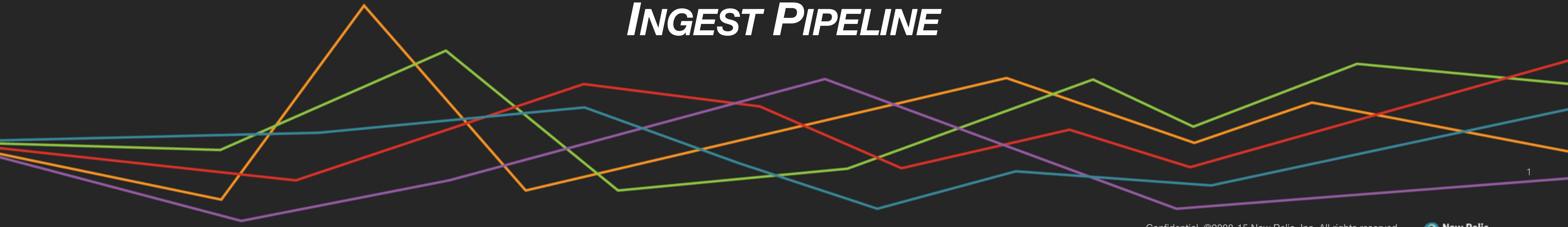


EVALUATING STREAMING FRAMEWORK PERFORMANCE FOR A LARGE-SCALE AGGREGATION PIPELINE

RON CROCKER

(rcrocker@newrelic.com)

*PRINCIPAL ENGINEER & ARCHITECT
INGEST PIPELINE*



This document and the information herein (including any information that may be incorporated by reference) is provided for informational purposes only and should not be construed as an offer, commitment, promise or obligation on behalf of New Relic, Inc. (“New Relic”) to sell securities or deliver any product, material, code, functionality, or other feature. Any information provided hereby is proprietary to New Relic and may not be replicated or disclosed without New Relic’s express written permission.

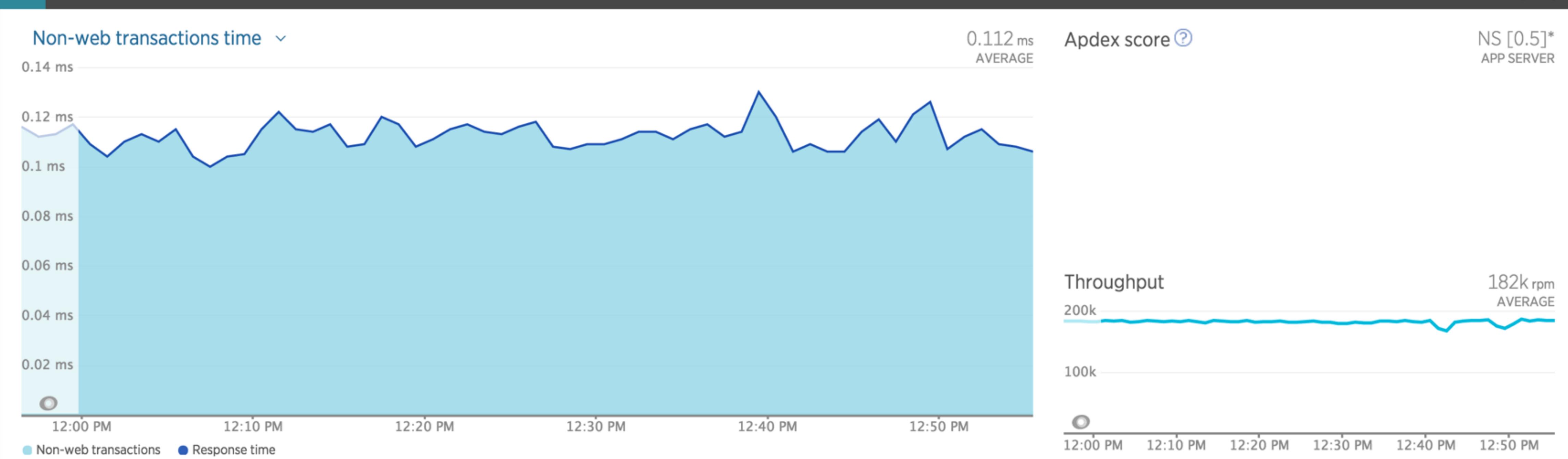
Such information may contain forward-looking statements within the meaning of federal securities laws. Any statement that is not a historical fact or refers to expectations, projections, future plans, objectives, estimates, goals, or other characterizations of future events is a forward-looking statement. These forward-looking statements can often be identified as such because the context of the statement will include words such as “believes,” “anticipates,” “expects” or words of similar import.

Actual results may differ materially from those expressed in these forward-looking statements, which speak only as of the date hereof, and are subject to change at any time without notice. Existing and prospective investors, customers and other third parties transacting business with New Relic are cautioned not to place undue reliance on this forward-looking information. The achievement or success of the matters covered by such forward-looking statements are based on New Relic’s current assumptions, expectations, and beliefs and are subject to substantial risks, uncertainties, assumptions, and changes in circumstances that may cause the actual results, performance, or achievements to differ materially from those expressed or implied in any forward-looking statement. Further information on factors that could affect such forward-looking statements is included in the filings we make with the SEC from time to time. Copies of these documents may be obtained by visiting New Relic’s Investor Relations website at ir.newrelic.com or the SEC’s website at www.sec.gov.

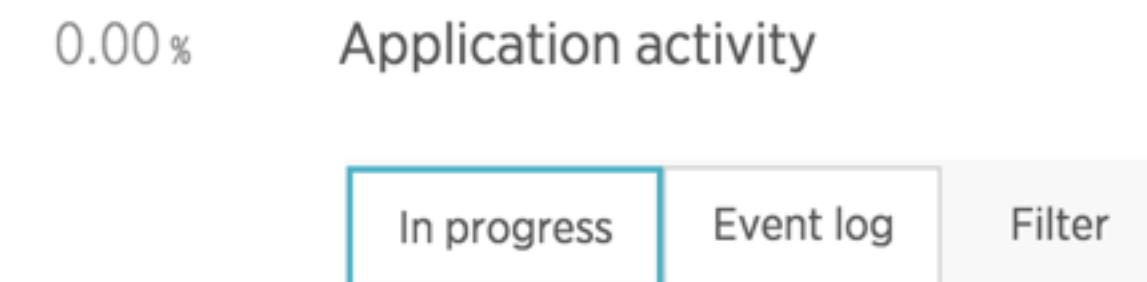
New Relic assumes no obligation and does not intend to update these forward-looking statements, except as required by law. New Relic makes no warranties, expressed or implied, in this document or otherwise, with respect to the information provided.

New Relic Software Analytics Cloud





Transactions	App server time
primaryTimesliceHarvest	609 ms
Transaction traces: 1 s 0.9 s 0.9 s	
publish	0.388 ms
Transaction traces: 0.5 s 0.5 s 0.4 s	
ResolvedTimeslices	0.144 ms
Transaction traces: 0.4 s 0.1 s 0.1 s	
publish	0.026 ms
Transaction traces: 0.1 s 0.1 s 0.1 s	
JettyMonitor	0.0134 ms
Transaction traces: n/a	



There are no violations in progress right now.

4 servers

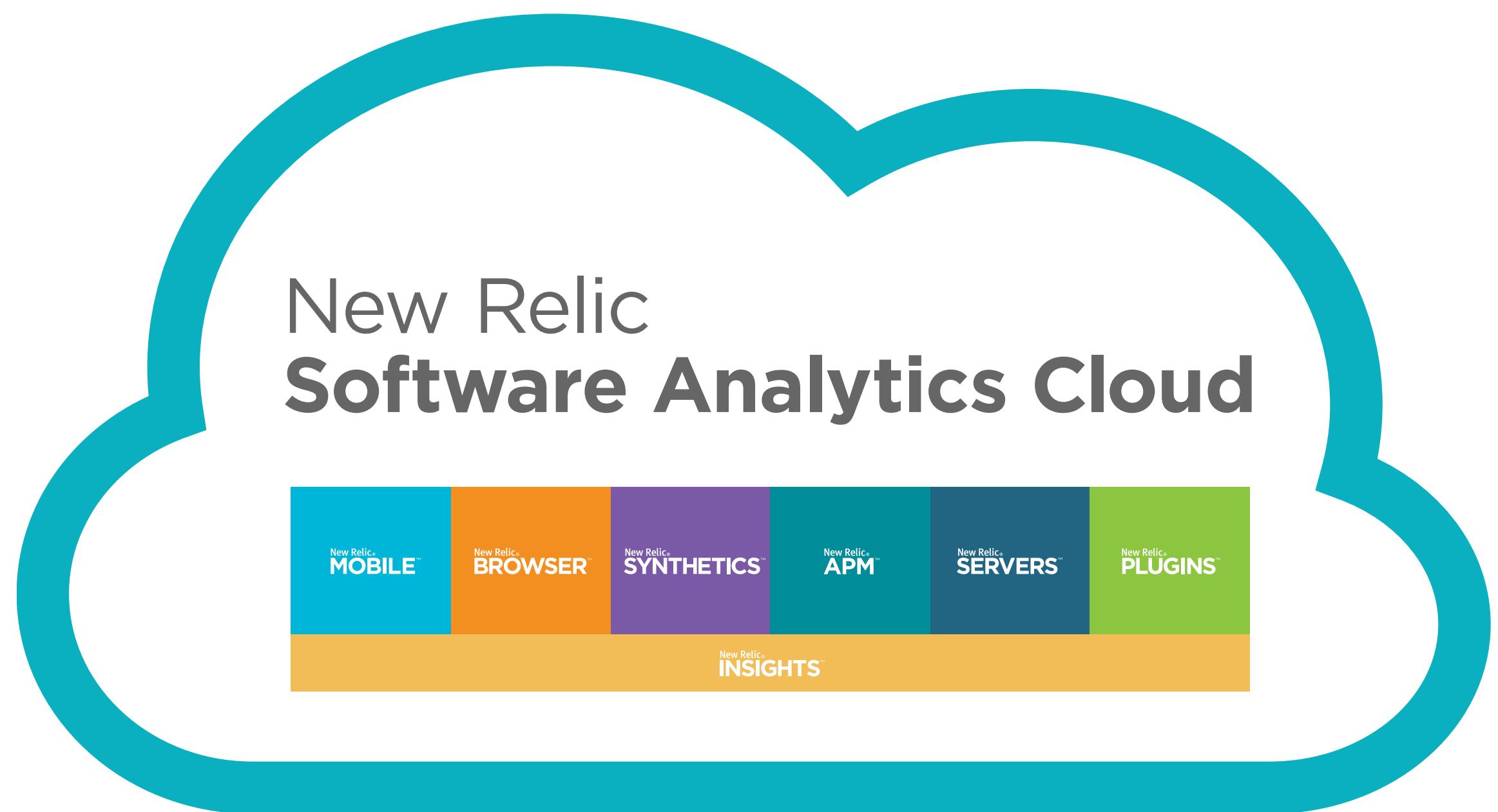


Server name

Apdex ▾ Resp. time ▾

Minute Timeslice Aggregator Throughput

0 ms 184k rpm 0.00 err% ⓘ CPU usage ▾ Memory ▾



accepts over **16M** requests

stores over **2M** analytic events

aggregates over **800M** metrics

queries over **3B** data points

EVERY MINUTE

New Relic **Software Analytics Cloud**



contains
over

200

different
services

more
than

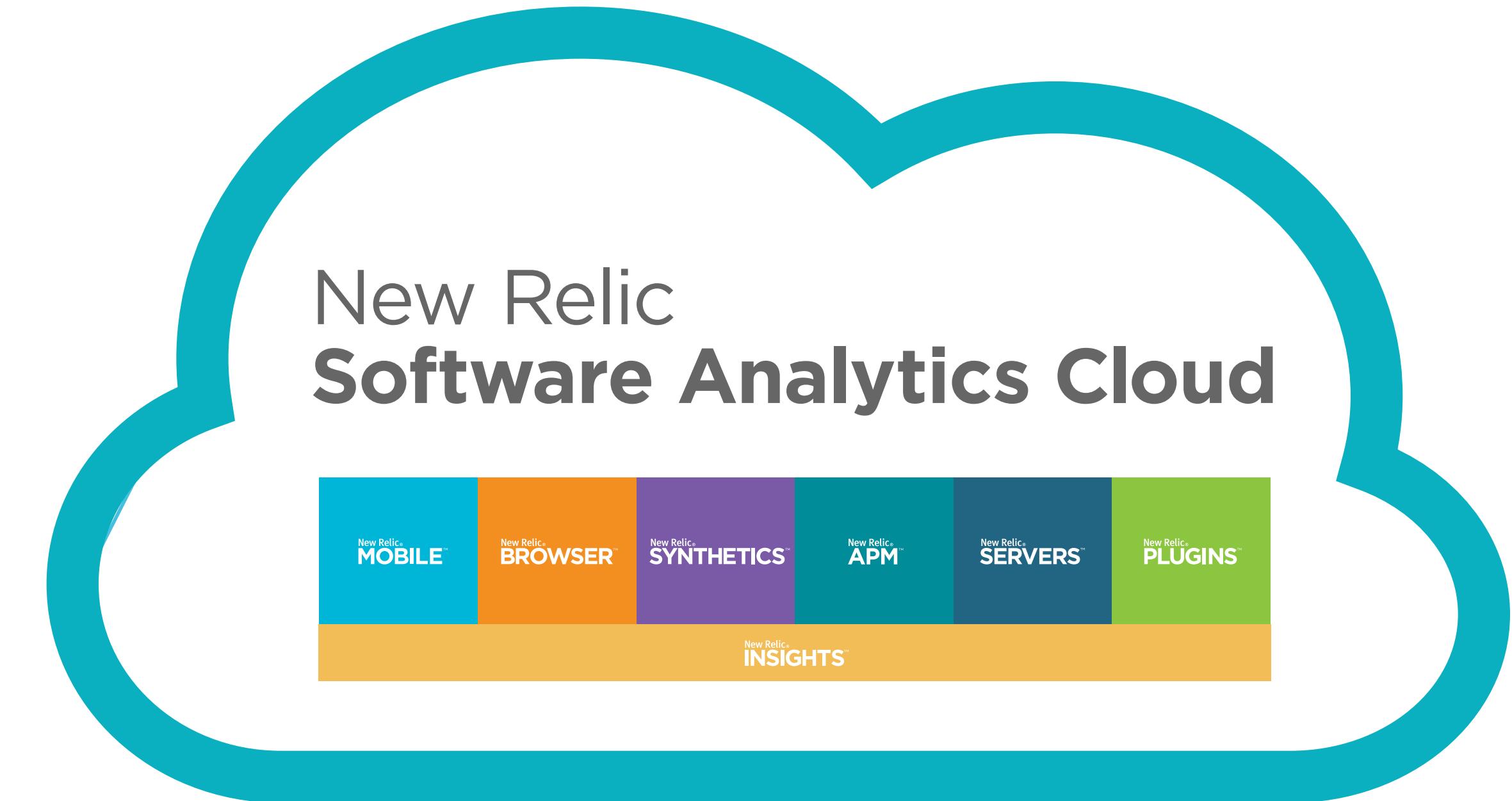
2.5
PETABYTES

SSD
storage

maintained/
built by

25+

engineering
teams





 New Relic®

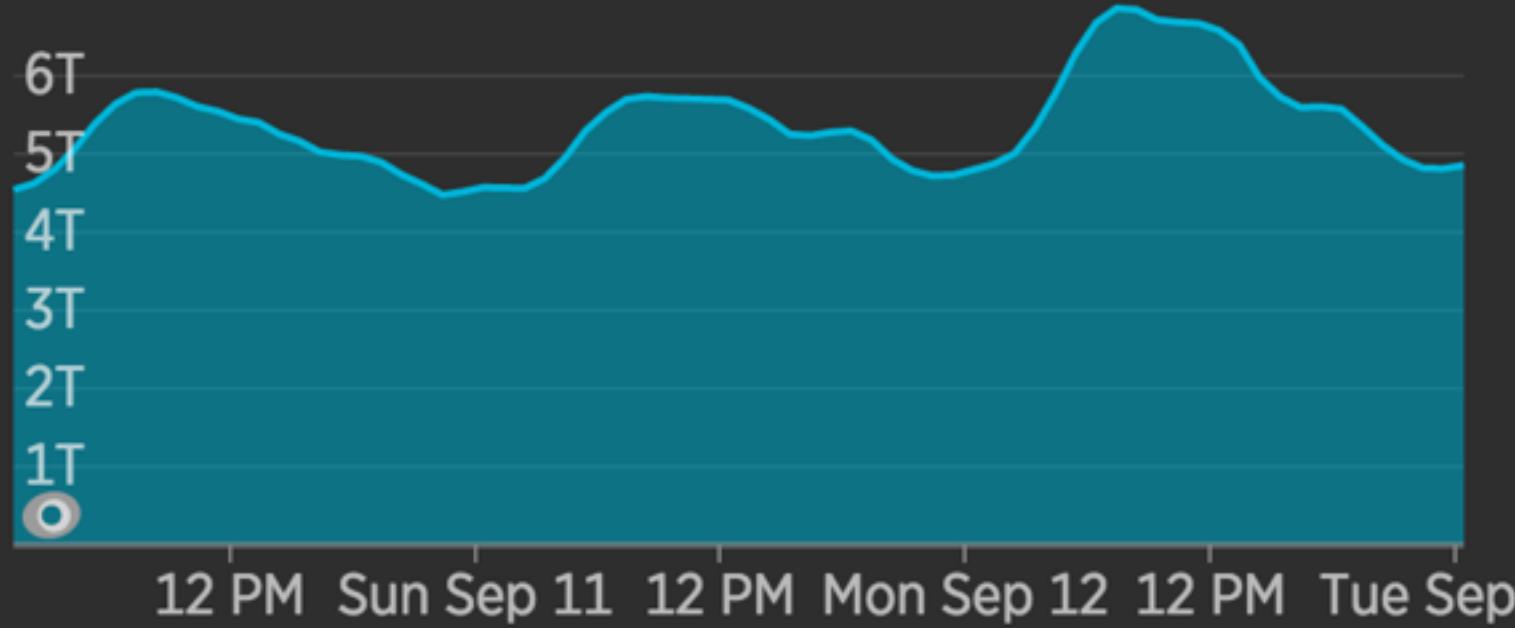
Production Kafka

Created by rcrocker@newrelic.com

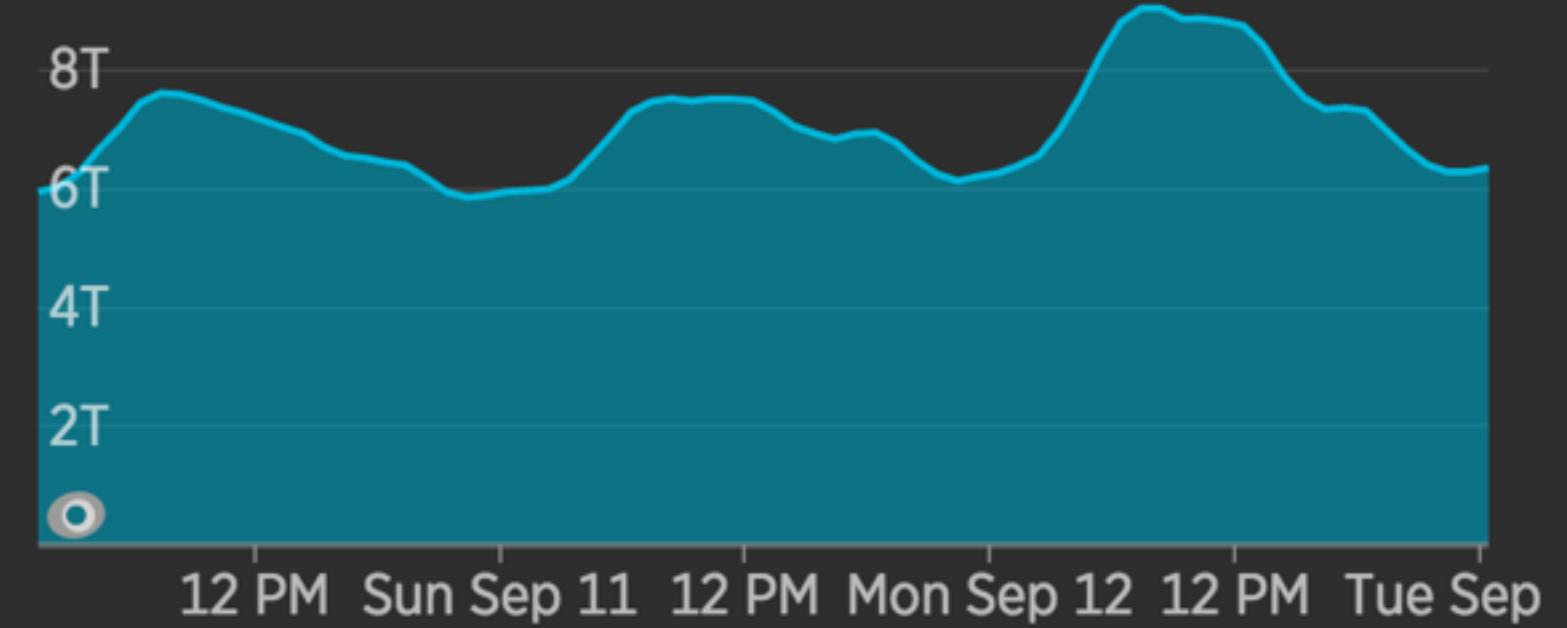
powered by **New Relic[®] INSIGHTS™** 

Last edited 9/12/16

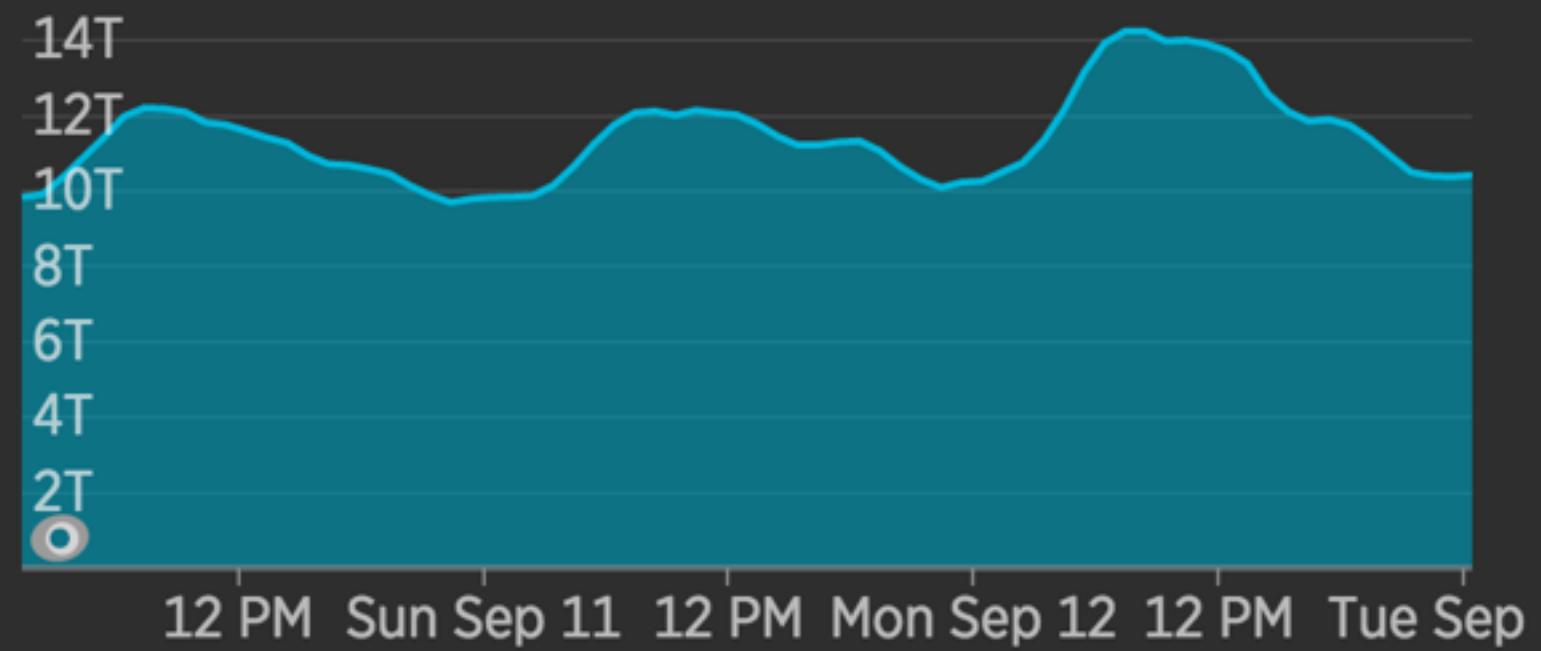
Kafka Producers Since 72 hours ago



Kafka Consumers Since 72 hours ago



Kafka replication Since 72 hours ago



Kafka Producers Since 24 hours ago

139

Total Byte Count (TB)

Kafka Consumers Since 24 hours ago

175

Total Byte Count (TB)

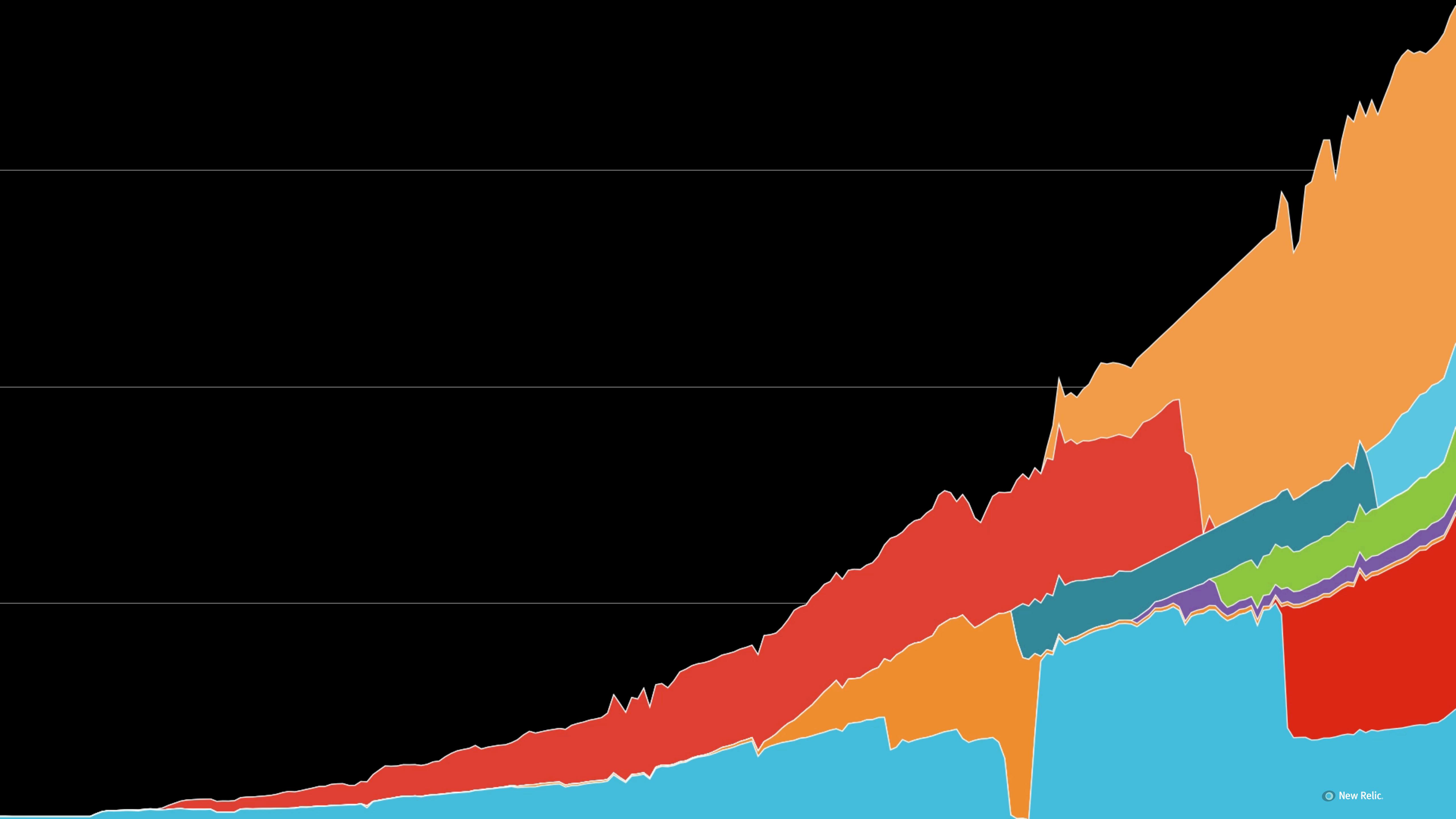
Kafka replication Since 24 hours ago

294

Total Byte Count (TB)

[NR Export](#)

[Permalink](#)

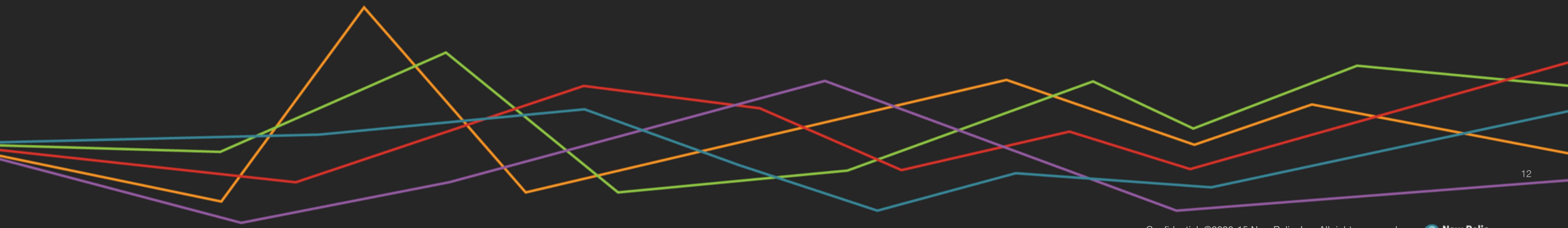


Goals for evaluating streaming systems

- Understand performance characteristics
- Understand operations characteristics

How New Relic works...

... the cartoon version



A₁

An instance of your application running on a host

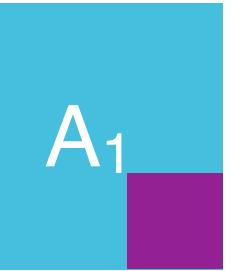
A₂

Another instance of your application running on another host

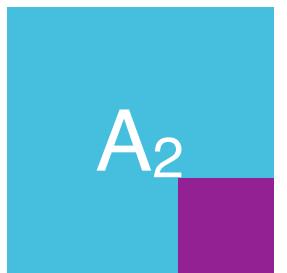
...

A_n

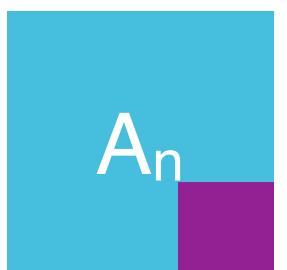
More instances of your application running on more hosts

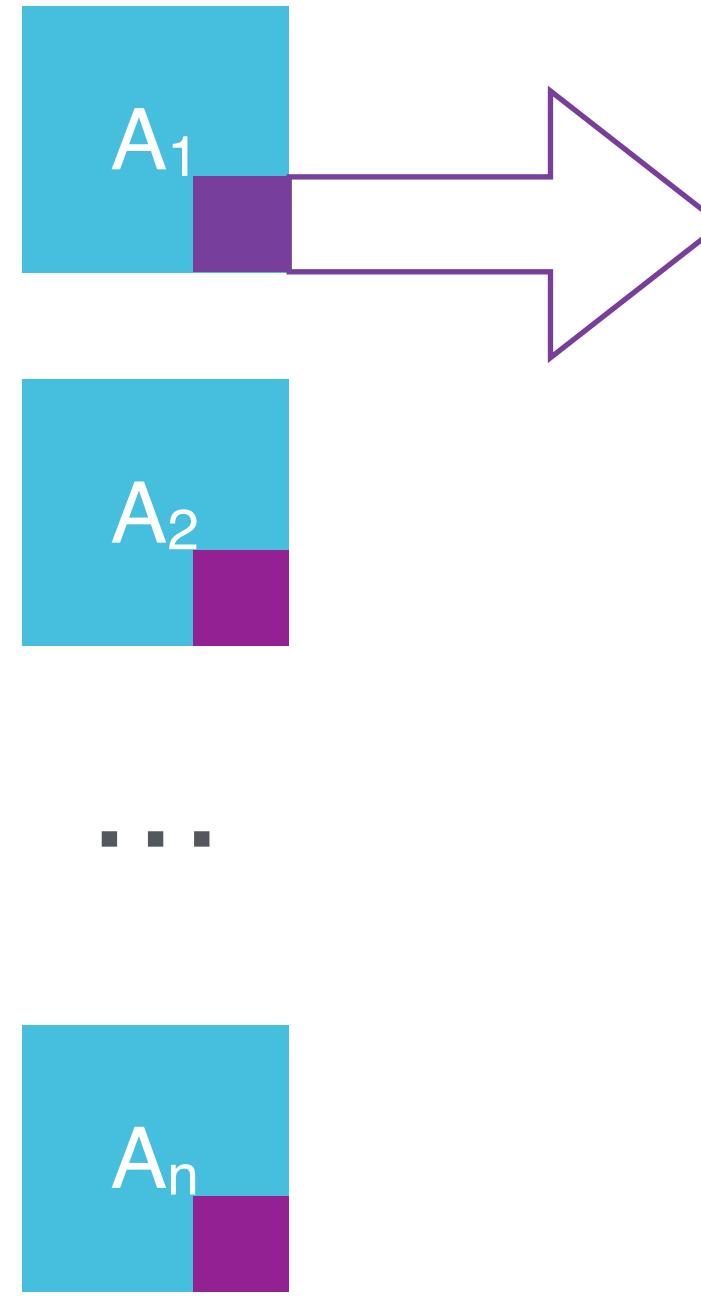


New Relic Agent reports data to New Relic

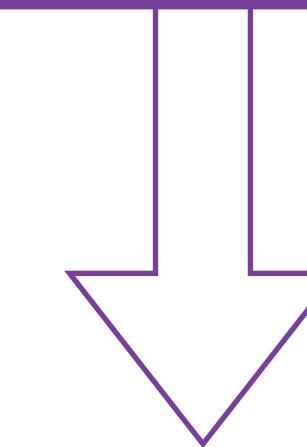


...





- Agent Token (\approx account ID, agent ID)
- Duration (time-period covered)
- Timeslices: Each timeslice contains
 - Metric name
 - Metric stats
 - Count, total time, exclusive time, min, max, sum of squares



HTTP post to <something>.newrelic.com

APPS
Minute Timeslice Aggreg...TIME PICKER
Last 60 minutes ending nowJVMS
All JVMs

MONITORING

Overview

Service maps

App map

Transactions

External services

JVMs

Threads

EVENTS

Error analytics

Errors

Violations

Dependencies

Threads

REPORTS

SLAs

Availability

Scalability

Web transactions

Database

Non-web transactions time

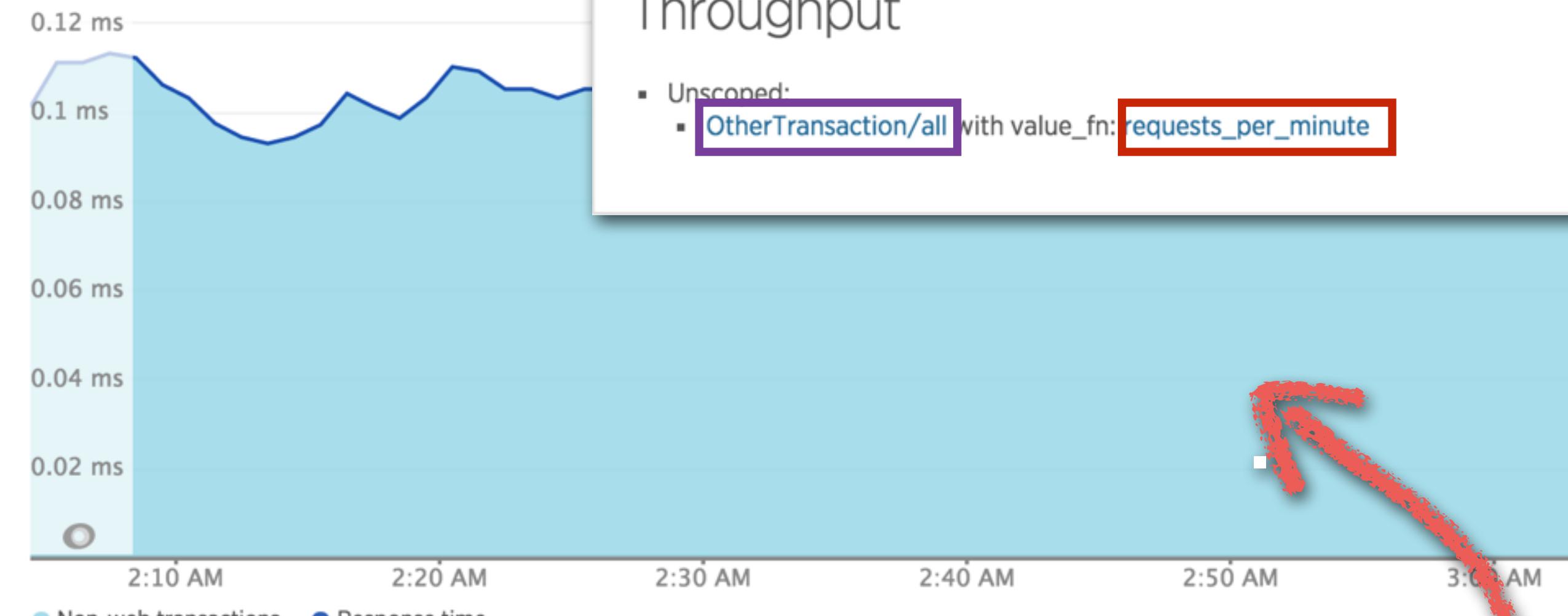
0.104 ms
AVERAGE

Apdex score

NS [0.5]*
APP SERVER

Throughput

- Unscoped:
 - OtherTransaction/all with value_fn: requests_per_minute



Throughput

200k

150k

100k

50k

0k

0.00 %

Application activity

In progress

Event log

Filter



There are no violations in progress right now.

Average response time from non-web transactions, broken down by measure

- Unscoped:
 - OtherTransaction/all with value_fn: average_exclusive_time
 - OtherTransaction/all with value_fn: average_exclusive_time

ResolvedTimeslices 0.159 ms

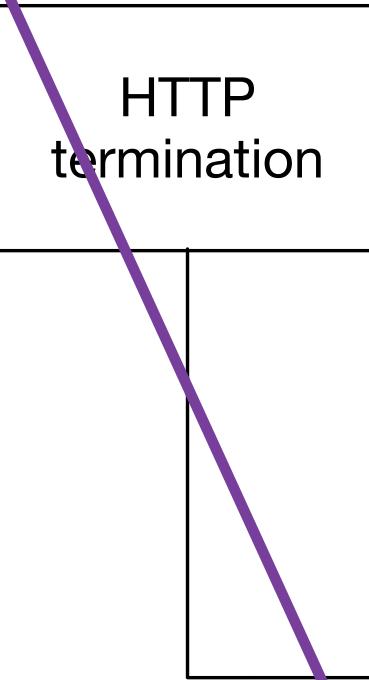
Transaction traces: 0.3 s 0.3 s 0.2 s

publish 0.022 ms

Transaction traces: 0.2 s 0.1 s 0.1 s

0 ms 162k rpm 0.00 err%

A₁

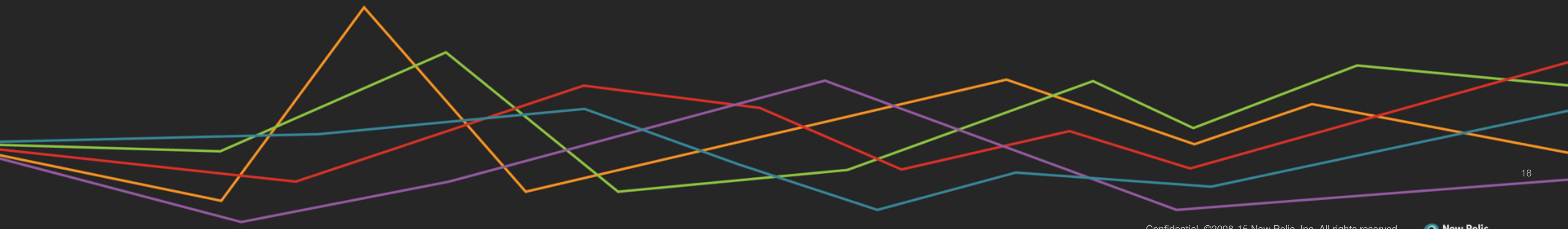


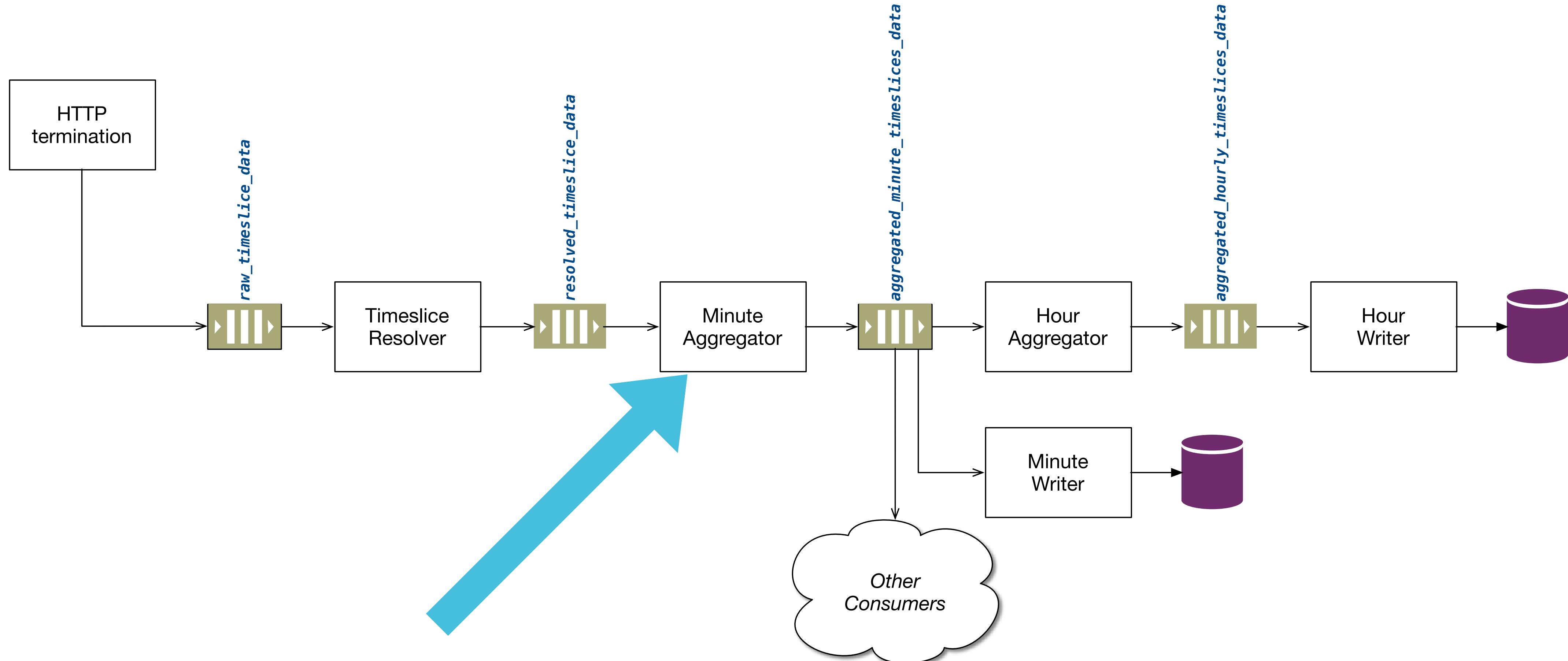
- Account ID
- Agent ID
- Duration
- Timeslices
- Metric ID
- Metric stats
 - **Start time**
 - **Duration (time-period covered)**
 - **Count, total time, exclusive time, min, max, sum of squares**

Metric ID

- Timeslices
 - Metric stats
 - Count, total time, exclusive time, min, max, sum of squares
 - Metric stats
 - Count, total time, exclusive time, min, max, sum of squares

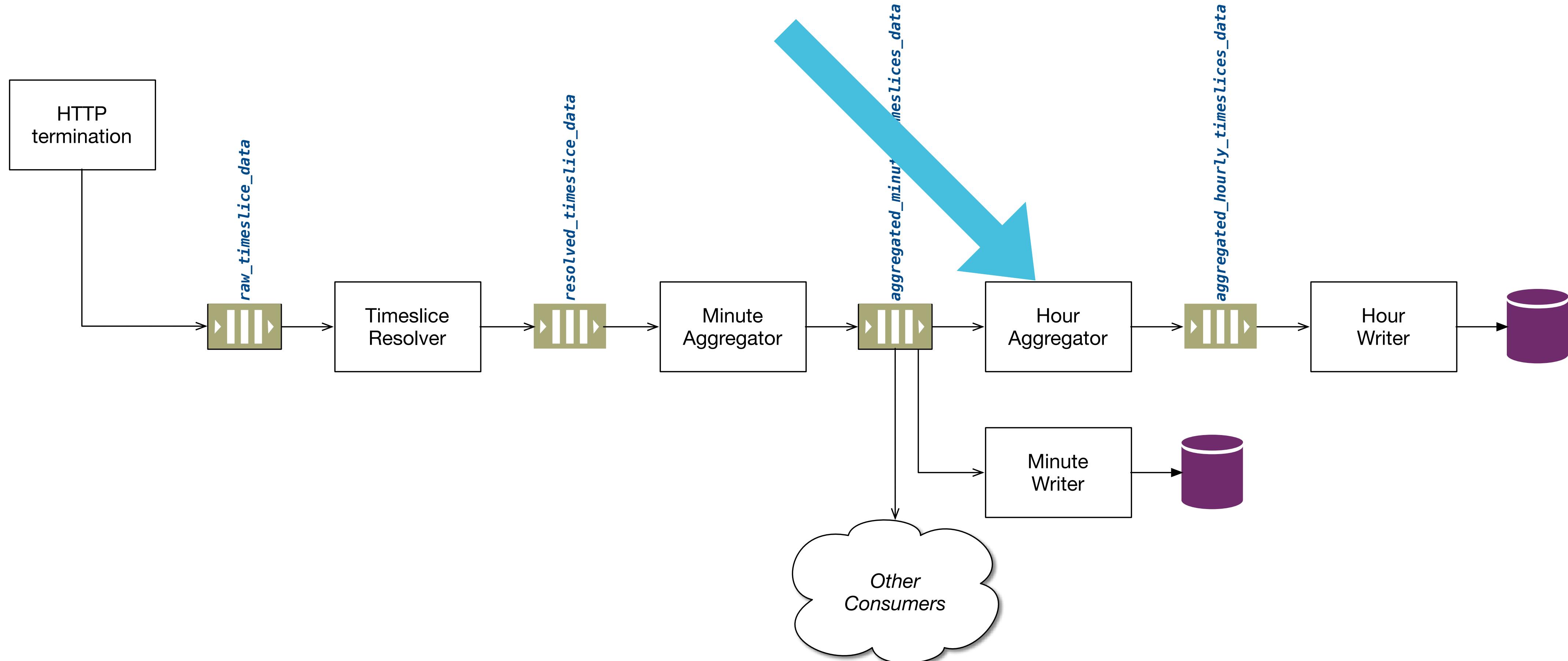
The Experiment





Why *Minute Aggregator*?

- **No external dependencies**
 - Performance comparisons solely focused on processing
- **Repeatable**
 - We can compare across technologies without needing to normalize
- **Important to our business**
 - Provides aggregation across instances of your application
 - We could have benchmarked something else, like Yahoo benchmark or word count, but would it have mattered?



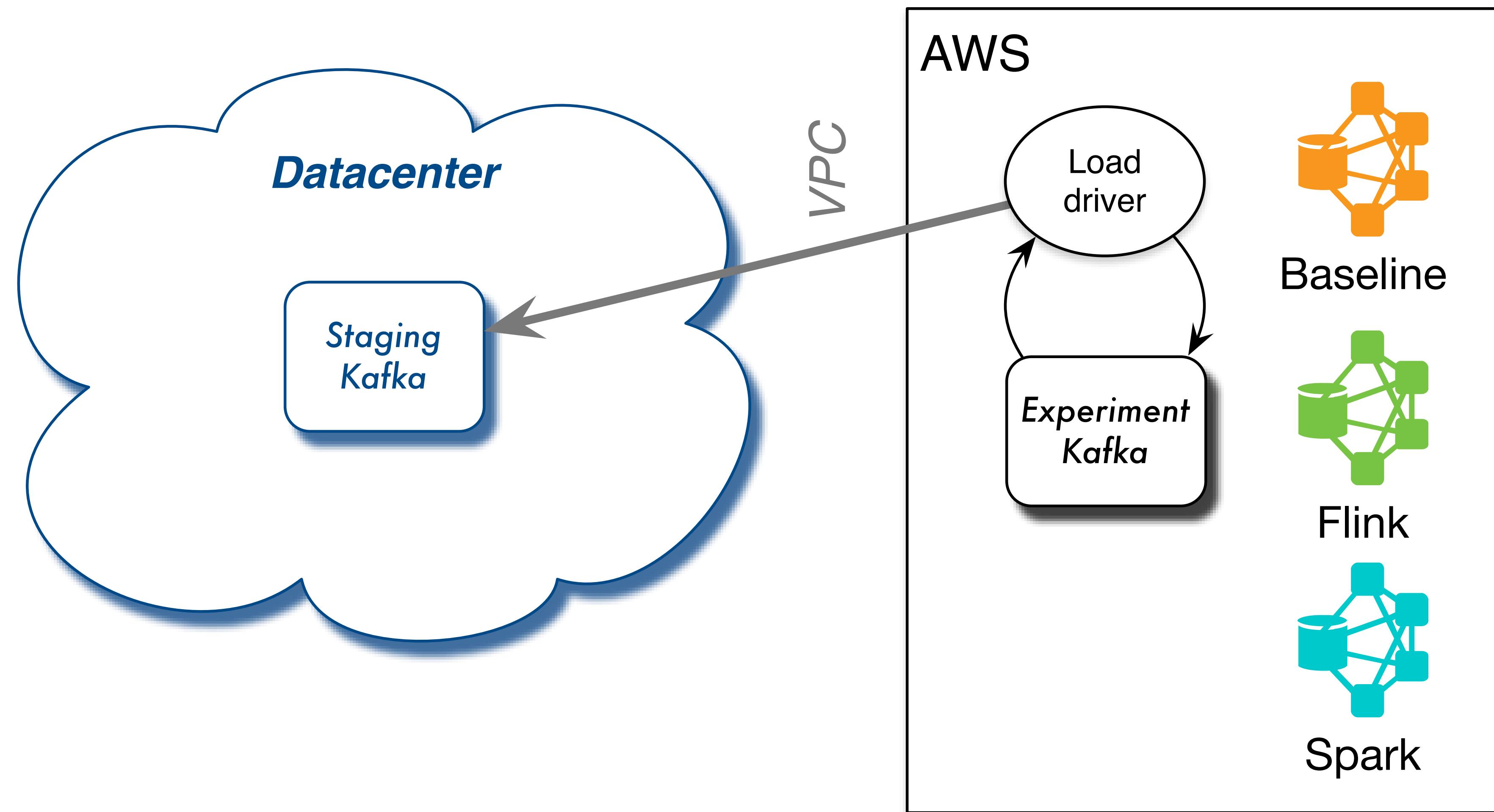
What about *Hour Aggregator*?

- **Similar to *Minute Aggregator***
 - *No external dependencies, Repeatable, Important to the business*
- **Needs to run for several hours to understand performance**
 - ... and I'm that patient
- **Extra credit: Integrate into stream implementations**

Goals for evaluating streaming systems

- **Understand performance characteristics**
 - Performance at different arrival rates:
 - 100%
 - 6000%
 - *To infinity and beyond*
- **Understand operations characteristics**
 - *No explicit goal*

Evaluation Framework



AWS Configurations

- **Kafka + ZK**
 - 3 i2.8xlarge hosts
- **Baseline**
 - 3 m4.4xlarge hosts
- **Flink**
 - 4 m4.4xlarge hosts
- **Spark**
 - EMR - 1 master , 3 workers, all m4.4xlarge

<i>AWS Configuration</i>	i2.8xlarge	m4.4xlarge
Cores	32	16
RAM	244GB	64GB
Network Bandwidth	10Gbps	2Gbps

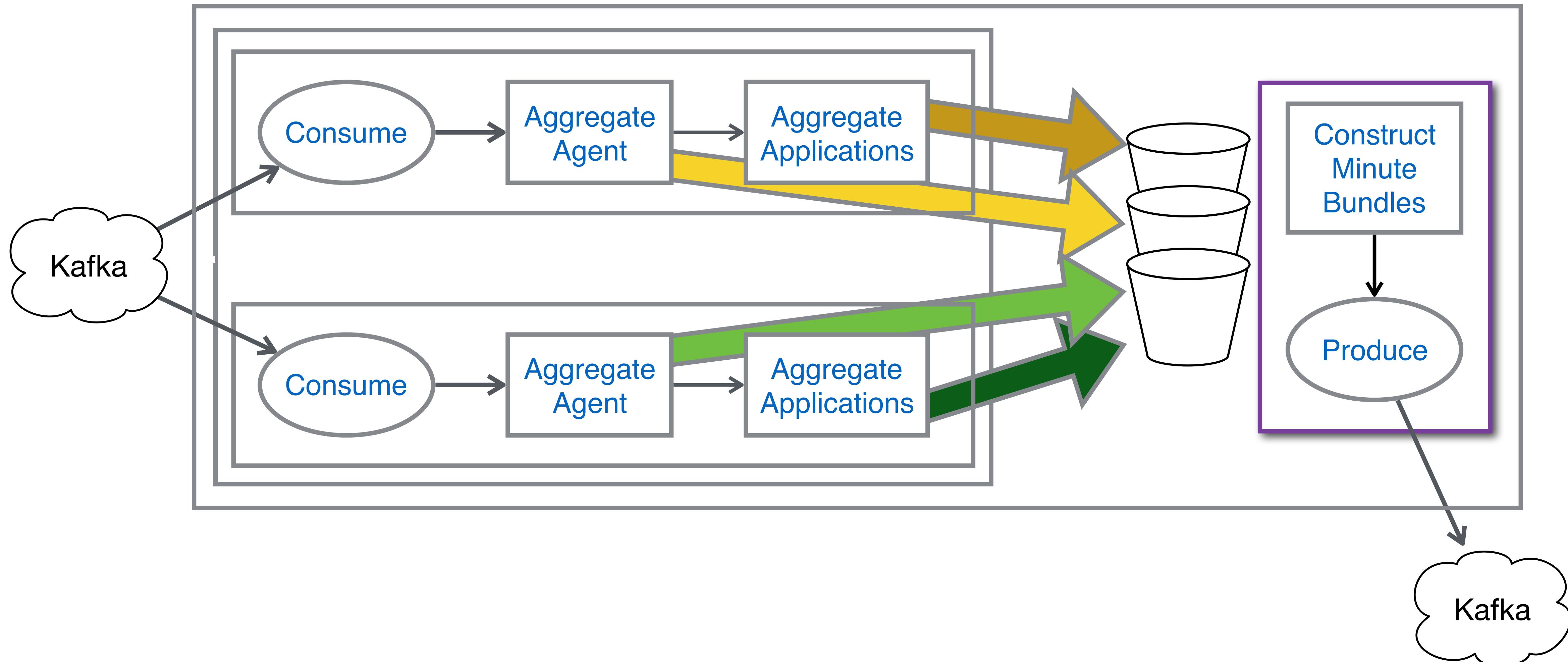
Experimental Kafka system

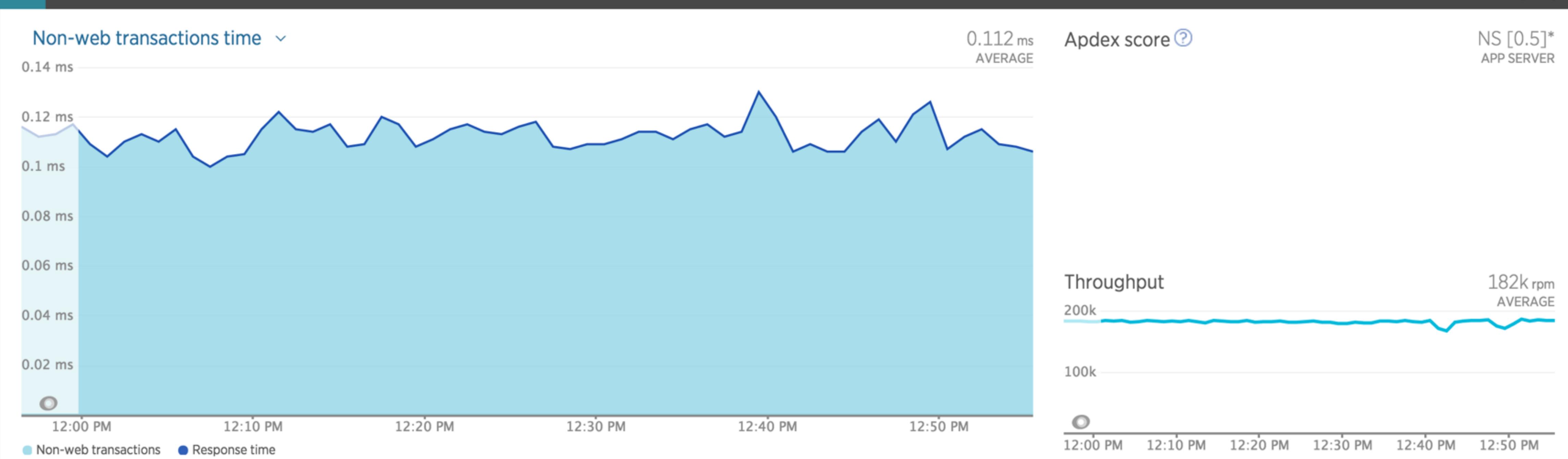
- **Kafka 0.8.2.2**
 - NR fork, includes back ports of some 0.9 features
- **# partitions: 16**
 - It's possible that this is too few partitions for the Baseline system

Load Driver

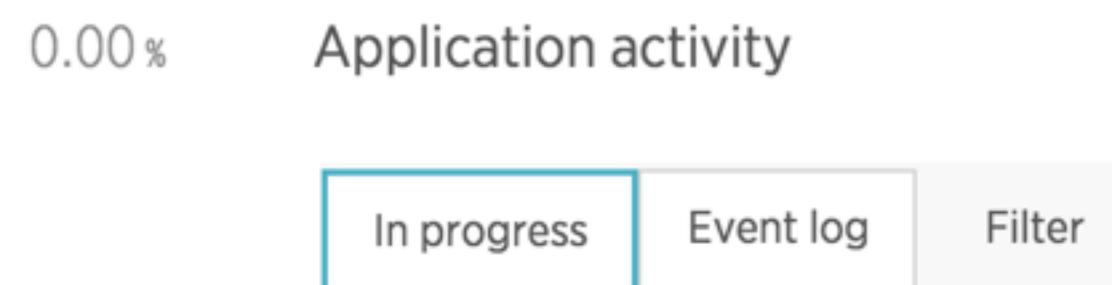
- Generates simple synthetic load *based on real traffic*
 - *Real traffic* = output of Timeslice Resolver
 - *Load generated based on repeated messages*
- Synthesizing *interesting* load is challenging:
 - Un-bundle timeslices
 - Generate re-bundled with new IDs - Agent, Account and/or Metrics
 - Repeat as necessary to get to load point

Baseline system - Our incumbent *Minute Aggregator*





Transactions	App server time
primaryTimesliceHarvest	609 ms
Transaction traces: 1 s 0.9 s 0.9 s	
publish	0.388 ms
Transaction traces: 0.5 s 0.5 s 0.4 s	
ResolvedTimeslices	0.144 ms
Transaction traces: 0.4 s 0.1 s 0.1 s	
publish	0.026 ms
Transaction traces: 0.1 s 0.1 s 0.1 s	
JettyMonitor	0.0134 ms
Transaction traces: n/a	



There are no violations in progress right now.

4 servers

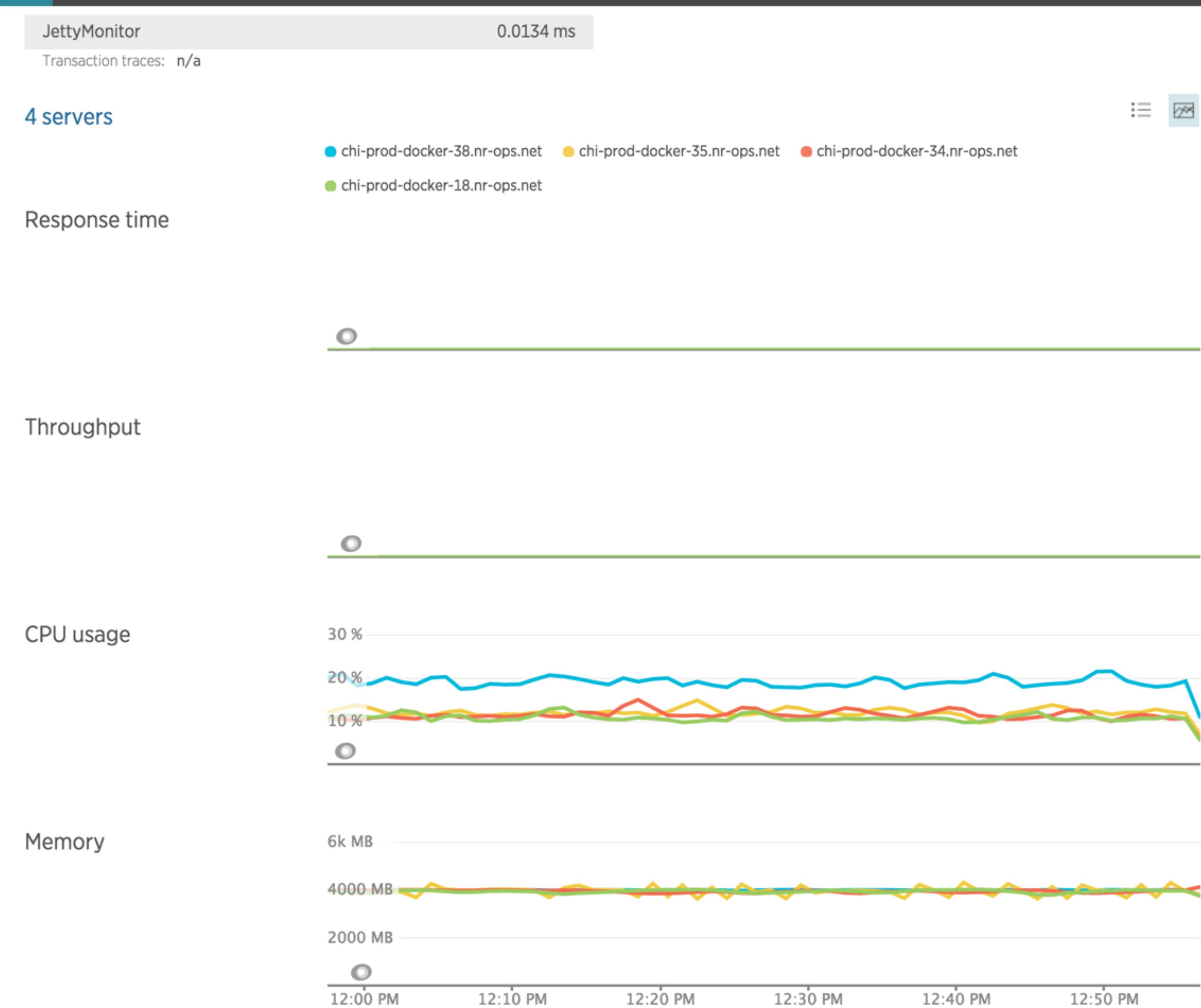


Server name

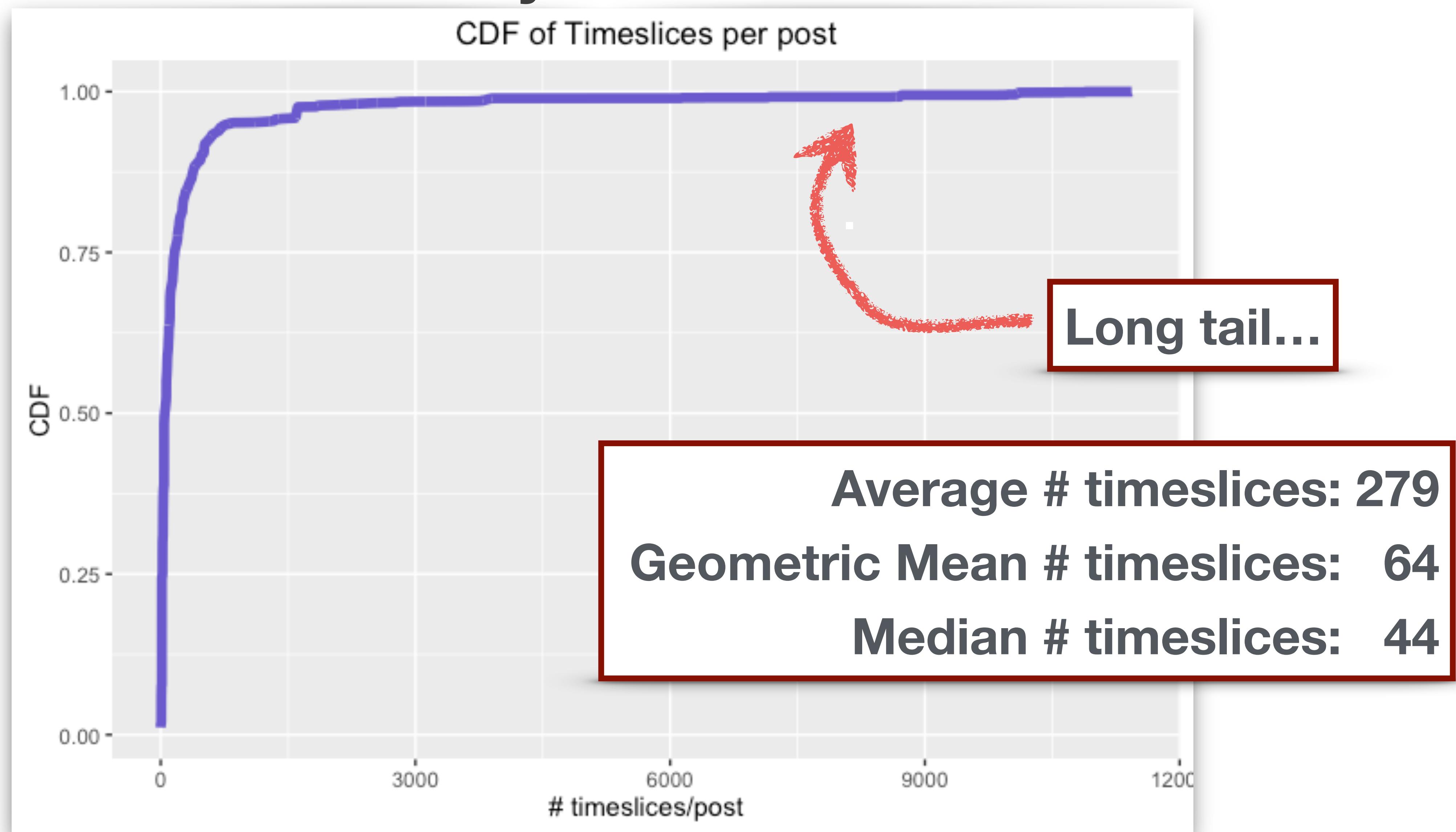
Apdex ▾ Resp. time ▾

Minute Timeslice Aggregator
Throughput ▾ Error Rate ▾

0 ms 184k rpm 0.00 err% ? CPU usage ▾ Memory ▾



Distributions are not friendly...



Flink configuration

Task Manager
(16 slots)

Job Manager

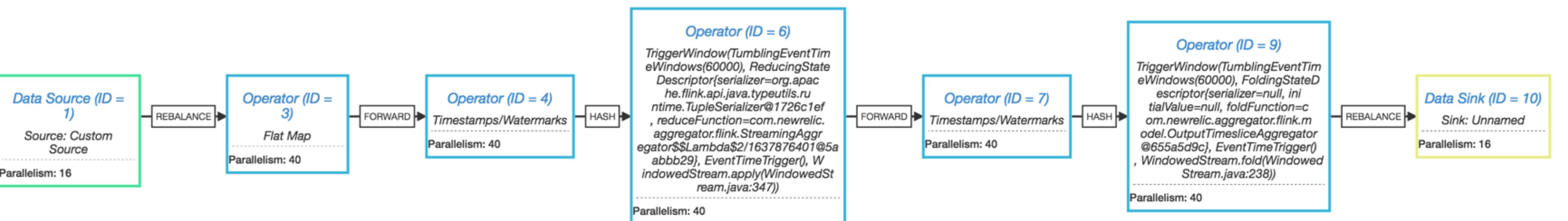
Task Manager
(16 slots)

Task Manager
(16 slots)

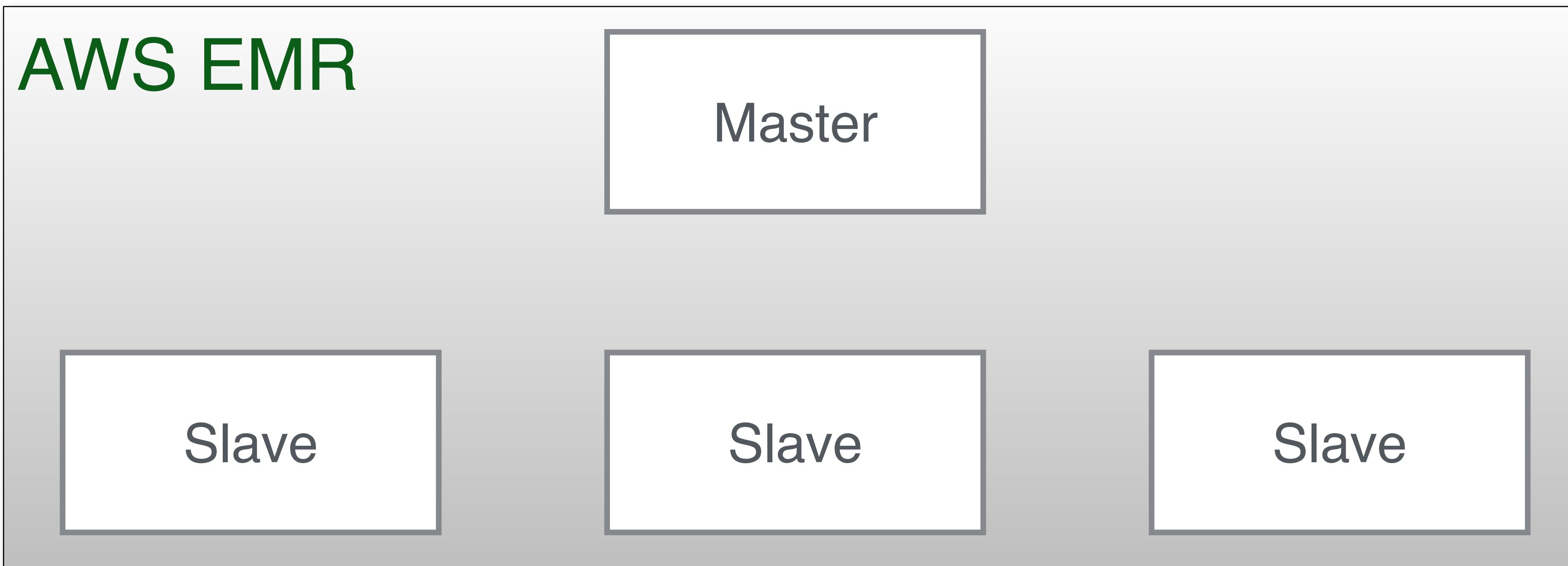
```
41 public static void main(String[] args) throws Exception {
42     final String confDir = System.getProperty("app.conf.dir");
43     final String environment = System.getProperty("newrelic.environment");
44     final DynamicConfiguration configuration = Configurations.configure(confDir, environment);
45
46     Properties kafkaConsumerProperties = getKafkaConsumerProperties(configuration);
47
48     logger.info("kafkaConsumerProperties: {}", kafkaConsumerProperties);
49
50     StreamExecutionEnvironment see = StreamExecutionEnvironment.getExecutionEnvironment();
51     see.setStreamTimeCharacteristic(TimeCharacteristic.EventTime);
52     see.getConfig().setAutoWatermarkInterval(1000);
53     see.setParallelism(40);
54
55     DataStream<TimesliceData> dataStream = see
56         .addSource(new FlinkKafkaConsumer08<>(
57             timesliceConstants.RESOLVED_TIMESLICE_DATA_KAFKA_TOPIC_NAME,
58             new TimesliceDeserializer(),
59             kafkaConsumerProperties)).setParallelism(16);
50
56
57
58
59
60
61 // aggregatedByMinuteStream contains events that are the aggregated timeslices by [MetricId,AgentId,AccountId]
62 DataStream<TimesliceTuple> aggregatedByMinuteStream = dataStream
63     .rebalance()
64     .flatMap(new TimesliceTupleMapper())
65     .assignTimestampsAndWatermarks(new TimesliceTupleTimestampExtractor())
66     .keyBy(TimesliceTuple.WIDE_METRIC_ID, TimesliceTuple.AGENT_ID, TimesliceTuple.ACCOUNT_ID)
67     .timeWindow(Time.seconds(60))
68     .reduce(TimesliceTuple::aggregateWith);
69
70 // packagedMinuteTimeslices are all the minute-aggregated timeslices for a given [AgentId,AccountId]
71 DataStream<OutputTimesliceData> packagedMinuteTimeslices = aggregatedByMinuteStream
72     .assignTimestampsAndWatermarks(new TimesliceTupleTimestampExtractor())
73     .keyBy(TimesliceTuple.AGENT_ID, TimesliceTuple.ACCOUNT_ID)
74     .timeWindow(Time.seconds(60))
75     .fold(null, new OutputTimesliceAggregator());
76
77 packagedMinuteTimeslices.addSink(new FlinkKafkaProducer08<OutputTimesliceData>(
78     AGGREGATED_MINUTE_TIMESLICE_DATA,
79     new OutputTimesliceDataSerializer(),
80     getKafkaProducerProperties(configuration))).setParallelism(16);
81
82 System.out.println(see.getExecutionPlan());
83
84 see.execute();
85 }
```



FLINK PLAN VISUALIZER

[Zoom In](#)[Zoom Out](#)

Spark configuration



```

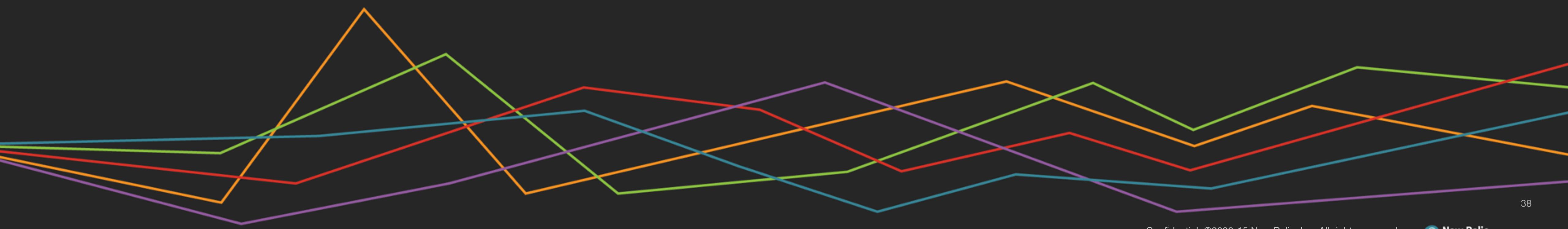
52     |     SparkConf conf = new SparkConf().setAppName("SparkTimesliceAggregator");
53     |     JavaStreamingContext jssc = new JavaStreamingContext(conf, Durations.seconds(1));
54
55     |     JavaPairInputDStream<byte[], TimesliceData> resolvedTimesliceStream = KafkaUtils.createDirectStream(jssc,
56     |             byte[].class, TimesliceData.class,
57     |             DefaultDecoder.class, TimesliceDeserializer.class,
58     |             getKafkaConsumerProperties(configuration),
59     |             (Set) ImmutableSet.of(RESOLVED_TIMESLICE_TOPIC));
60
61     |     JavaDStream<AggregatableTimeslice> resolvedTimeslicesFlattenedStream = resolvedTimesliceStream
62     |             .flatMap(new TimesliceTupleMapper());
63
64     |     JavaPairDStream<Tuple3<Long, Integer, Integer>, AggregatableTimeslice>
65     |             minuteAggregatedTimesliceStream = resolvedTimeslicesFlattenedStream
66     |             .mapToPair(new AggregatableTimesliceKeyMapper())
67     |             .reduceByKeyAndWindow(new TimesliceReducer(), Durations.minutes(1));
68
69     |     JavaMapWithStateDStream<Tuple2<Integer, Integer>, AggregatableTimeslice, OutputTimesliceData, OutputTimesliceData>
70     |             packagedAggregatedMinuteTimeslices = minuteAggregatedTimesliceStream.mapToPair(new PairFunction<Tuple2<Tuple3<Long,
71     |             Integer, Integer>, AggregatableTimeslice>, Tuple2<Integer, Integer>, AggregatableTimeslice>() {
72     |         ↑
73     |         @Override
74     |         public Tuple2<Tuple2<Integer, Integer>, AggregatableTimeslice> call(Tuple2<Tuple3<Long, Integer, Integer>,
75     |             AggregatableTimeslice> keyValuePair) throws Exception {
76     |             return new Tuple2<Tuple2<Integer, Integer>, AggregatableTimeslice>(
77     |                 new Tuple2<Integer, Integer>(keyValuePair._2.accountId, keyValuePair._2.agentId),
78     |                 keyValuePair._2);
79     |         }
79     |         ↑
80     |         }.mapWithState(StateSpec.function(new Function4<Time, Tuple2<Integer, Integer>, Optional<AggregatableTimeslice>,
81     |             State<OutputTimesliceData>, Optional<OutputTimesliceData>)()
82     |             @Override
83     |             public Optional<OutputTimesliceData> call(Time time, Tuple2<Integer, Integer> tuple, Optional<AggregatableTimeslice>
84     |                 wrappedValue, State<OutputTimesliceData> stateObject) throws Exception {
85     |                 AggregatableTimeslice value = wrappedValue.get();
86     |                 OutputTimesliceData newState = stateObject.exists() ?
87     |                     stateObject.get() :
88     |                     new OutputTimesliceData(value.getAccountID(),
89     |                         value.getAgentID(),
90     |                         value.getBeginTimeInSeconds(),
91     |                         value.getEndTimeInSeconds(),
92     |                         value.getMetricStoragePolicy());
93     |                 newState.appendTimeslice(value);
94     |                 stateObject.update(newState);
95     |                 return Optional.of(newState);
96     |             }
97     |         });
98
99     |         packagedAggregatedMinuteTimeslices.foreachRDD(outputTimesliceDataJavaRDD -> produceToKafka(outputTimesliceDataJavaRDD));

```

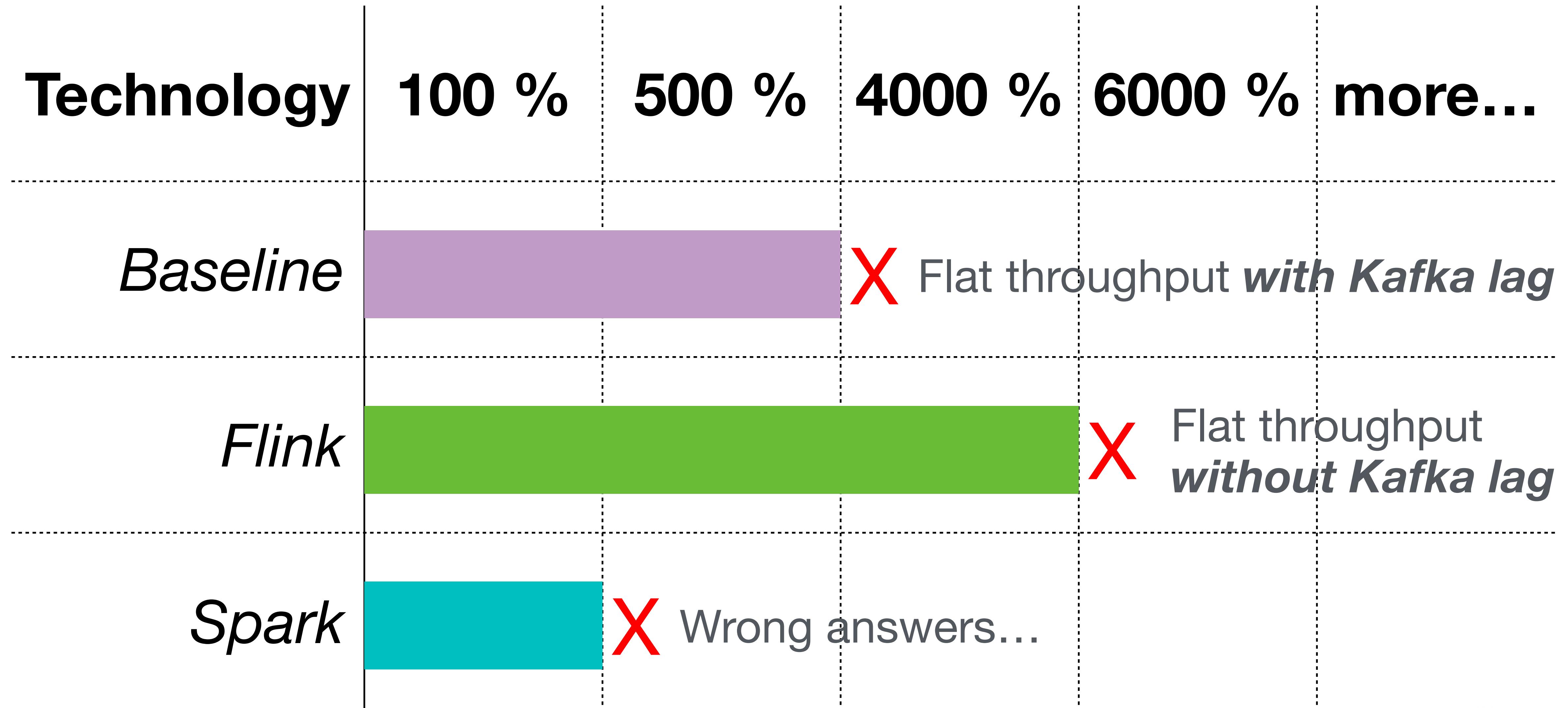
But the Spark Streaming solution generates **WRONG** results

- ... because there is no Event Time windowing
- ... leading to me abandoning Spark Streaming

Results



Results



Opportunities to improve the experiment

- **MORE BANDWIDTH**

- I don't know the limit of the Flink implementation

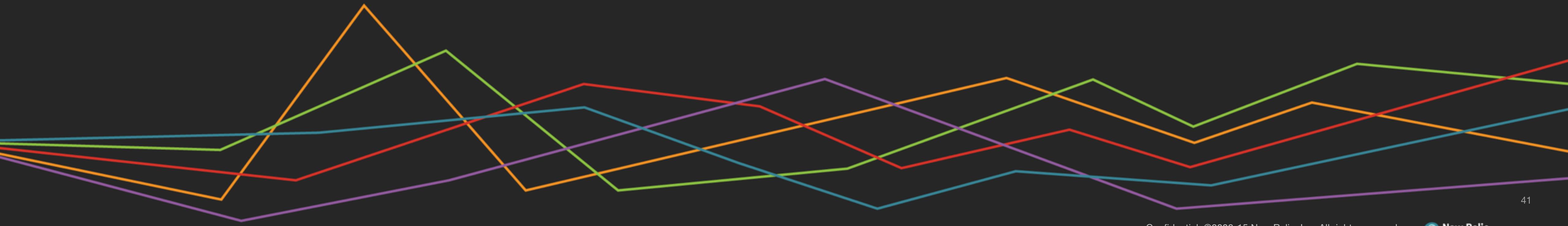
- **Key space domain expansion [All]**

- Scaling in rate domain only, with the same set of keys
 - This is *too easy* on the key-based systems [**Flink**, **Spark**]
 - This may be hard on the baseline system as well

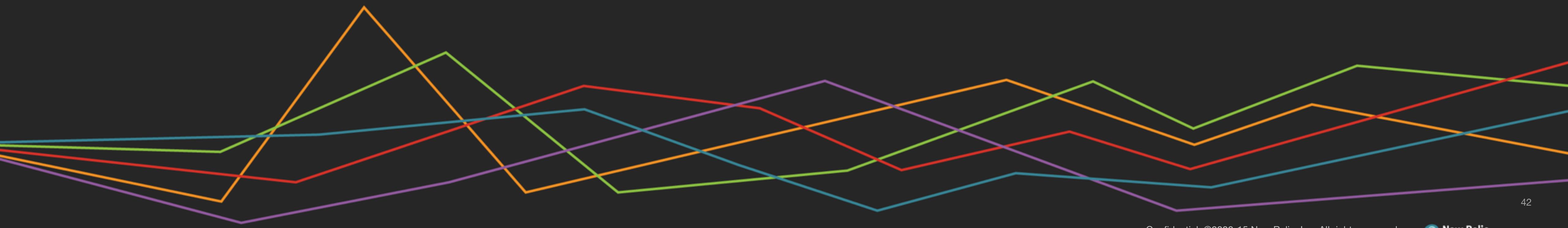
- **Inclusion of Database sinks [Flink, Spark]**

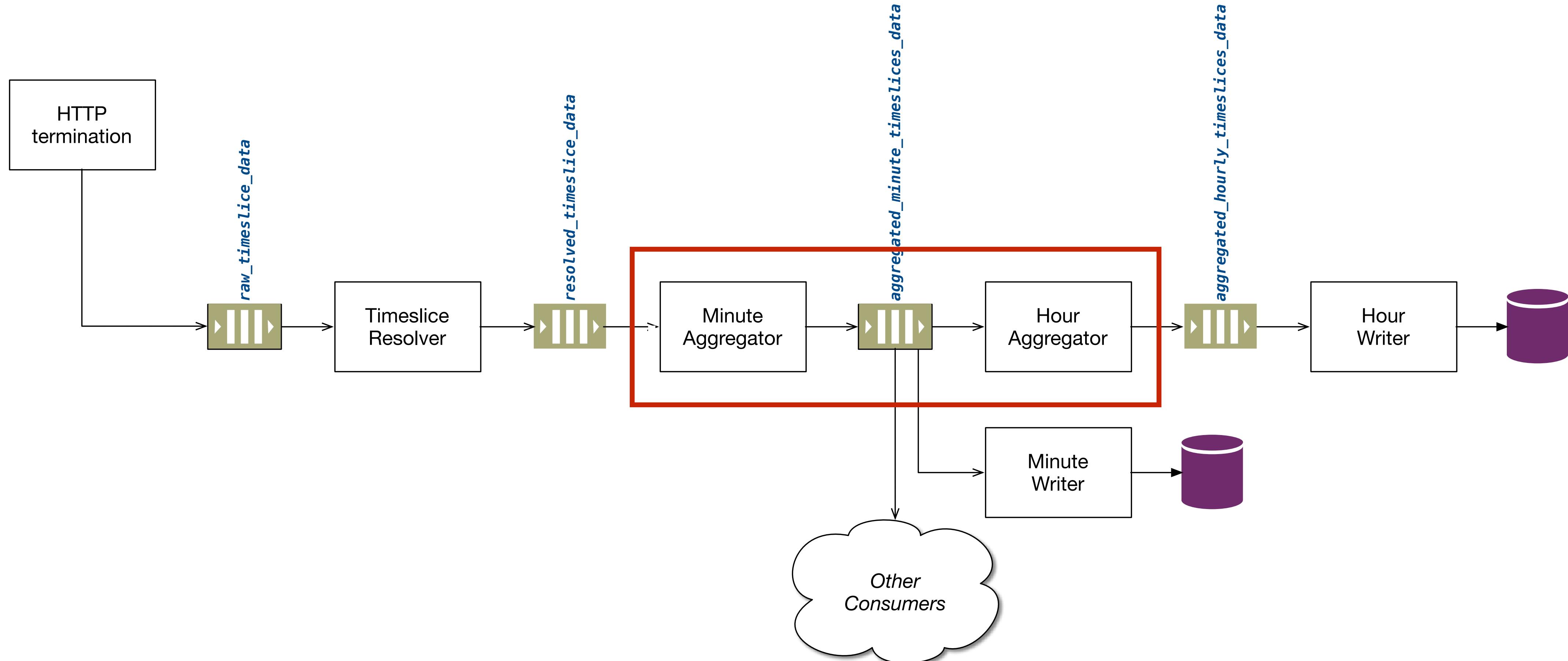
- Kafka sinks are still needed for downstream functions

Thank you



Extra credit

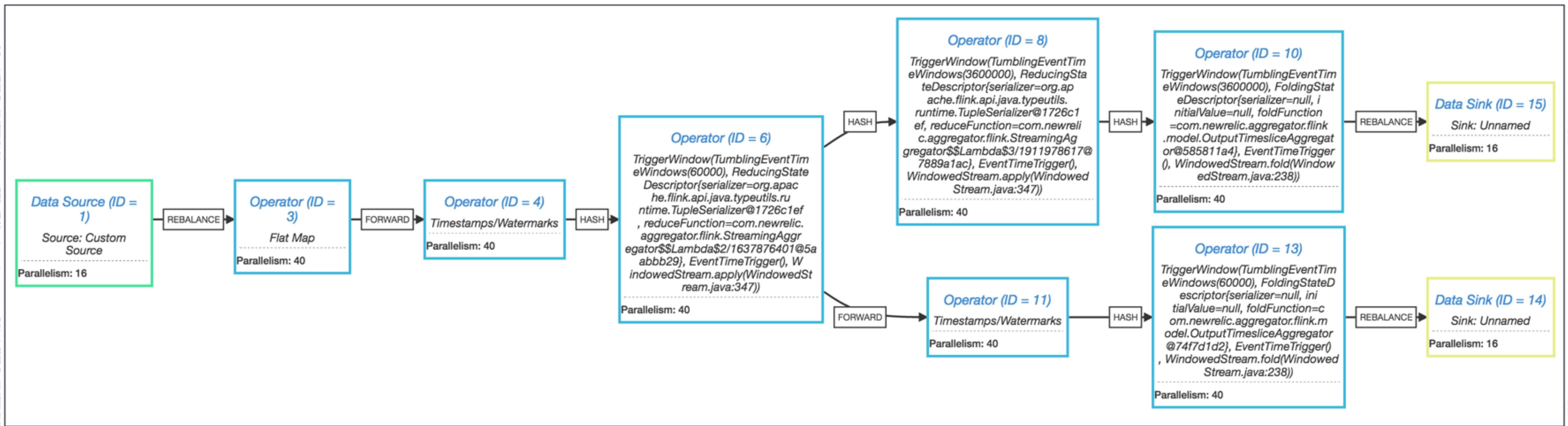




```
63
64    // aggregatedByMinuteStream contains events that are the aggregated timeslices by [MetricId,AgentId,AccountId]
65    DataStream<TimesliceTuple> aggregatedByMinuteStream = dataStream
66        .rebalance()
67        .flatMap(new TimesliceTupleMapper())
68        .assignTimestampsAndWatermarks(new TimesliceTupleTimestampExtractor())
69        .keyBy(TimesliceTuple.WIDE_METRIC_ID, TimesliceTuple.AGENT_ID, TimesliceTuple.ACCOUNT_ID)
70        .timeWindow(Time.seconds(60))
71        .reduce(TimesliceTuple::aggregateWith);
72
73    // packagedMinuteTimeslices are all the minute-aggregated timeslices for a given [AgentId,AccountId]
74    DataStream<OutputTimesliceData> packagedMinuteTimeslices = aggregatedByMinuteStream
75        .assignTimestampsAndWatermarks(new TimesliceTupleTimestampExtractor())
76        .keyBy(TimesliceTuple.AGENT_ID, TimesliceTuple.ACCOUNT_ID)
77        .timeWindow(Time.seconds(60))
78        .fold(null, new OutputTimesliceAggregator());
79
80    packagedMinuteTimeslices.addSink(new FlinkKafkaProducer08<OutputTimesliceData>(
81        AGGREGATED_MINUTE_TIMESLICE_DATA,
82        new OutputTimesliceDataSerializer(),
83        getKafkaProducerProperties(configuration))).setParallelism(16);
84
85    DataStream<TimesliceTuple> aggregatedByHourStream = aggregatedByMinuteStream
86        .keyBy(TimesliceTuple.WIDE_METRIC_ID, TimesliceTuple.AGENT_ID, TimesliceTuple.ACCOUNT_ID)
87        .timeWindow(Time.minutes(60))
88        .reduce(TimesliceTuple::aggregateWith);
89
90    DataStream<OutputTimesliceData> packagedHourTimeslices = aggregatedByHourStream
91        .keyBy(TimesliceTuple.ACCOUNT_ID, TimesliceTuple.AGENT_ID)
92        .timeWindow(Time.minutes(60))
93        .fold(null, new OutputTimesliceAggregator());
94
95    packagedHourTimeslices.addSink(new FlinkKafkaProducer08<OutputTimesliceData>(
96        AGGREGATED_HOUR_TIMESLICE_DATA,
97        new OutputTimesliceDataSerializer(),
98        getKafkaProducerProperties(configuration))).setParallelism(16);
99
```



FLINK PLAN VISUALIZER

[Zoom In](#)[Zoom Out](#)

Thank you

