

Heroku Operations

Noah Zoschke

noah@heroku.com

heroku status

current status and incident report

Widespread Application Outage

Production **5h** **8h** Development **5h** **8h**

Update

We have lost connectivity to some of our infrastructure. Our engineers are working to restore affected systems. We've disabled API access while we work through the issues.

Posted about a year ago, Jun 30, 2012 03:49 UTC

Issue

We're currently experiencing a widespread application outage. We've disabled API access while engineers work on resolving the issues.

Posted about a year ago, Jun 30, 2012 03:25 UTC

Investigating

Our automated systems have detected potential platform errors. We are investigating.

Posted about a year ago, Jun 30, 2012 03:18 UTC

[← Current Status](#)

Problem: Cloud Services

Lost EBS + EC2

API + Customer Databases

30% App Servers

Internal Ops Apps + Instances

AWS API Offline

Can't get new capacity

PagerDuty Crushed

No alerts coming through

Heroku API Offline

Customers helpless

Solution: HA Architecture



Solution: HA Architecture

Solution: People &
Operational Culture



Oh sh*t, I think the pager is blowing up...

A Personal Account...

Problem: Culture

Feature Culture → Too Much Software

No inventory of what's up or down

Surprising dependencies

Hacker Culture → Poorly Written Software

Feature rich, not fool proof

Lots of “beta” services with production workloads

Rockstar Culture → Individual Ownership

Lots of low bus factor

Implicit Culture → Unclear Expectations

Can I escalate?

How do we prioritize services and customers?

“Feature Culture” Side Effect:
Legacy Services

Problem: Legacy Services

Two Routing Services

Router → Nginx/Varnish → Dyno (Bamboo)

Router → Dyno (Cedar)

Two Database Services

Shiny New Dedicated Databases (Heroku Postgres)

Years-old Legacy Shared Databases

Five Metrics Services, etc...

Solution: Sunsetting Culture

Treat Sunsetting as First Class Product and
Engineering Work

Meticulously Catalog Running Services

Celebrate Success When Shutting One Down

Recipe: Lifecycle Board

Prototype → Development → Production →
Deprecated → Deactivated → Sunset

Follow Checklists to Advance

Reflect Explicit Owners

Engineering LifecycleHerokuOrg Visible

Prototyped

pickle

sokoban

test-runner

slug-compiler

ribo

daedalus

kafka-http

Add a card...

In Development

log-shuttle

cloudwatcher

m2met

tempo

manager-api

event-manager-api

ringmaster

megaphone

addons-ss0

scooter

clisaurus

anvil

shh

arsenal

watchtower-tng

direwolf

Production: platform apps

canary

canary-v2

identity

scheduler

vixie

deployhooks

buildkits

codon

direct-to

s3pository

vacuole

reservoir

izanami

maas

umpire

bob

Production: kernel apps

rallgun

rendezvous

hermes

wombat

maestro

logplex

psmgr

api

qitproxy

bastion

ion

darwin

naqios

NTP

splunk

deploymaster

Production: other stuff

stacks

toolbelt

heroku-sudo

heroku-repo

heroku-pq-extras

api-admin

ssl:endpoint

ion-client

logplex-client

keyblaster

elbping

pq_logfebe

logplexc

WAL-E

pq_logplexcollector

htcat

keppa

Deprecated

wolfpack

squid

graphite

ssl:ip

elbssl

ssl:hostname

collectd

backstop

apt-mirror (andromeda)

qit-mirror (andromeda)

face

balrog

sfdc-core-loader

payments-prod

Add a card...

Deactivated

torch

kerosene

taps

ssl:snl

Add a card...

Sunset

apollo

ssl-ip-check

exprd

re/pdump

aorta

pulse

slinky

amitinker

izanami-client

mitt

l2met

emailgateway

Add a card...



“Hacker Culture” Side Effect: Inoperable Software

Recipe: Production Checklist

Code is visible on GitHub

Has operations docs with executable instructions for common tasks

Has a high-fidelity staging setup with production parity

Alerts a human if it is down

Uses structured logging

Enforces SSL access

Any credentials and their rotation procedures are added to “cred rolls” list

Send a launch email to engineering@ describing the new component

Move to Production on the Engineering Lifecycle board

Auto-scaled to maintain the needed number of instances

Set up to terminate unhealthy instances

“Implicit Culture” Side Effect: Platform and Pager Chaos

Problem: Implicit

Do I need to fix these warnings this week? Or put it off?

Can I escalate this alert? Should I?

Should I update the status site? Will someone else?

Recipe: PagerDuty Discipline

Everyone engineer is on-call

Every page is visible in HipChat

Monkey – Everyone should help ack pages in HipChat during work hours

Level 1 – Explicit expectation of first responder after hours

Level 2 – Explicit value that the team has each other's back

Engineering Manager – Explicit accountability for the whole team and its body of work service

Incident Commander – Experts trained in explicit procedures around updating the status site, opening up AWS tickets, paging extra engineers, etc.

Escalation Policies

+ New Escalation Policy

1-1 of 1

Per Page: 25

↑ Runtime



When an incident is created by a service, notify the following on-call:

- | | | |
|---|------------------------------|--------------------------|
| 1 | Runtime Monkey | Escalate after 2 minutes |
| 2 | Runtime Primary | Escalate after 5 minutes |
| 3 | Runtime Secondary | Escalate after 5 minutes |
| 4 | Noah Zoschke | Escalate after 5 minutes |
| 5 | Incident Commander Primary | Escalate after 5 minutes |
| 6 | Incident Commander Secondary | |

Used by 4 services

- Runtime
- Runtime API
- Runtime Pingdom

1-1 of 1

Per Page: 25

Recipe: Pager Metrics

Measure everything and review weekly

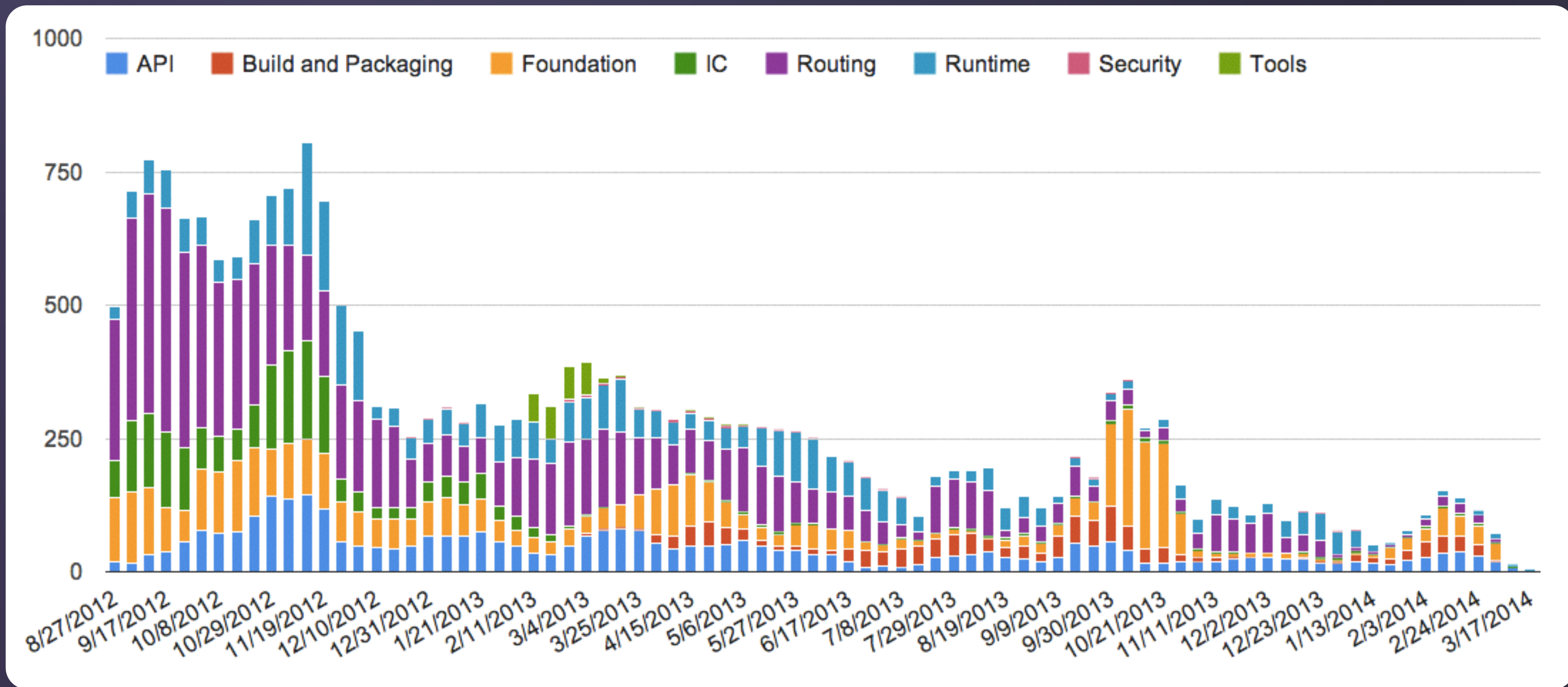
After hours pages are detrimental to engineering health and well being

Engineers deserve weeks with no pages

Engineers have power to improve the operator experience

Engineering Managers are responsible for managing balance between operations and feature work

Service Reliability Engineering (SRE) team is accountable for overall pager burden program



Weekly Pages By Team

