# Shifting Left on Production Excellence with Observability

DevOps Enterprise Summit
October 5-7, 2021

featuring illustrations
by @emilywithcurls!

# Your business wants to move faster

# Velocity requires tight feedback loops

@lizthegrey + @shelbyspees

# We need to shift production left.

@lizthegrey + @shelbyspees

# Liz Fong-Jones

Principal Developer Advocate at Honeycomb.io

lizthegrey.com

@lizthegrey



# Shelby Spees

Site Reliability Engineer at Equinix Metal

spees.dev

@shelbyspees

# We've come a long way

@lizthegrey + @shelbyspees

# What's the next step?



* BEFUDDLED *

# Should developers be on-call?

@lizthegrey + @shelbyspees

# Production is increasingly complex

# Production feels intimidating

# Traditional tools are inscrutable

# Prod is always encountering new issues

**emergent failure modes**

    small, unrelated failures cascading together to degrade or take down a system

see also: **how.complexsystems.fail**

# 42%

of developer time is spent dealing with bad code and tech debt

# Teams can't make forward progress

**@lizthegrey + @shelbyspees**

# Our heroes are exhausted

# We shouldn't stop here

# All teams need production excellence

@lizthegrey + @shelbyspees

# Can't we just buy it?

# Invest in people, culture, and process

# Production Excellence is Business Excellence
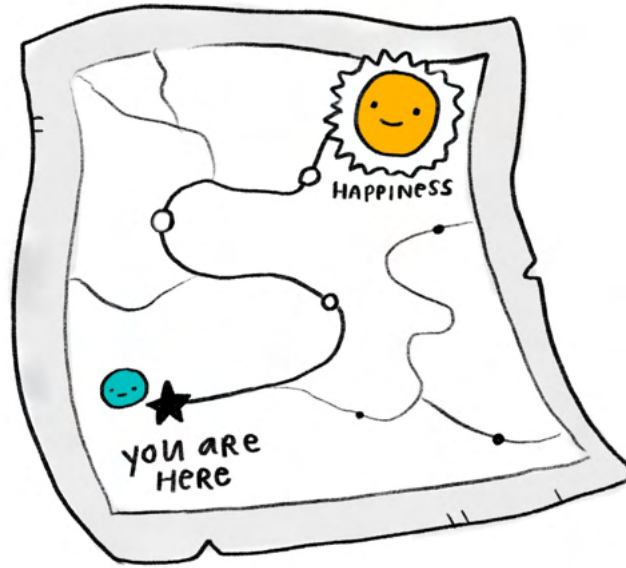
# Start with observability

# What is observability?

The ability to inspect and understand a system's internal state using the telemetry data it's already outputting.

(Even if your system is in a state you didn't know was possible!)

# Hard-to-debug problems

**Distributed Systems**

small change causing downstream effects?

**Poor Performance**

what is worth optimizing?

**Subset of Traffic**

only some users are complaining?

**Regulatory Requirements**

can't touch the box because of SOX?

# Observability finds answers.

**Distributed Systems**

Not a problem, we have traces.

**Poor Performance**

Identify the bottleneck.

**Subset of Traffic**

Slice and dice with ad-hoc queries.

**Regulatory Requirements**

Read-only access.

**@lizthegrey + @shelbyspees**
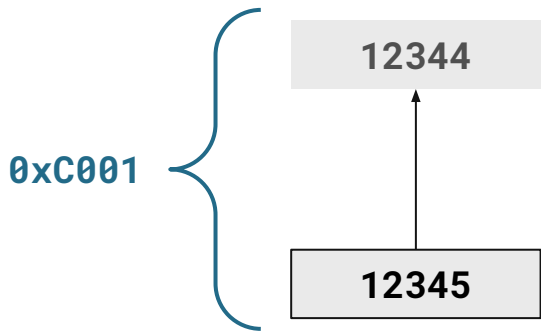
# Example event

```
event = {
    "duration_ms" = 1198,
    "request.id" = 12345,
    "request.method" = "GET",
    "request.path" = "/search",
    "request.query_string" = "category=decor&price<=50",
    "response.status_code" = 200
}
```

@lizthegrey + @shelbyspees

# Example event

```
event = {
    "trace_id" = 0xC001,
    "span_id" = 12345,
    "parent_id" = 12344,
    "duration_ms" = 1198,
    "request.method" = "GET",
    "request.path" = "/search",
    "request.query_string" = "category=decor&price<=50",
    "response.status_code" = 200
}
```

12344

0xC001

12345

@lizthegrey + @shelbyspees

# Examples of dimensions to capture

infrastructure:
- build ID
- kubernetes pod
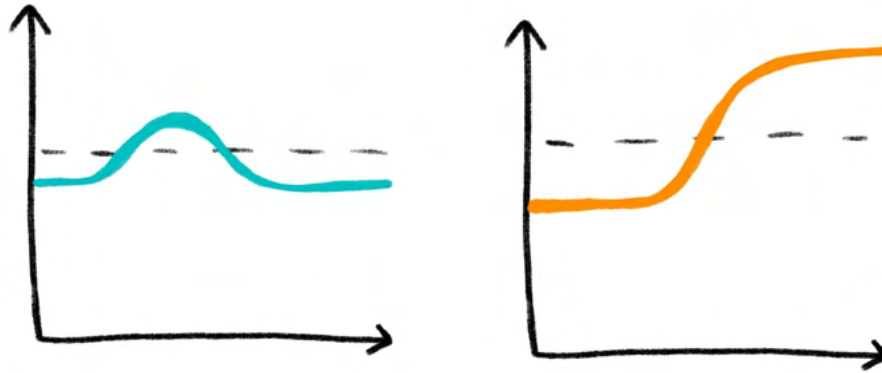- kafka broker

application:
- library version
- user-agent
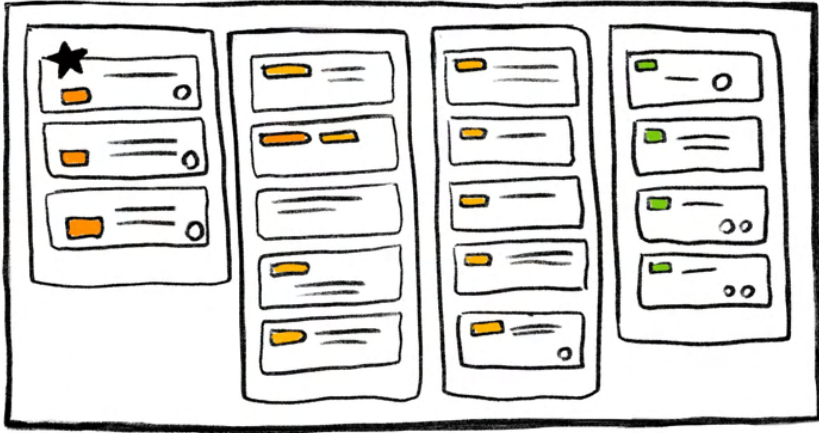- sanitized SQL query string

domain:
- number of items in cart
- coupon code
- payment processor
  (e.g. PayPal vs. Stripe)

# Technical Decisions
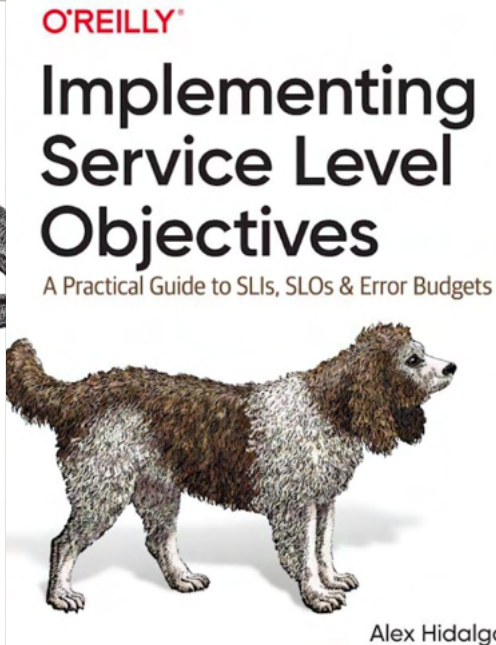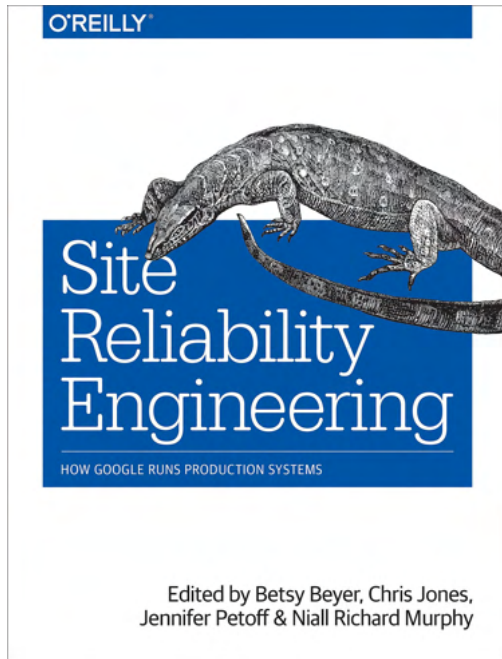# are Business Decisions

# How do we map effort to impact?

@lizthegrey + @shelbyspees

# "How does this drive the business forward?"

@lizthegrey + @shelbyspees

# Service Level Objectives (SLOs)

Common language for engineers and business stakeholders



O'REILLY®

Site Reliability Engineering

HOW GOOGLE RUNS PRODUCTION SYSTEMS

Edited by Betsy Beyer, Chris Jones,
Jennifer Petoff & Niall Richard Murphy



O'REILLY®

Implementing Service Level Objectives

A Practical Guide to SLIs, SLOs & Error Budgets

Alex Hidalgo

# Think in terms of events in context

# SLI: for each event, good or bad?

@lizthegrey + @shelbyspees

# Use a window and target percentage

window is usually 30, 60, or 90 days

# 99.9%

of app Home Page loads
over the past 30 days
were "successful" and "fast enough"

# Historical SLO compliance

# A good SLO barely keeps users happy

# Error budget: allowed unavailability

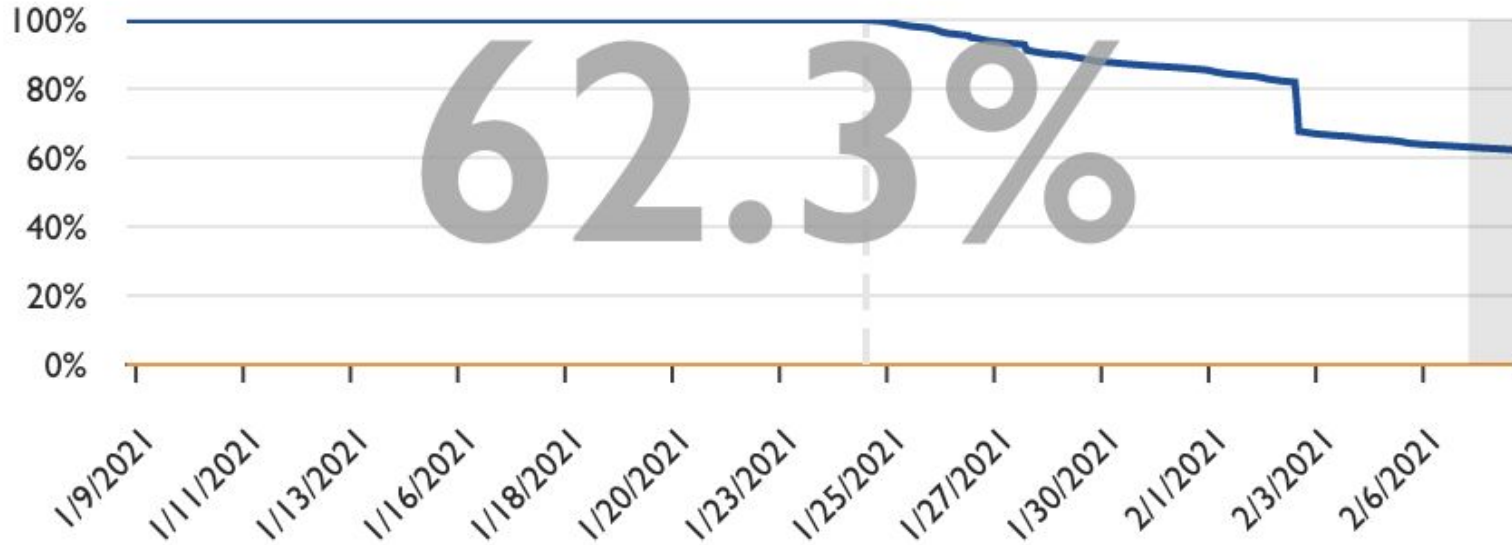$$1 - x$$

99.9% SLO Target → 0.1% Error Budget

@lizthegrey + @shelbyspees

# Example error budget

99.9% SLO Target → 0.1% Error Budget

$$
\begin{array}{r}
1{,}000{,}000 \quad \text{requests/month} \\
\times \qquad\qquad 0.1\% \\
\hline
1{,}000 \quad \text{requests/month}
\end{array}
$$

We're allowed 1000 "bad" requests/month

# Error budget remaining

# When is it okay to take risks?



RISKS

$$\frac{\text{PERCENTAGE of USERS IMPACTED} * \text{DURATION of OUTAGE}}{\text{TIME BETWEEN FAILURES}} = \text{TOTAL RISK}$$

@lizthegrey + @shelbyspees
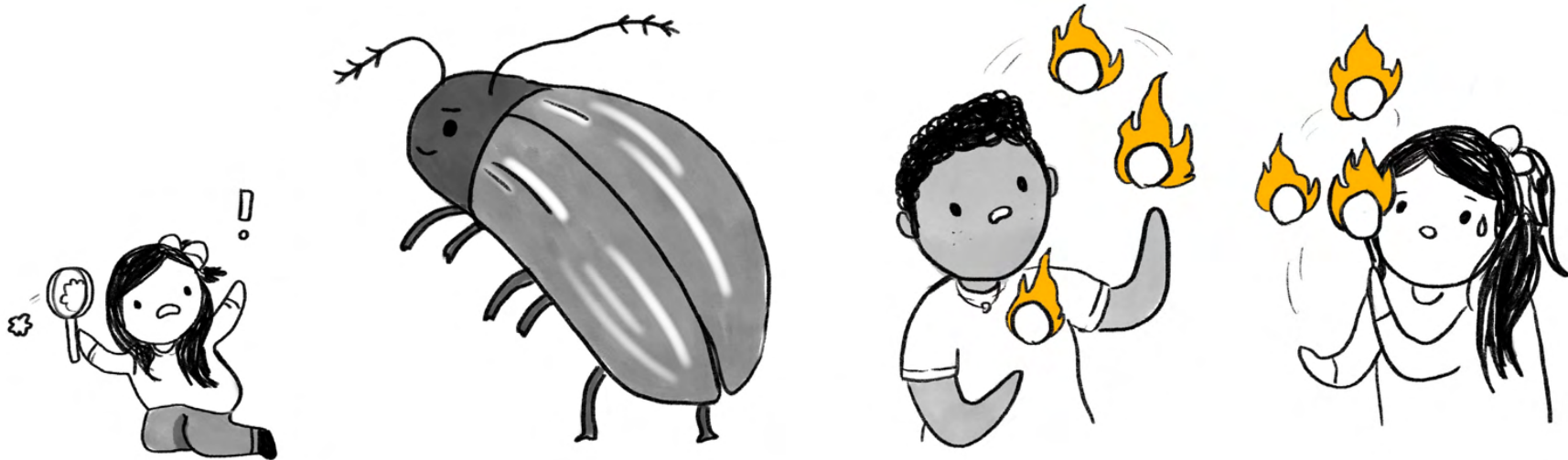
# When is it *not* okay?

# Alerting on SLOs
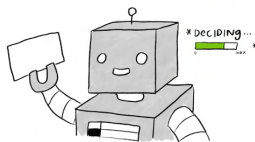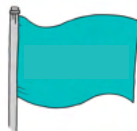
Alert on error budget **burndown rate**

   "How soon are we going to run out of budget?"

NOT: every **potential cause** of an alert
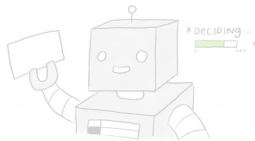
   "Disk is at 90.05%!!!" (does it matter?)

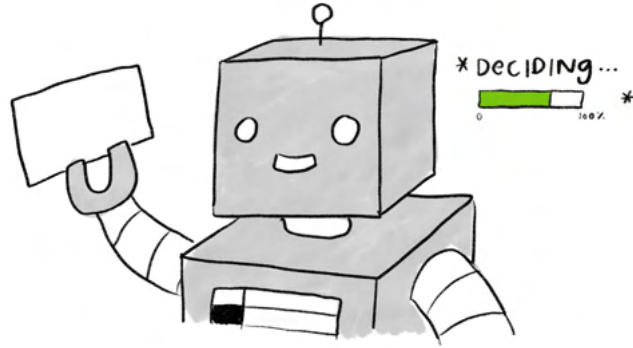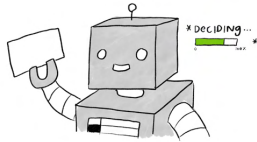# Observability & Continuous Delivery

@lizthegrey + @shelbyspees

# Continuous Delivery (at Honeycomb)

**@lizthegrey + @shelbyspees**

# Instrument as we code.

# Ensure fast integration.

# Observe behavior in prod.

# Bring production excellence to your org

# Invest in the right places

# Instrument your builds



@lizthegrey + @shelbyspees

# Find a champion

# Learn and iterate

**SLOs aren't one-and-done**

- Adjusting SLO targets
- Improving instrumentation
- Updating alerting methods

see also: **go.hny.co/biz-goals-slos**

# Make it scale to your org

Prioritize:
- developer experience
- on-call onboarding

Invest in:
- config as code
- custom instrumentation libraries

@lizthegrey + @shelbyspees

# Case Study: Observability at Vanguard

Observability champion(s) ✔️

Prioritizing knowledge transfer ✔️

OpenTelemetry ✔️

Service Level Objectives ✔️

# 18 million developers

It's on us to grow them into production engineers.

# Help we're looking for

CNCF Observability
Technical Advisory Group (TAG)
- CNCF Slack
  **#tag-observability** channel

OpenTelemetry
- CNCF Slack **#opentelemetry** channel
- OpenTelemetry.io

OpenSLO
- OpenSLO Slack
- OpenSLO.com

Schedule 30 minutes
- hny.co/meet/liz

Find us on Twitter
- twitter.com/lizthegrey
- twitter.com/shelbyspees