

The logo consists of a central teal circle with three lines extending from it to three smaller teal circles, one above and two below, forming a stylized network or data flow icon.

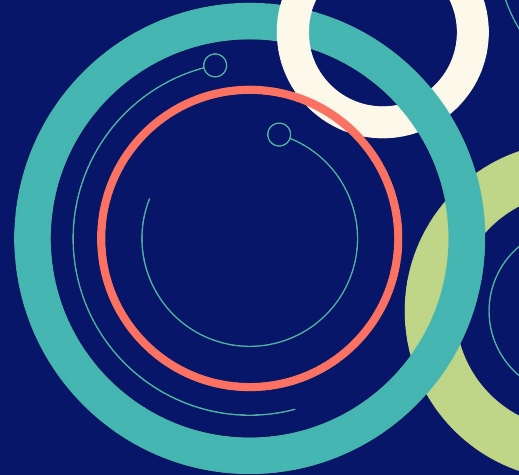
Observability Summit North America

The Speed of Metrics, The Fidelity of Traces:

Architecting Post Collection Aggregation

Zack Owens

Principal Software Engineer, New Relic



Telemetry Fidelity is a Choice

What to measure is essential to observability practice

When you're in an incident or want to do long range analysis, these choices matter!

Today's problem: the choice you make at collection dictates what questions your observability practice can answer



Metrics

Cost effective to store

(Usually) fast to query or alert

Problem: lacks the dimensions necessary to answer any observability question

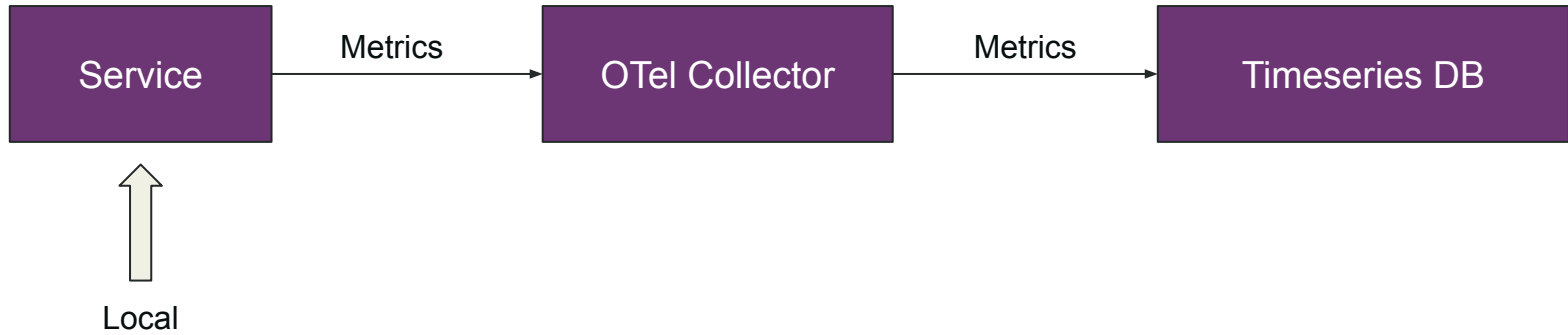
“Raw Data” (Events/Logs/Spans/Traces)

Rich context to answer ad-hoc questions

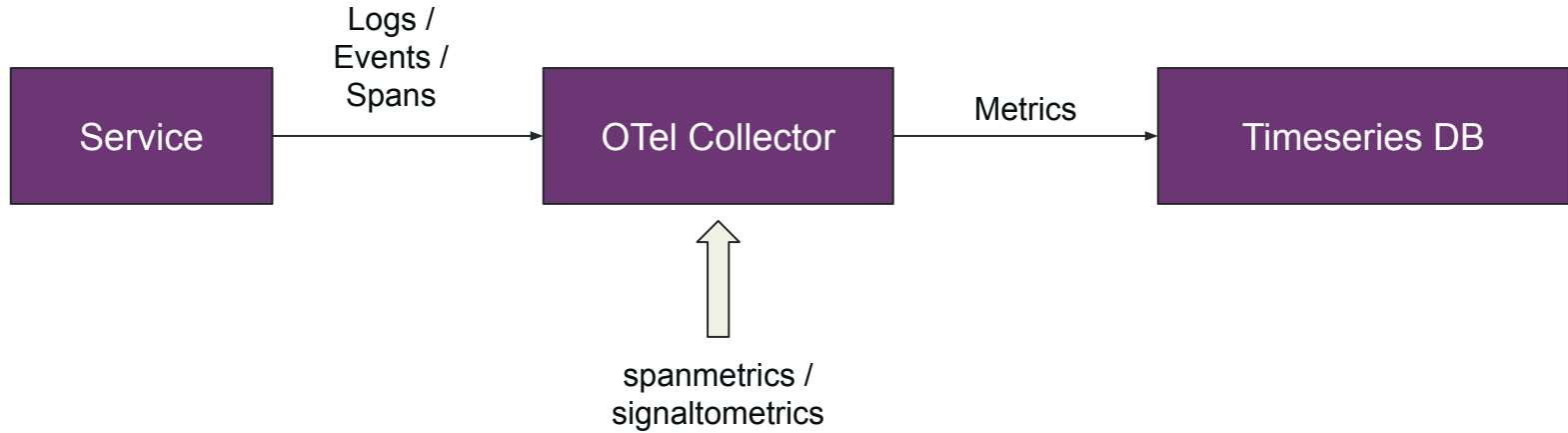
Problem: SLOW to query with large volumes of data



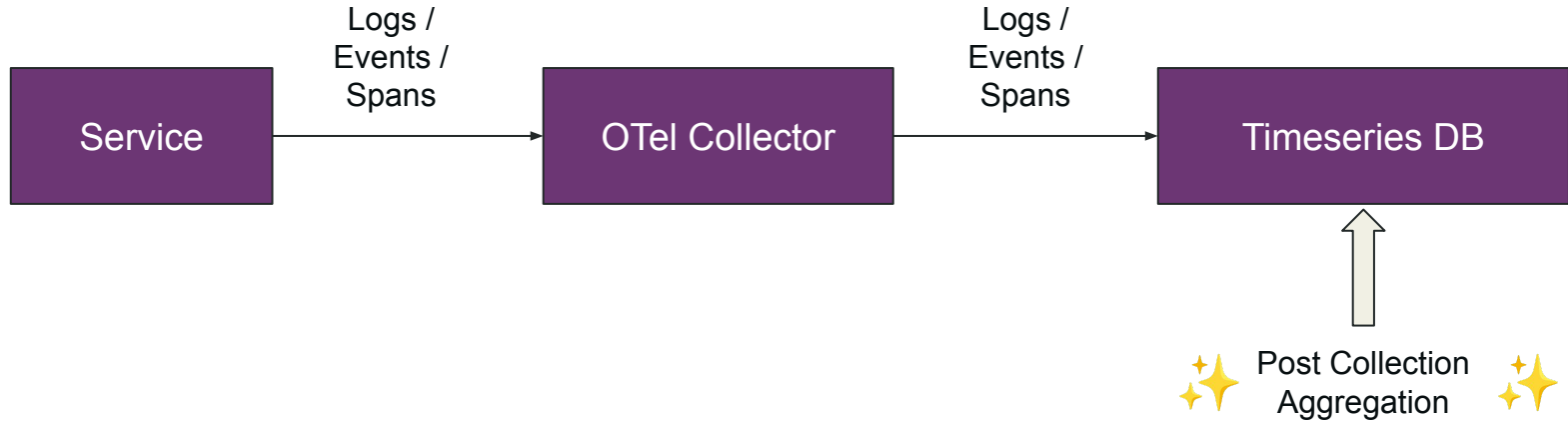
Where do we aggregate?



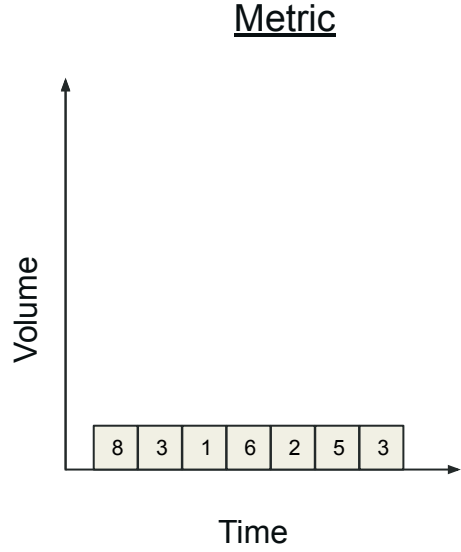
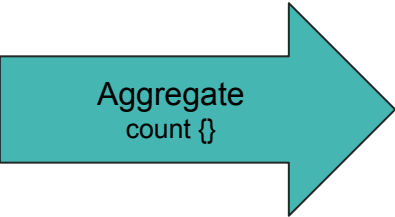
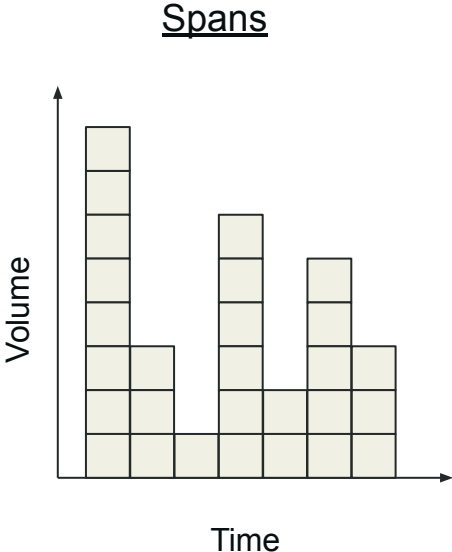
Where do we aggregate?



Where do we aggregate?

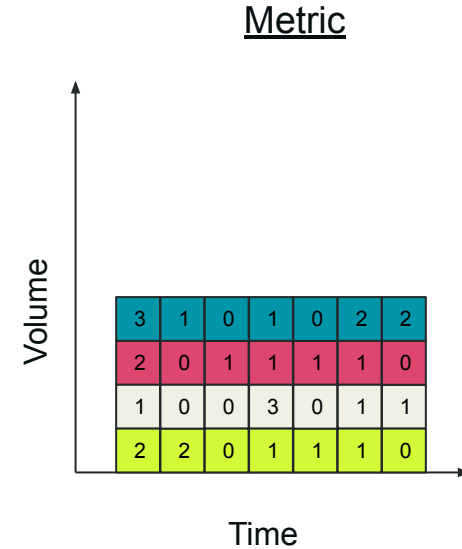
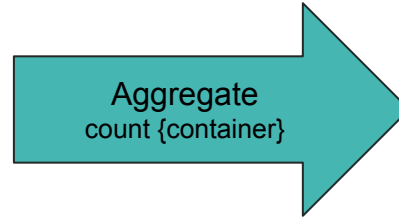
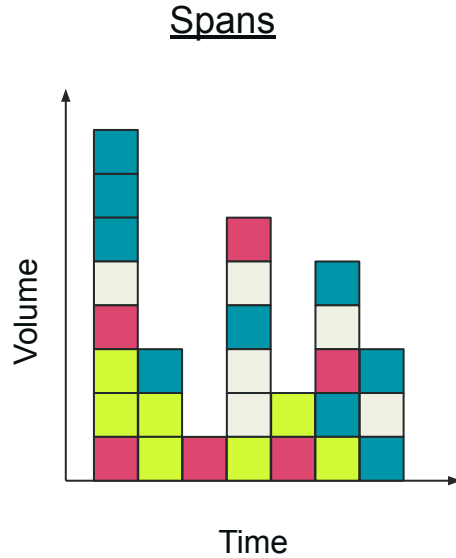


What is aggregation anyway?



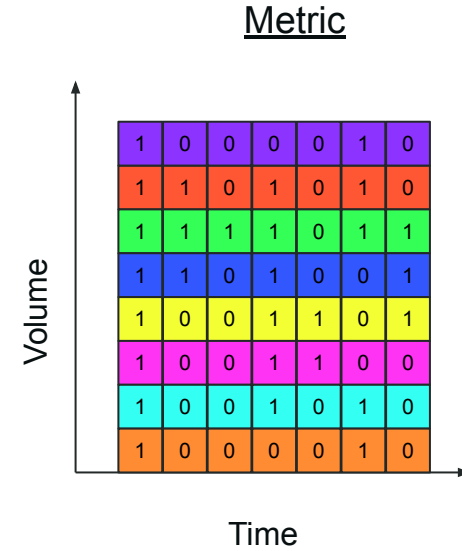
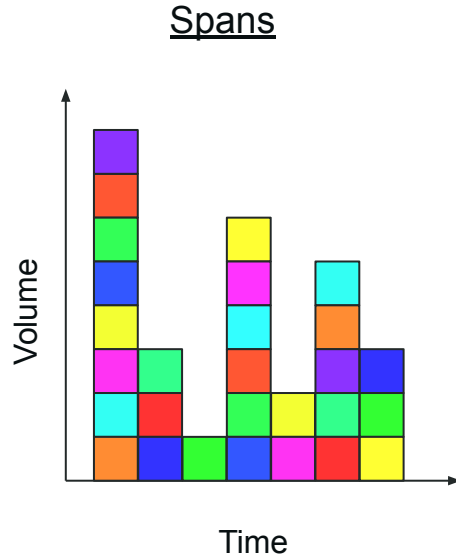
What is aggregation anyway?

With Dimensions



What is aggregation anyway?

High Cardinality



High Cardinality has a price (storage + memory to aggregate)

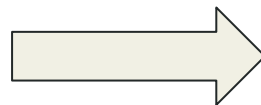
Might be better off storing raw data instead

Most cases: cardinality reduction for speed!



Before

End-to-end aggregation with
reduction in dimensionality



After

Aggregation as a performance
optimization over underlying data,
keep high cardinality data



Advantages of Post Collection

Keep high fidelity raw data for “interesting” observability questions

Measure cardinality *before* layering on aggregation

Leading Timeseries DBs and Observability Vendors automatically apply aggregation



“Events to Metrics”

Create metrics behind the scenes by aggregating raw telemetry into metrics

User decides at query time to use aggregates instead of raw data

Examples:

Clickhouse: Incremental Materialized Views

New Relic: Events to Metrics

Automatic Materialized View

Similar to indexes, aggregations build behind the scenes on raw data table

Database decides to use the aggregates at query time (no query syntax change!)

Examples:

Clickhouse: Projections

New Relic: Accelerated Views



“Events to Metrics”

Use when your retention of aggregates exceeds retention of raw data

Can be useful for metric standardization

Automatic Materialized View

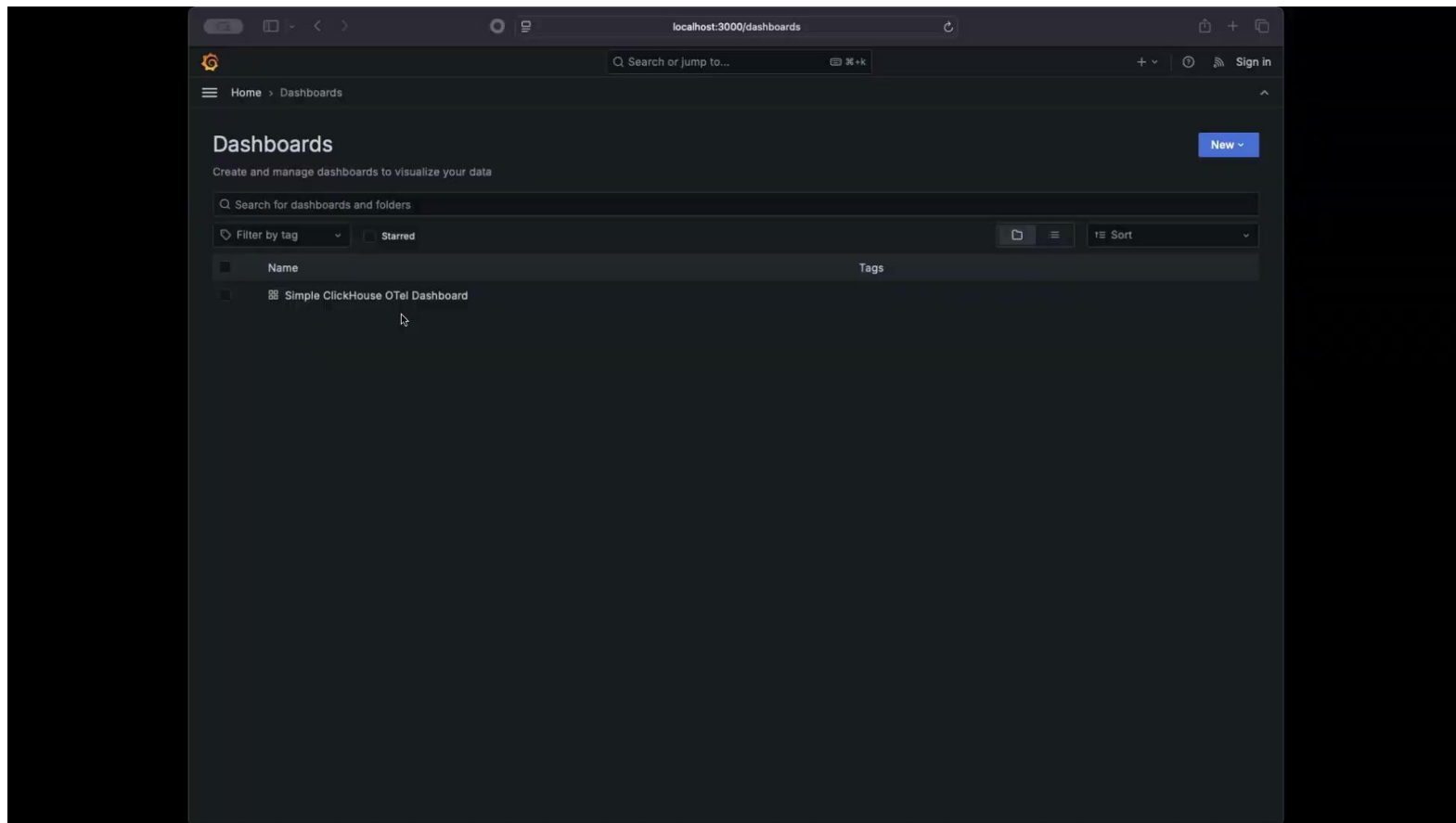
Default choice to accelerate common dashboards/alerts

Starting point to “drill into” data without switching contexts between aggregates and raw



Ex: Clickhouse Projections

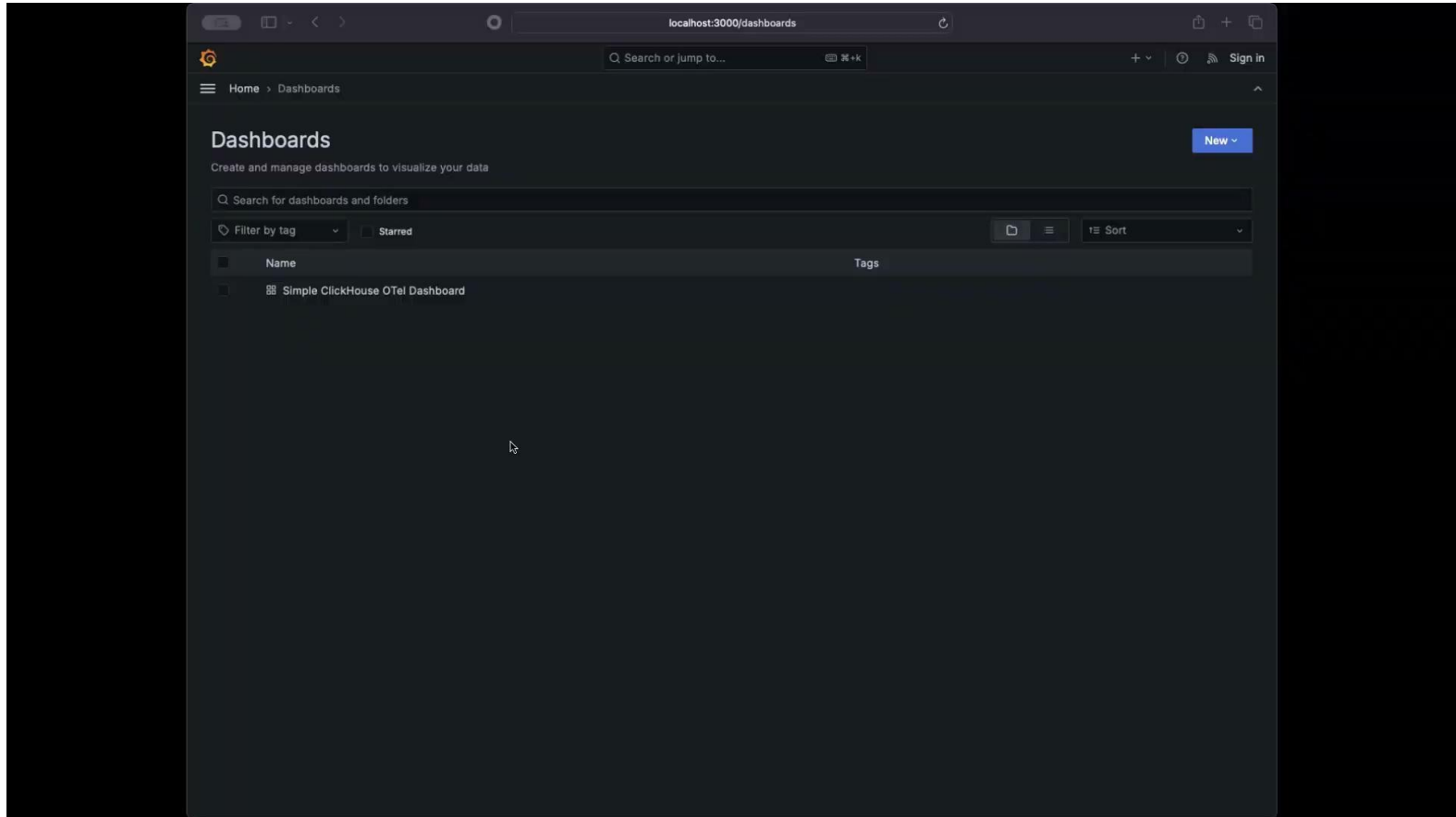
```
1 CREATE TABLE IF NOT EXISTS otel_traces
2 (
3     Timestamp      DateTime64(9)          CODEC(Delta, ZSTD(1)),
4     TraceId        String              CODEC(ZSTD(1)),
5     SpanId         String              CODEC(ZSTD(1)),
6     ParentSpanId  String              CODEC(ZSTD(1)),
7     TraceState    String              CODEC(ZSTD(1)),
8     SpanName       LowCardinality(String) CODEC(ZSTD(1)),
9     SpanKind       LowCardinality(String) CODEC(ZSTD(1)),
10    ServiceName    LowCardinality(String) CODEC(ZSTD(1)),
11    ResourceAttributes Map(LowCardinality(String), String) CODEC(ZSTD(1)),
12    ScopeName      String              CODEC(ZSTD(1)),
13    ScopeVersion   String              CODEC(ZSTD(1)),
14    SpanAttributes  Map(LowCardinality(String), String) CODEC(ZSTD(1)),
15    Duration       Int64               CODEC(ZSTD(1)),
16    StatusCode     LowCardinality(String) CODEC(ZSTD(1)),
17    StatusMessage  String              CODEC(ZSTD(1)),
18
19    Events Nested (Timestamp DateTime64(9), Name LowCardinality(String), Attributes Map(LowCardinality(String), String)) CODEC(ZSTD(1)),
20
21    Links Nested (TraceId String, SpanId String, TraceState String, Attributes Map(LowCardinality(String), String)) CODEC(ZSTD(1)),
22
23    INDEX idx_trace_id      TraceId              TYPE bloom_filter(0.001) GRANULARITY 1,
24    INDEX idx_res_attr_key  mapKeys(ResourceAttributes) TYPE bloom_filter(0.01)  GRANULARITY 1,
25    INDEX idx_res_attr_value mapValues(ResourceAttributes) TYPE bloom_filter(0.01)  GRANULARITY 1,
26    INDEX idx_span_attr_key mapKeys(SpanAttributes)   TYPE bloom_filter(0.01)  GRANULARITY 1,
27    INDEX idx_span_attr_value mapValues(SpanAttributes)   TYPE bloom_filter(0.01)  GRANULARITY 1,
28    INDEX idx_duration      Duration              TYPE minmax              GRANULARITY 1
29 )
30 ENGINE = MergeTree()
31 PARTITION BY toDate(Timestamp)
32 ORDER BY (ServiceName, SpanName, toUnixTimestamp(Timestamp), TraceId)
33 TTL toDate(Timestamp) + INTERVAL 30 DAY
34 SETTINGS index_granularity = 8192, ttl_only_drop_parts = 1;
35
```



Ex: Clickhouse Projections

```
1 -- add "minute" column to the table (simplifies aggregation)
2 ALTER TABLE otel_traces ADD COLUMN minute DateTime MATERIALIZED toStartOfMinute(Timestamp);
3 ALTER TABLE otel_traces MATERIALIZE COLUMN minute;
4
5 -- Add projection with count(), sum(Duration), p50, p95, p99
6 ALTER TABLE otel_traces ADD PROJECTION p_svc_minute
7 (
8     SELECT
9         ServiceName,
10        StatusCode,
11        minute,
12        count(),
13        sum(Duration),
14        quantile(0.5)(Duration),
15        quantile(0.95)(Duration),
16        quantile(0.99)(Duration)
17    GROUP BY
18        ServiceName, StatusCode, minute
19 );
20
21 ALTER TABLE otel_traces MATERIALIZE PROJECTION p_svc_minute;
```





```
EXPLAIN
SELECT
  minute AS time,
  ServiceName,
  quantile(0.99)(Duration) / 1000000. AS p99_ms
FROM otel_traces
WHERE minute >= (now() - toIntervalYear(6))
GROUP BY
  time,
  ServiceName
ORDER BY time ASC
```

Query id: 9bc83f04-330b-4a54-8be9-ba3c01bale10

```
explain
1. Expression ((Project names + (Before ORDER BY + Projection) [lifted up part]))
2.   Sorting (Sorting for ORDER BY)
3.     Expression ((Before ORDER BY + Projection))
4.       Aggregating
5.         Expression (Before GROUP BY)
6.           Expression
7.             ReadFromMergeTree (default.otel_traces)
```



```
EXPLAIN
SELECT
    minute AS time,
    ServiceName,
    quantile(0.99)(Duration) / 1000000. AS p99_ms
FROM otel_traces
WHERE minute >= (now() - toIntervalYear(6))
GROUP BY
    time,
    ServiceName
ORDER BY time ASC
```

Query id: b8873535-9a10-4d7d-8a0c-d6ea4092b8ec

```
explain
1. Expression ((Project names + (Before ORDER BY + Projection) [lifted up part]))
2.   Sorting (Sorting for ORDER BY)
3.     Expression ((Before ORDER BY + Projection))
4.       Aggregating
5.         Filter
6.           ReadFromMergeTree (p_svc_minute)
```



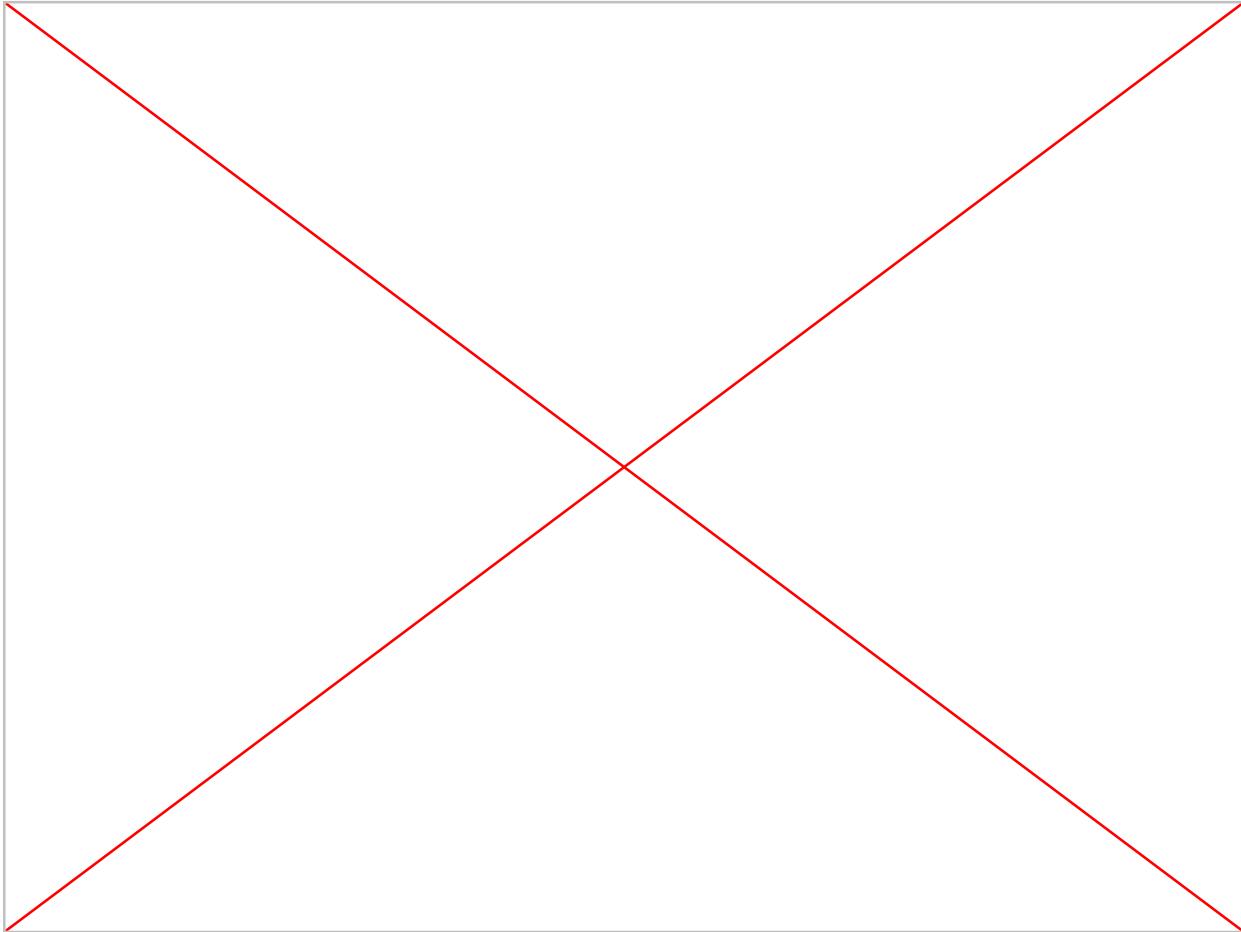
Before Projections

```
300 rows in set. Elapsed: 5.183 sec. Processed 360.00 million rows, 4.68 GB (69.45 million rows/s., 902.87 MB/s.)  
Peak memory usage: 81.51 MB.
```

After Projections

```
300 rows in set. Elapsed: 0.091 sec.
```





Measure Cardinality

```
0828bcd6562a :) SELECT
count() AS raw_rows,

-- Cardinality: ServiceName + StatusCode
uniq(ServiceName, StatusCode) AS k_svc_status,
k_svc_status * 360 AS proj_rows_svc_status,
round(raw_rows / proj_rows_svc_status) AS shrink_svc_status,

-- Cardinality: ServiceName + StatusCode + customer ID
uniq(ServiceName, StatusCode, SpanAttributes['customer.id']) AS k_plus_customer,
k_plus_customer * 360 AS proj_rows_plus_customer,
round(raw_rows / proj_rows_plus_customer) AS shrink_plus_customer,

-- Cardinality: ServiceName + StatusCode + user ID
uniq(ServiceName, StatusCode, SpanAttributes['user.id']) AS k_plus_user,
k_plus_user * 360 AS proj_rows_plus_user,
round(raw_rows / proj_rows_plus_user) AS shrink_plus_user

FROMotel_traces
FORMAT Vertical;
```



Raw Rows*: 3,600,000

```
Row 1:  
raw_rows: 360000000  
k_svc_status: 10  
proj_rows_svc_status: 3600  
shrink_svc_status: 100000  
k_plus_customer: 10000  
proj_rows_plus_customer: 3600000  
shrink_plus_customer: 100  
k_plus_user: 997416  
proj_rows_plus_user: 359069760  
shrink_plus_user: 1
```

Dimensions	Unique Groups	Projection Rows*	Reduction*
<i>ServiceName, Statuscode</i>	10	3,600	100,000:1
<i>ServiceName, Statuscode, Customer ID</i>	10,000	3,6600,000	100:1
<i>ServiceName, Statuscode, User ID</i>	~1,000,000	~360,000,000	1:1

*assuming 6 hour timeframe



Small Cost for Speed

Small write time performance cost of aggregation

Additional storage space for aggregates

Still paying cost to send & store raw data

Query matching can be temperamental (check EXPLAIN)



Recommendations:

- Add as much detail as possible to your spans/logs
- Ingest/store everything as long as useful
- For long range, known analysis create aggregates after collection

Future: automatic identification and creation of aggregates based on query patterns!



Q/A

async: zowens@newrelic.com

