicatedReplacingMergeTree

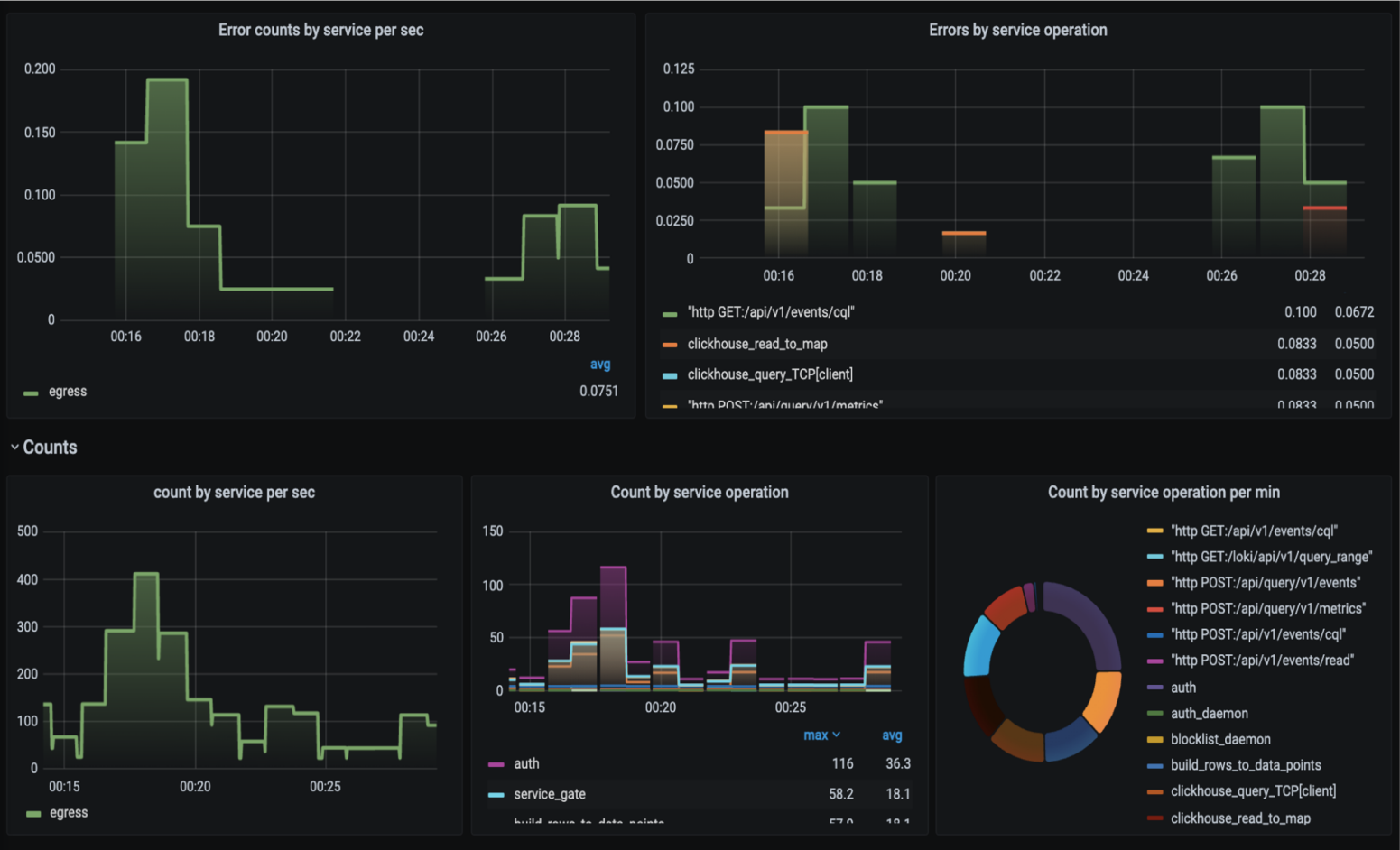
```
CREATE TABLE sherlockio.traces
(
  `svc_name` LowCardinality(String),
  `operation` String,
  `traceid` String,
  `spanid` String,
  `pspanid` String,
  `start_time` DateTime,
  `status` LowCardinality(String),
  `httpstatus` LowCardinality(String),
  `log` Array(String) CODEC(ZSTD(1)),
  `duration` UInt64,
  `offset` UInt64,
  `attr_key` Array(String),
  `attr_value` Array(String),
  `attr_type` Array(UInt8),
  INDEX traceid traceid TYPE bloom_filter GRANULARITY 5,
  INDEX operation operation TYPE bloom_filter GRANULARITY 5
)
ENGINE = ReplicatedMergeTree('/clickhouse/{installation}/{cluster}/tables/{shard}/sherlockio/traces', '{replica}')
PARTITION BY toYYYYMMDDhhmmss(toStartOfHour(start_time))
ORDER BY (svc_name, traceid, start_time)
SETTINGS index_granularity = 8192
```

name	compressed	uncompressed
svc_name	24.05 MiB	4.99 GiB
status	438.64 MiB	89.52 GiB
start_time	4.41 GiB	19.90 GiB
operation	7.17 GiB	185.39 GiB
log	13.25 GiB	206.65 GiB
offset	22.18 GiB	39.81 GiB
traceid	23.19 GiB	164.20 GiB
attr_type	25.07 GiB	125.24 GiB
duration	25.34 GiB	39.81 GiB
pspanid	52.65 GiB	72.03 GiB
attr_key	66.01 GiB	1.46 TiB
spanid	83.97 GiB	84.59 GiB
attr_value	194.59 GiB	7.64 TiB

RED Metrics from Traces

Auto
Instrumentation

- Requests
- Errors
- Duration
- HTTP status



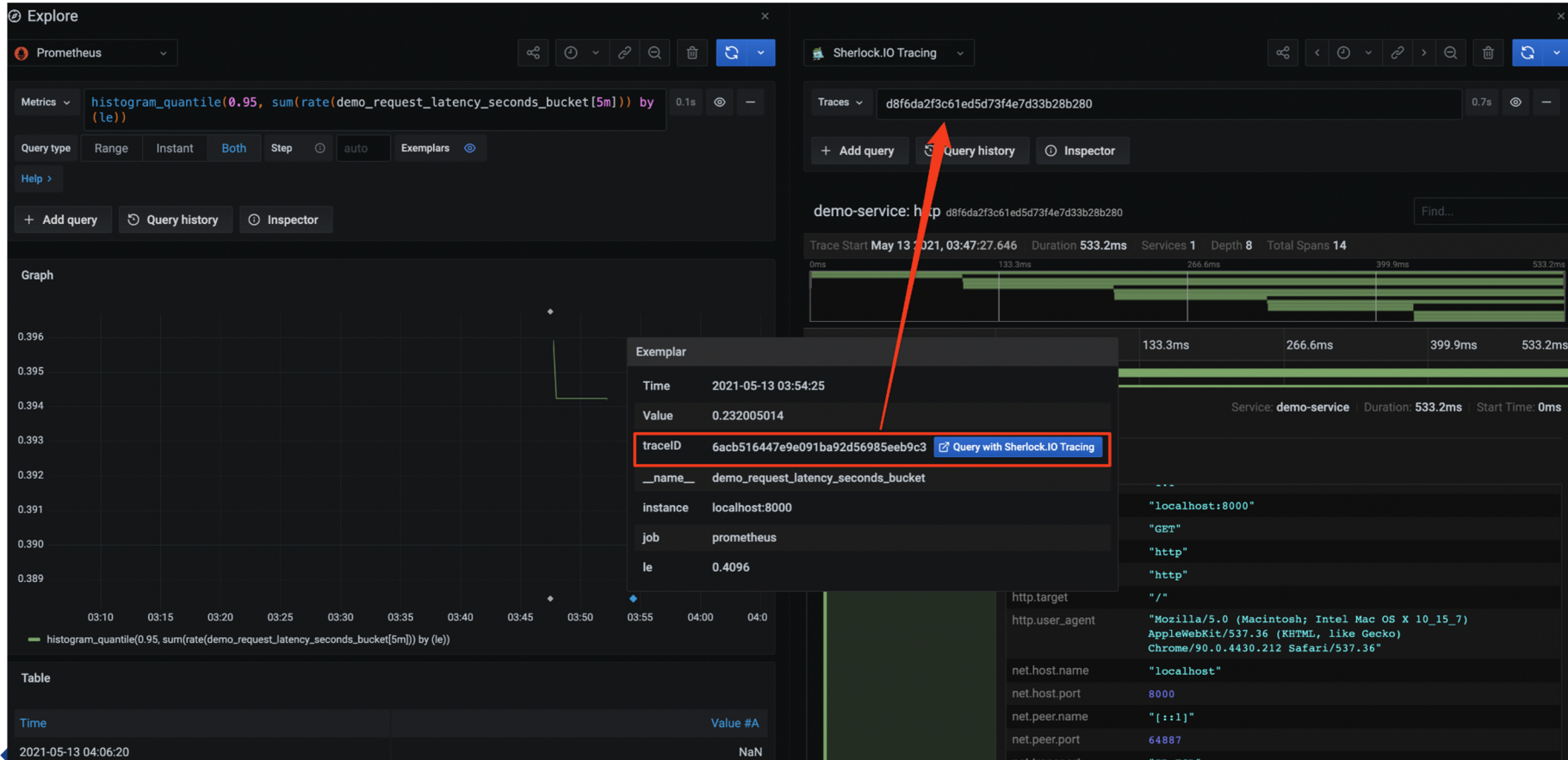
Sampling

- Decision on whether to process/export a span or not
- Tail based sampling
- Sampling based on
 - *Latency*
 - *Error*
 - *HTTP response status codes*

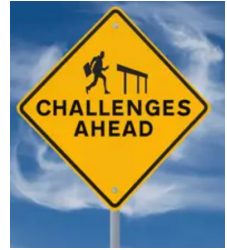




Metrics Traces



Challenges



- Adoption
 - *Frmk integration*
- Coverage
 - *Service Mesh*
- Sampling

eBay is Hiring

- ❖ Globally across all major cities and countries.
- ❖ Remote Work is an potential option in current times.
- ❖ Areas
 - RCGs
 - Research and Engineering
 - Data Science
 - Product Management
- ❖ Explore opportunities at <https://careers.ebayinc.com>
- ❖ Reach out to amvaidya@ebay.com or sudekumar@ebay.com



Thank You!