

izettle



public class SRE implements DevOps



Who am I?

Who am I?

- Have come from more of an Ops background
- Am really interested in how and why systems breaks
- Am interested in how people respond to incidents



Why am I giving this
talk?

The background is a dark blue gradient. Scattered across the slide are several small, solid-colored circles in white, light blue, orange, green, and purple, creating a starry or abstract pattern.

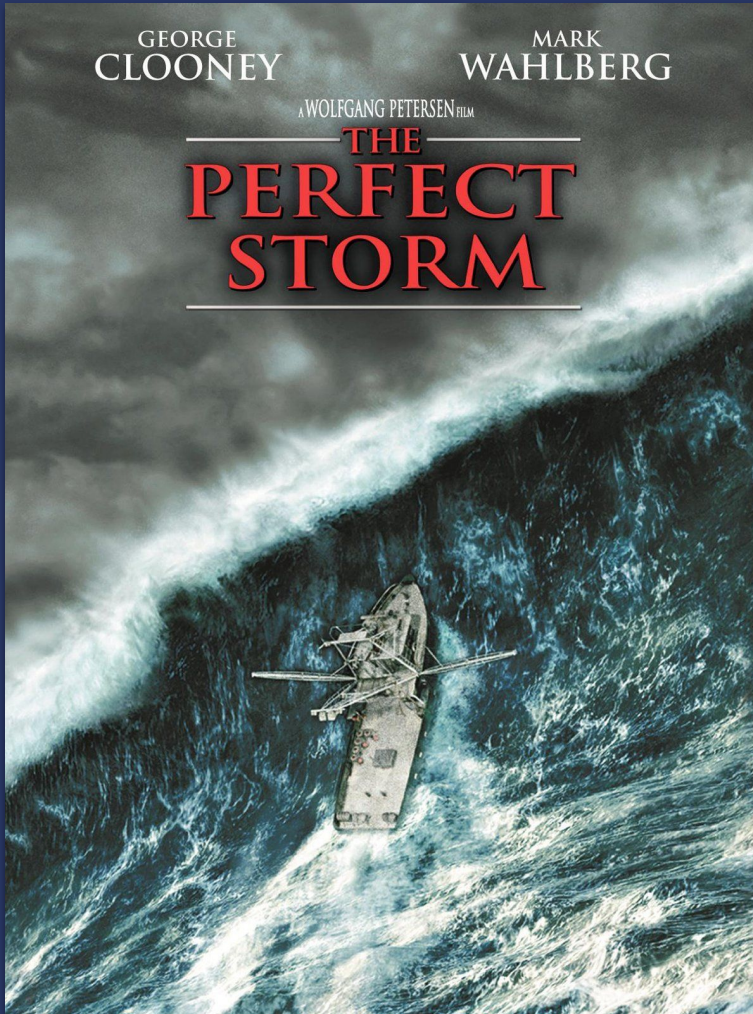
**Over the next 12 months
we are going to have a lot
of incidents**

GEORGE
CLOONEY

MARK
WAHLBERG

A WOLFGANG PETERSEN FILM

THE PERFECT STORM





Roughly 70%¹ of outages are due to changes

*¹ Analysis of Google internal data 2011 - 2018
Solving reliability fears with SRE Cloud Next '18*

The background is a solid dark blue. Scattered across the slide are several small, semi-transparent colored dots in shades of white, light blue, orange, green, and purple.

How we manage those
incidents will define our
success

**SRE provides a clean
outline of how to manage
incidents and failure**

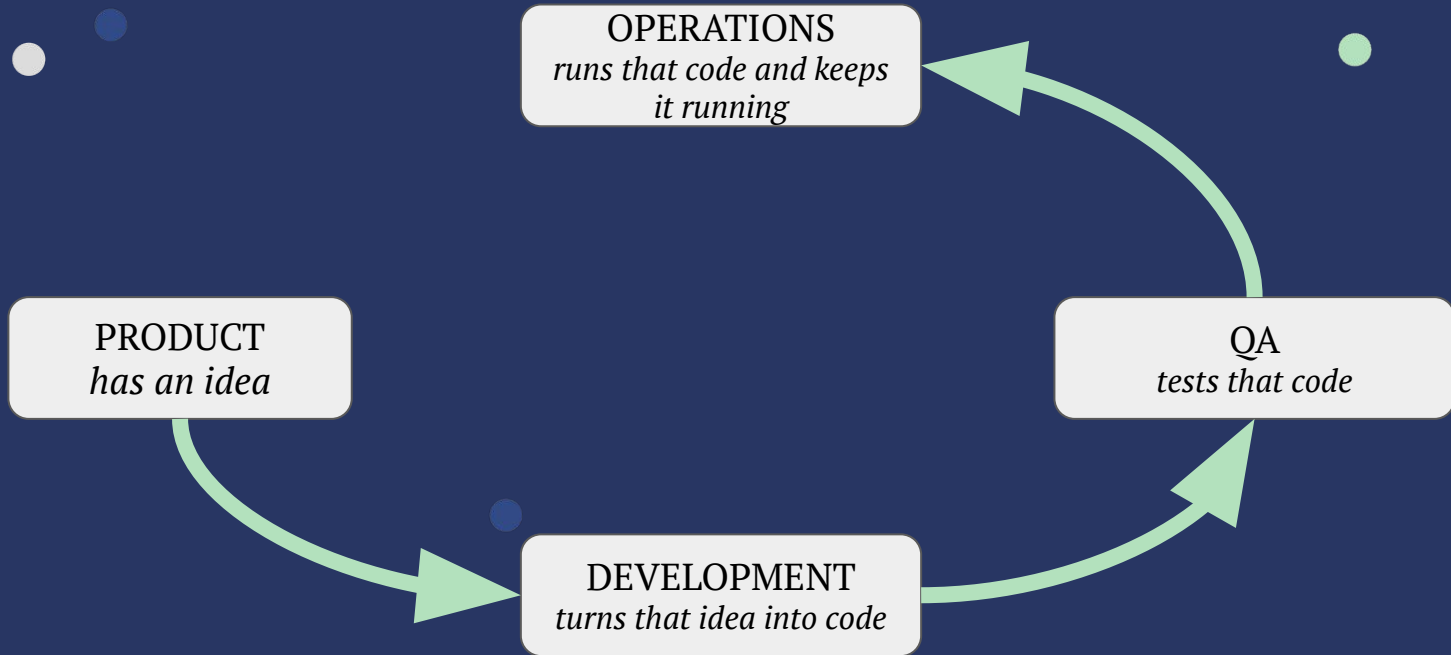
What is DevOps?



Software Service Operations is hard



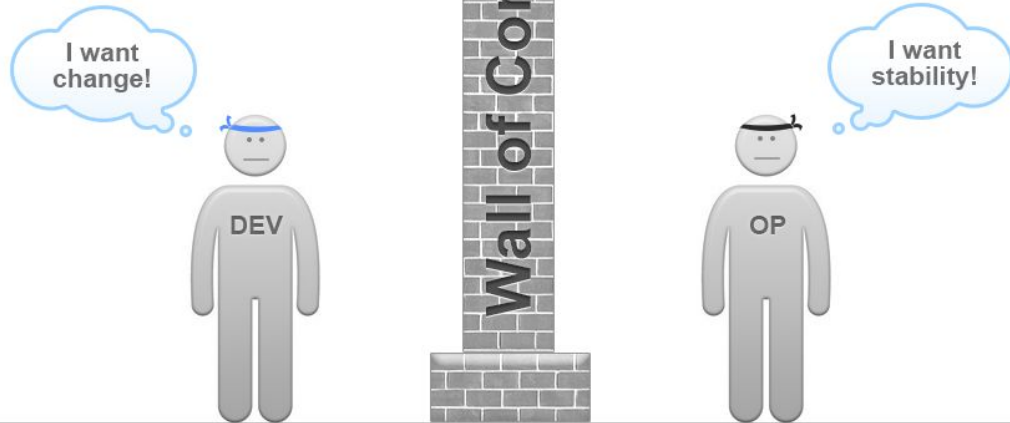
**Traditionally Operations
was treated as a Cost
Centre**

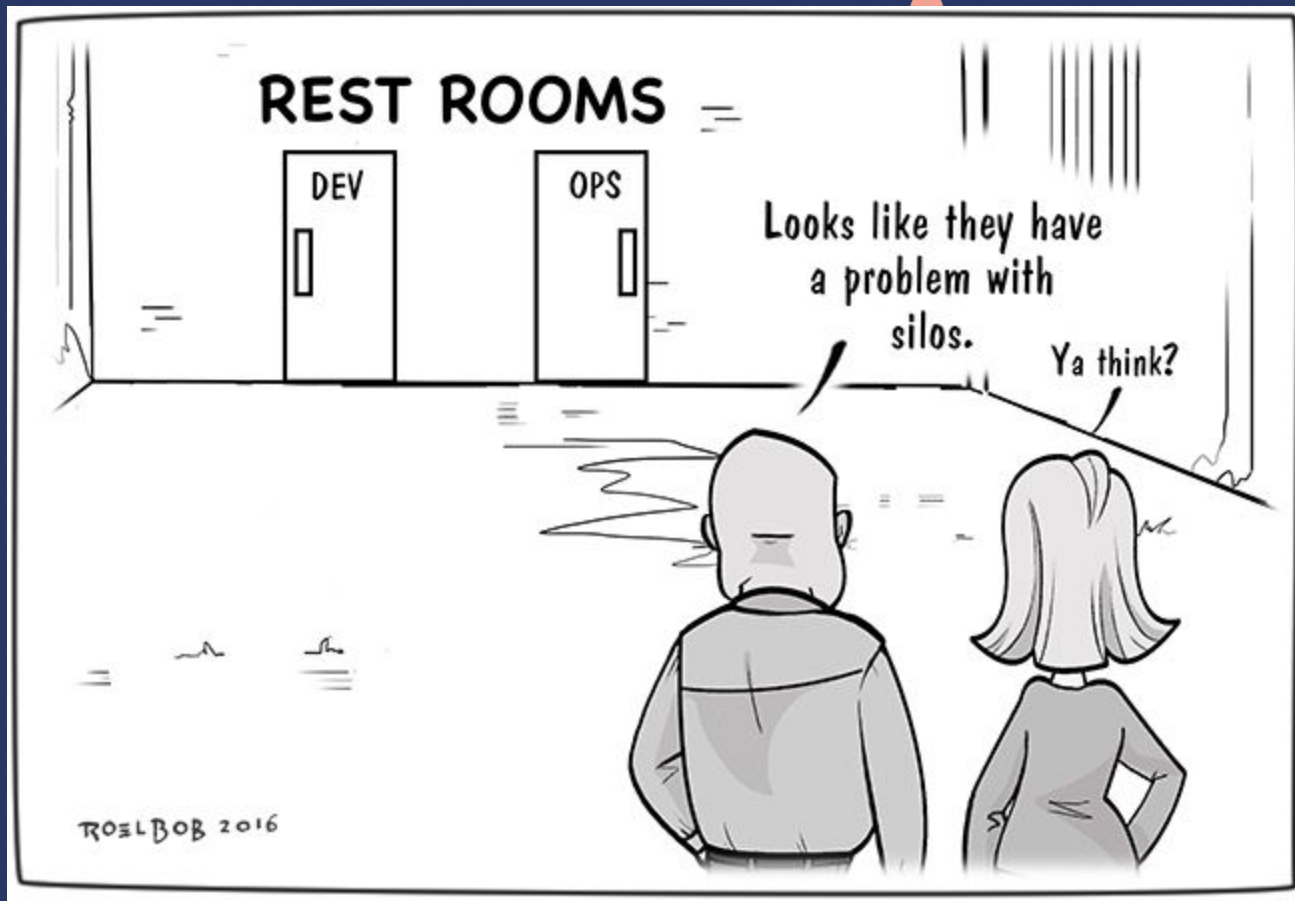


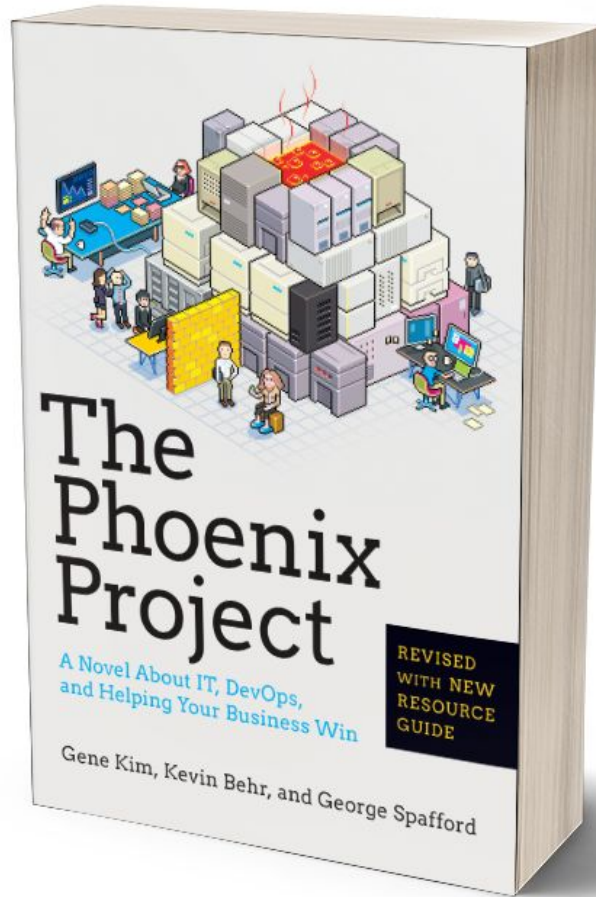
WORKED FINE IN DEV

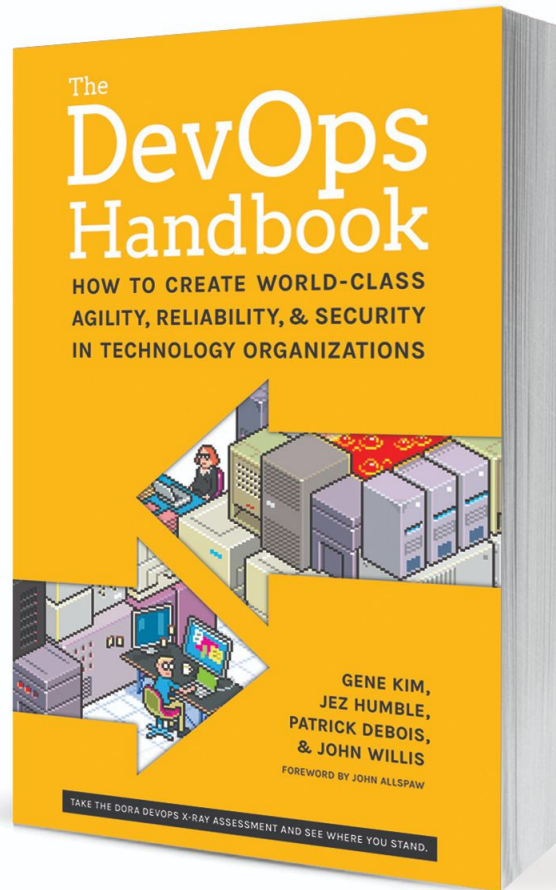
OPS PROBLEM NOW

makeameme.org









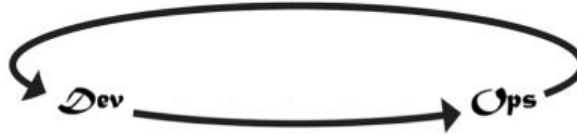
The First Way: Systems Thinking



The First Way emphasizes the performance of the entire system, as opposed to the performance of a specific silo of work or department..

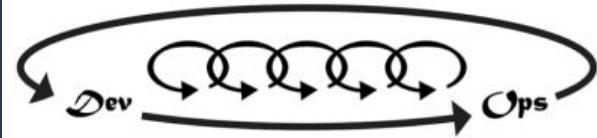
The focus is on all business value streams that are enabled by IT.

The Second Way: Amplify Feedback Loops



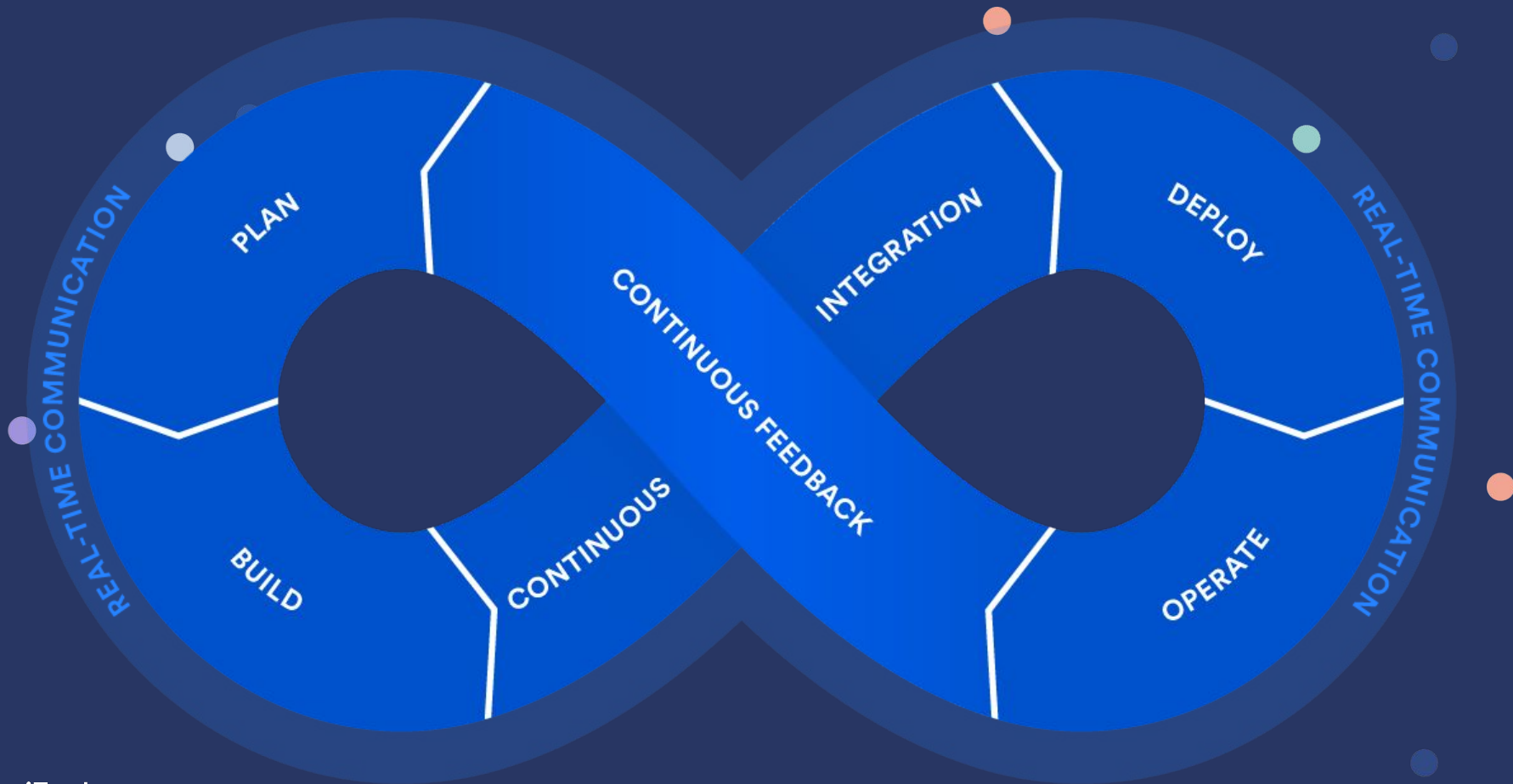
The Second Way is about creating the right to left feedback loops. The goal of almost any process improvement initiative is to shorten and amplify feedback loops so necessary corrections can be continually made.

The Third Way: Culture Of Continual Experimentation And Learning



The Third Way is about creating a culture that fosters two things: continual experimentation, taking risks and learning from failure; and understanding that repetition and practice is the prerequisite to mastery.

We need both of these equally. Experimentation and taking risks are what ensures that we keep pushing to improve, even if it means going deeper into the danger zone than we've ever gone. And we need mastery of the skills that can help us retreat out of the danger zone when we've gone too far.



**ACCEPT FAILURE
AS NORMAL**

**IMPLEMENT
GRADUAL
CHANGE**

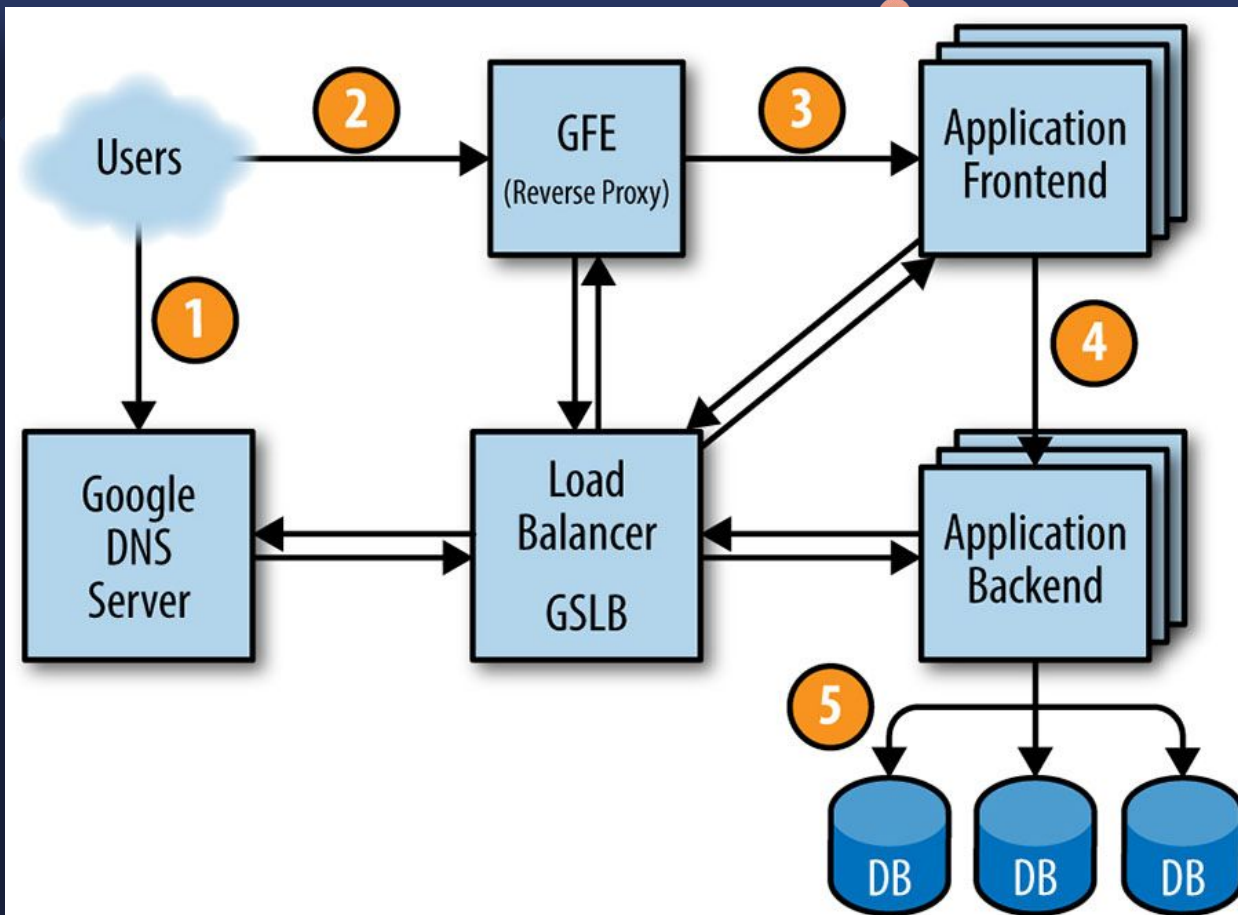
**REDUCE
ORGANISATIONAL
SILOS**

**LEVERAGE
AUTOMATION**

**MEASURE
EVERYTHING**

What about SRE?

History



SRE Principles

- Use data to guide decisions
- Embracing Risk through Error Budgets
- Service Level Objectives that have consequences
- Eliminating Toil so we have time to make tomorrow better

Use Data to Guide Decisions



Embracing Risk

The background is a solid dark blue color. Scattered across the slide are approximately 15 small, semi-transparent dots in various colors including white, light blue, orange, green, and purple. These dots are positioned at various angles and distances from the center, creating a starry or abstract pattern.

100% uptime is not feasible



Most laptops are ~70%
reliable

The difference between
99.9% and 99.999% is
indistinguishable for
most users

The background is a solid dark blue color. It is decorated with several small, semi-transparent colored dots in white, light blue, orange, green, and purple, scattered across the slide.

Reliability and Innovation and competing goals

Error Budget

Error Budget

- Given an SLO of 99.9% what downtime can we have?

$$\frac{100 - 99.9}{100} \times 365 \text{ days / year} \times 24 \text{ hours / day} \times 60 \text{ mins / hour} = 525 \text{ mins / year}$$

| Availability Level | Allowed Unavailability Window | | |
|--------------------|-------------------------------|---------------|--------------|
| | per year | per quarter | per 30 days |
| 90% | 36.5 days | 9 days | 3 days |
| 95% | 18.25 days | 4.5 days | 1.5 days |
| 99% | 3.65 days | 21.6 hours | 7.2 hours |
| 99.5% | 1.83 days | 10.6 hours | 3.6 hours |
| 99.9% | 8.76 hours | 2.16 hours | 43.2 minutes |
| 99.95% | 4.38 hours | 1.08 hours | 21.6 minutes |
| 99.99% | 52.6 minutes | 12.96 minutes | 4.32 minutes |
| 99.999% | 5.26 minutes | 1.30 minutes | 25.9 seconds |

Error Budget

- Given an SLO of 99.9%, after a 20 minute outage you still have 23 minutes of budget available for the remainder of the month.

Exceeding Your Error Budget

- Deployment of new features not allowed
- Focus on reliability/performance improvements



Service Level Objectives
Service Level Indicators
Service Level Agreements

Service Level Indicator

- Something you can measure
 - Request latency
 - Error rates
- *95th Percentile latency of homepage requests over past 5 minutes less than 300ms*

Service Level Objective

- Informed by SLIs
- Target level or range that states if a service is normal
- Specifically tracking customer experience
- If SLO is being met then customers are happy
- *95th percentile homepage latency SLI will succeed 99.9% of the time within a month*

Service Level Agreement

- Contractual obligation
- If breached then need to pay someone £££
- *Users will receive service credits if 95th percentile homepage latency SLI succeeds less than 99.9% over the trailing month*



● **SLIs drive SLOs which inform SLAs**

Error Budgets and SLOs

- Common incentive for Dev and Ops
- Dev teams decide how they want to spend their budget
- Unrealistic reliability goals become unattractive
- Builds shared responsibility between teams

Eliminating Toil

Toil

- SRE puts a cap of 50% on the aggregate “ops” work.
 - Being on-call
 - Tickets
 - Manual tasks

Toil

- Toil is work that is:
 - **Manual**
 - **Repetitive**
 - **Automatable**
 - Devoid of **long-term value**

Toil

- What toil is not but is overhead
 - **Meetings**
 - **Emails**
 - **Travelling**
 - **Expenses**



Pursue Maximum Change Velocity Without Violating a Service's SLO

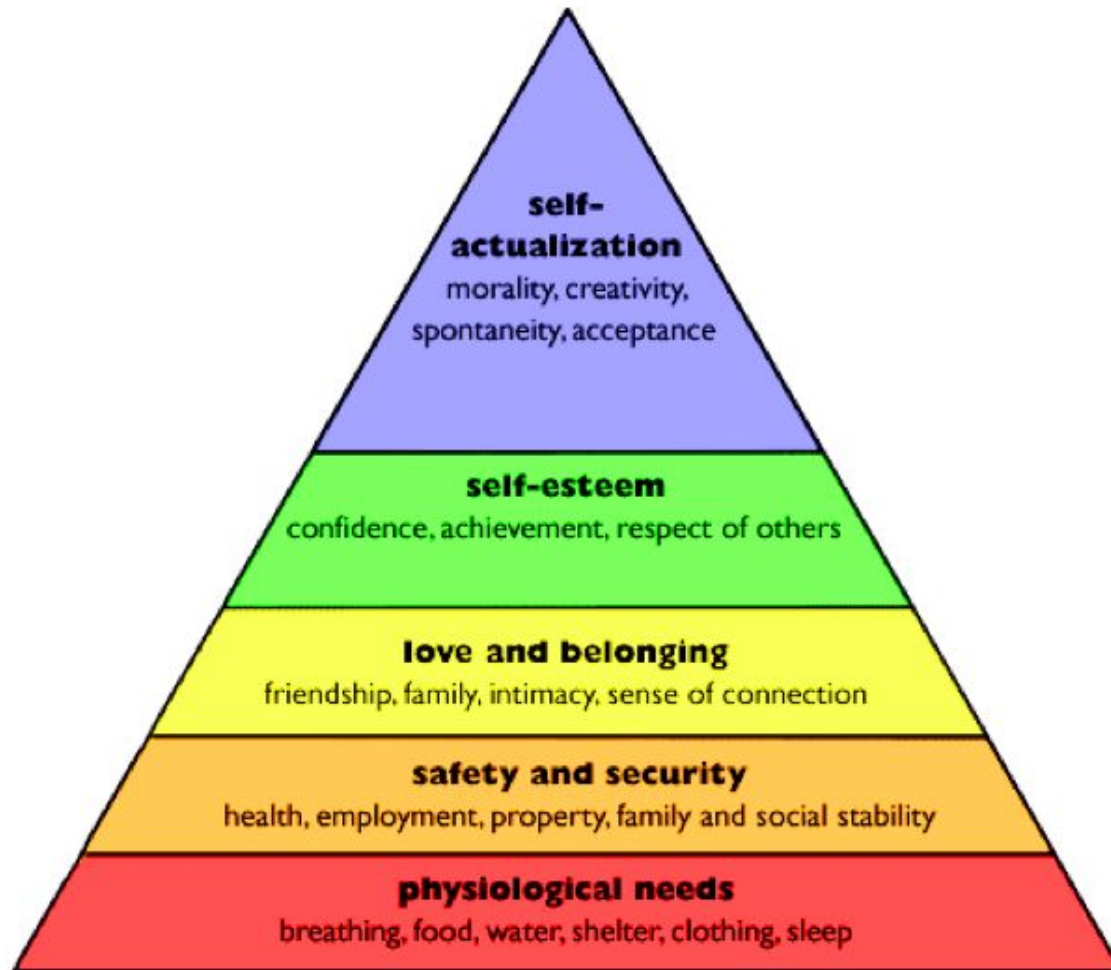


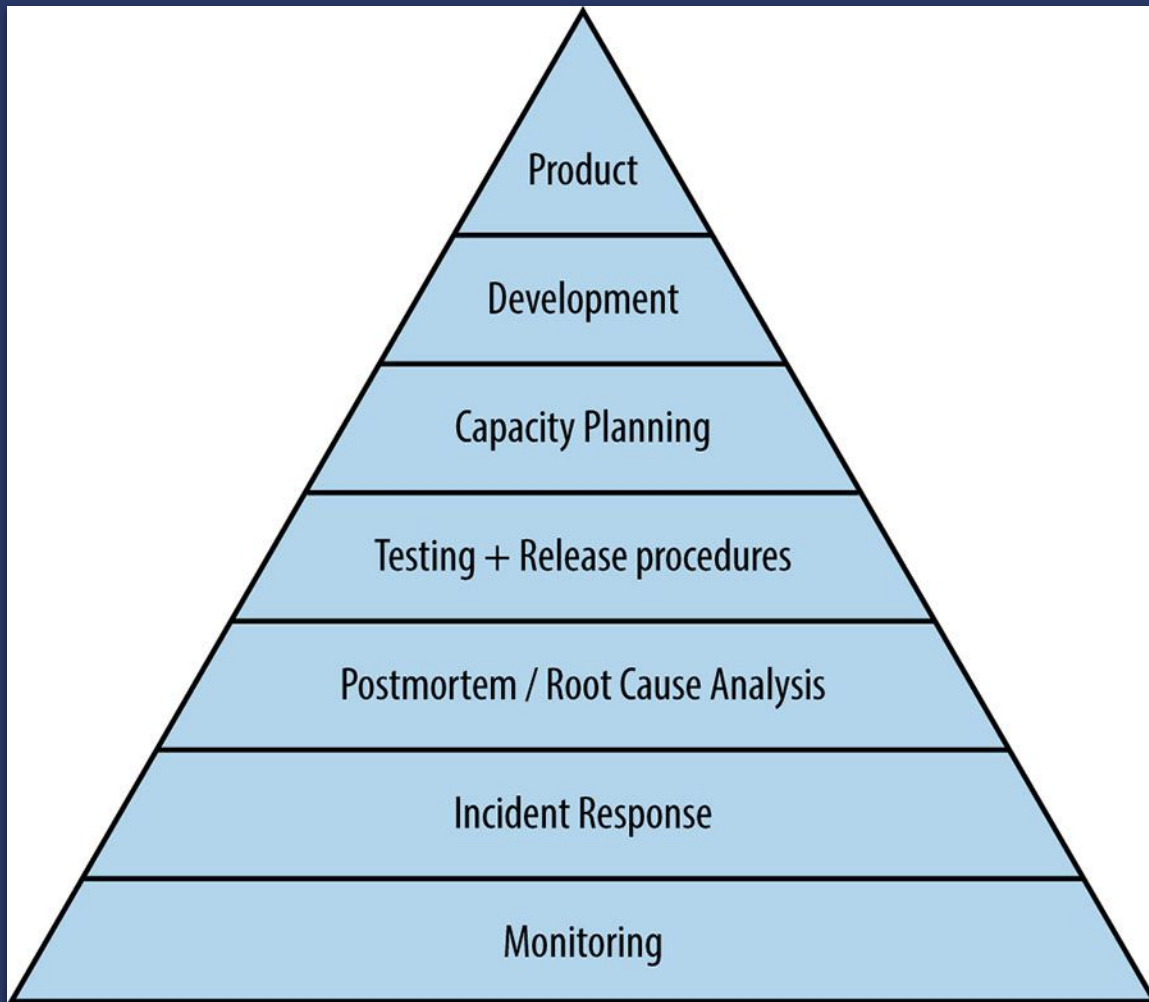
The background is a solid dark blue color. It is decorated with several small, semi-transparent colored dots in white, light blue, orange, green, and purple, scattered across the slide.

Reducing Risky Behaviour with Canary Releases



The Hierarchy of Reliability







Postmortem / Root Cause Analysis

Incident Response

Monitoring



Monitoring is necessary to:

- **Alert** on conditions that require attention
- **Investigate** and **diagnose** those issues
- Gain **insight** into **trends** in resource usage or service health for **long-term planning**
- **Compare** the behaviour of the system before and after a change, or between groups in an **experiment**

Data Sources

1. Logs

Append-only record of events

2. Metrics

Numerical measurements representing attributes and events, gathered at regular intervals

3. Distributed Traces

Individual request level observability

Data Sources

1. Logs

More granular, high cardinality but delayed delivery

2. Metrics

Less granular, low cardinality, near real time and service level

3. Traces

Very granular, high cardinality and focussed on individual requests

The background is a solid dark blue. Scattered across the slide are several small, semi-transparent colored dots in shades of orange, light blue, green, and purple. The text is centered and reads:

Monitoring is there to
give you the context you
need to fix failing
services

The background is a solid dark blue color. Scattered across the background are several small, semi-transparent dots in various colors: white, light blue, orange, light green, and purple. These dots are positioned at various angles and distances from the center, creating a starry or abstract pattern.

What should you observe?

The Four **Golden** Signals

1. **Latency**

How long does it take to service a request

2. **Traffic**

Measure of how much demand is being placed on your system

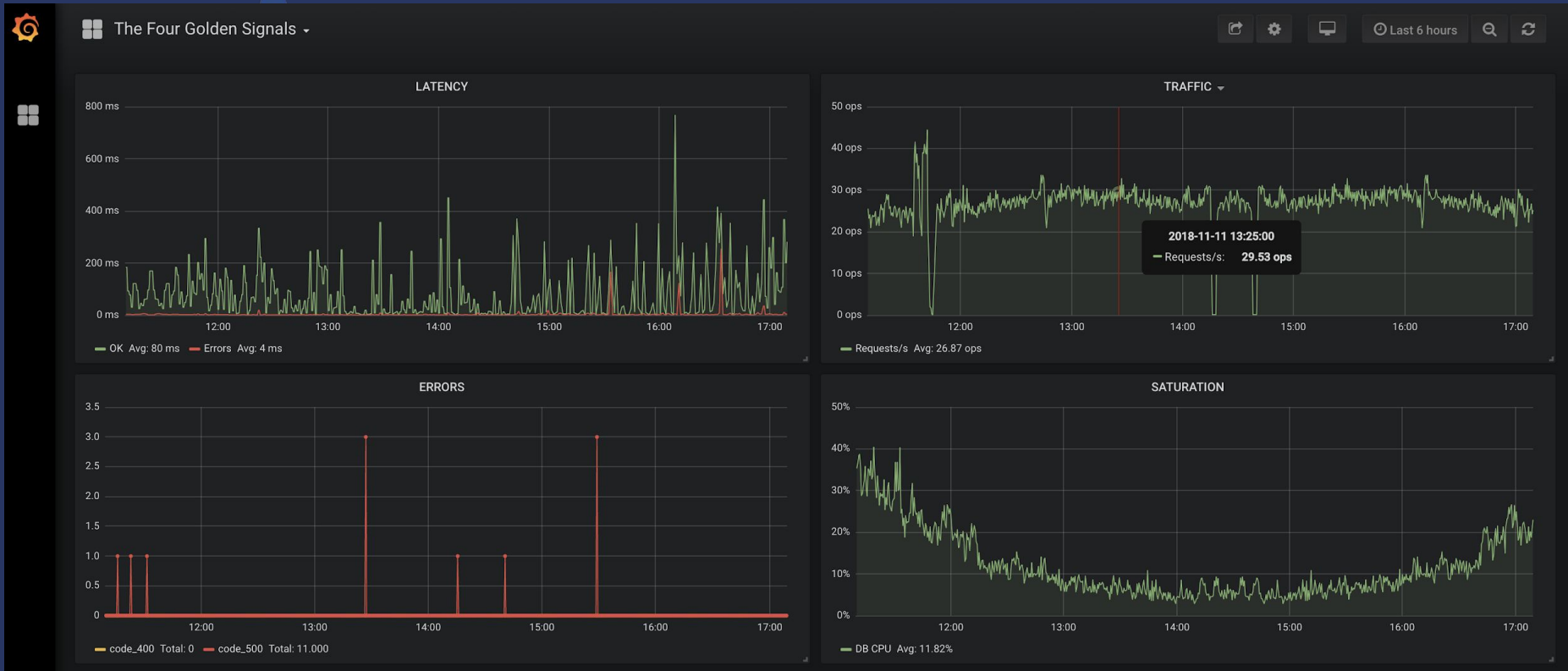
3. **Errors**

The rate of requests that fail

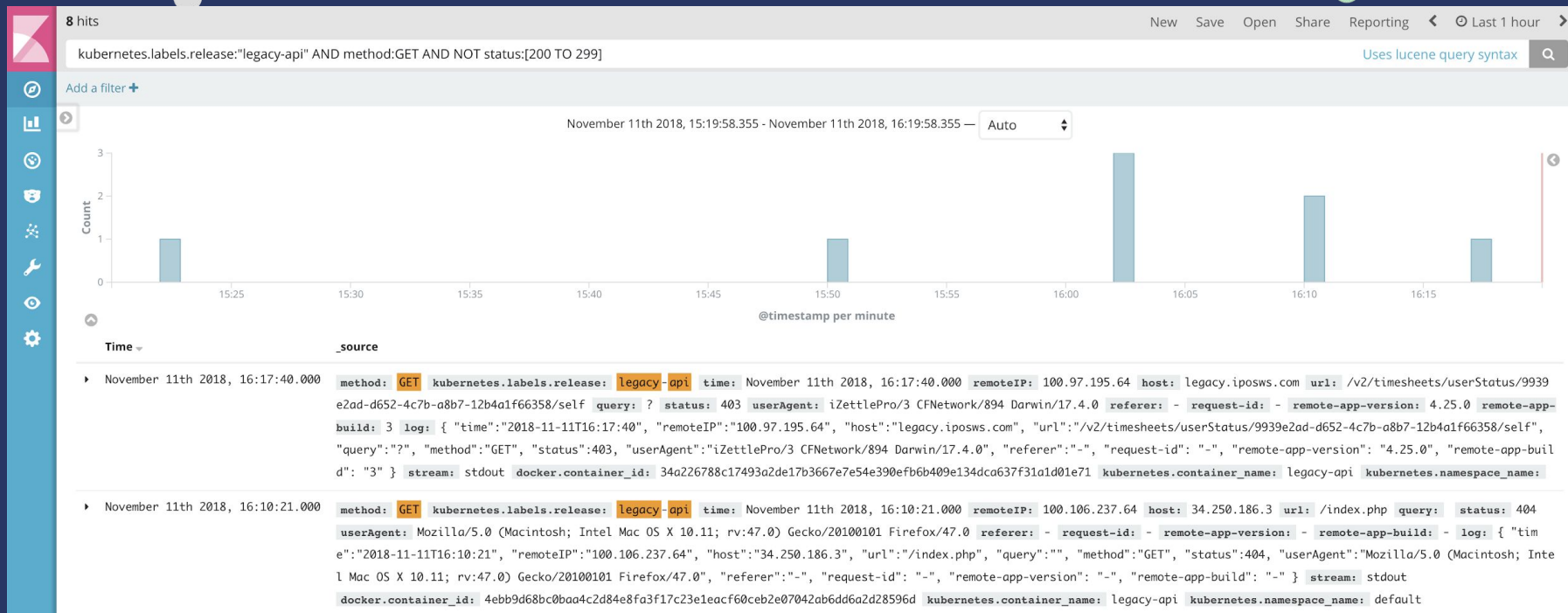
4. **Saturation**

How “full” your service is

Dashboards for Starters



Log Querying for Deep Diving



Distributed Tracing for Following Users Experience

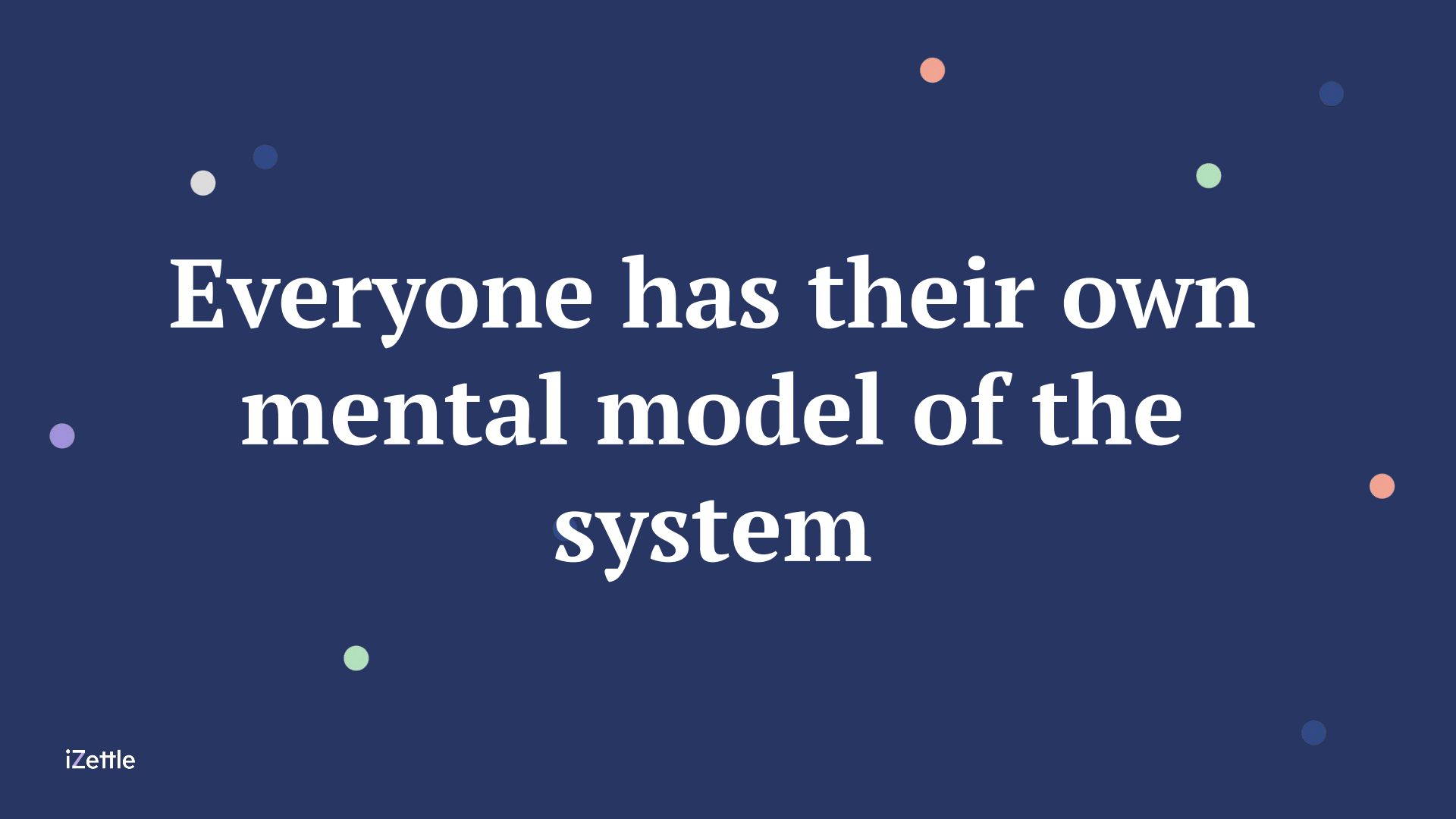
| Services | | 24.980ms | 49.959ms |
|------------------|---|-----------|--------------------------|
| - edge-server | . | 124.898ms | http:/api/accounts/10000 |
| - edge-server | . | 107.662ms | http:/api/accounts/10000 |
| - accountservice | . | 80.536ms | getaccount |
| accountservice | . | 33μ | queryaccount |
| accountservice | . | 82.461ms | getquote |
| - accountservice | . | 7μ | messaging |



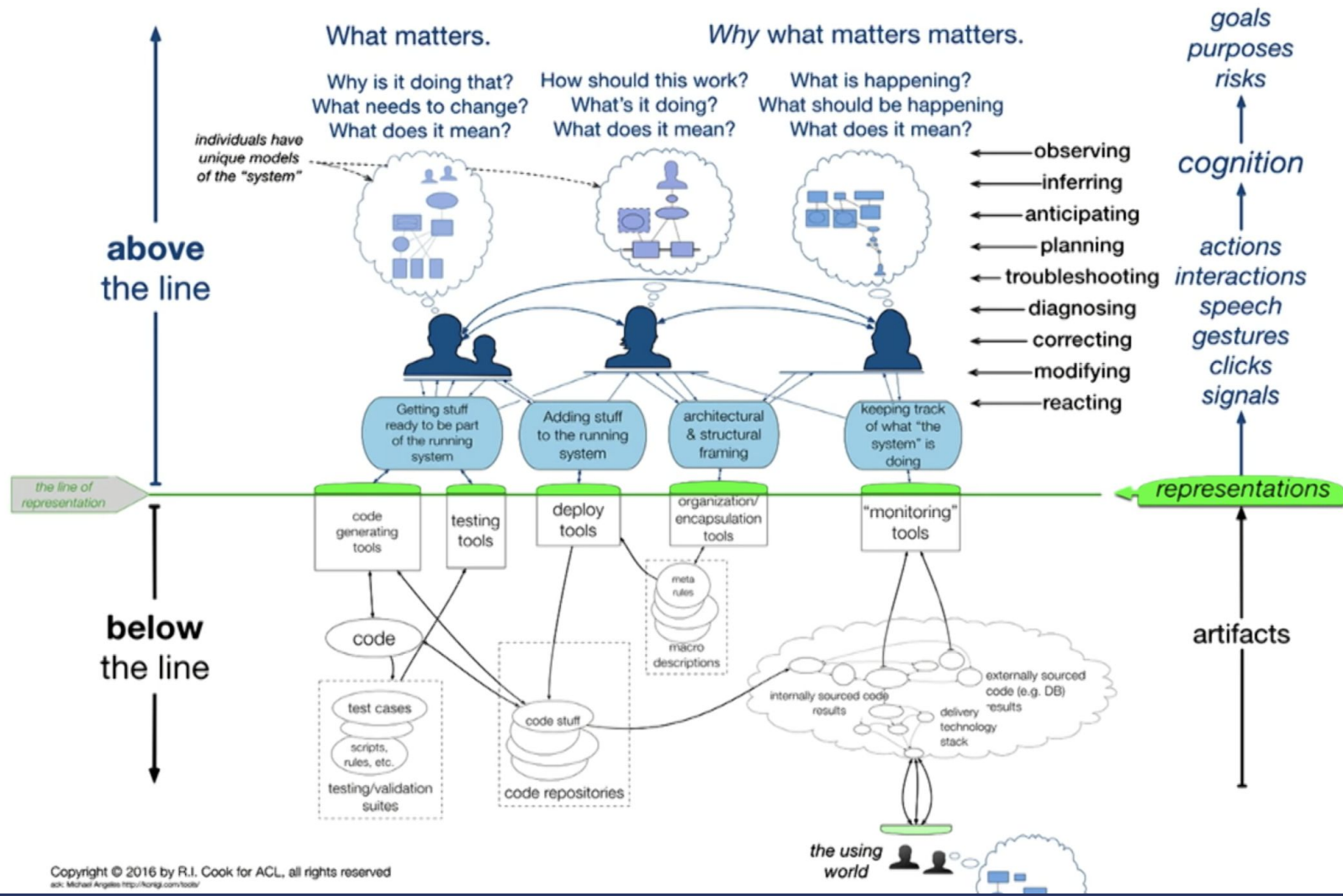
Incident Response



You Can't See Code

The background is a solid dark blue. It is decorated with several small, semi-transparent colored dots in white, light blue, orange, green, and purple, scattered across the slide.

Everyone has their own
mental model of the
system



Track Outages

What was the Time to Detection?

What was the Time to Recovery?

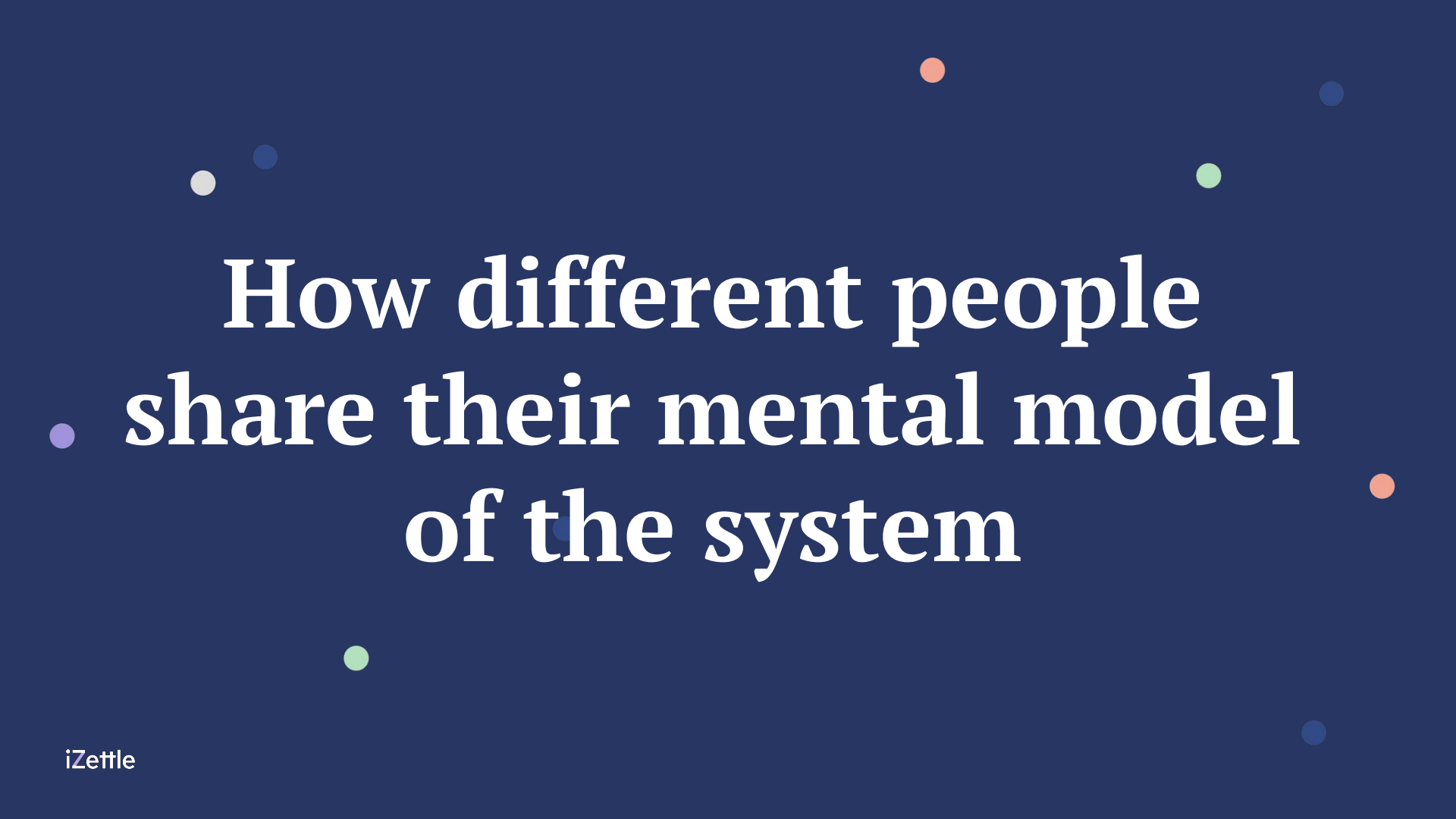
How many users were affected?

How many times in a period did outage happen?

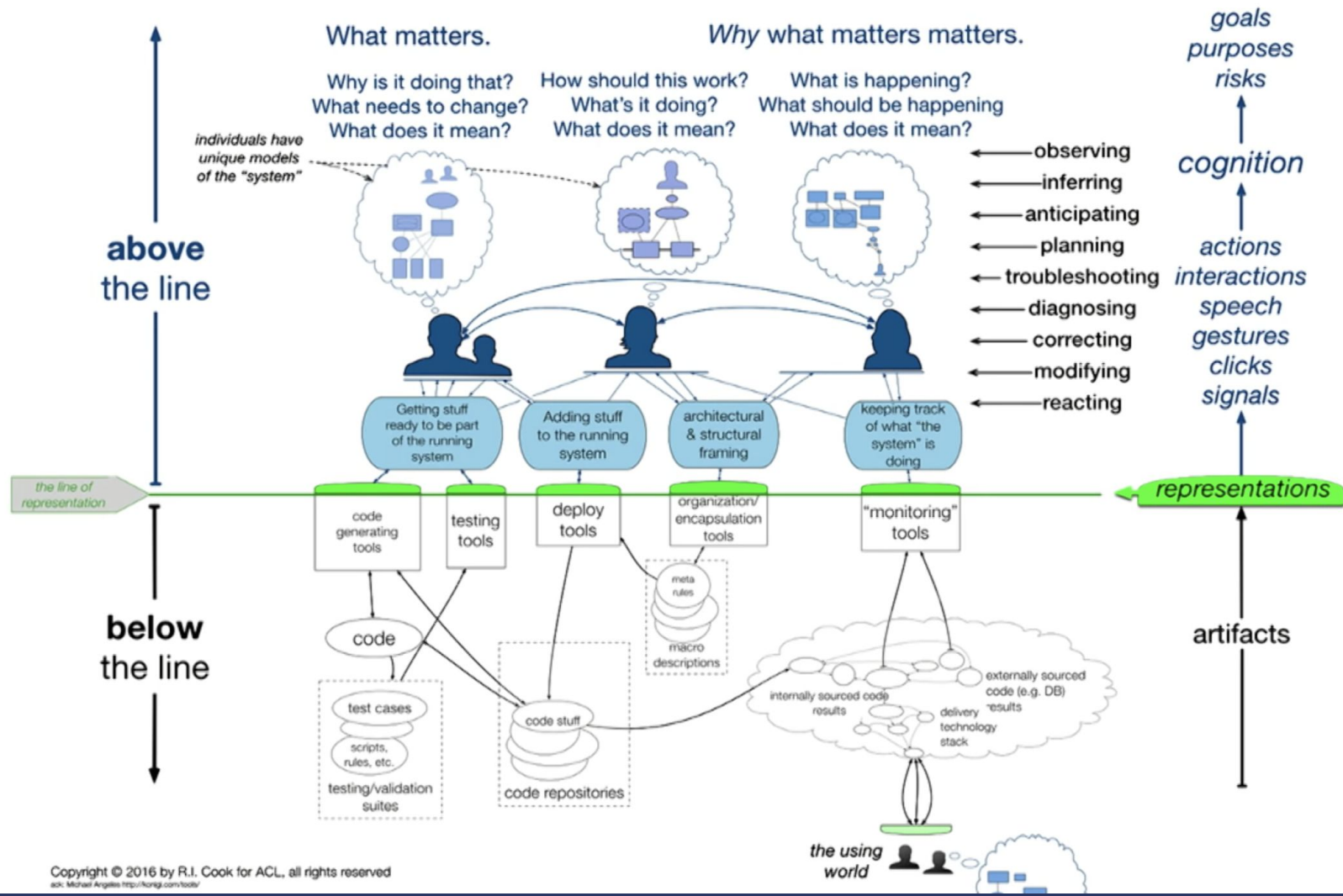
Post Incident Reviews

Post Incident Review

- Ensuring the incident is documented
- All contributing root causes are understood
- Effective preventive actions are put in place to reduce recurrence

The background is a solid dark blue. It is decorated with several small, semi-transparent colored dots in white, light blue, orange, green, and purple, scattered across the slide.

How different people
share their mental model
of the system



STELLA

Report from the SNAFUcatchers Workshop on Coping With Complexity

Brooklyn NY, March 14-16, 2017



Winter storm STELLA

Woods' Theorem: *As the complexity of a system increases, the accuracy of any single agent's own model of that system decreases rapidly.*

Blameless PostMortems and a Just Culture



Posted by **John Allspaw** on May 22, 2012

Last week, Owen Thomas wrote a flattering [article over at Business Insider](#) on how we handle errors and mistakes at Etsy. I thought I might give some detail on how that actually happens, and why.

Anyone who's worked with technology at any scale is familiar with failure. Failure cares not about the architecture designs you slave over, the code you write and review, or the alerts and metrics you meticulously pore through.

So: failure happens. This is a foregone conclusion when working with complex systems. But what about those failures that have resulted due to the actions (or lack of action, in some cases) of individuals? What do you do with those careless humans who caused everyone to have a bad day?

...participants should be able to give
a detailed account *without fear of
punishment or retribution.*

The background is a dark blue gradient. Scattered across the slide are several small, solid-colored circles in white, light blue, orange, green, and purple, creating a starry or abstract pattern.

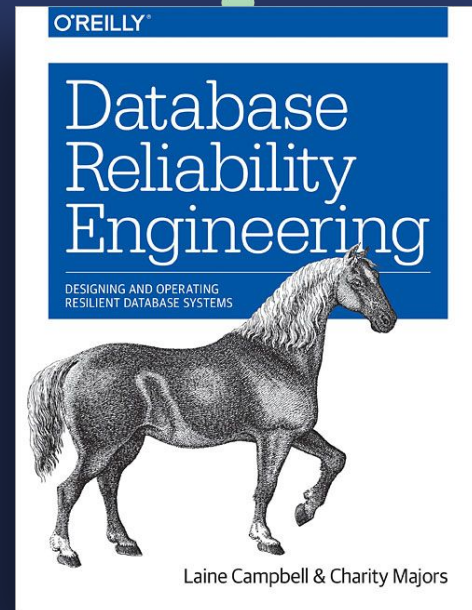
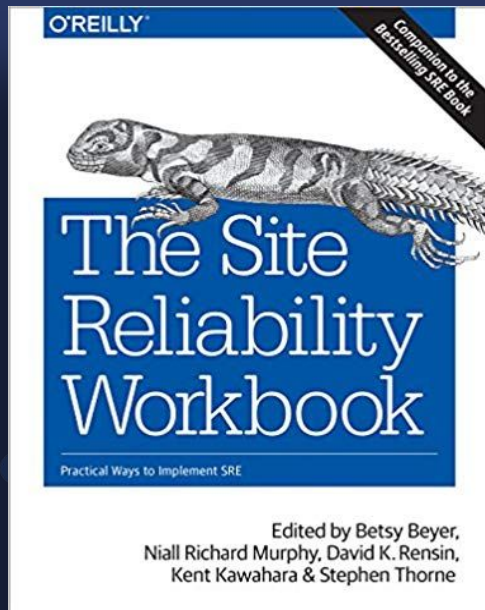
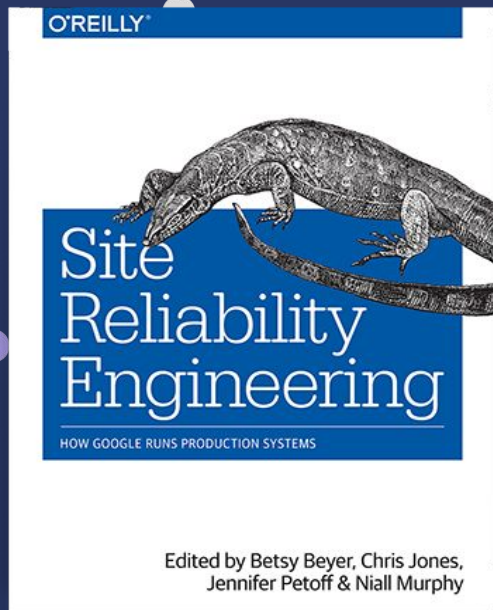
Share Reviews So Other Teams Can Learn From Your Incidents

WHEEL of MISFORTUNE



Summary

1. Failure happens
2. How we manage failure is key
3. What can we do to make tomorrow better than today?



izettle

