**Daniele Polencic — @danielepolencic@hachyderm.io**
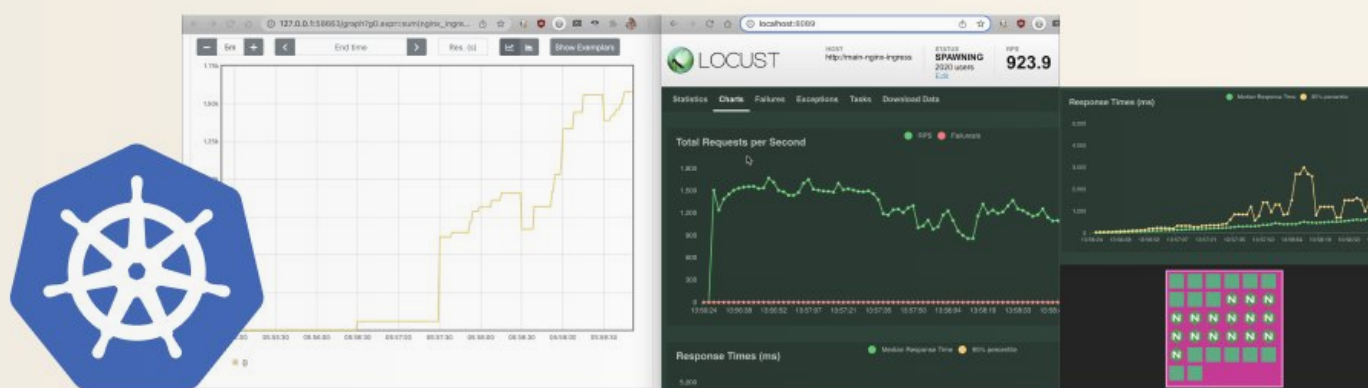@danielepolencic

How do you deal with peaks of traffic in Kubernetes?

You can use an autoscaler, but how should you configure and test it?

Let's dive into it.



AUTOSCALING
INGRESS CONTROLLERS
in KUBERNETES

8:09 PM · Apr 17, 2023

285 Likes

**Daniele Polencic — @danielepolencic@hachyderm.io**
@danielepolencic

1/

To autoscale the Ingress controller based on incoming requests, you need:

① Metrics (e.g. the requests per second)
② A metrics collector (to store the metrics)
③ An autoscaler (to act on the data)

**SCALING BASED ON NUMBER OF HTTP REQUESTS**

**1 Expose metrics** → prometheus exporters

**2 Collect & store**
→ prometheus
→ metrics server
→ KEDA

**3 Autoscaler**
→ horizontal pod autoscaler
→ vertical pod autoscaler
→ cluster autoscaler

2 Likes

**Daniele Polencic — @danielepolencic@hachyderm.io**
@danielepolencic

2/

Let's start with metrics

The nginx-ingress can be configured to expose Prometheus metrics

The official documentation has a page dedicated to it
kubernetes.github.io/ingress-nginx/...

You can use `nginx_connections_active` to count the number of active requests

**Stub status metrics**

| Name | Type | Description |
|------|------|-------------|
| nginx_connections_accepted | Counter | Accepted client connections. |
| nginx_connections_active | Gauge | Active client connections. |
| nginx_connections_handled | Counter | Handled client connections. |
| nginx_connections_reading | Gauge | Connections where NGINX is rea header. |
| nginx_connections_waiting | Gauge | Idle client connections. |
| nginx_connections_writing | Gauge | Connections where NGINX is wri back to the client. |

**①**

https://kubernetes.github.io/ingress-nginx/user-guide/monitoring/
https://github.com/nginxinc/nginx-prometheus-exporter

8:10 PM · Apr 17, 2023

3 Likes

3/

Next, you need a way to scrape the metrics

As you've already guessed, you can install Prometheus to do so

Since Nginx-ingress uses annotations for Prometheus, I installed the server without the Kubernetes operator

Community helm chart without CRDs

```
$ helm install prometheus prometheus-community/prometheus
NAME: prometheus
NAMESPACE: default
STATUS: deployed
REVISION: 1
TEST SUITE: None
```
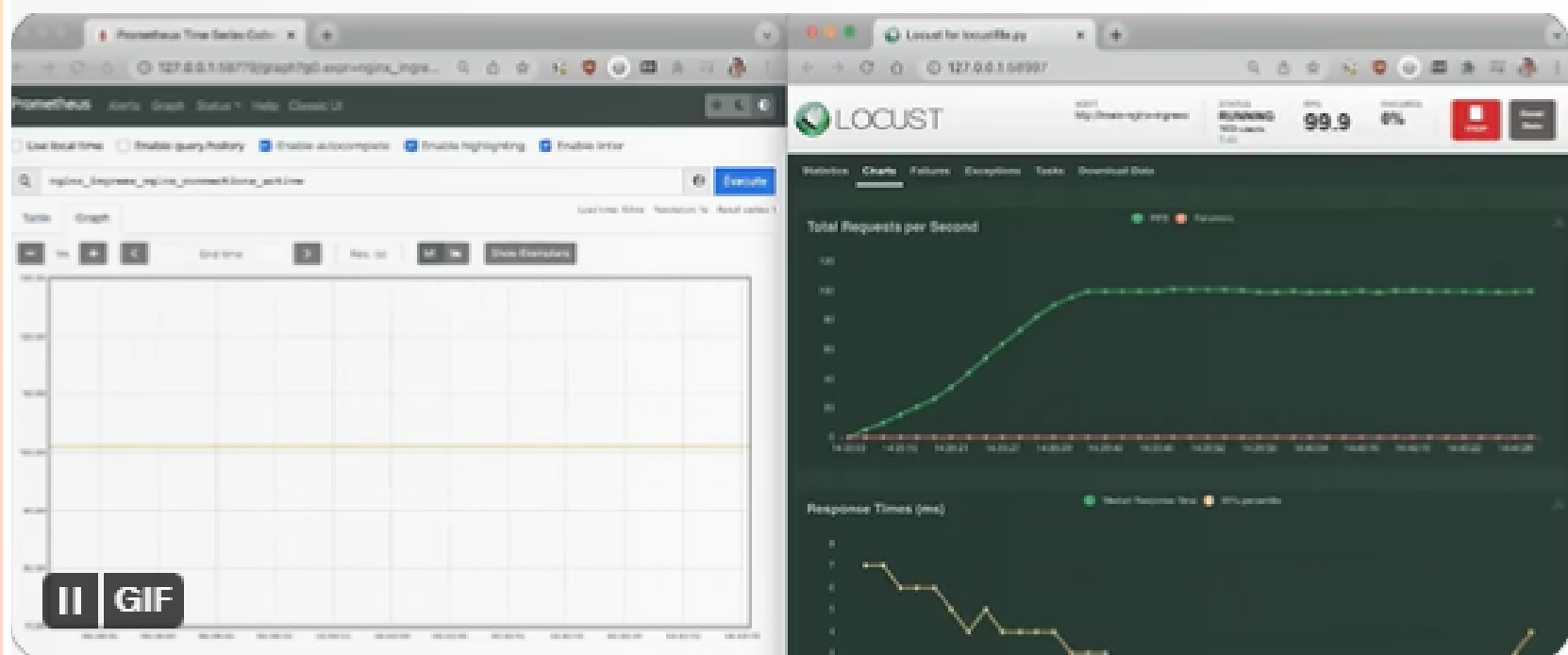
8:10 PM · Apr 17, 2023

1 Like

4/

I used Locust to generate some traffic to the Ingress to check that everything was running smoothly

With the Prometheus dashboard open, I checked that the metrics increased as more traffic hit the controller
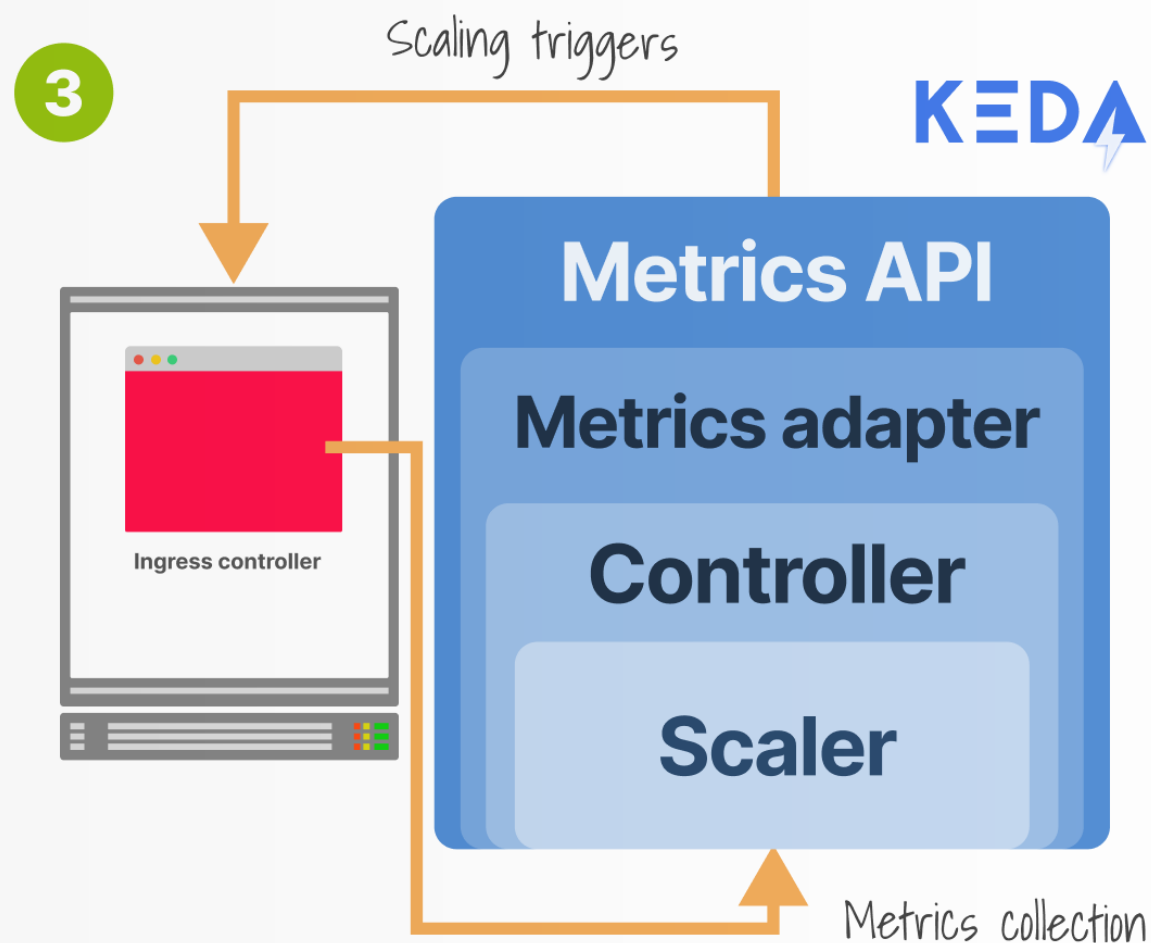


8:10 PM · Apr 17, 2023

4 Likes

5/

The last piece of the puzzle is the autoscaler

I decided to go with KEDA because:

① It's an autoscaler with a metrics server (so I don't need to install 2 different tools)
② It's easier to configure than the Prometheus adapter
③ I can use the HPA with PromQL

**Daniele Polencic — @danielepolencic@hachyderm.io**
@danielepolencic

6/

Once I installed KEDA, I only had to create a ScaledObject, configure the source of the metrics (Prometheus), and scale the Pods (with a PromQL query)

KEDA connects the dots and automatically creates the HPA for me

The ScaledObject creates the Horizontal Pod Autoscaler resource!

```yaml
apiVersion: keda.sh/v1alpha1
kind: ScaledObject
metadata:
 name: nginx-scale
spec:
 scaleTargetRef:
   kind: Deployment
   name: main-nginx-ingress
 minReplicaCount: 1
 maxReplicaCount: 20
 cooldownPeriod: 30
 pollingInterval: 1
 triggers:
 - type: prometheus
   metadata:
     serverAddress: http://prometheus-server
     metricName: nginx_connections_active_keda
     query: |
       sum(avg_over_time(nginx_ingress_nginx_connections_active{app="main-nginx-ingress"}[1m]))
     threshold: "100"
```

Ingress deployment

Prometheus scaler

PromQL

8:11 PM · Apr 17, 2023

3 Likes

**Daniele Polencic — @danielepolencic@hachyderm.io**
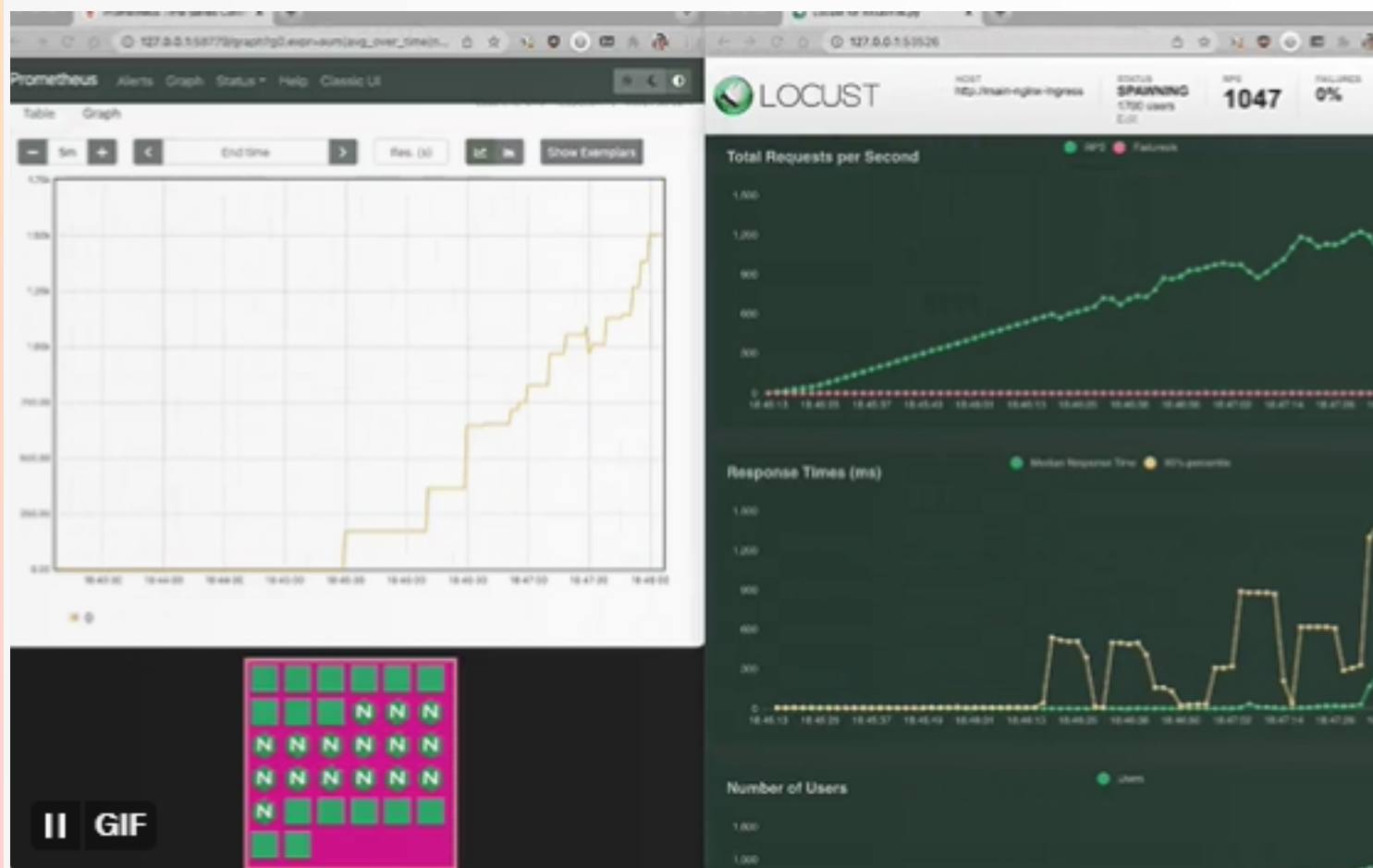@danielepolencic

7/

I repeated the tests with Locust and watched the replicas increase as more traffic hit the Nginx Ingress controller!

Can this pattern be extended to any other app?

Can you autoscale all microservices on the number of requests received?



8:11 PM · Apr 17, 2023

2 Likes

**Daniele Polencic — @danielepolencic@hachyderm.io**
@danielepolencic

8/

Unless they expose the metrics, the answer is no

However, there's a workaround

KEDA ships with an HTTP add-on to enable HTTP scaling
github.com/kedacore/http-...

How does it work!?

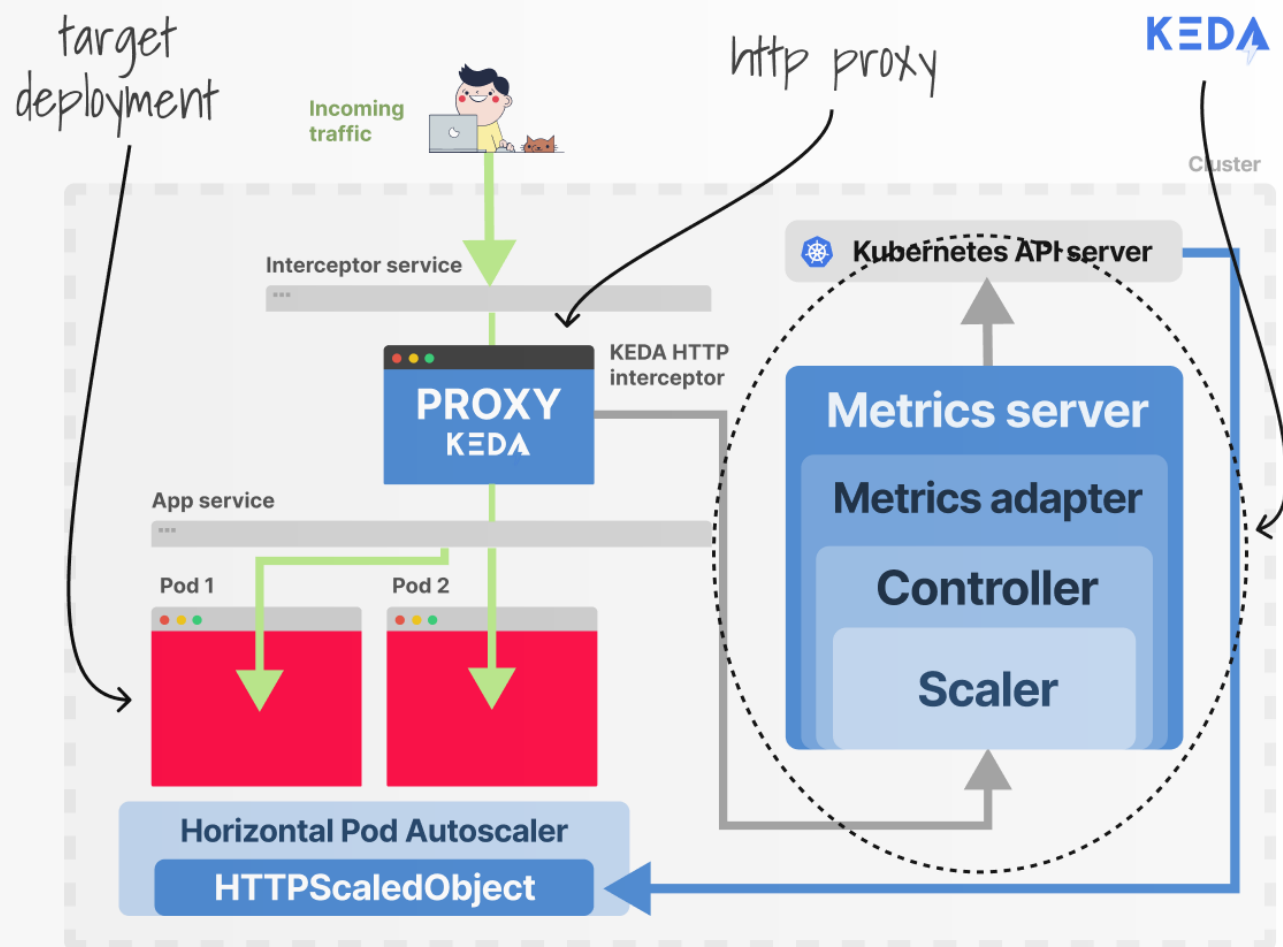8:11 PM · Apr 17, 2023

2 Likes

9/

KEDA injects a sidecar proxy in your pod so that all the HTTP traffic is routed first

Then it measures the number of requests and exposes the metrics

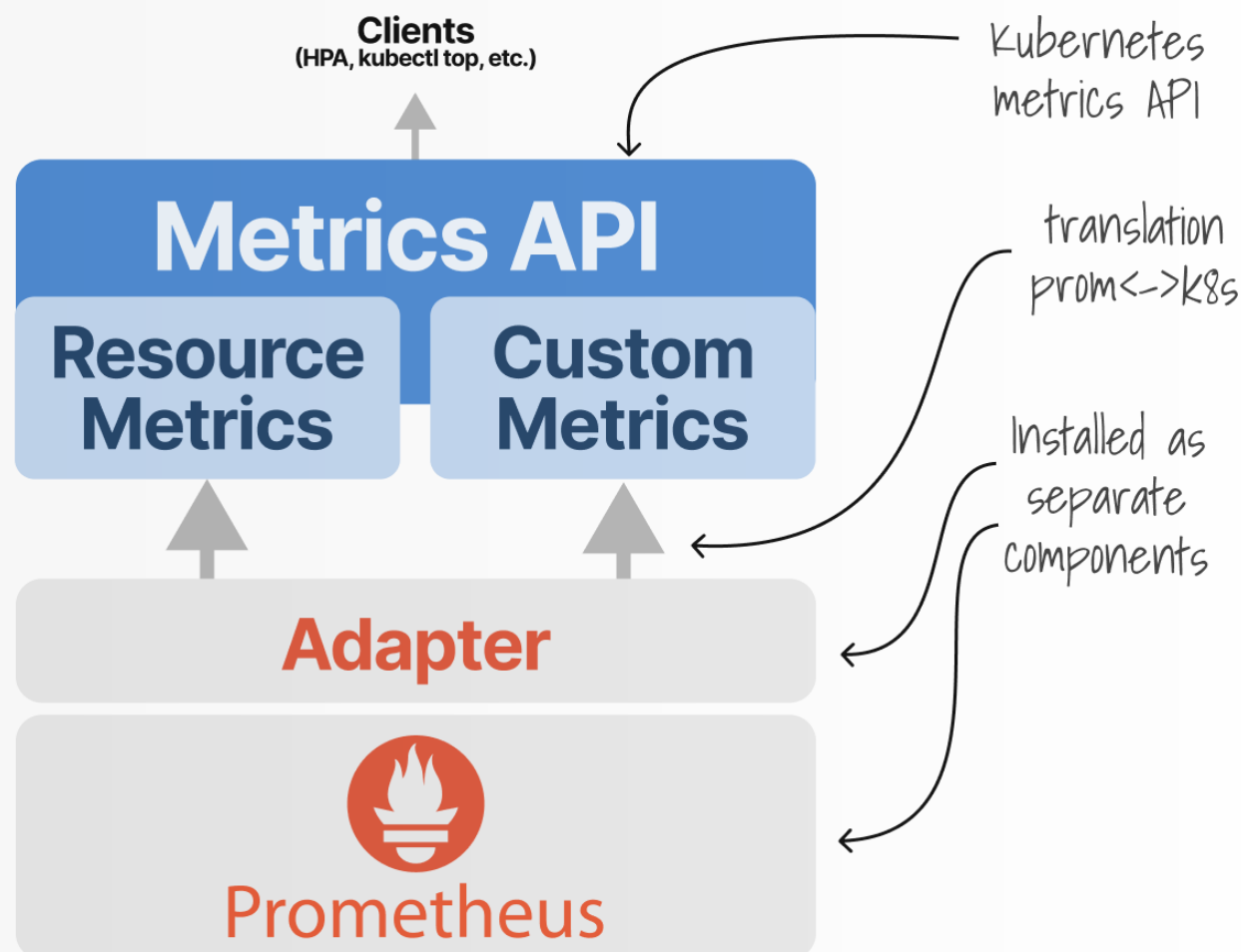With that data at hand, you can trigger the autoscaler finally

10/

KEDA is not the only option, though

You could install the Prometheus Adapter

The metrics will flow from Nginx to Prometheus, then the Adapter will make them available to Kubernetes

From there, they are consumed by the HPA



Kubernetes metrics API

translation prom<->k8s

Installed as separate components

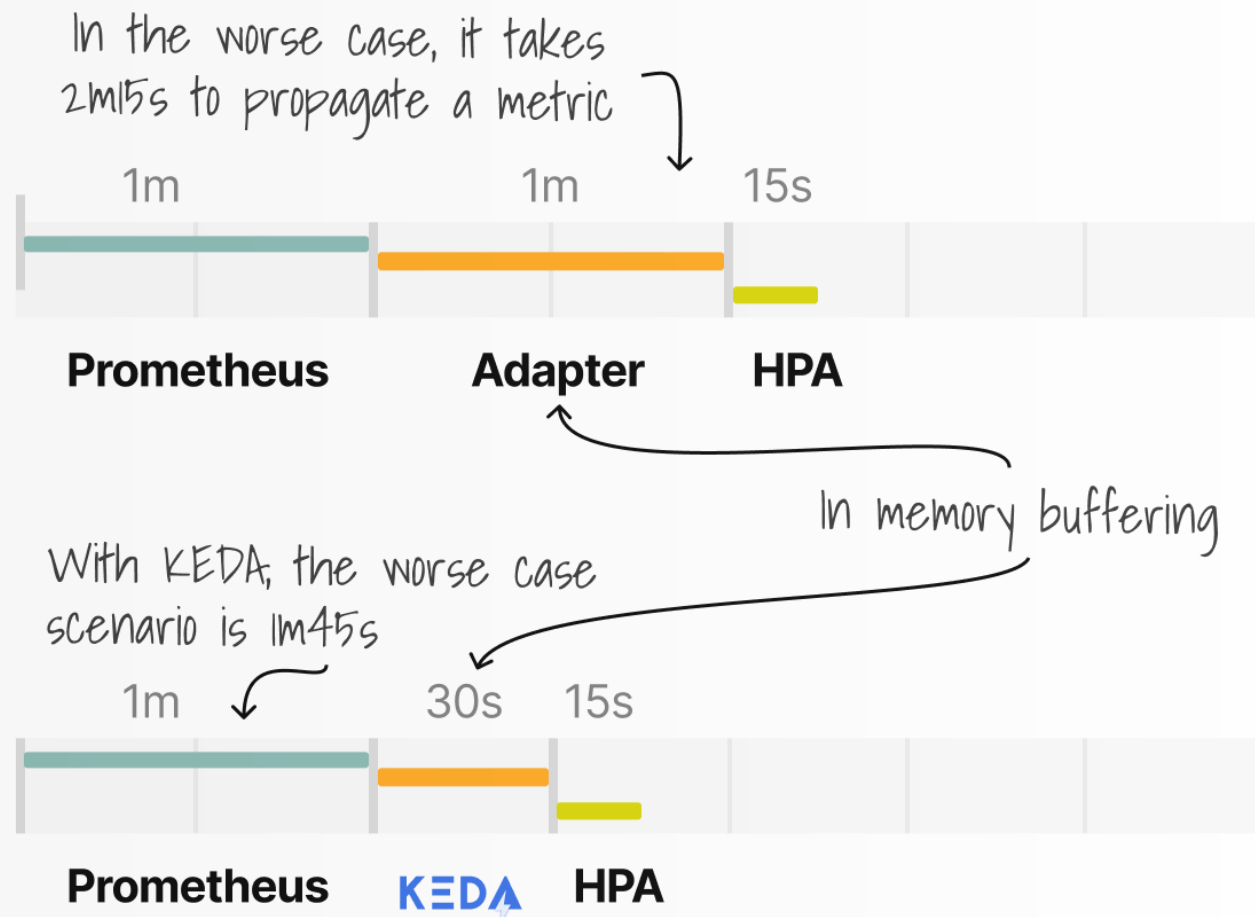**Daniele Polencic — @danielepolencic@hachyderm.io**
@danielepolencic

11/

Is this better than KEDA?

They are similar, as both have to query and buffer metrics from Prometheus

However, KEDA is pluggable and the Adapter works exclusively with Prometheus

In the worse case, it takes 2m15s to propagate a metric

| 1m | 1m | 15s |
|---|---|---|

**Prometheus** **Adapter** **HPA**

In memory buffering

With KEDA, the worse case scenario is 1m45s

| 1m | 30s | 15s |
|---|---|---|

**Prometheus** **KEDA** **HPA**

8:12 PM · Apr 17, 2023

1 Like

**Daniele Polencic — @danielepolencic@hachyderm.io**
@danielepolencic

12/

Is there a competitor to KEDA?

A promising project called the Custom Pod Autoscaler aims to make the pod autoscaler pluggable

However, the project focuses more on how those pods should be scaled (i.e. algorithm) than the metrics collection

custom-pod-autoscaler.readthedocs.io/en/latest/

8:12 PM · Apr 17, 2023

1 Like

**Daniele Polencic — @danielepolencic@hachyderm.io**
@danielepolencic

13/

During my research, I found these links helpful:

- keda.sh/docs/2.10/scal...
- sysdig.com/blog/kubernete...
- github.com/nginxinc/nginx...
- learnk8s.io/scaling-celery...
-

8:12 PM · Apr 17, 2023

5 Likes

**Daniele Polencic — @danielepolencic@hachyderm.io**
@danielepolencic

14/

And finally, if you've enjoyed this thread, you might also like:

- The Kubernetes workshops that we run at Learnk8s
learnk8s.io/training
- This collection of past threads twitter.com/danielepolenci...
- The Kubernetes newsletter I publish every week learnk8s.io/learn-kubernet...

8:12 PM · Apr 17, 2023

5 Likes