# Building a Culture of Reliability

SRECON EMEA 2017

Arup Chakrabarti
Director of Engineering, PagerDuty

@arupchak

pagerduty

# Disclaimers

pagerduty

I work with ~~smrt~~ smart people

pagerduty

# You are not PagerDuty

pagerduty

# We get this wrong too

@arupchak

pagerduty

# Definitions

pagerduty

# Reliability

pagerduty

# Probability that your software works*

pagerduty

# What every CTO claims they want because numbers

pagerduty

# Culture

@arupchak

**pagerduty**

# Social behavior and norms for a group of people

@arupchak

pagerduty

A way to get your colleagues to behave the way you want them to without staring at them all the time

pagerduty

# Metrics

@arupchak

**pagerduty**

*"Show me the business impact"*
-Your Pointy Haired Manager

**pagerduty**

*"Here is a graph of open File Descriptors going through the roof"*
-Frustrated Engineer

@arupchak

**pagerduty**

*"What the $%#! is a File Descriptor?"*
-Your Pointy Haired Manager

@arupchak

pagerduty

# Business Metrics Managers Care About

@arupchak

pagerduty

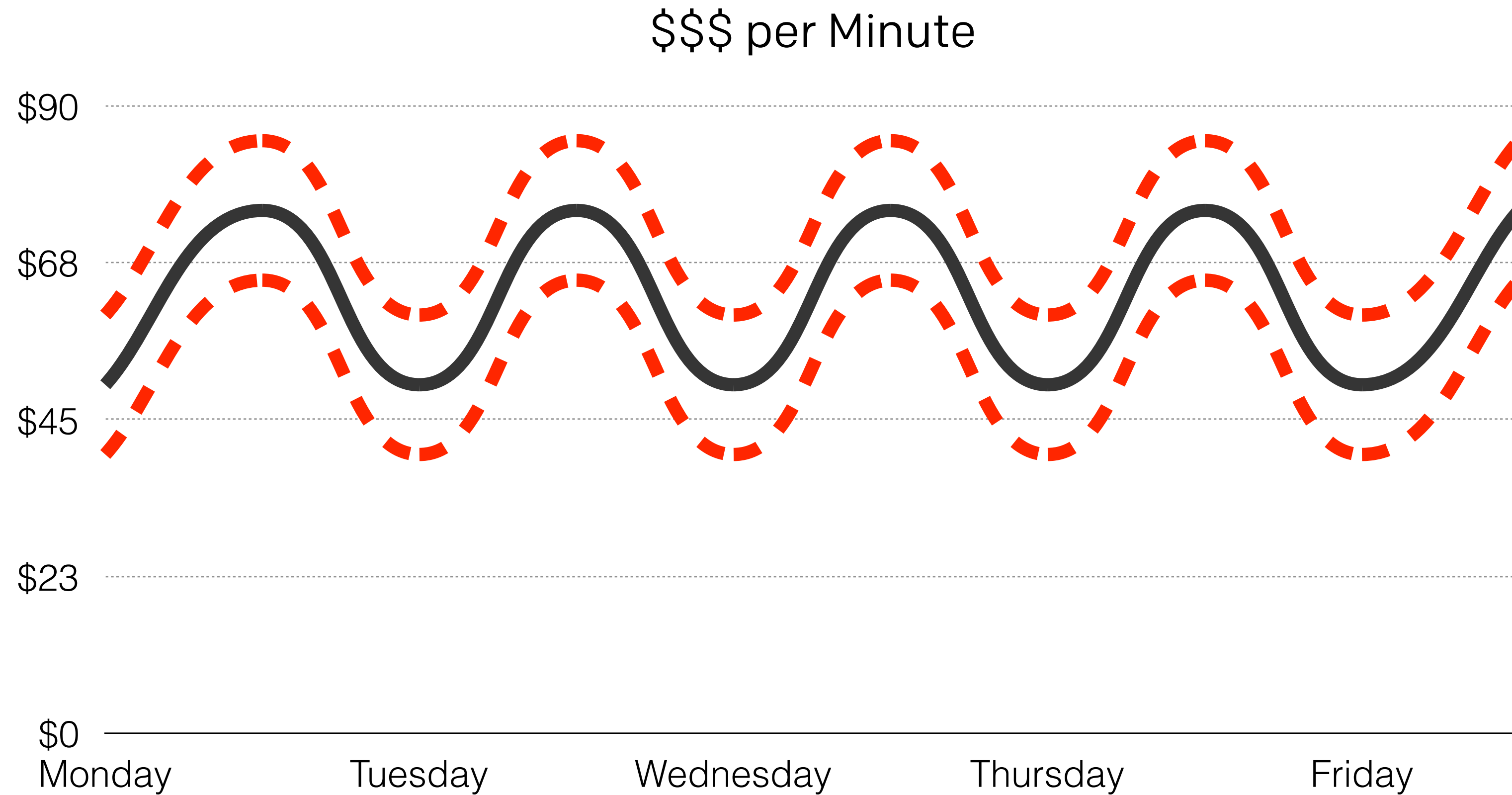# Metrics Your Customers Care About

@arupchak

pagerduty

# Two Types of Online Businesses

- Individual Transaction Businesses

- Subscription Businesses

@arupchak

**pagerduty**

# Individual Transaction Business

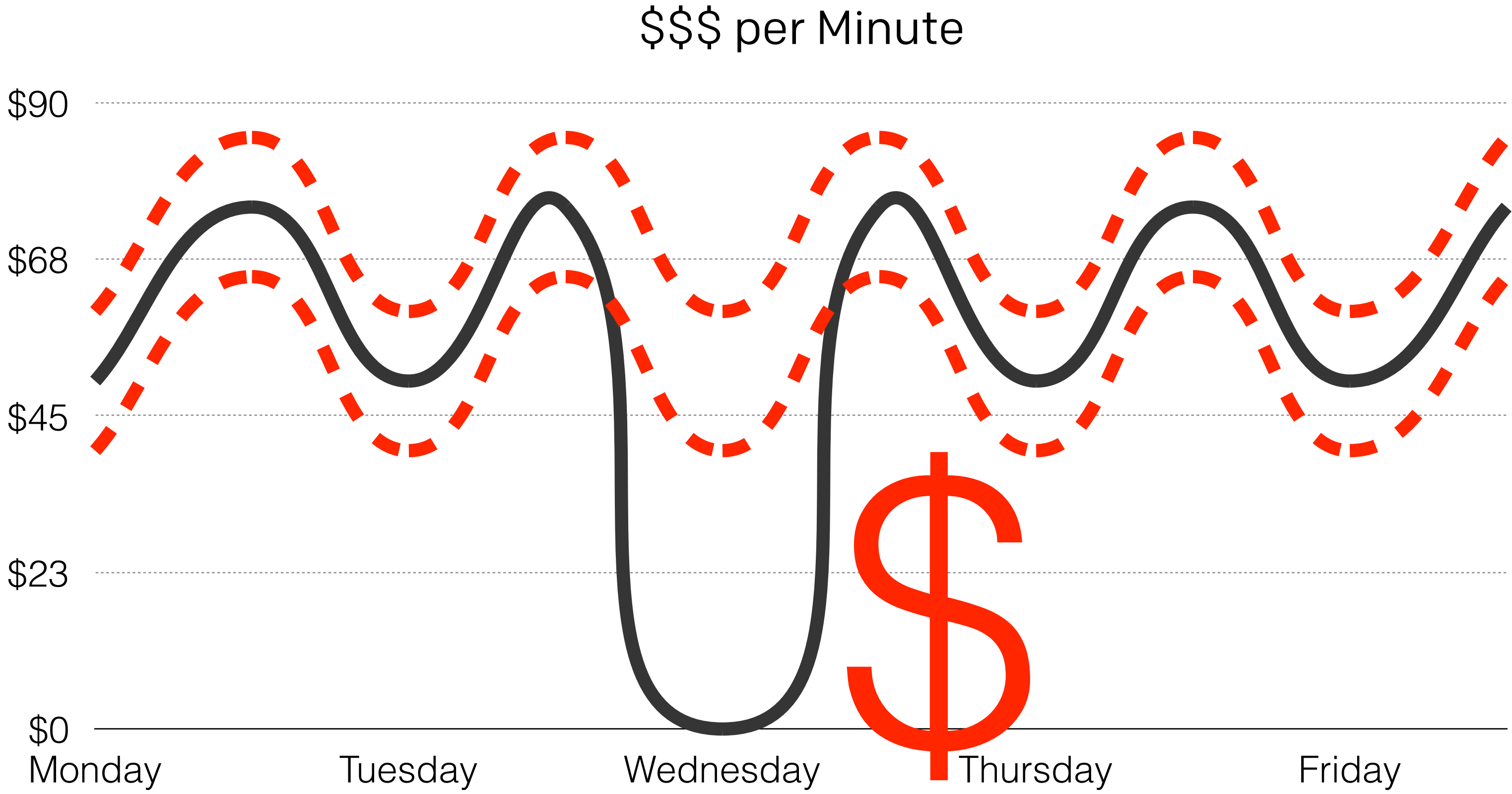## $$$ per Minute

$90

$68

$45

$23

$0

Monday          Tuesday          Wednesday          Thursday          Friday

@arupchak

**pagerduty**

# Individual Transaction Business

## $$$ per Minute



@arupchak

pagerduty

# Individual Transaction Business

## $$$ per Minute

| | Monday | Tuesday | Wednesday | Thursday | Friday |
|---|---|---|---|---|---|
| $90 | | | | | |
| $68 | | | | | |
| $45 | | | | | |
| $23 | | | | | |
| $0 | | | | | |

@arupchak

pagerduty

# Individual Transaction Business

## $$$ per Minute



@arupchak

pagerduty

# Individual Transaction Business



@arupchak

pagerduty

# Individual Transaction Business
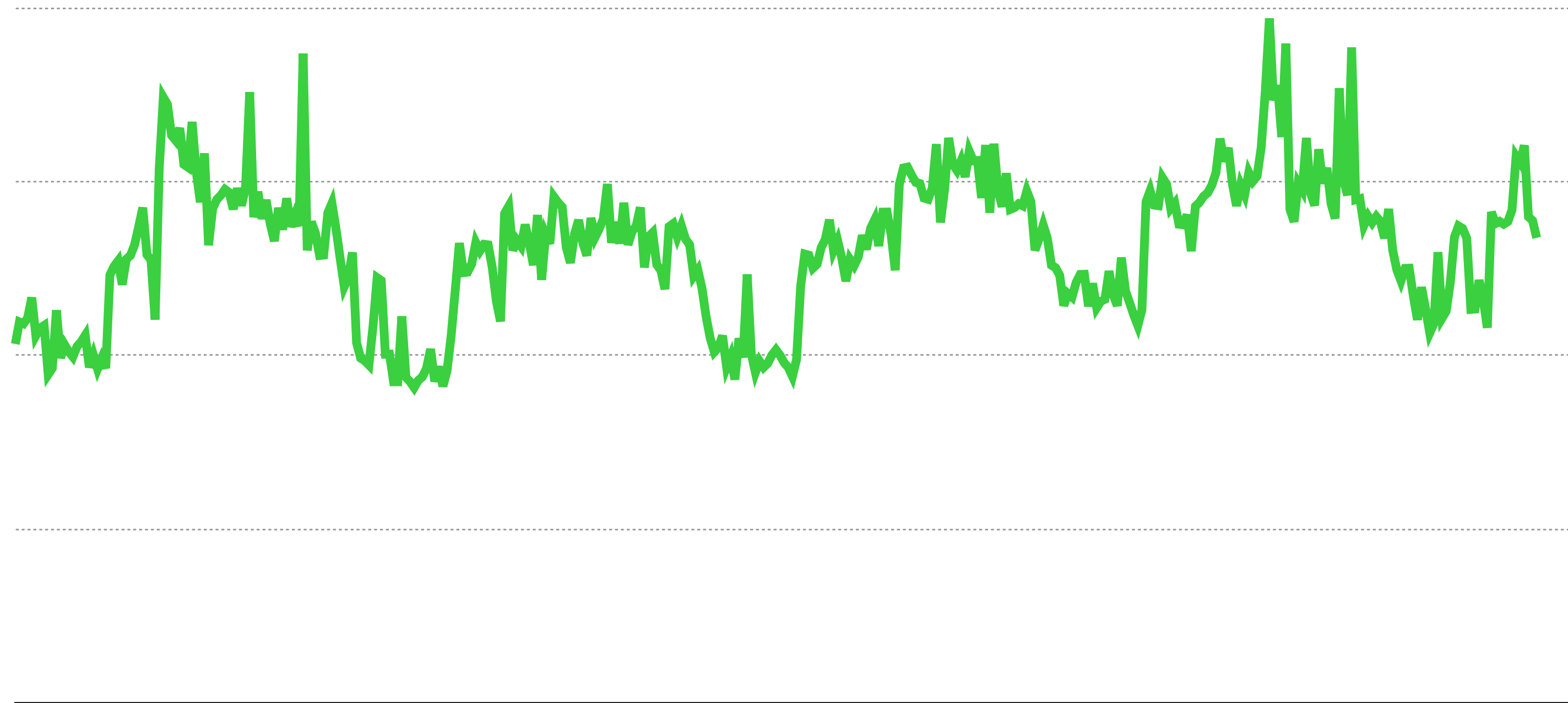
$$$ per Minute

@arupchak

# Subscription Businesses

- Cannot solely measure when you make money

- Poor Reliability erodes trust and will cause you lose revenue

- Need to find something between how money is made and what customers care about

@arupchak

**pagerduty**

# Subscription Businesses

Incidents Resolved per Hour - July 2017



@arupchak

pagerduty

# Finding the right metrics is hard

@arupchak

pagerduty

# But still worth looking for

@arupchak

**pagerduty**

# More People On-Call

# Customers do not care who gets paged

**pagerduty**

# Customers just want to use your service
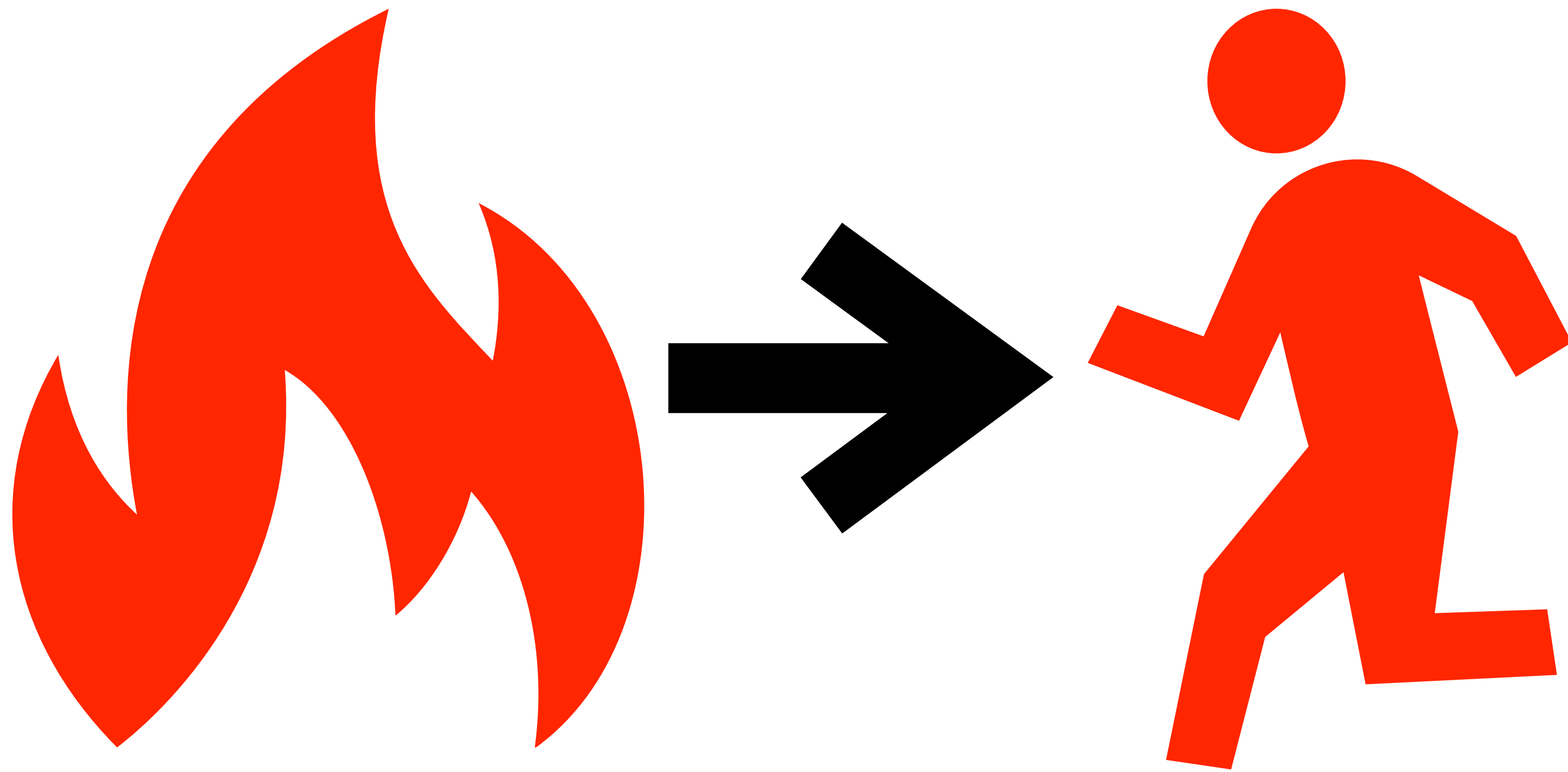
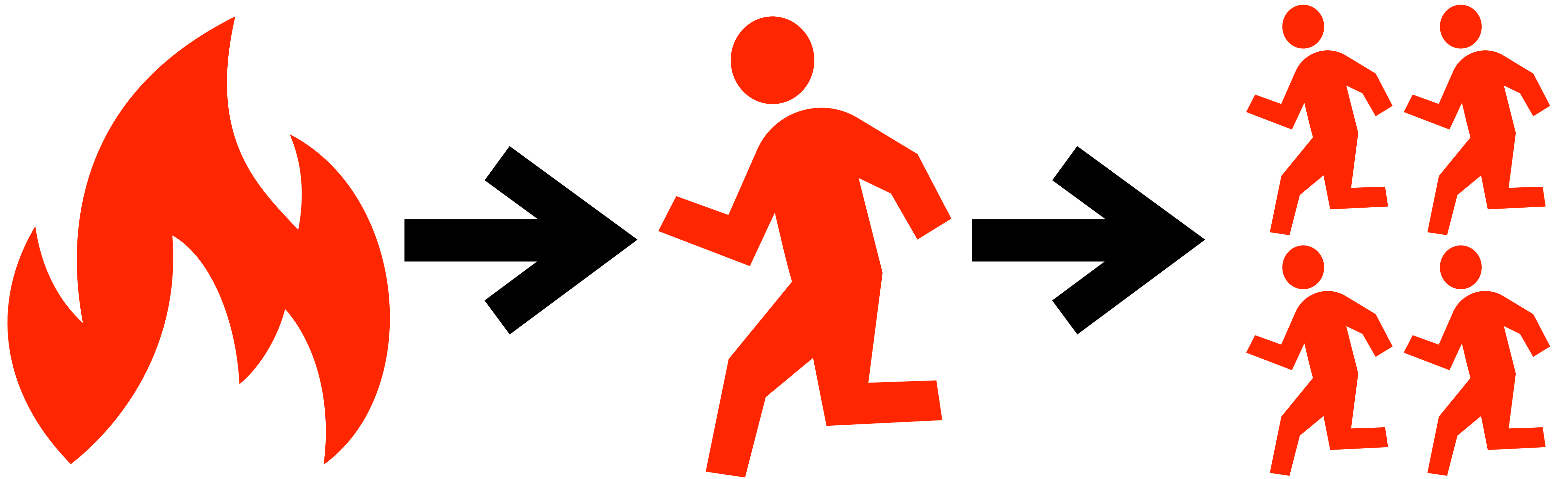@arupchak

pagerduty

# Centralized Operations Engineering Org

@arupchak

pagerduty
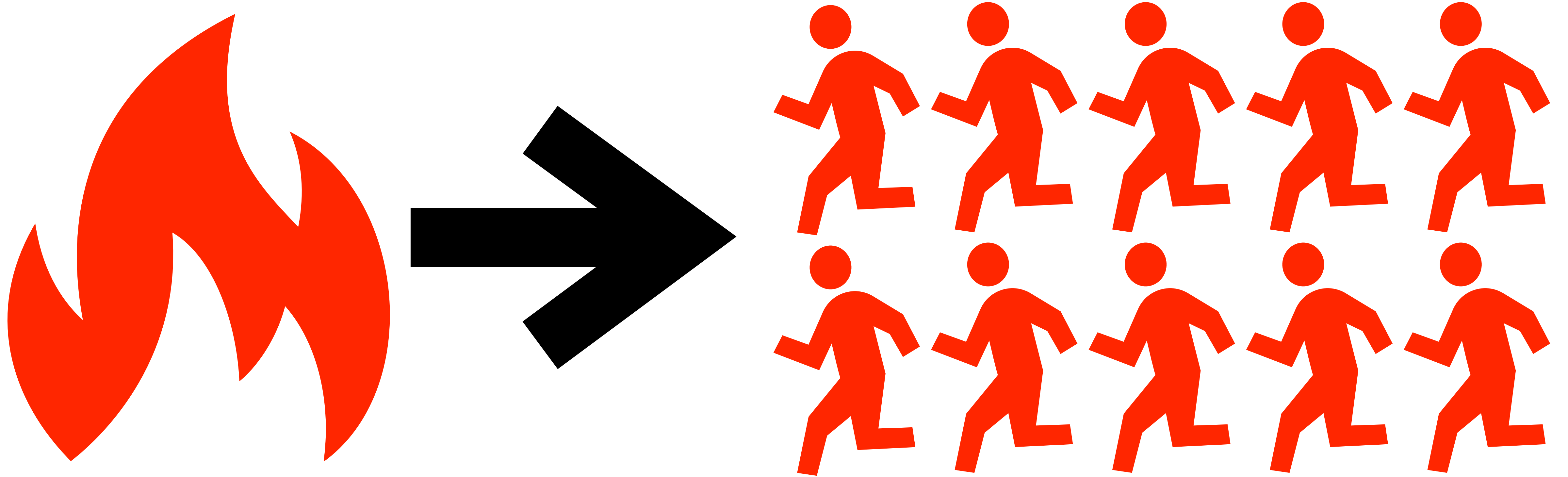
# Centralized Operations Engineering Org



@arupchak

**pagerduty**

# Centralized Operations Engineering Org



@arupchak

pagerduty

# Distributed Operations Engineering Org

@arupchak

pagerduty

# Distributed Operations Engineering Org

# Distributed Operations ~~Engineering~~ Org



Eng

Product

HR

UX

Marketing

Execs

@arupchak

pagerduty

# Distributed Operations Org

- Sets expectations around availability of people

- More small incidents over single major incident

- Builds empathy and why Reliability is hard

@arupchak

**pagerduty**

# Tooling and Processes

pagerduty

*"If we just install Nagios, everything will be fine and all of our problems will be solved"*
-Arup in 2002

@arupchak

pagerduty

*"We humans co-evolve with our tools. We change the tools, and the tools change us, and that cycle repeats."*
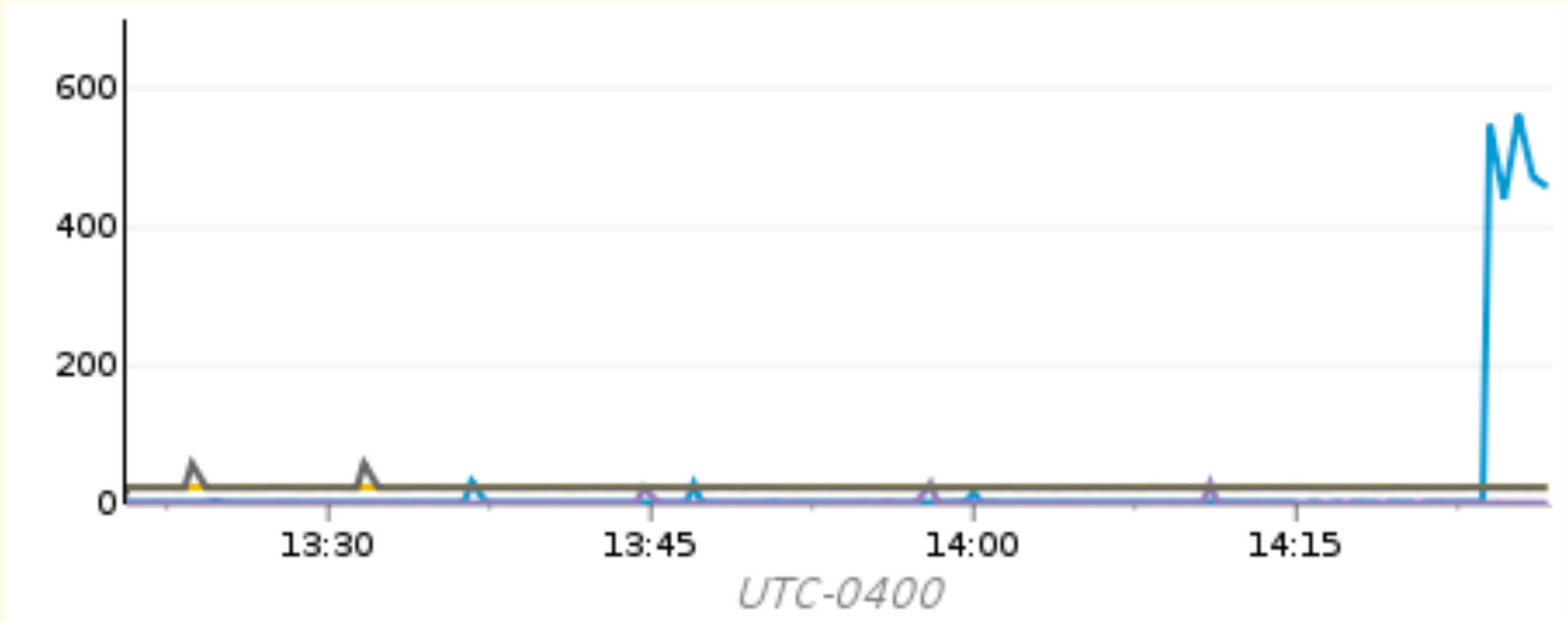-Jeff Bezos

# Started Small

pagerduty

| Arup Chakrabarti | cass08 is still being marked as up | Oct–18 11:25 AM |
| John Laban | | Oct–18 11:27 AM |



| John Laban | cass10 latency | Oct–18 11:27 AM |
| John Laban | @hipchat–Failure_Friday yep, cass08 looks slow *Via Datadog* | Oct–18 11:27 AM |
| Arup Chakrabarti | cass09 can haz  latency now | Oct–18 11:29 AM |
| Arup Chakrabarti | both are still part of the ring | Oct–18 11:29 AM |

@arupchak

**pagerduty**

# Got Bigger and Smarter

pagerduty

**Chaos Cat** APP 12:55 PM
!status

**Officer URL** APP 12:55 PM

> Status: NORMAL

**Chaos Cat** APP 12:55 PM
!ff flaky-network-roulette production

**Officer URL** APP 12:55 PM
flaky-network-roulette chose prod-gemini████ as the victim. Network latency/loss will be added, and automatically removed in 7 minutes. To remove it early, run `!ff unlatency-loss-node production prod-gemini███`

Igor job #33594 created

Starting job #33594 [repo: smoothie, ref: master, cmd: ff_latency_loss, args: {"env"=>"production", "hostlist"=>"prod-gemini███"}]

@chaoscat: ✅ Completed job #33594 [repo: smoothie, ref: master, cmd: ff_latency_loss, args: {"env"=>"production", "hostlist"=>"prod-gemini███"}]

@arupcʜaк

pagerduty

**?**

**Chaos Cat** APP 12:55 PM
!status

**Officer URL** APP 12:55 PM

> Status: NORMAL

**Chaos Cat** APP 12:55 PM
!ff flaky-network-roulette production

**Officer URL** APP 12:55 PM
flaky-network-roulette chose prod-gemini▓▓▓▓ as the victim. Network latency/loss will be added, and automatically removed in 7 minutes. To remove it early, run `!ff unlatency-loss-node production prod-gemini▓▓▓`

Igor job #33594 created

Starting job #33594 [repo: smoothie, ref: master, cmd: ff_latency_loss, args: {"env"=>"production", "hostlist"=>"prod-gemini▓▓▓▓"}]

@chaoscat: ✅ Completed job #33594 [repo: smoothie, ref: master, cmd: ff_latency_loss, args: {"env"=>"production", "hostlist"=>"prod-gemini▓▓▓▓"}]

**Chaos Cat** APP 12:55 PM
!status

**Officer URL** APP 12:55 PM

```
Status: NORMAL
```

**Chaos Cat** APP 12:55 PM
!ff flaky-network-roulette production

**Officer URL** APP 12:55 PM
flaky-network-roulette chose prod-gemini▓▓▓ as the victim. Network latency/loss will be added, and automatically removed in 7 minutes. To remove it early, run `!ff unlatency-loss-node production prod-gemini▓▓▓`

Igor job #33594 created

Starting job #33594 [repo: smoothie, ref: master, cmd: ff_latency_loss, args: {"env"=>"production", "hostlist"=>"prod-gemini▓▓▓"}]

@chaoscat: ✅ Completed job #33594 [repo: smoothie, ref: master, cmd: ff_latency_loss, args: {"env"=>"production", "hostlist"=>"prod-gemini▓▓▓"}]

@arupcнaк

pagerduty

**Chaos Cat** APP 12:55 PM
!status

**Officer URL** APP 12:55 PM
> Status: NORMAL

**Chaos Cat** APP 12:55 PM
!ff flaky-network-roulette production

**Officer URL** APP 12:55 PM
flaky-network-roulette chose prod-gemini▒▒▒▒ as the victim. Network latency/loss will be added, and automatically removed in 7 minutes. To remove it early, run `!ff unlatency-loss-node production prod-gemini▒▒▒▒`

Igor job #33594 created

Starting job #33594 [repo: smoothie, ref: master, cmd: ff_latency_loss, args: {"env"=>"production", "hostlist"=>"prod-gemini▒▒▒▒"}]

@chaoscat: ✅ Completed job #33594 [repo: smoothie, ref: master, cmd: ff_latency_loss, args: {"env"=>"production", "hostlist"=>"prod-gemini▒▒▒▒"}]

@arupchaк

pagerduty

# Reboot Roulette (Tool)

pagerduty

**Ka** 11:50 AM
!ff reboot-roulette production

**Officer URL** APP 11:50 AM
reboot-roulette chose prod-permissions ▓▓▓ as the victim

Igor job #32138 created

Starting job #32138 [repo: smoothie, ref: master, cmd: ff_reboot, args: {"env"=>"production", "hostlist"=>"prod-permissions▓▓▓"}]

✅ Completed job #32138 [repo: smoothie, ref: master, cmd: ff_reboot, args: {"env"=>"production", "hostlist"=>"prod-permissions▓▓▓"}]

@arupchak

pagerduty

# Major Incident Response (Process and Tooling)

@arupchak

pagerduty

# Started Really Poorly

@arupchak

pagerduty

# Got A Little Better Each Time

@arupchak

pagerduty

| Command | Description |
| --- | --- |
| `!ic who` | Displays who the current primary IC is.<br>(Also useful for when people join the call to see who the assigned IC is) |
| `!ic who backup` | Same as above but for the backup. |
| `!ic page` | "Pages" the current primary IC by issuing a Twilio call to them. Also @ mention's them on Slack. |
| `!ic page backup` | Same as above, but for the backup. |

@arupchak

pagerduty

# Still Not Perfect

@arupchak

pagerduty

**PagerDuty Incident
Response
Documentation**
Source on GitHub

DOWNLOAD
STARS 410

Home

Being On-Call

Before an Incident

During an Incident
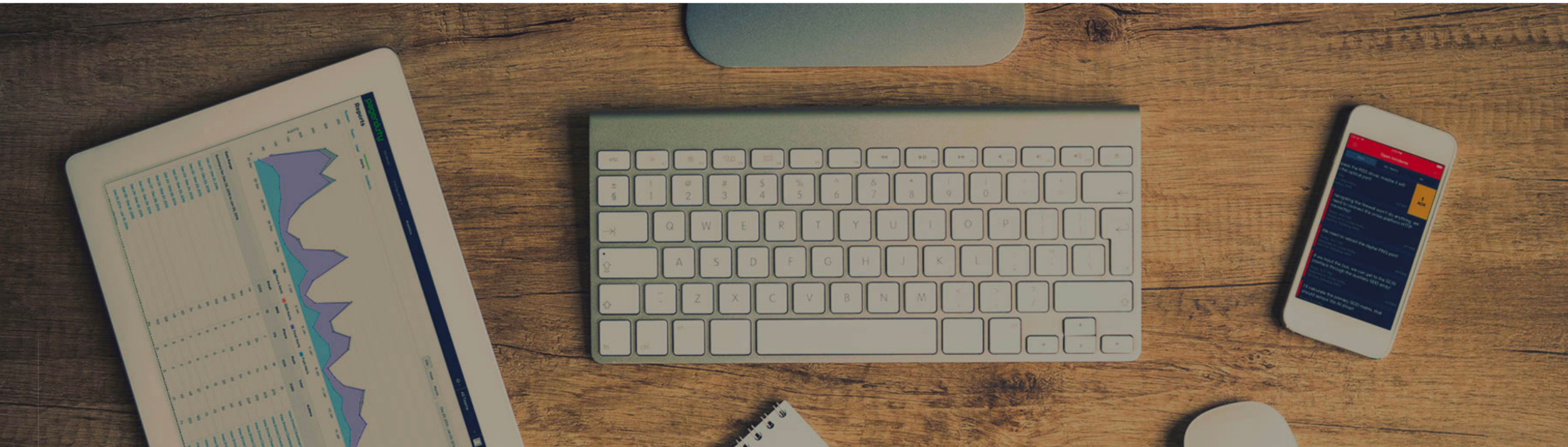
After an Incident

Training

Additional Resources

On-Call

Being On-Call

Alerting Principles

Before an Incident

What is an Incident?

This documentation covers parts of the PagerDuty Incident Response process. It is a cut-down version of our internal documentation, used at PagerDuty for any major incidents, and to prepare new employees for on-call responsibilities. It provides information not only on preparing for an incident, but also what to do during and after. It is intended to be used by on-call practitioners and those involved in an operational incident response process (or those wishing to enact a formal incident response process). See the about page for more information on what this documentation is and why it exists.

## Being On-Call

If you've never been on-call before, you might be wondering what it's all about. These pages describe what the expectations of being on-call are, along with some resources to help you.

> Being On-Call - *A guide to being on-call, both what your responsibilities are, and what they are not.*

> Alerting Principles - *The principles we use to determine what things page an engineer, and what time of day they page.*

**pagerduty**

# Internal Liaison Role (Process)

@arupchak

pagerduty

# Over-communicate during Major Incidents

@arupchak

pagerduty

**demitri** 🌴 1:06 PM
@here There has been an issue causing a brief period of degraded service in ▓▓▓▓ ▓▓▓▓▓▓ ▓▓▓▓▓. We have recovered from SEV-2 and monitoring the situation. (edited)

   💬 1 reply   17 days ago

**demitri** 🌴 1:11 PM
@here At this point we are fairly confident that customer impact is zero. SRE is continuing to monitor and investigate.

@arupchak

pagerduty

Improving Reliability means constantly failing, constantly recovering, and constantly learning

@arupchak

pagerduty

# Yes, it can be exhausting, but it is worth it

pagerduty

Improving Culture means constantly failing, constantly recovering, and constantly learning

@arupchak

pagerduty

Yes, it can be even more exhausting, but it is really really really worth it

@arupchak

pagerduty

# Thank You

Arup Chakrabarti
Director of Engineering, PagerDuty

@arupchak

**pagerduty**

# Related Reading

- https://response.pagerduty.com/

- https://www.pagerduty.com/blog/intern-insights-on-call-experience/

- https://www.pagerduty.com/blog/failure-fridays-four-years/

- https://speakerdeck.com/arupchak

@arupchak

**pagerduty**