

# Scaling Telemetry Systems

---

Reliability best practices for streaming data

**Liz Fong-Jones & Terra Field**

Field CTO  
@lizthegrey

Staff Platform Engineer  
@rainofterra



# Scaling Telemetry Systems

---

Reliability lessons for streaming data

**Liz Fong-Jones & Terra Field**

Field CTO  
@lizthegrey

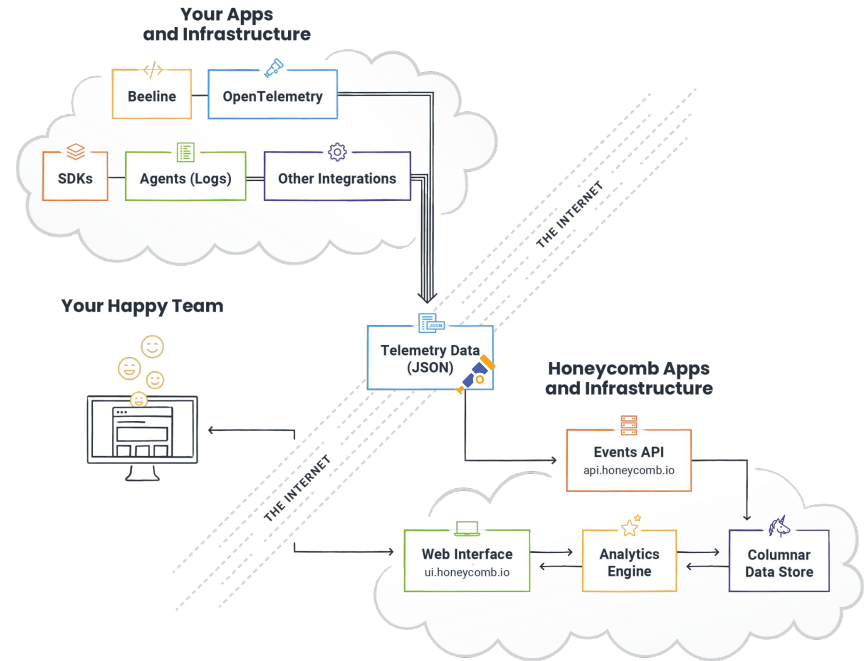
Staff Platform Engineer  
@rainofterra



# Our data ingest, storage, & analytics use case

## What Honeycomb does

- Ingests customer's telemetry
- Indexes on every column
- Enables near-real-time querying on newly ingested data



# Our data ingest, storage, & analytics use case

---

What Honeycomb *really* does



“

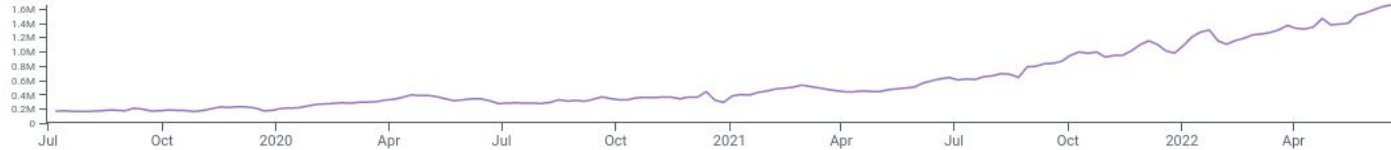
**Kafka is the beating heart of Honeycomb.**

Ben Hartshorne, first employee at Honeycomb

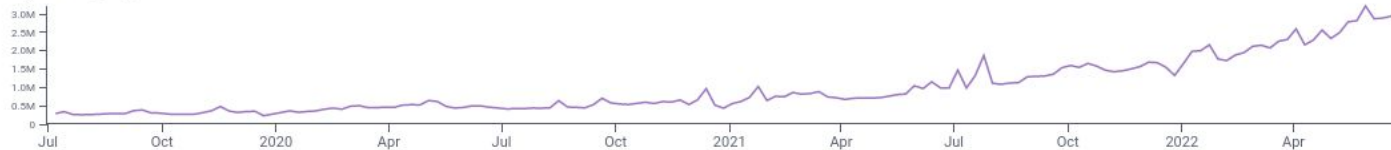
# 10x growth in three years

Jul 1 2019, 12:00 AM – Jun 26 2022, 12:00 AM (Granularity: 1 day)

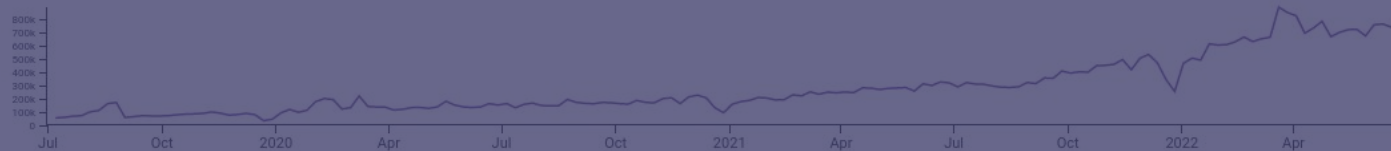
AVG(avg\_ingest)



MAX(peak\_ingest)



SUM(human\_queries)



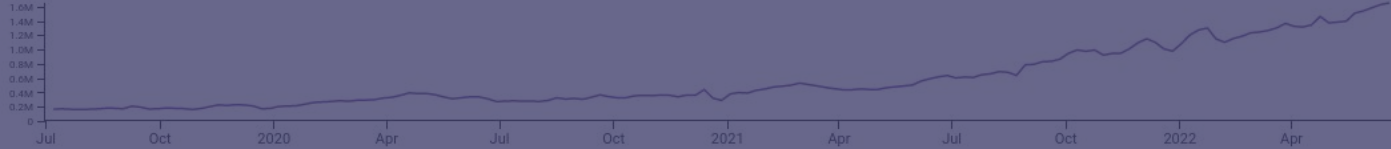
SUM(triggers\_per\_minute)



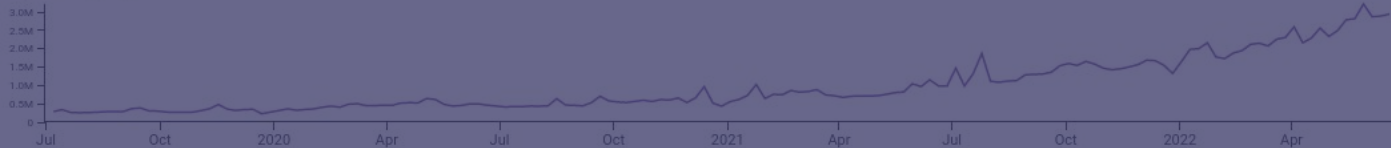
# 10x growth in three years

Jul 1 2019, 12:00 AM – Jun 26 2022, 12:00 AM (Granularity: 1 day)

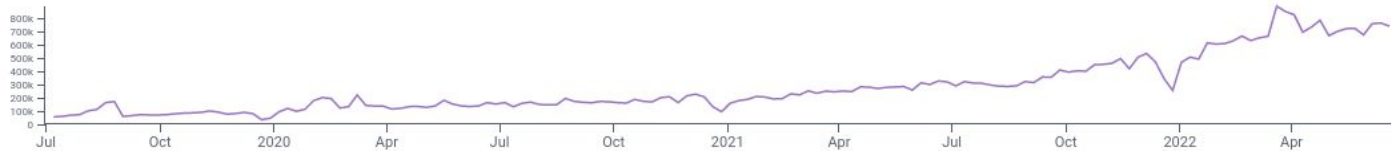
AVG(avg\_ingest)



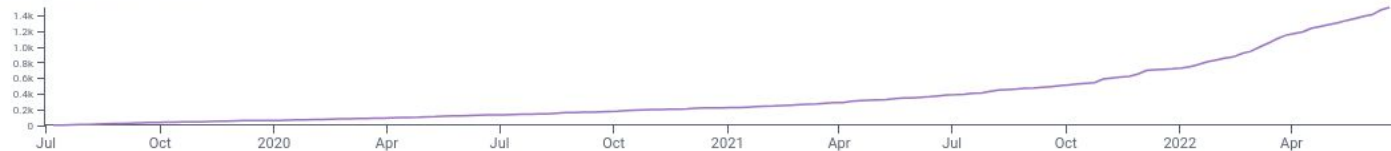
MAX(peak\_ingest)



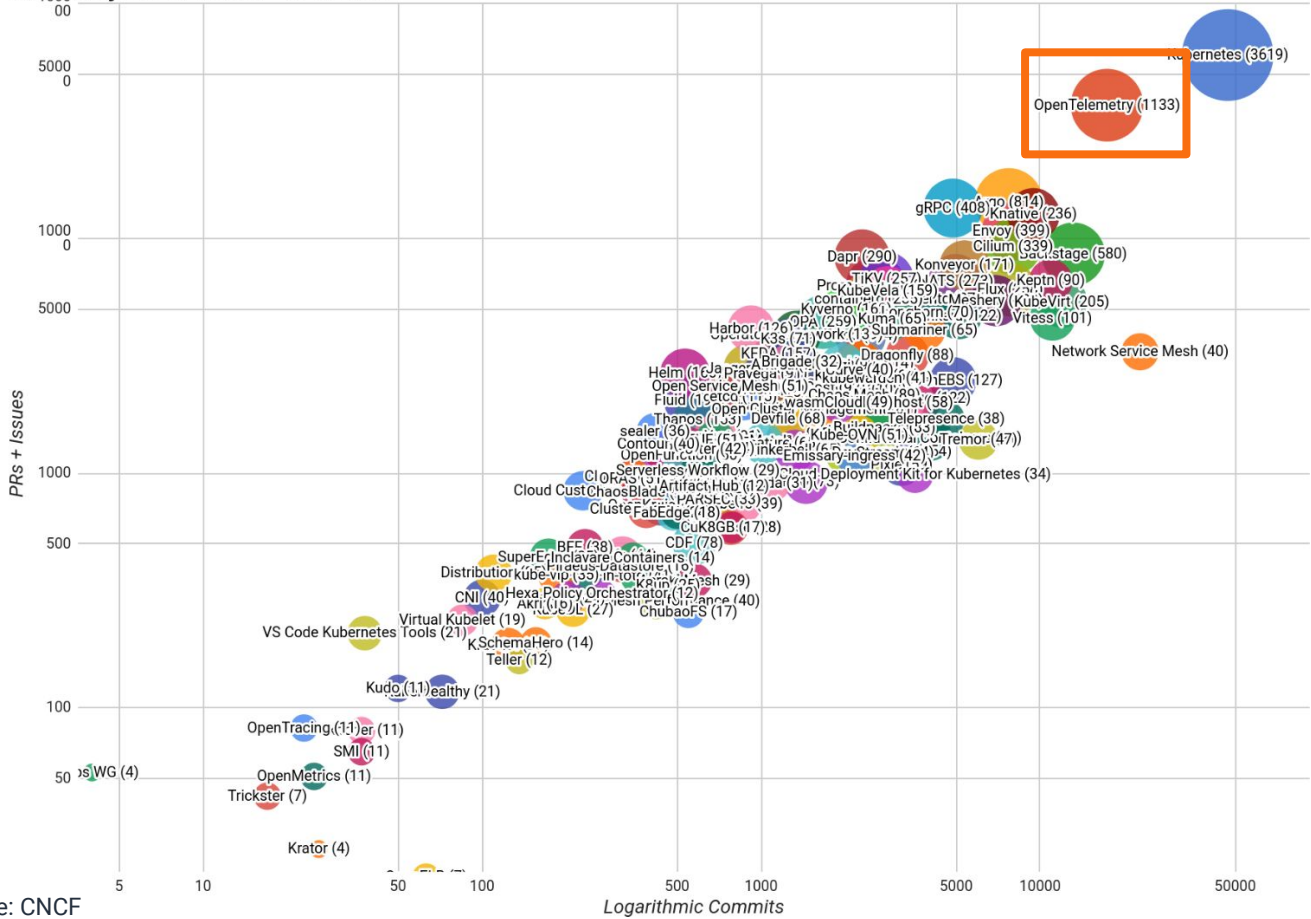
SUM(human\_queries)



SUM(triggers\_per\_minute)



# CNCF Projects 8/1/2021 - 8/1/2022



- Kubernetes (kubernetes.io) 3619 authors
- OpenTelemetry (opentelemetry.io) 1133 authors
- Argo (argoproj.github.io) 814 authors
- Backstage (backstage.io) 580 authors
- Prometheus (prometheus.io) 424 authors
- gRPC (grpc.io) 408 authors
- Envoy (www.envoyproxy.io) 399 authors
- Cilium (cilium.io) 339 authors
- Dapr (dapr.io) 290 authors
- Fluentd (fluentd.org) 275 authors
- NATS (nats.io) 273 authors
- OPA (openpolicyagent.org) 259 authors
- Flux (github.com/fluxcd) 258 authors
- TiKV (tikv.org) 257 authors
- containerd (containerd.io) 255 authors
- Knative (knative.dev) 236 authors
- Meshery (layer5.io/meshery) 209 authors
- KubeVirt (kubevirt.io) 205 authors
- Falco (falco.org) 194 authors
- Fluid (github.com/fluid-cloudnative/...) 194 authors
- Konveyor (konveyor.io) 171 authors
- Kyverno (kyverno.io) 162 authors
- Helm (helm.sh) 160 authors
- KubeVela (kubevela.io) 159 authors
- KEDA (keda.sh) 157 authors
- External Secrets Operator (external-secrets.io) 157 authors
- Thanos (thanos.io) 153 authors

Source: CNCF





# DATA GENERATION

# STREAMING INGEST

# STORAGE + PROCESSING

# VISUAL ANALYSIS LOOP

## Logging Data

Amazon Relational Database Service (Amazon RDS)	AWS Elastic Load Balancing
AWS Elastic Beanstalk	Amazon Elastic Kubernetes Service (Amazon EKS)

## Metrics Data

Host Metrics	App Metrics
Amazon CloudWatch	Prometheus

## Trace Span Data

JavaScript	Go
Jaeger	Java
Ruby	Python
.NET	Front-End

Realtime Ingest  
*Auth and validation*

**Honeycomb API**

OpenTelemetry

**Refinery**  
*User-controlled dynamic tail sampling*

Unpacked Into Columns

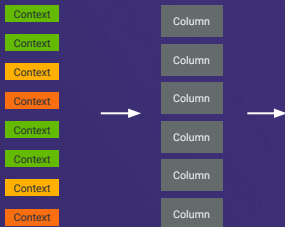
## Query Engine

Unlimited users can store thousands of dimensions at no additional cost and query any arbitrary combinations without pre-aggregation.

Proprietary distributed computing and parallelized processing returns query results in < 3 seconds across billions of rows of data.

## Columnar Datastore

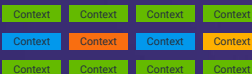
High-cardinality column file      Segment with time range



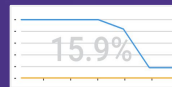
Amazon Simple Storage Service (Amazon S3)



Events contain many fields and distinct values



"Wide events" packed with as much context as you need for debugging



**Graph Rendering**  
Click through visuals based on granular data



**Query Builder UI**  
Intuitive GUI with multiple group-by



**Query Result History**  
Shared team intelligence



# Agenda

---

## **When to use streaming systems**

e.g. Honeycomb's use case (and yours?)

## **Bottlenecks to streaming**

Producers? Brokers? Consumers?

## **Patterns for scaling**

How to address each bottleneck

## **Patterns for observing**

Because you can't run it unless you can see it

**TL;DR: tricky to get right, but more supportable than request/response.**



## Liz Fong-Jones

Field CTO, honeycomb.io



## Terra Field

Staff Platform Engineer, honeycomb.io



# **What streaming is good for**

---

How streaming powers Honeycomb's telemetry pipeline

# Streaming data decouples systems

---

## Separation of stateless & stateful

Update producers & consumers  
on-demand without dropping data.  
(mostly! SLO != 100%)

Keep one single record of truth.

## Multiple fan-out on the event bus

This need came later, but was incredibly helpful.

- Originally: one producer, one consumer.
- Now: two producers, three consumers



# Now all the state lives in one place...

---

## This is both good...

- Rolling restarts of everything else.
- Replay in case of incorrect consumer behaviour

## ... and scary.

- If it breaks, everything breaks.
- Running third party software = harder to debug/understand



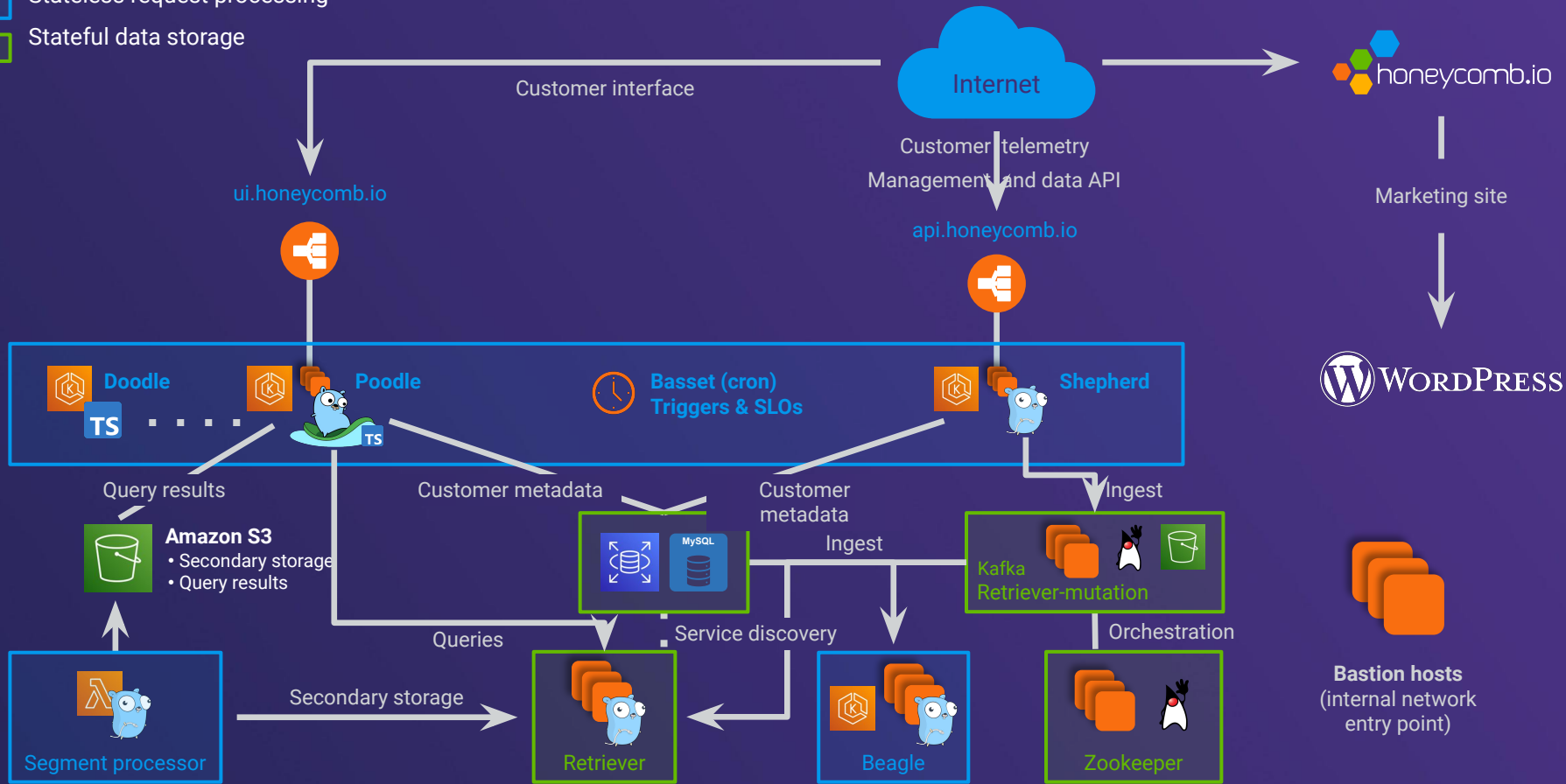
# But here's how our use case is atypical

---

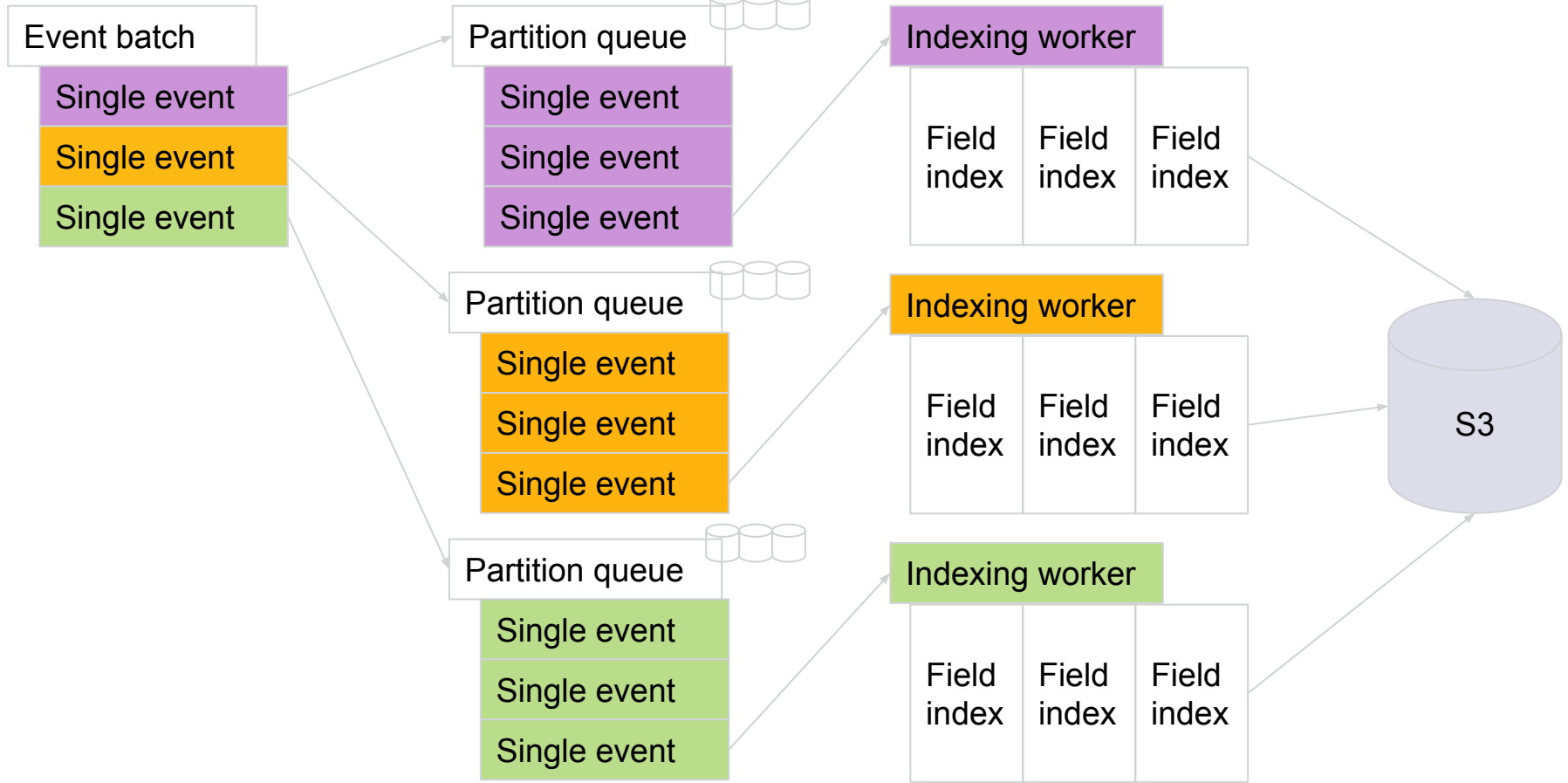
- We don't keep weeks or months of history
- We only typically read the last hour
- 1-2 main topics, not hundreds
- Self-managed partition allocation (~100 partitions)
- High throughput per partition (50k msgs/sec/partition)
- No ksql
- No librdkafka; pure Go Shopify Sarama (Shopify team, we owe you a drink)

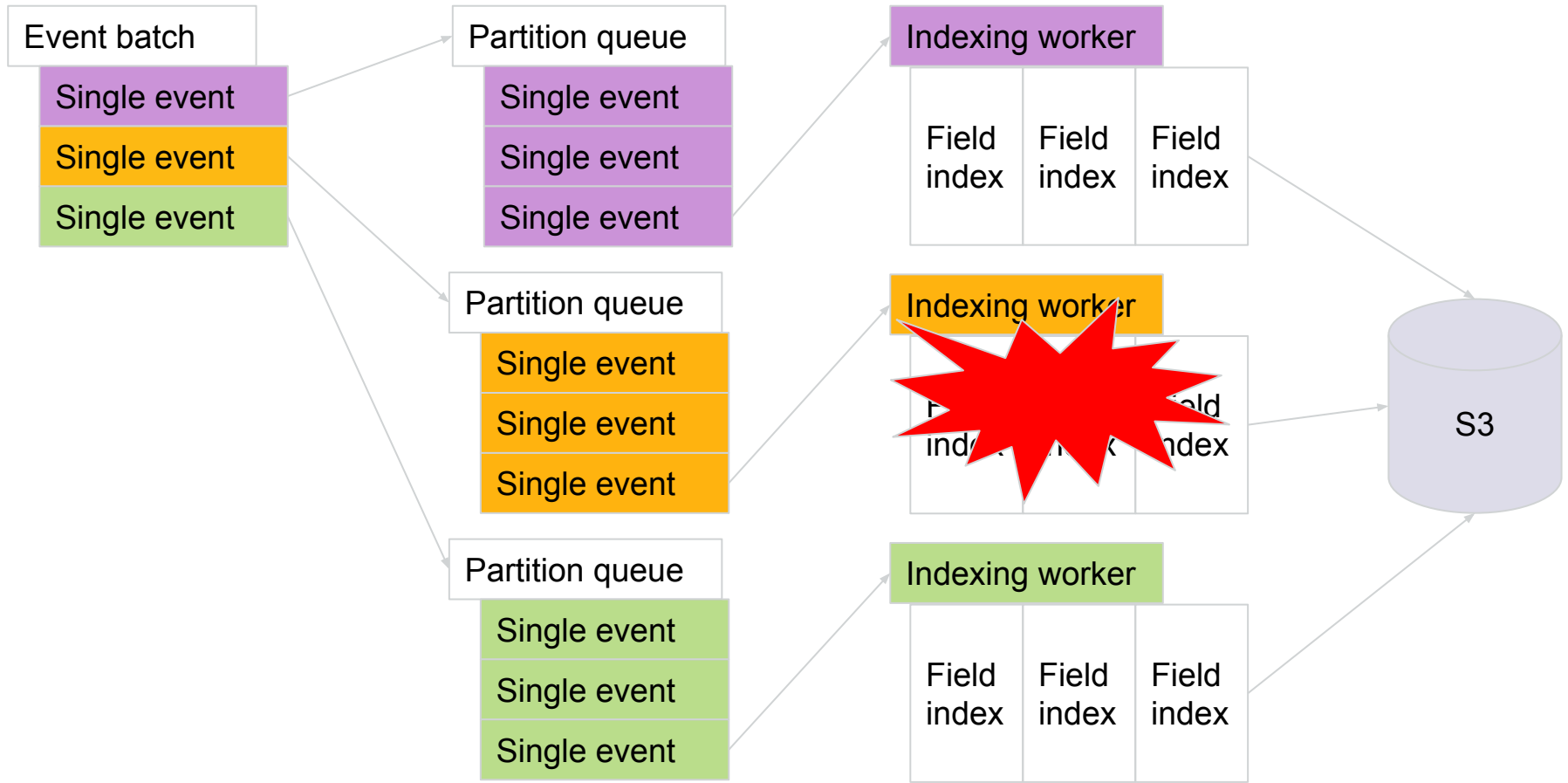


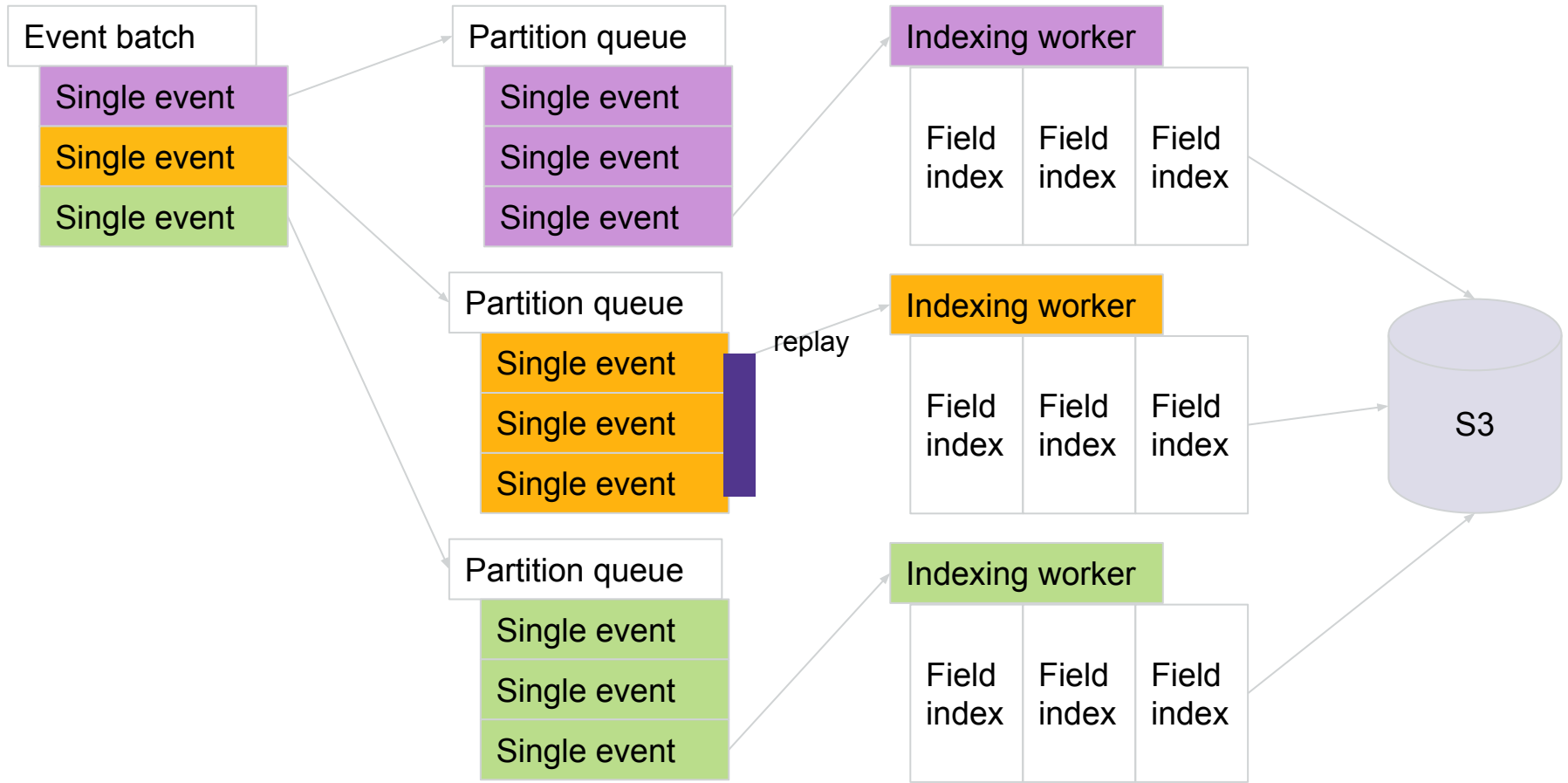
- Stateless request processing
- Stateful data storage











# Identifying & solving bottlenecks

---

What happens when the system exceeds constraints?

# Per-producer limits

---

- Least outstanding requests LB for ingest
- Tune batch sizes (MB, seconds), queue depths, etc
- Guard against OOMs
- Avoid persistently bad partitions
- Allocate load between partitions



# Broker limits

---

- CPU
  - Use Graviton (ideally im4gn/is4gen)
- On-disk storage
  - Use smaller NVMe for predictable latency
  - Use tiered storage for bulk, less frequently accessed
    - DIY tiering w/ writeback cache does not work. We tried it.
  - We do *not* recommend use of EBS volumes for scaling out (latency, cost)
- Auto-balancing
- Horizontal scale-out (if needed)



# Per-consumer limits

---

- Annoying: Kafka limits consumer-broker BW regardless of distinct partitions.
- Run more brokers (or map consumers:partitions 1:1) if all else fails.
- Watch out for consumer group rebalancing
  - Rolling restarts are a good SRE practice everywhere EXCEPT here
  - (but this advice may change with sticky consumer group assignment)
- You are only as good as your offset commit



# Optimizations to consider

---

How to get the most out of your cluster



# Consider not running brokers yourself

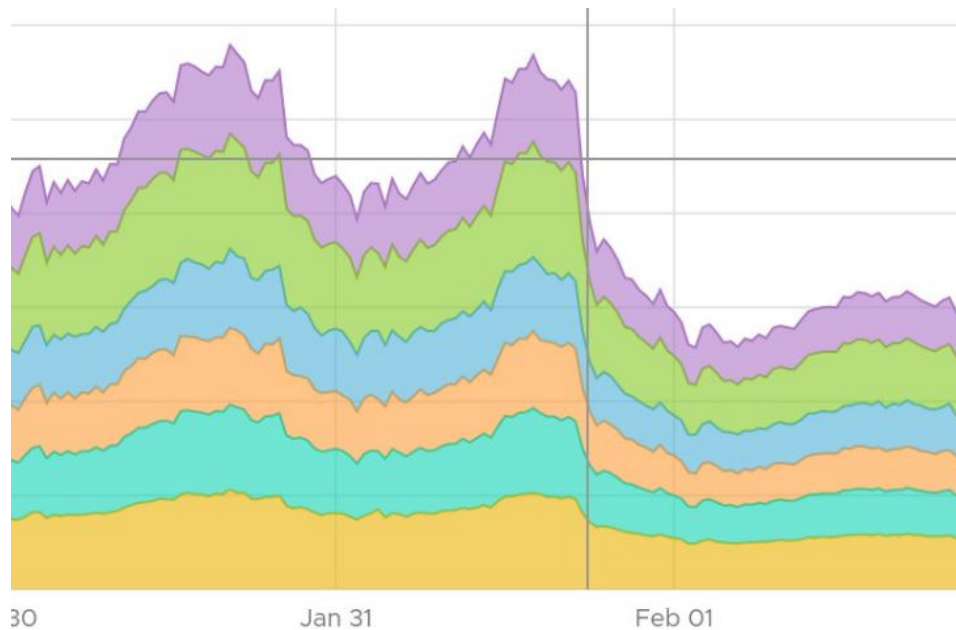
---

- Confluent Cloud, Google Cloud Pub/Sub, Amazon Kinesis, etc.
- **Run. Less. Software.**



# Use zstd. Seriously.

- CPU is ~cheaper than network.
- 20%+ savings on bandwidth vs snappy



# Fully utilise Kafka's replication

---

- DTAZ = \$\$\$\$
- Don't re-copy data between AZs; read from followers
- Don't pay for more durability than you need (Kafka already provides R=3)



# Use the most efficient base you can

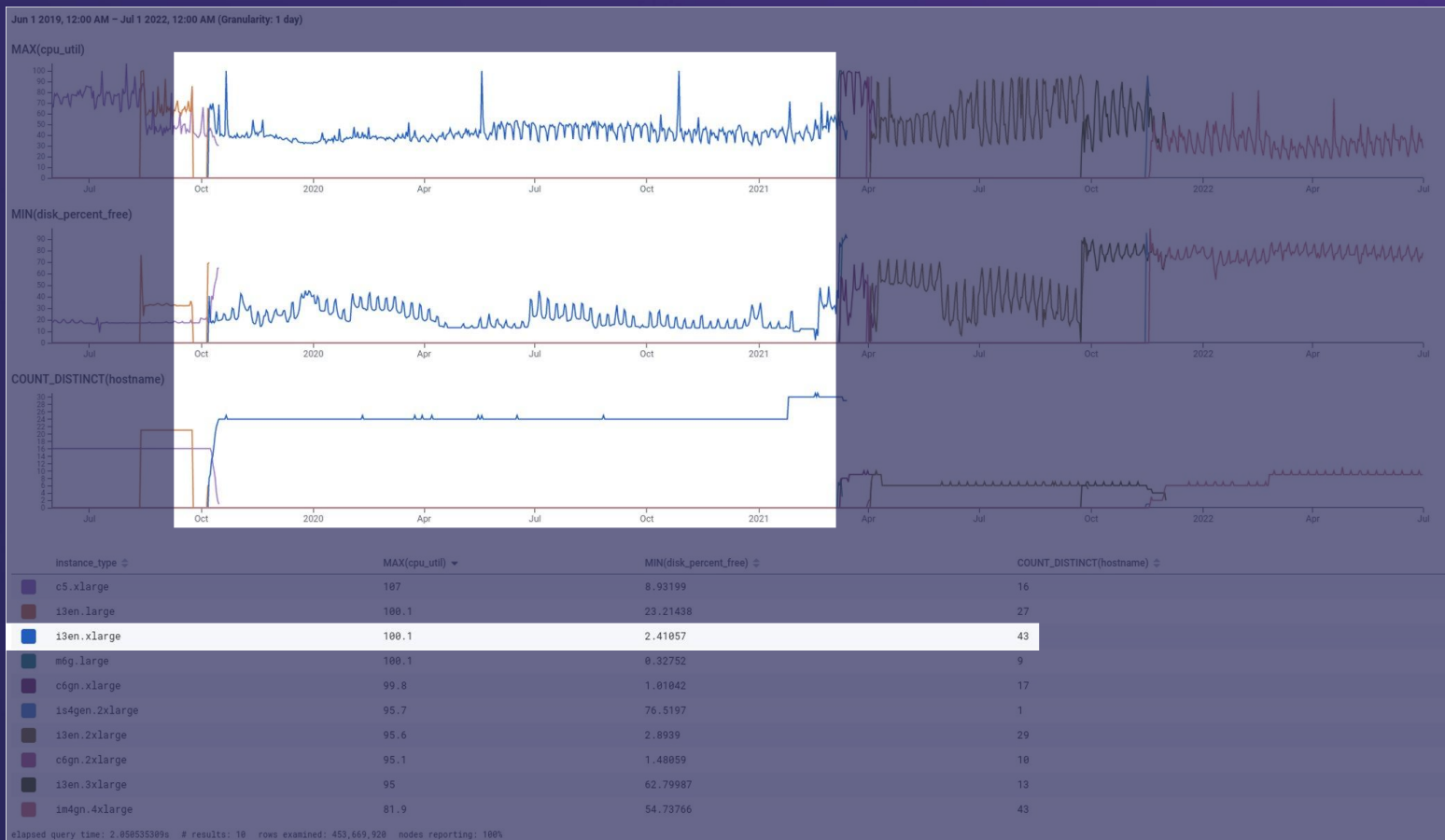
---

## Hardware

- Be careful of unknown unknowns
  - There may be dimensions you're unaware of.
- Burst balances create metastable systems
  - There may be dimensions you're unaware of.



# Finding the right way to migrate Kafka



# Our month of Kafka pain

---

Longtime Confluent Kafka users

First to use Kafka on Graviton2 at scale

## Changed multiple variables at once

- move to tiered storage
- i3en → c6gn
- AWS Nitro



Read more: [go.hny.co/kafka-lessons](https://go.hny.co/kafka-lessons)



# Unexpected constraints

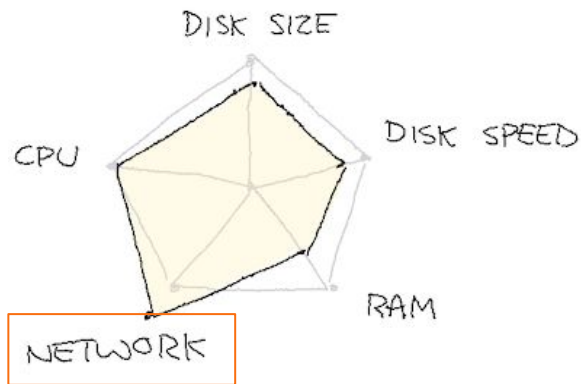
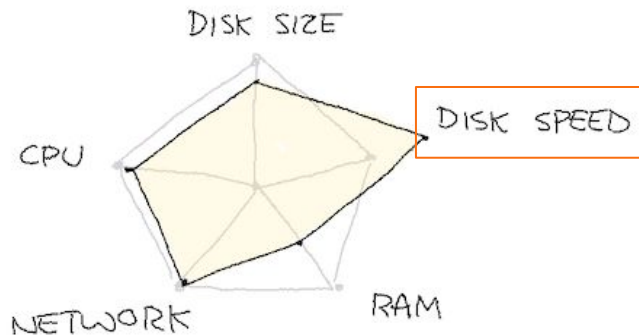
We thrashed multiple dimensions.

We tickled hypervisor bugs.

We tickled EBS bugs.

Burning our people out wasn't worth it.

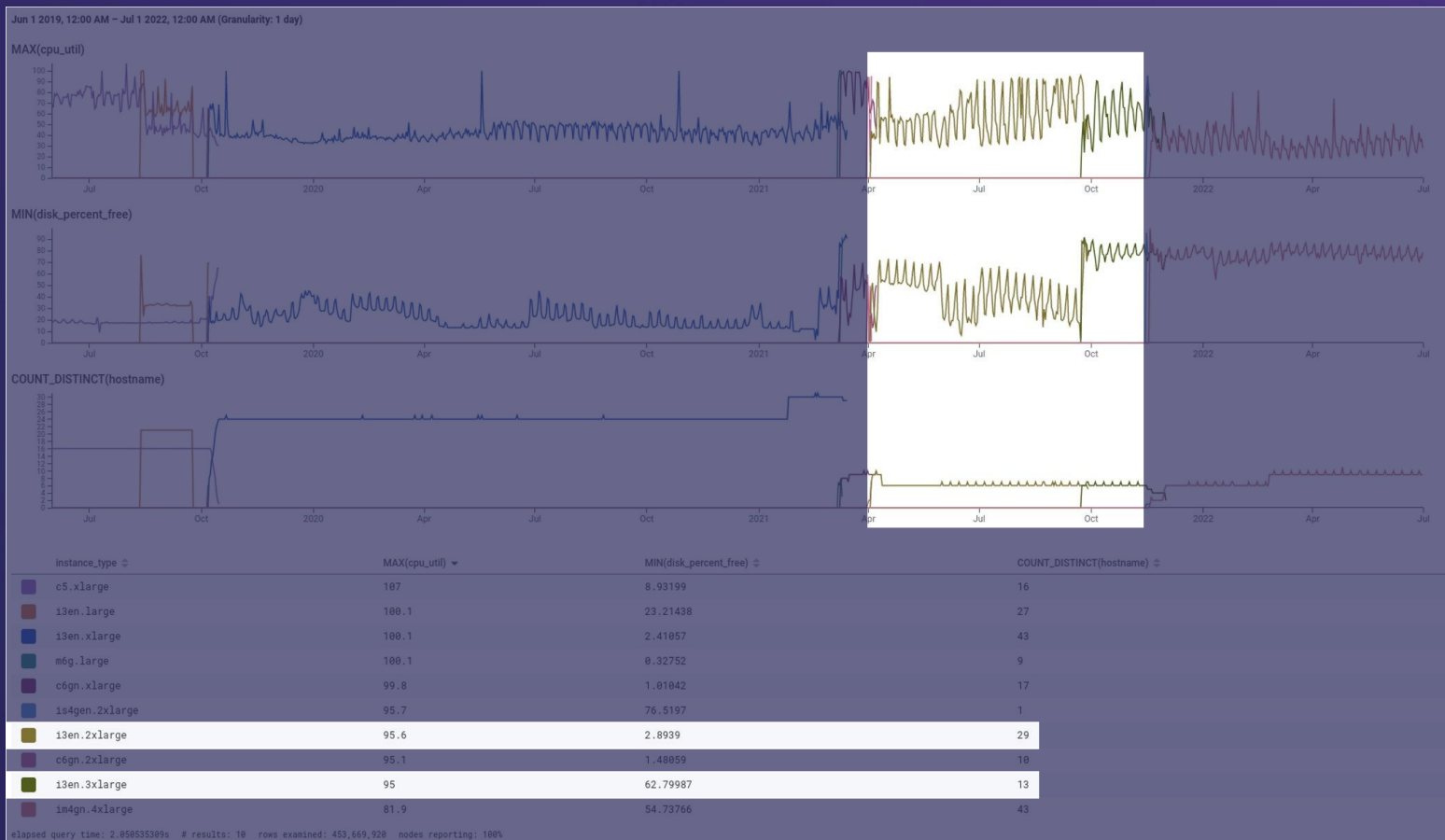
***But we were finally able to move forward in Dec 2021 with im4gn!***



Read more: [go.hny.co/kafka-lessons](https://go.hny.co/kafka-lessons)

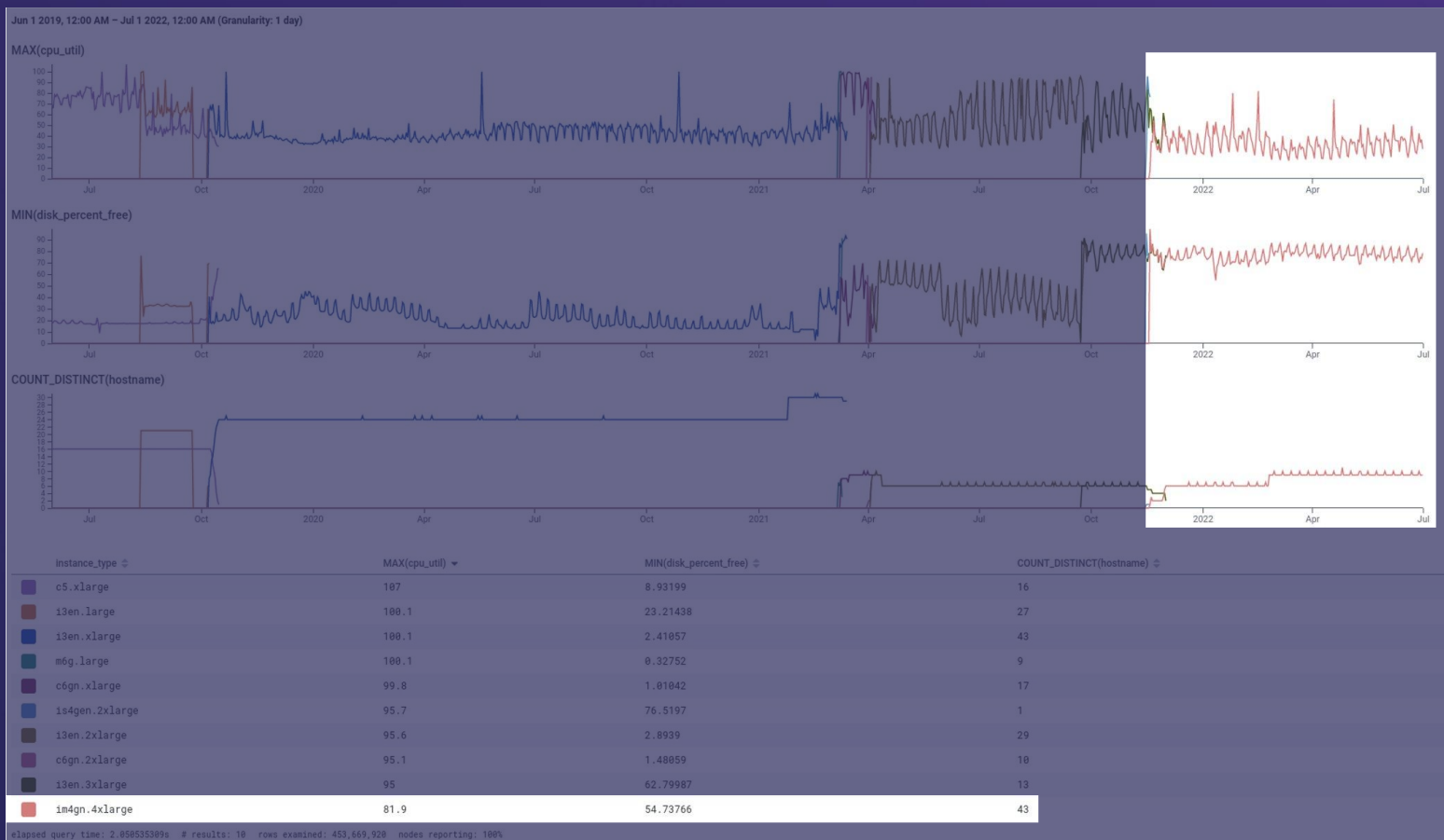


# Finding the right way to migrate Kafka





# Finding the right way to migrate Kafka



# Use the most efficient base you can

---

## Hardware

- Be careful of unknown unknowns
  - There may be dimensions you're unaware of.
- Burst balances create metastable systems
  - There may be dimensions you're unaware of.

## Software

- Profile, profile, profile.
- Use Corretto JVM, not GetOpenJDK
- Use a well-tuned GC algorithm
- Upgrade your JNI deps (eg Zstd)
- Replace Java crypto libraries with AWS Corretto Crypto Provider



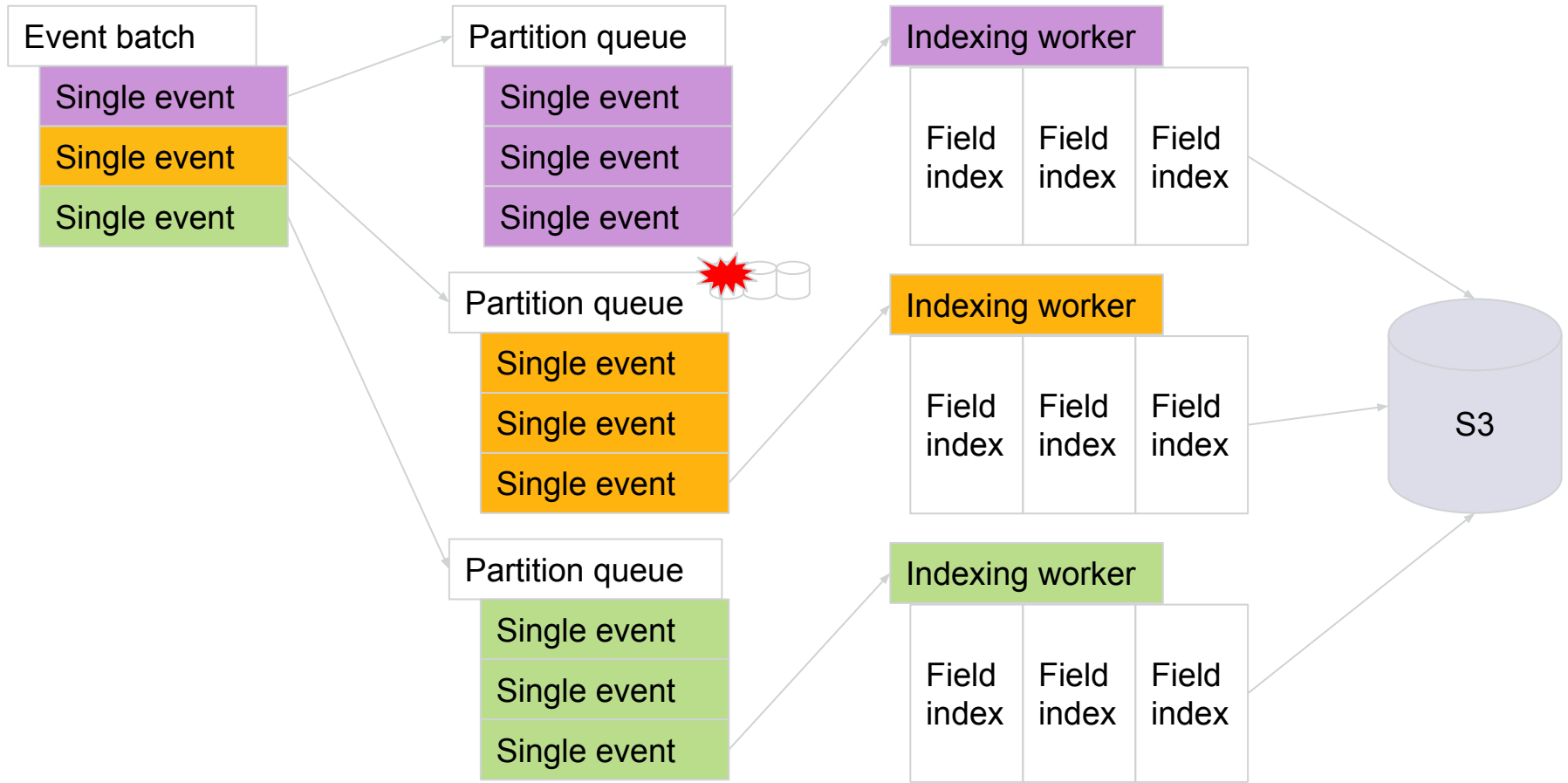


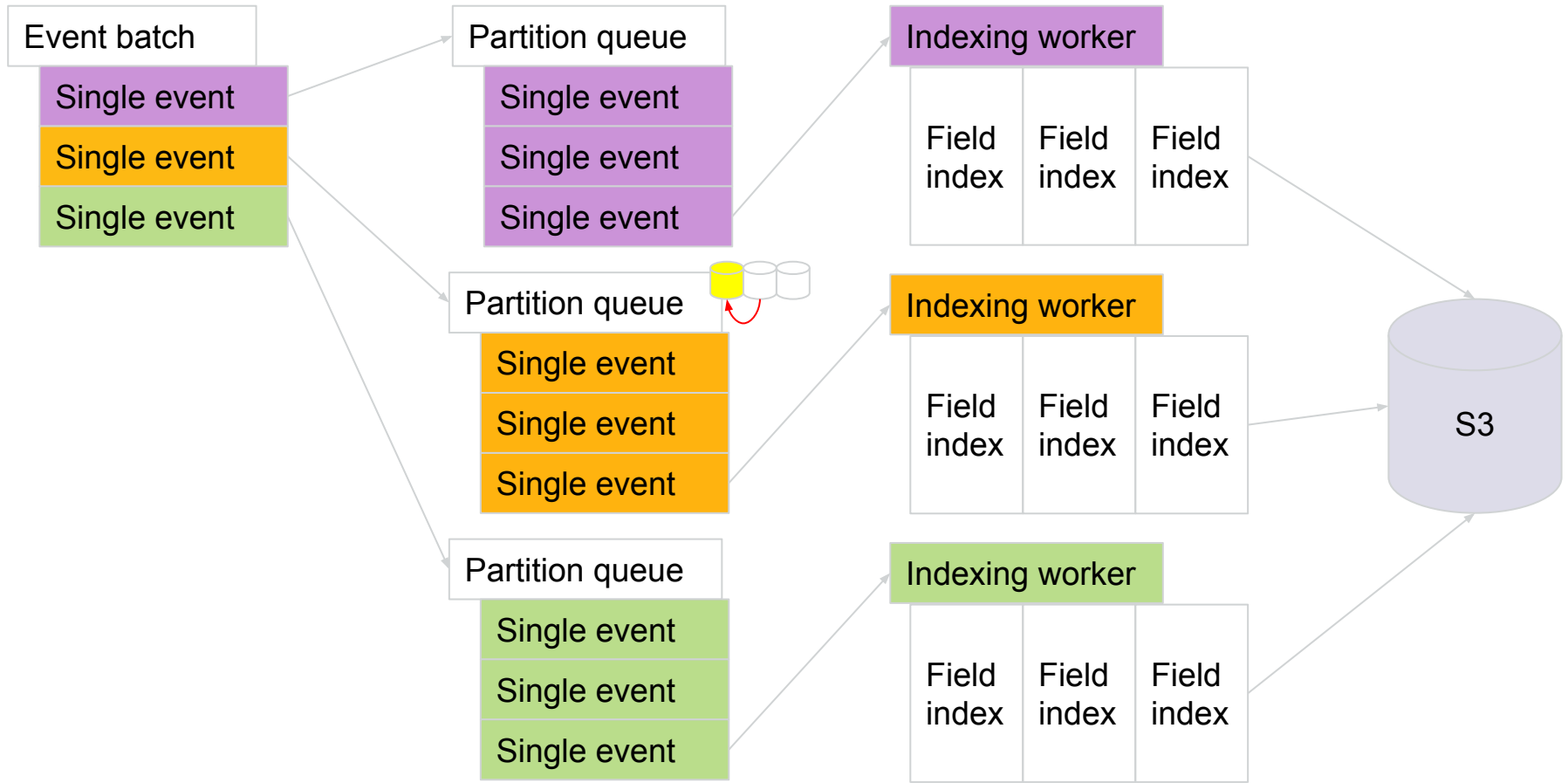
# Continuously chaos test your DR strategy

---

- Weekly consumer and broker replacement to verify cold start
- Remember: an untested backup is *not* a restore.
- Leave plenty of room for unexpected scenarios
- In streaming, headroom = *time* before pear-shaped







# Observing streaming systems

---

How we make sure everything is working correctly

# Application vs system observability

---

## App level: trace spans & links

- Periodic trace spans per consumer ("tick")
- Heavily sampled produce requests
- Trace links between consume & produce

## System level: broker metrics

- Basics: Msgs/sec, CPU, URP, Disk
- Advanced: GC, Network, Controller, Rebalancing

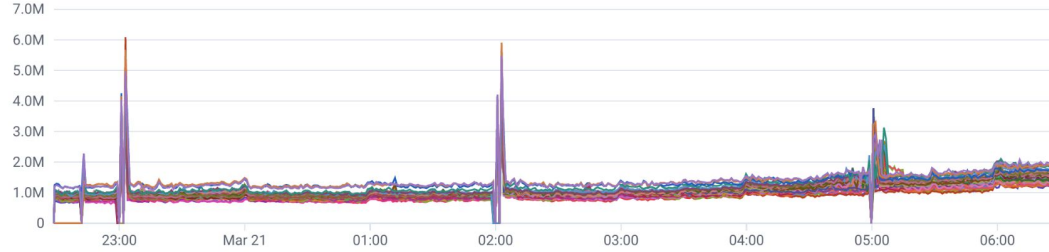




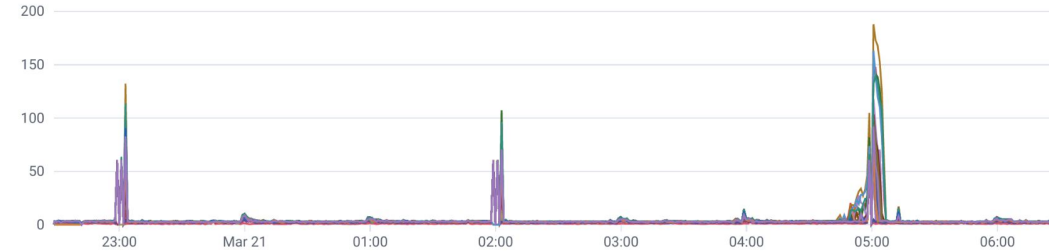


Mar 20 2023 22:28:44 – Mar 21 2023 06:28:44 UTC-07:00 (Granularity: 1 min)

SUM(consumer.messages.count)



MAX(delay.duration\_sec)



meta.local_hostname	SUM(consumer.messages.count)	MAX(delay.duration_sec)
beagle-667bc7b6c5-z9q4m	675,889,923	92.03093
beagle-667bc7b6c5-fb7rr	658,901,116	76.23326





← Trace 3224b822936253f854e89a84bdf7bd4e at 2023-03-21 03:27:40

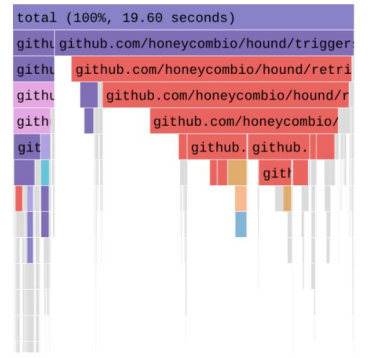
Rerun

beagle >  
tick

Search spans  0 spans with errors

name	service.name	0s	20s	40s	60.75s
1K+ tick	beagle	60.753s			
1 getSLOsForDataset	beagle	1.900ms			
11 getSLOsForDataset	beagle	12.97ms			
15 getSLOsForDataset	beagle	17.96ms			
worker	beagle	58.113s			
5K+ getSLOsForDataset	beagle	644.1ms			
worker	beagle	58.095s			
worker	beagle	57.442s			
worker	beagle	57.415s			
worker	beagle	57.414s			
worker	beagle	57.404s			
worker	beagle	57.396s			
worker	beagle	57.388s			
worker	beagle	57.380s			
worker	beagle	57.354s			
worker	beagle	57.354s			
worker	beagle	57.354s			

Latencies Profile Self Profile Children

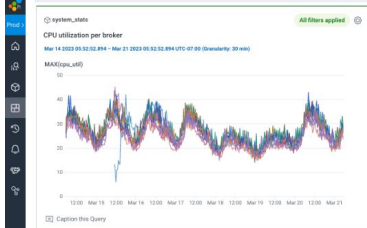


Fields Span events (128) Links (0)

Filter fields and values in span

- Timestamp 2023-03-21T10:27:40.010006877Z
- consumer.generation\_duration\_sec





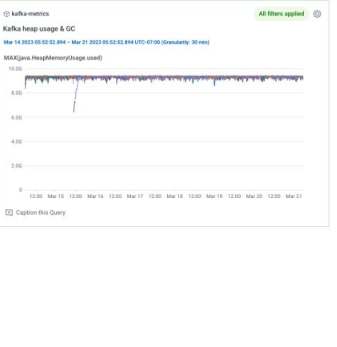
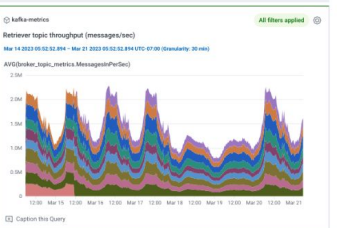
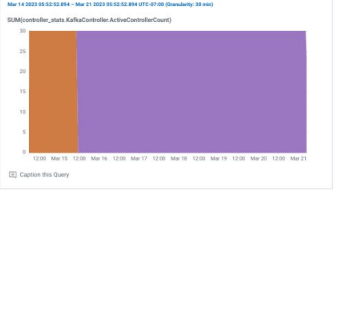
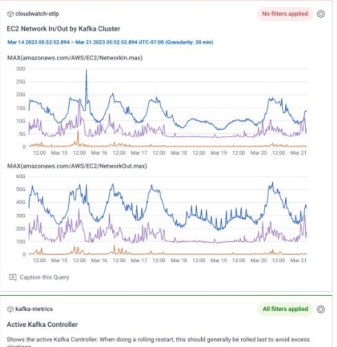
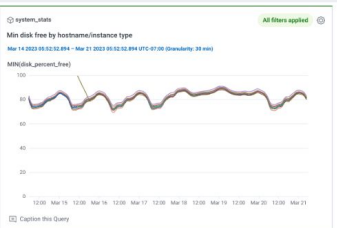
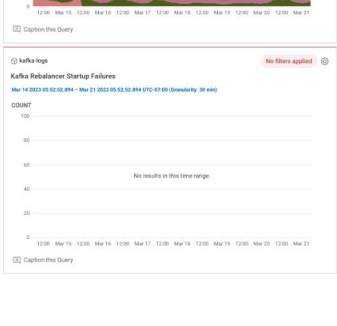
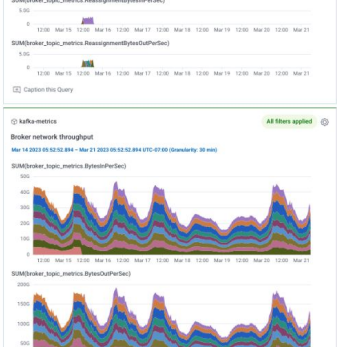
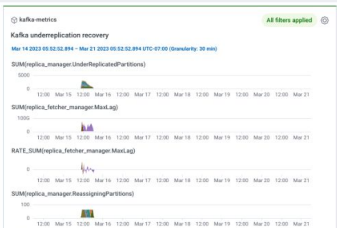
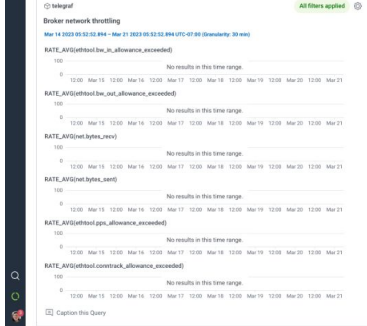
### kafka-metrics

Kafka broker metadata

Mar 14 2023 05:52:52.894 - Mar 21 2023 05:52:52.894 UTC-07:00 (Sunday: 30 min)

kafka_broker_id	host	env	env_instance	env	MAX(overridden)
1365	kafka-00174764a6d6384	us-east-1	saegm-4c1arge	production	7
1363	kafka-nd3rF86k0a7nF2	us-east-1	saegm-4c1arge	production	7
1364	kafka-088174763293926	us-east-1	saegm-4c1arge	production	7
1363	kafka-8ba251997164821	us-east-1	saegm-4c1arge	production	7
1362	kafka-020772c0d74b5ca	us-east-1	saegm-4c1arge	production	7
1361	kafka-8f5784c44627275	us-east-1	saegm-4c1arge	production	7
1368	kafka-5a08576a5c109857	us-east-1	saegm-4c1arge	production	7
1358	kafka-0513878271191568	us-east-1	saegm-4c1arge	production	7
1357	kafka-061786546039260	us-east-1	saegm-4c1arge	production	7
1356	kafka-0c5764e0d870d112	us-east-1	saegm-4c1arge	production	7

📄 Caption this Query



# Application vs system SLOs

---

## App level: many SLOs

- SLOs on producer write success/latency
- SLOs on consumer freshness per-consumer
- Implied SLO for durability (~never lose data)

## System level: No SLOs.

- Because we use the producer/consumer views instead!



# Future work

---

Where we want to go next

# tl;dr better balancing & auto-healing

---

Broker self-balancing sometimes sticks

Leadership imbalance causes CPU anomalies/perf pain

Long-term partition imbalance causes operational pain

Short-term spikes on specific partitions from new customers cause pain

**Kafka is the beating heart, but should not produce toil.**

*Non-goal: k8s. None of these problems are things k8s would solve for us.*





# Visit our booth!

---

[hny.co/srecon23-americas](https://hny.co/srecon23-americas)

@lizthegrey & @rainofterra



# Observability Engineering

Get our new book, free!

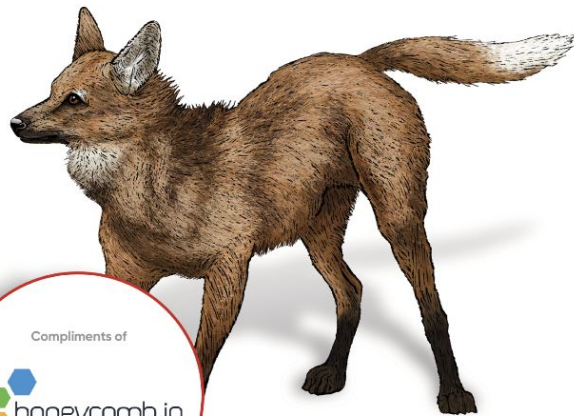


@lizthegrey

O'REILLY®

## Observability Engineering

Achieving Production Excellence



Compliments of



Charity Majors,  
Liz Fong-Jones  
& George Miranda





# Questions?

---



[www.honeycomb.io](http://www.honeycomb.io)