# Observing *from* Incidents

Cory Watson
SREcon20 Americas

Surprises, incidents, and mistakes are a fact of life.

# Cory Watson
# Reliability at Stripe

Past: **SRE** @ Twitter, **Observability** @ Twitter & Stripe, **Director** @ SignalFx, **Advocacy and Product** @ Splunk

# Things We **Do** From Incidents

Add an alert for that

   Page someone

   Email someone

   Slack someone

Add a chart for that

Add a dashboard for that

Add a metric for that

Add that to the runbook

Filter those shapes out

Catch this in code review

Add tests for that

Test this in (QA, Staging, Pre-Prod, Blue/Green)

Fix that one thing nobody wants to touch

Org wide "someone needs to fix QA"

Rewrite it all

Learning from Incidents

Photo by Ivo Rainha from Pexels

# Information Overload

"Perhaps a better way is to make memory unnecessary: **put the required information into the world**"
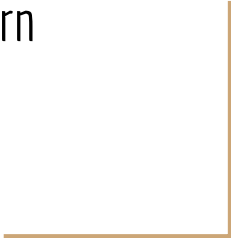Don Norman, *The Design of Everyday Things*

Example: Labels

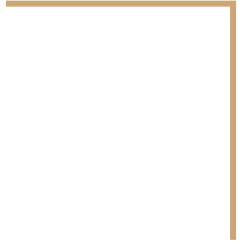How We Interact With Our Systems

# Observing
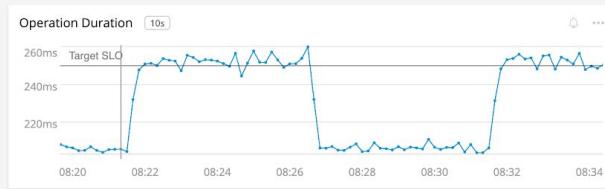## *from*
# Incidents

Putting what we learn
into the world.

# Goal: Add What We Learn To The Tools

- Dashboards, mostly

- Minimizing the need for memory / recall

- Quick, "good enough"

- Awareness of boundaries

- Feedback

# Dashboards

Models of Our Systems

# Problems we're looking for

Surprise

Confusion

Dead ends

Frustration

Edges

Misunderstanding

Boundaries

Inefficiency

Timeliness

Staleness

Lack of confidence

General gripes

# Incidents can inform how we need to improve our tools

# When someone asks a question in an incident, take note of it

# Link, Describe

# Declutter

# Combine



Backend-ops-02
## 100%

**Memory / CPU** ⌄

2020-11-11 05:40:50
— cpu:     20.00%

— memory                                    — cpu

**Statler** Oct 20th, 2:49 pm

I can't find it, hang on. there is so much here.

👍 1

Participant Follow-up

# Divide

# Divide (cont)



cory 12:46
where's the stuff for the darn node 😡 (edited)

# Timeliness & Rate of Change

# Units

Latency                    10s

*Chart description*



10k

5k

0

14:28
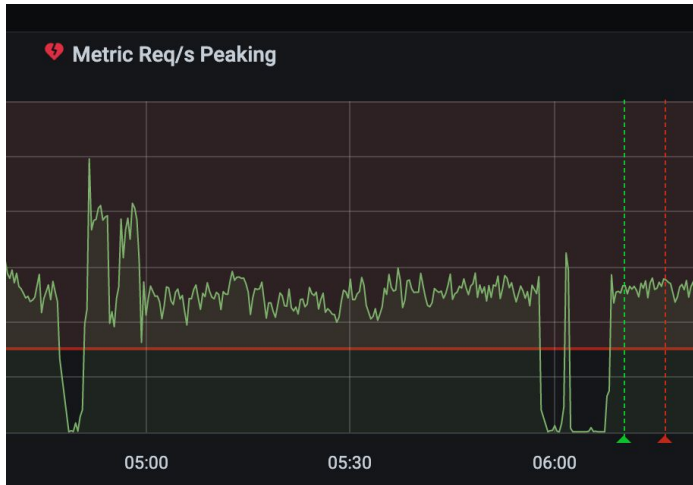
**cory**  06:43
5k what!? where is the code for this

"...the focus should be on control of behavior by **making the boundaries explicit and known**..."
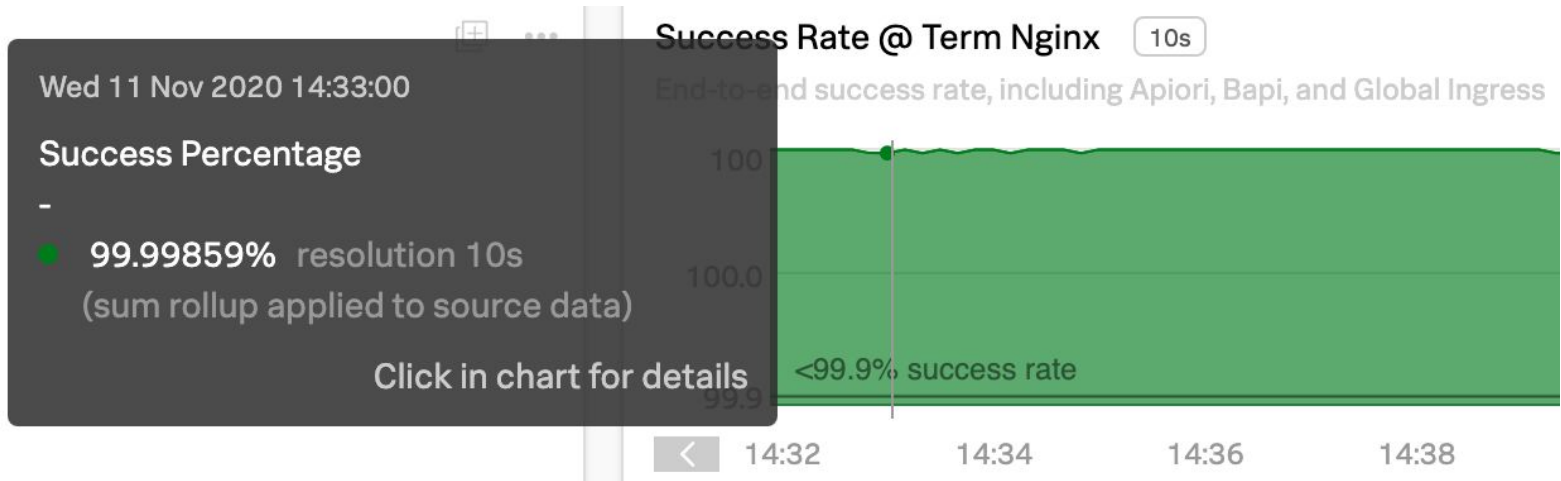Rasmussen, *Risk management in a dynamic society*

# Colocate Signals

# Communicate Goals

# Give Context

Load Shedding    [10s]    📈 🏔 📊 📋 ☰ 4 ⊞ ⚠ 📝

Previous day in gray

```
20k ┤
15k ┤  SLA Violation
10k ┤  Concerning
 5k ┤
  0 ┤
     └─────────┬──────────┬──────────┬──────────┬──────────
            14:45      14:46      14:47      14:48      14:49
```

Wed 11 Nov 2020 14:49:00

Previous Day
-

● **3,629**   resolution 10s
(sum rollup applied to source data)

Click in chart for details

**cory**  07:00
what is normal for this?

# Work With Automation

# Go Analytic

# Stay Perceptive / Mental Models

# Spot Contributing Factors

# Help with Situation Awareness

Open Incidents

cory 07:15
wait, what's happening to <other service>?

# Encode Common Workflows



cory 07:25
how many transactions failed?

cory 07:24
when did the 500s start?

cory 07:29
can we get a list of everyone who was affected?

cory 07:14
how much did this affect <big customer>?

# Summary

# Learn
## From Your
# ~~Incidents~~ Opportunities

# No neutral in design

Easy to fall into

Organic growth

Examine usage, collect signals

Sit with users

Share your work

# Not (just) your vendor / tool problem

What's your model?

Collecting more feedback

Use what you have

# Spot the signs
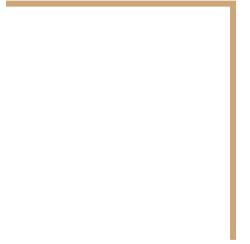
- Dashboards as a tool: need a goal, get state, way to affect state, model

- Mental models, avoid dead ends, allow side to side, analytic

- Events to aid awareness (stuff happening, alert state, using alerts as an "info" channel)

- Bounds, limits, norms

- Automation actions, decisions, and goals

# Would you like to know more?

- *The Design of Everyday Things*, Don Norman
- *Above the Line, Below the Line*, Dr. Richard Cook
- *Risk management in a dynamic society*, Rasmussen
- Ecological Interface Design

[DESIRE TO KNOW MORE INTENSIFIES]

# Thanks!

Twitter: @gphat
Email: cory@onemogin.com