

*please*

**Ask questions  
through the app**



*Rate Session*

*Thank you!*



# Four Years in – How Prometheus Revolutionized Monitoring at SoundCloud



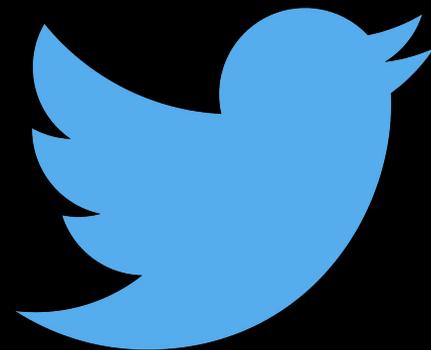
GOTO Berlin – 2017-11-16

Björn “Beorn” Rabenstein, Production Engineer, SoundCloud Ltd.



# Level 0

No monitoring.



## Community impact

Source	Amount
Emails	4
Tweets	3
Facebook posts	
Help Community	1 post (96 views)
Total	8 direct (visibility to over 100 users)

# Level 1

Nagios.

**Nagios<sup>®</sup>**

## Current Network Status

Last Updated: Fri Nov 10 11:24:13 GMT 2017 - Update is PAUSED [continue] 

Icinga 1.8.4 - Logged in as *admin*

- ▶ View Service Status Detail For All Hosts
- ▶ View Host Status Detail For All Hosts

### Commands for checked host(s)

Select command

Submit

### Host Status Details For Host/Services matching '10.33.84.91'

Page 1 of 1 Results: 50

Host	Status	Last Check	Duration	Attempt	Status Information	
ip-10-42-84-91	UP	2017-11-10 11:24:01	66d 17h 45m 39s	1/3	OK - 10.42.84.91: rta 0.098ms, lost 0%	<input type="checkbox"/>

Displaying Result 1 - 1 of 1 Matching Hosts

### Commands for checked services

Select command

Submit

### Service Status Details For Host/Services matching '10.33.84.91'

Host	Service	Status	Last Check	Duration	Attempt	Status Information	
ip-10-42-84-91	NRPE available	OK	2017-11-10 11:21:39	66d 21h 49m 50s	1/3	NRPE checks working	<input type="checkbox"/>
	check-statsd-local	OK	2017-11-10 11:19:59	142d 17h 23m 48s	1/3	OK: run: statsd: (pid 3786) 5780956s: run: log: (pid 3778) 5780956s	<input type="checkbox"/>
	nfs_mounts	OK	2017-11-10 11:20:42	388d 20h 25m 30s	1/3	NFS OK: 1 mount points avg of 0.00161 secs, max 0.00161 secs.	<input type="checkbox"/>

Page 1 of 1 Results: 50

Displaying Result 1 - 3 of 3 Matching Services

## Current Network Status

Last Updated: Fri Nov 10 11:28:51 GMT 2017 - Update in 13 seconds [pause]

Icinga 1.8.4 - Logged in as *admin*

- ▶ View Service Status Detail For All Hosts
- ▶ View Host Status Detail For All Hosts

### Commands for checked host(s)

Select command

Submit

### Host Status Details For Host/Services matching 'go-link-redirector'

Page 1 of 1 Results: 50

Host	Status	Last Check	Duration	Attempt	Status Information
go-link-redirector	PENDING	N/A	939d 21h 41m 38s	1/3	Host check scheduled for Fri Nov 10 11:29:33 GMT 2017

Displaying Result 1 - 1 of 1 Matching Hosts

### Commands for checked services

Select command

Submit

### Service Status Details For Host/Services matching 'go-link-redirector'

Page 1 of 1 Results: 50

Host	Service	Status	Last Check	Duration	Attempt	Status Information
go-link-redirector	health	OK	2017-11-10 11:27:45	226d 20h 15m 21s	1/3	HTTP OK: Status line output matched "200" - 1870 bytes in 0.019 second response time

Displaying Result 1 - 1 of 1 Matching Services

O'REILLY®



# Site Reliability Engineering

HOW GOOGLE RUNS PRODUCTION SYSTEMS

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

# Symptom-based alerting

Symptoms → pages. Causes → tickets.

Black-box monitoring FTW?

We combine *heavy use of white-box* monitoring with *modest but critical uses of black-box* monitoring. The simplest way to think about black-box monitoring versus white-box monitoring is that black-box monitoring is *symptom-oriented and represents active—not predicted—problems*. [...]

For paging, black-box monitoring has the key benefit of forcing discipline to only nag a human when a problem is both *already ongoing and contributing to real symptoms*. On the other hand, for *not-yet-occurring but imminent* problems, black-box monitoring is *fairly useless*.

# Why black-box monitoring is not sufficient:

1. Your probe is not real user traffic.
2. Long-tail latency matters.
3. You still need to investigate causes.

### Critical: Page Failure

[Availability Test for API Mobile Track \[107111\]](#)

<https://api-mobile.soundcloud.id=4711424248273487234>

TYPE	Web	REPORT	11/07/2017 09:30:00
MONITOR	Object	ALERT	11/07/2017 09:34:22

RUN THRESHOLD	1
TIME THRESHOLD	default

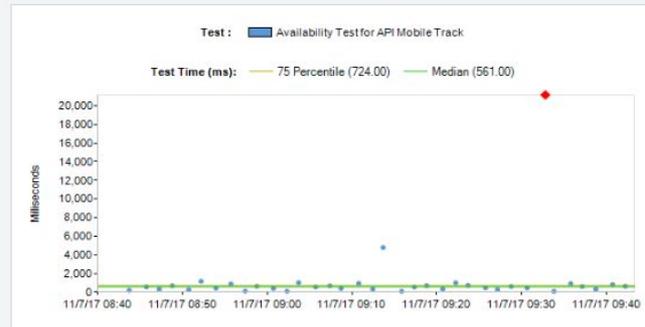
[Performance](#)
[Scatterplot](#)
[Waterfall](#)

### Alert Details

NODE	CAUSE	HOST IP
<b>Dallas - Level3</b> 4.30.68.66	Unknown	52.85.210.186
0 x 1 [2]		

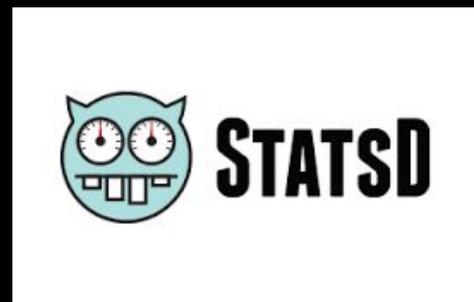
Triggered on 1 of 3 runs, expected 0 more runs in the alert interval.

### Analysis



# Level 2

StatsD and Graphite.



# Monitoring is many things...

- Observe. (“Shiny dashboards.”)
- Explore. (“What caused this outage?”)
- Alert. (“The machines wake me up when they need me.”)

# Container orchestration

Bazooka (R.I.P.) and...



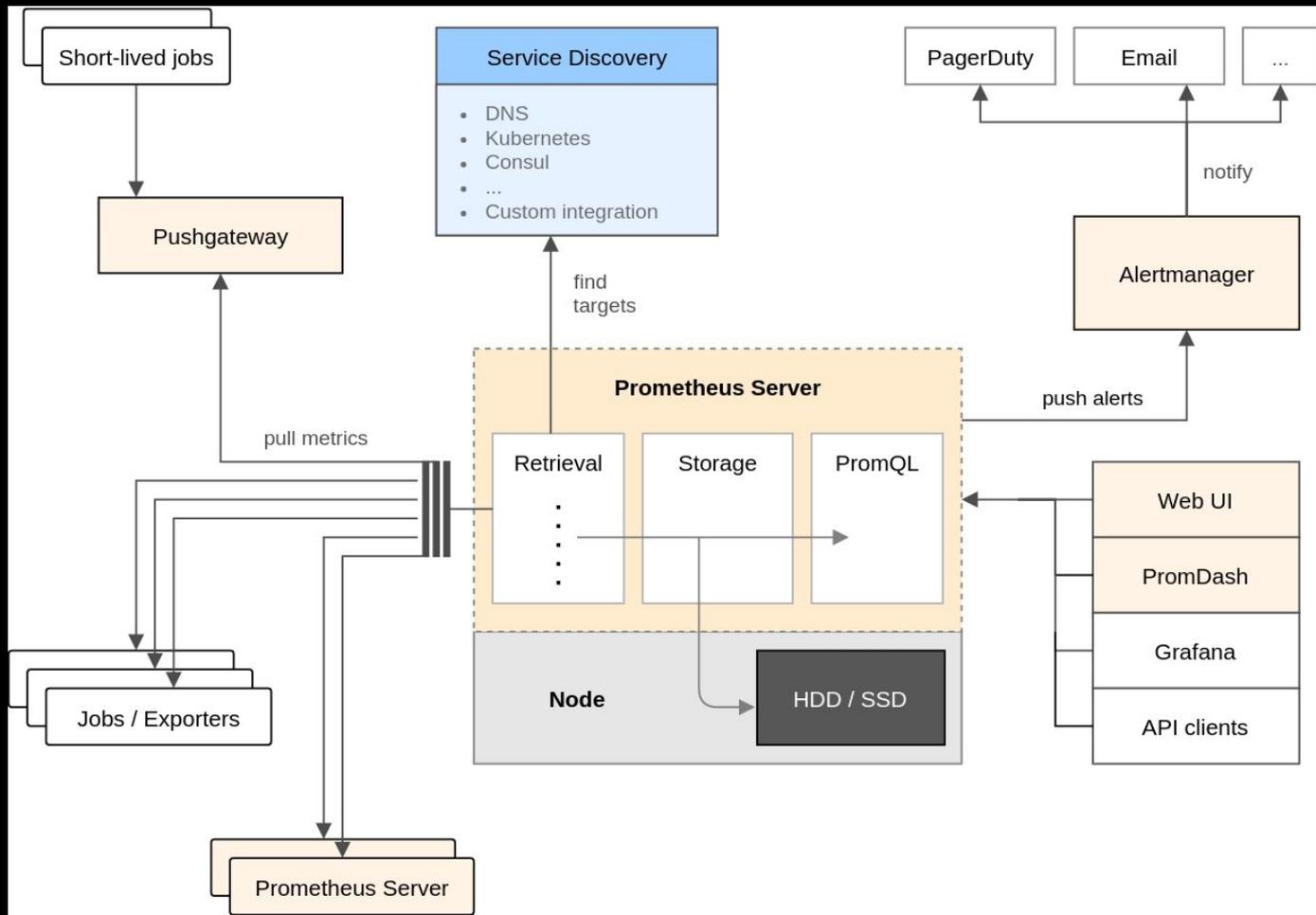
We need monitoring systems that allow us to alert for high-level service objectives, but retain the granularity to inspect individual components as needed.

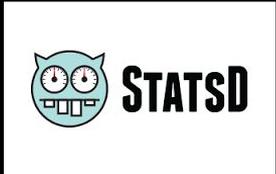
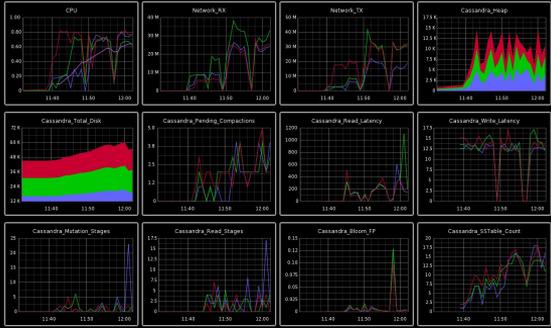
*Chapter 10: Practical Alerting from Time-Series Data*

# Level 3

Prometheus.







Application

Container

Orchestration

Host (OS, Hardware)

Network





Application

Container



Orchestration

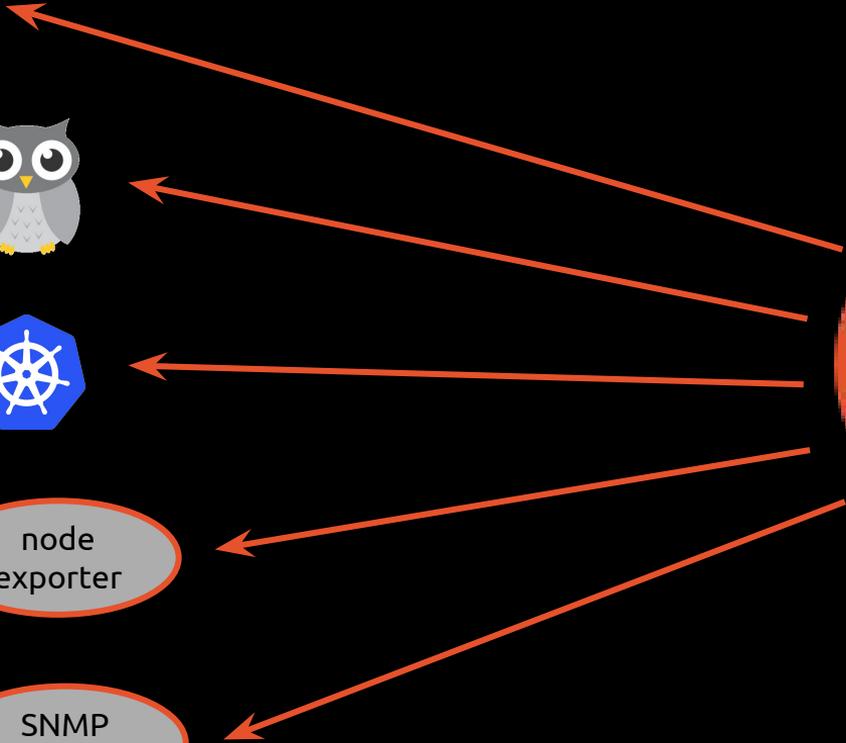


Host (OS, Hardware)

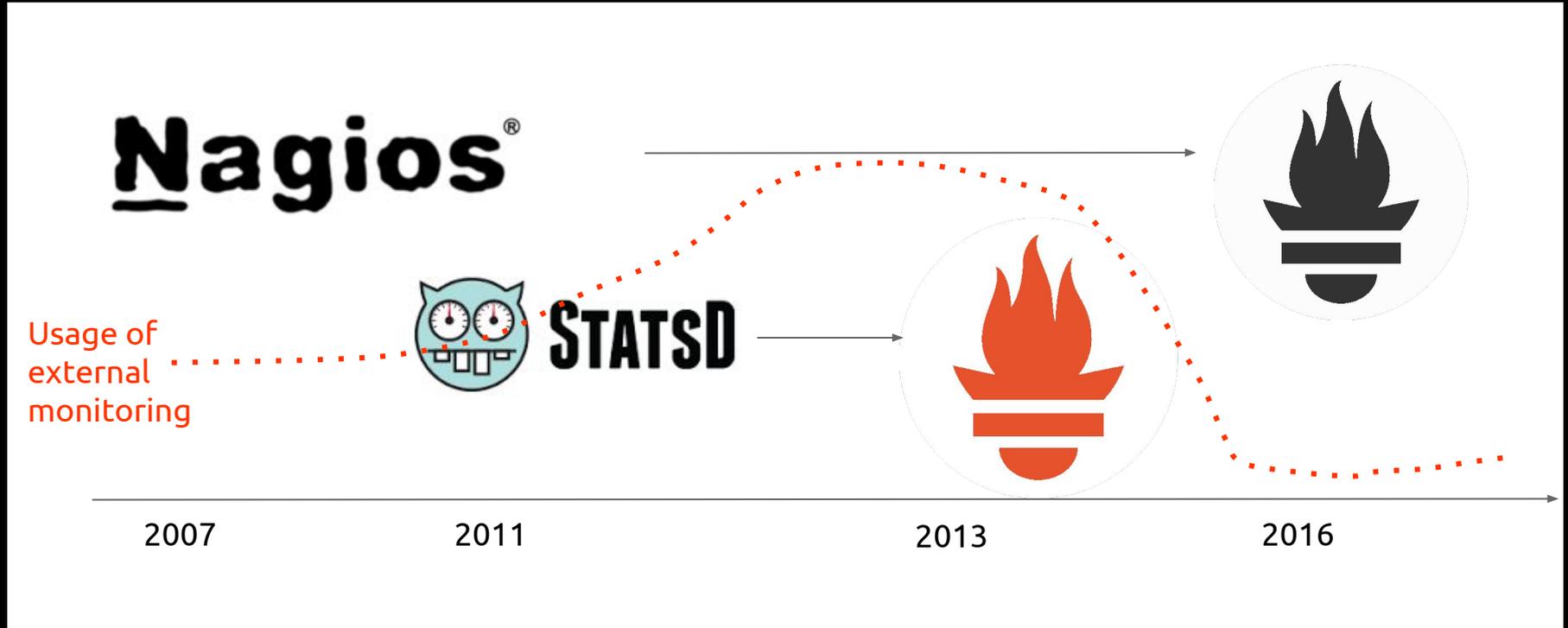
node exporter

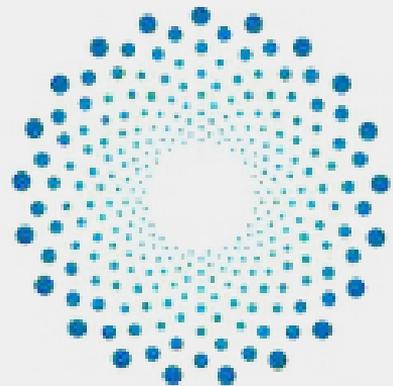
Network

SNMP exporter



# A word about external monitoring providers





**catchpoint®**

**L**

# Off-premises insight for your on-premises Prometheus servers.

Latency.at measures performance and availability of your sites and services from multiple global locations and provides the results as Prometheus metrics.

# Instrumentation

# JVMKit

```
func ExampleHandler(w http.ResponseWriter, r *http.Request) {
    status := http.StatusOK
    defer func(begun time.Time) {
        duration.Observe(time.Since(begun).Seconds())
        counter.With(prometheus.Labels{
            "status": fmt.Sprint(status),
        }).Inc()
    }(time.Now())
    result, err := computeResult()
    switch err {
    case nil:
        fmt.Fprintln(w, result)
        return
    case errTimeout:
        status = http.StatusServiceUnavailable
```



# Collection

```
'prometheus' => {  
  'storage' => {  
    'retention' => '360h',  
  },  
  'owner' => 'api-mobile',  
  'scrape_interval' => '30s',  
  'jobs' => [  
    {  
      'name' => 'api-mobile',  
      'k8s' => {  
        'cluster' => 'ab',  
        'system' => 'api-mobile',  
      },  
    },  
  ],  
  # ...  
}
```

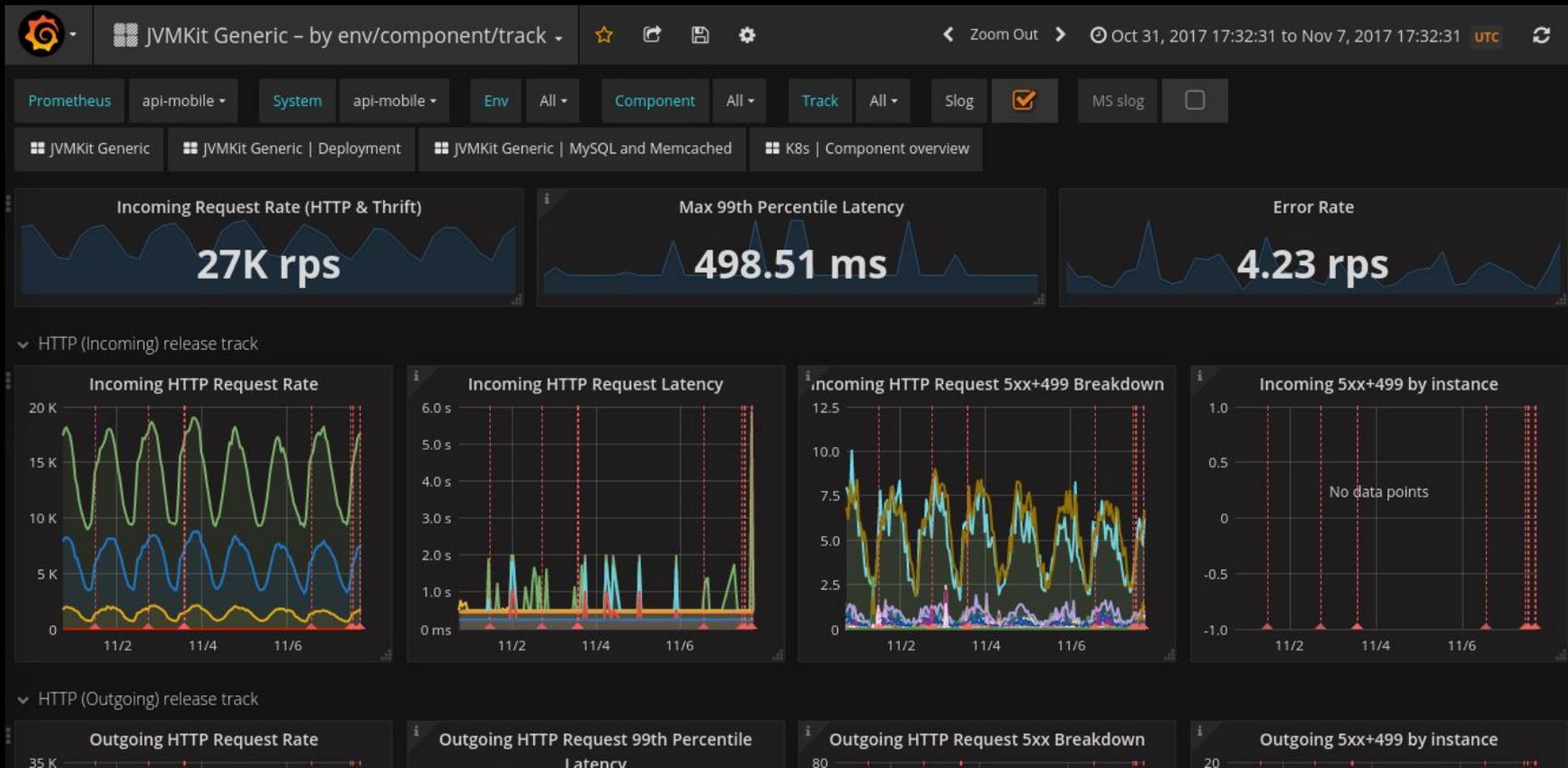
```
- job_name: api-mobile  
  scrape_interval: 30s  
  scrape_timeout: 10s  
  metrics_path: /metrics  
  scheme: http  
  kubernetes_sd_configs:  
  - api_server: https://api.k2.k8s.s-cloud.net  
    role: pod  
    tls_config:  
      ca_file: /mnt/prometheus/k8s-certificates/ab/ca.crt  
      cert_file: /mnt/prometheus/k8s-certificates/ab/client.crt  
      key_file: /mnt/prometheus/k8s-certificates/ab/client.key  
      insecure_skip_verify: false  
    namespaces:  
      names: []  
  relabel_configs:  
  - source_labels: [__meta_kubernetes_namespace, __meta_kubernetes_pod_annotation_prometheus_io_scheme]  
    separator: ;  
    regex: .+;(?:[0-9]+,?)+|;  
    replacement: $1  
    action: keep  
  - source_labels: [__meta_kubernetes_namespace, __meta_kubernetes_pod_label_system]  
    separator: ;  
    regex: .+;api-mobile|;  
    replacement: $1  
    action: keep  
  - source_labels: [__meta_kubernetes_namespace, __meta_kubernetes_pod_label_env]  
    separator: ;  
    regex: .+;production|;  
    replacement: $1  
    action: keep  
  - source_labels: [__meta_kubernetes_pod_annotation_prometheus_io_scheme]  
    separator: ;  
    regex: (https?)  
    target_label: __scheme__  
    replacement: $1
```

# Targets

## api-mobile (200/200 up)

Endpoint	State	Labels	Last Scrape	Error
http://10.146.40.231:9150/metrics	UP	<code>cluster="ab"</code> <code>component="memcached"</code> <code>env="production"</code> <code>instance="api-mobile-memcached-3708501969-8d414"</code> <code>namespace="api-mobile"</code> <code>system="api-mobile"</code> <code>track="release"</code> <code>version="14-3-3069a38"</code>	24.932s ago	
http://10.146.0.131:9150/metrics	UP	<code>cluster="ab"</code> <code>component="memcached"</code> <code>env="production"</code> <code>instance="api-mobile-memcached-3708501969-9v8d8"</code> <code>namespace="api-mobile"</code> <code>system="api-mobile"</code> <code>track="release"</code> <code>version="14-3-3069a38"</code>	3.785s ago	
http://10.145.220.181:9150/metrics	UP	<code>cluster="ab"</code> <code>component="memcached"</code> <code>env="production"</code> <code>instance="api-mobile-memcached-3708501969-cnss6"</code> <code>namespace="api-mobile"</code> <code>system="api-mobile"</code> <code>track="release"</code> <code>version="14-3-3069a38"</code>	27.397s ago	
http://10.144.16.234:9150/metrics	UP	<code>cluster="ab"</code> <code>component="api"</code> <code>env="production"</code> <code>instance="api-mobile-api-3708501969-dfnf0"</code> <code>namespace="api-mobile"</code> <code>system="api-mobile"</code> <code>track="release"</code> <code>version="14-3-3069a38"</code>	8.999s ago	

# Dashboards



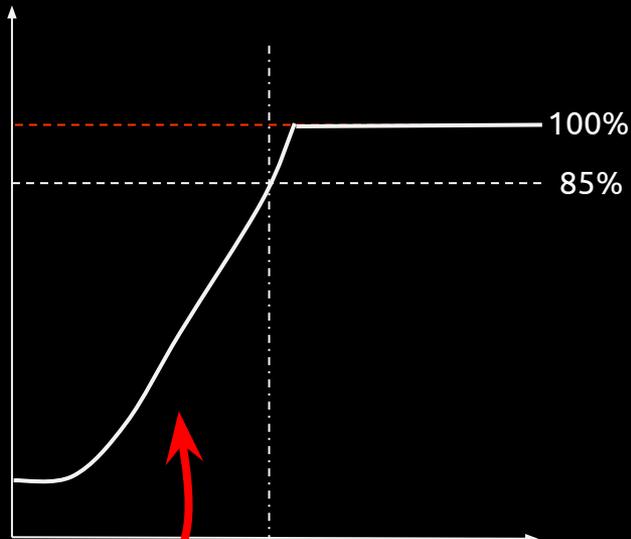


# Alert creation

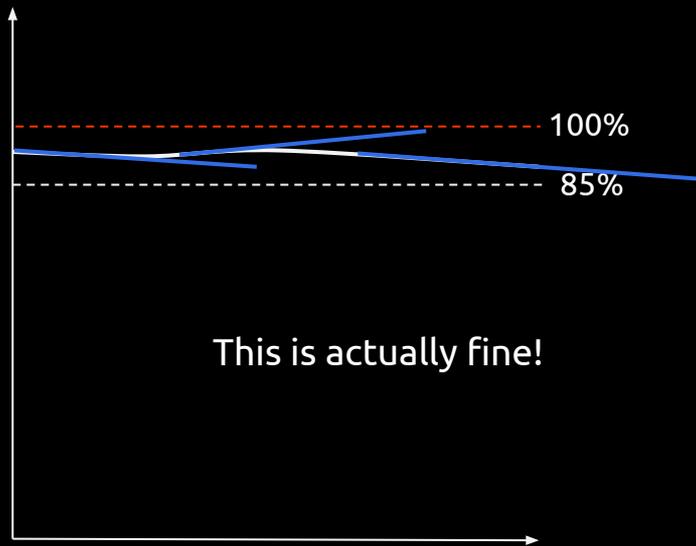
```
ALERT PrometheusRuleEvaluationSlow
  IF
    rate(prometheus_evaluator_iterations_missed_total[1h])
    /
    rate(prometheus_evaluator_iterations_total[1h])
    * 100
    > 5
  FOR 4h
  LABELS {
    severity = "warning",
  }
  ANNOTATIONS {
    summary = "{{${labels.job}} is evaluating rules too slowly",
    description = "In the last hour, {{${labels.job}} at {{${labels.instance}} has ...",
    runbook = "http://doc/runbooks/prometheus/#prometheusruleevaluationslow",
    dashboard = 'http://grafana/dashboard/db/prometheus-srv?var-job={{${labels.job}}...',
  }
```



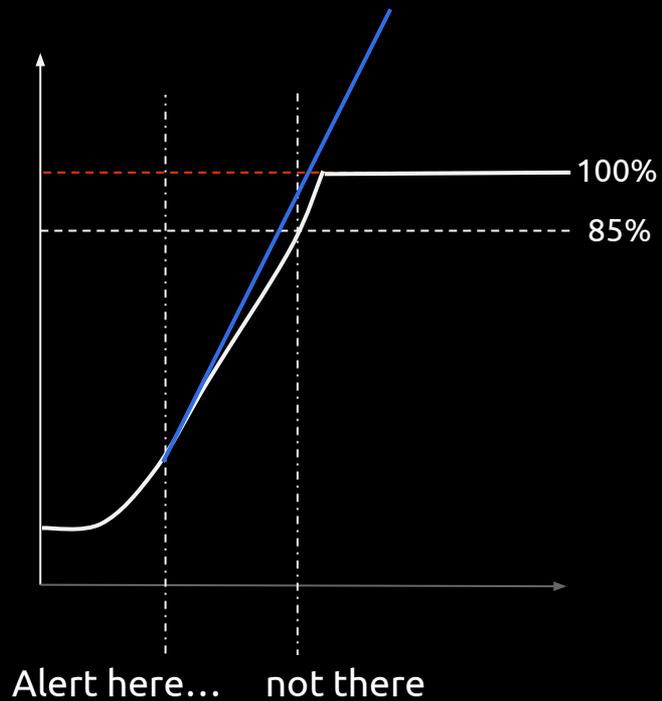
This is fine!?!



Alert here.



This is actually fine!



```
ALERT NodeFilesystemSpaceFillingUp
  IF
    predict_linear(node_filesystem_avail{fstype=~"ext."}[1h], 4*60*60) < 0
  AND
    node_filesystem_avail{fstype=~"ext."} / node_filesystem_size{fstype=~"ext."} < 0.2
  AND
    node_filesystem_readonly{fstype=~"ext."} == 0
  FOR 10m
  LABELS {
    severity = "critical",
  }
  ANNOTATIONS {
    summary = "Filesystem space is filling up",
    description = 'Filesystem on {{$labels.device}} at {{ $labels.instance }} is \
      predicted to run out of space within the next 4 hours.',
    runbook = "http://doc/runbooks/node/#nodefilesystemspacefillingup",
    dashboard = 'http://grafana/dashboard/db/node?var-node={{ $labels.instance ...}',
  }
```



# Alert delivery

pagerduty Incidents Alerts Configuration Analytics NEW ?

INCIDENTS > INCIDENT #230307

**[FIRING:1] NodeFilesystemFilesFillingUp node db (/dev/disk/by-uuid/04fd58f2-675c-4fba-9cfb-726774219b8c xfs ip-10-42-88-55.ab3.nyc5.s-cloud.net:7700 node / data critical)** [Edit](#)

STATUS **Resolved** DURATION **00h 00m**

More Actions ▾

**Status** Resolved **Urgency** High

**Incident Times** Open from Nov 13, 2017 at 1:12 PM to Nov 13, 2017 at 1:12 PM (for seconds) **Incident Key** ffb802f4094aaf3a64cd1877e7f55d36ae89cf16ef39b0db8273d5b151

**Impacted Service** [Alertmanager Data Team](#) **Integration** Generic API

**Responders** 0 **Notes** 0

**Details** **Timeline**

**Name**

[FIRING:1] NodeFilesystemFilesFillingUp node db (/dev/disk/by-uuid/04fd58f2-675c-4fba-9cfb-726774219b8c xfs ip-10-42-88-55.ab3.nyc5.s-cloud.net:7700 node / data critical)

**Custom Details**

**firing** [HIDE DETAILS](#)

**Labels:**

- alertname = NodeFilesystemFilesFillingUp
- device = /dev/disk/by-uuid/04fd58f2-675c-4fba-9cfb-726774219b8c
- fstype = xfs
- instance = ip-10-42-88-55.ab3.nyc5.s-cloud.net:7700
- job = node
- mountpoint = /
- owner = data
- severity = critical

**Notes** 0

[+ Add Note](#)

There are no recent notes.



**AlertManager** APP 7:24 PM ☆

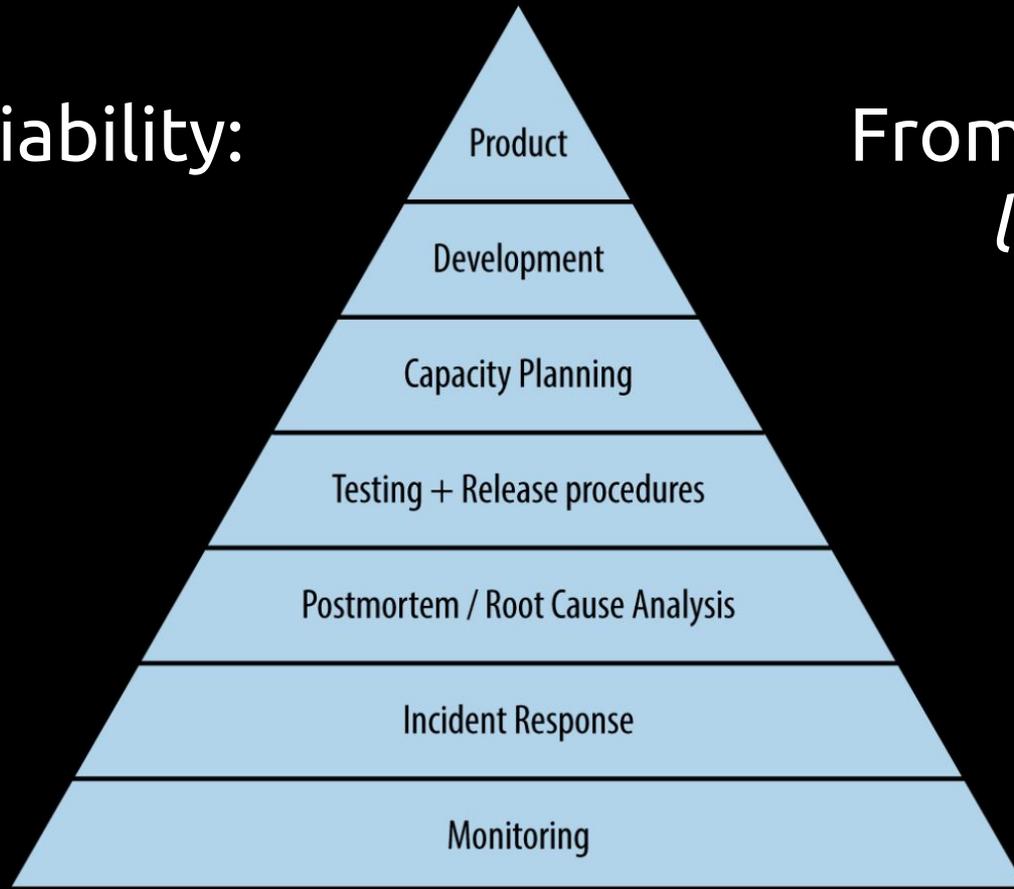


**[FIRING:1] prometheus-2 is evaluating rules too slowly prometheus-2 prometheus (PrometheusRuleEvaluationSlow k2 ip-10-64-70..5.s-cloud.net:9090 prodeng warning db)**

In the last hour, prometheus-2 at [ip-10-42-64-70.ab3.nyc5.s-cloud.net:9090](#) has skipped 11.0% of rule evaluation cycles because previous cycles took too long.

[Runbook Dashboard Source](#)

Service reliability:



From *unknown* to  
*low to high* in  
four years.

Dickerson's hierarchy of service reliability

*Site Reliability Engineering – How Google Runs Production Systems*, B. Beyer et al. (ed.), O'Reilly 2016, p104

Slides will be linked at

<https://github.com/beorn7/talks>

*Please*

**Remember to  
rate this session**

*Thank you!*

