# Using Serverless Functions for Real-time Observability

SRECon, 16 March 2022

Liz Fong-Jones @lizthegrey
Jessica Kerr @jessitron

honeycomb.io

**Liz Fong-Jones**

Principal Developer Advocate

@lizthegrey

**Jessica Kerr**

Principal Developer Advocate

@jessitron

# Today

How serverless is useful for on-demand compute

How serverless is painful for on-demand compute

How to experiment with serverless in your environment

# What is Lambda for?

Let's talk use cases of serverless

# What is ~~Lambda~~ for?

We'd like to optimize our custom datastore, Retriever

@jessitron

# What is Retriever for?

It's a distributed column store for **real-time event aggregation**

# What is ~~Retriever~~ for?

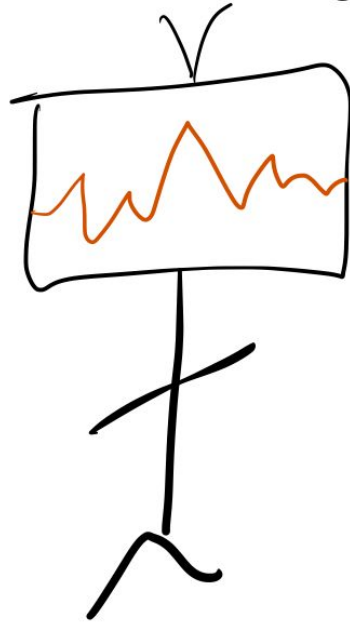**Real-time event aggregation** for interactive querying over traces

# What is Honeycomb for?

Observability: finding out what is going on
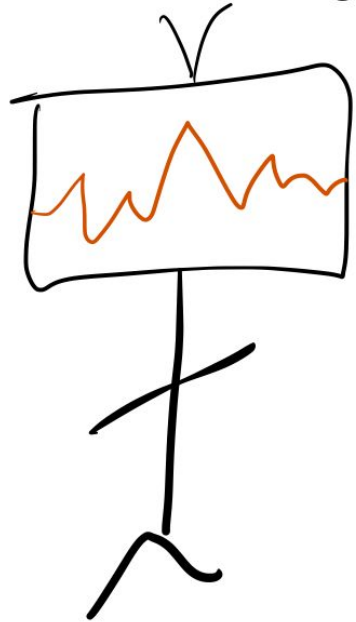(by querying traces!)
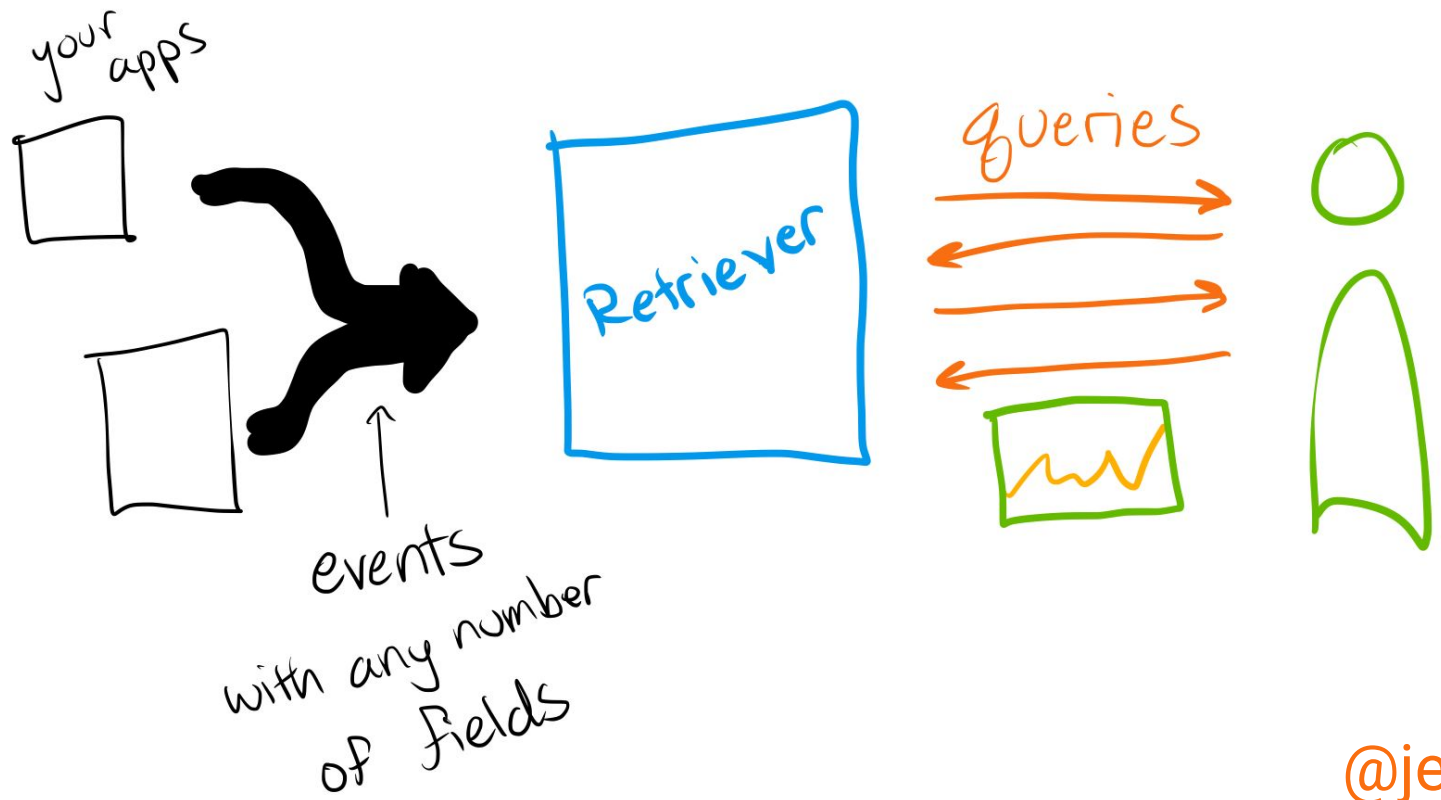
@jessitron

@jessitron

# Interactive investigation of production behavior

We run fast queries across any combination of fields.
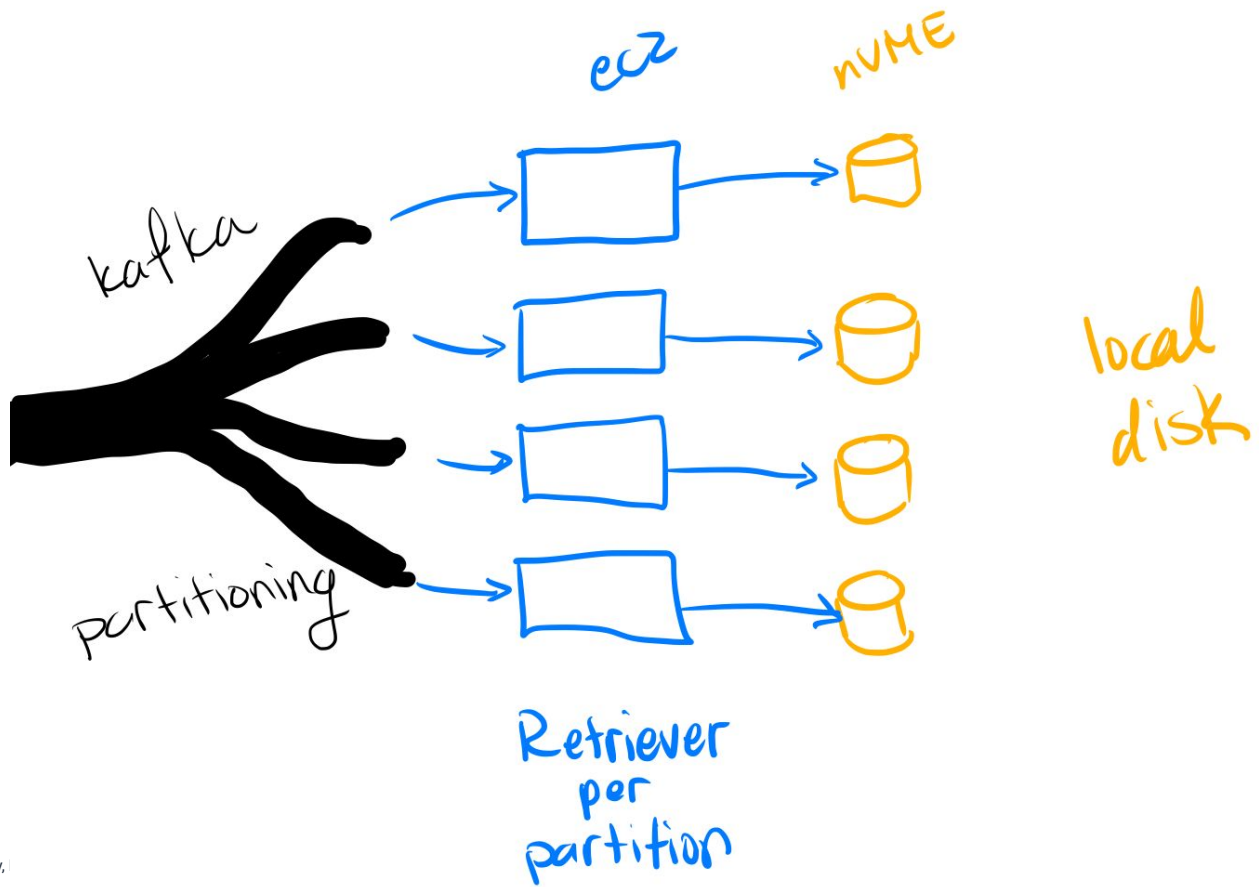
@jessitron

# Emphasis: *interactive*.

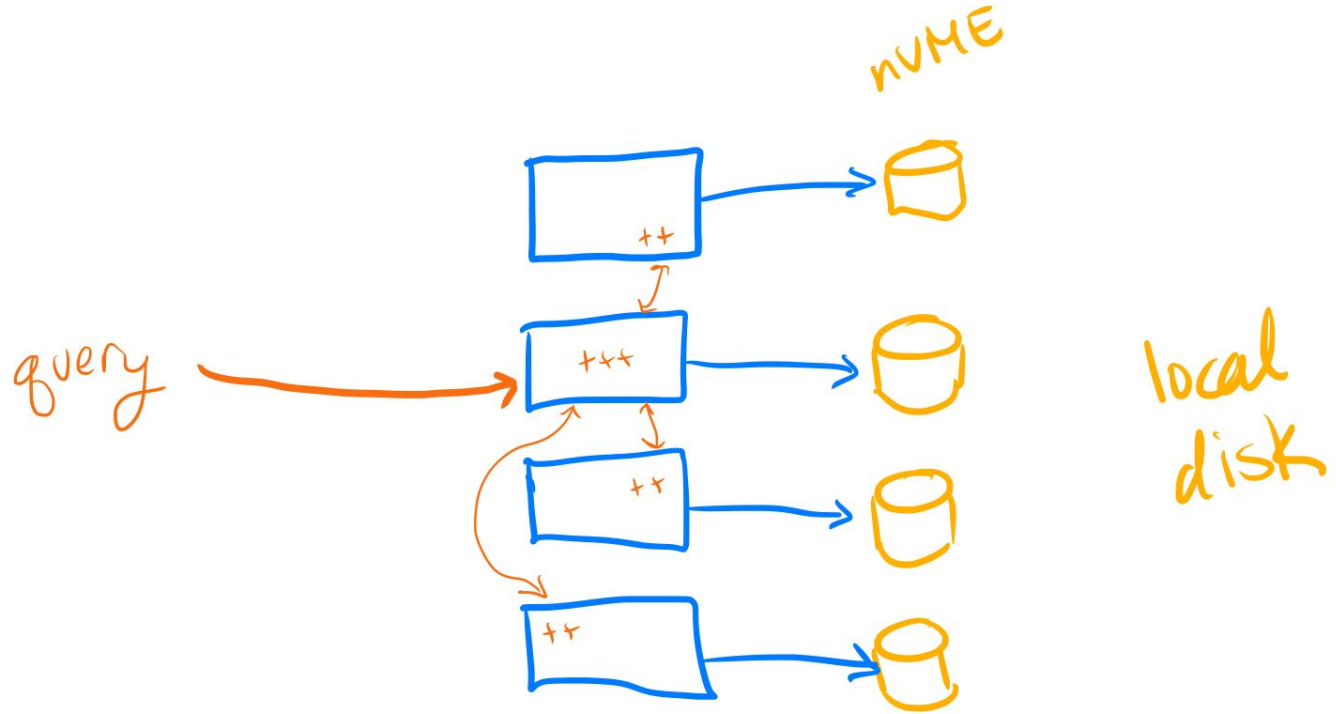100ms is fast. 1000ms is ok. 10sec is slow. 100sec is unacceptable.

# Retriever stores all your event data
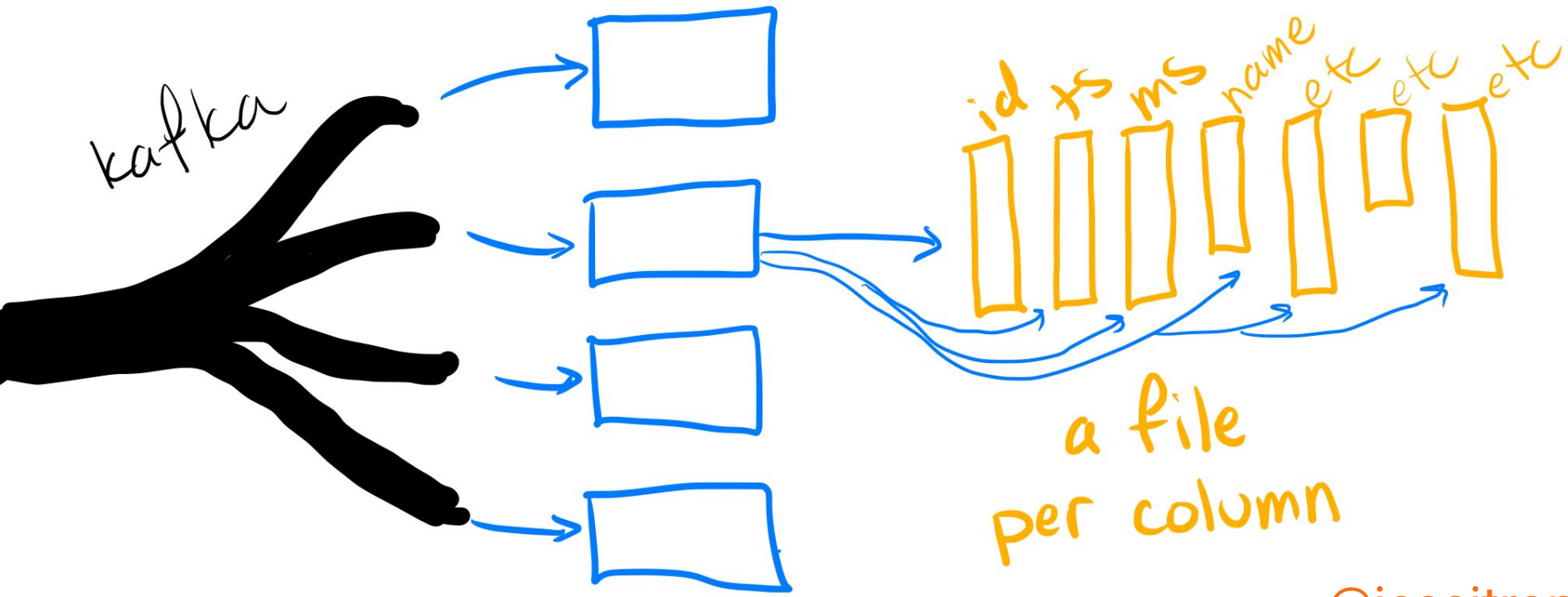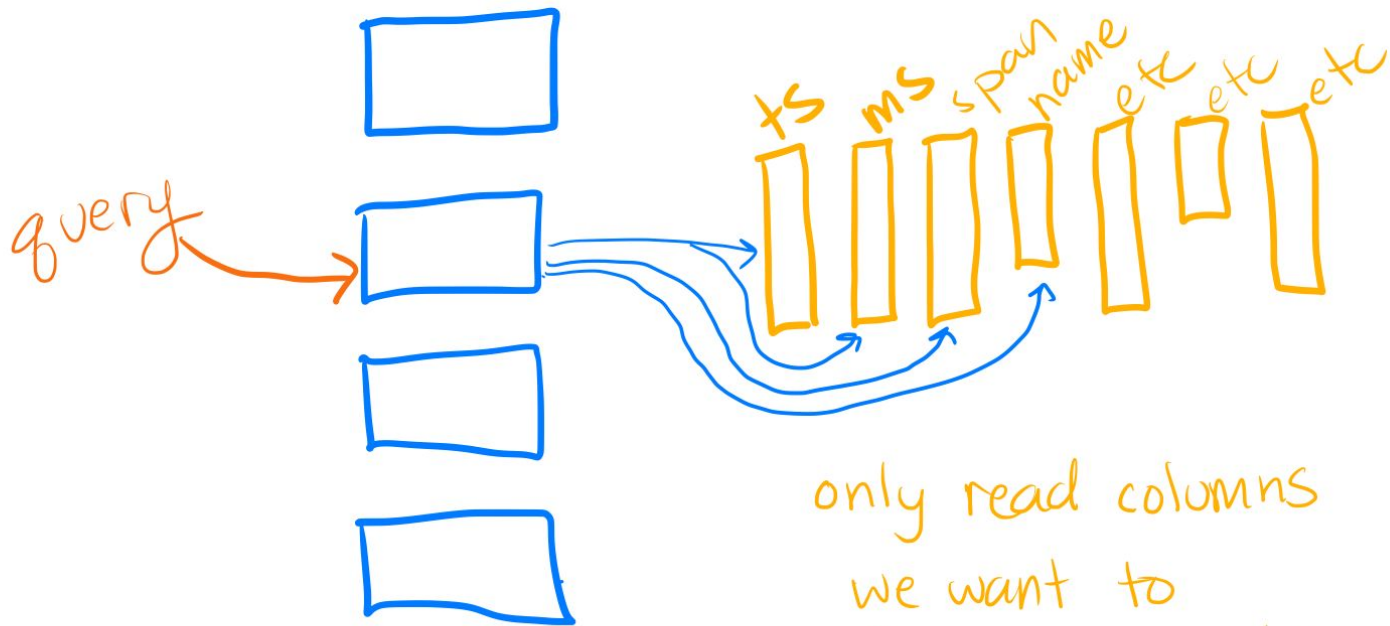
# Retriever is a distributed datastore.



ec2

nVME

kafka

partitioning

local disk

Retriever per partition

@jessitron

# Retriever is a distributed datastore.



query

nVME

local disk

each Retriever
reads & aggregates

@jessitron

# Retriever is a distributed column store.



@jessitron

# Retriever is a distributed column store.



query

ts ms span name etc etc etc

only read columns
we want to
filter, aggregate,
or group by.

@jessitron

# Retriever indexes segments by timestamp.

kafka

events
are
ordered
by arrival

segment by arrival.
break at a million, or 1Gb,
or 12hrs
Record the timestamp range
per segment.

@jessitron

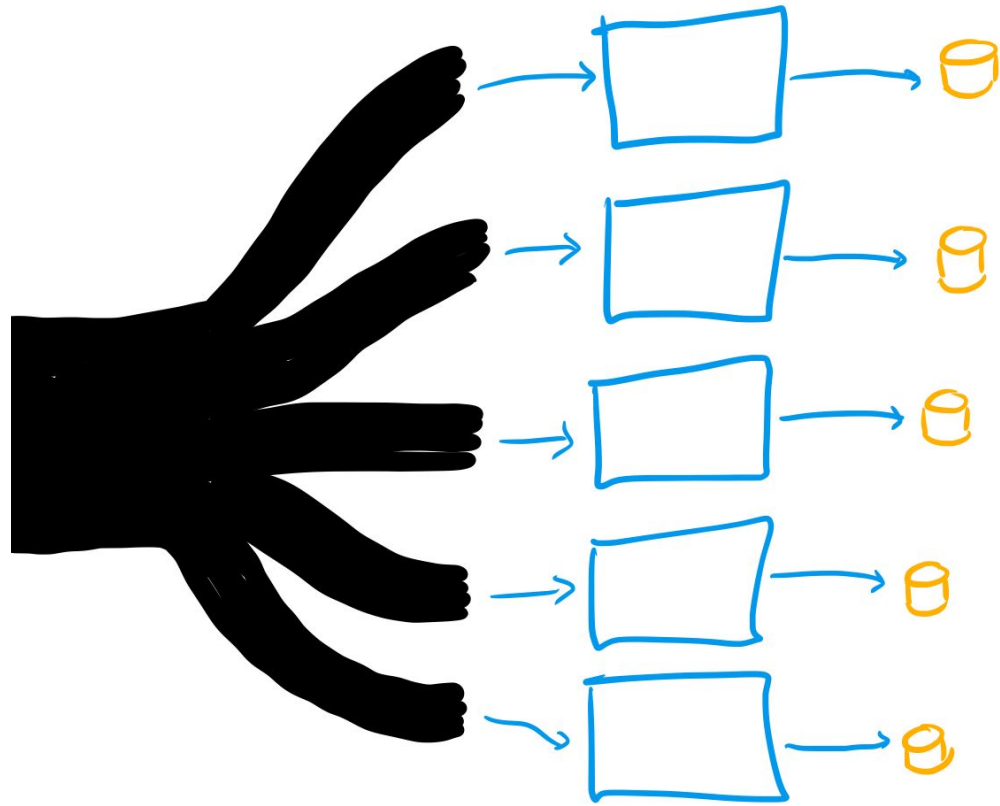# Retriever indexes segments by timestamp.



@jessitron

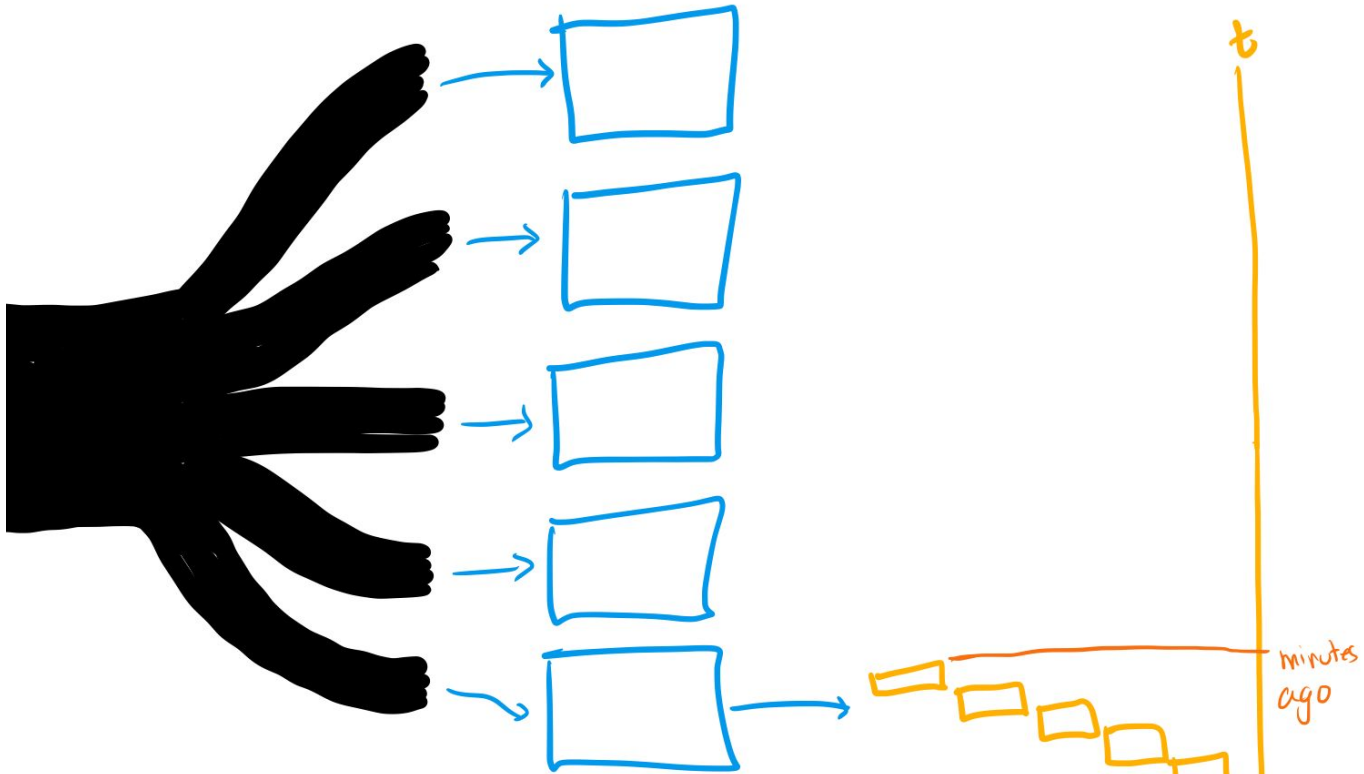# Retriever indexes segments by timestamp.

# Dynamic aggregation of any fields across any time range
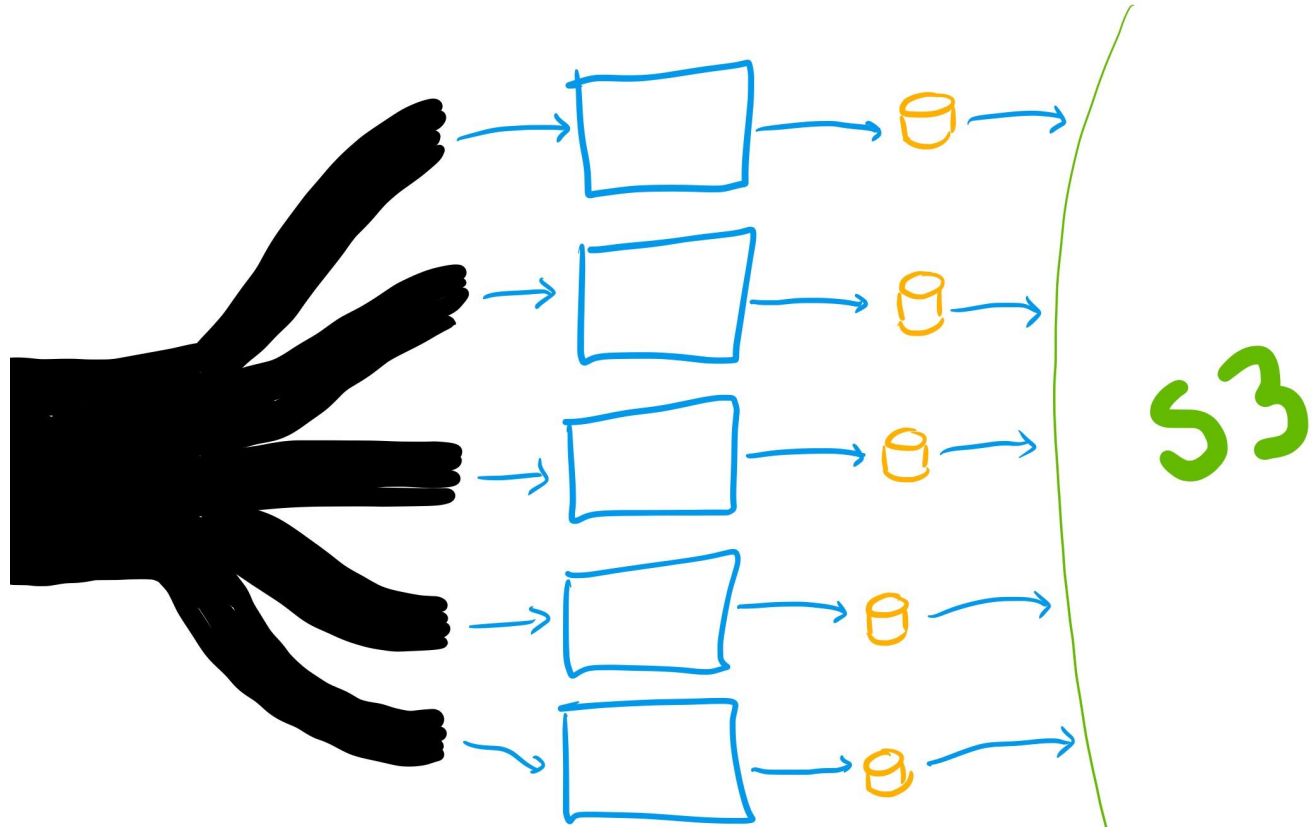
A custom datastore, carefully suited, continually optimized.
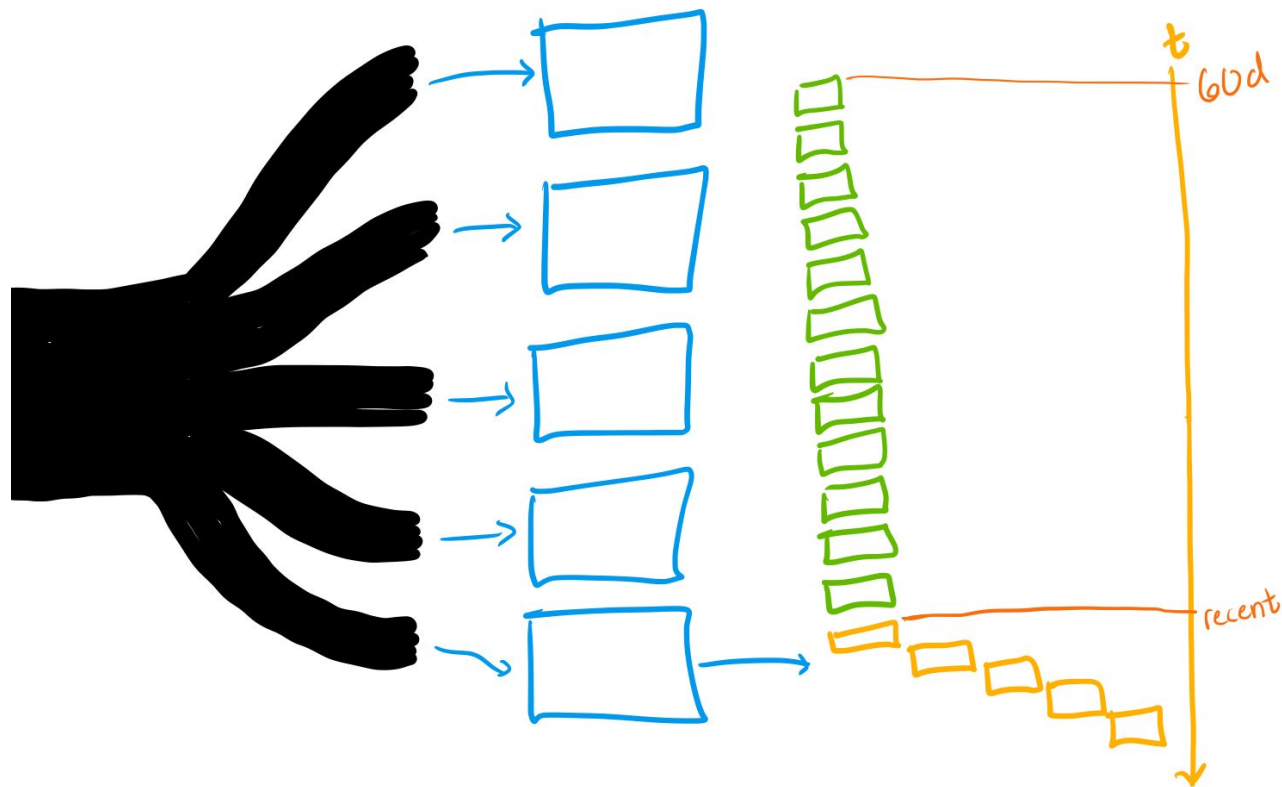
@jessitron

# Bigger customers, more data coming in.

@jessitron

# Segments hold a smaller time range.

@jessitron

# Solution: MOAR storage

@jessitron

# Now we can keep data for a fixed time range!

@jessitron

# Retrievers grab data back from S3 at need.

© 2021 Hound Technology, In

@jessitron

# Now people can run queries over 60 days 😯

@jessitron

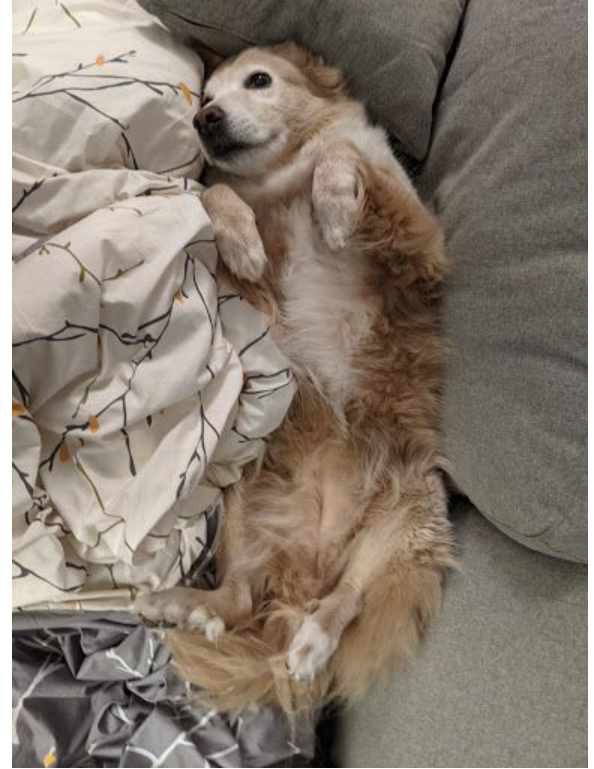# Now people can run queries over 60 days 😲

@jessitron

" Lots more compute to play with, pretty please!
but only if I want to play!

Retrievers

@jessitron

# MOAR compute, on demand.



@jessitron

# Problem: too much data for one retriever...

60-day query

t

query

@jessitron

# Solution: more compute, on demand.

@jessitron

# Increase in query time is sublinear

@jessitron

34

# Buy compute in ~~100ms~~1ms units

Compute scales with time range, so response time doesn't have to.

@jessitron

# Lambda scales* up our compute

**50ms**

median* startup time

**90%**

of ours return* within 1.5s

**3–4x**

as expensive* as EC2

@jessitron

# Considerations

Lambda *is* on-demand compute, but they didn't build it for this.

# Lambda scales up our compute

**50ms**

median startup time

**90%**

of ours return within 1.5s

**3–4x**

as expensive as EC2

@lizthegrey

# Lambda scales... within limits



@lizthegrey

# Lambda scales... within limits



absolute limit

increment

burst limit

AWS concurrency limit

requires 1 min sustained load

lambdas we try to run

@lizthegrey

# Observability helps: concurrency

# Lambda scales... within limits

Study your limits:

https://docs.aws.amazon.com/lambda/latest/dg/gettingstarted-limits.html

Change the SDK retry parameters

Observability helps 😉

Talk to your account reps

@lizthegrey

# Lambda scales up our compute

## 50ms
median startup time

## 90%
of ours return within 1.5s

## 3-4x
as expensive as EC2

@lizthegrey

# Functions start up… when they do

CONCURRENCY

Lambda invocations running

CONCURRENCY

Lambdas running/sleeping

@lizthegrey

# Lambda scales up our compute

**50ms**

median startup time

**90%**

of ours return within 1.5s

**3-4x**

as expensive as EC2

@lizthegrey

# Functions return... usually

Sep 28 2021, 9:14 PM − Sep 29 2021, 9:14 PM (Granularity: 1 min)

HEATMAP(duration_ms)

@lizthegrey

# Functions accept… JSON

Put the data in S3 and send a link.

@lizthegrey

# Functions return... up to 6Mb

**Sep 28 2021, 9:22 PM – Sep 29 2021, 9:22 PM (Granularity: 1 min)**

HEATMAP(app.response_size)



Put the data in S3 and send a link.

@lizthegrey

# Lambda scales up our compute

**50ms**

median startup time

**90%**

of ours return within 1.5s

**3-4x**

as expensive as EC2

@lizthegrey

# Functions cost... something

**Query run every 1440 minutes**   Define the calculation to perform and any relevant filters

| VISUALIZE | WHERE | AND ⌄ | GROUP BY |
|---|---|---|---|
| ✕ SUM(lambda_cost) | ✕ dataset_id exists | | ✕ dataset_id |

Triggering Queries are constrained to one calculation, and as many filters as you'd like.

Below, we show the `SUM(lambda_cost)` trends for each dataset_id (`where dataset_id exists`) for **the last 16** 1440-minute intervals.

The markers indicate the last 16 points at which the trigger would have run.

**Sep 15 2021, 8:15:43 PM – Oct 1 2021, 8:15:43 PM (Granularity: 1 day)**

SUM(lambda_cost)

Sep 30, 2021 7:00 PM

## Threshold

Trigger notification if returned `SUM(lambda_cost)` for any dataset_id is   `>= ⌄`   `300`

# Functions cost... let's make it less?

## AWS Lambda Functions Powered by AWS Graviton2 Processor – Run Your Functions on Arm and Get Up to 34% Better Price Performance

by Danilo Poccia | on 29 SEP 2021 | in AWS Lambda, Compute, Graviton, Serverless | Permalink | 💬 Comments | ➤ Share

@lizthegrey

> **M6g instances are superior to C5 in every aspect—they cost less, have more RAM, exhibit lower median and significantly narrower tail latency, and run cooler with the same proportional workload per host. Converting our entire ingest worker fleet has allowed us to run 30% fewer instances, and each instance costs 10% less.**

Yours Truly

# Observability helps!



P99(duration_ms)

P50(duration_ms)

| arch | COUNT | HEATMAP(Log_Duration) | P99(duration_ms) | P50(duration_ms) |
|---|---|---|---|---|
| amd64 | 262,988 | | 1,168.09377 | 139.24663 |
| arm64 | 161,394 | | 2,677.62006 | 175.50275 |

| COUNT | arch exists<br>name = processSegment | arch | Run a few<br>seconds ago |
|---|---|---|---|

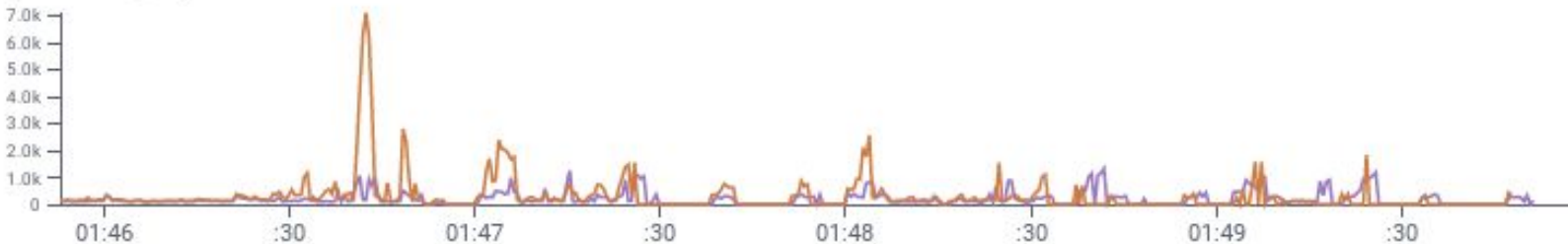| ORDER BY | LIMIT | HAVING |
|---|---|---|
| COUNT desc | None | None; include all results |

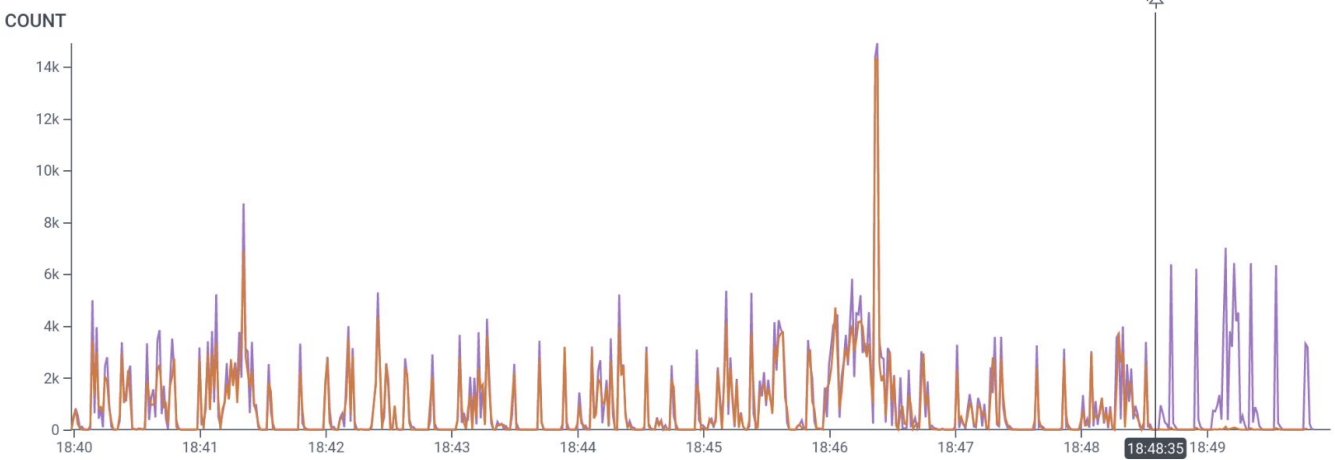Results    BubbleUp    Metrics    Traces    Raw Data          ☐ Compare to    [ 10 minutes prior ⌄ ]    ⚙ Graph Settings

**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 |

@lizthegrey

# Why so slow?

- AWS capacity constraints
- Go register calling convention
- lz4 library asm optimization

@lizthegrey

# 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`
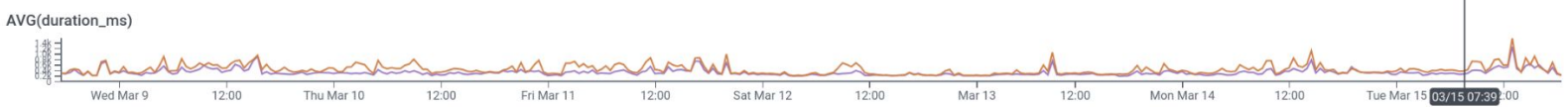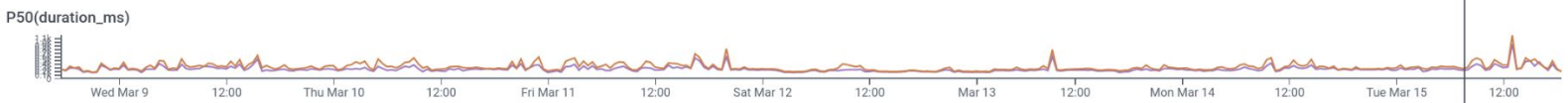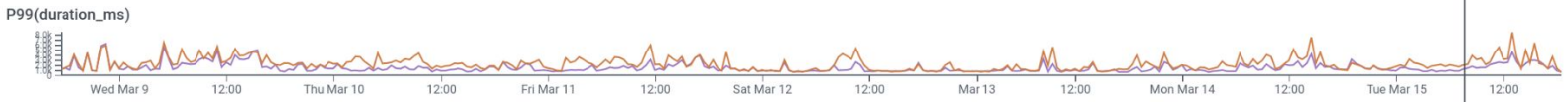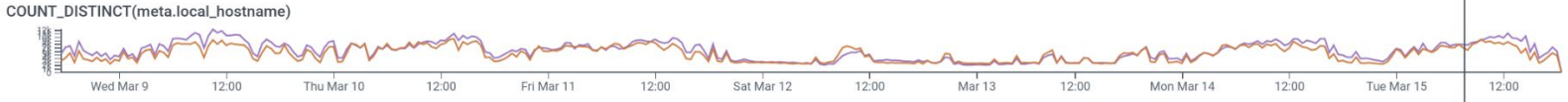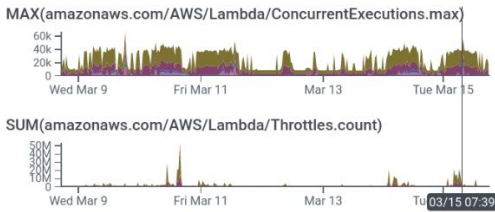- Update default variation to **serve** `1% ARM`

**LaunchDarkly** `APP` 11:15 AM

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

- Turned the flag off

@lizthegrey

MAX(cpu_util)

MAX(amazonaws.com/AWS/Lambda/ConcurrentExecutions.max)

SUM(amazonaws.com/AWS/Lambda/Throttles.count)

COUNT

HEATMAP(Log_Duration)

COUNT_DISTINCT(meta.local_hostname)

P99(duration_ms)

P50(duration_ms)

AVG(duration_ms)

| arch | global.go_runtime_shortversion | COUNT | HEATMAP(Log_Duration) | COUNT_DISTINCT(meta.local_hostname) | P99(duration_ms) | P50(duration_ms) | AVG(duration_ms) |
|---|---|---|---|---|---|---|---|
| amd64 | 1.17 | 2,436,396,531 | | 16,194 | 1,940.61384 | 242.42897 | 325.01033 |
| arm64 | 1.18 | 953,247,976 | | 16,071 | 2,869.15901 | 307.28266 | 477.1254 |

elapsed query time: 8.285680504s

# Yes*, do this at home!

# Most realtime bulk workloads benefit

- **Move** state from local machines onto object storage
- **Shard** list of objects into work units
- **Parallelize** object processing
- **Reduce** results outside Lambda afterwards

@lizthegrey

# Just beware the dragons

- **Avoid latency-insensitive** batch workloads (cost)
- **Avoid tiny** workloads (set-up latency)
- Check **cloud provider limits**, state your intentions (capacity planning)

@lizthegrey

# Do this before scaling out

- Ensure it's **tuned properly** (items/invoke, CPU/RAM ratio)
- Ensure your code is **optimized properly** (esp if multi-arch)
- Ensure you use **observability layers** (e.g. OTel layer)
- Measure **metrics** carefully (esp cost)

@lizthegrey

**Remember: nothing matters
unless users (developers) are happy**

# Observability Engineering

Explore preview chapters
from our new book

@lizthegrey

@honeycombio   @jessitron   @lizthegrey



www.honeycomb.io

https://www.honeycomb.io/blog/speeding-things-up-so-your-queries-can-bee-faster/

AWS Lambda Instrumentation | Honeycomb