# 可观测数据融合处理平台探索
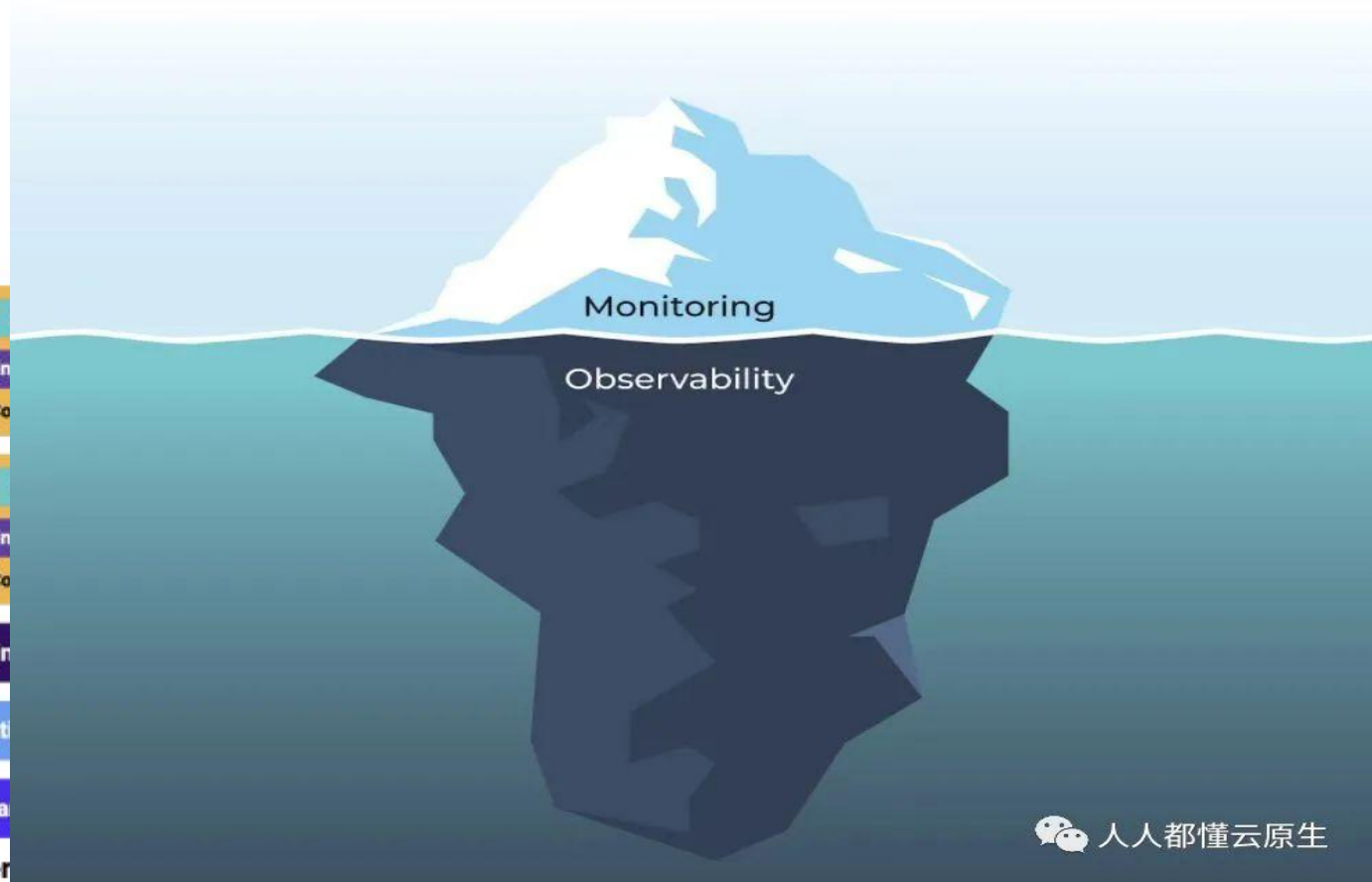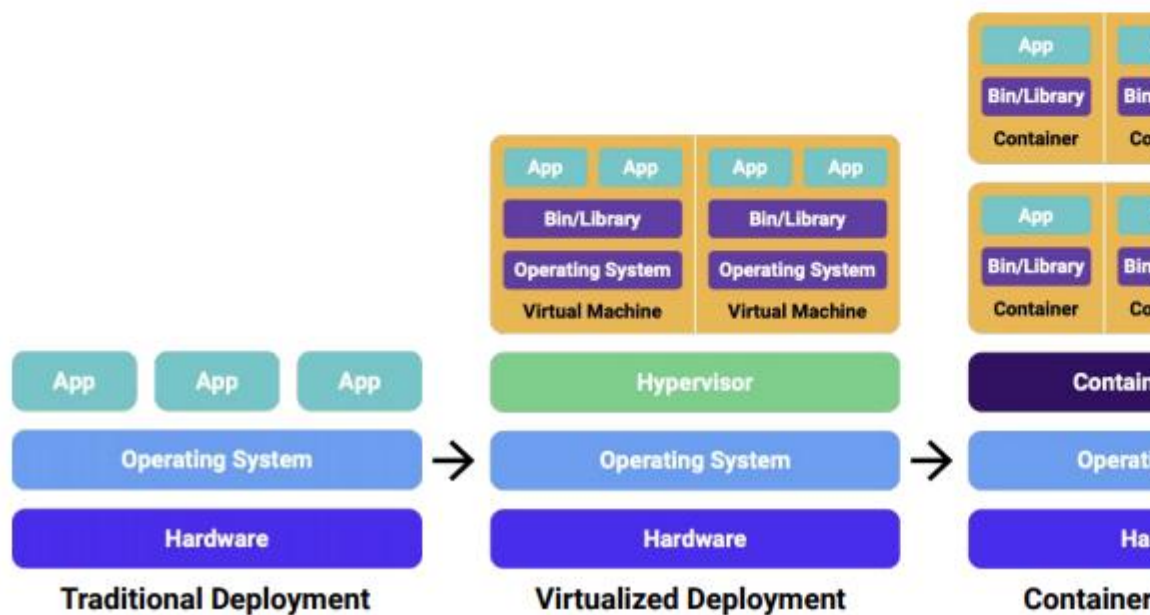
唐恒建

OPPO 高级后端工程师

- 1、可观测要素
- 2、opentelemetry现状
- 3、各类可观测方案对比
- 4、可观测数据融合平台
- 5、总结

# Increasingly complex software deployments

# 可观测要素
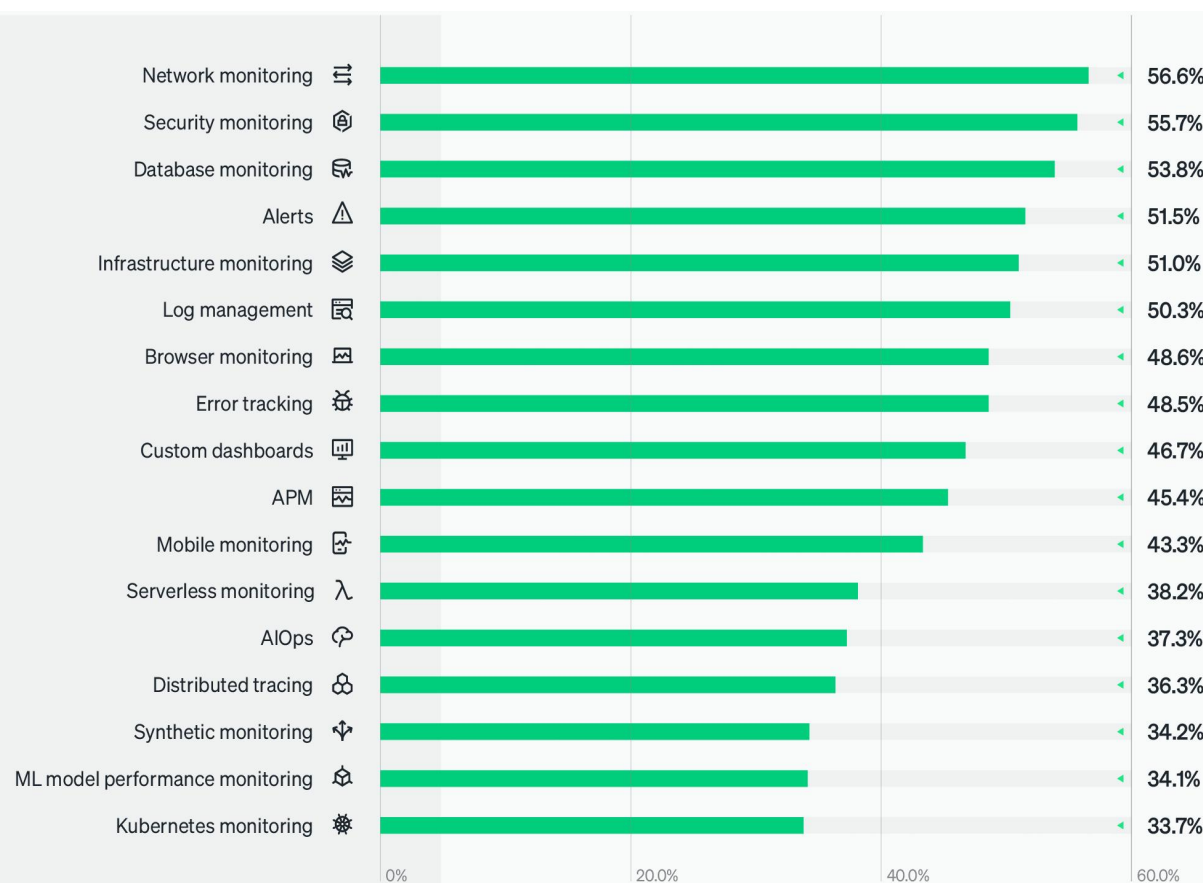


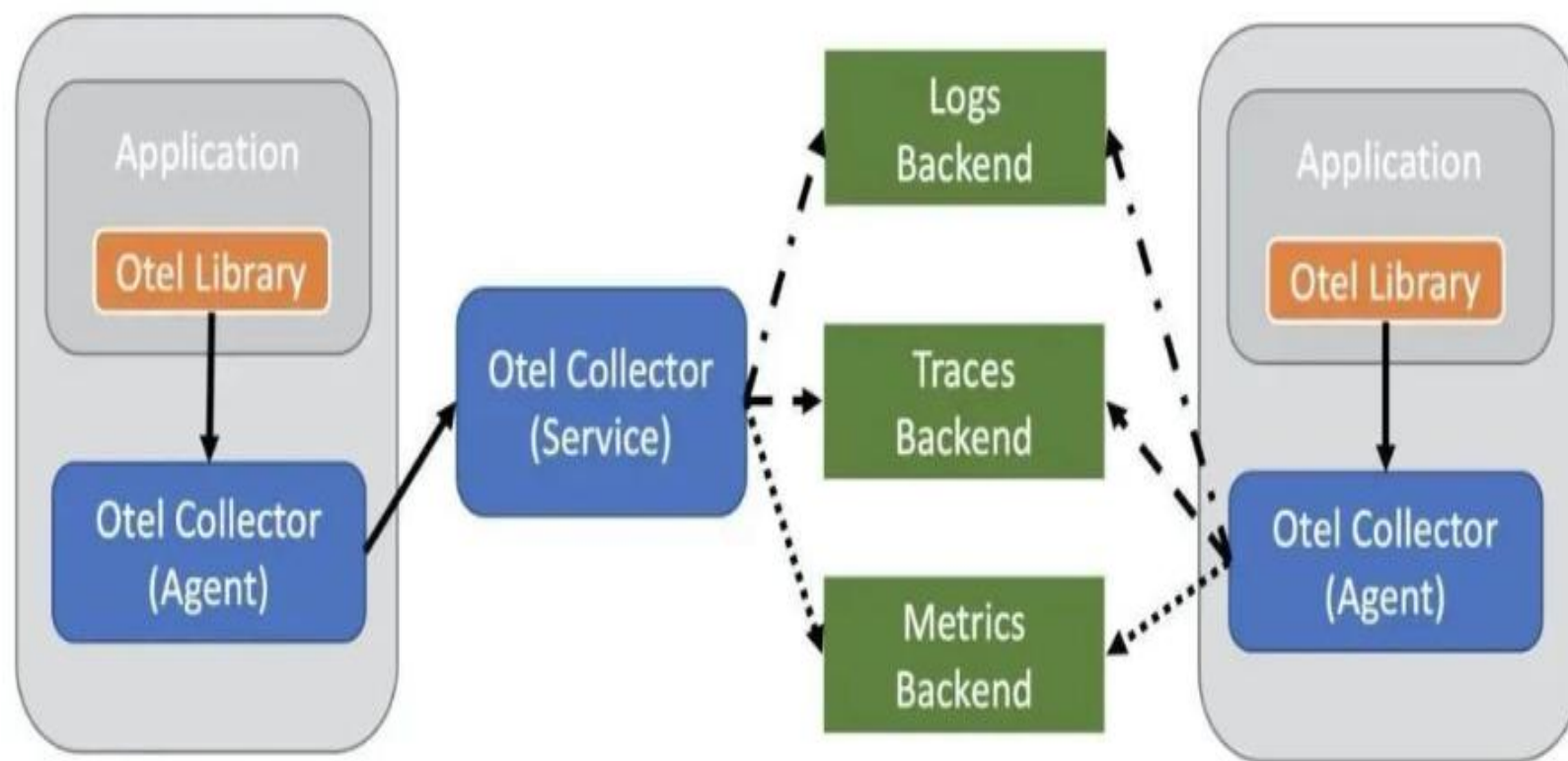| | |
|---|---|
| Network monitoring ⇄ | 56.6% |
| Security monitoring 🛡 | 55.7% |
| Database monitoring | 53.8% |
| Alerts ⚠ | 51.5% |
| Infrastructure monitoring | 51.0% |
| Log management | 50.3% |
| Browser monitoring | 48.6% |
| Error tracking | 48.5% |
| Custom dashboards | 46.7% |
| APM | 45.4% |
| Mobile monitoring | 43.3% |
| Serverless monitoring λ | 38.2% |
| AIOps | 37.3% |
| Distributed tracing | 36.3% |
| Synthetic monitoring | 34.2% |
| ML model performance monitoring | 34.1% |
| Kubernetes monitoring ⚙ | 33.7% |



Low volume — Request-scoped metrics — **Metrics** Aggregatable

**Tracing** Request scoped

Aggregatable events e.g. rollups

Request-scoped, aggregatable events

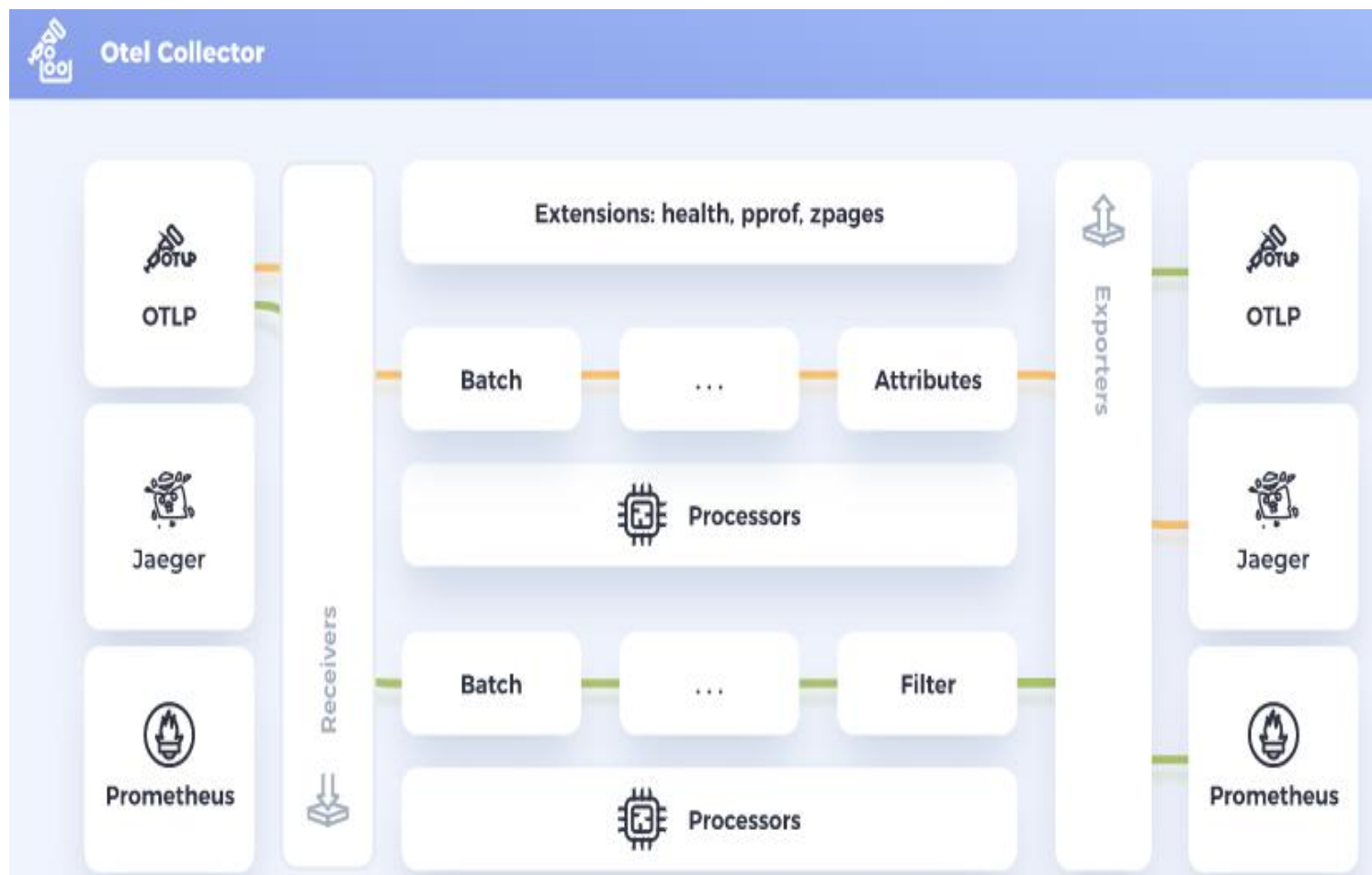**Logging** Events

High volume — Request-scoped events

# opentelemetry现状

- 维度：metric、trace、log
- 语言：JIT：Java、Rust、Python
  非JIT类：C++、Go

- agent、collect（batch）

- Multi protocol：receiver、exporter

# opentelemetry现状

```yaml
receivers:
  otlp:
    protocols:
      grpc:
      http:

  prometheus:
    config:
      scrape_configs:
        - job_name: 'app'
          scrape_interval: 10s
          static_configs:
            - targets: ['app:8080']

exporters:
  otlp:
    endpoint: ███████4317
    tls:
      insecure: true

  prometheusremotewrite:
    endpoint: http://██████.8080/api/v1/push
    tls:
      insecure: true
    headers:
      X-Scope-OrgID: demo
```
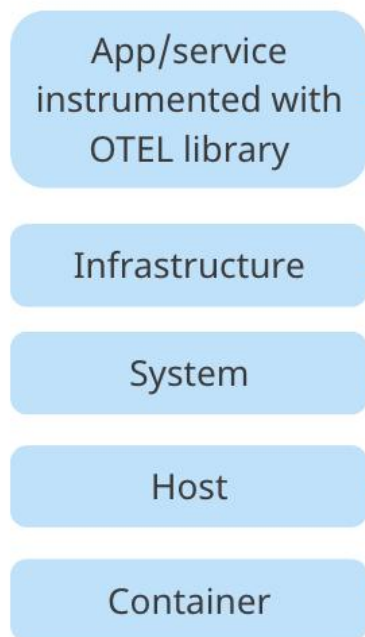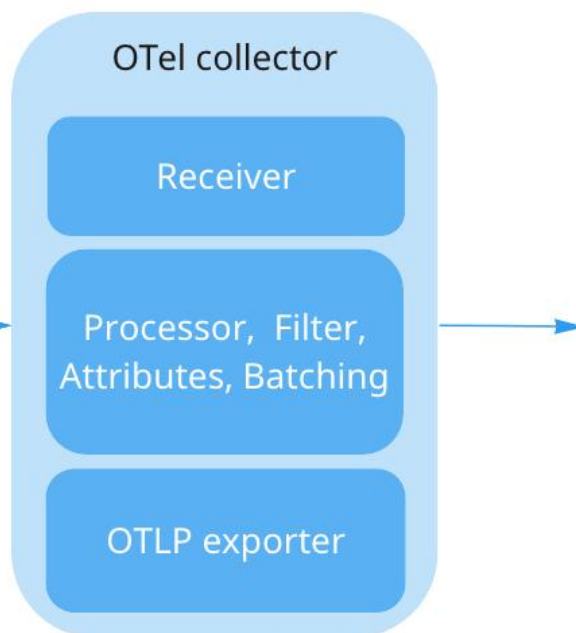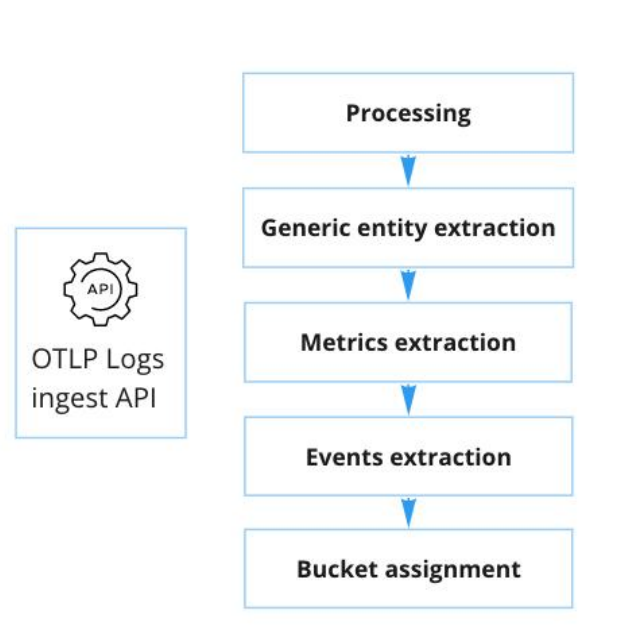
# 各类可观测方案对比

- Dynatrace

- Dynatrace

- Datadog



| | | | | | | | 2017 | 2018 | 2019 | 2020 | 2021... |
|---|---|---|---|---|---|---|---|---|---|---|---|

**2010** — Real-Time Unified Data Platform

**2012** — Infrastructure Monitoring Hosts / Clouds / VMs / Containers / Processes / IoT

**2017** — APM, Distributed Tracing

**2018** — Log Management, Logging without Limits™, Watchdog Alerts

**2019** — Serverless Monitoring, Tracing without Limits™, Synthetic Monitoring, Network Performance Monitoring

**2020** — Continuous Profiler, Deployment Tracking, Incident Management, Real User Monitoring, Error Tracking, Cloud SIEM, Mobile RUM

**2021...** — Session Replay, Network Device Monitoring, Cloud Security Posture Management, Cloud Workload Security, Database Monitoring, CI Visibility, Watchdog Root Cause Analysis, Watchdog Insights

**Founded Datadog to break down silos**

**Deployed everywhere, used by everyone**
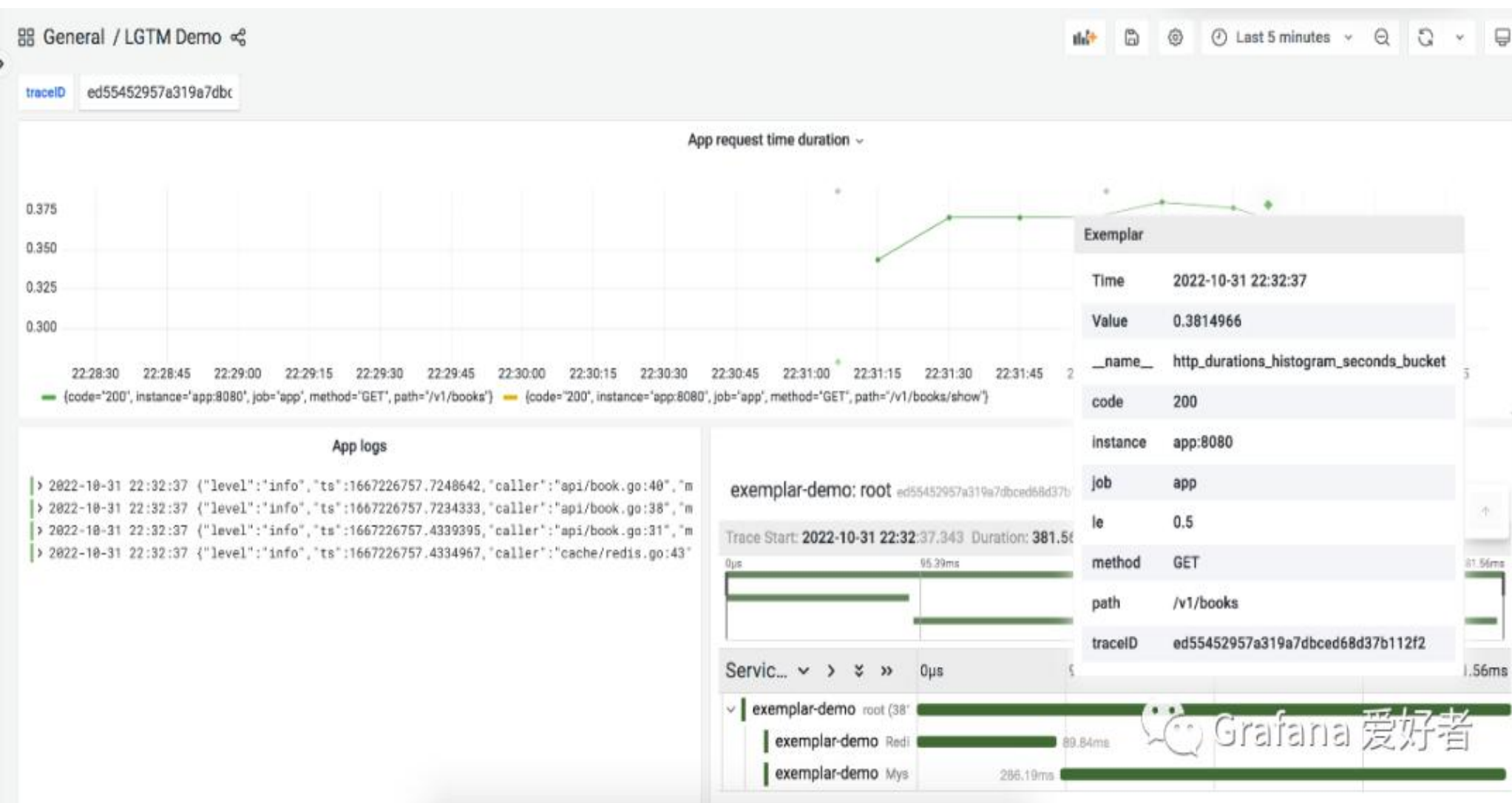
- Datadog

- Grafana

- Grafana



- 以exemplar为连接
- 联动metrics、log、trace
- Metrics–>logs
- metrics–>trace

# • Victoriametrics

**Clients**

vmselect fully supports PromQL and can be used as Prometheus datasource in Grafana

**Stateless**

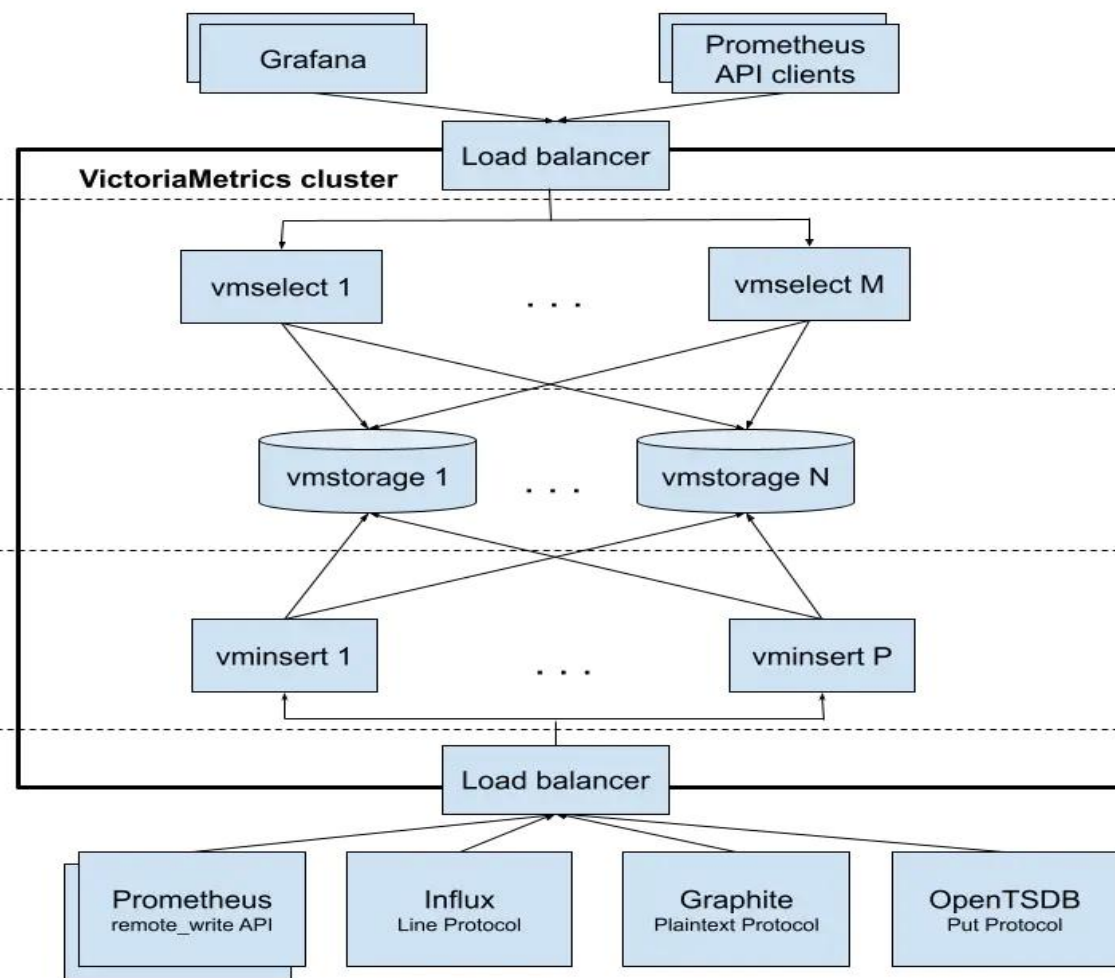vmselect fetches and merges data from vmstorage during queries

**Stateful**

vmstorage stores time series data

**Stateless**

vminsert spreads time series across available vmstorage nodes

**Writers**

Multiple Prometheus instances may write data to VictoriaMetrics cluster
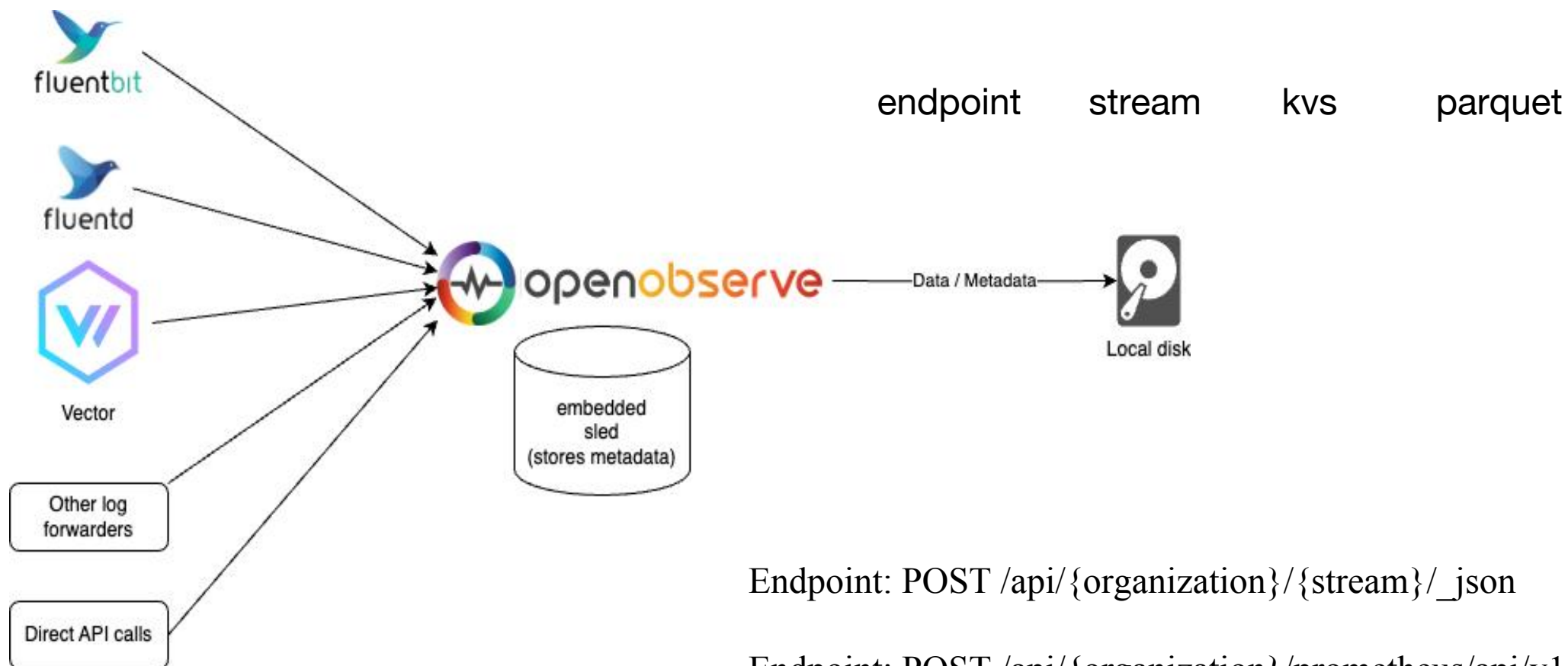There is support for other ingestion protocols

• Less CPU Usage

• Less Disk Usage

- Victoriametrics

VictoriaLogs --- Preview stage

- inspired by ClickHouse architecture
- It uses bloom filters
- encoding and compression for fields with different data types
- logs for the same log stream close to each other.
- maintains sparse index for log timestamps
- high cardinalitymetrics added.

- Openobserve
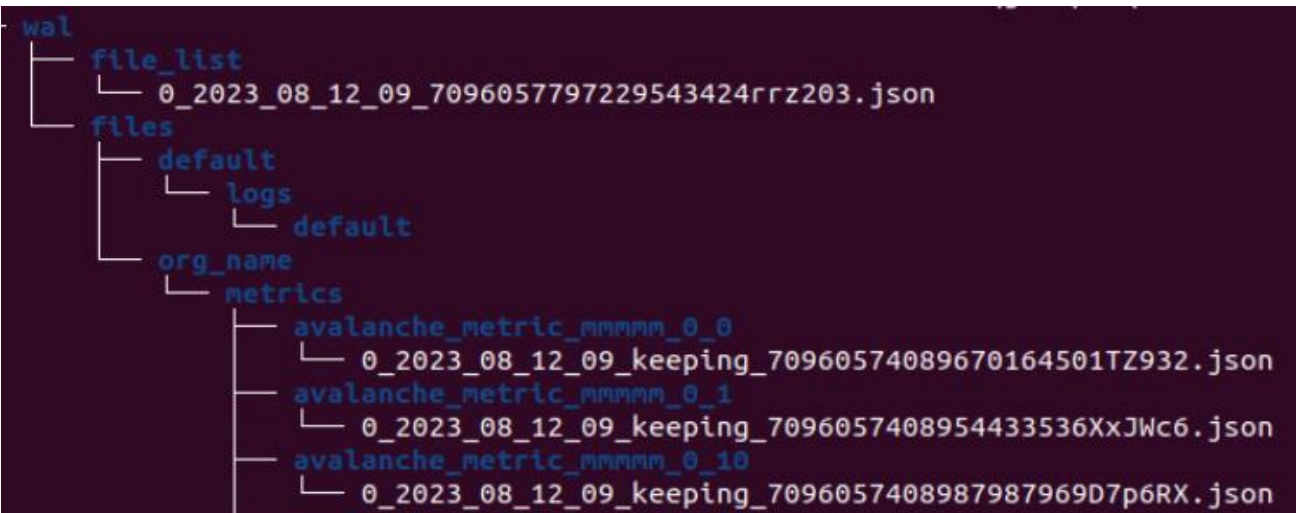


endpoint        stream        kvs        parquet

Endpoint: POST /api/{organization}/{stream}/_json

Endpoint: POST /api/{organization}/prometheus/api/v1/write

Endpoint: POST /api/{organization}/traces

# • Openobserve



WAL

blocks

- Openobserve

# 可观测数据融合平台

夯实metrics能力

项目启动

基于多副本的分布式存储、实现 scraper、ruler 高可靠、多租户。

基于对象存储的分级存储、集中压缩去重、自动降精度。

DB分片和并行压缩、乱序/历史数据写入、exempalr、多协议支持。

2019年　　　　　v1.0　　　　　v2.0　　　　　v3.0

抽象平台层

分布式、高可用、高基数

数据融合平台

扩展log、trace底座能力,作为数据存储层

# 夯实metrics能力

## v1.0 - 基于多副本实现存储高可靠

# v1.0 – 小结

实现：

1. 基于 Prometheus TSDB 构建多租户的分布式存储服务。

2. 租户之间可以配置不同的读写QPS 限制和存储周期。

3. 分布式 scraper 和 ruler 实现数据抓取和 rule 评估的高可靠。

不足：

1. 长时间存储集群运维压力大，扩容和坏盘修复较复杂。

2. 读写没有分离，相互影响，尤其长时间范围查询。

3. 本地存储多副本，对于冷数据（一个月前）没有进行去重，造成存储浪费。

# v2.0 - 基于对象存储的分级存储

# v2.0 - 集中压缩和自动降精度

# v2.0 – 小结

实现：

1. 集中压缩，数据去重，降低成本。
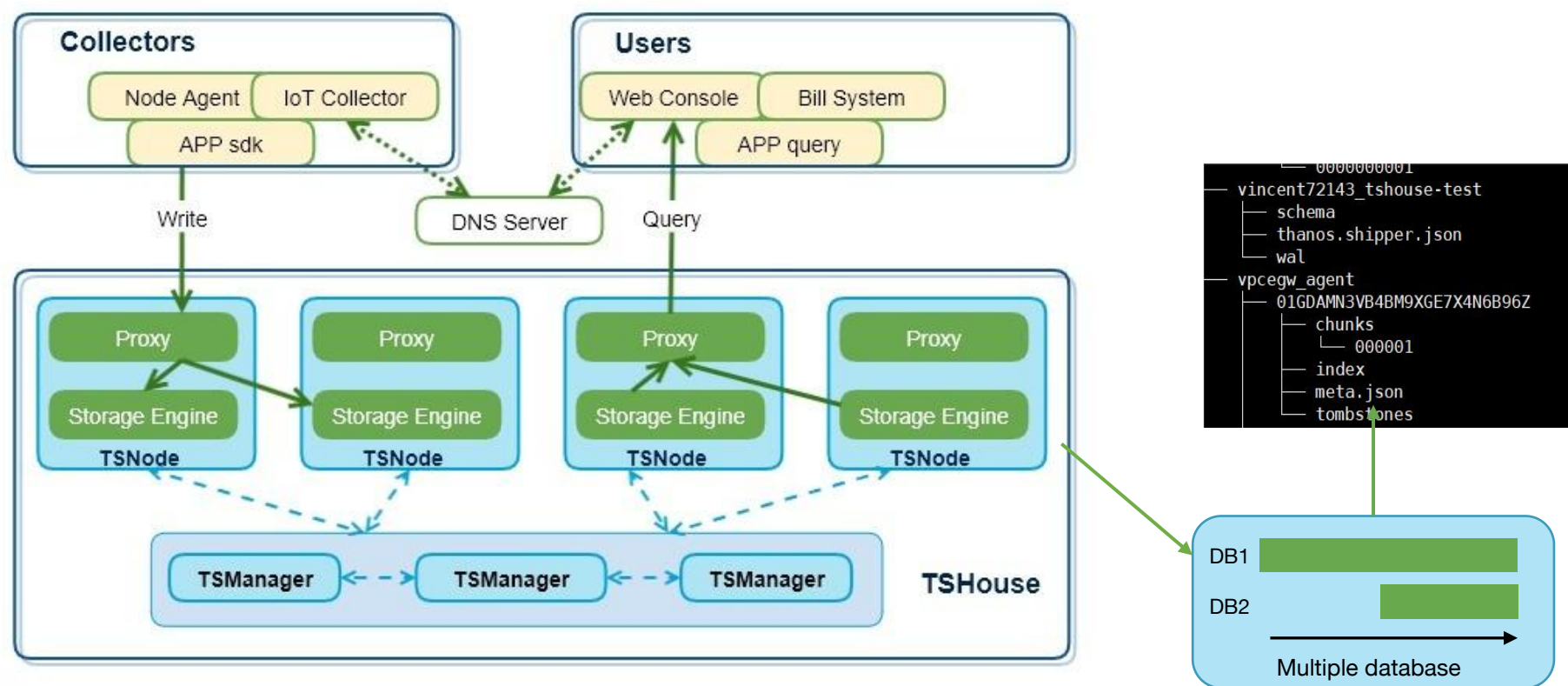
2. 自动降精度，长时间范围查询自动转化，提升查询效率。

3. 依赖对象存储实现长时间数据集中存储。

不足：

1. 超大规模数据压缩时间长，占用磁盘和内存大，受限 TSDB 索引中 symbols 长度 64GB 限制，高 level 压缩会导致失败。

2. 历史数据和乱序数据写入支持较差。

3. 仅支持 OpenMetrics/Prometheus 数据格式写入。

# v3.0 - 高基数和大规模数据存储挑战



```
523     }
524
525     func (w *Writer) finishSymbols() error {
526  +      symbolTableSize := w.f.pos - w.toc.Symbols - 4
527  +      // The symbol table's <len> part is 4 bytes. So the total symbol table size must be less than or equal to 2^32-1
528  +      if symbolTableSize > 4294967295 {
529  +          return errors.Errorf("symbol table size exceeds 4 bytes: %d", symbolTableSize)
530  +      }
531  +
532        // Write out the length and symbol count.
533        w.buf1.Reset()
534  +      w.buf1.PutBE32int(int(symbolTableSize))
535        w.buf1.PutBE32int(int(w.numSymbols))
536        if err := w.writeAt(w.buf1.Get(), w.toc.Symbols); err != nil {
537            return err
```
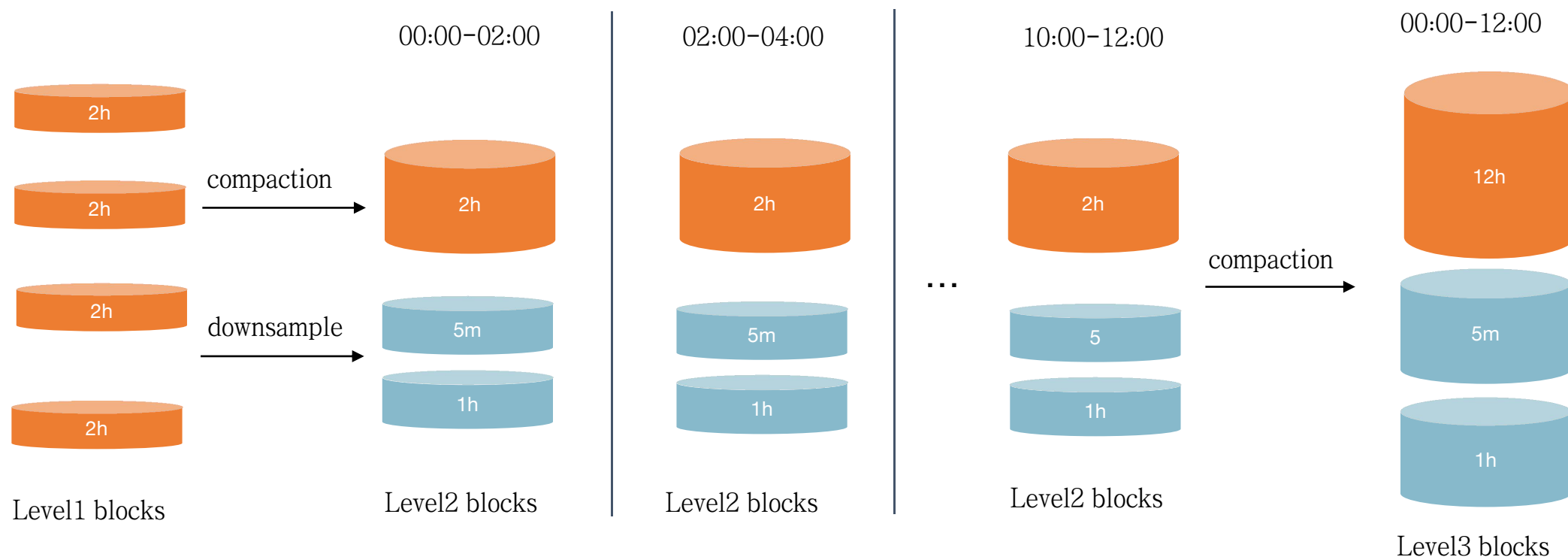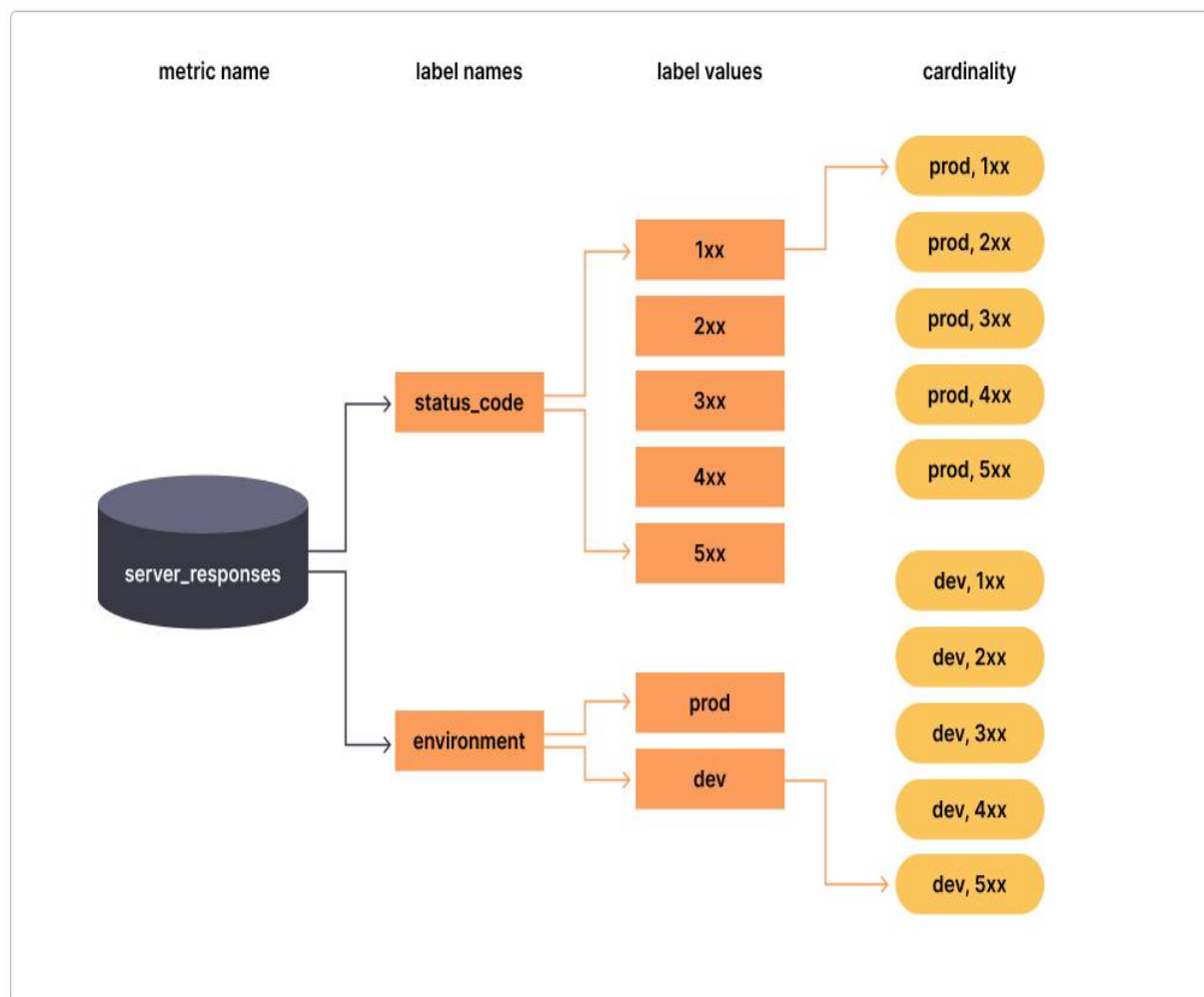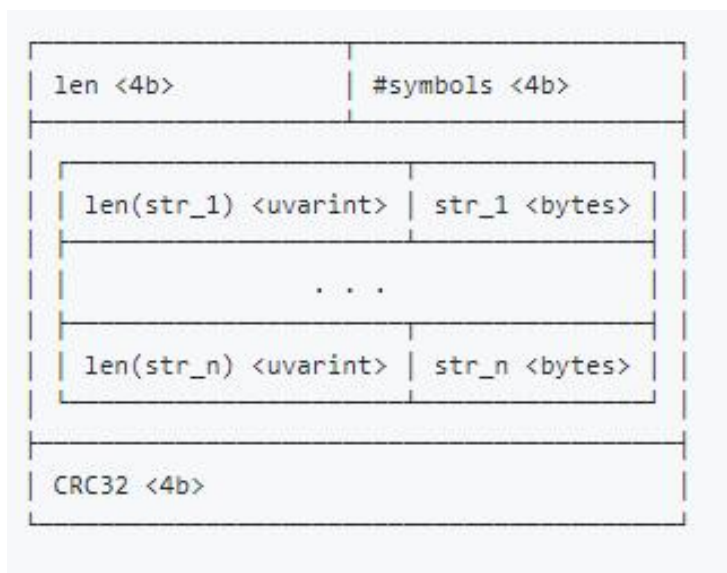
https://github.com/prometheus/prometheus/pull/9104



Level 5 compaction

# v3.0 - 高基数和大规模数据存储挑战



```
- sourceLabels: [__name__]
  action: drop
  regex: 'node_(nf_conntrack_statl
```

# v3.0 - DB 分片与并行压缩

# v3.0 - DB 分片与并行压缩

# v3.0 - DB 分片与并行压缩

# v3.0 - DB 分片与并行压缩

Path:

Name
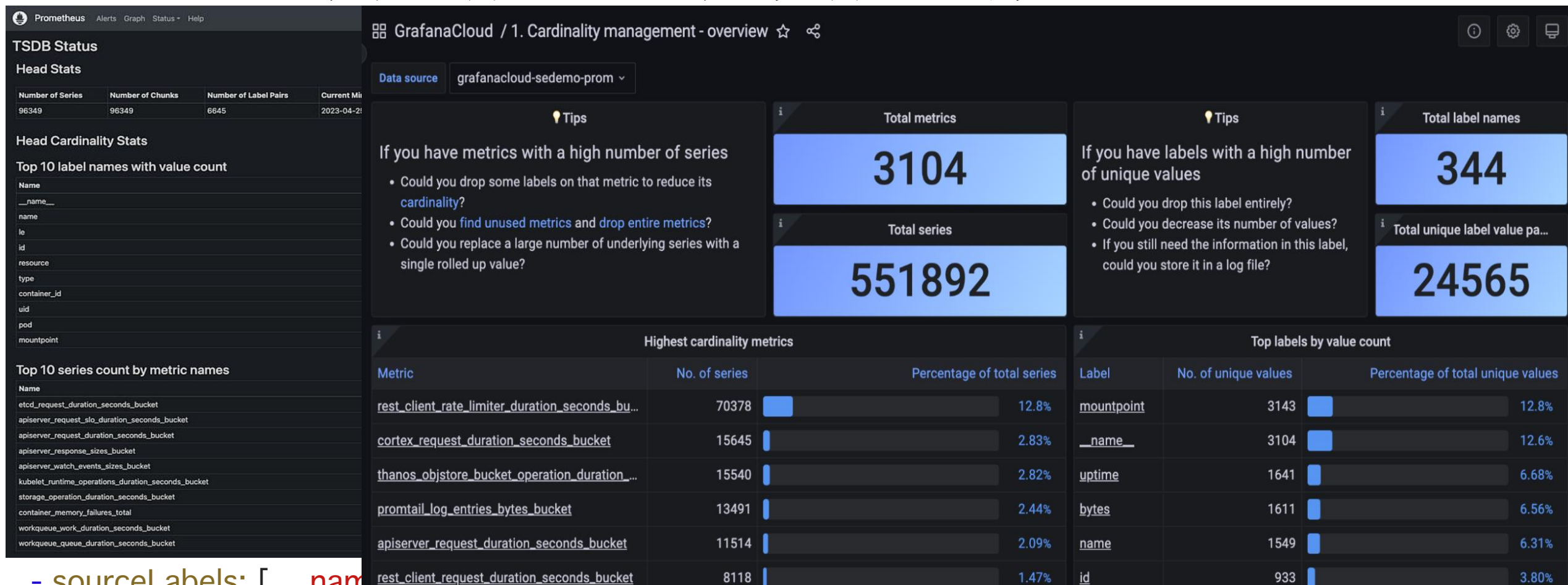
📁 01H6ZNQJ2AP1XJXKTFH38GFYSB/
📁 01H70017RX6BDD5NHP78V3FCJP/
📁 01H7038RZ7WEZFM8Z6YK2NB8KR/

Path:

Name

📁 01H70A4G6YQAMQRQBV4EJCFWE7/
📁 01H70H07J3FEJ2FEPYEGP0WCSG/
📁 01H70QVYPZ45JMN6THVWMV21YT/
📁 01H70YQNZ0YWWK9KE4JCTTD568/
📁 01H715KDD8YV7P58GYMW40M0GN/
📁 01H71CF4F1D97M586661DM6QHQ/

"thanos": {
        "labels": {
    "receive_cluster": "eu1",
    "receive_replica": "0",
    "tenant_id": "ad",
    "tenant-split": "2",
        },

"thanos": {
        "labels": {
    "receive_cluster": "eu1",
    "receive_replica": "0",
    "tenant_id": "ad",
        "tenant-shard": "1",
        "tenant-split": "2",
        },

# v3.0 - OpenTelemetry 协议支持



OpenTelemetry Metrics → Prometheus Metrics

- Gauge
- Sum(cumulative, non-monotonic) → Gauge
- Sum(cumulative, monotonic)
- Sum(delta, monotonic) → Counter (convert to cumulative)
- Histogram(cumulative)
- Histogram(delta) → Histogram (convert to cumulative)
- Summary → Summary

```
Metric #4
Descriptor:
     -> Name: http_durations_histogram_seconds
     -> Description: Http latency distributions.
     -> Unit:
     -> DataType: Histogram
     -> AggregationTemporality: Cumulative
HistogramDataPoints #0
Data point attributes:
     -> code: Str(200)
     -> method: Str(GET)
     -> path: Str(/v1/books/show)
StartTimestamp: 2022-10-19 14:38:52.634 +0000 UTC
Timestamp: 2022-10-19 14:39:22.597 +0000 UTC
Count: 4938
Sum: 94.990601
ExplicitBounds #0: 0.050000
ExplicitBounds #1: 0.100000
ExplicitBounds #2: 0.250000
ExplicitBounds #3: 0.500000
ExplicitBounds #4: 1.000000
ExplicitBounds #5: 2.000000
Buckets #0, Count: 4845
Buckets #1, Count: 26
Buckets #2, Count: 37
Buckets #3, Count: 30
Buckets #4, Count: 0
Buckets #5, Count: 0
Buckets #6, Count: 0


# HELP http_durations_histogram_seconds Http latency distributions.
# TYPE http_durations_histogram_seconds histogram
http_durations_histogram_seconds_bucket{code="200",method="GET",path="/v1/books/show",le="0.05"} 4845
http_durations_histogram_seconds_bucket{code="200",method="GET",path="/v1/books/show",le="0.1"} 4871
http_durations_histogram_seconds_bucket{code="200",method="GET",path="/v1/books/show",le="0.25"} 4908
http_durations_histogram_seconds_bucket{code="200",method="GET",path="/v1/books/show",le="0.5"} 4938
http_durations_histogram_seconds_bucket{code="200",method="GET",path="/v1/books/show",le="1"} 4938
http_durations_histogram_seconds_bucket{code="200",method="GET",path="/v1/books/show",le="2"} 4938
http_durations_histogram_seconds_bucket{code="200",method="GET",path="/v1/books/show",le="+Inf"} 4938
http_durations_histogram_seconds_sum{code="200",method="GET",path="/v1/books/show"} 94.9906009000003
http_durations_histogram_seconds_count{code="200",method="GET",path="/v1/books/show"} 4938
```
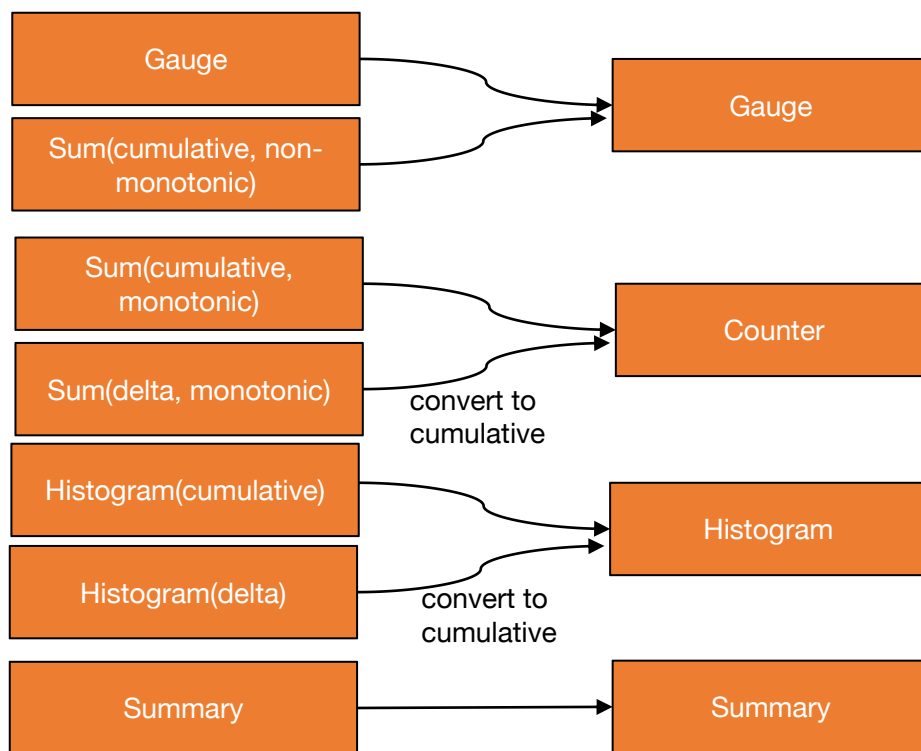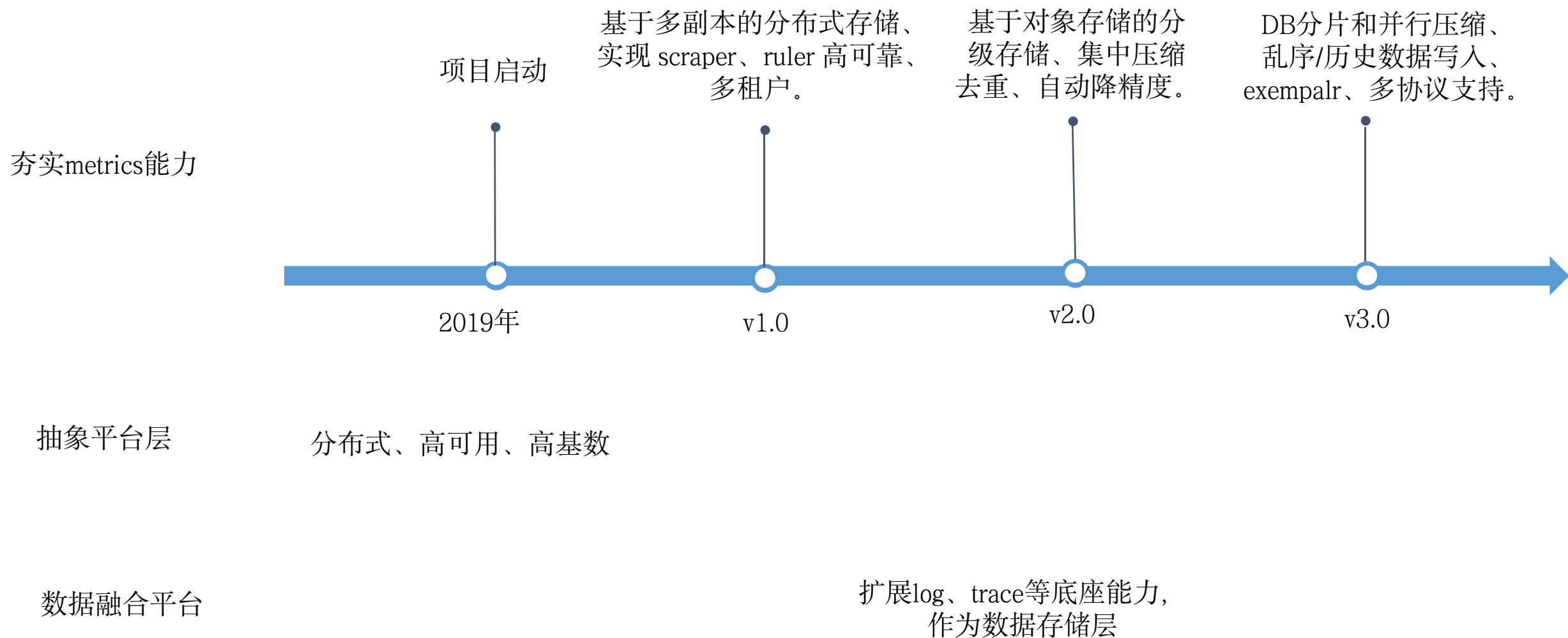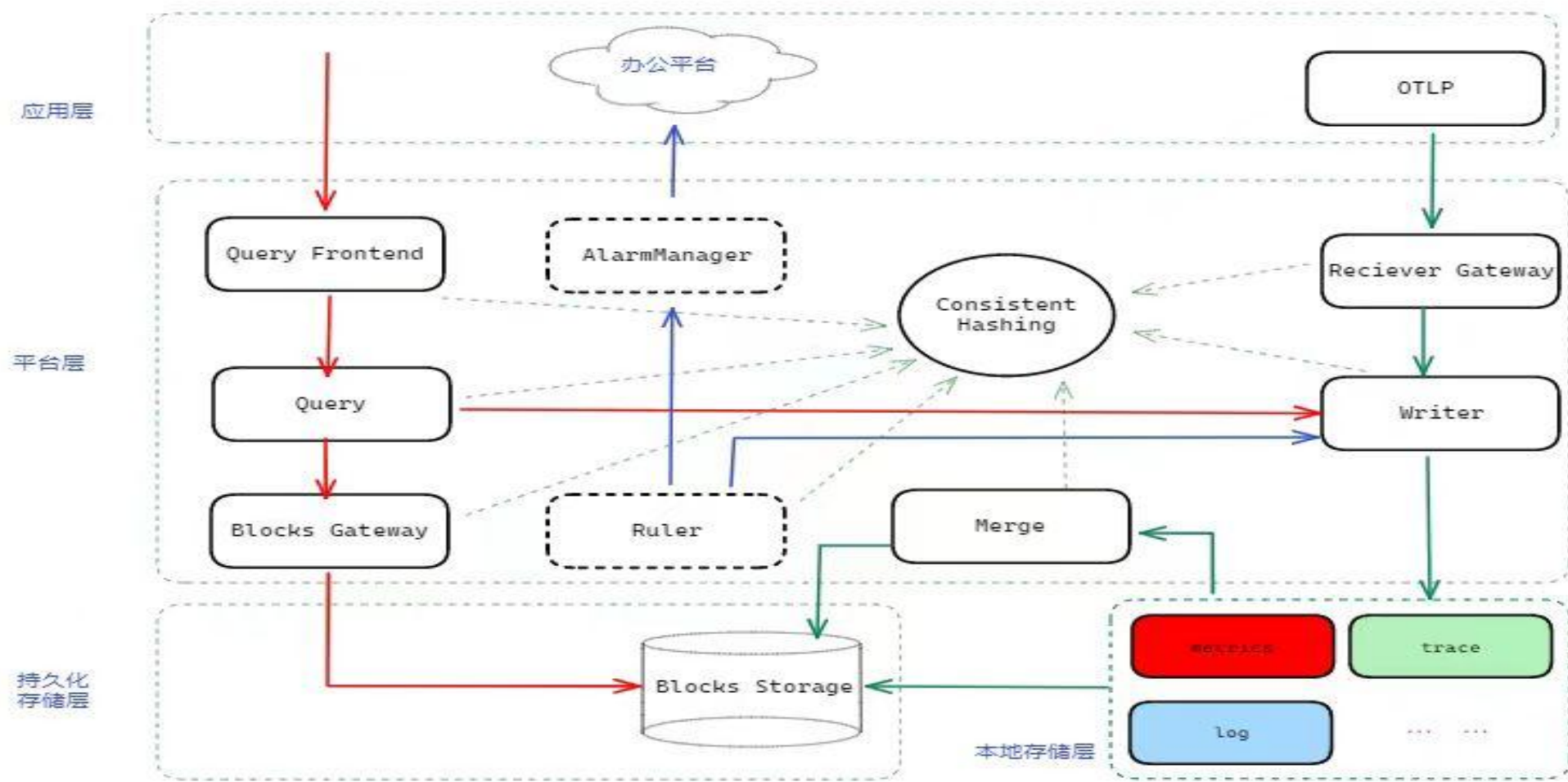
# 可观测数据融合平台

基于多副本的分布式存储、实现 scraper、ruler 高可靠、多租户。

基于对象存储的分级存储、集中压缩去重、自动降精度。

DB分片和并行压缩、乱序/历史数据写入、exempalr、多协议支持。

项目启动

夯实metrics能力

2019年      v1.0      v2.0      v3.0

抽象平台层

分布式、高可用、高基数

数据融合平台

扩展log、trace等底座能力，作为数据存储层

# 可观测数据融合平台
抽象平台层

# 可观测数据融合平台

扩展log、trace能力

## TSDB format

XOR chunk data

- Index
- Chunks
- Head Chunks
- Tombstones
- Wal
- Memory Snapshot

| num_samples <uint16> | ts_0 <varint> | v_0 <float64> | ts_1_delta <uvarint> | v_1_xor <varbit_xor> | ts_2_dod <varbit_ts> |
|---|---|---|---|---|---|

log chunk data

| num_samples <uint16> | ts_0 <int64> | data_len <int64> | log_data <varint> | ts_1 <int64> |
|---|---|---|---|---|

# 扩展log 能力

```go
func (a *xorAppender) Append(t int64, v float64) {
    var tDelta uint64
    num := binary.BigEndian.Uint16(a.b.bytes())
    switch num {
    case 0:
        buf := make([]byte, binary.MaxVarintLen64)
        for _, b := range buf[:binary.PutVarint(buf, t)] {
            a.b.writeByte(b)
        }
        a.b.writeBits(math.Float64bits(v), nbits: 64)
    case 1:
        tDelta = uint64(t - a.t)

        buf := make([]byte, binary.MaxVarintLen64)
        for _, b := range buf[:binary.PutUvarint(buf, tDelta)] {
            a.b.writeByte(b)
        }

        a.writeVDelta(v)

    default:
        tDelta = uint64(t - a.t)
```

```go
func (a *logAppender) AppendLog(t int64, v string) {

    /**
    chunk format:
    t int64;datalen int64;logdata []byte;// default lz4 compression
    t int64;datalen int64;logdata []byte;
    ...
     **/

    // 1. write time section
    buf := make([]byte, 8)
    binary.BigEndian.PutUint64(buf, uint64(t))
    for _, b := range buf {
        a.b.writeByte(b)
    }


    // 2. write datalen section
    lenbuf := make([]byte, 8)
    fmt.Println( a...: "en(v)",len(v))
    binary.BigEndian.PutUint64(lenbuf, uint64(len(v)))
    for _, b := range lenbuf {
        a.b.writeByte(b)
    }
    fmt.Println( a...: "2--a.b.stream:",a.b.stream)
```
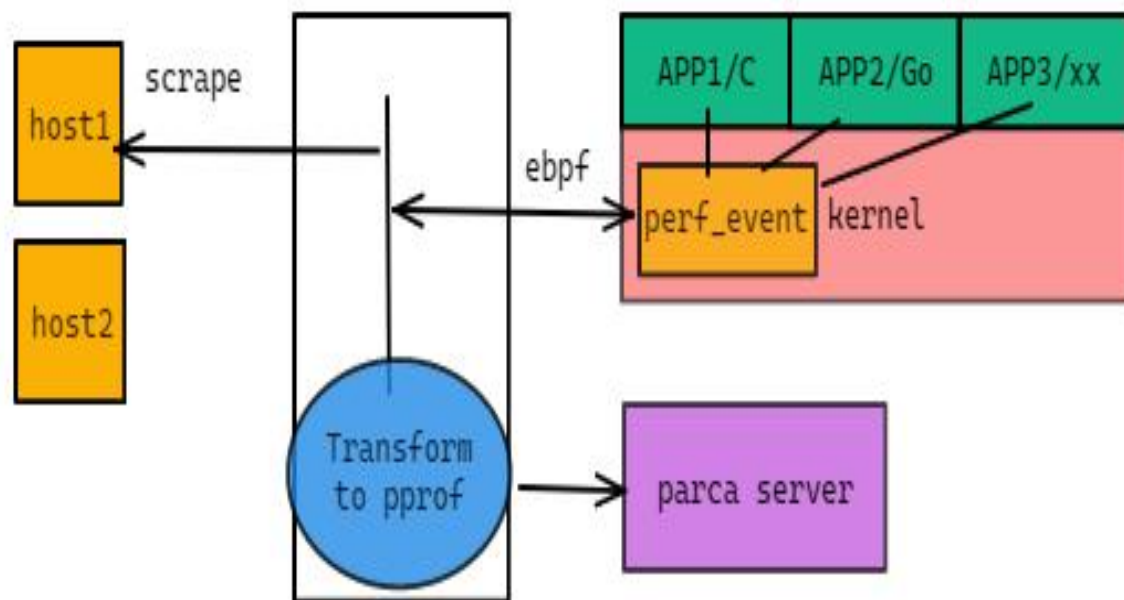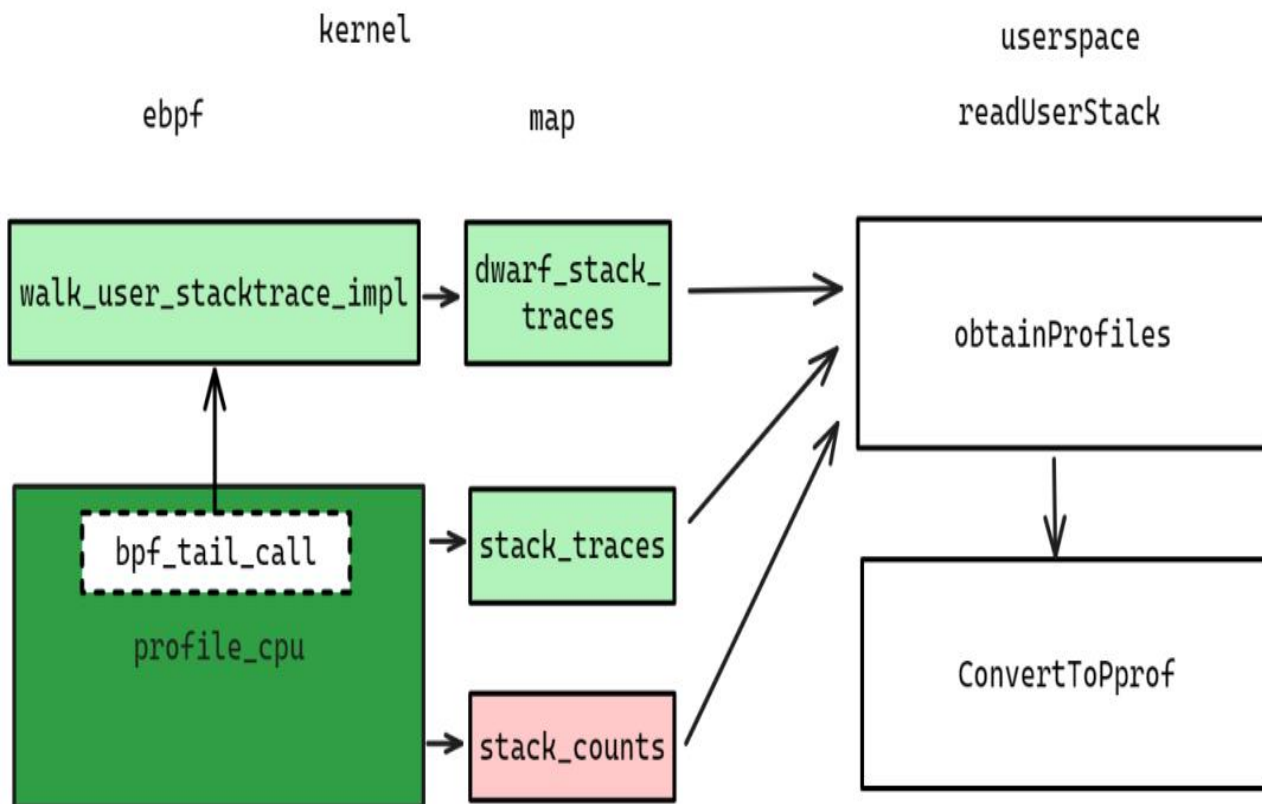
# eBPF

Parca





[bpf profile_cpu] "has_unwind_information" works rather than "has_fp" even though a porgrame(C) compiled with frame pointers
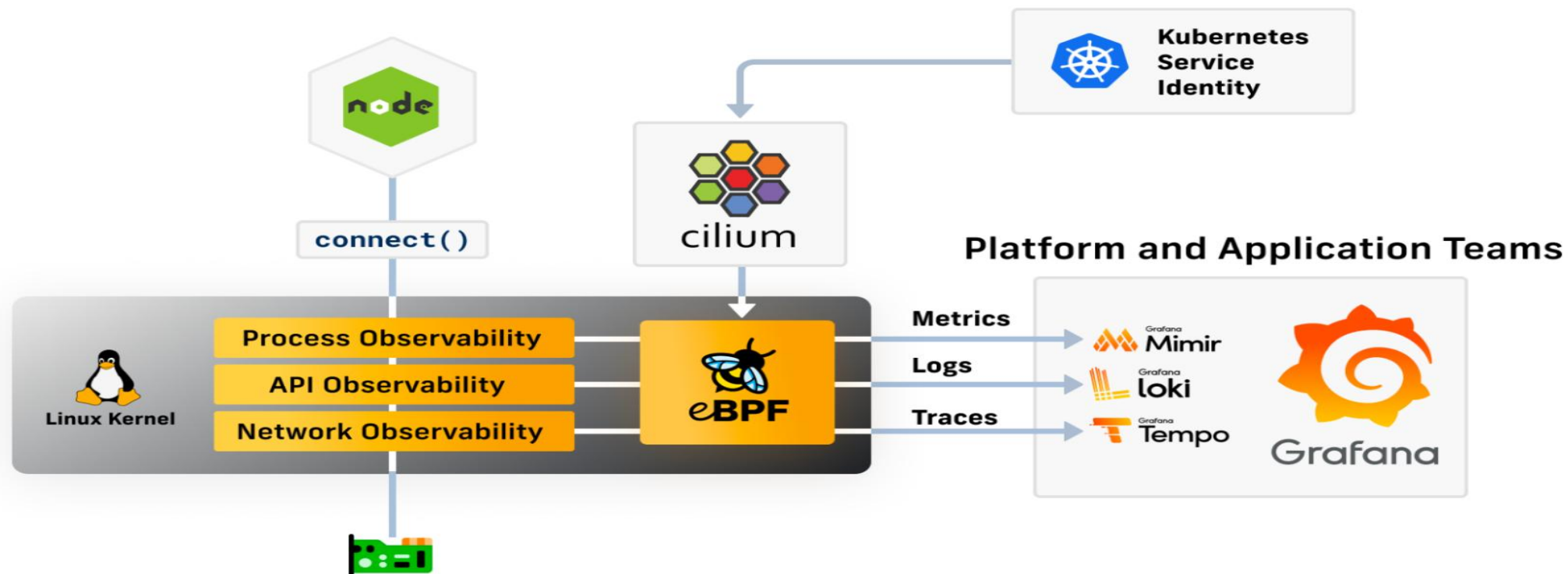https://github.com/parca-dev/parca-agent/issues/1657

OpenTelemetry Auto Instrumentation using eBPF
https://github.com/open-telemetry/opentelemetry-go-instrumentation/pull/149

# eBPF

# 总结

- 1、各类方案百花齐放，多调研和灰度验证
- 2、平台聚焦服务业务为主
- 3、找到适合业务的技术方案
- 4、慎重叠加解决方案
- 5、尽量抽象平台能力，扩展底层能力

谢谢

欢迎关注公众号~



Grafana fans

欢迎关注公众号~