# Driving Service Ownership with Distributed Tracing

Daniel "Spoons" Spoonhower, CTO and Co-founder

Lightstep

# Who am I?

**spoons (aka Daniel Spoonhower)**
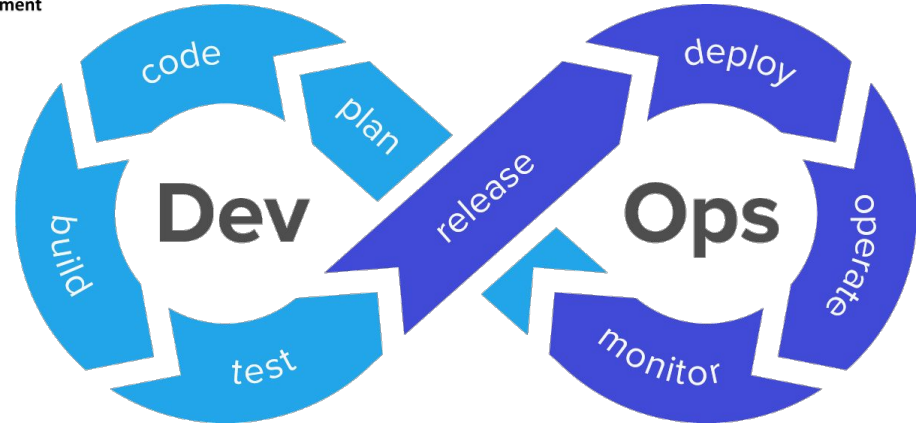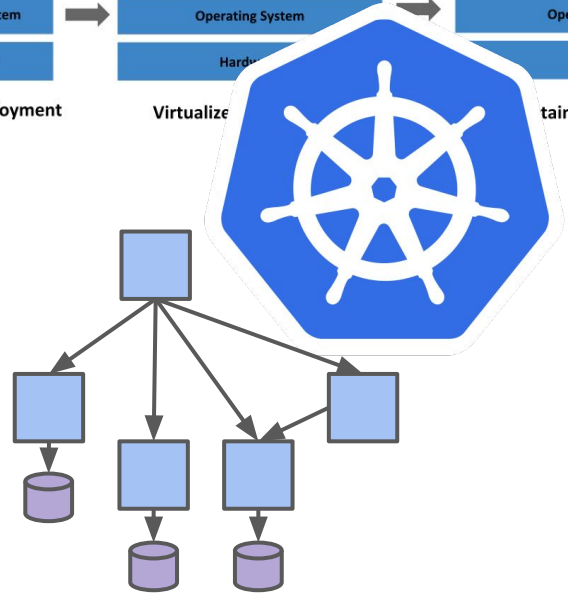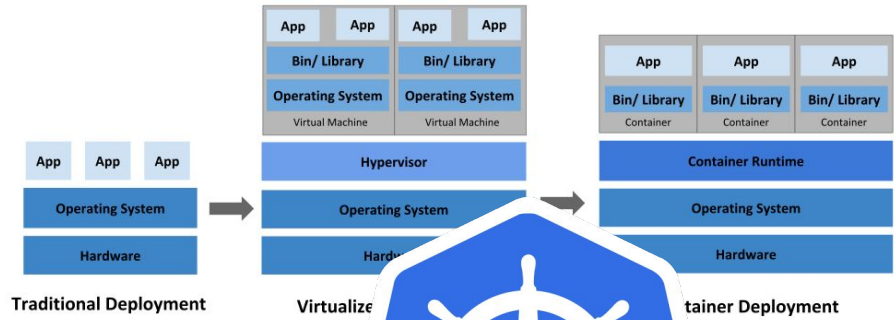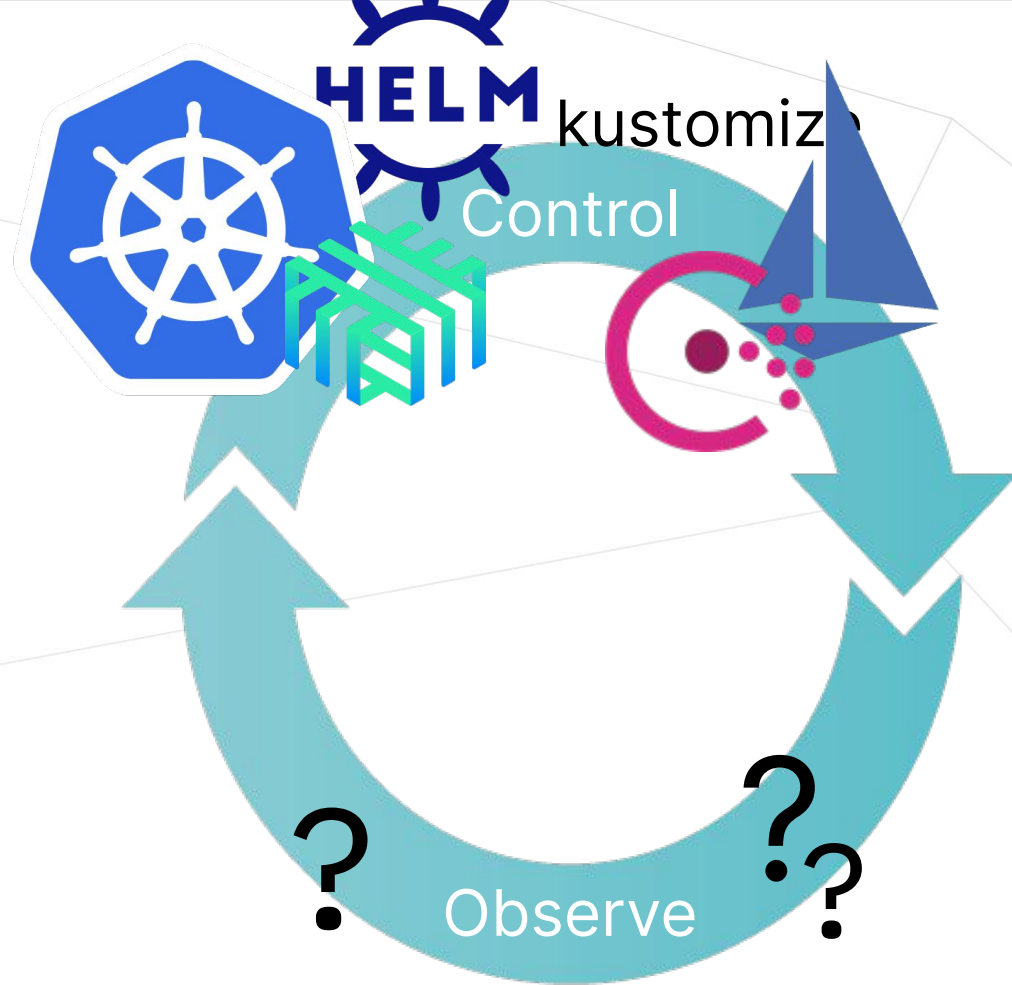*CTO and Co-founder, Lightstep*

@save_spoons

spoons@lightstep.com

# What Changed?

*Cloud Native, Microservices & DevOps*

HELM kustomize

Control

Observe

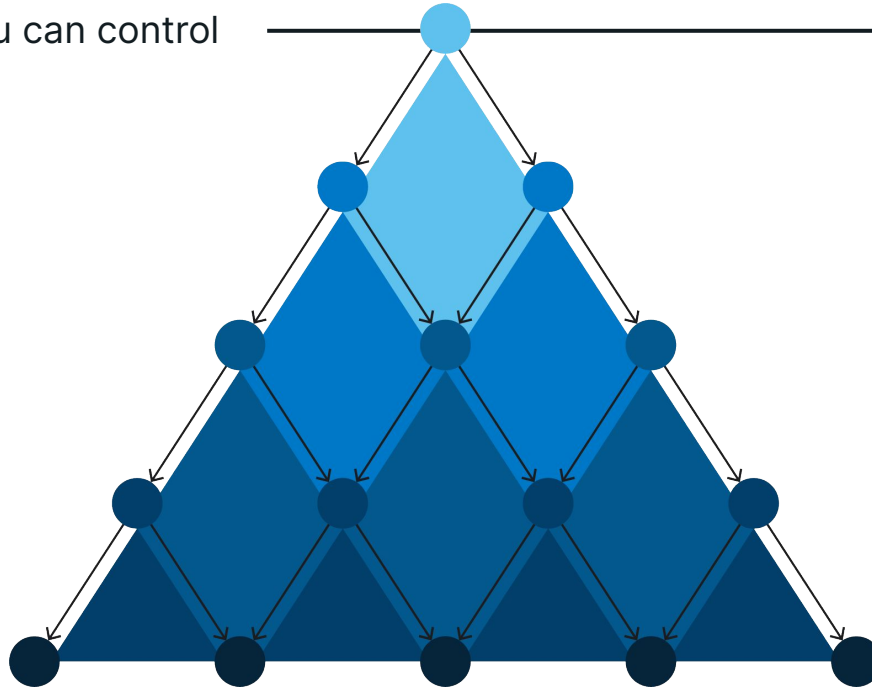# *Stress (n):* responsibility without control



What you can control

What you are responsible for

# Closing the gap between control and responsibility

Responsibility for delivering **performance** and **reliability**

Like many problems, the solution requires:

- Having the right data
- Setting the right goals
- Giving teams ownership

*Distributed tracing is essential to closing this gap*

Ownership means...

**Accountability**

+
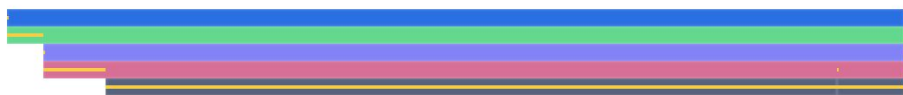
**Agency**

Loosely coupled services require...

**Distributed Tracing**

# Distributed Tracing

12:24:53
02/01

106μs   0.3s   0.6s   0.9s   1.2s   1.35s

Operation
/lightstep.database.ProjectService/GetProjects

Service
tracea

apptrace   Traces

Trace ba4782064cbf6ec4

Type to filter...

OPERATION

cache_struct.get_p

5

/lightstep.database

4

/lightstep.database

3

project_service.get

2

sql/select_one                    bento: gorp client

sql/select_one                    bento: gorp client

def            def (124ms)
SQL
SQL            SQL (1ms)
repo           repo (19ms)
SQL
SQL
SQL            SQL (1ms)
SQL: perms     SQL: perms (1ms)
SQL: perms
repo.commit    repo.commit (20ms)
SQL            SQL (1ms)
SQL            SQL (1ms)
SQL
SQL: perms
SQL: perms     SQL: perms (1ms)
               (6ms)

SQL            SQL (1ms)
repo.refresh-vcs-data   repo.refresh-vcs-data (10ms)
29.816ms : call-tchannel
12.153ms : call_in_request_context     8699Z"
9.712ms : endpoint

Jaeger UI   Lookup by Trace ID...   Search   Dependencies

⌄ frontend: HTTP GET /dispatch
Trace Start: July 20, 2018 2:48 PM   Duration: 732.11ms   Services: 6   Depth: 5   Total Spans: 51

Service & Operation                0ms        183.03ms      366.06ms      549.08ms

⌄ frontend HTTP GET /dispatch
  ⌄ frontend HTTP GET: /customer                                          305.89ms
    ⌄ frontend HTTP GET /customer                                         305.88ms
      ⌄ customer HTTP GET /customer                                       304.46ms
        ▸ mysql SQL SELECT                                                304ms
  ⌄ frontend Driver.findNearest                            211.19ms
    ⌄ driver Driver.findNearest
      ● redis FindDriverIDs
      ● redis GetDriver
      ● redis GetDriver
      ● redis GetDriver
      ● redis GetDriver
      ● redis GetDriver
      ● redis GetDriver
      ● redis GetDriver
      ● redis GetDriver
      ● redis GetDriver
      ● redis GetDriver
  ⌄ frontend HTTP GET /route
    ⌄ frontend HTTP GET
      ⌄ route HTTP GET /route
  ⌄ frontend HTTP GET /route
    ⌄ frontend HTTP GET
      ⌄ route HTTP GET /route
  ⌄ frontend HTTP GET /route
    ⌄ frontend HTTP GET
      ⌄ route HTTP GET /route
  ⌄ frontend HTTP GET /route
    ⌄ frontend HTTP GET
      ⌄ route HTTP GET /route

Duration: 209.323ms   Services: 5   Depth: 7   Total Spans: 24

Expand All   Collapse All   Filter Service Se...  ▾

client x4   flask-server x10   missing-service-name x2   tchannel-server x2   tornado-server x11

Services                      41.864ms        83.729ms
⌄ client               -181.126ms : client-calls-server-via-get
  ⌄ flask-server       -180.527ms : get
    flask-server              605μ : mysqldb:connect
    flask-server              54.152ms : mysqldb:select
    flask-server
    flask-server                        394μ : mysqldb:c
    flask-server                        46μ : mysqldb:t
    flask-server                        40.910ms : my
    tornado-server
    tornado-server
    tornado-server
    tornado-server
    tornado-server
    tornado-server
    tornado-server
    tornado-server
    tornado-server
    tchannel-server

100ms

# Relationships matter

calls

O'REILLY®

# Distributed Tracing in Practice

Instrumenting, Analyzing, and Debugging Microservices

Austin Parker,
Daniel Spoonhower,
Jonathan Mace
& Rebecca Isaacs
Foreword by Ben Sigelman

Google Technical Report dapper-2010-1, April 2010

## Dapper, a Large-Scale Distributed Systems Tracing Infrastructure

Benjamin H. Sigelman, Luiz André Barroso, Mike Burrows, Pat Stephenson,
Manoj Plakal, Donald Beaver, Saul Jaspan, Chandan Shanbhag

Tra                                tween callers and callees

# Traces are the raw material, not the finished product

Distributed *traces* – basically just structs

Distributed *tracing* – **the art and science of deriving value from traces**

# Service Ownership

*Benefits & Risks*

# Service ownership, defined

Ownership means teams are responsible for delivery of their software and services

Includes responsibilities like:

- Incident response
- Cost management
- Fixing bugs

Compare to DevOps...

- As a engineering culture
- As a set of tools
- As a feedback loop between developers and their users

# Benefits

Team are more independent $\Rightarrow$ higher developer velocity

Engineering team performance tied to real business metrics

- Application developers
- Platform engineers

# Risks

Teams are more independent $\Rightarrow$ more frameworks & tools

- Higher vendor costs
- More training new team members & internal transfers
- Harder to get a "big picture" view of your application

# Managing trade-offs around benefits and risks

Allowing for independence

Define clear responsibilities and goals

...while ensuring consistency

Measure progress toward goals and hold teams accountable

# Driving Toward Ownership

# 3 pieces to the puzzle

Docs

Oncall

SLOs

# Centralized documentation



Start with *expertise* then *ownership*

Make it easy to find related...

- Telemetry & dashboards
- Alert definitions
- Playbooks

Use a template!

Track last-modified dates

- Require periodic audits & updates

# ~~Centralized~~ documentation

## MACHINE-READABLE

```
       | Services

192      - name: kafka
193        team: data_ingest
194        properties:
195          build_type: helm
196          k8s_type: statefulset
197          skip_custom_values: true
198        environments:
199          staging:
```

**Purpose**

Messenger is designed to provide an abstraction layer...
Lightstep web application.

**Architecture**

Make documentation *machine-readable*

Use it to generate:

- Dashboard config
- Escalation policy config
- Deployment pipeline config
- ...

Make documentation *necessary* for day-to-day work

# ~~Centralized~~ documentation

*MACHINE-READABLE + DYNAMIC!*

It's hard to keep service dependencies up to date *manually...*

**So don't!**

- Use telemetry from the application

# Why is documentation important?

Record who is accountable

Automate many mundane tasks

Train new team members

Build confidence

# Oncall

Oncall is (often) responsible for...

- Incident response
- Communicating status internally & externally
- Production change management
  - Deploying new code
  - Pushing infrastructure changes
- Monitoring dashboards
- Low-urgency alert triage
- Customer requests
  - And other interrupt driven work
- Shift handoffs
- Writing postmortems

# Handling alerts



How to improve incident response:

- Make pages more actionable
- Deliver pages to the right teams
- Reduce the number of pages

# More context means faster root cause analysis

# Dynamic alert delivery

Send alerts directly
to the teams that
are responsible for
taking action!



**WALKING THROUGH A TRACE**

1. Starting at the span which was defined as the signal - **place_order**

2. Inspect every child span's tags

3. Follow path with **error=true**

4. Rinse and repeat until no more children

Are We All on the Same Page? Let's Fix That
Luis Mineiro @voidmaze SRE @ Zalando, SREcon EMEA 2019
https://www.usenix.org/conference/srecon19emea/presentation/mineiro

# Improving postmortems

"How will we do better next time?"

- Make sure root causes are fixed
- Improve responses for novel issues

Make sure postmortems are blameless

- By establishing what happened *objectively* using traces

# Why is improving oncall important?

Direct impact on customer experience (revenue, reputation, etc.)

Time spent handing pages, writing postmortems, handling interrupts is... time **not** spent building new features, proactive optimization

Stress of oncall has major impact on job satisfaction

# SLOs

Service Level Objectives

- Promises service owners make to their customers
  - Both external and internal customers
- Stated in a way that can be **measured** on short time scales
  - Days, hours, or even minutes

Common indicators

- Latency (p50, p99, etc.)
- Error rate
- Availability

threshold

**p99 latency < 5s** over the last **5 minutes**

service level indicator (SLI)

measurement window

"instantaneous" latency



| — p50 **2.9s** | — p95 **3.19s** | — p99 **3.22s** | — p99.9 **3.23s** |

latency over 5 minute window

# Determining SLOs



Ask:

- What do your customers expect?
- What can you provide today?
- How do you expect that to change?

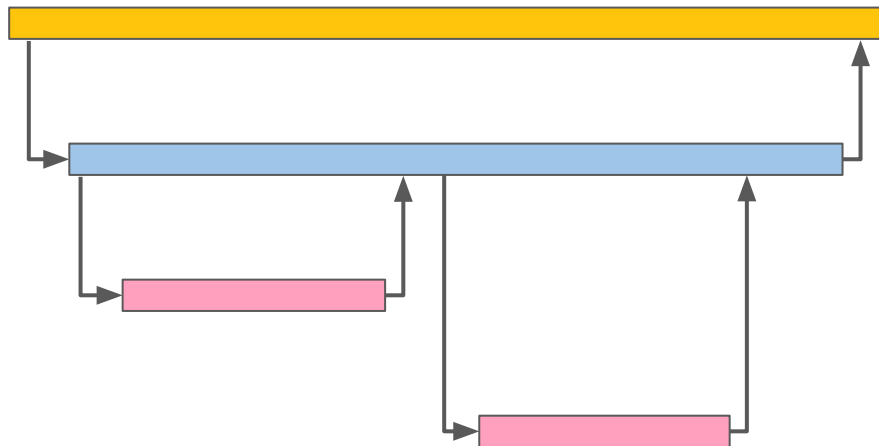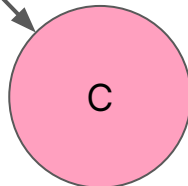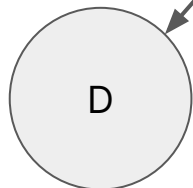# Derive internal SLOs using tracing

A: `p99 latency < 5s`

B: `p99 latency < 5s`

~p99.5 latency < 2.5s

C: `p99 latency < 2.5s`

# Why are SLOs important?

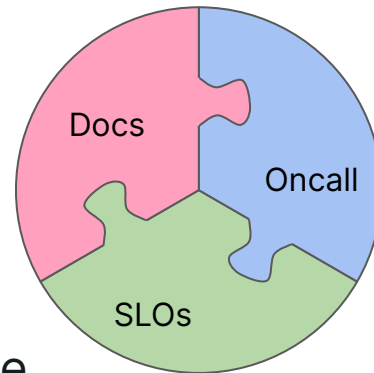They measure success in delivering service

Teams use them as a guide to prioritize work

Consistency across your org

- Hold teams accountable in a consistent way
- Do it transparently

# 3-piece puzzle review

Documentation

- Establishes ownership: who you will hold accountable
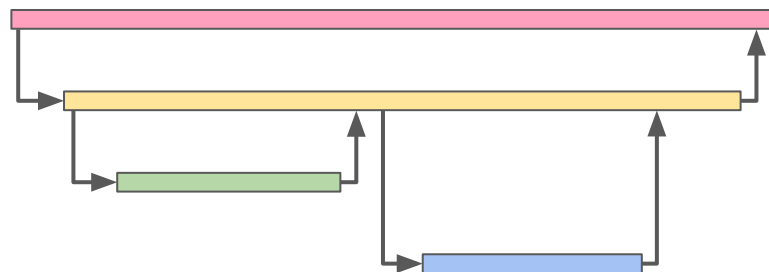- Also provides critical information for building confidence

Oncall

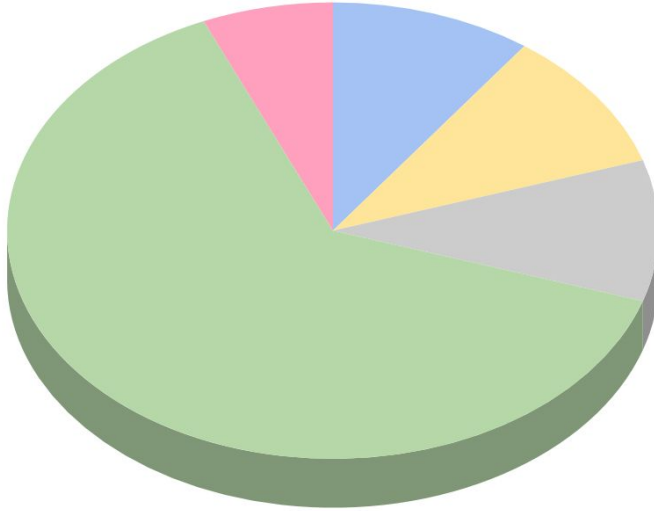- Oncall (and service ownership) are more than incident response

SLOs

- How you measure success!

Tracing provides key information for each

# Budgeting for ownership

Giving teams agency to improve reliability

- Will help them hit their goals
- Lower their stress

But having agency requires

- The right information
- The time to act on it

Ownership doesn't come for free!

# Next steps...

# Making changes

Rolling out new processes/tools in a DevOps org is hard

- Process/tools must provide **value to app dev teams**
- Ideally, they are **necessary** parts of their day-to-day work

To establish and maintain service ownership

- Use a combination of docs, oncall process, and SLOs
- *Manufacture* a need for those process/tools where necessary

# Ownership = Accountability + Agency

Accountability

- Set deliverables and goals for service owners
- Judge their performance based on those deliverables and goals

Agency

- Offer the information, confidence, and budget to improve

# Thank you

🐦 @save_spoons

Daniel "Spoons" Spoonhower, CTO and Co-founder

lightstep.com

Lightstep