# high reliability infrastructure migrations

Julia Evans
@b0rk

stripe

# about me

infrastructure
engineer
@stripe

payments company
billions of dollars /year

our challenges:
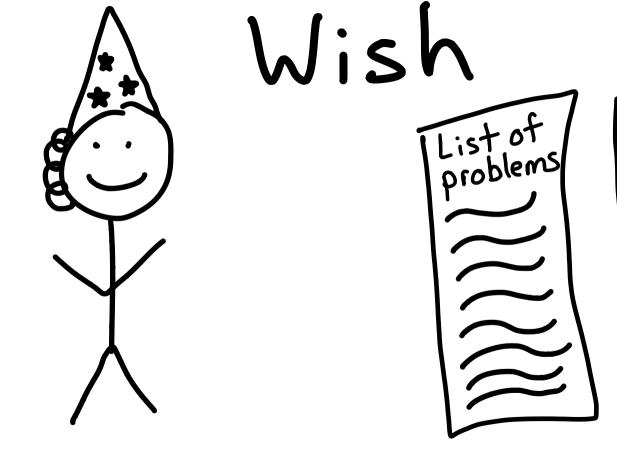
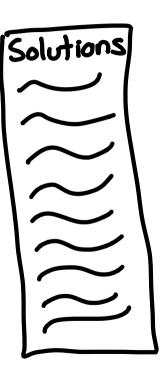~~10 million QPS~~

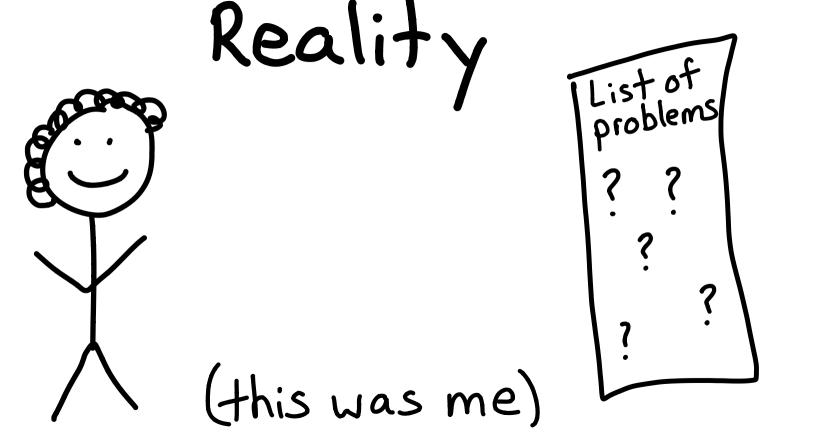~~sub-millisecond latency~~
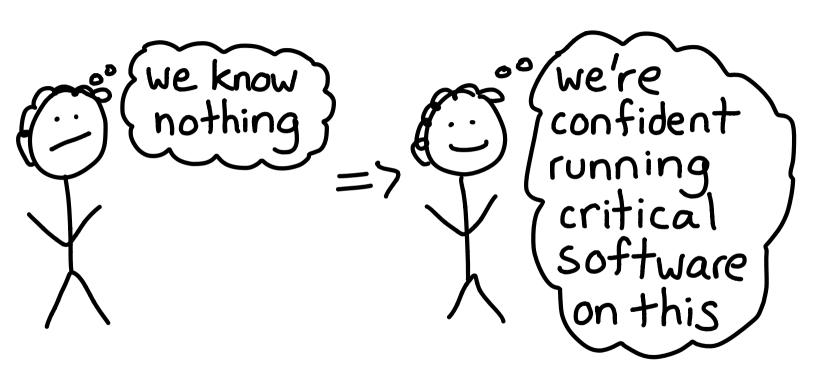
RELIABILITY    SECURITY

# 99.99%

# 1 minute / week

# We made 2 changes

→ move some workloads
to kubernetes

→ use Envoy for all
service-to-service networking

# Wish

# the goal



√ FIXED
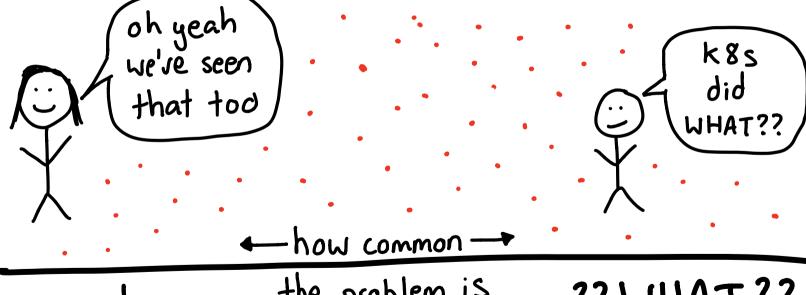
99.99%

normal                    ??WHAT??

# how to get there

- Understand the design
- run gamedays
- classify your failures
- have incidents only once
- make incremental changes
- always have a rollback plan

# Understand kubernetes' design

# k8s design

etcd — API server — kubelet

controllers

everything else

ignore (most)
new software

kubernetes   Envoy

that's it

# theory isn't enough

how can kubernetes break 🔍

learn how _your_ system breaks

# cause problems on purpose

_culture of learning & rigor_

"gamedays"
"DiRT"

_weekly cadence_

# gamedays:

-> test how your system behaves
   under known failures

-> let you learn without duress

-> share knowledge

# Run gamedays

- terminate an etcd instance
- push invalid configuration
- destroy all apiserver instances (or just 1)
- container registry outage
- take down Envoy control plane

Run these in QA, but also in production!

Kubernetes terminated
every running pod
in the cluster "pod eviction"

We fixed & then tested the
fix

classify your
failure modes

~all our failure modes

→ containers don't start
→ permissions errors
→ networking issues

# Reasons pods don't start:

- IAM rate limiting
- scheduler bug
- etcd is down
- lots more

so many reasons

classification
=> monitoring

♥ heartbeat job ♥

Have every
incident
★ only once ★

If you don't understand your incidents, they <u>will</u> happen again

# Your problem space
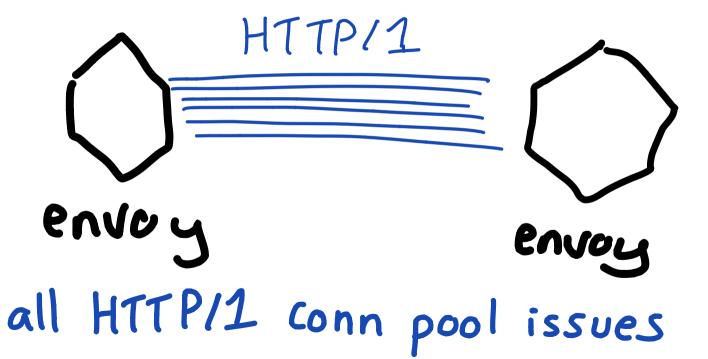
fix
these

99.99%

normal                    ??WHAT??

① Find a problem
② Find causes
③ Implement remediations
④ Problem never comes back*

* usually

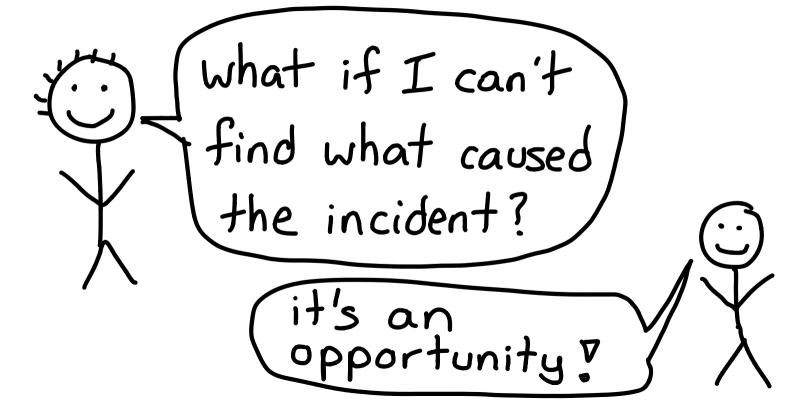# Fix categories of incidents

# Some Envoy issues

request timeouts

connection timeouts

slow request

thundering herd

HTTP/1

envoy                    envoy

all HTTP/1 conn pool issues

# solution:

## use HTTP/2

(Envoy is designed for HTTP/2)

HTTP/2

envoy                                    envoy

# tell your coworkers what you learned

♥ incident reports ♥

example: etcd EBS issue
throttling => leader elections
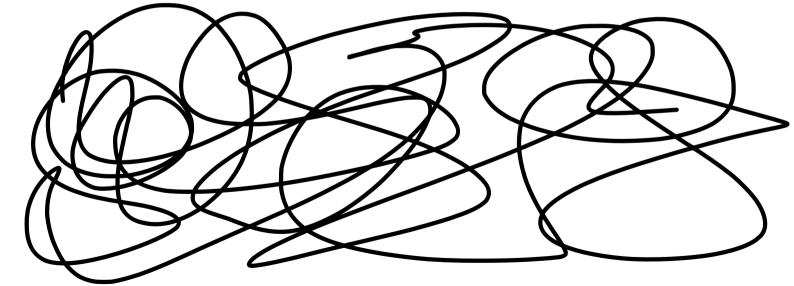
incidents teach you
how to build a
reliable system

make incremental changes

- 5% of traffic
- 1 host
- a non-critical service

establish an
interface boundary

make incremental changes

our deploy system

circa 2017

make incremental changes

client

server

make incremental changes

new client

server

new
client

new
server

"no haunted forests"

# don't expose kubernetes to developers

① reduce cognitive load
② reduce support burden

# escape from YAML:

* skycfg *

skycfg. fun

YAML

```yaml
name: missing-review-finder
owner: risk
schedule: 30 0 * * *
disabled: false
command:
- ruby
- scripts/cron/risk-missing-review-finder
```

- what other attributes are supported?
- what k8s config does it generate?

```
return stripe_service(
    image = default_image,
    command = einhorn(henson_service = "home-srv",
        script = "home/srv.rb",
        workers = 8,
        port = 9768,
    ),
    iam_role = "homesrv.kube.%s.%s" % (
        ctx.vars["stripe.cluster"],
        ctx.vars["stripe.environment"],
    ),
    replicas = 3,
    cpu = kube.cores(4),
    mem = kube.gigabytes(16),
    block_egress = False,
)
```
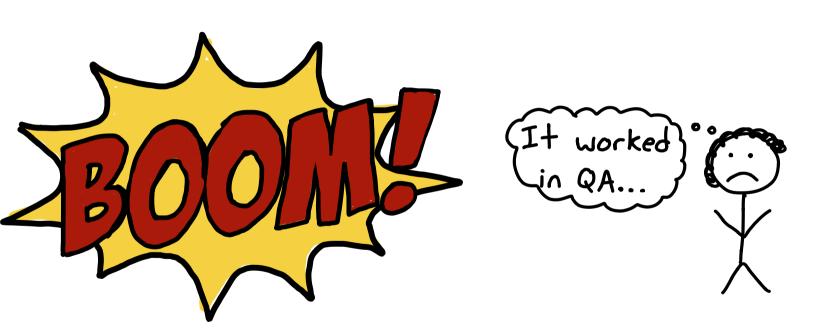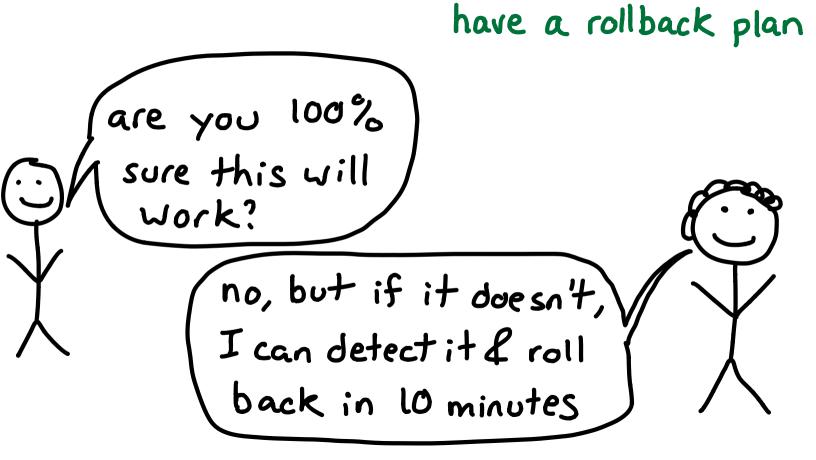
- subset of Python
- typechecked
- sandboxed

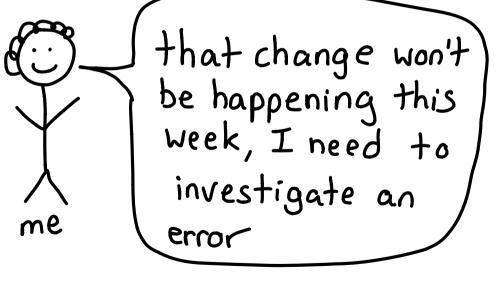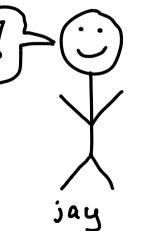github.com/stripe/skycfg

skycfg.fun

always have a

♥ rollback plan ♥

# * playbook *

- Understand the design
- run gamedays
- classify your failures
- have incidents only once
- make incremental changes
- always have a rollback plan

# culture & leadership

- it's ok to start out not being an expert
  **but you need to become one**
- build an engine of learning
- building that expertise takes time

# thanks

a lot!

ps: we're hiring ˇ in Seattle!

stripe.com/jobs

intro 0-4

gamedays 4-7

modes 7-9

only once 9 - 14:30

incremental 14:30 - 18

rollback 18 - 19

end 19 - 22