

Humane On-Call

Alerting doesn't have to be painful

Bad alerting adds friction to operating production software

**Bad alerting adds friction to operating
production software**

We can do better!

On-Call in a Nutshell



Being *available* to handle an incident



Constantly monitor production

As part of a *rotation*

Platform Primary

Export ▾ ⚙ ▾

On-Call Now
Mario Fernandez from Oct 14, 2021 at 00:00 to Oct 18, 2021 at 00:00

Your Next On-Call
Oct 14, 2021 at 00:00 to Oct 18, 2021 at 00:00

Teams using this Schedule
Platform Services

Escalation Policies using this Schedule
Platform

November

10/31	M 11/1	T 11/2	W 11/3	T 11/4	F 11/5	S 11/6	S 11/7	M 11/8	T 11/9	W 11/10	T 11/11	F 11/12	S 11/13
Another Developer				Mario Fernandez					Another Developer				

Why do this to
ourselves?

Reason 1

Users expect constant availability

Google

facebook

NETFLIX

Reason 2

You build it, you run it

stevesmith.tech/blog/you-build-it-you-run-it/

Alerting 101

**Detect problems in production before
your customers do**

Monitoring + Notification

The Four Golden Signals

Latency

Traffic

Errors

Saturation

sre.google/sre-book/monitoring-distributed-systems/



Constant Interruptions

User	Time on Call	High-Urgency Incidents	Acknowledged	Timeout Escalations	Manual Escalations	Reassignments	MTTA
[REDACTED]	168h 0m	2	100% (2)	50% (1)			57m
[REDACTED]	168h 0m	3	100% (3)			33% (1)	2m
[REDACTED]	168h 0m	1					
[REDACTED]	168h 0m	10	80% (8)	10% (1)			4m
[REDACTED]	123h 0m	1	100% (1)	100% (1)			13m
Mario Fernandez	96h 0m	20	40% (8)	30% (6)			5m
[REDACTED]	96h 0m	10	50% (5)	40% (4)			12m
[REDACTED]	72h 0m	17	71% (12)	6% (1)			2m
[REDACTED]	45h 0m						
[REDACTED]	45h 0m						
[REDACTED]		37		19% (7)			

Bad Night(s)

Today 08:03

All Missed Edit

Recents

PagerDuty Outgoing... mobile	01:02	(i)
PagerDuty Outgoi... mobile	00:05	(i)

Entgangener Anruf:
██████████
am 29.09. um 01:03 Uhr
Ihre Telekom

Entgangener Anruf:
██████████
am 29.09. um 01:51 Uhr
Ihre Telekom



[P1] [Recovered on {resource_name:get_/v1/get/min-layout}] service prod / Error Rate [Change%]

created_by:terraform env:prod monitor resource_name:get_/v1/get/min-layout

Error Rate is at 9.577% for the route[get_/v1/get/min-layout]



Show More ▾

Fri Oct 08 2021 08:33:05 CEST (2 days and 3 hours ago)

Flakiness



[P1] [Triggered on {resource_name:get_/v1/get/min-layout}] service prod / Error Rate [Change%]

created_by:terraform env:prod monitor resource_name:get_/v1/get/min-layout

Error Rate is at 13.551% for the route[get_/v1/get/min-layout]



Show More ▾

Fri Oct 08 2021 08:31:05 CEST (2 days and 3 hours ago)



Alerting Dysfunctions

Mixed abstraction levels

Lack of automation

Noisy alerts

Inadequate tools

Mismatched tuning

Mixed Abstraction Levels

What are you monitoring?

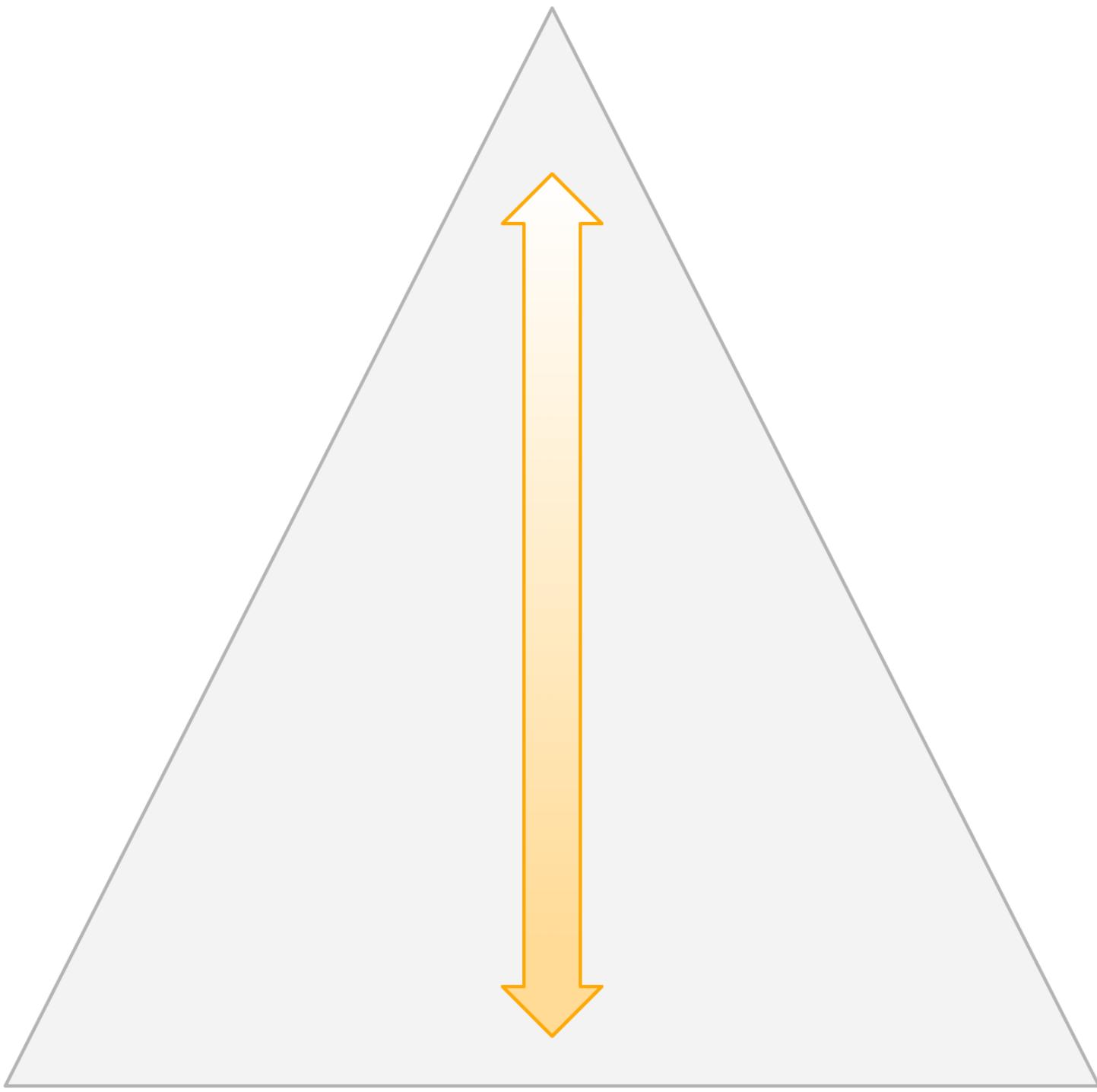
What are you monitoring?

High-Level: Close to the user

What are you monitoring?

***High-Level:* Close to the user**

***Low-Level:* Close to the infrastructure**





OK Last ran 1 min ago (Oct 18, 2021, 5:24 pm)

Navigation Available

STATUSOK**STEPS**

5 / 5

DURATION

29s

LOCATION

Frankfurt (AWS)

DEVICE

Laptop Large

BROWSER

edge 95.0.1020.9

TIME RAN

3 mins ago (Oct 18, 2021, 5:19 pm)

RUN TYPE

Scheduled

	SCREENSHOT	ACTION	DURATION
0 ✓		.Navigate to start URL Started at about:blank	2 s
1 ✓		.Login > Expand 10 steps	26.2 s
2 ✓		✓ Navigation Panel present	207 ms
3 ✓		✓ User displayed Should not be empty	194 ms
4 ✓		✓ "Help & Support" is present	211 ms

CPU load is very high on {{host.name}}

Properties

Metric Monitor
ID: 11872480QUERY `avg(last_5m):100 - avg:system.cpu.idle{{host}} by {host} > 60`

TAGS

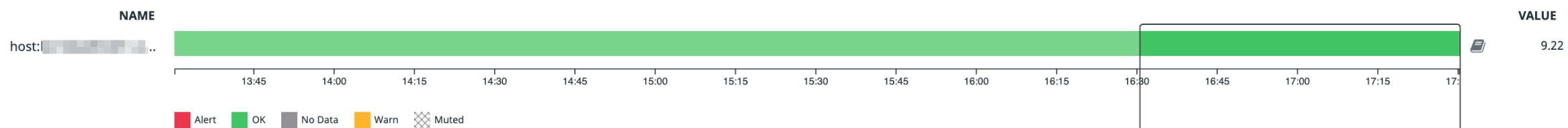
PRIORITY Not Defined

Filter monitor groups and their events

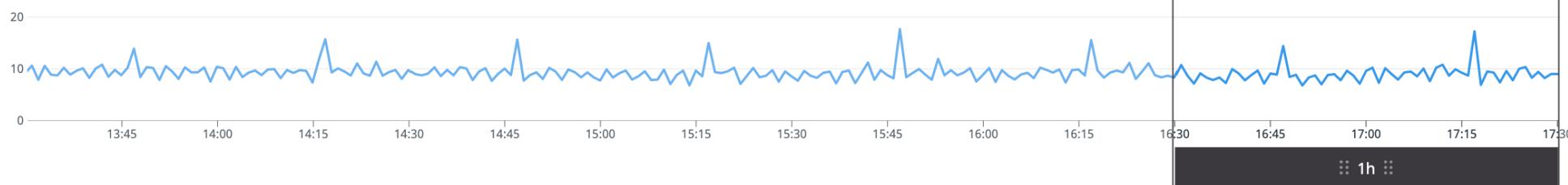
 Alert 0 Warn 0 No Data 0 OK 1 | 1 of 1 groups

Status & History

GROUP STATUS

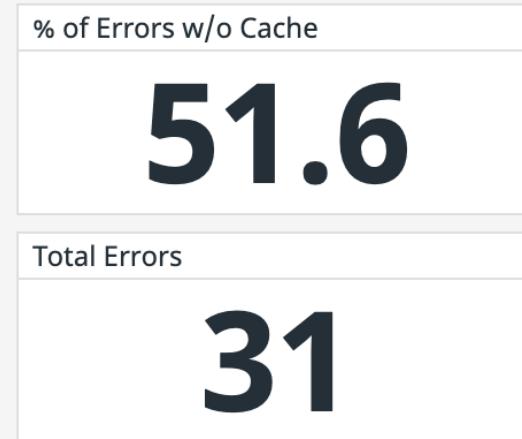
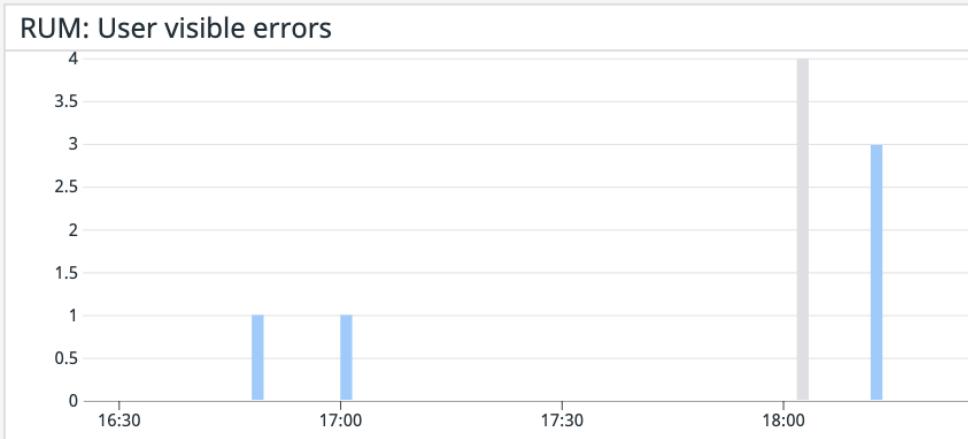


HISTORY

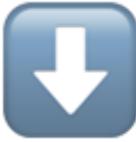


**Choose the alert based on what you
want to monitor**

User Impact



Learning 1



Use the right alert

Lack of Automation

Manual work is not reproducible

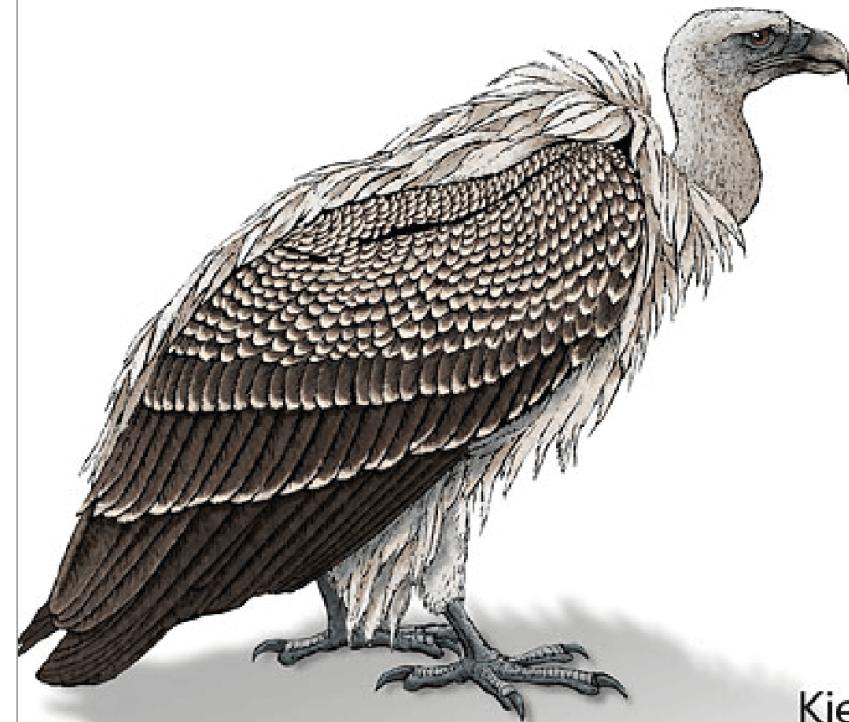
Automation is crucial in alerting

O'REILLY®

Second
Edition

Infrastructure as Code

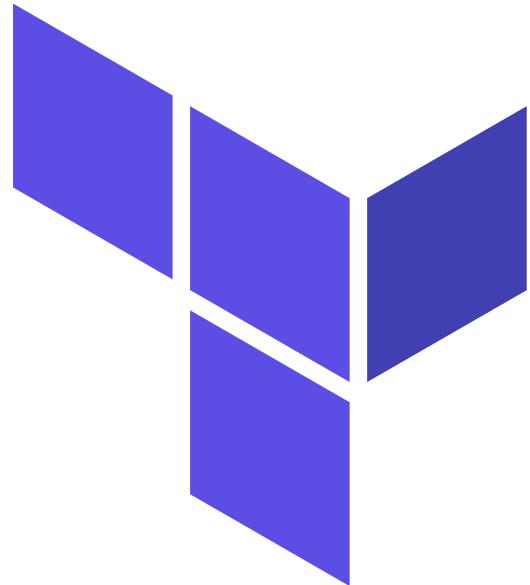
Dynamic Systems for the Cloud Age



Kief Morris

Reduce maintenance

**Reduce maintenance
Spread best practices**



HashiCorp
Terraform

<https://registry.terraform.io/browse/providers>

```
resource "datadog_monitor" "monitor" {
  count = try(var.enabled, false) == true ? 1 : 0

  name          = var.name
  type          = var.type
  query         = var.query
  message       = var.message

  monitor_thresholds {
    warning = lookup(var.monitor_thresholds, "warning", null)
    ok      = lookup(var.monitor_thresholds, "ok", null)
    critical = lookup(var.monitor_thresholds, "critical", null)
  }

  priority      = var.priority
}
```

```
resource "pagerduty_team" "team" {
  name = "Service Team"
}

resource "pagerduty_schedule" "schedule" {
  name      = "Weekly Engineering Rotation"
  time_zone = "Europe/Berlin"

  layer {
    name          = "Night Shift"
    start         = "2015-11-06T20:00:00-05:00"
    rotation_virtual_start = "2015-11-06T20:00:00-05:00"
    rotation_turn_length_seconds = 86400
    users          = [pagerduty_user.team.id]
  }
}

teams = [pagerduty_team.example.id]
}
```

**What about more complex stuff?
Multiwindow, Multi-Burn-Rate Alerts**

Error Rate SLO



Fast Burn Rate



Medium Burn Rate



Slow Burn Rate



```
module "error_rate_slo_burn_rate_long_window" {
  source  = "../../monitor"
  enabled = var.enabled

  name      = <<-EOT
  ${var.service_name} ${var.environment} -
  Error Rate SLO [Long Window, ${var.burn_rate}x Burn rate]"
EOT

  query = <<-EOT
sum(${local.period}):(
  sum:trace.${local.metric}.request.errors{env:prod,service:${var.service_name}}
    .rollup(sum, ${var.window_in_seconds}) /
  sum:trace.${local.metric}.request.hits{env:prod,service:${var.service_name}}
    .rollup(sum, ${var.window_in_seconds})
) * 100 > ${local.threshold}
EOT

  monitor_thresholds = {
    critical = local.threshold
  }
}
```

```
module "error_rate_slo_quick_burn" {
  source  = "../error_rate_slo_burn_rate"
  enabled = var.enabled

  service_name = var.service_name
  environment = var.environment

  slo           = var.slo
  burn_rate     = 14.4
  window_in_seconds = 3600

  language = var.language
}
```

```
module "error_rate_slo" {
  source  = "../../monitor"
  enabled = var.enabled

  name      = "${var.service_name} - Error Rate SLO [${var.slo}%]"
  type     = "composite"

  query = <<-EOT
(module.error_rate_slo_quick_burn.long_window_id &&
  module.error_rate_slo_quick_burn.short_window_id)
 ||
(module.error_rate_slo_slow_burn.long_window_id &&
  module.error_rate_slo_slow_burn.short_window_id)
EOT

  message = templatefile("${path.module}/message.md", {
    service      = var.service_name
    environment = var.environment
    notify       = local.notify
  })
}
```

Scale alerts inside a team

Scale alerts across teams

Learning 2



Leverage automation at scale

Noisy Alerts

**When the signal is drowned by the
noise**

Alert fatigue leads to ignored alerts

atlassian.com/incident-management/on-call/alert-fatigue

A lot of the noise is *accidental*

Service [REDACTED] has low pod count

Integration Monitor

QUERY `avg(last_5m):sum:kubernetes_state.deployment.replicas_available{env:prod} < 1`

TAGS env:prod

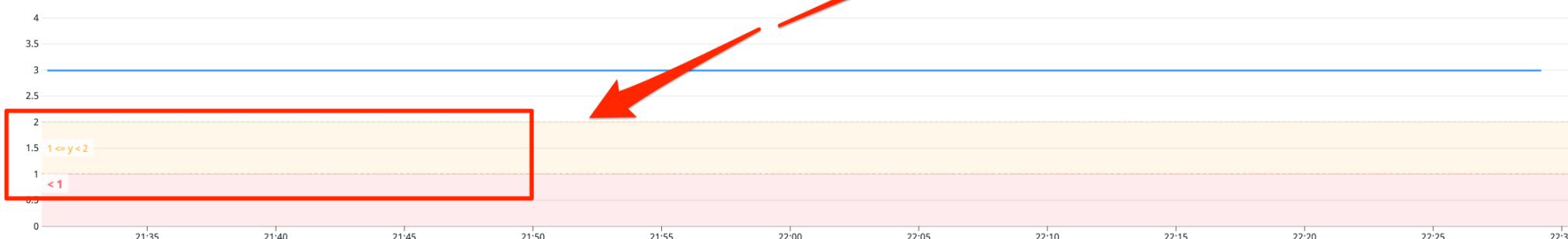
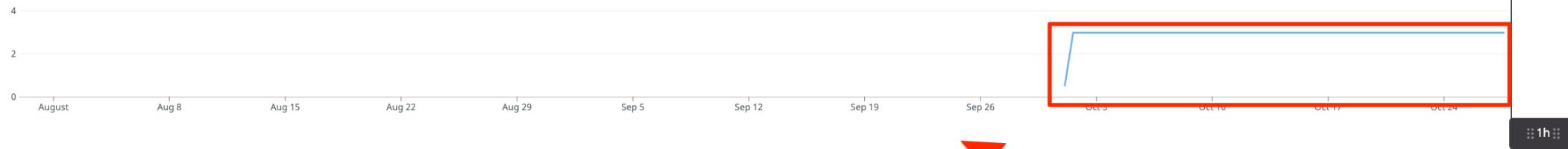
PRIORITY P1 (Critical)

Status & History

GROUP STATUS



HISTORY

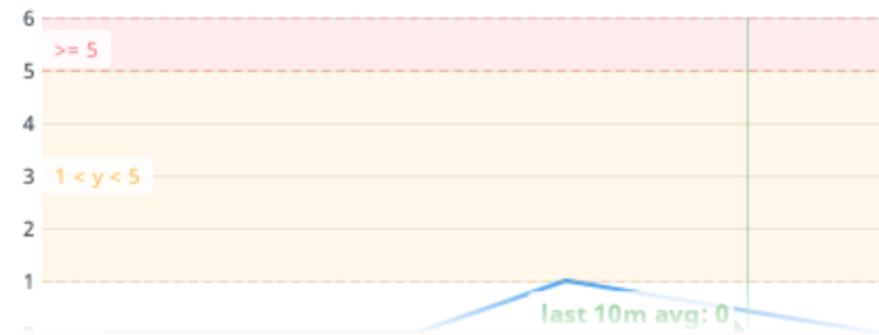


▽ Events & Watchdog

Events Watchdog



[P4] [Recovered] eDocs Signing Failed ?



[Tue Oct 19 2021 12:30:50 CEST \(9 hours and 12 minutes ago\)](#)

What should I do?

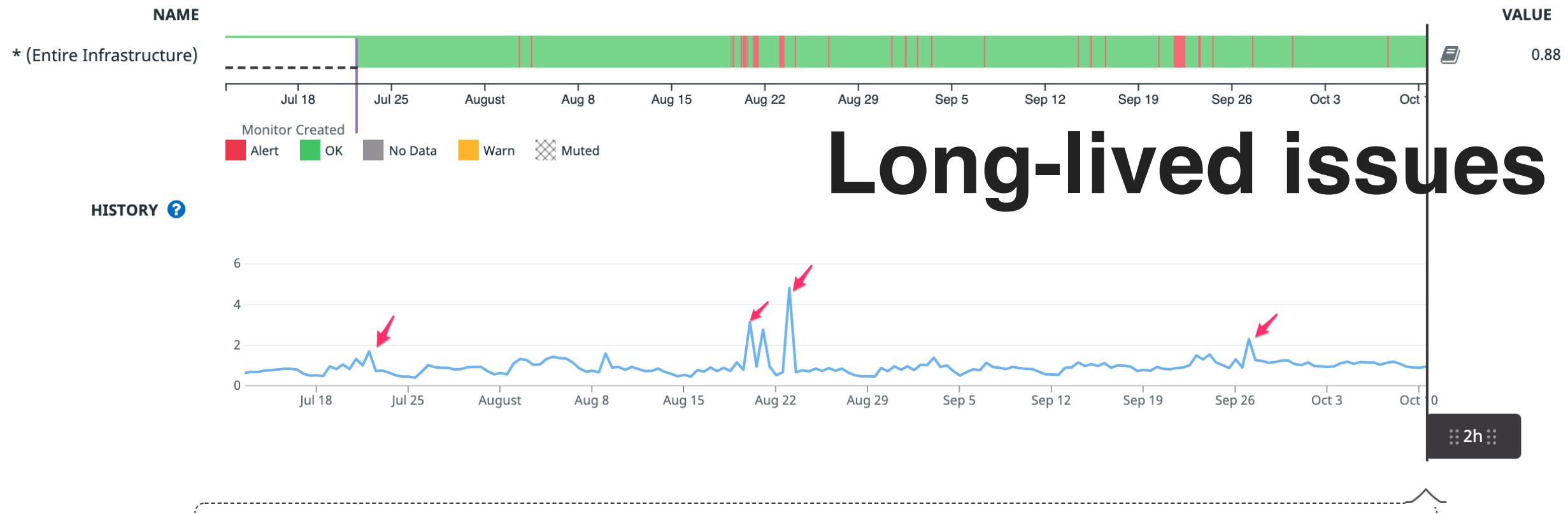
Warnings are evil

Make them failures or ignore them

*When that alert wakes me up I just snooze it
and keep sleeping*

OK Monitor status since 22 hours and 2 minutes ago (9 Oct, 16:13:59 CEST)

▼ Status & History



2h EVALUATION GRAPH

Each point in the graph reflects an evaluation window of 10min ?

Delete things you don't need!

Like, seriously

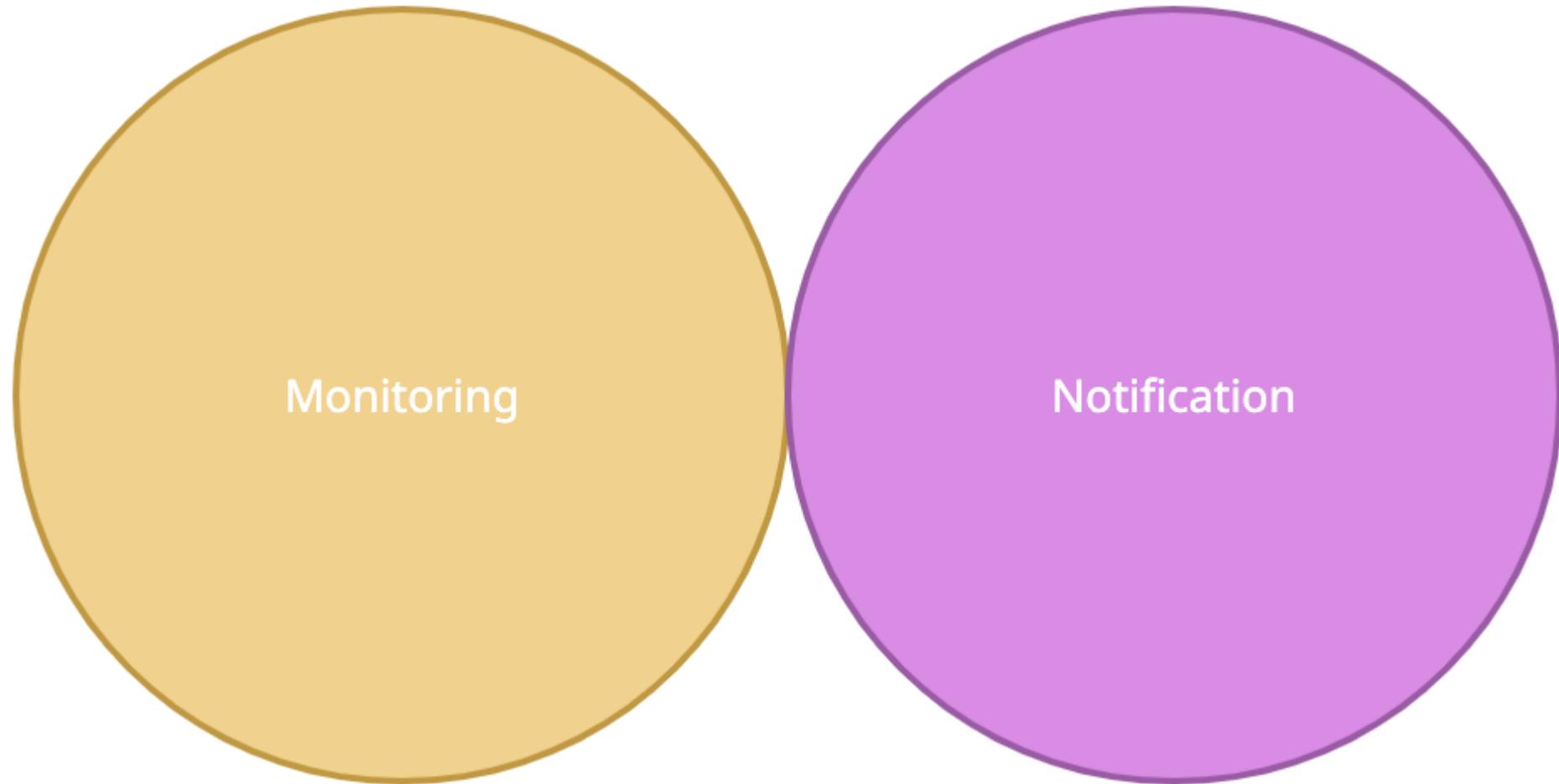
Learning 3



Reduce noise ruthlessly

Inadequate Tools

Crappy tools are extra infuriating at 3 in the morning



The market is full of high-quality tools



DATADOG

PagerDuty

Do you take good  for granted?

I've learned not to 

Project at a German client

Ticketing system without features such as

Auto closing

Auto closing

Auto grouping

Auto closing

Auto grouping

Conditional routing

**Auto closing
Auto grouping
Conditional routing**



Learning 4



Adopt tools that support you

Mismatched Tuning

Minimize
false *positives*
false *negatives*

Use ratios for your metrics

Properties

 Real User Monitoring

QUERY `rum("@type:error -@error.source:network @application.id: [REDACTED] env:prod").rollup("count").last("5m") >= 20`

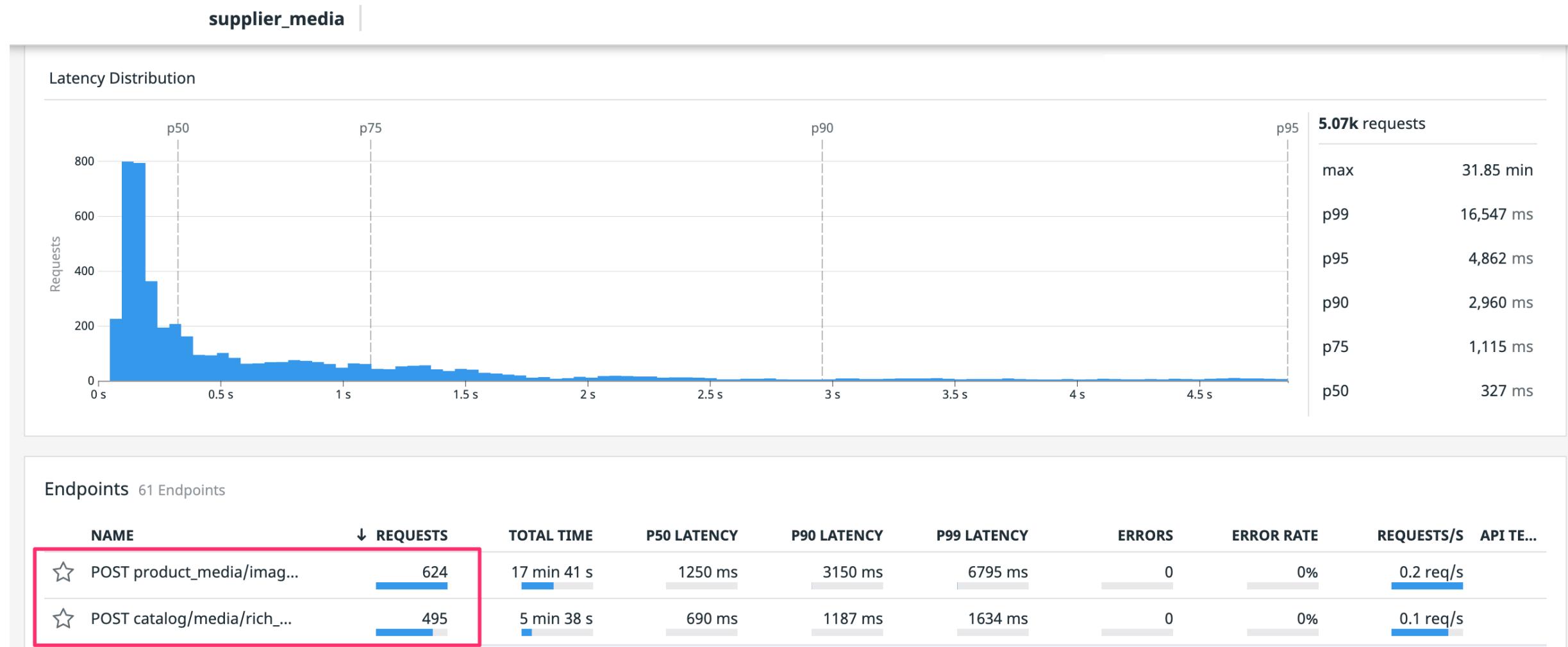
MESSAGE The App seeing high number of code errors

1 notified

TAGS  service:some-service

PRIORITY Not Defined

Get enough data



Multiple small alerts

<input type="checkbox"/>	PRIORITY	STATUS ↑	MUTED LEFT	NAME	
<input type="checkbox"/>	P1	OK		navigation-service prod / P1 - Error Rate SLO [99,9%]	2
<input type="checkbox"/>	P1	OK		navigation-service prod - Error Rate SLO [Short Window, 6x Burn rate]	
<input type="checkbox"/>	P1	OK		navigation-service prod - Error Rate SLO [Long Window, 14.4x Burn rate]	
<input type="checkbox"/>	P1	OK		navigation-service prod - Error Rate SLO [Short Window, 14.4x Burn rate]	
<input type="checkbox"/>	P1	OK		navigation-service prod - Error Rate SLO [Long Window, 6x Burn rate]	

Reliable alerts are hard!

hceris.com/monitoring-alerts-that-dont-suck/

Learning 5



Tune alerts often

Dysfunctions and Learnings Side by Side

Mixed abstraction levels

Lack of automation

Noisy alerts

Inadequate tools

Mismatched tuning

Use the right alert

Leverage automation at scale

Reduce noise ruthlessly

Adopt tools that support you

Tune alerts often

O'REILLY®

Building Secure & Reliable Systems

Best Practices for Designing, Implementing
and Maintaining Systems



Heather Adkins, Betsy Beyer,
Paul Blankinship, Piotr Lewandowski,
Ana Oprea & Adam Stubblefield

O'REILLY®

The Site Reliability Workbook

Practical Ways to Implement SRE

Companion to the
Bestselling SRE Book



Edited by Betsy Beyer,
Niall Richard Murphy, David K. Rensin,
Kent Kawahara & Stephen Thorne

O'REILLY®

Site Reliability Engineering

HOW GOOGLE RUNS PRODUCTION SYSTEMS

Edited by Betsy Beyer, Chris Jones,
Jennifer Petoff & Niall Murphy

Mario Fernandez
Staff Engineer
Wayfair

