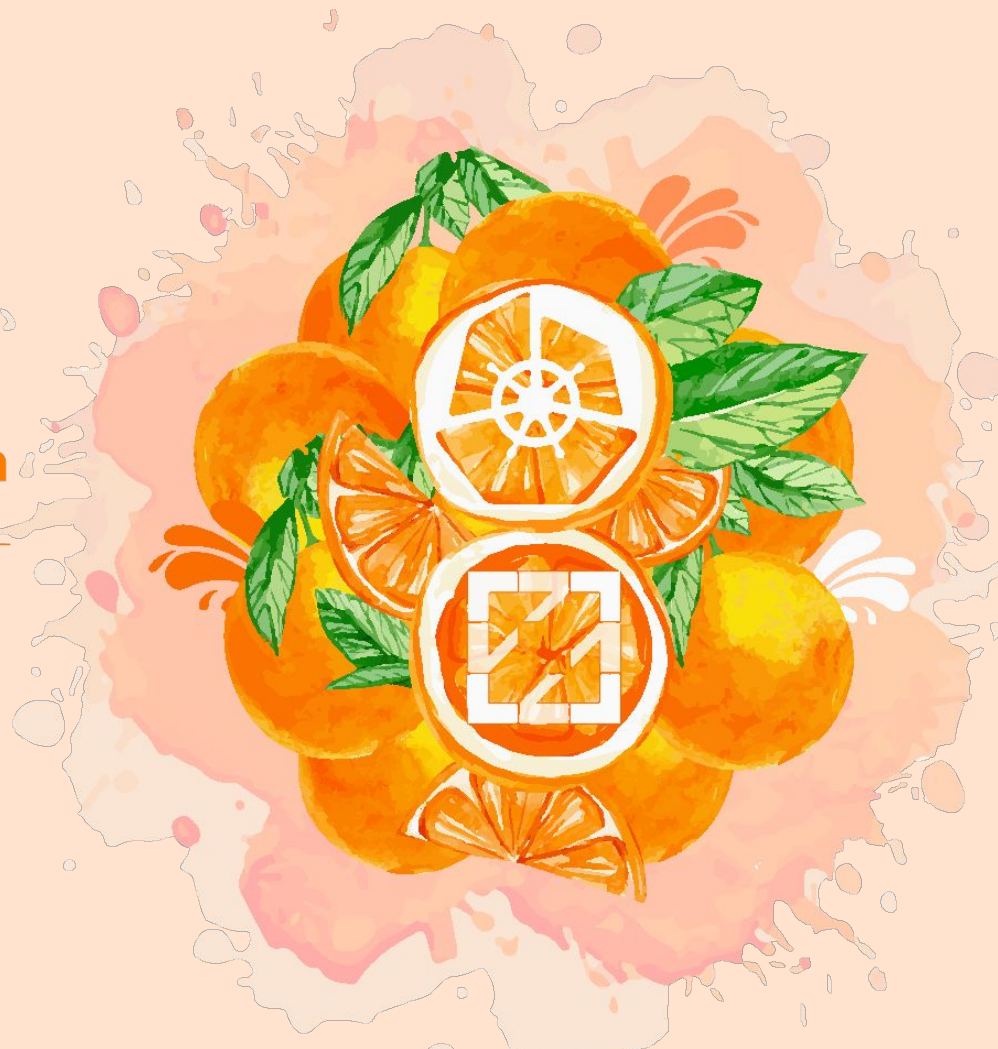KubeCon | CloudNativeCon

Europe 2022

WELCOME TO VALENCIA

# Autoscaling Kubernetes Deployments:
# A (Mostly) Practical Guide

Natalie Serrino, New Relic

# About

💻 Software engineer / TLM at New Relic

💥 Contributor to Pixie

💜 Observability and performance

📊 Formerly worked in the data space

**Natalie Serrino**

# Content

🤔 What is Kubernetes autoscaling (and why is it useful)?

🕛 What knobs does Kubernetes autoscaling give us?

📈 Selecting an autoscaling metric for your application

🤓 A Turing-complete autoscaler (?!)

# What is Kubernetes autoscaling?

# Resource sizing in Kubernetes

How do you select the following values?

💰 # of nodes in your cluster
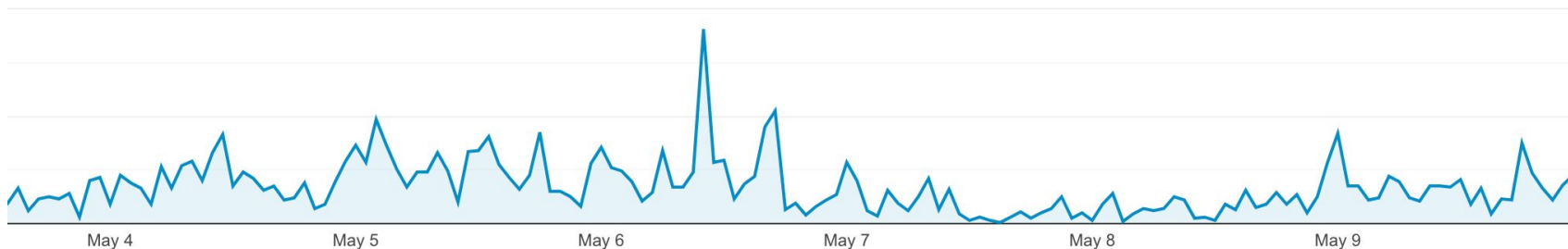🚀 # of pods in a deployment
📊 Amount of resources to give a pod

Methods…

● Random guess
● Copypasta
● Proactive iteration
● Reactive iteration

# Why is autoscaling useful?

- Ideal resource allocation depends on workload
- Workloads are often spikey and unpredictable



- Too few resources → bad user experience, latency, outages 😱
- Too many resources → wasteful, expensive 💸💸💸

# The types of autoscaling in Kubernetes

| Kubernetes autoscaler | …adds/removes… | …to your… | …based on… |
|---|---|---|---|
| `Cluster Autoscaler` | Nodes | Cluster | Resource utilization |
| `VerticalPodAutoscaler` | Resources like CPU/Memory | Existing replicas | Resource utilization |
| `HorizontalPodAutoscaler` | Replicas | Workload | Resource utilization or user-defined metrics |

# Cluster Autoscaler

💡 Set pod resource requests and limits

💡 Make sure resource requests reflect actual usage

💡 Specify PodDisruptionBudgets

💡 Compare limits with available quota from your cloud provider (if applicable)

⚠️ Kubernetes contributors tested for <=1000 nodes with 30 pods each



*Image Credit: Understanding Kubernetes Cluster Autoscaling by Ajay Tripathy*

# VerticalPodAutoscaler (VPA)

💡 Can still set resource caps with VPA

💡 Updates may result in container restarts or pod rescheduling

💡 Use in cluster autoscaler to avoid VPA recommending more than available resources

⚠️ Can't use with HorizontalPodAutoscaler on the same application for CPU/memory yet

⚠️ Has not yet been tested on large clusters



*Image Credit: Red Hat, How Full is my Cluster, Raffaele Spazzoli*

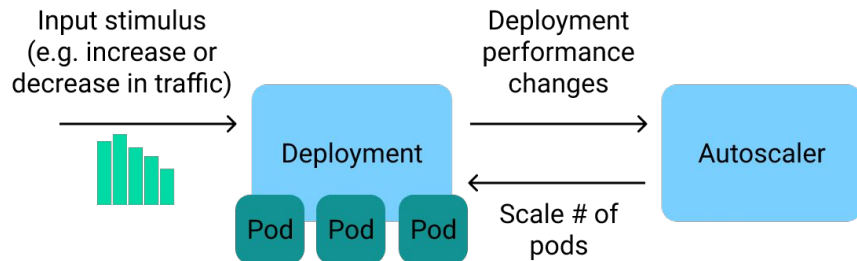# HorizontalPodAutoscaler (HPA)

💡 Lots of flexibility for metric selection

💡 Check your service client affinity policies to ensure even load distribution

💡 When scaling on CPU/memory, make sure to set resource requests

⚠️ Can't use with VerticalPodAutoscaler on same application on CPU/memory yet

Input stimulus
(e.g. increase or
decrease in traffic)

Deployment
performance
changes

Deployment

Pod  Pod  Pod

Autoscaler

Scale # of
pods

# Demo: Horizontal pod autoscaling on CPU

# HPA autoscaling equation

```
outputReplicas = ceil(
    currentReplicas * ( currentMetricValue / desiredMetricValue )
)
```

What knobs does Kubernetes autoscaling give us?

# Scaling up and down

- How to set the minimum and maximum number of replicas?

- How often to look for changes in the metric?

- How quickly to add or remove pods?
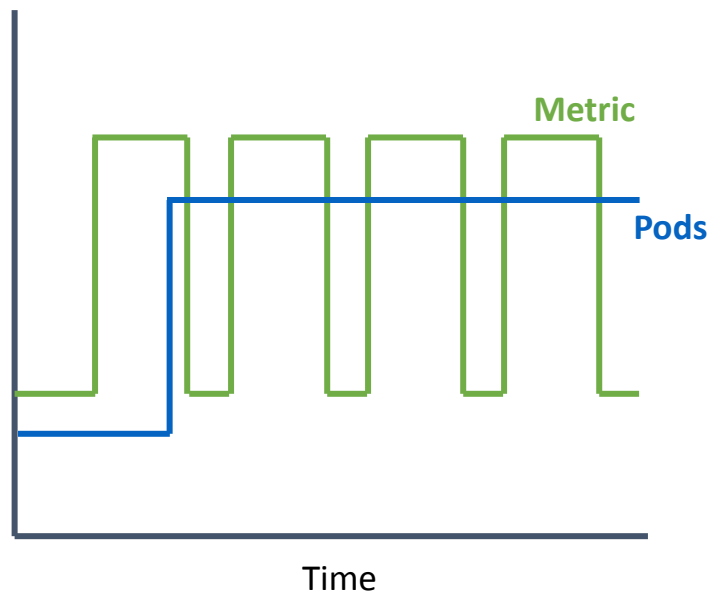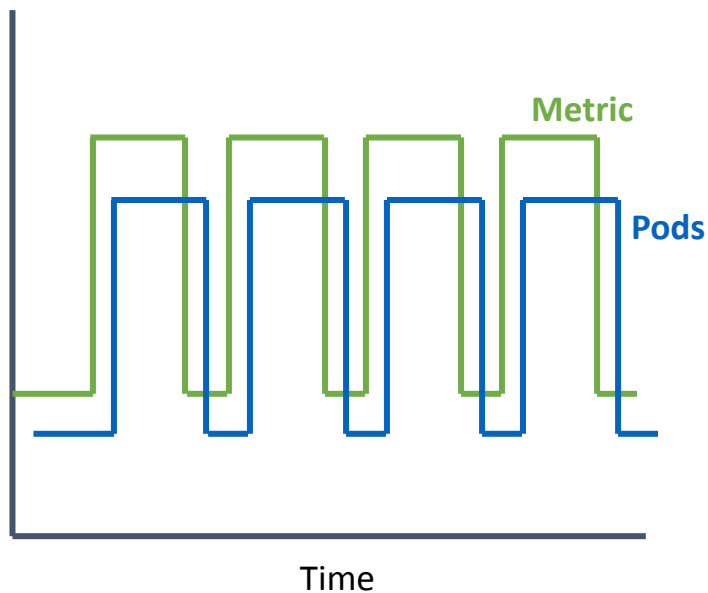
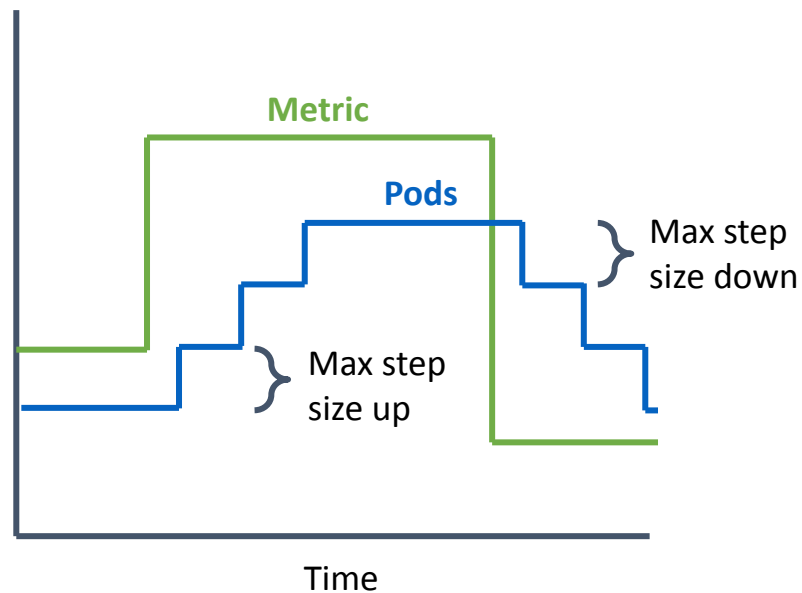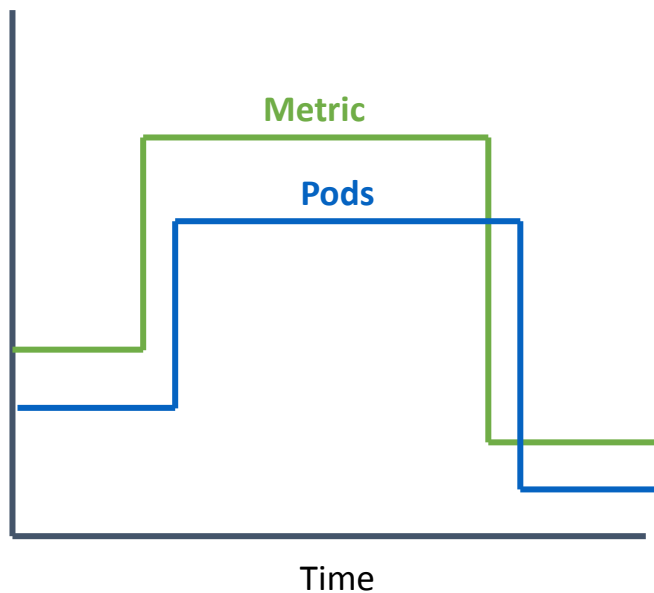- How many pods to add/remove in a single period?
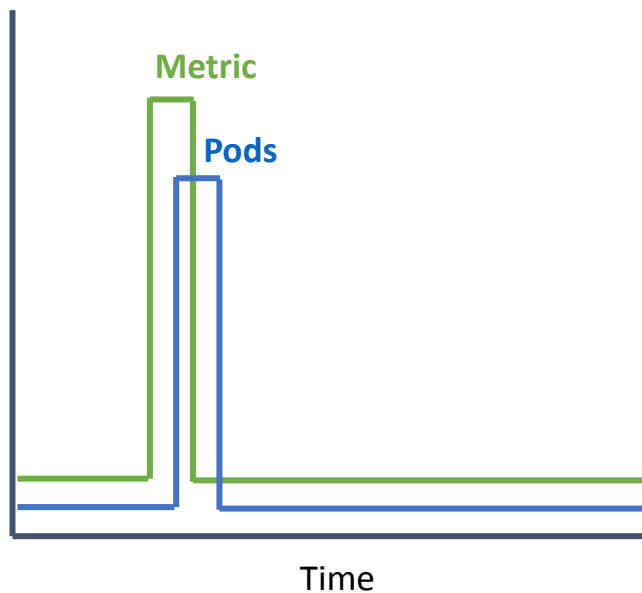
# Capping min/max pods

# Stabilization period

# Stabilization period reduces pod churn

# Capping max pods to add/remove per period

# Capping max step size reduces pod churn

# Selecting an autoscaling metric for your application

# HPA Metric Types

| Category | Type | Description | Examples |
|----------|------|-------------|----------|
| `Resource Metrics` | Built-in | Resource utilization metrics for pods and nodes only | Currently limited to CPU and memory |
| `Custom Metrics` | User-defined | Custom metrics about Kubernetes resources | Latency, throughput, queue depth |
| `External Metrics` | User-defined | Custom metrics NOT about Kubernetes resources | # of customers using website |

# Possible bottlenecks in an application

- CPU
- Memory
- Network
- # of worker threads
- # of outbound connections
- Downstream dependencies
- Queue depth
- …Many more

The best metric to scale on depends on your workload!

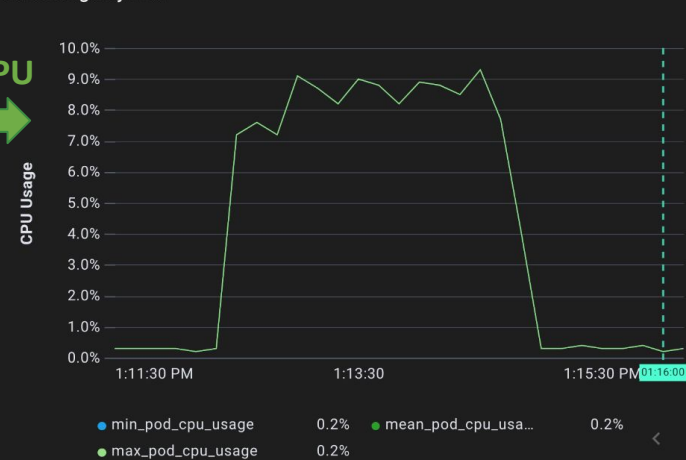A Turing-complete autoscaler?

# Turing machine

Turing machines are capable of any computation, given enough time and tape.

## Input program

```
Instruction 1
Instruction 2
Instruction 3
…
Instruction N
```

*Next instruction*

Turing Machine

*Write value*

## Output tape

| 3 | 0 | 7 | 1 | 4 | 3 | 3 | 2 | 5 | 2 | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# Subleq

"One-instruction set computer", sufficient for Turing completeness.

```
Instruction subleq a, b, c
    Mem[b] = Mem[b] - Mem[a]
    if (Mem[b] ≤ 0)
        goto c
```

*Subleq pseudocode (credit: Wikipedia)*

# HPA subleq Turing machine

Input program

```
subleq 3, 0, 3
subleq 5, 3, 2
subleq 2, 1, 1
...
subleq 4,-1, 0
```

*Execute 1 instruction per HPA interval*

Custom HPA metric provider

*Set # pods*

Number of pods over time

| 3 | 1 | 7 | 1 | 4 | 3 | 3 | 2 | 5 | 2 | | | | |

# Input program: deployment name

```
9x-1x3x10x-1x6x0x0x-1x72x105x0
```

*Split on "x"*

```
9,-1,3,10,-1,6,0,0,-1,72,105,0
```

*3 input args per subleq*

```
subleq  9, -1, 3
subleq 10, -1, 6
subleq  0,  0,-1
subleq 72,105, 0
```
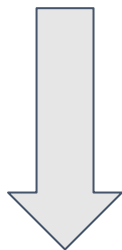
# Setting a certain # of output pods

```
outputReplicas = ceil(
    currentReplicas * ( currentMetricValue / targetMetricValue )
)
```



*Backwards calulating
"current metric value"*

```
currentMetricValue = (
    targetMetricValue * ( outputReplicas / currentReplicas )
)
```

# Demo: Turing-complete autoscaler

SIG instrumentation

[github.com/kubernetes-sigs/custom-metrics-apiserver](github.com/kubernetes-sigs/custom-metrics-apiserver)

CNCF Sandbox Project

[github.com/pixie-io/pixie](github.com/pixie-io/pixie)

[blog.px.dev](blog.px.dev)

Load generator

[github.com/rakyll/hey](github.com/rakyll/hey)

Thanks!