



KubeCon



CloudNativeCon

— North America 2018 —

SIG Autoscaling

Solly Ross (Google) and Michael Hausenblas (Red Hat)



Autoscaling?



KubeCon



CloudNativeCon

North America 2018

- On an abstract level:
 - Calculate resources to cover demand
 - Demand measured by metrics
 - Metrics must be collected, stored and queryable
- Ultimately to fulfill:
 - Service Level Objectives (SLO) ...
 - of Service Level Agreements (SLA) ...
 - through Service Level Indicators (SLI)

Types of autoscaling (in Kubernetes)



KubeCon



CloudNativeCon

North America 2018

- Cluster-level
- App-level
 - Horizontal
 - Vertical

Cluster Autoscaling

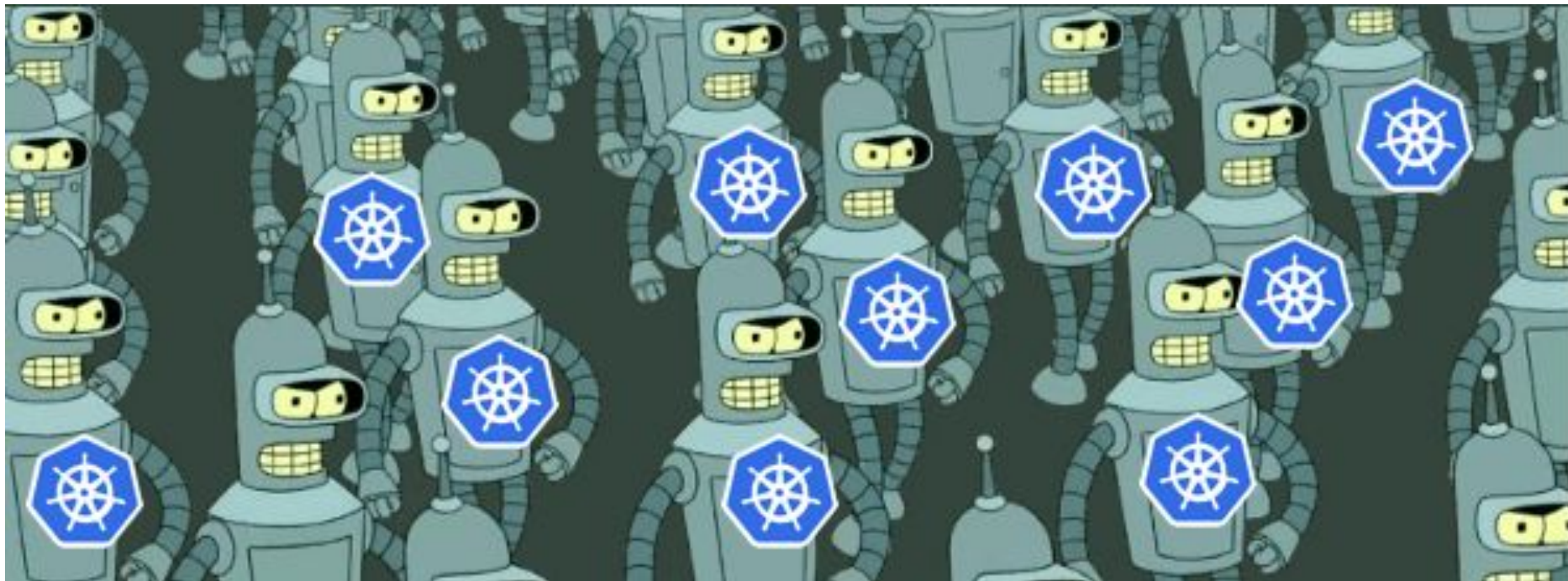


KubeCon



CloudNativeCon

North America 2018



https://cdn-images-1.medium.com/max/1600/1*P8Frww4Nvhrv_LwPnu5-w.png

Horizontal autoscaling



KubeCon



CloudNativeCon

North America 2018

- Horizontal pod autoscaler
- Resource: replicas
- “Increasing replicas when necessary”
- Requires application to be designed to scale horizontally



Vertical autoscaling



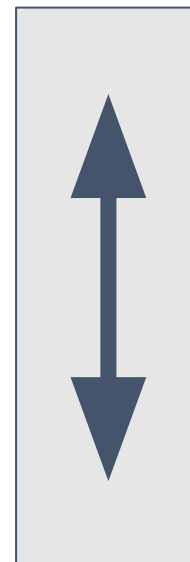
KubeCon



CloudNativeCon

North America 2018

- Vertical pod autoscaler
- Resource: CPU/Memory
- “Increasing CPU/Memory when necessary”
- Less complicated to design for resource increase
- Harder to autoscale



History of Autoscaling on Kubernetes

- Autoscaling used to heavily rely on Heapster
 - Heapster collects metrics and writes to time-series database
 - Metrics collection via cAdvisor (container + custom-metrics)
- We could autoscale!





KubeCon

CloudNativeCon

North America 2018

Custom Metrics API



Resource & Custom Metrics APIs



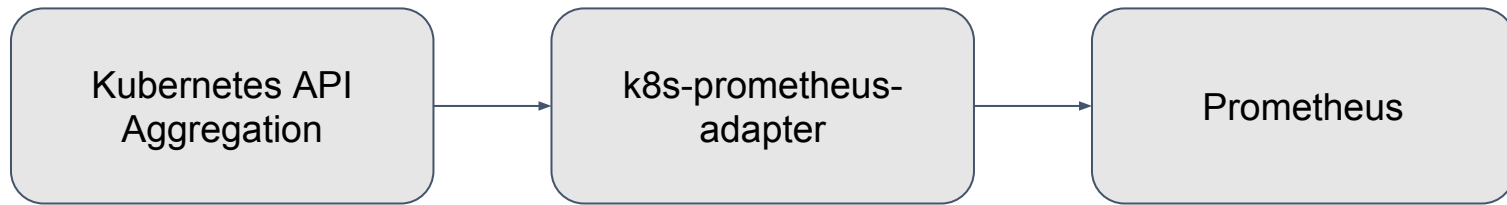
KubeCon



CloudNativeCon

North America 2018

- Well defined APIs:
 - Not an implementation, an API spec
 - Implemented and maintained by vendors
 - Returns single value





KubeCon

CloudNativeCon

North America 2018

Vertical Pod Autoscaling



Background & terminology



KubeCon



CloudNativeCon

North America 2018

```
15  apiVersion: extensions/v1beta1
16  kind: Deployment
17  metadata:
18    name: hamster
19    namespace: default
20  spec:
21    replicas: 2
22    template:
23      metadata:
24        labels:
25          app: hamster
26      spec:
27        containers:
28          - name: hamster
29            image: k8s.gcr.io/ubuntu-slim:0.1
30            resources:
31              requests:
32                cpu: 100m
33                memory: 50Mi
34            command: ["/bin/sh"]
35            args: ["-c", "while true; do timeout 0.5s yes >/dev/null; sleep 0.5s; done"]
```

Background & terminology



KubeCon



CloudNativeCon

North America 2018

- Scheduling
 - nodes offer resources
 - pods consume resources
 - scheduler matches needs of pods based on requests
- Types of resources (compressible/incompressible)
- Quality-of-Service (QoS)
 - Guaranteed: $\text{limit} == \text{request}$
 - Burstable: $\text{limit} > \text{request} > 0$
 - Best-Effort: \nexists (limit, request)

Motivation



KubeCon



CloudNativeCon

North America 2018

Unfortunately, Kubernetes has not yet implemented dynamic resource management, which is why we have to set resource limits for our containers. I imagine that at some point Kubernetes will start implementing a less manual way to manage resources, but this is all we have for now.

Ben Visser, 12/2016

[Kubernetes—Understanding Resources](#)

Kubernetes doesn't have dynamic resource allocation, which means that requests and limits have to be determined and set by the user. When these numbers are not known precisely for a service, a good approach is to start it with overestimated resources requests and no limit, then let it run under normal production load for a certain time.

Antoine Cotten, 05/2016

[1 year, lessons learned from a 0 to Kubernetes transition](#)

- Automating configuration of resource requirements
 - manually setting requests is brittle & hard so people don't do it
 - no requests set → QoS is *best effort* :(
- Improving utilization
 - can better bin pack
 - impact on other functionality such as [out of resource handling](#) or an (aspirational) optimizing scheduler

Use Cases



KubeCon



CloudNativeCon

North America 2018

- For stateful apps, for example Wordpress or single-node databases
- Can help on-boarding of "legacy" apps, that is, non-horizontally scalable ones



Interlude: API server

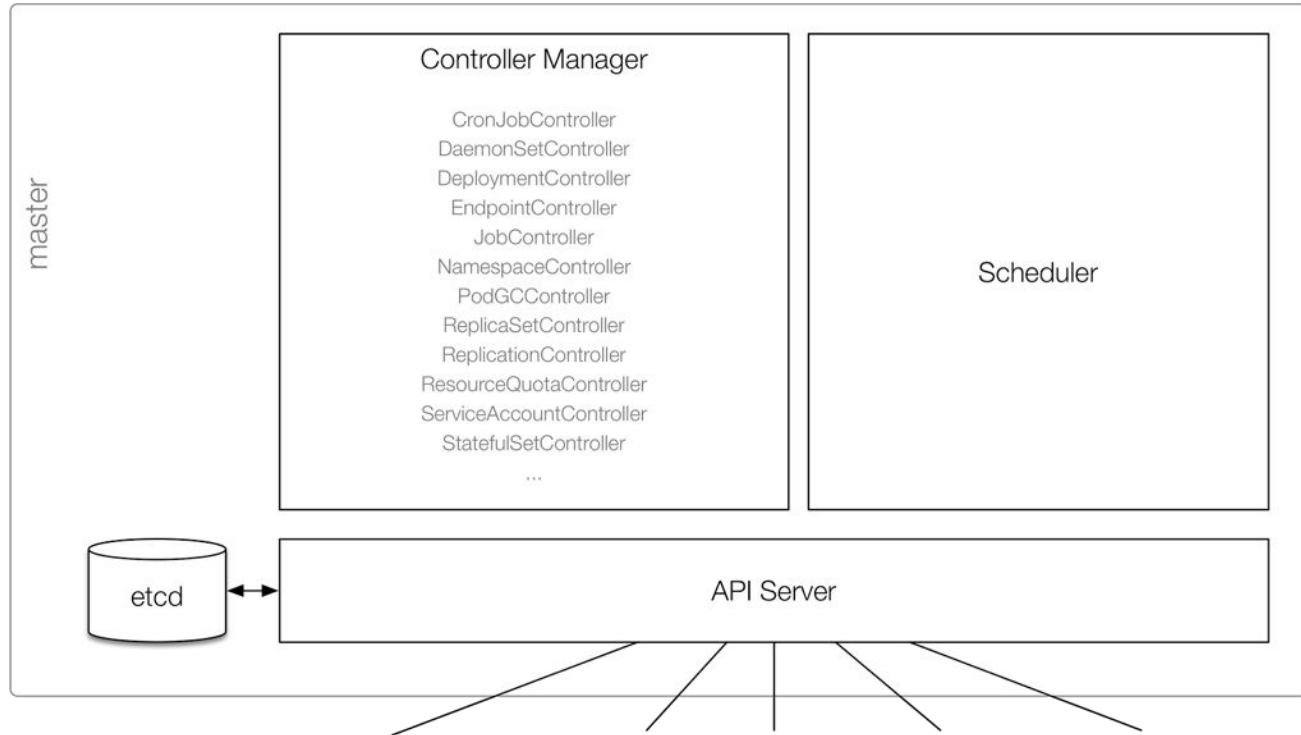


KubeCon



CloudNativeCon

North America 2018



Interlude: API server

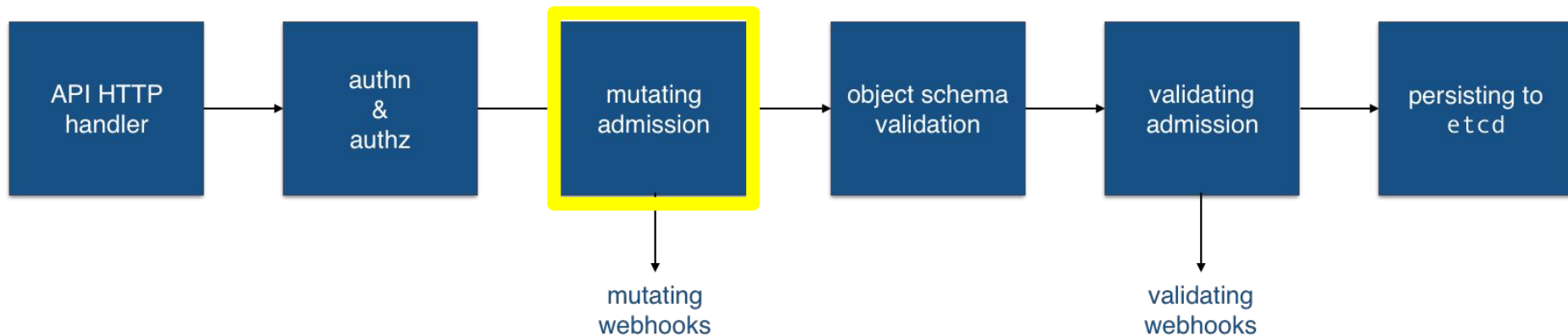


KubeCon



CloudNativeCon

North America 2018



Basic idea



KubeCon



CloudNativeCon

North America 2018

- observe resource consumption of all pods
- build up historic profile (recommender)
- apply to pods on an opt-in basis via labels (updater)

VPA architecture

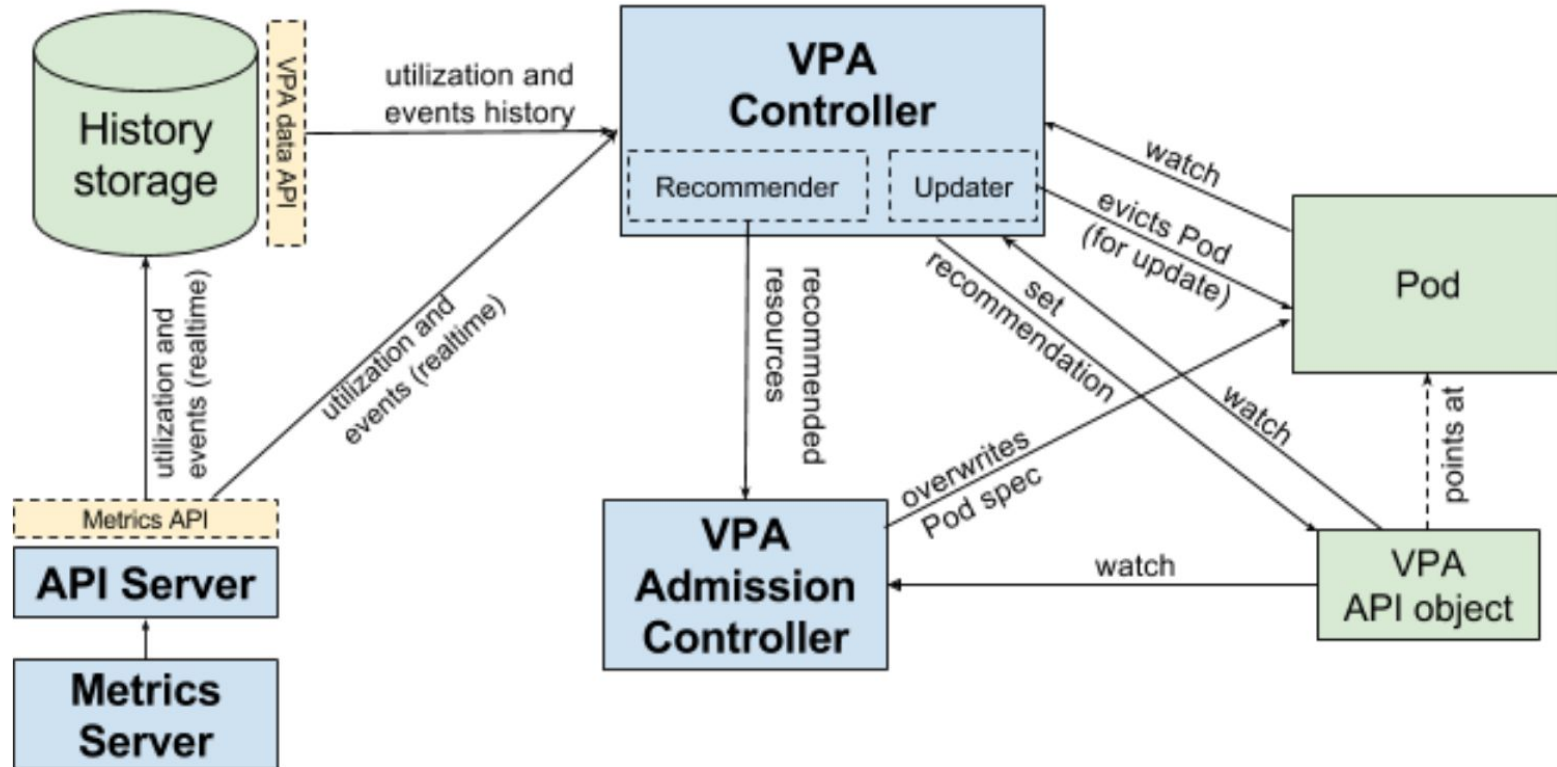


KubeCon



CloudNativeCon

North America 2018



Limitations



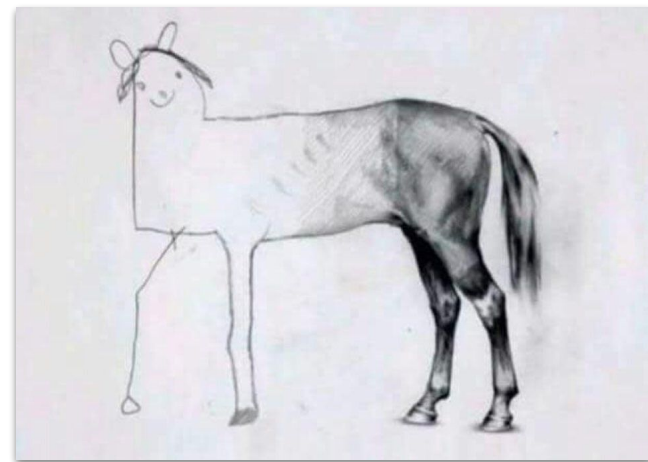
KubeCon



CloudNativeCon

North America 2018

- pre-alpha, so need testing and tease out edge-cases
- [in-place updates](#) (requires support from container runtime)
- usage spikes—how to deal with it best?



VPA demo



KubeCon



CloudNativeCon

North America 2018

The screenshot displays the OpenShift console interface for a pod named 'hamster-6db596f5b4-mrp2c'. The left sidebar shows navigation options like Overview, Applications, Builds, Resources, Storage, Monitoring, and Catalog. The main content area shows the pod's status as 'Running', its IP address (18.0.18.16), and its node (fedora-kube-single.localdomain). Below this, the container 'hamster' is shown with its state as 'Running since Aug 6, 2018 2:11:20 PM'. The 'Template' section shows the pod's configuration, including the image 'k8s.gcr.io/ubuntu-slim0.1', the command to run a loop, and resource requests for CPU (300m) and Memory (319572800).

```
Update Policy:
Update Mode: Auto
Events: <none>
demo@: ~[go:autoscaler]/vertical-pod-autoscaler $ kubectl describe vpa
Name: hamster-vpa
Namespace: vpa-demo
Labels: <none>
Annotations: <none>
API Version: poc.autoscaling.k8s.io/v1alpha1
Kind: VerticalPodAutoscaler
Metadata:
  Creation Timestamp: 2018-08-06T18:09:42Z
  Generation: 1
  Resource Version: 866836
  Self Link: /apis/poc.autoscaling.k8s.io/v1alpha1/namespaces/vpa-demo/verticalpodautoscalers/hamster-vpa
  UID: e47d5eaa-99a3-11e8-be76-5254002b8d24
Spec:
  Selector:
    Match Labels:
      App: hamster
  Update Policy:
    Update Mode: Auto
Status:
  Conditions:
    Last Transition Time: 2018-08-06T18:10:28Z
    Status: True
    Type: RecommendationProvided
  Recommendation:
    Container Recommendations:
      Container Name: hamster
      Lower Bound:
        Cpu: 159m
        Memory: 254761479
      Target:
        Cpu: 200m
        Memory: 319572800
      Upper Bound:
        Cpu: 24200m
        Memory: 38668308800
Events: <none>
demo@: ~[go:autoscaler]/vertical-pod-autoscaler $
```

Resources & what's next?



- VPA issue [10782](#)
- VPA [design](#)
- Test, provide feedback
- [SIG Autoscaling](#)—come and join us on #sig-autoscaling
 - or weekly online meetings on Monday
- SIG Instrumentation and SIG Autoscaling work towards a historical metrics API—get involved there!