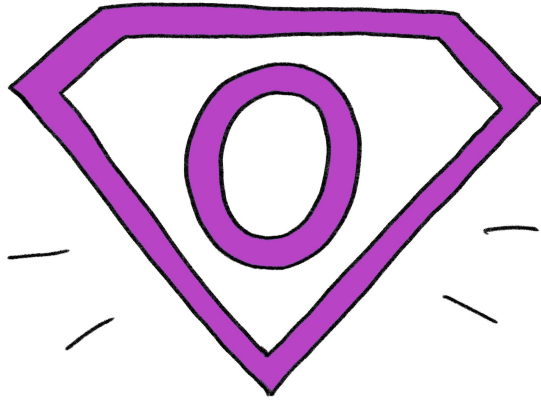


Managing to your SLO (amidst chaos!)

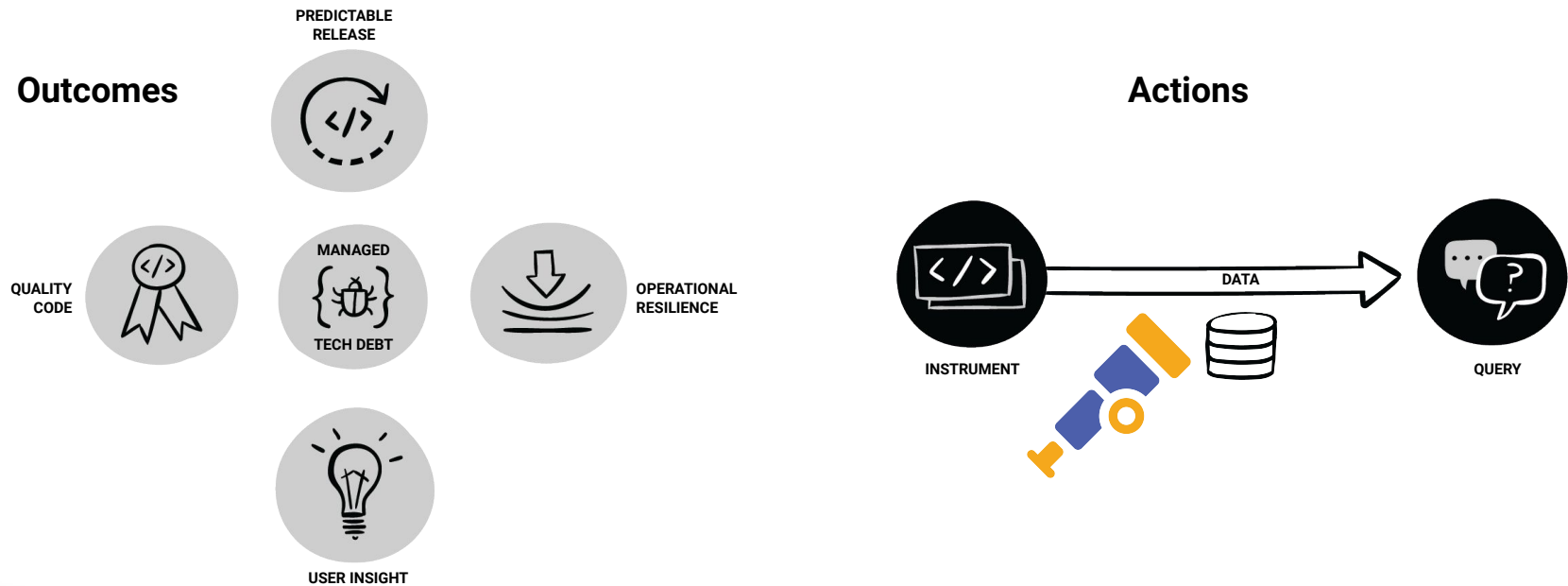
Liz Fong-Jones
Field CTO, Honeycomb.io
YOW! 2022



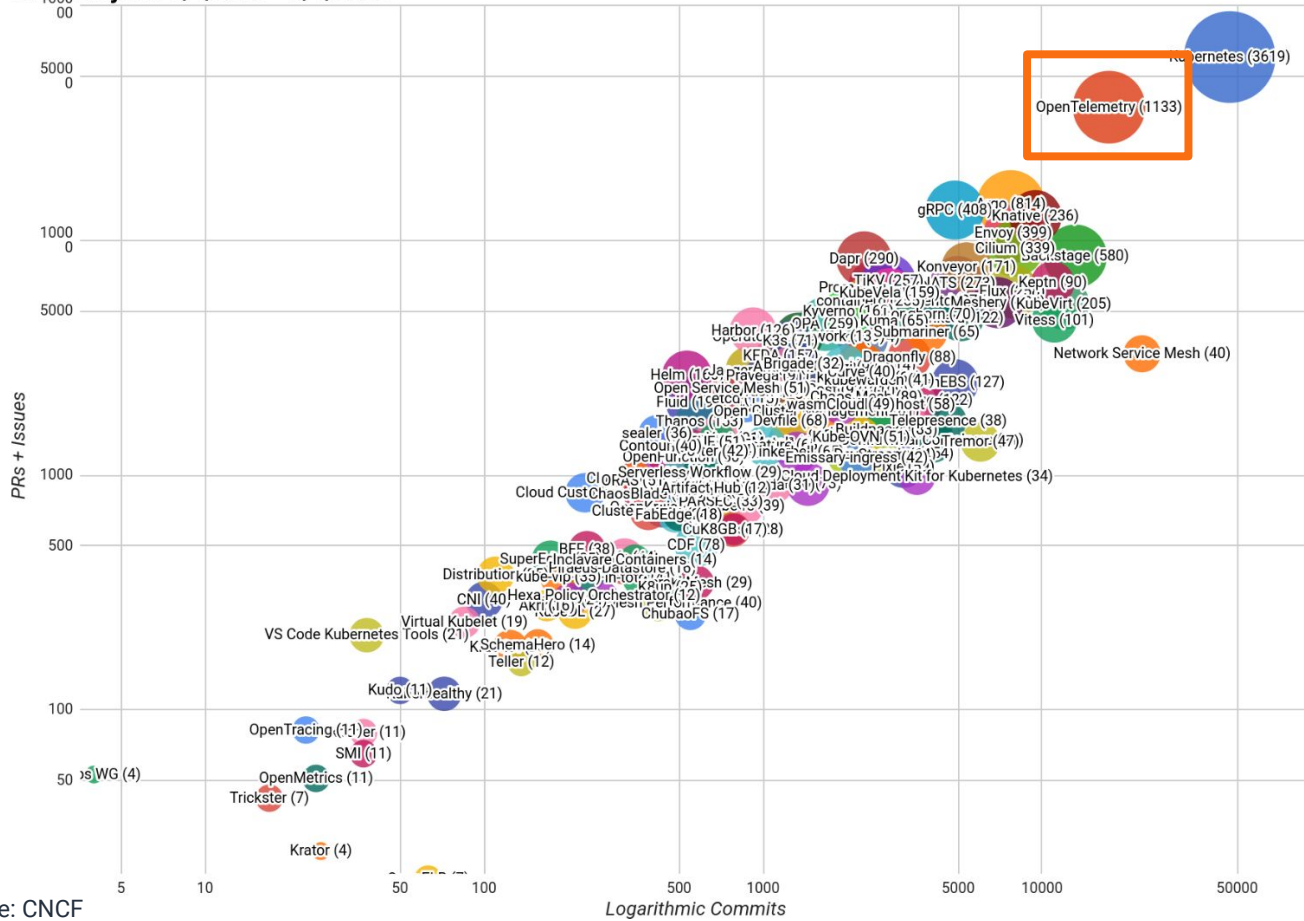


Observability is evolving quickly.

And the problem space is complex.



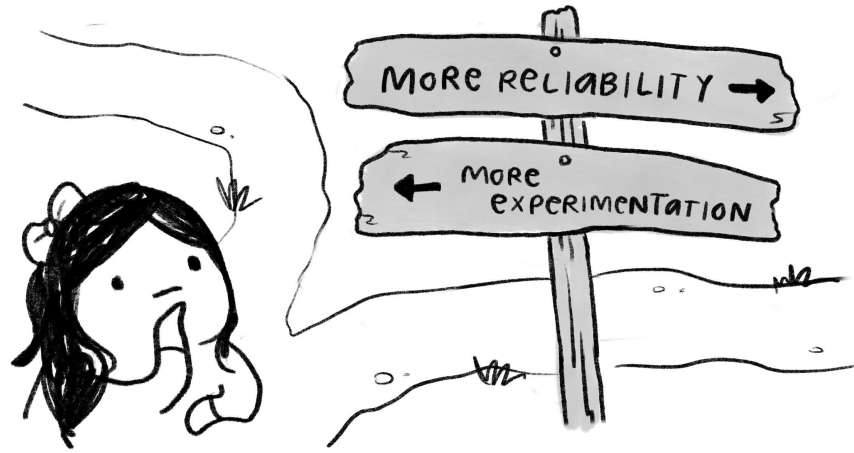
CNCF Projects 8/1/2021 - 8/1/2022



- 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
- 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
- 111 more

Source: CNCF

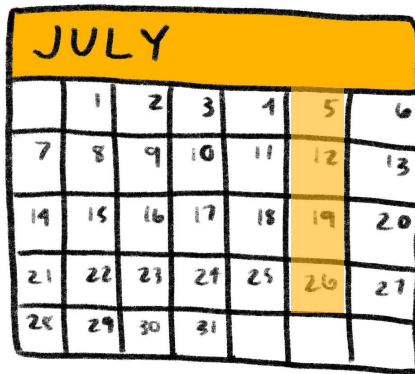




Practitioners need **velocity**, **reliability**, & **scalability**.



A small but **growing** team builds Honeycomb.




We **deploy** with confidence.

| | | | | |
|---|--|--|-----|---------------------------------------|
| VISUALIZE COUNT_DISTINCT(global.build_id) | WHERE trace.parent_id does-not-exist | GROUP BY None; don't segment | ... | Run Query Run 2 minutes ago |
| + ORDER BY | + LIMIT | + HAVING | | |

Results BubbleUp Metrics Traces Raw Data

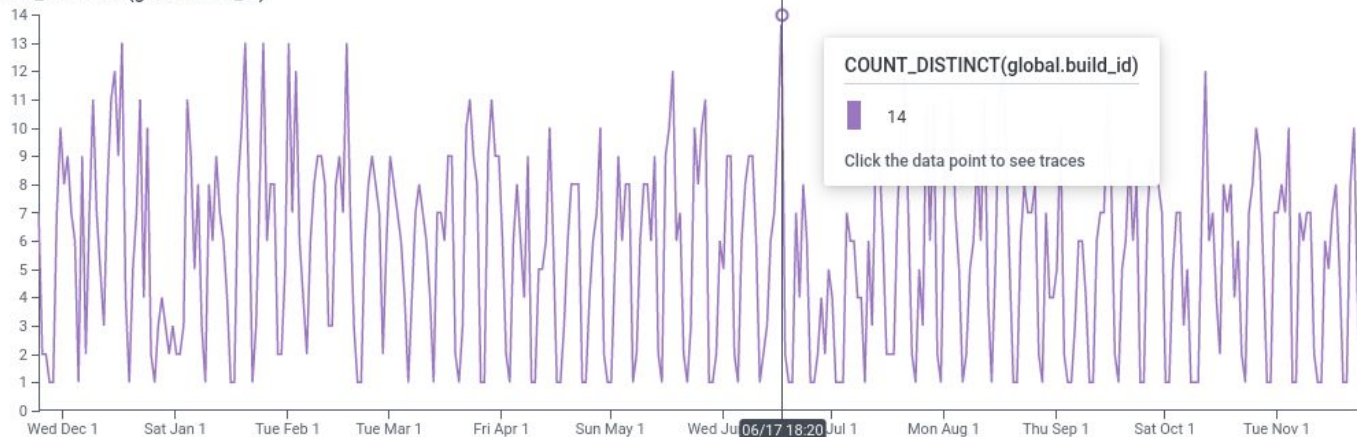
☐ Compare to

Previous time range

 Graph Settings

Nov 24 2021 19:13:24 – Nov 24 2022 19:13:24 UTC+11:00 (Granularity: 1 day)

COUNT_DISTINCT(global.build_id)



COUNT_DISTINCT(global.build_id) 

 1,564

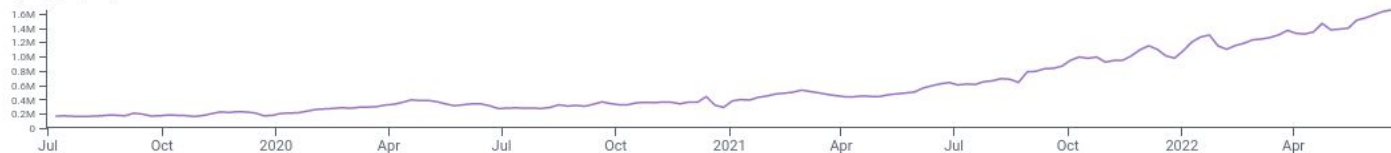
elapsed query time: 4.951435582s rows examined: 13,381,219,529 nodes reporting: 100%



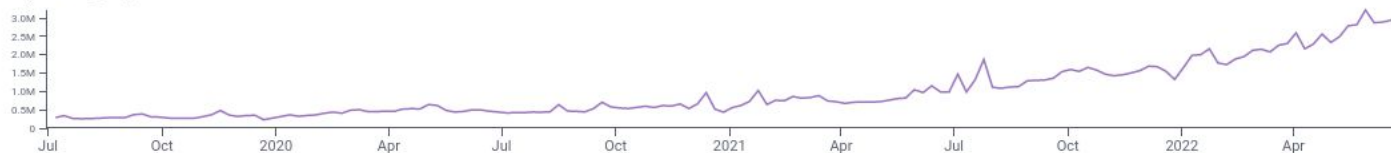
10x growth in three years

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

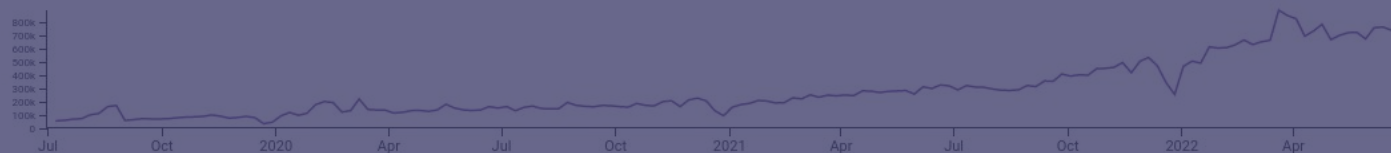
AVG(avg_ingest)



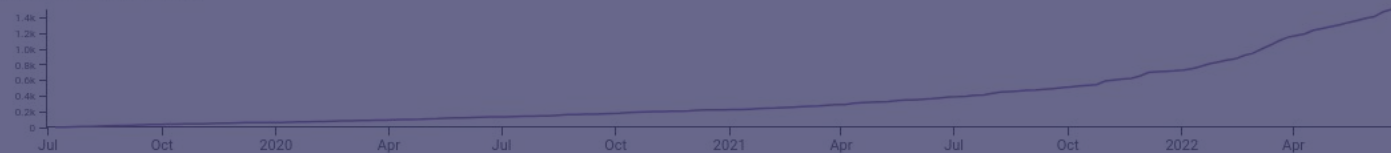
MAX(peak_ingest)



SUM(human_queries)



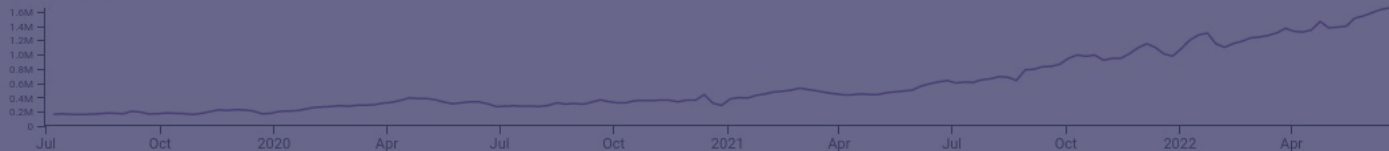
SUM(triggers_per_minute)



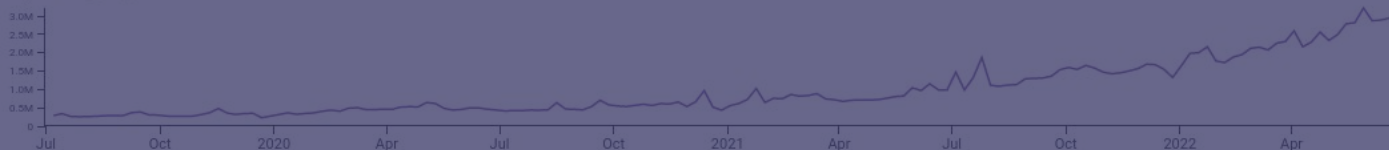
10x growth in three years

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

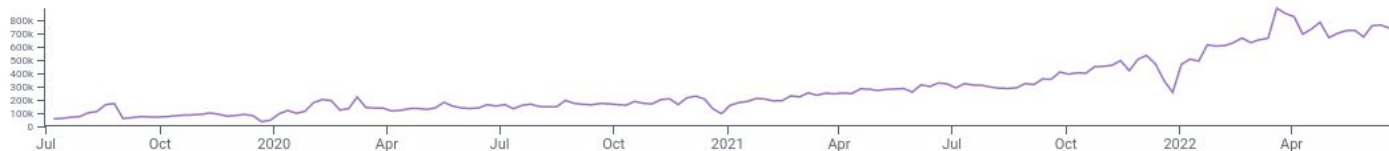
AVG(avg_ingest)



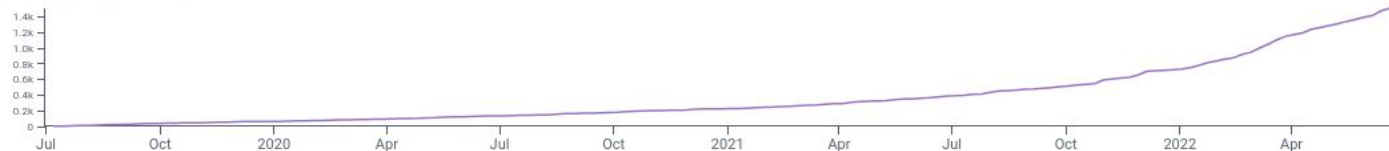
MAX(peak_ingest)



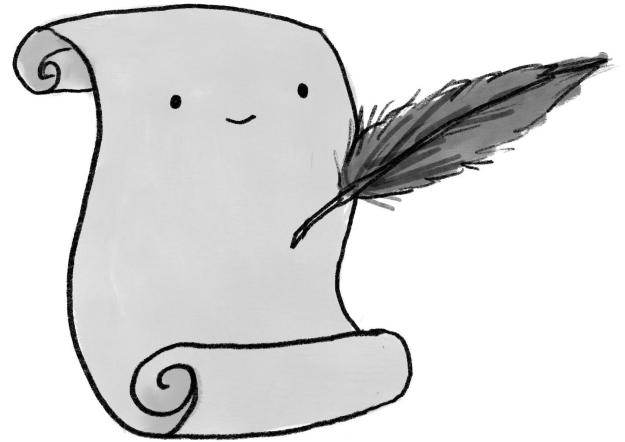
SUM(human_queries)



SUM(triggers_per_minute)



Our **confidence** recipe:





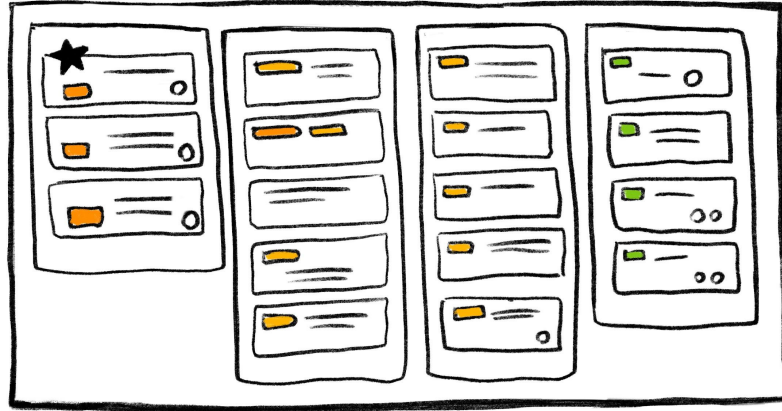
Quantify reliability.



Identify potential areas of risk.



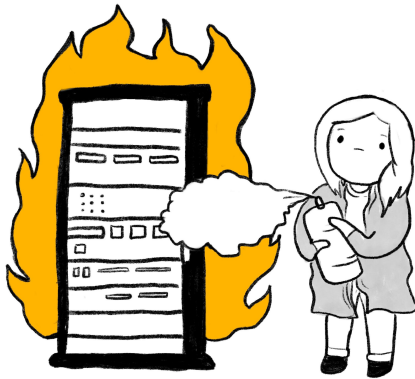
Design **experiments** to probe risk.



Prioritize **addressing risks.**

Measuring **reliability**:



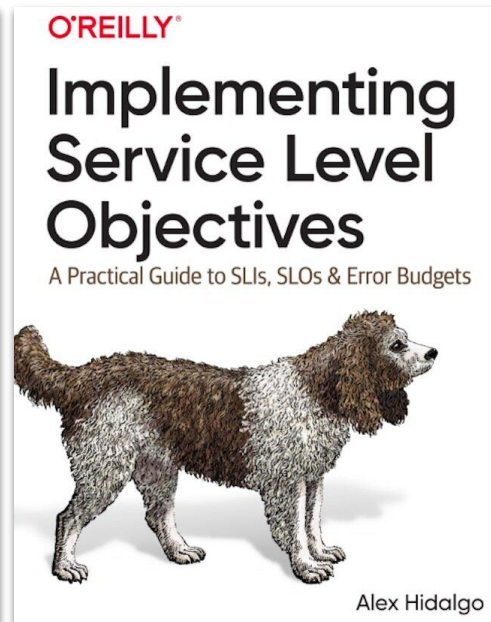
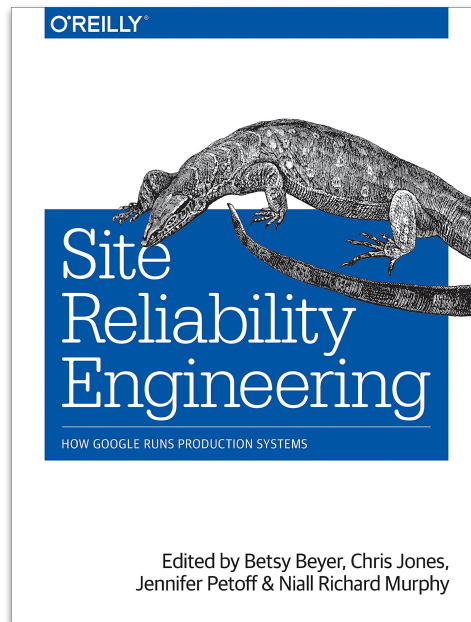


How broken is “**too broken**”?

Service Level Objectives (SLOs)

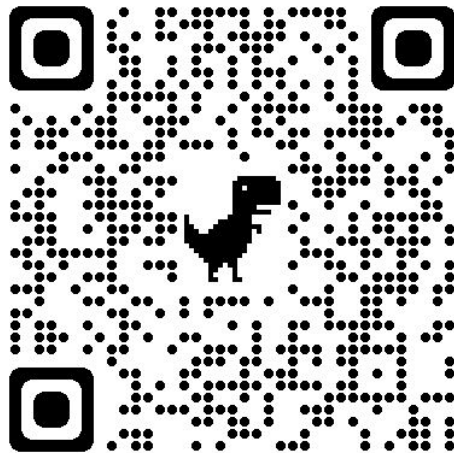
Define and measure success!

Popularized by Google, widely adopted now!



Service-Level Objectives

- Described in Chapter 4 of the Google SRE Book
- Based on what users actually want, not what's in a contract
- Don't plan for 100% uptime.
- SLOs are an opportunity to negotiate!
- See also: #YOW19 talk by yours truly



Service-Level Objectives

- Negotiation:
 - What's important to your customers?
 - What can everyone agree on? If it's not important, maybe it's not worth measuring in an SLO?
- Paging only on SLOs helps reduce pager fatigue.



DATA GENERATION

STREAMING INGEST

STORAGE + PROCESSING

VISUAL ANALYSIS LOOP

Logging Data

Amazon Relational
Database Service
(Amazon RDS)AWS Elastic
Load BalancingAWS Elastic
BeanstalkAmazon Elastic
Kubernetes Service
(Amazon EKS)

Realtime Ingest

Auth and validation

Metrics Data

Host Metrics

App Metrics

Amazon
CloudWatch

Prometheus

OpenTelemetry

Honeycomb API

Unpacked
Into
Columns

Trace Span Data

JavaScript

Go

Jaeger

Java

Ruby

Python

.NET

Front-End

OpenTelemetry

Refinery

User-controlled
dynamic tail sampling

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 fileSegment with
time rangeAmazon Simple Storage
Service (Amazon S3)Context
Context
Context
Context
Context
Context
ContextColumn
Column
Column
Column
Column
Column
ColumnSegment
Segment
Segment
Segment
Segment
Segment
Segment

AWS Lambda

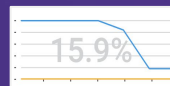
ui.honeycomb.io

ui.honeycomb.io

ui.honeycomb.io

Events contain many
fields and distinct
valuesContext Context Context Context
Context Context Context Context
Context Context Context Context

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



Graph Rendering

Click through
visuals based on
granular dataQuery
Builder UIIntuitive GUI with
multiple group-byQuery Result
HistoryShared team
intelligence

SLOs are user flows

Honeycomb's SLOs

- home page loads quickly
- user-run queries are fast
- customer data gets ingested fast



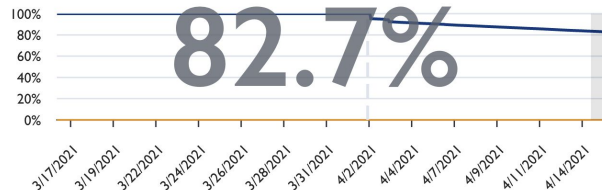
Latency per-event

Shepherd ingestion latency should be below 5ms per event within a batch. We ignore values from user-triggered issues, deprecated endpoints we won't support as extensively as the main ones, and also ignore values coming from collectd which historically was a misbehaving client whose API we don't control.

99.99% of eligible events from the **shepherd** column **sli** will succeed over a period of 30 days.

Budget Burndown


How much of the error budget remains after the last 30 days. Starts at 100% and burns down.



Service-Level Objectives

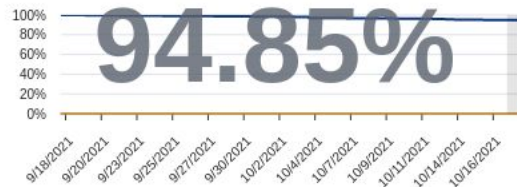
- Example Service-Level Indicators:
 - 99.99% of events succeed without error over a period of 30 days.
 - 99.99% of events are processed in 5ms over a period of 30 days.
- Translates to about 4.5 minutes of violation in a month.



99.99% of eligible events from the  shepherd column
 error_sli will succeed over a period of 30 days.

Budget Burndown

How much of the error budget remains after the last 30 days. Starts at 100% and burns down.



Reset

Historical SLO Compliance

For each day of the past 30, how often this SLI has succeeded over the preceding 30 days.



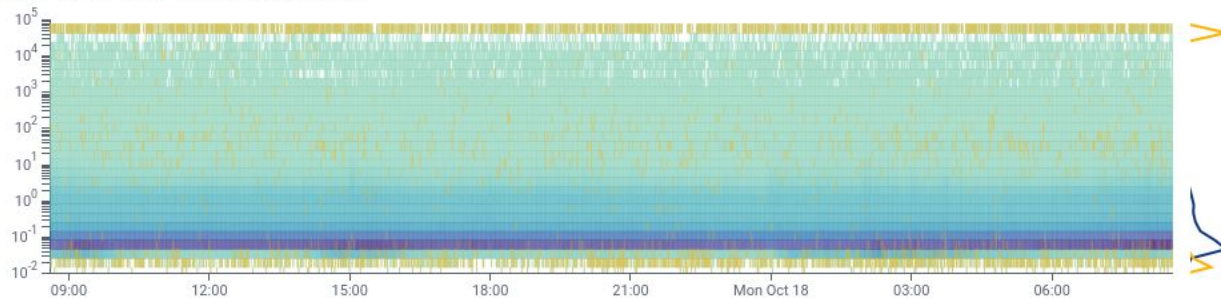
Distribution of Events failing SLI by

duration_ms

Make Default





Last 24 hours

Oct 17 2021, 8:36:05 AM – Oct 18 2021, 8:36:05 AM

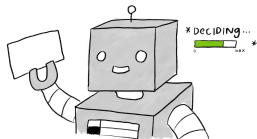
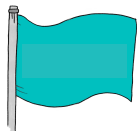


How to stay within SLO

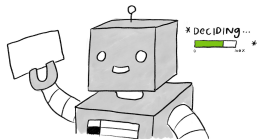
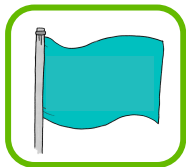
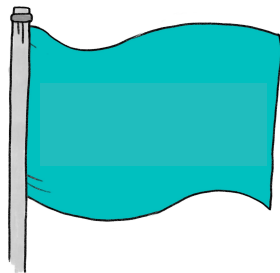
Minimise time to mitigate & recover.

| Software delivery performance metric | Elite | High | Medium | Low |
|--|--------------------------------------|--|--|--------------------------------|
|  Deployment frequency For the primary application or service you work on, how often does your organization deploy code to production or release it to end users? | On-demand (multiple deploys per day) | Between once per week and once per month | Between once per month and once every 6 months | Fewer than once per six months |
|  Lead time for changes For the primary application or service you work on, what is your lead time for changes (i.e., how long does it take to go from code committed to code successfully running in production)? | Less than one hour | Between one day and one week | Between one month and six months | More than six months |
|  Time to restore service For the primary application or service you work on, how long does it generally take to restore service when a service incident or a defect that impacts users occurs (e.g., unplanned outage or service impairment)? | Less than one hour | Less than one day | Between one day and one week | More than six months |
|  Change failure rate For the primary application or service you work on, what percentage of changes to production or released to users result in degraded service (e.g., lead to service impairment or service outage) and subsequently require remediation (e.g., require a hotfix, rollback, fix forward, patch)? | 0%-15% | 16%-30% | 16%-30% | 16%-30% |

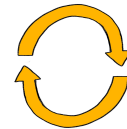
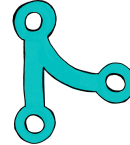
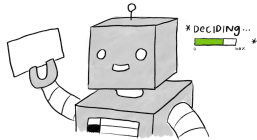
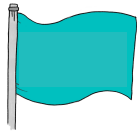
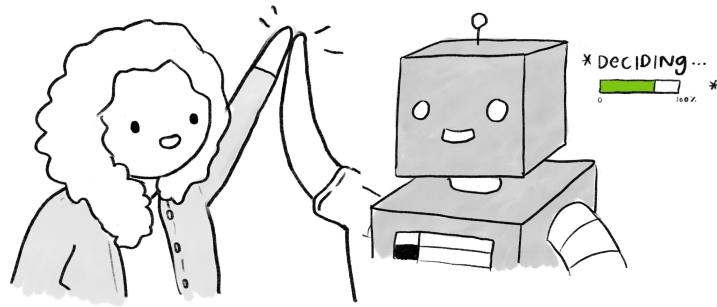
What's our recipe?



Design for feature flag deployment.



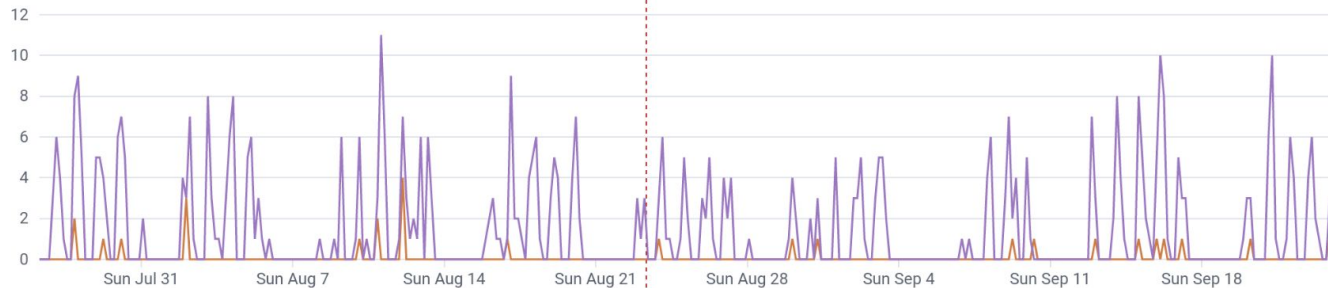
Automated integration.



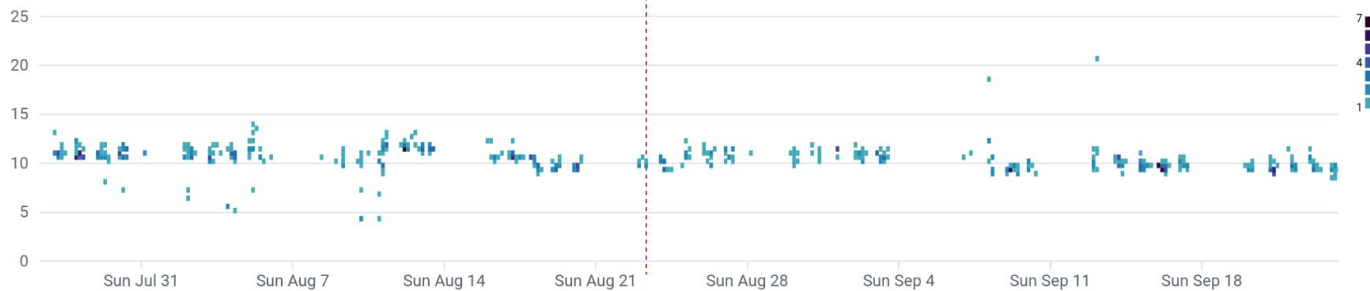


Jul 26 2022 07:07:58 – Sep 24 2022 07:07:58 UTC+10:00 (Granularity: 4 hr)

COUNT



HEATMAP(duration_min)



status ▾

COUNT ▾

HEATMAP(duration_min) ▾

success

521



failed

29





CI/C >

VISUALIZE

HEATMAP(duration_min)

WHERE

trace.parent_id does-not-exist
 branch = main
 status = success
 repo = git@github.com:honeycombio/hound.git

GROUP BY

None; don't segment

Run Query

Run a few seconds ago



+ ORDER BY

+ LIMIT

+ HAVING

Results BubbleUp Metrics Traces Raw Data

☐ Compare to

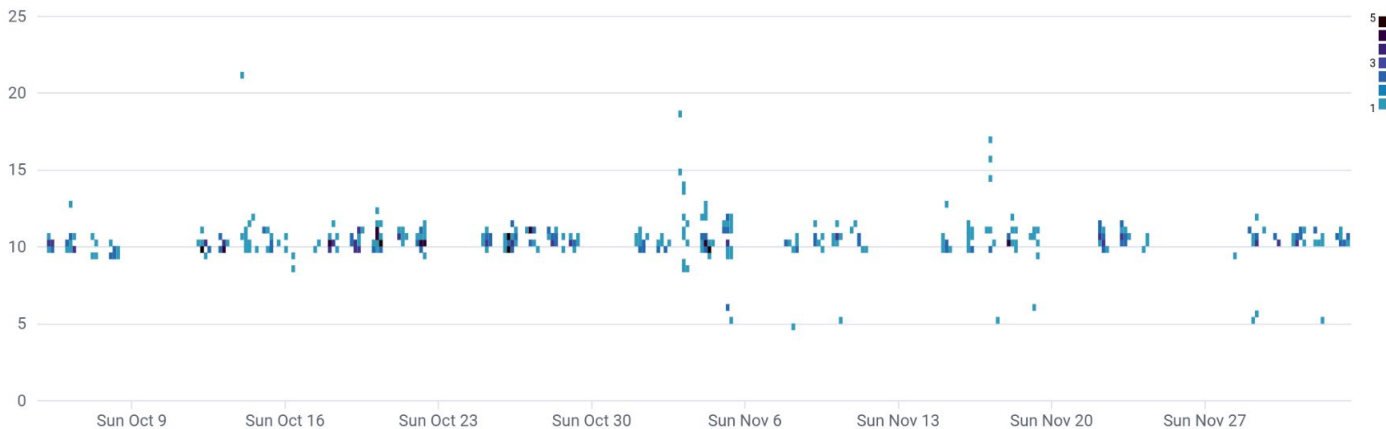
Previous time range



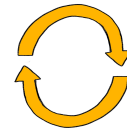
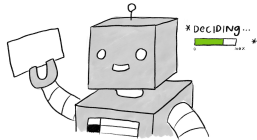
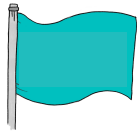
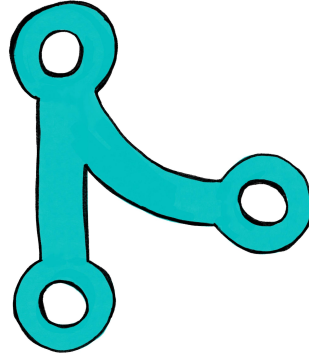
Graph Settings

Oct 4 2022 16:31:20 – Dec 3 2022 16:31:20 UTC+11:00 (Granularity: 4 hr)

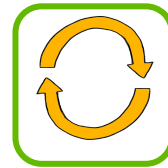
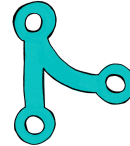
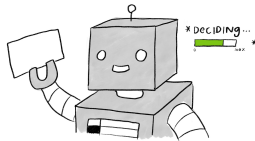
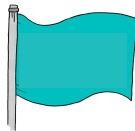
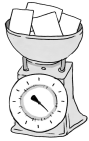
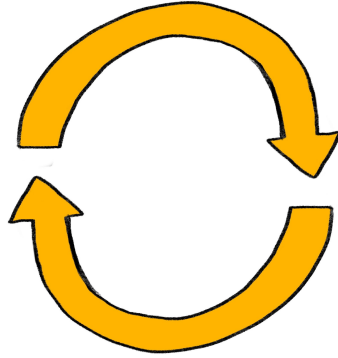
HEATMAP(duration_min)



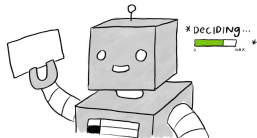
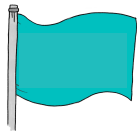
Green button merge.



Auto-updates, rollbacks, & pins.

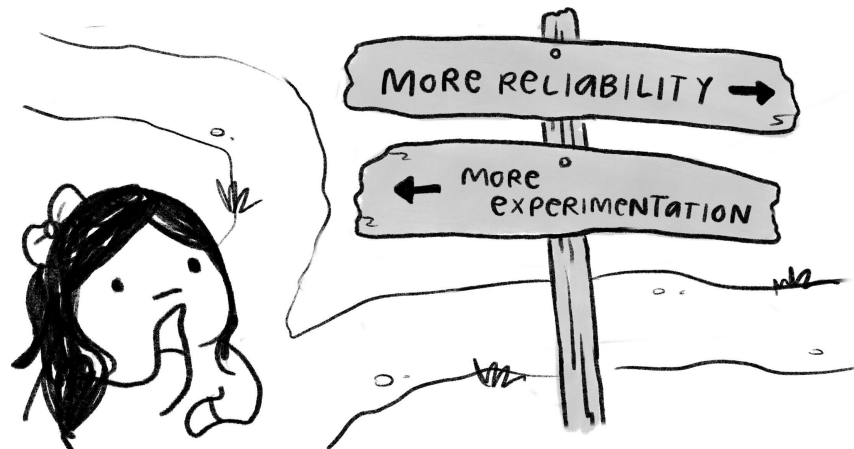


Observe behavior in prod.

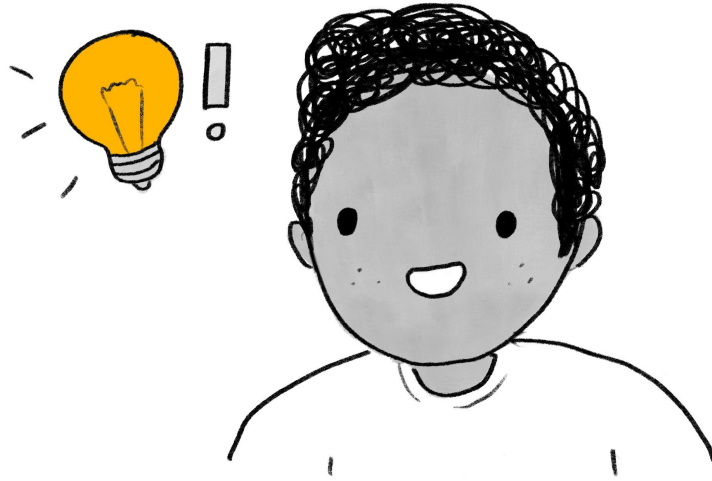


Validating our expectations

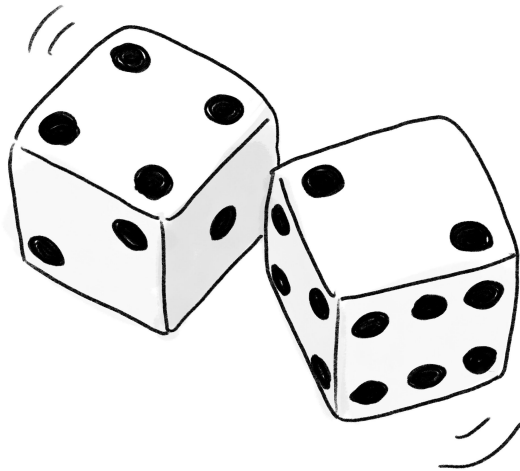




Experiment using **error budgets**.
(only when there's error budget left)




Keep learning objectives in mind.



Always ensure **safety.**

Targeting ☒

Prerequisites 

Target individual users 

Target users who match these rules 

IF

team_id



is one of



SERVE

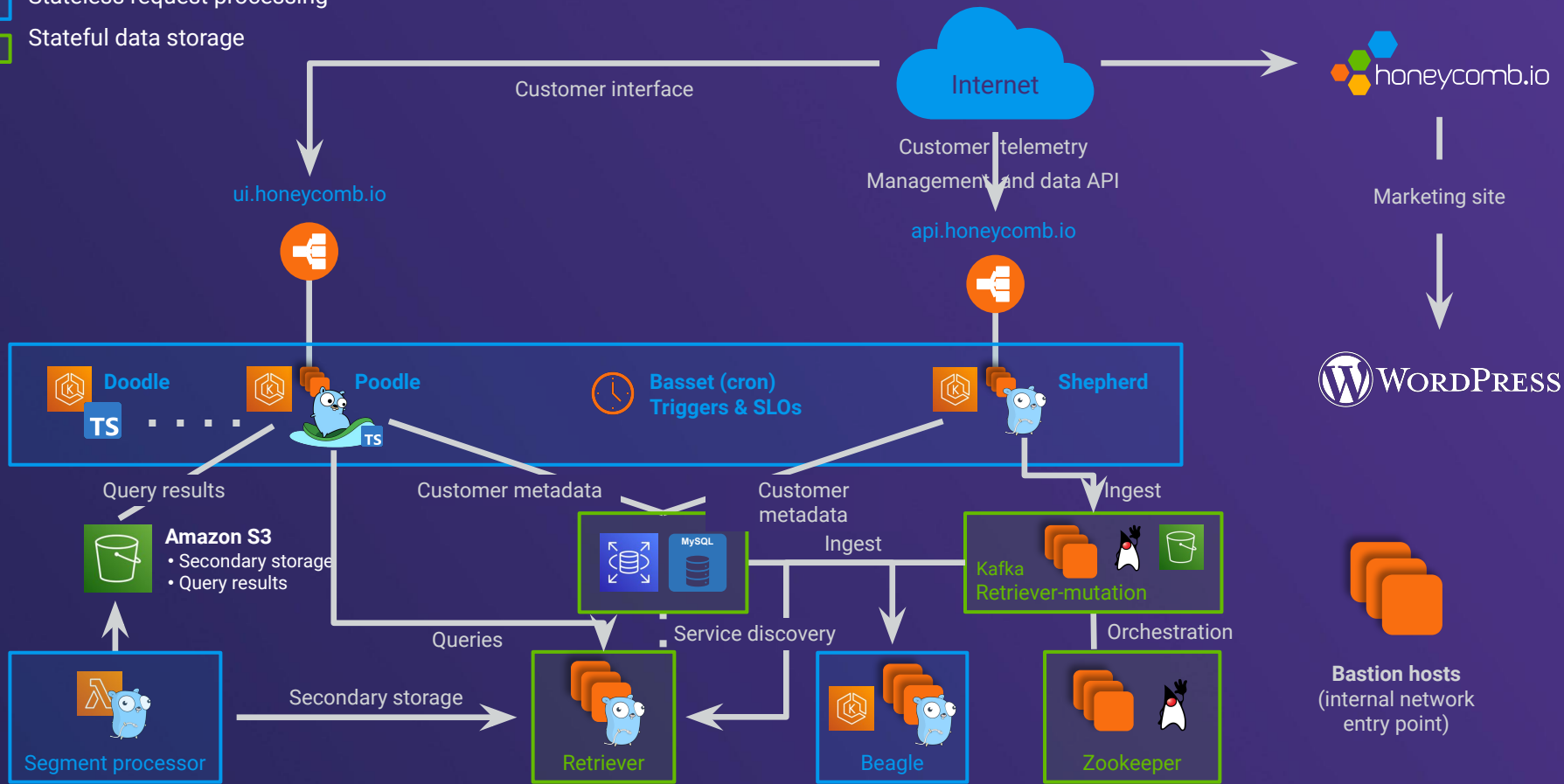
☒ true

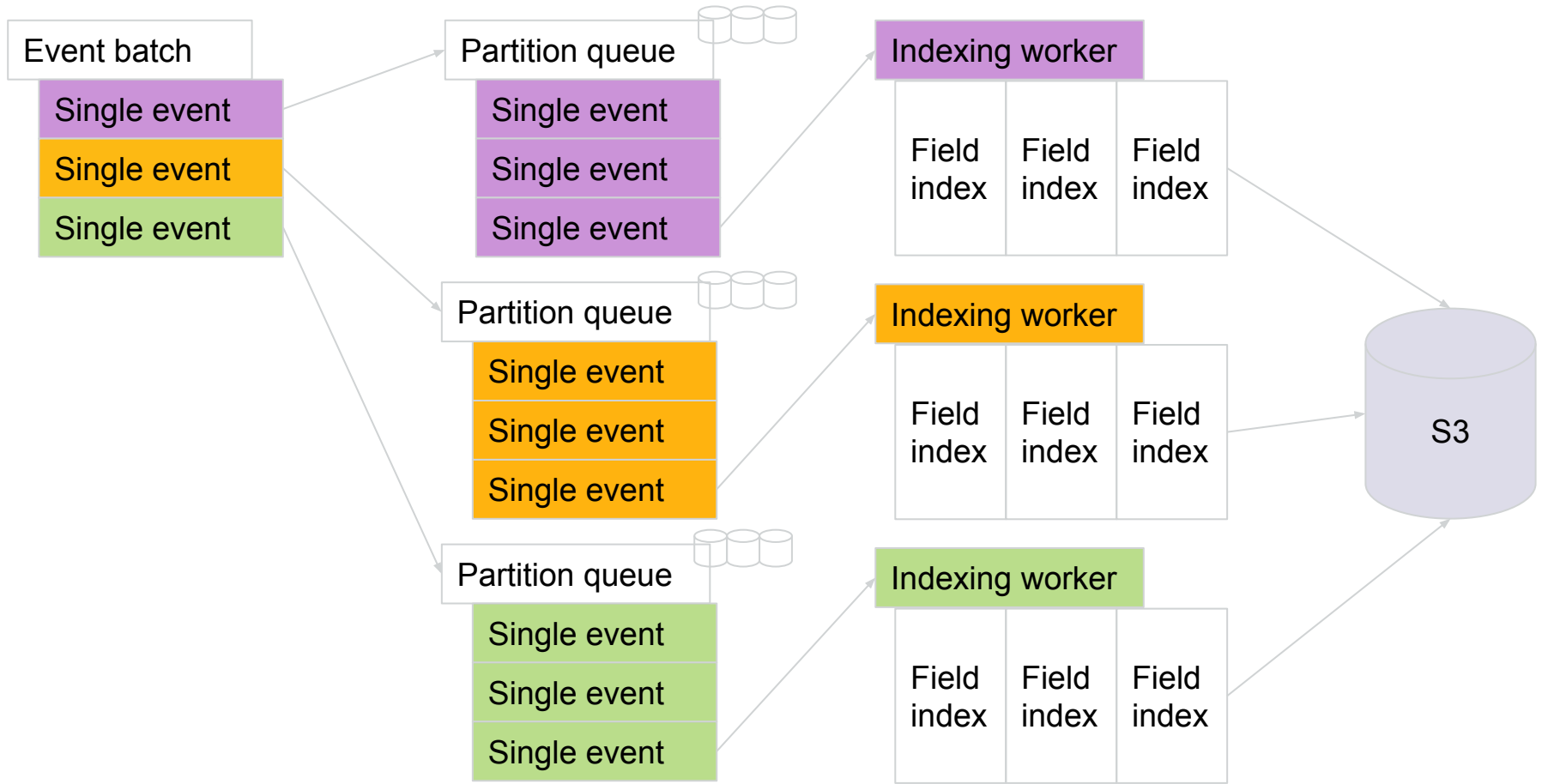


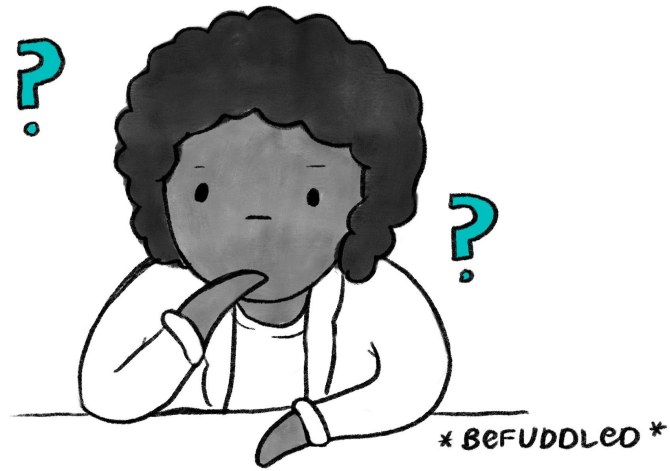


Data **persistence** is tricky.

- Stateless request processing
- Stateful data storage



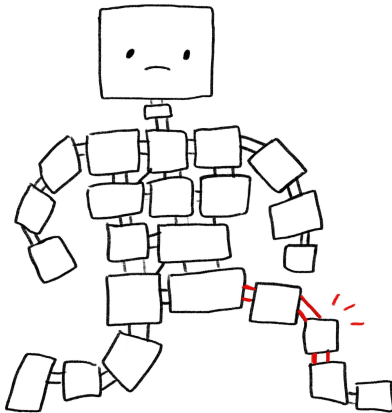




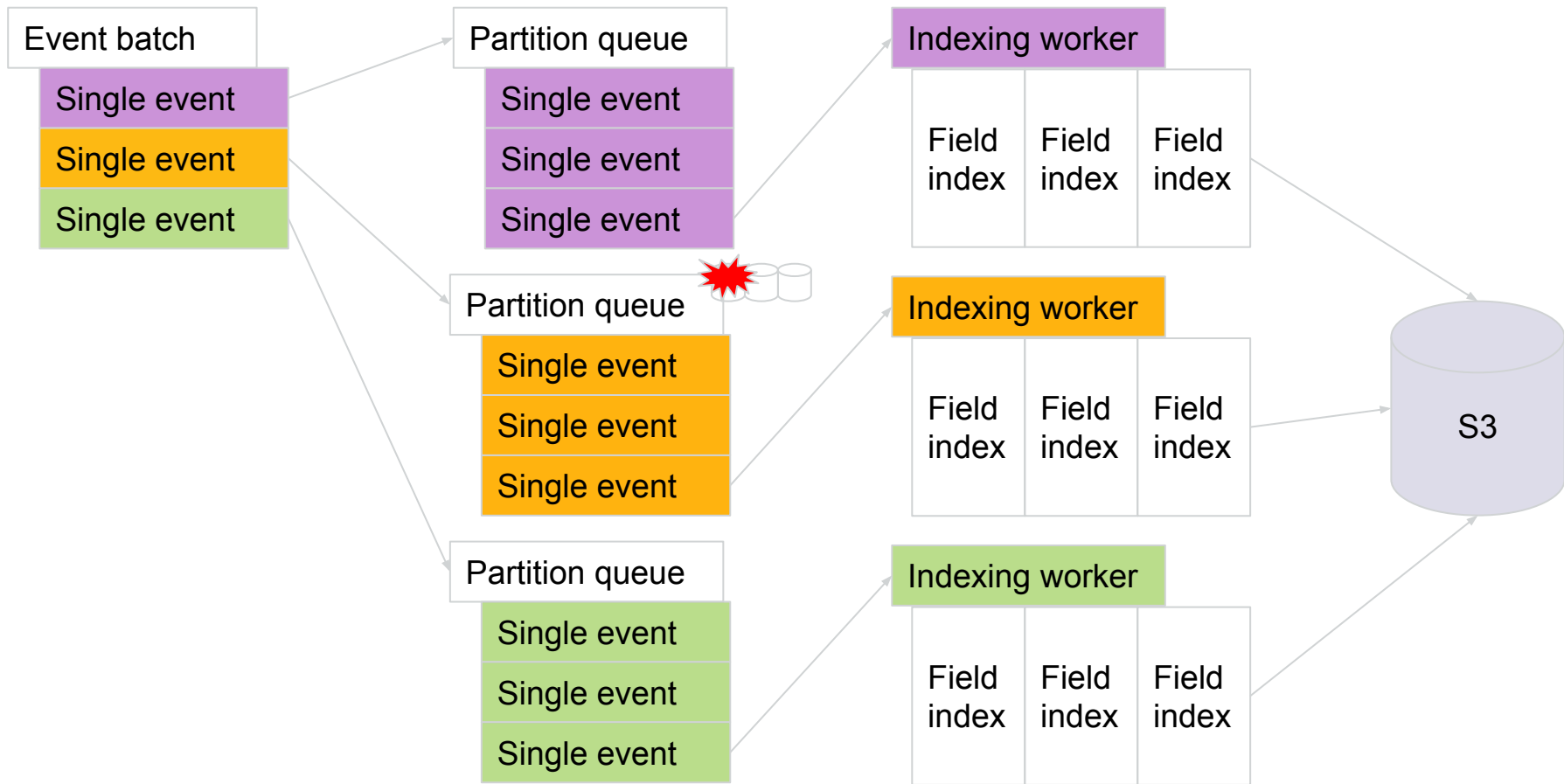
Infrequent changes.

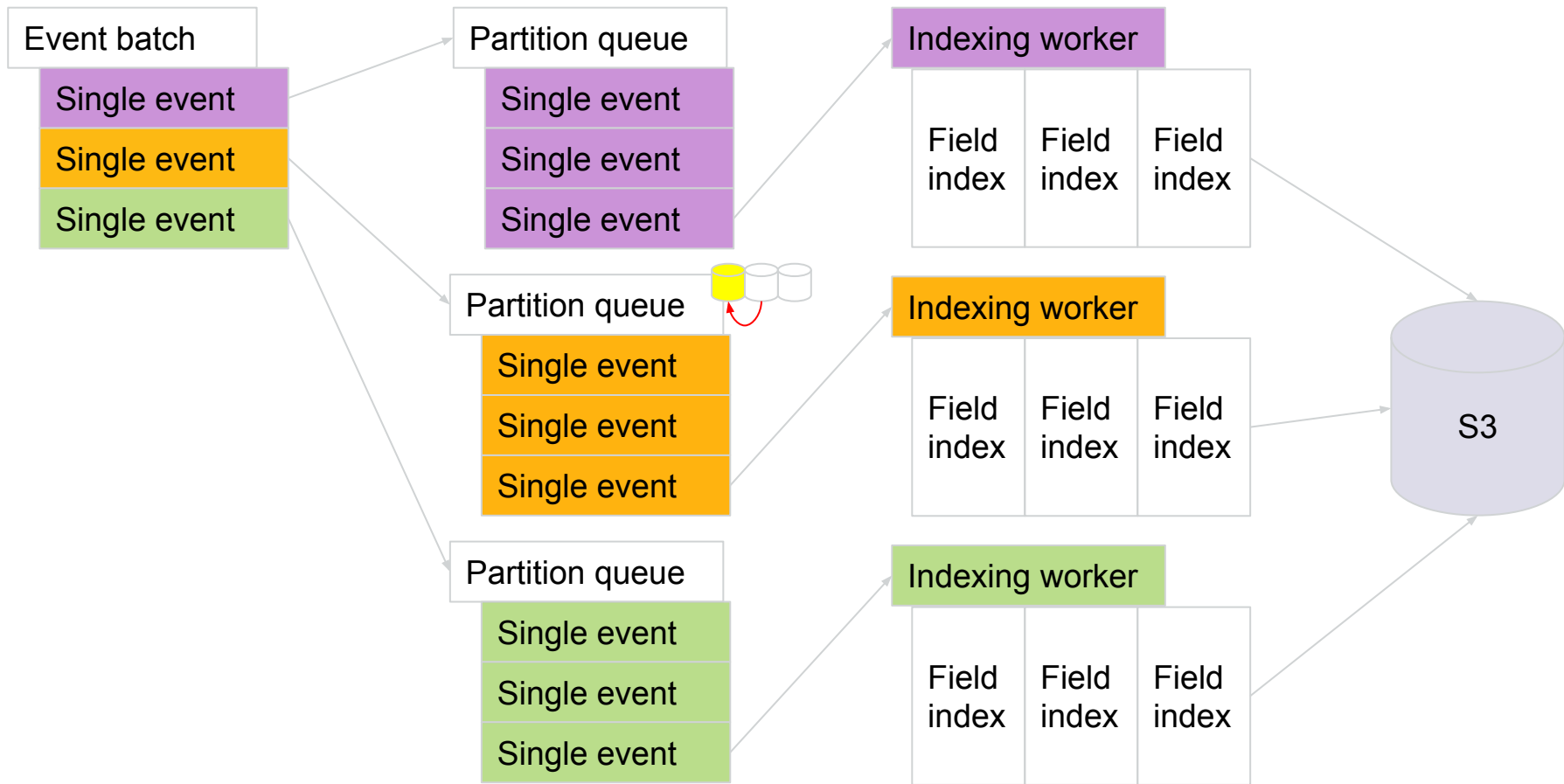


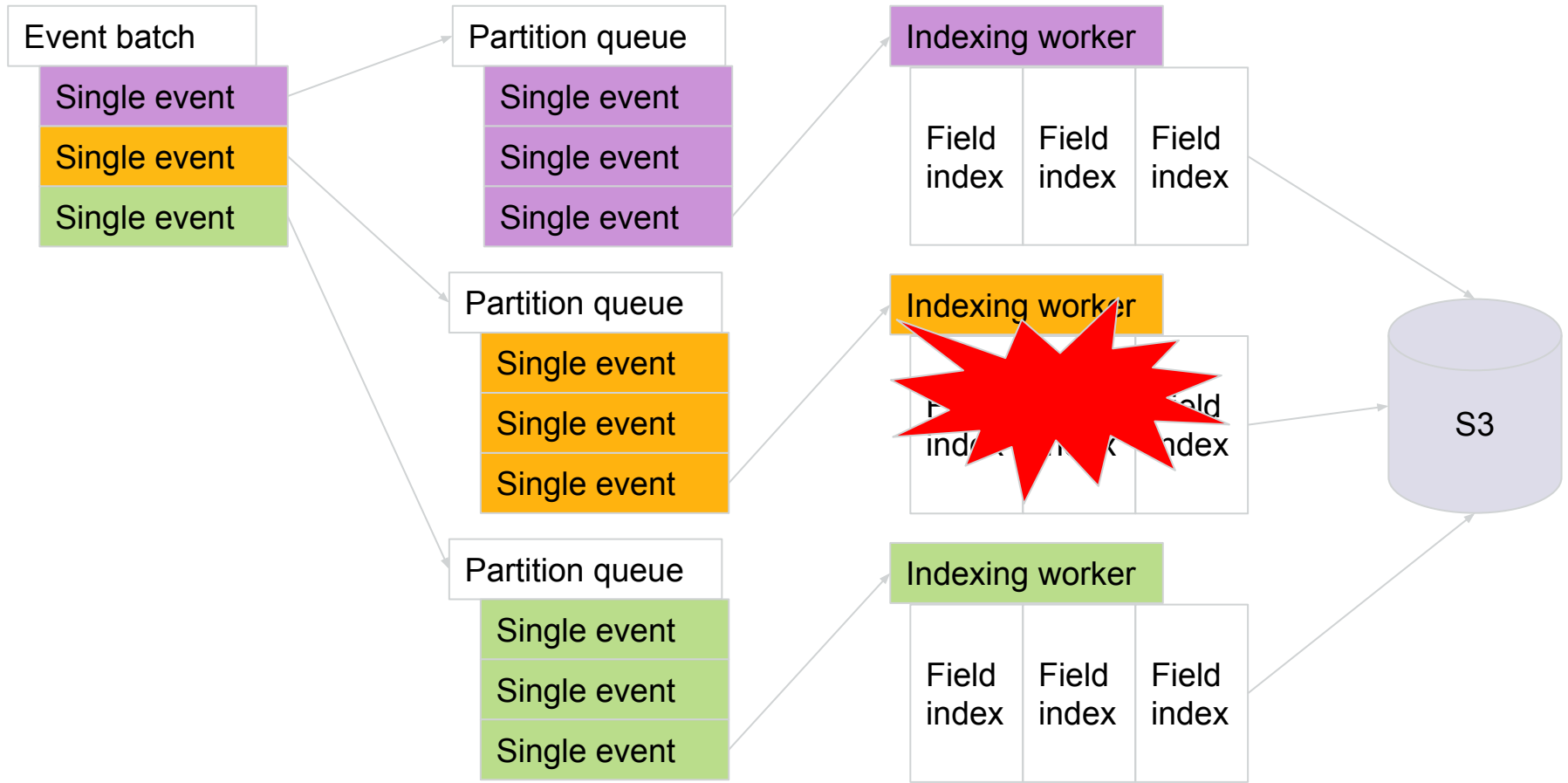
Data integrity and consistency.

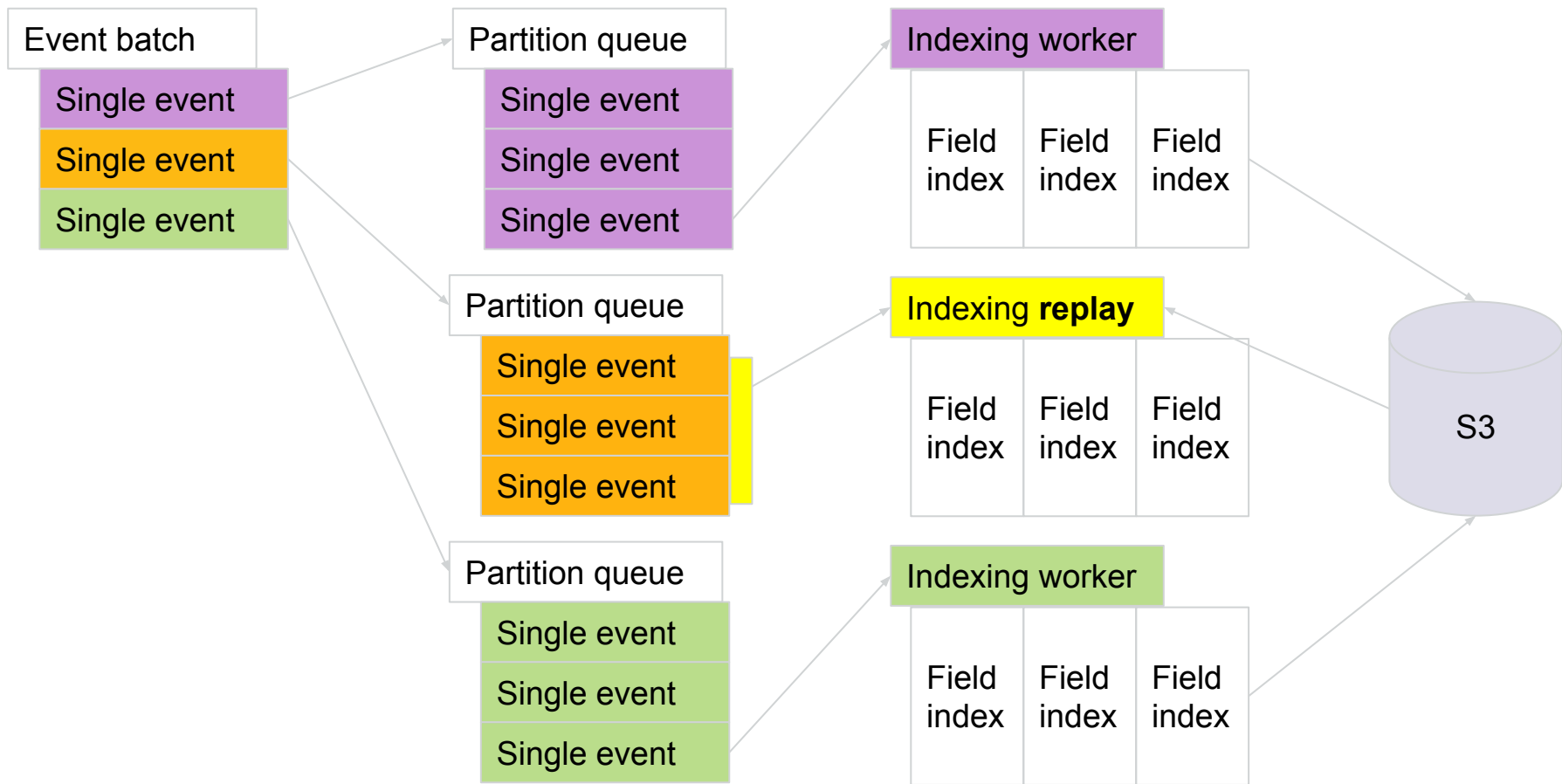


Delicate failover dances

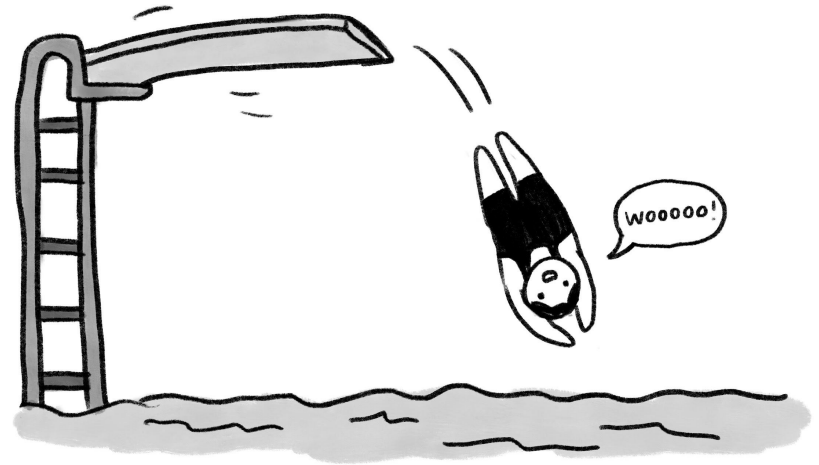




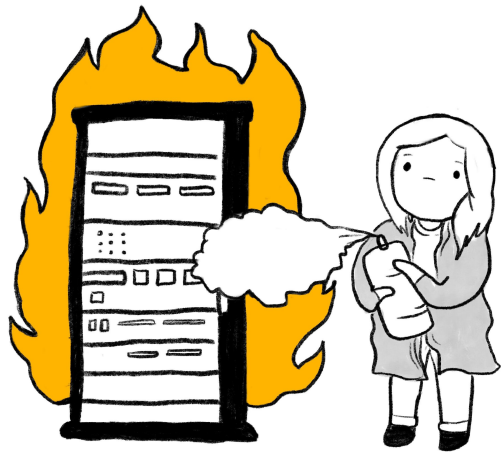




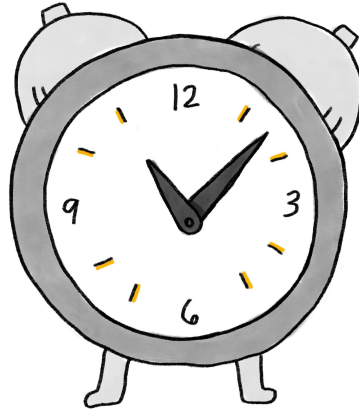
Experimenting in prod



@lizthegrey at #YOW22



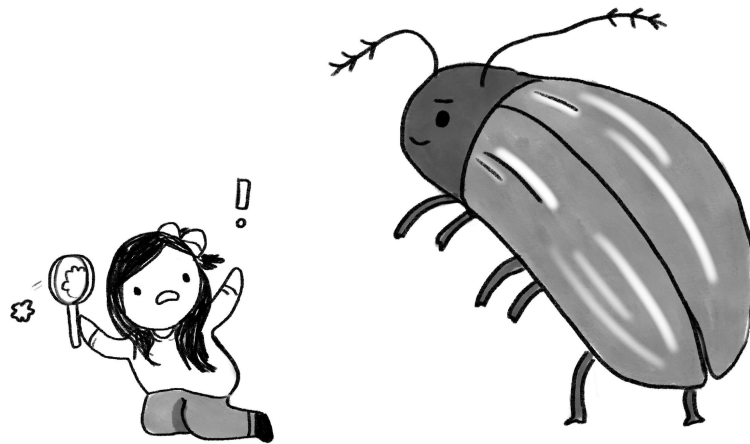
Restart **one server & service** at a time.



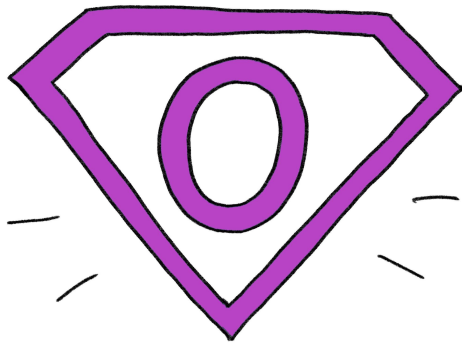
At 3pm, not at 3am.



"Bugs are shallow with more eyes."



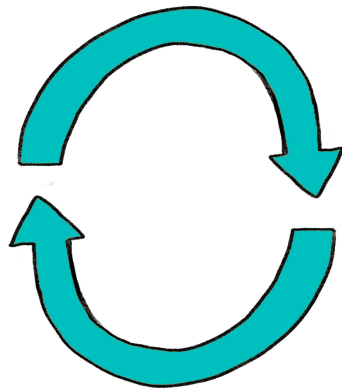
Monitor for changes using SLIs.



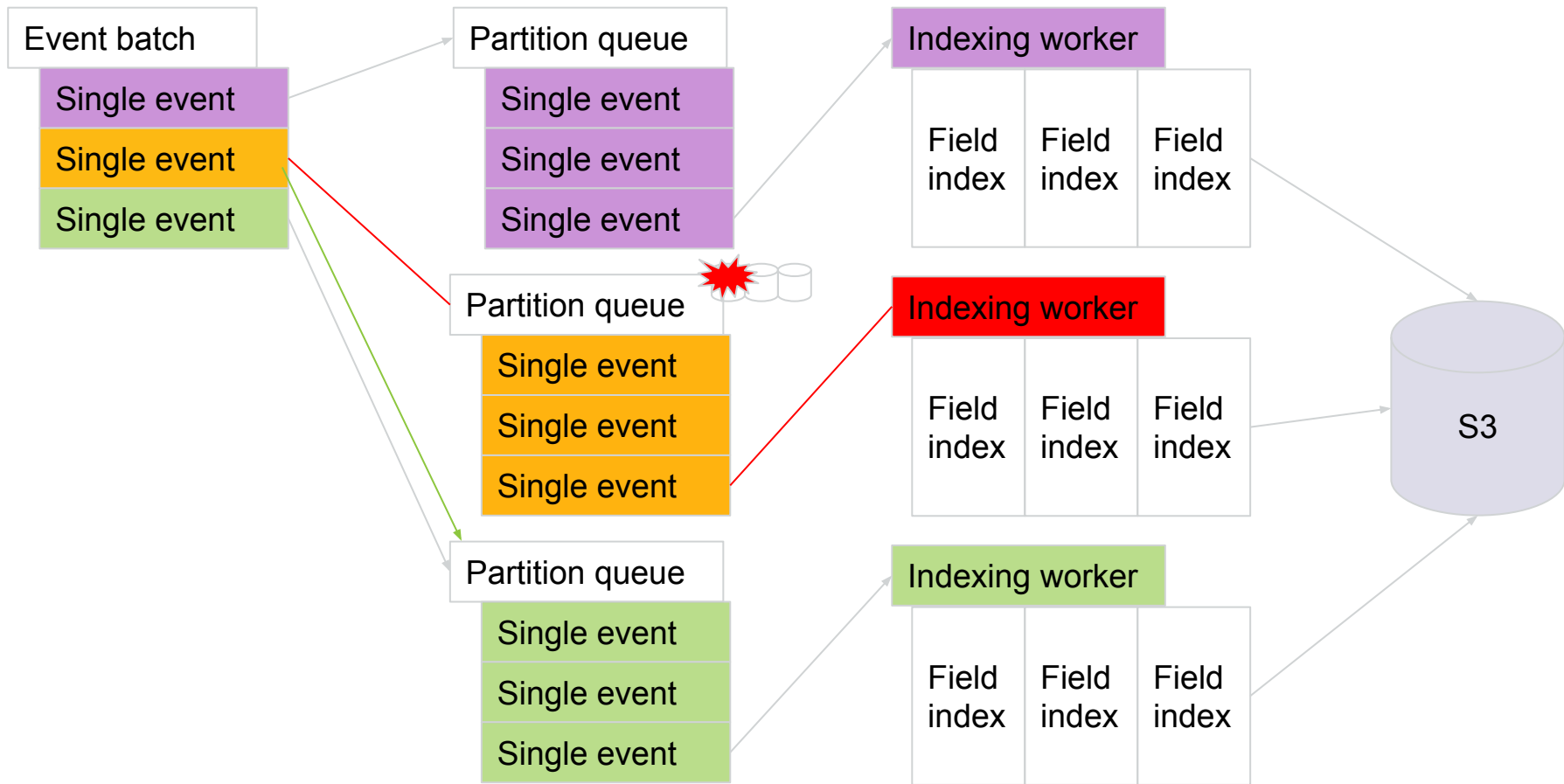
Debug with observability.

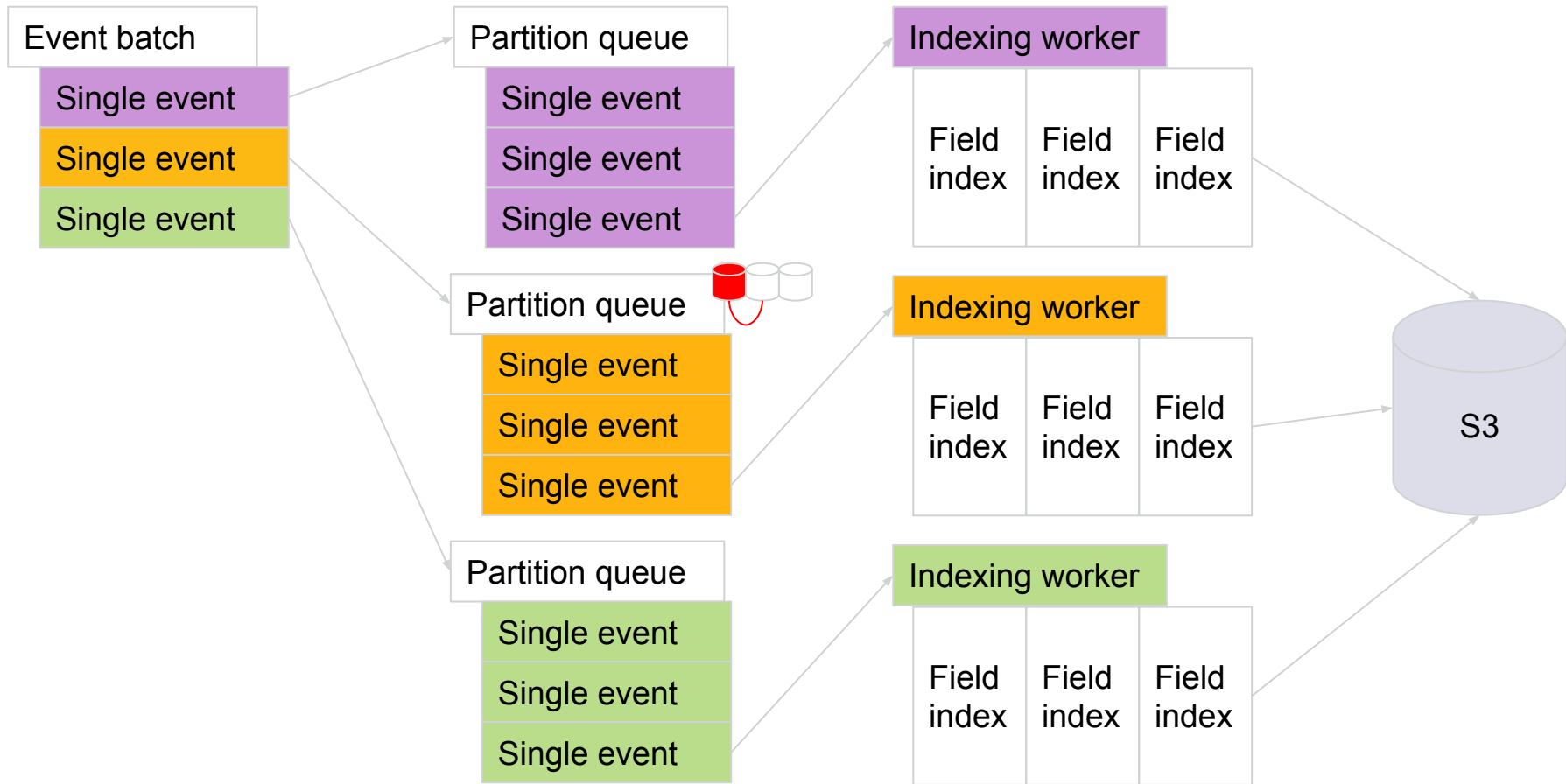


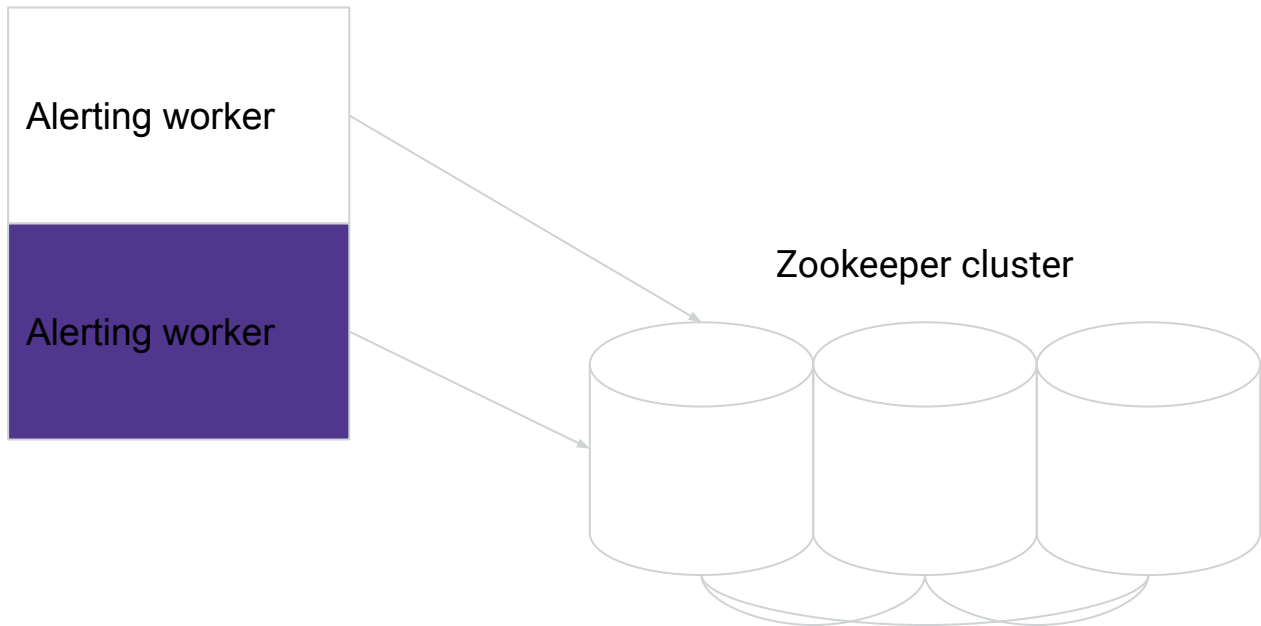
Test the **telemetry** too!

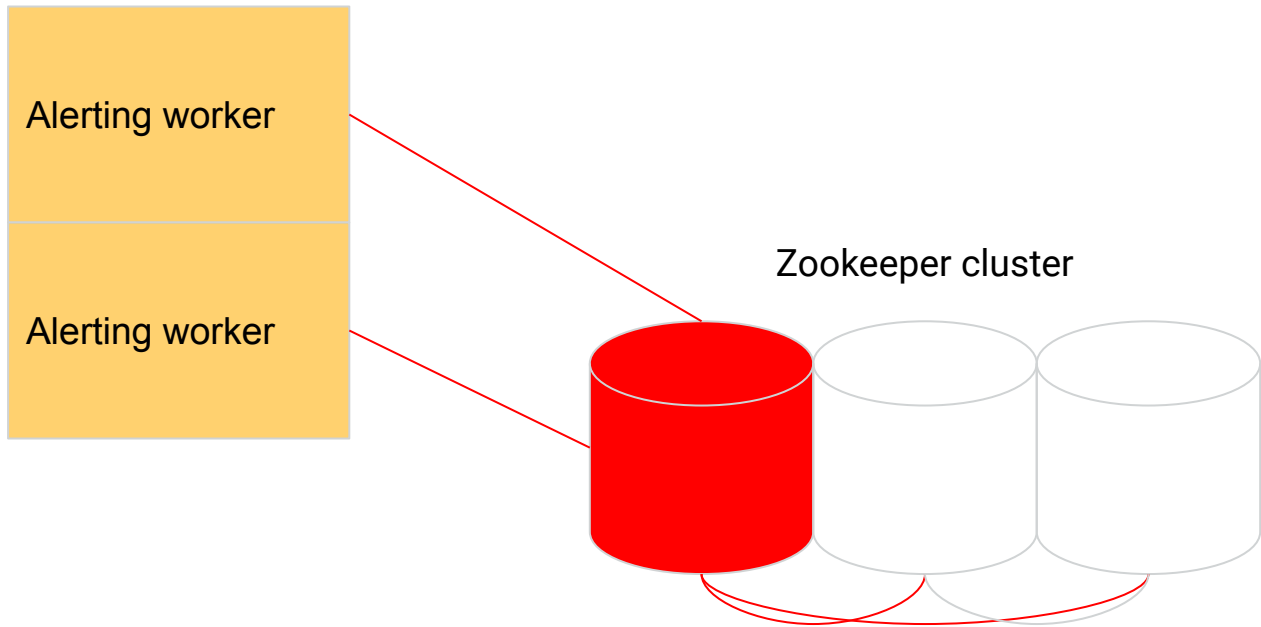


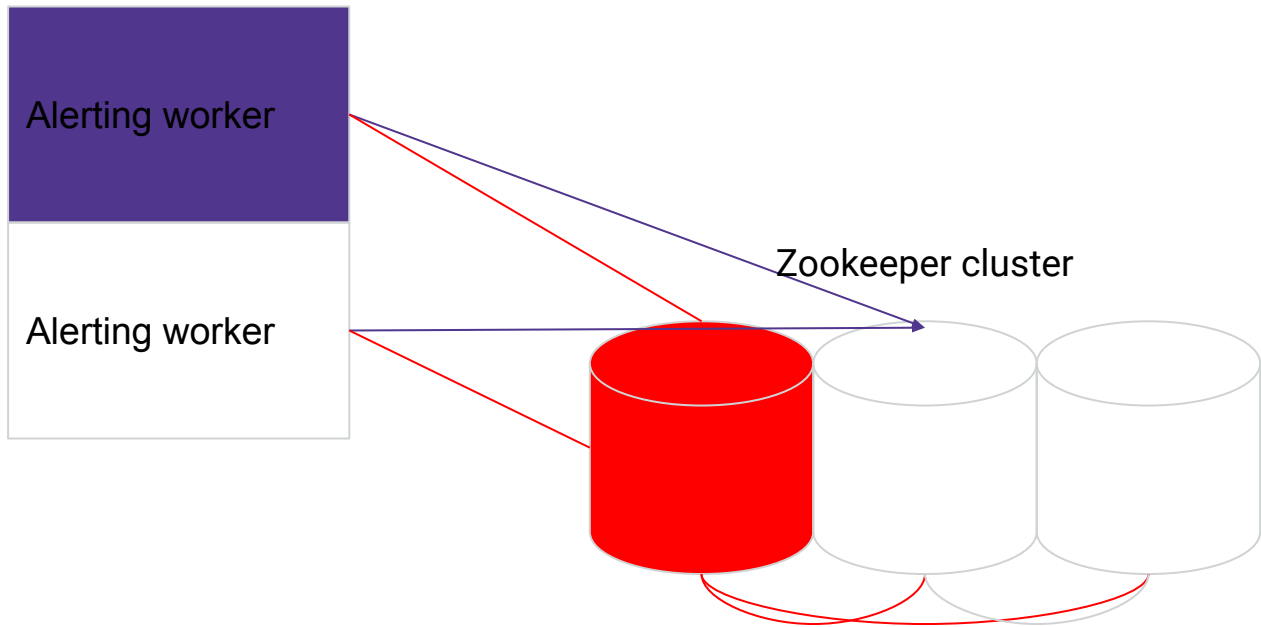
Verify fixes by repeating.

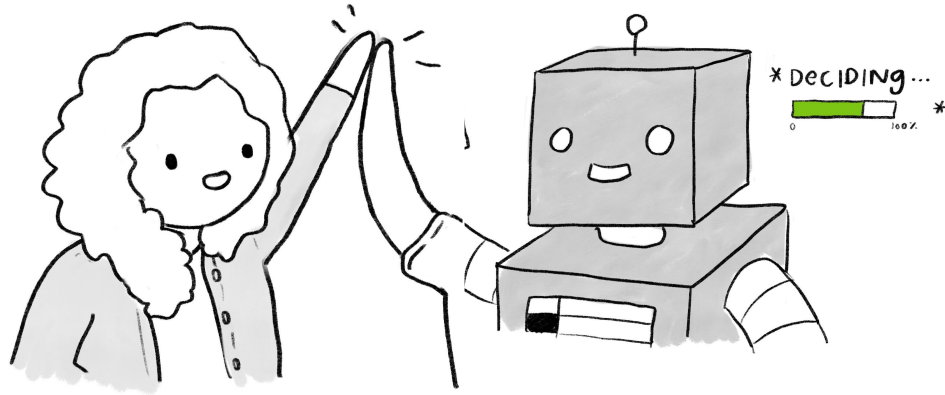




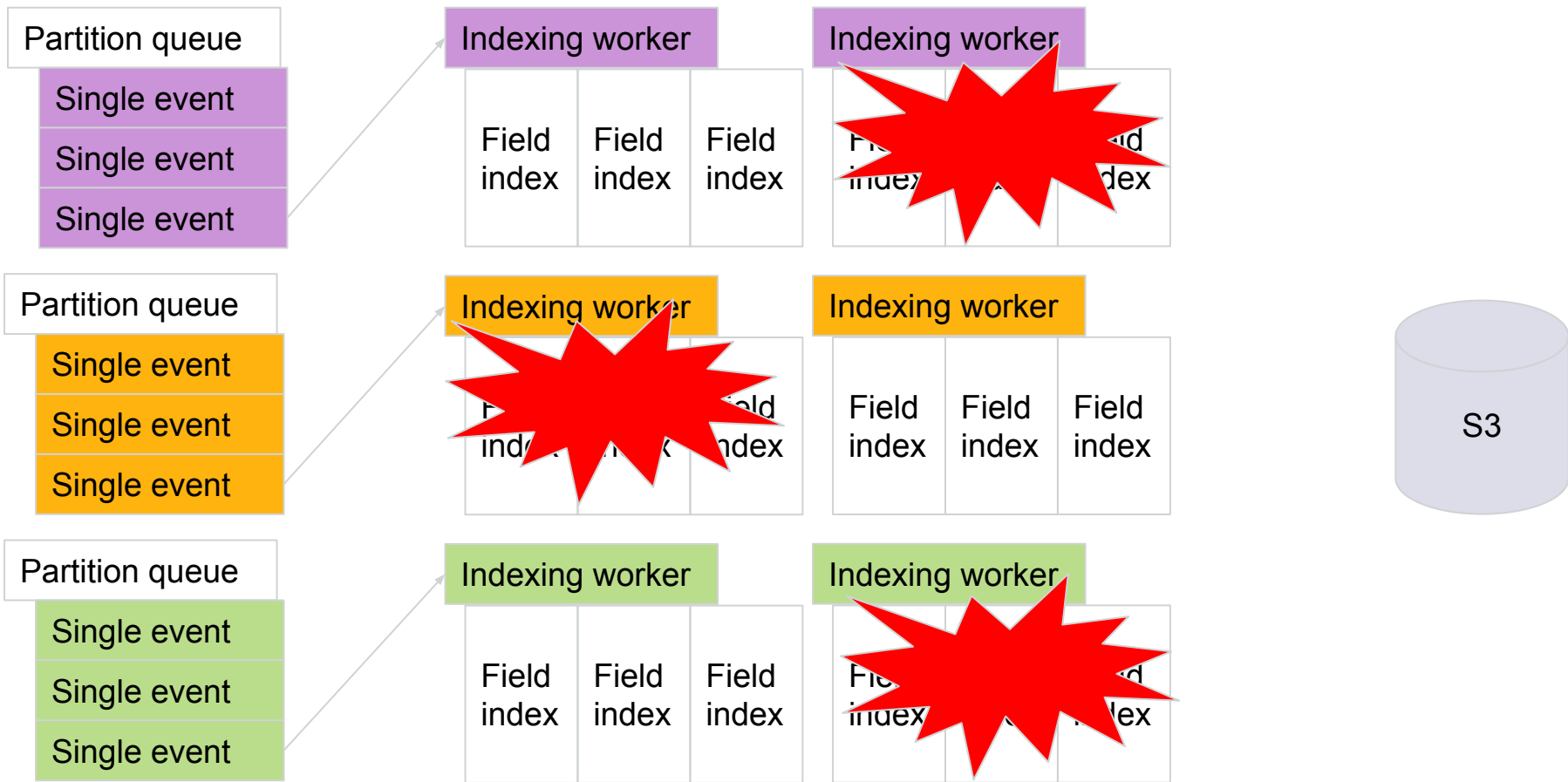


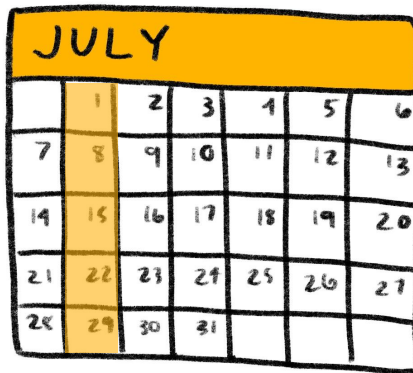






De-risk with **design & automation.**





Continuously verify to stop regression.



Spot and Kubernetes cause (un)planned chaos.



Chaos testing creates stronger platforms.

Not every experiment succeeds.

But you can mitigate the risks.

Three case studies of failure

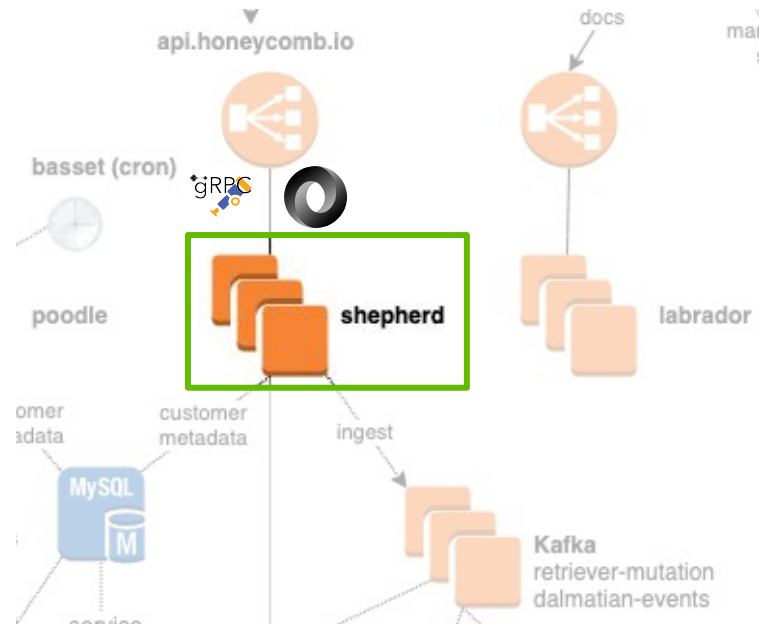
- Ingest service crash
- Kafka instability
- Query performance degradation

and what we learned from each.

1) Shepherd: ingest API service

Shepherd is the gateway to all ingest

- highest-traffic service
- stateless service
- cares about throughput first, latency close second
- used compressed JSON
- gRPC was needed.



Honeycomb Ingest Outage

- In November 2020, we were working on OTLP and gRPC ingest support
- Let a commit deploy that attempted to bind to a privileged port
- Stopped the deploy in time, but scale-ups were trying to use the new build
- Latency shot up, took more than 10 minutes to remediate, blew our SLO

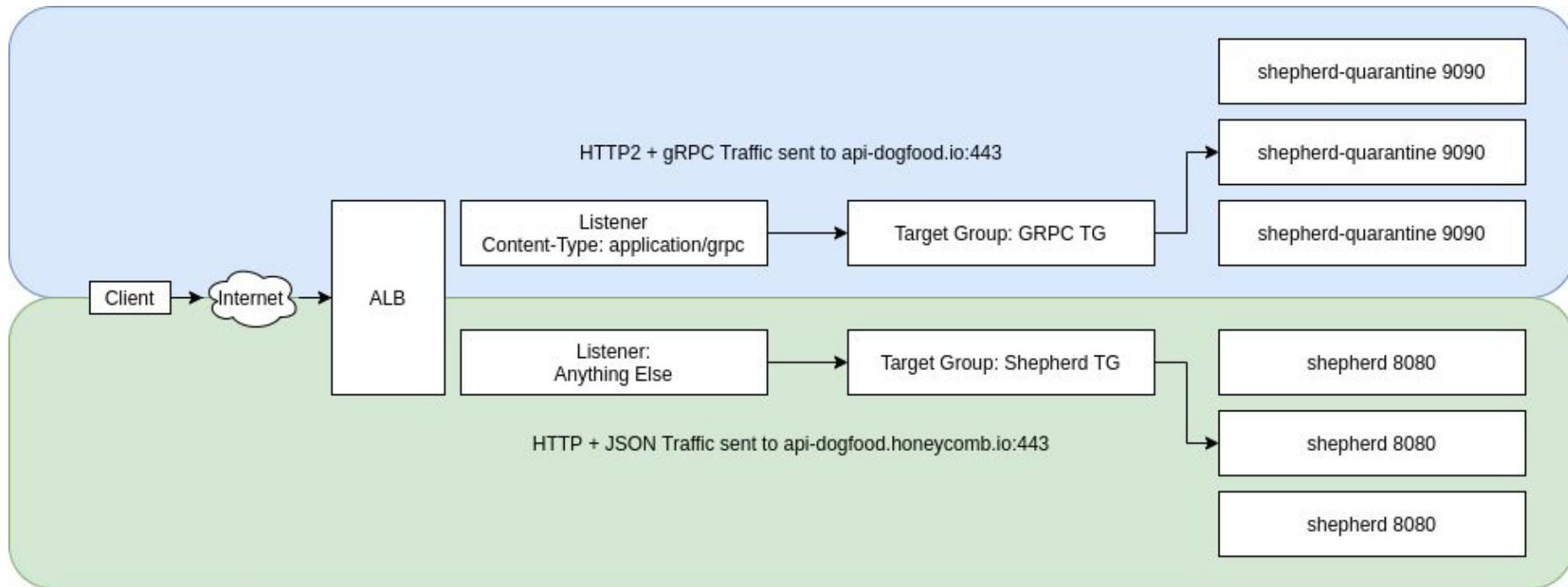


Now what?

- We could freeze deploys (oh no, don't do this!)
- Delay the launch? We considered this...
- Get creative!



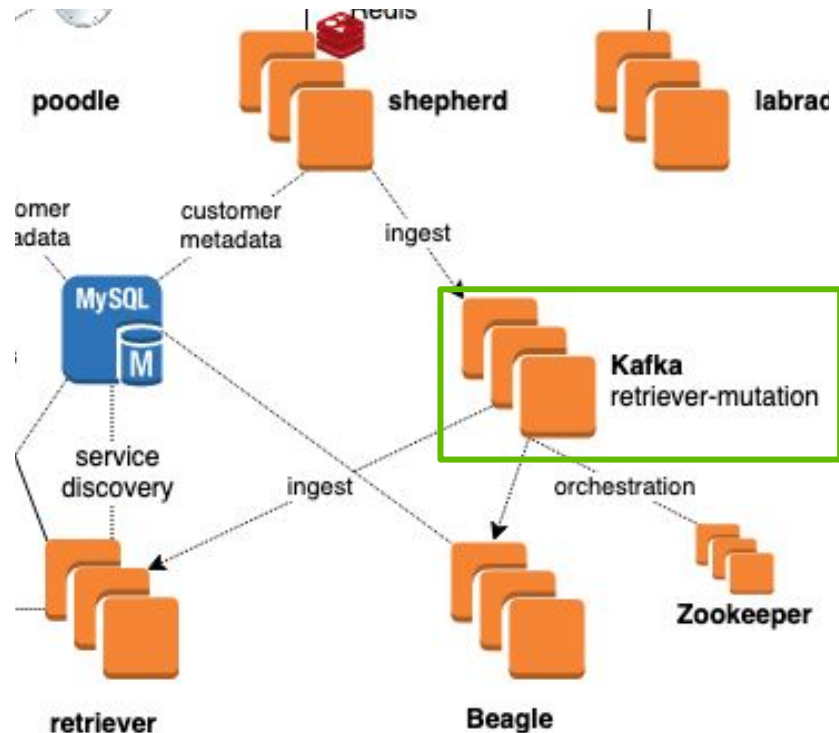
Reduce Risk



2) Kafka: data bus

Kafka provides durability

- Decoupling components provides safety.
- But introduces new dependencies.
- And things that can go wrong.



Our month of Kafka pain

Longtime Confluent Kafka users

First to use Kafka on Graviton2 at scale

Changed multiple variables at once

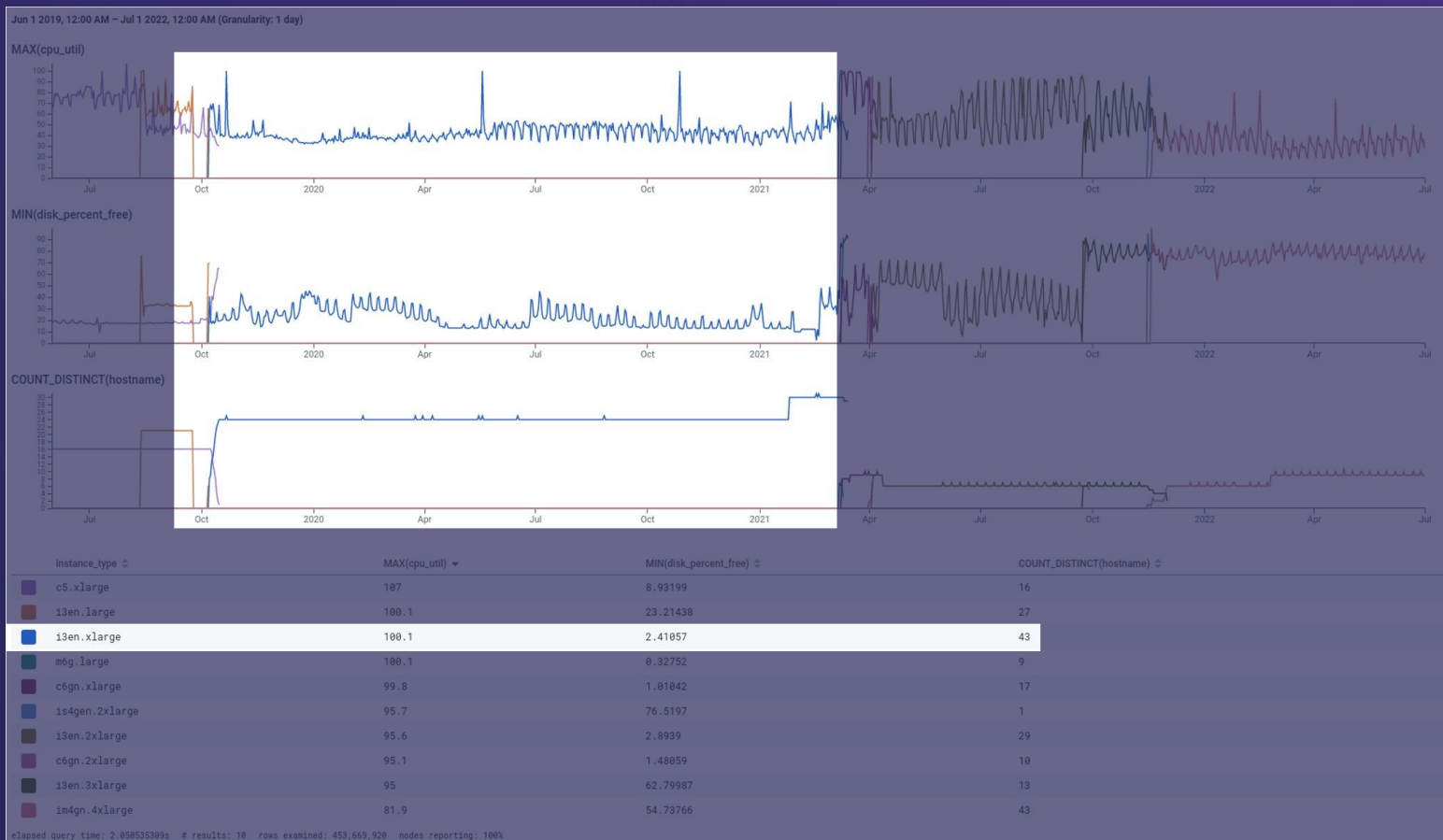
- move to tiered storage
- i3en → c6gn
- AWS Nitro



Read more: go.hny.co/kafka-lessons



Finding the right way to migrate Kafka



Unexpected constraints

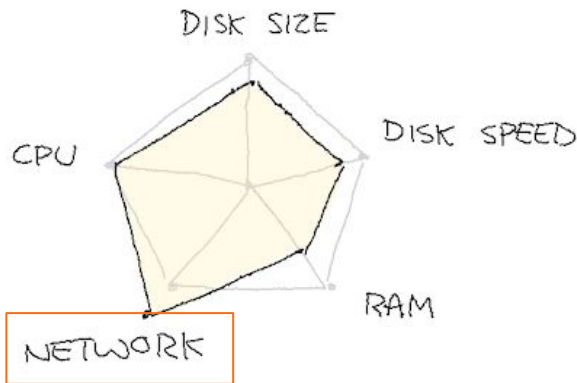
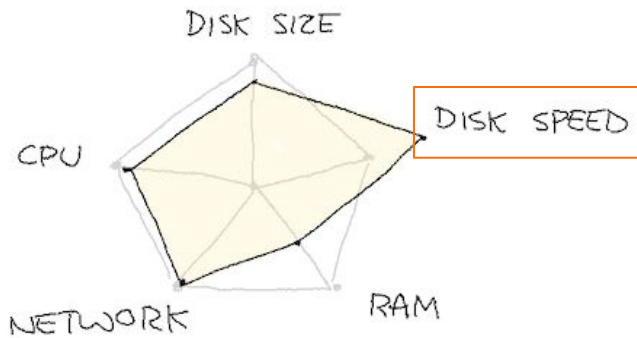
We thrashed multiple dimensions.

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!

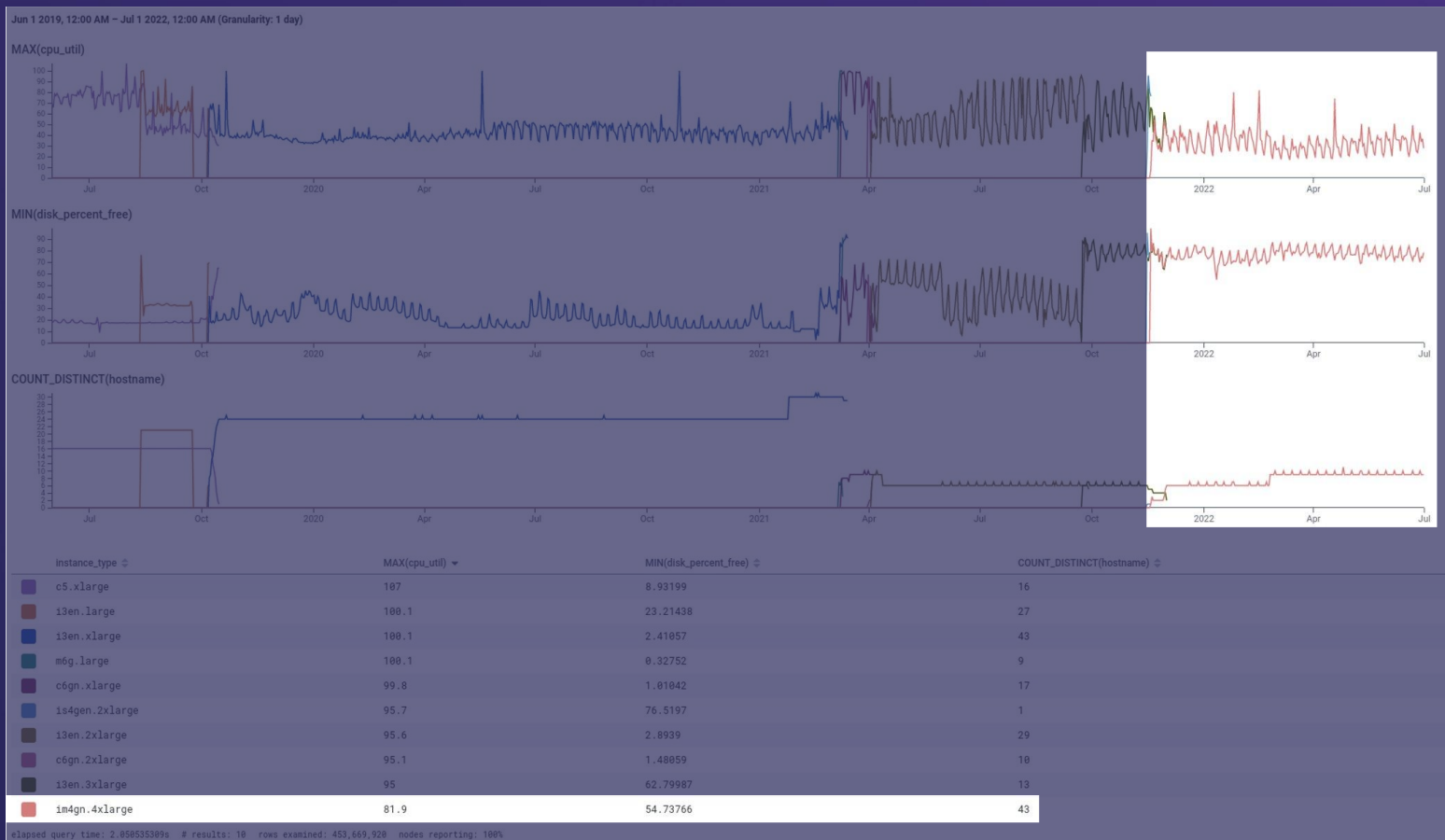


Read more: go.hny.co/kafka-lessons

Finding the right way to migrate Kafka



Finding the right way to migrate Kafka



Take care of your people

Existing incident response practices

- Escalate when you need a break / hand-off
- Remind (or enforce) time off work to make up for off-hours incident response

Official Honeycomb policy

- Incident responders are encouraged to expense meals for themselves and family during an incident



We hire adults.

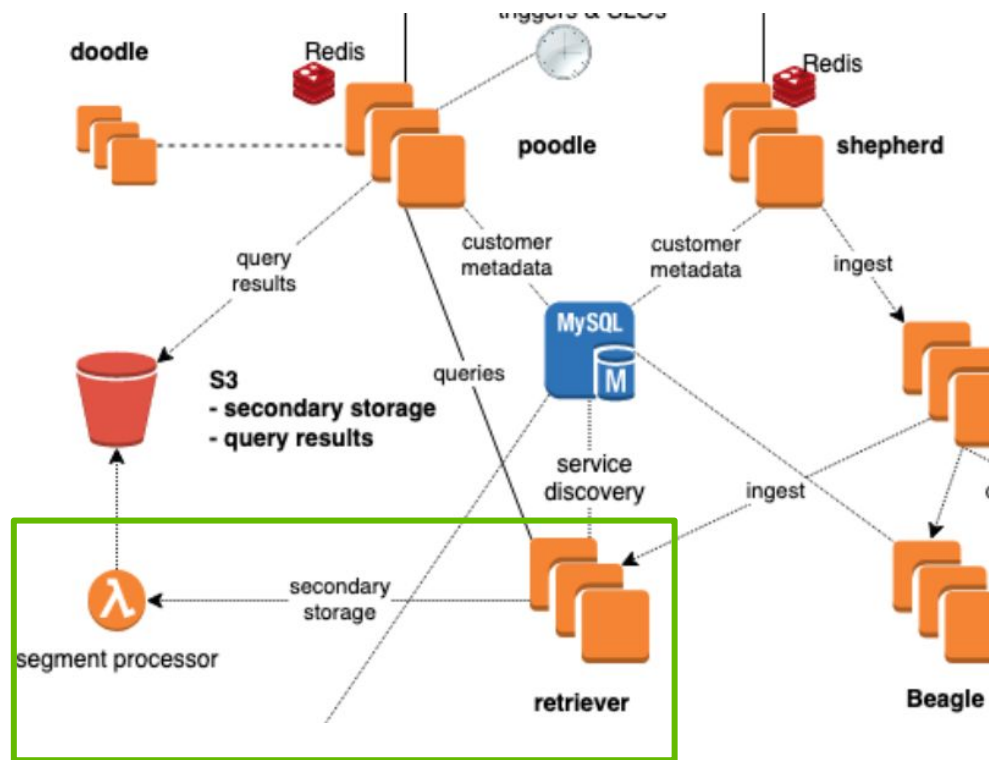
Pay attention to your mind and body so you can give and get help. All of us wobble, and being transparent about that means we can support each other. Participate fully in collaboration, coaching and management. If any group of us were together in a car on a long road trip, there would be no need for a dividing line in the back seat to keep people from hitting each other.



3) Retriever: query service

Retriever is performance-critical

- It calls to Lambda for parallel compute
- Lambda use exploded.
- Could we address performance & cost?
- Maybe.





If **infra is code, we can use CI & flags!**



LaunchDarkly APP 6:48 PM

Liz Fong-Jones updated the flag **Retriever Lambda ARM Percentage**

- Added the variation **1% ARM**

Liz Fong-Jones updated the flag **Retriever Lambda ARM Percentage** in **Production**

- Changed the default variation from **50% ARM** to **1% ARM**



lizf 🌙 6:49 PM

reverting ARM experiment, just keeping a trickle on 1% for validation of non-breakage/dogfooding of the lambda layer on both archs. it was 20% slower at p50 and 100% slower at p99, so we need to roll back.



1



1



1



1 reply 17 days ago



COUNT

arch exists

arch

Run a few
seconds agoORDER BY

COUNT desc

LIMIT

None

HAVING

None; include all results



Results BubbleUp Metrics Traces Raw Data

☐ Compare to

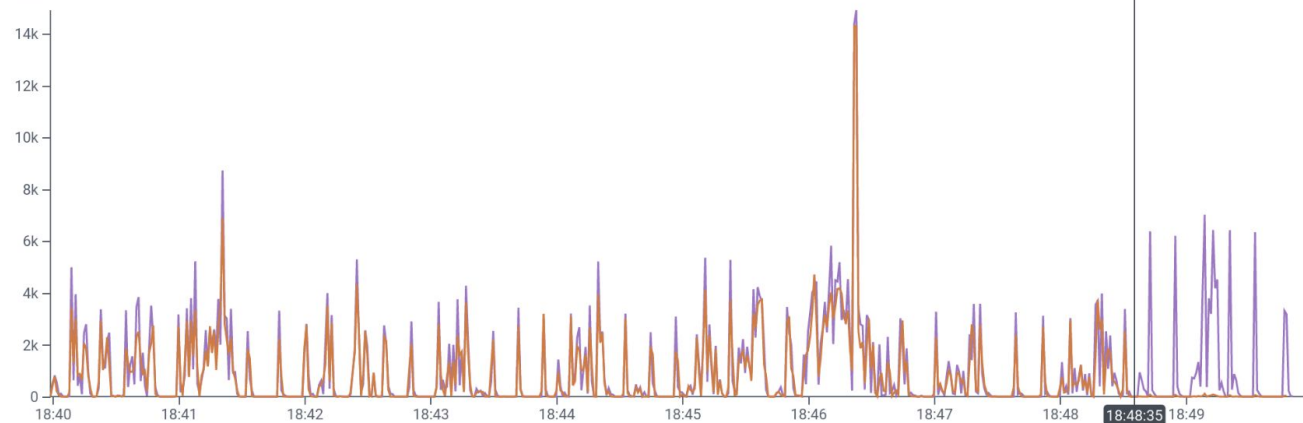
10 minutes prior



Graph Settings

Oct 1 2021, 6:39:59 PM – Oct 1 2021, 6:49:59 PM (Granularity: 1 sec)

COUNT



arch

COUNT



amd64

583,172



arm64

455,704



Making progress carefully



LaunchDarkly APP 11:06 AM

Liz Fong-Jones turned on the flag **Profile Lambda Percent** in **Production**

Liz Fong-Jones scheduled changes for the flag **Profile Lambda Percent** in **Production**

- Changes will occur on **Sat, 16 Oct 2021 18:15:00 UTC**
- Turn off the flag

Liz Fong-Jones scheduled changes for the flag **Retriever Lambda ARM Percentage** in **Production**

- Changes will occur on **Sat, 16 Oct 2021 18:20:00 UTC**
- Update default variation to **serve 1% ARM**



LaunchDarkly APP 11:15 AM

Completed scheduled changes to the flag **Profile Lambda Percent** in **Production** (via API)

- Turned the flag off



Today we're 99%+ ARM lambda



VISUALIZE

SUM(Fraction)

WHERE

None; include all events

GROUP BY

Resource

Run Query

Run a few
seconds ago

ORDER BY

SUM(Fraction) desc

LIMIT

None

HAVING

None; include all results

Results BubbleUp Traces Raw Data

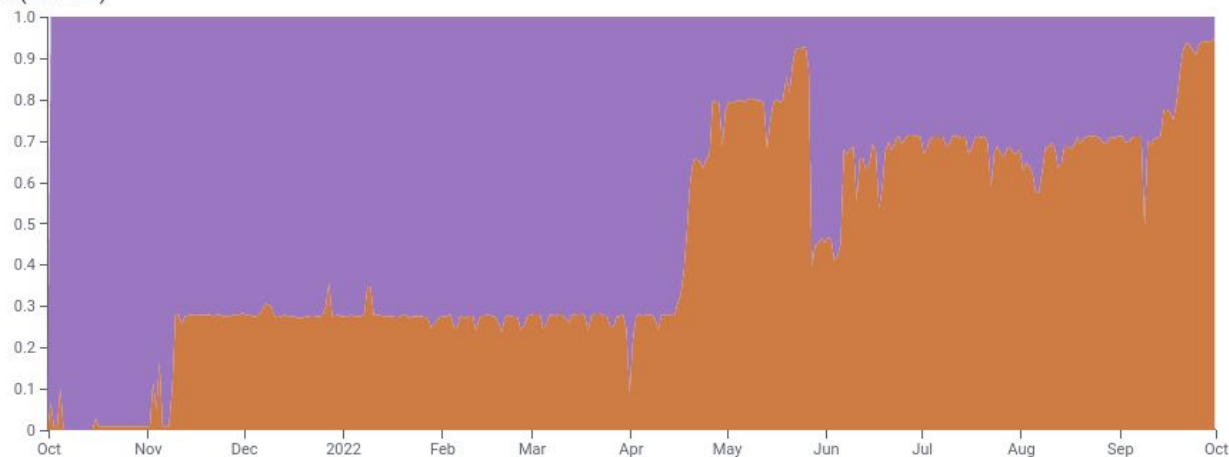
☐ Compare to

Previous time range

Graph Settings

Oct 1 2021, 12:00 AM – Oct 1 2022, 12:00 AM (Granularity: 1 day)

SUM(Fraction)



Fast and reliable: pick both!

Go faster, safely.

Takeaways

- Design for reliability through full lifecycle.
- Feature flags can keep us within SLO, most of the time.
- But even when they can't, find other ways to mitigate risk.
- Discovering & spreading out risk improves customer experiences.
- Black swans happen; SLOs are a guideline, not a rule.



Acknowledge hidden risks

Examples of hidden risks

- Operational complexity
- Existing tech debt
- Vendor code and architecture
- Unexpected dependencies
- **SSL certificates**
- **DNS**

Discover early and often through testing.

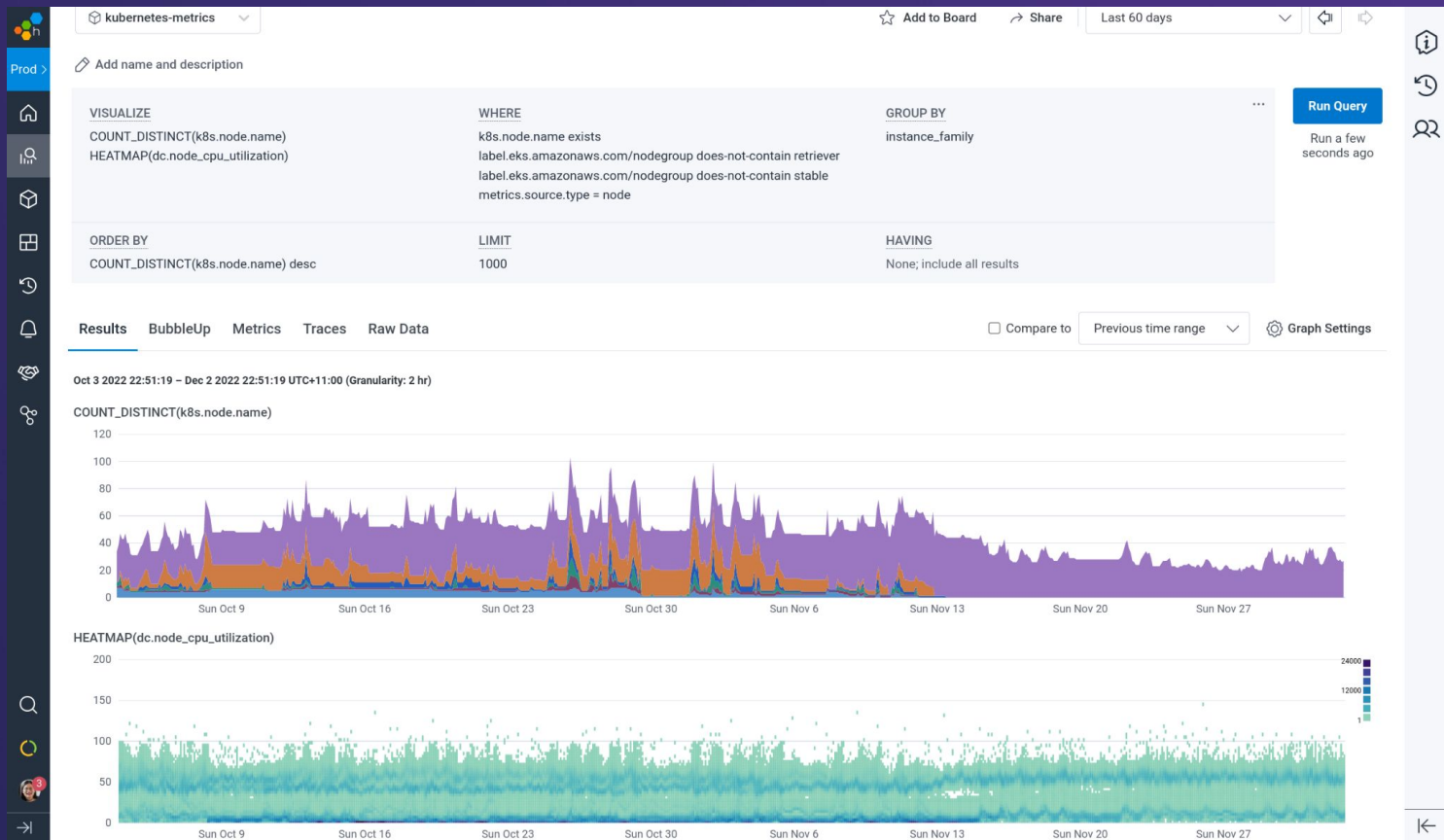


Takeaways

- We are part of sociotechnical systems. Customers, Engineers, Stakeholders.
- Outages and failed experiments are unscheduled learning opportunities.
- Nothing happens without discussions between different people and teams.
- DevOps is just talking to each other! Figuring out how to put customers first.



PS: we're now 100% Graviton3 k8s nodes.



Observability Engineering

Get our new book, free!

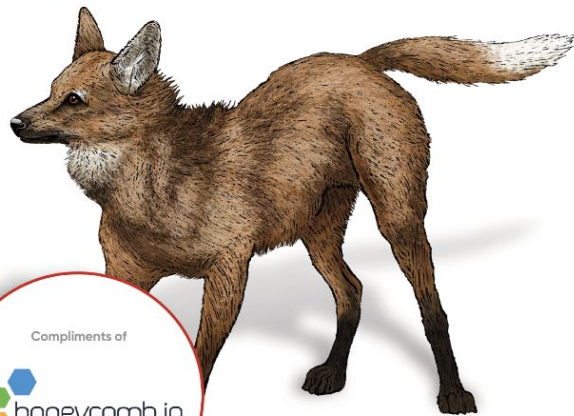


@lizthegrey

O'REILLY®

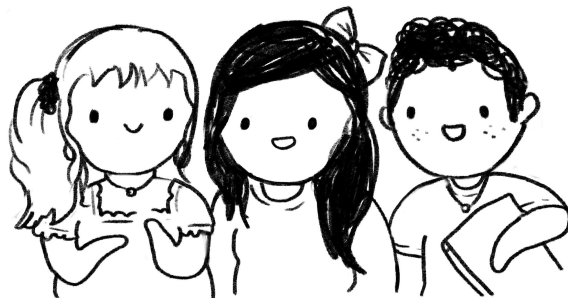
Observability Engineering

Achieving Production Excellence



Charity Majors,
Liz Fong-Jones
& George Miranda

Understand & control production.



Go faster on stable infra.
Manage risk and iterate.



lizthegrey.com; [@lizthegrey](https://twitter.com/lizthegrey)





www.honeycomb.io