



Monitoring Kubernetes with Prometheus

Tom Wilkie, July 2017



Tom Wilkie *VP Product, Grafana Labs*

Previously: Kausal, Weaveworks, Google, Acunu, Xensource

Twitter: @tom_wilkie Email: tom@grafana.com

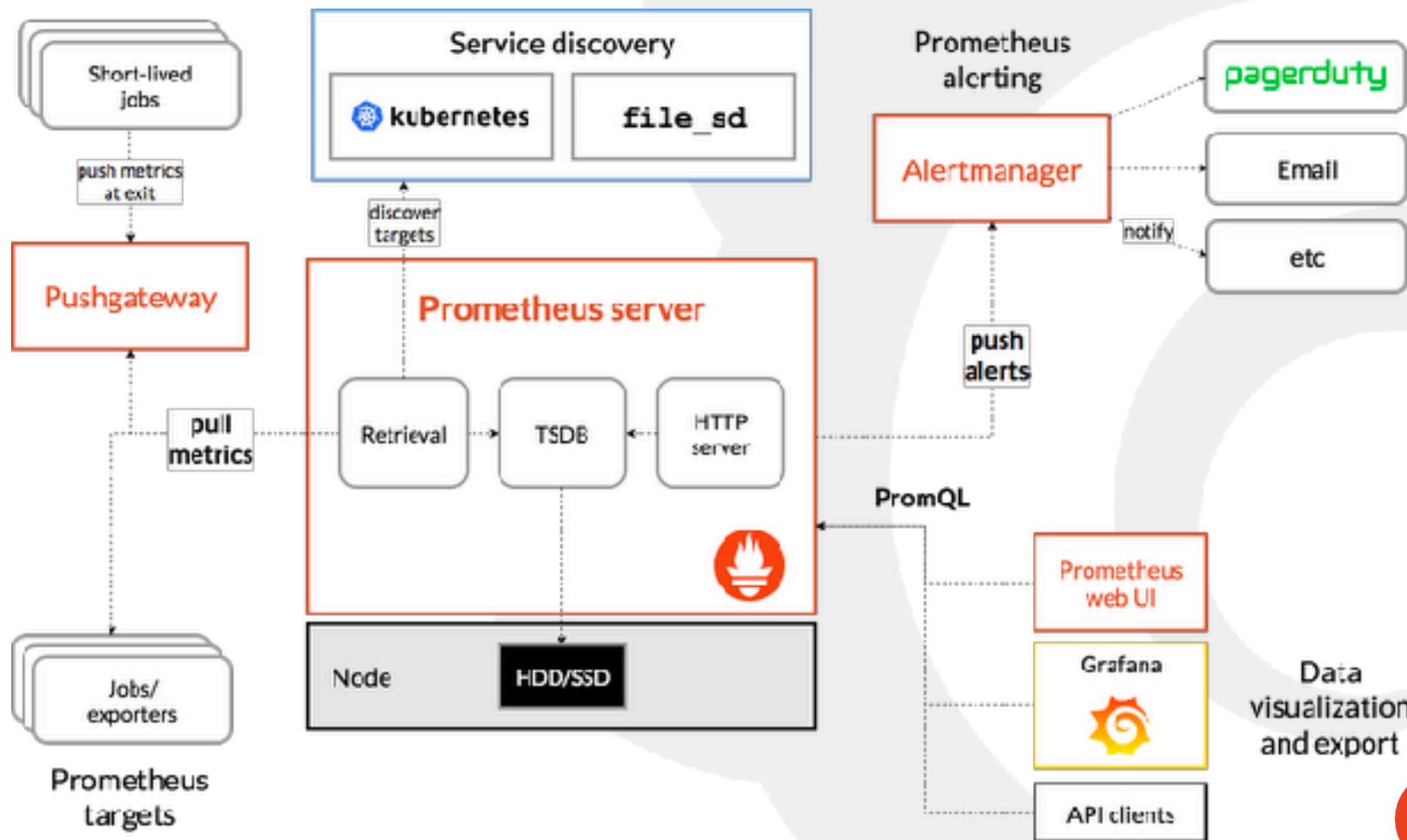


Prometheus Kubernetes Monitoring & Alerting Getting Started

Prometheus

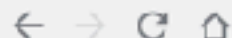


- A monitoring & alerting system.
 - Inspired by Google's BorgMon
 - Originally built by SoundCloud in 2012
 - Open Source, now part of the CNCF
-
- Simple text-based metrics format
 - Multidimensional datamodel
 - Rich, concise query language





localhost:8080/metrics



localhost:8080/metrics



```
07513
# HELP cortex_ring_ingester_ownership_percent The percent ownership of the ring by ingester
# TYPE cortex_ring_ingester_ownership_percent gauge
cortex_ring_ingester_ownership_percent{ingester="ingester-7677ff7756-6nm68"} 0.139281139974315
cortex_ring_ingester_ownership_percent{ingester="ingester-7677ff7756-7hq5d"} 0.14205375247217103
cortex_ring_ingester_ownership_percent{ingester="ingester-7677ff7756-845ng"} 0.1458392255347779
cortex_ring_ingester_ownership_percent{ingester="ingester-7677ff7756-8h25r"} 0.13680296045187929
cortex_ring_ingester_ownership_percent{ingester="ingester-7677ff7756-bpxc4"} 0.14248830665426523
cortex_ring_ingester_ownership_percent{ingester="ingester-7677ff7756-mlktg"} 0.1536875248778815
cortex_ring_ingester_ownership_percent{ingester="ingester-7677ff7756-rmp28"} 0.13984709003471002
# HELP cortex_ring_ingesters Number of ingesters in the ring
# TYPE cortex_ring_ingesters gauge
cortex_ring_ingesters{state="ACTIVE"} 7
cortex_ring_ingesters{state="JOINING"} 0
cortex_ring_ingesters{state="LEAVING"} 0
cortex_ring_ingesters{state="PENDING"} 0
cortex_ring_ingesters{state="Unhealthy"} 0
# HELP cortex_ring_tokens Number of tokens in the ring
# TYPE cortex_ring_tokens gauge
cortex_ring_tokens 3584
# HELP go_gc_duration_seconds A summary of the GC invocation durations.
# TYPE go_gc_duration_seconds summary
go_gc_duration_seconds{quantile="0"} 0.000143842
go_gc_duration_seconds{quantile="0.25"} 0.000270938
go_gc_duration_seconds{quantile="0.5"} 0.000337269
go_gc_duration_seconds{quantile="0.75"} 0.000486132
```



Prometheus' data model is very simple:

`<identifier> → [(t0, v0), (t1, v1), ...]`

Timestamps are millisecond int64, values are float64



Prometheus identifiers

```
http_requests_total{job="nginx", instances="1.2.3.4:80", path="/home", status="200"}  
http_requests_total{job="nginx", instances="1.2.3.4:80", path="/home", status="500"}  
http_requests_total{job="nginx", instances="1.2.3.4:80", path="/settings", status="200"}  
http_requests_total{job="nginx", instances="1.2.3.4:80", path="/settings", status="502"}
```

Prometheus series selector

```
http_requests_total{job="nginx", status=~"5.."}
```



Building queries usually starts with a selector

PromQL: `http_requests_total{job="nginx", status=~"5.."}`

`{job="nginx", instances="1.2.3.4:80", path="/home", status="500"}` 34

`{job="nginx", instances="1.2.3.4:80", path="/settings", status="502"}` 56

`{job="nginx", instances="2.3.4.5:80", path="/home", status="500"}` 76

`{job="nginx", instances="2.3.4.5:80", path="/setting", status="502"}` 96

...



Can select vectors of values...

PromQL: `http_requests_total{job="nginx", status=~"502"} [1m]`

`{job="nginx", instances="1.2.3.4:80", path="/home", status="500"} [30, 31, 32, 34]`

`{job="nginx", instances="1.2.3.4:80", path="/settings", status="500"} [4, 24, 56, 56]`

`{job="nginx", instances="2.3.4.5:80", path="/home", status="500"} [76, 76, 76, 76]`

`{job="nginx", instances="2.3.4.5:80", path="/setting", status="500"} [56, 106, 5, 96]`

...



And apply functions...

PromQL: **rate**(http_requests_total{job="nginx", status=~"502"}[1m])

```
{job="nginx", instances="1.2.3.4:80", path="/home", status="500"}      0.0666  
{job="nginx", instances="1.2.3.4:80", path="/settings", status="500"} 0.866  
{job="nginx", instances="2.3.4.5:80", path="/home", status="500"}      0.0  
{job="nginx", instances="2.3.4.5:80", path="/settings", status="500"} 2.43  
...
```



And aggregate by a dimension...

```
PromQL: sum by (path) (rate(http_requests_total{job="nginx", status=~"502"}[1m]))
```

```
{path="/home"}      0.0666
```

```
{path="/settings"} 3.3
```

```
...
```



Do binary operations...

```
PromQL: sum by (path) (rate(http_requests_total{job="nginx", status=~"502"}[1m]))  
/  
      sum by (path) (rate(http_requests_total{job="nginx"}[1m]))
```

```
{path="/home"}      0.001  
{path="/settings"} 1.0  
...
```

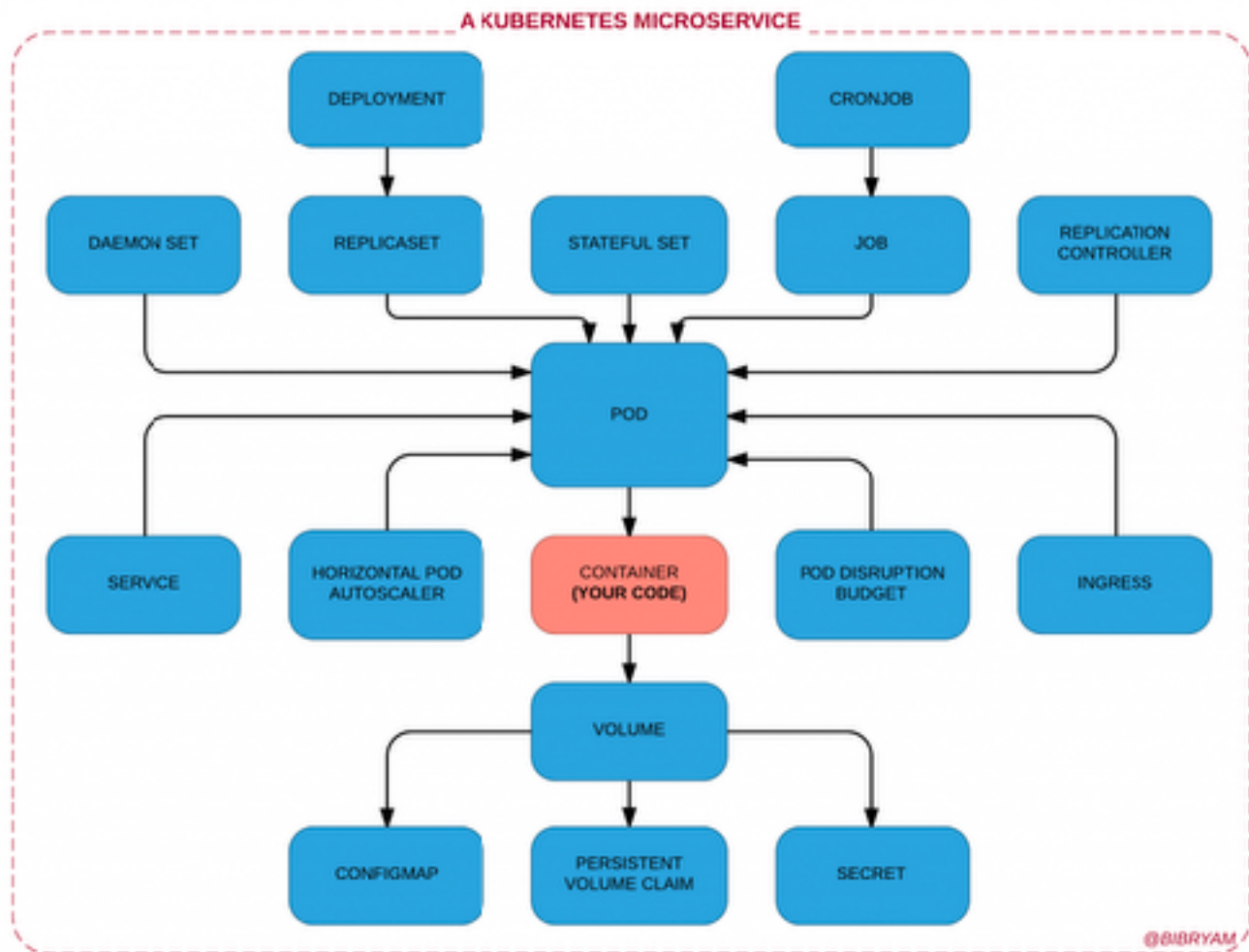


Kubernetes

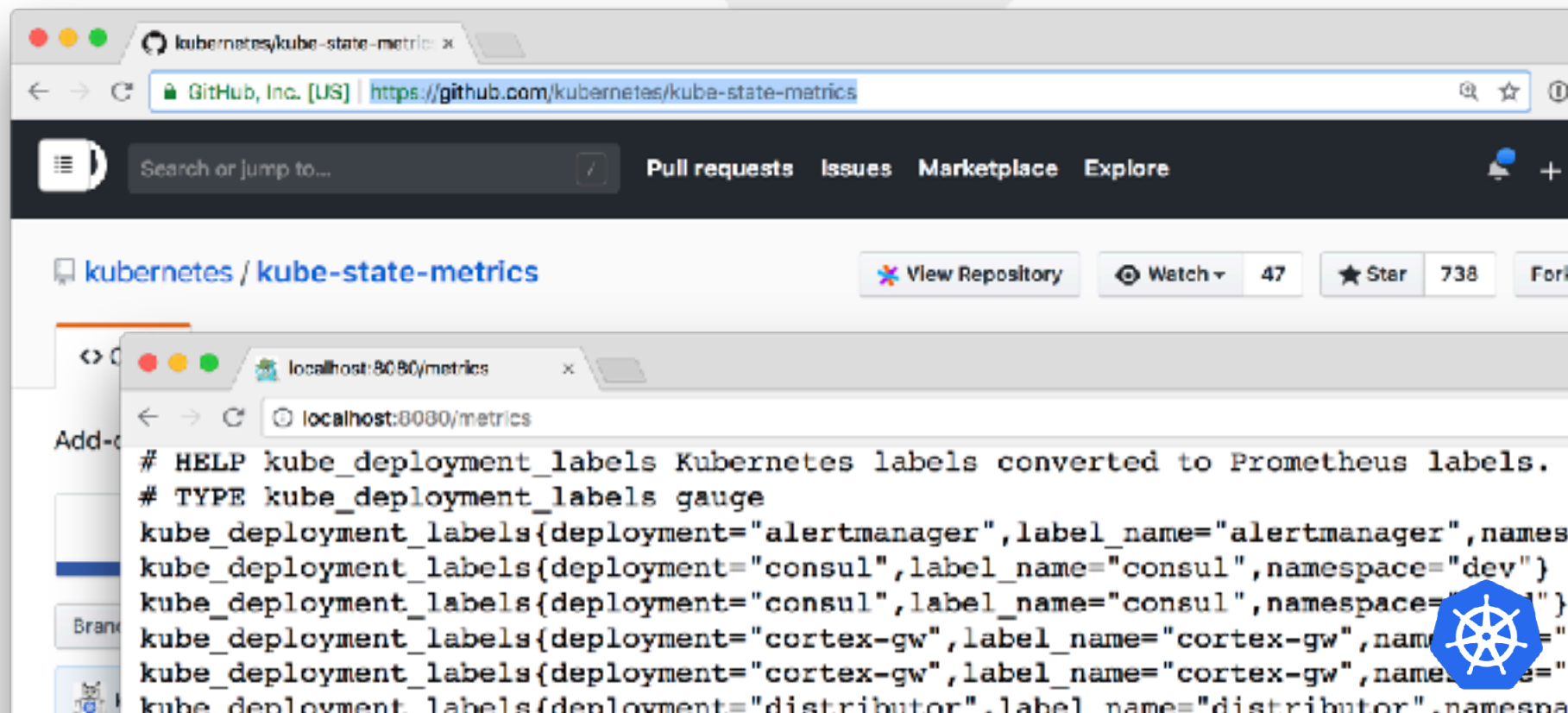


- Platform for managing containerized workloads and services
- “operating system for you datacenter”
- Inspired by Google’s Borg
- Also part of the CNCF
- Distributed, fault tolerant architecture
- Rich object model for you applications





kube-state-metrics



The image shows two overlapping browser windows. The background window displays the GitHub repository for `kubernetes/kube-state-metrics`. The address bar shows the URL `https://github.com/kubernetes/kube-state-metrics`. The repository page includes a search bar, navigation links for Pull requests, Issues, Marketplace, and Explore, and repository statistics: 47 watches, 738 stars, and a fork button. The foreground window shows a local Prometheus metrics endpoint at `localhost:8080/metrics`. The metrics output is in text format, starting with a help message and followed by several gauge metrics for Kubernetes deployments.

Background Window: GitHub Repository

Address bar: `https://github.com/kubernetes/kube-state-metrics`

Repository: `kubernetes / kube-state-metrics`

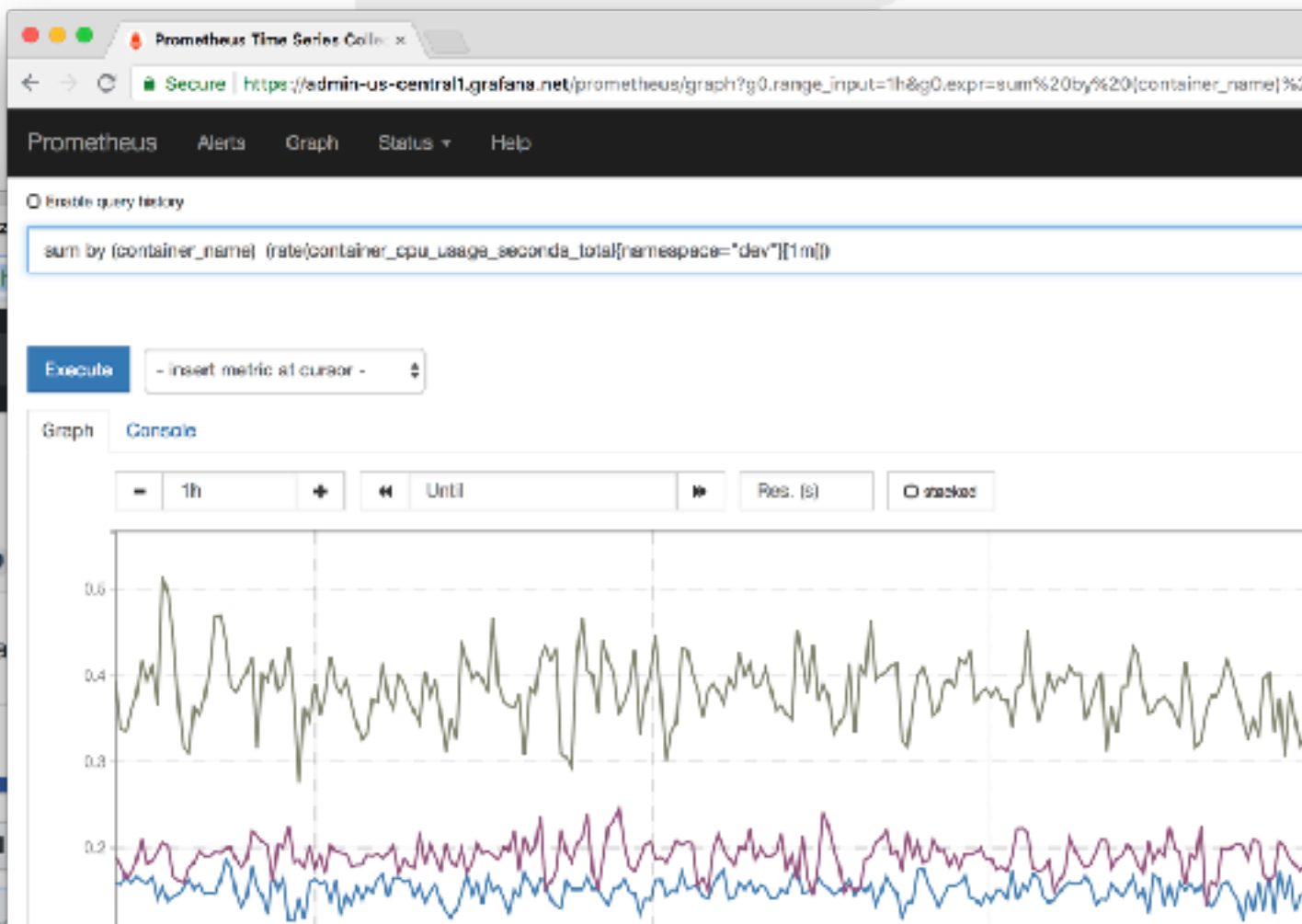
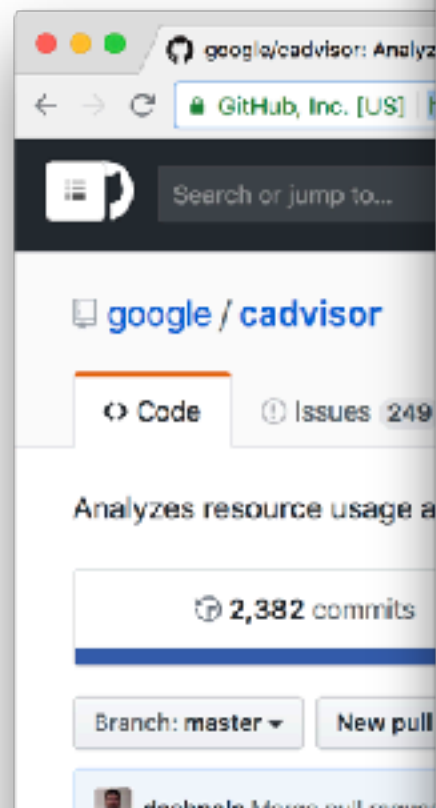
Statistics: 47 Watch, 738 Star, Fork

Foreground Window: Local Metrics

Address bar: `localhost:8080/metrics`

```
# HELP kube_deployment_labels Kubernetes labels converted to Prometheus labels.
# TYPE kube_deployment_labels gauge
kube_deployment_labels{deployment="alertmanager",label_name="alertmanager",namespace="default"} 1
kube_deployment_labels{deployment="consul",label_name="consul",namespace="dev"} 1
kube_deployment_labels{deployment="consul",label_name="consul",namespace="default"} 1
kube_deployment_labels{deployment="cortex-gw",label_name="cortex-gw",namespace="default"} 1
kube_deployment_labels{deployment="cortex-gw",label_name="cortex-gw",namespace="dev"} 1
kube_deployment_labels{deployment="distributor",label_name="distributor",namespace="default"} 1
```

cAdvisor



Monitoring & Alerting

What should I monitor?

USE Method

- Utilisation, Saturation, Errors...

RED Method

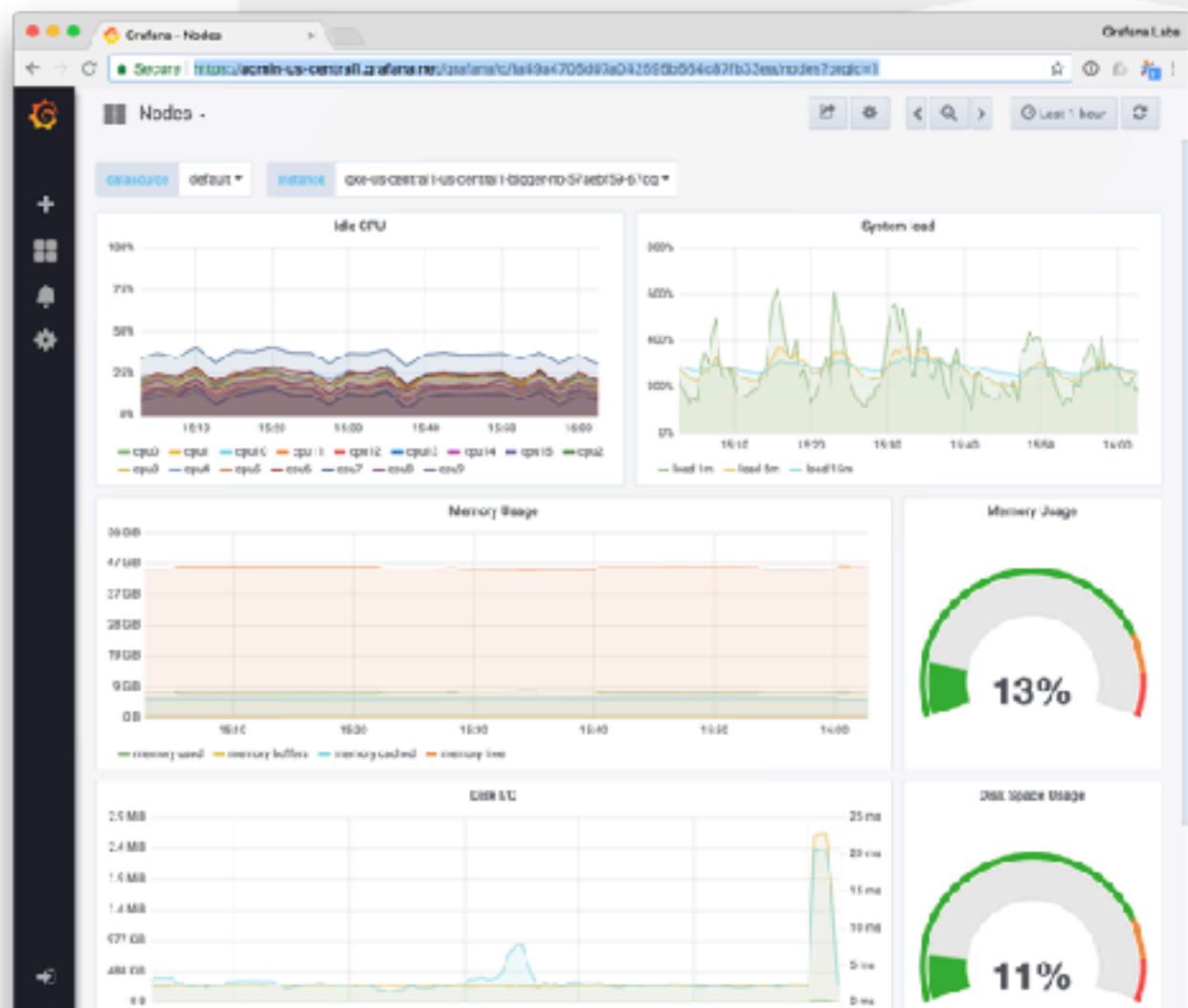
- Requests, Errors, Duration...

??? Method

- Expected system state...

USE Method

- cluster and node level metrics
- node_exporter run as a daemonset



USE Method

CPU Utilisation:

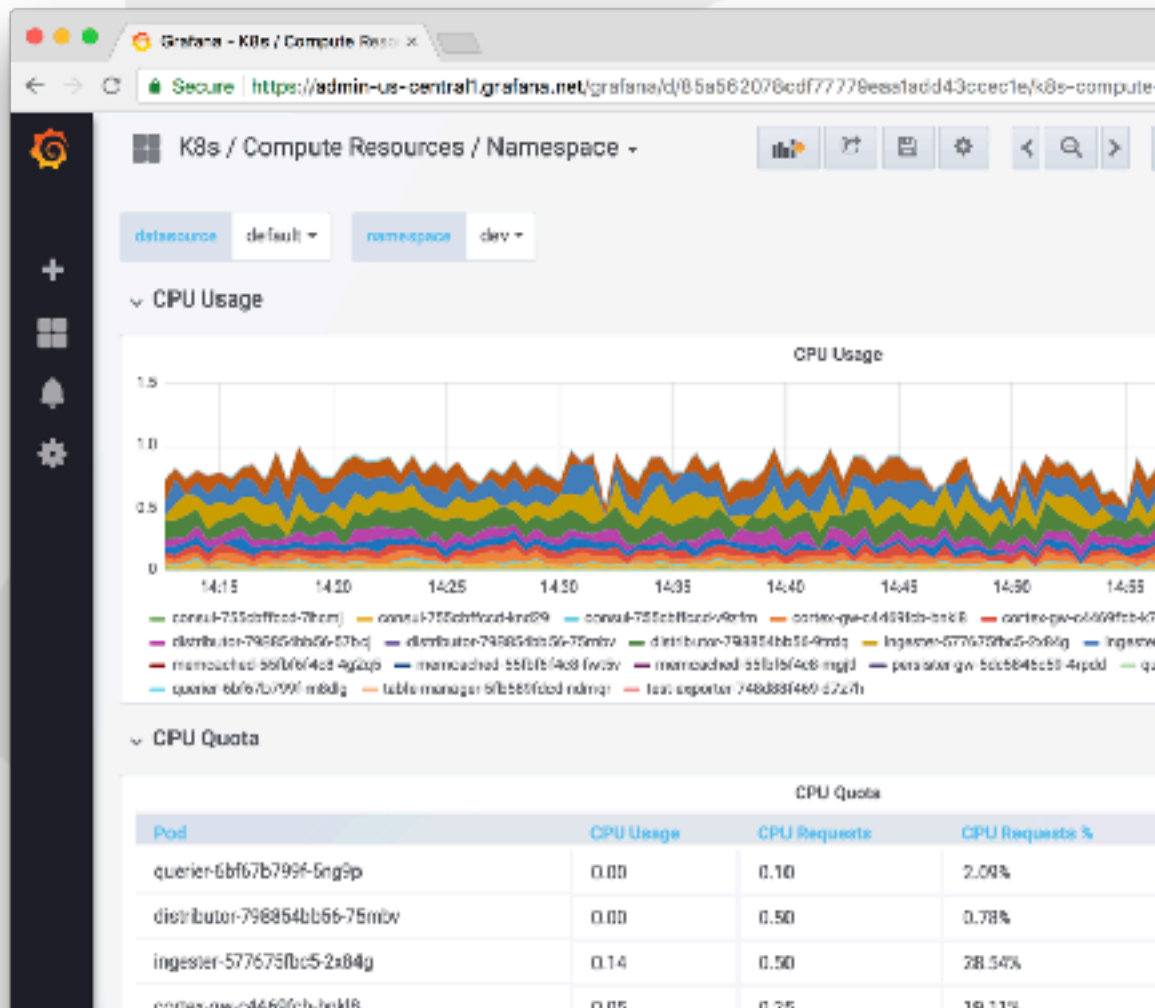
```
1 - avg(rate(node_cpu{mode="idle"}[1m]))
```

CPU Saturation:

```
sum(node_load1) / sum(node:num_cpu:sum)
```

USE Method

- Can also look at container level metrics from cAdvisor...
- ...and combine them with metadata from kube-state-metrics.



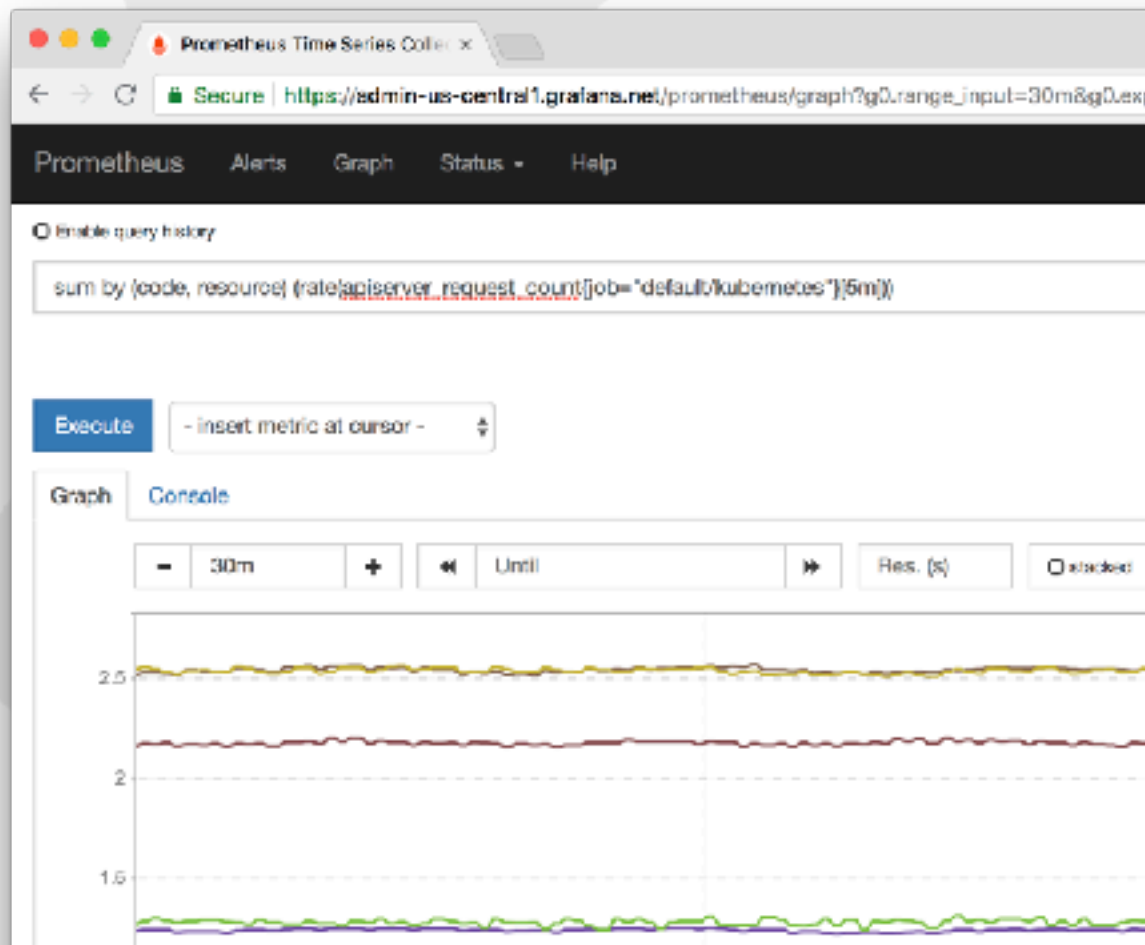
USE Method

Container CPU usage by “app” label

```
sum by (namespace, label_name) (  
    sum(rate(container_cpu_usage_seconds_total[5m])) by (pod_name)  
    * on (pod_name) group_left(label_name)  
    label_join(kube_pod_labels, "pod_name", ",", "pod")  
)
```


RED Method

- Metrics exposed by components for RED-style monitoring



RED Method

Most useful alert I've found:

```
100 * sum by(instance, job) (  
    rate(rest_client_requests_total{code!~"2.."} [5m])  
)  
/  
sum by(instance, job) (  
    rate(rest_client_requests_total[5m])  
)
```

??? Method

Alert expressions are invariants that describe a healthy system

```
kube_deployment_spec_replicas !=  
kube_deployment_status_replicas_available  
  
rate(kube_pod_container_status_restarts_total[15m]) > 0
```

??? Method

Alert expressions are invariants that describe a healthy system

```
(kube_pod_status_phase{phase!~"Running|Succeeded"}) > 0
```

```
sum(kube_pod_container_resource_requests_cpu_cores)
  / sum(node:node_num_cpu:sum)
  >
(count(node:node_num_cpu:sum) - 1)
  / count(node:node_num_cpu:sum)
```

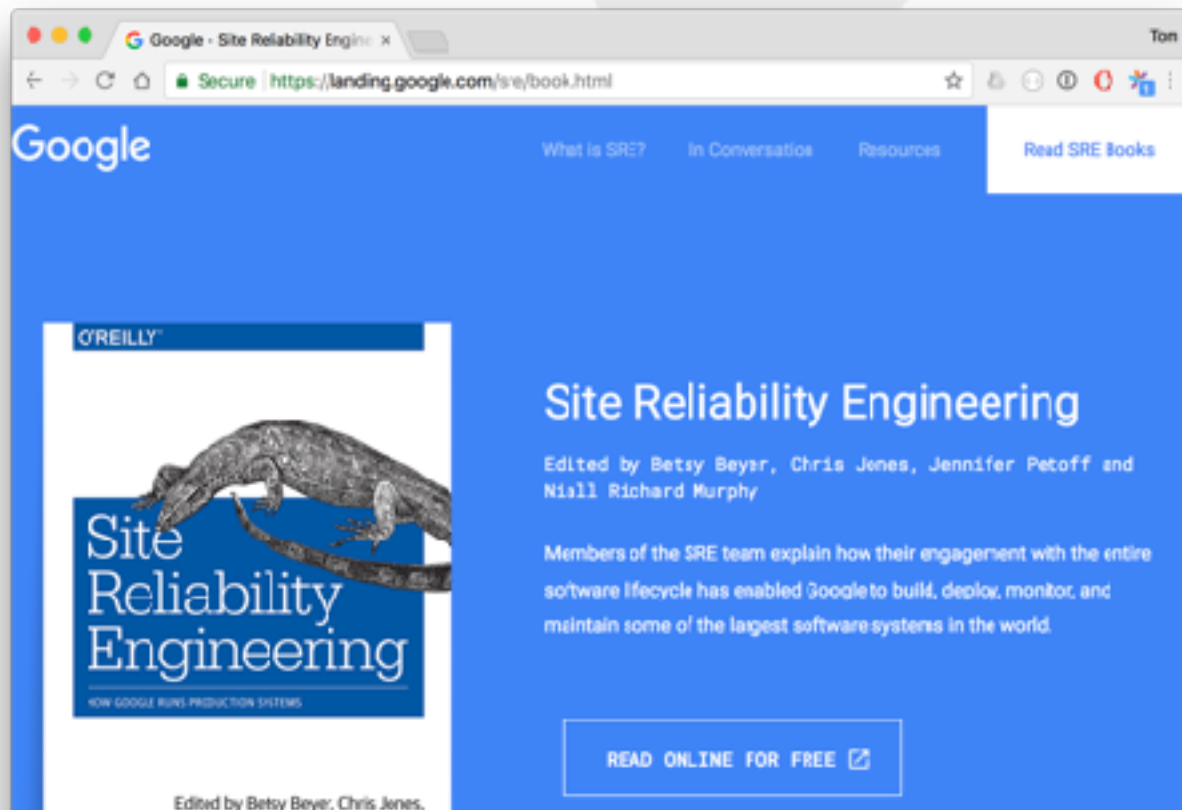
Getting Started

Getting setup

- github.com/coreos/prometheus-operator - Job to look after running Prometheus on Kubernetes
- github.com/coreos/kube-prometheus - Set of configs for running all there other things you need.
- github.com/kausalko/public/tree/master/prometheus-ksonnet - My configs for running Prometheus, Alertmanager, Grafana etc
- github.com/kubernetes-monitoring/kubernetes-mixin - Joint project to unify and improve common alerts for Kubernetes.

More reading...

<https://landing.google.com/sre/book.html>



The screenshot shows a web browser window with the Google Site Reliability Engineering (SRE) book landing page. The browser's address bar displays the URL <https://landing.google.com/sre/book.html>. The page has a blue header with the Google logo on the left and navigation links: "What is SRE?", "In Conversation", "Resources", and "Read SRE books". The main content area features the book cover for "Site Reliability Engineering: How Google Runs Production Systems" by O'Reilly, which includes an image of a gecko. To the right of the cover, the title "Site Reliability Engineering" is displayed in large white text, followed by the editors: "Edited by Betsy Beyer, Chris Jones, Jennifer Petoff and Niall Richard Murphy". Below this, a paragraph states: "Members of the SRE team explain how their engagement with the entire software lifecycle has enabled Google to build, deploy, monitor, and maintain some of the largest software systems in the world." At the bottom right, there is a white button with the text "READ ONLINE FOR FREE" and an external link icon.

Google


What is SRE? In Conversation Resources **Read SRE books**

O'REILLY

Site Reliability Engineering
HOW GOOGLE RUNS PRODUCTION SYSTEMS

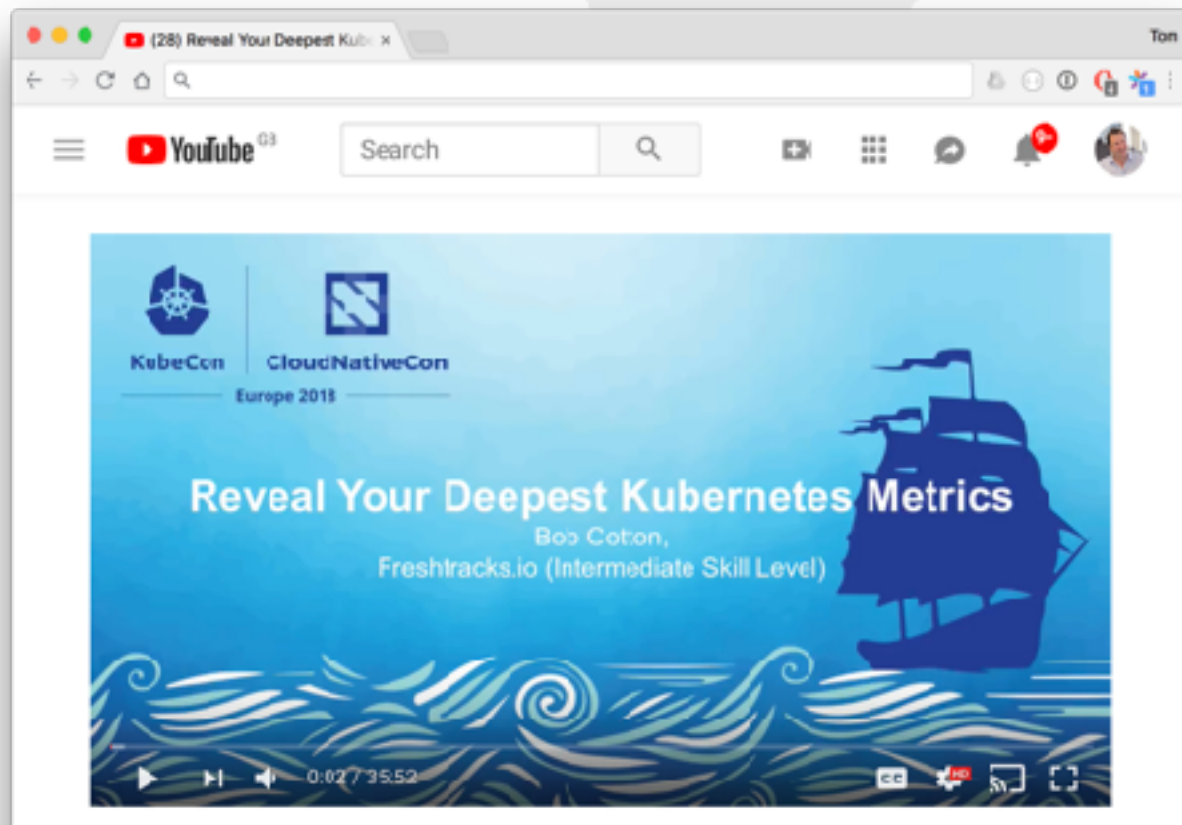
Edited by Betsy Beyer, Chris Jones, Jennifer Petoff and Niall Richard Murphy

Members of the SRE team explain how their engagement with the entire software lifecycle has enabled Google to build, deploy, monitor, and maintain some of the largest software systems in the world.

READ ONLINE FOR FREE 

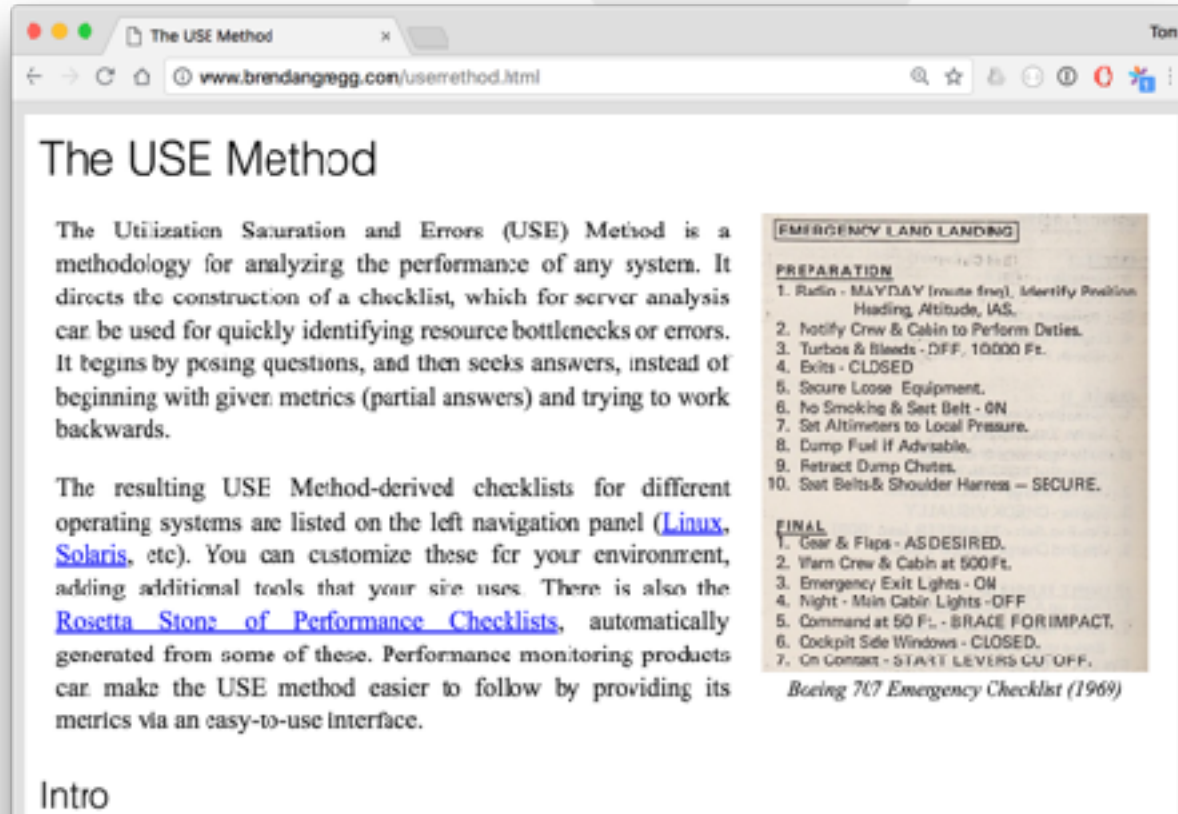
rafana Labs

<https://www.youtube.com/watch?v=1oJXMdVi0mM>



rafana Labs

<http://www.brendangregg.com/usemethod.html>



Questions?

