# Spot Instances

Philipp Kühn @pkse as 1st killed soldier
Alby Hernández @achetronic as 2nd killed soldier

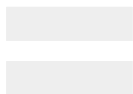## How we discovered they are not reliable

Docplanner

# Spot Instances: on-demand vs. spot

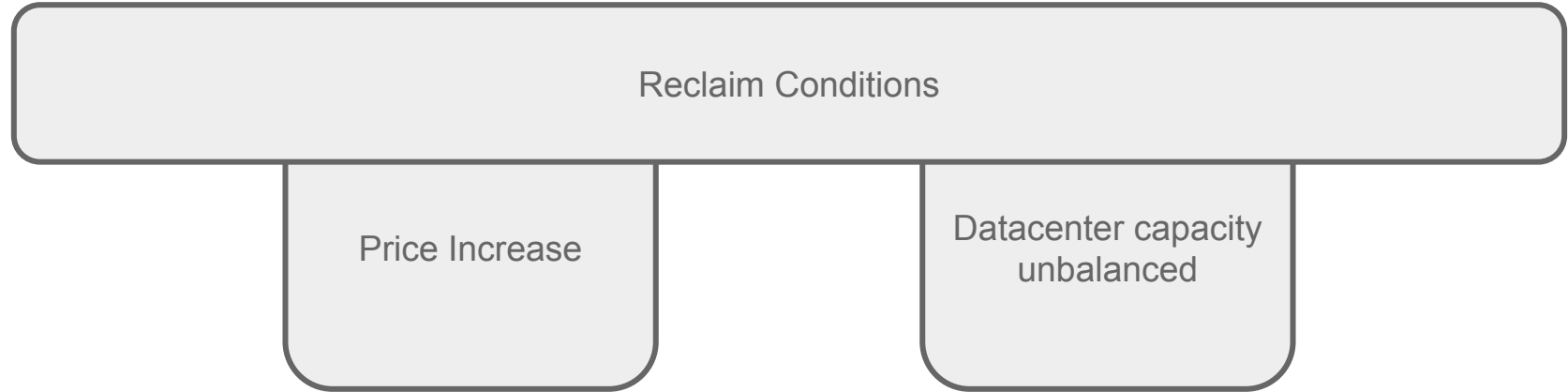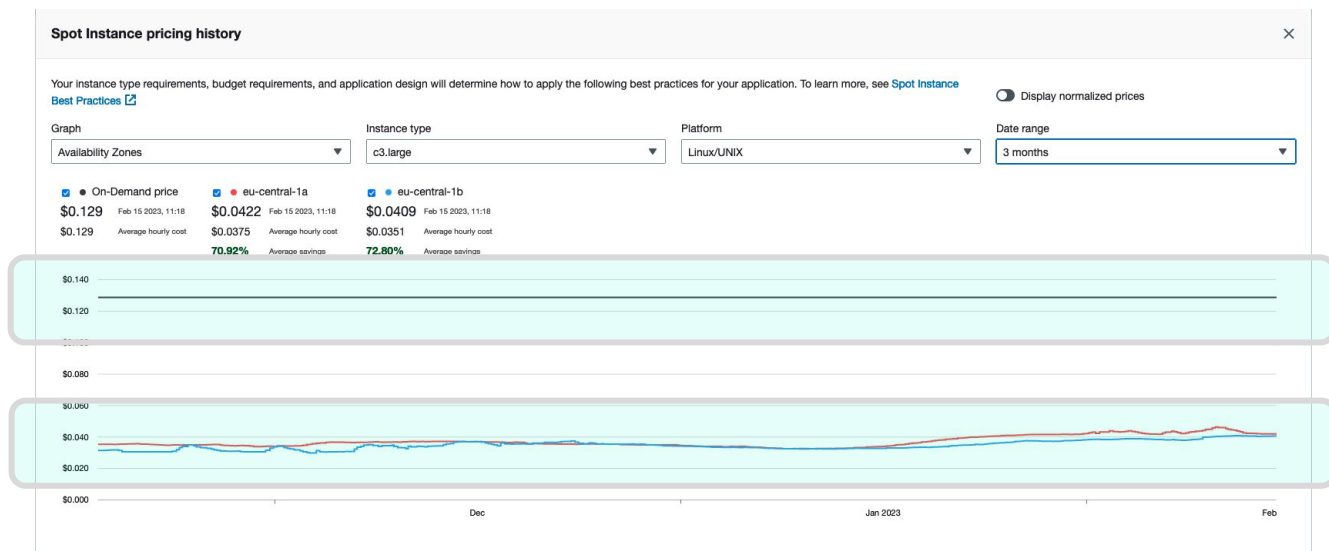Spot = On-demand + Reclaim Conditions

Docplanner

# Spot Instances: on-demand vs. spot

Reclaim Conditions = On some situations related to AWS, these instances can be destroyed

# Spot Instances: reclaim conditions
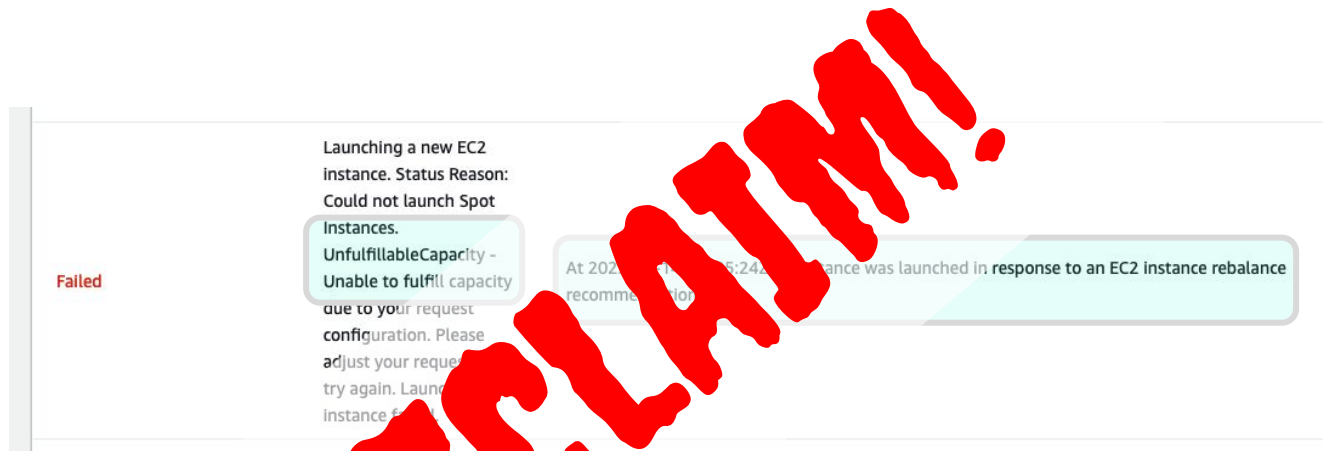
Reclaim Conditions

Price Increase

Datacenter capacity unbalanced

Docplanner

# Spot Instances: Price increased

# Spot Instances: reclaim conditions



...nted with old price?
how we loose money

RECLAIM!

Jan 2023

Docplanner

# Spot Instances: UNBALANCED DATA CENTER



**Failed** — Launching a new EC2 instance. Status Reason: Could not launch Spot Instances. UnfulfillableCapacity - Unable to fulfill capacity due to your request configuration. Please adjust your request and try again. Launching instance for...

At 202... 5:24... ...ance was launched in **response to an EC2 instance rebalance** recommendation

**RECLAIM!**

aws forgot to pre-balance data center's VMs

Docplanner

# The Pieces

# Overview



AWS

| ASG 1 | ASG 2 | ... | ASG n |

Kubernetes

| Cluster Autoscaler | Cluster Overprovision | Node Termination Handler | Tolerations & Affinity |

Docplanner

# The Pieces

## Auto Scaling Groups

# Spot Auto Scaling Group: Instance types



- As many as possible

- Same CPU & Memory

# Spot Auto Scaling Group: DECREASE THE RISK



- Proper strategy

- Allow AWS create before destroy

Docplanner

# The Pieces

## Cluster Autoscaler

Docplanner

# Cluster Autoscaler: Priority Expander



- Support Go regexp for matching ASGs

- Provide more control over ASGs we scale

Docplanner

# Cluster Autoscaler: Priority Expander

```yaml
apiVersion: v1
kind: ConfigMap
metadata:
  name: cluster-autoscaler-priority-expander
  namespace: kube-system
data:
  priorities: |-
    10:
      - .*t2\.large.*
      - .*t3\.large.*
    50:
      - .*m4\.4xlarge.*
```

```
23    autoscaler_priorities = <<EOT
24    1000:
25      - .*spot.*
26    5:
27      - .*
28    EOT
29
```

Example on docs

Our configuration

Docplanner

# The Pieces

## Node Termination Handler

# Node Termination Handler: What?

Able to
Cordon & drain
on

EC2 maintenance events

ASG AZ Rebalance

EC2 Spot interruptions

EC2 Instance Termination

ASG Scale-In

Docplanner

# Node Termination Handler: Why required for spots?

EC2 Spot interruptions

=

Rebalance Recommendation

+

2 Grace Minutes

Early indicator for certain instance

Do something…

Docplanner

# Node Termination Handler: Why required for spots?

EC2 maintenance events = No more emails about rebooting some instances

Docplanner

# The Pieces

Tolerations & Affinity

# Tolerations & Affinity: Let Kubernetes be kubernetes

```yaml
1  tolerations: &tolerations
2    - effect: NoSchedule
3      key: capacityType
4      operator: Equal
5      value: spot
6  affinity: &affinity
7    nodeAffinity:
8      requiredDuringSchedulingIgnoredDuringExecution:
9        nodeSelectorTerms:
10          - matchExpressions:
11            - key: eks.amazonaws.com/nodegroup
12              operator: In
13              values:
14              # ATTENTION: It is important to add the spo
15              # just in case spots are not available
16              - nodes-app-spot
17              - nodes-app
18        preferredDuringSchedulingIgnoredDuringExecution:
19          - weight: 1
20            preference:
21              matchExpressions:
22                - key: capacityType
23                  operator: In
24                  values:
25                    - spot
26
```

Hey, Kubernetes!
I can handle spot tainted nodes

Hey, Kubernetes!
I require nodes with these labels

Hey, Kubernetes!
I prefer spot

Docplanner

# Tolerations & Affinity: Let Kubernetes be kubernetes

```
27   # ATTENTION: Following lines are totally optional,
28   # when some services require to be scheduler avoid
29   # where other known application is already running
30   podAntiAffinity:
31     requiredDuringSchedulingIgnoredDuringExecution:
32     - labelSelector:
33         matchExpressions:
34         - key: app.kubernetes.io/instance
35           operator: In
36           values:
37             - one
38   topologyKey: "kubernetes.io/hostname"
39       # Remember that pods are namespaced, so implic
40       # so you have to specify the namespaces where
41       # Ref: https://kubernetes.io/docs/concepts/sch
42   namespaces:
43   - one-app
```

Hey, Kubernetes!
I hate Saas, I wanna be
allocated in other place

Remember labels belong to pods, so they
are namespaced too

Docplanner

# The Pieces

Cluster Overprovisioner

# C. Overprovisioner: Components

Cluster Overprovisioner = Placeholder Pods + Cluster Proportional Autoscaler

Docplanner

# C. Overprovisioner: Placeholder

Placeholder Pods

- Use pause Docker image

- Requests resources

- Lowest PriorityClass

# C. Overprovisioner: Placeholder

Placeholder Pods **=** Force schedule some unused resources

Docplanner

# C. Overprovisioner: Cluster Proportional Autoscaler

Cluster Proportional Autoscaler

- Watches a deployment

- Count number of certain labeled nodes

- Scale paused pod replicas

Docplanner

# C. Overprovisioner: Cluster Proportional Autoscaler

Cluster Proportional Autoscaler **=** Force schedule DYNAMIC amount of placeholders

Docplanner

# C. Overprovisioner: CPA. Our config

| | | |
|---|---|---|
| 0 - 5 Nodes | Automatically… | 4 Placeholders |
| 5 - 10 Nodes | | 6 Placeholders |
| 10+ Nodes | | 8 Placeholders |

Docplanner

# C. Overprovisioner: CPA. Our config

Placeholder
Pods

=

```
resources:
  requests:
    cpu: 4
    memory: 2000Mi
  limits:
    cpu: 5
    memory: 3000Mi
```

More than a Monolith

Docplanner

# Initial Results

Round 1

Docplanner

# Results: Everything started to work



Spots nodes over the time
Feb 12, 10:20am – Feb 13, 10:20am

60
50
40
30
20
10
0

Feb 12, 10:00am | Feb 12, 1:00pm | Feb 12, 4:00pm | Feb 12, 7:00pm | Feb 12, 10:00pm | Feb 13, 1:00am | Feb 13, 4:00am | Feb 13, 7:00am | Fel 10:(
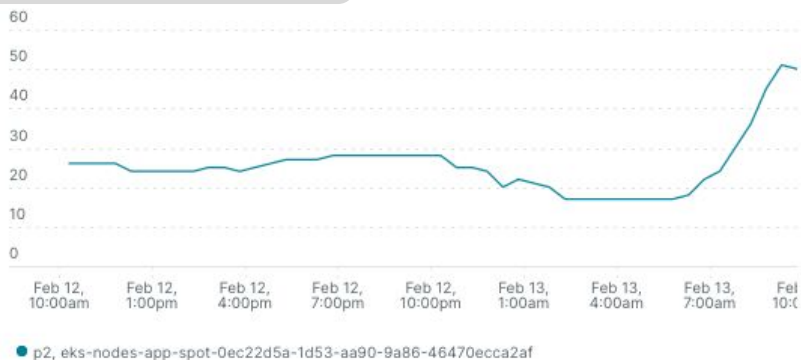
● p2, eks-nodes-app-spot-0ec22d5a-1d53-aa90-9a86-46470ecca2af

On-demand nodes over the time
Feb 12, 10:20am – Feb 13, 10:20am

18
16
14
12
10
8
6
4
2
0

Feb 12, 10:00am | Feb 12, 1:00pm | Feb 12, 4:00pm | Feb 12, 7:00pm | Feb 12, 10:00pm | Feb 13, 1:00am | Feb 13, 4:00am | Feb 13, 7:00am | Fel 10:(

● p2, eks-nodes-svc-02c22d55-95cf-6683-90f1-...   ● p2, eks-nodes-app-a2c22d55-9602-35a9-b85f...
● p2, eks-nodes-anonymization-e2c22d55-95fa-...   ● p2, eks-nodes-app-saas-66c22d55-95fe-7ac1-...

Docplanner

# Results: Everything broke



**Spots nodes over the time**
Feb 14, 9:11am – Feb 14, 12:11pm

● p2, eks-nodes-app-spot-0ec22d5a-1d53-aa90-9a86-46470ecca2af

**On-demand nodes over the time**
Feb 14, 9:11am – Feb 14, 12:11pm

● p2, eks-nodes-app-a2c22d55-9602-35a9-b85f-    ● p2, eks-nodes-svc-02c22d55-95cf-6683-90f1-
● p2, eks-nodes-anonymization-e2c22d55-95fa-    ● p2, eks-nodes-app-saas-66c22d55-95fe-7ac1-

Docplanner

# Results: Everything broke

# Break (10 minutes)

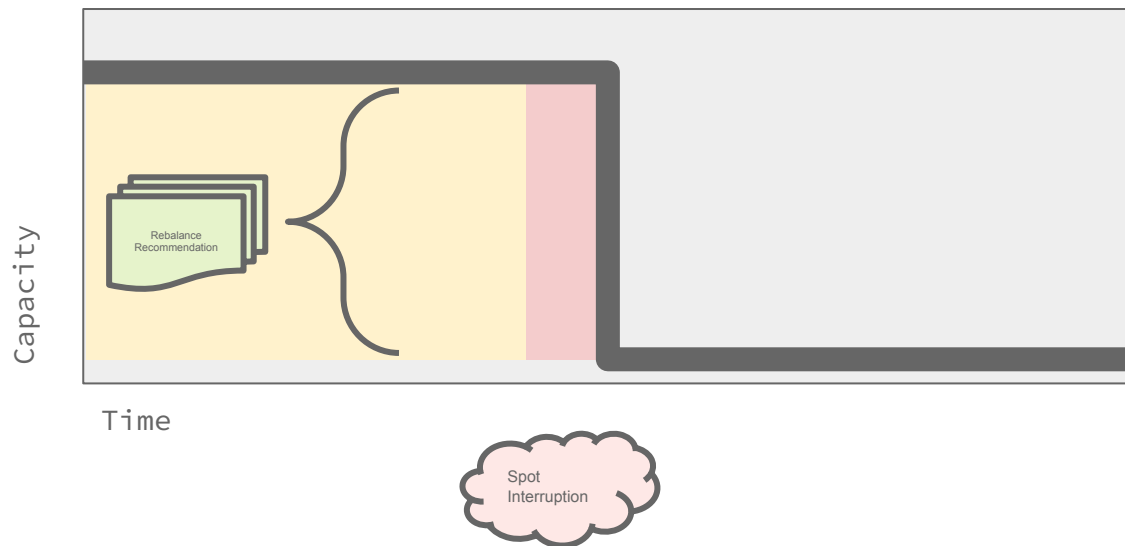Go for a coffee *cough cough w.c*

Docplanner

# Potential solution

Round 2

Docplanner

# Possible solutions: The problem

# Possible solutions: The problem

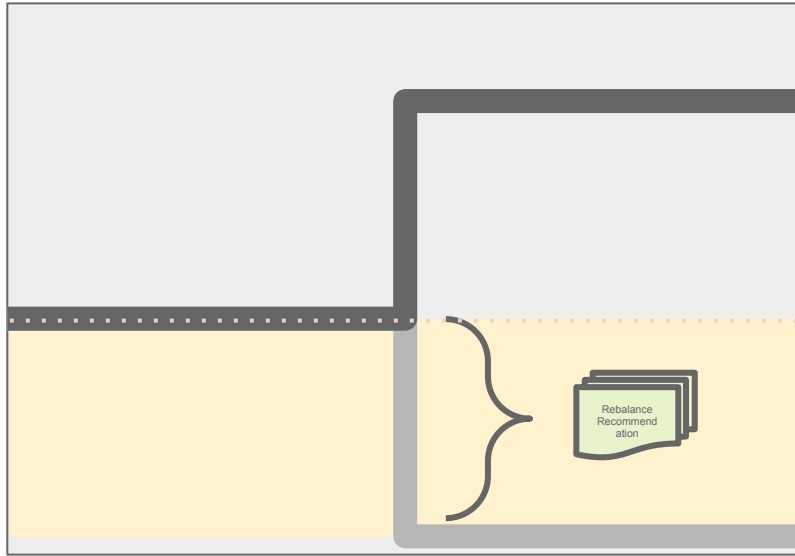# Possible solutions: Dig Deeper

10
Ready
Nodes

Rebalance
Recommend
ation

7 RR = 7 potential lost nodes

Docplanner

# Possible solutions: Dig Deeper

| 10 Ready Nodes | | 7 RR = 7 potential lost nodes | | 3 Sane Nodes |

Action time

![Docplanner logo]

# Possible solutions: Dig Deeper

**10 Ready Nodes**

**7 needed extra nodes**

Rebalance Recommendation

**7 potential lost nodes**

Docplanner

# Possible solutions: Solution

| | | | | | |
|---|---|---|---|---|---|
| **10 Ready Nodes** | **+** | **7 RR = 7 potential lost nodes** | **=** | **17 Nodes Needed** | |

Extended action time

# Possible solutions: Solution

10
Ready
Nodes

17
Ready
Nodes

10
Ready
Nodes

Rebalance
Recommend
ation

Docplanner

# Possible solutions: How to calculate a boost

Resources Calculations  **+**  Rebalance Recommendation

Docplanner

# Possible solutions: Calculations. Not that easy

| Cluster Autoscaler Configmap | ASG Tags | Rebalance Recommendations | Nodes |
|---|---|---|---|
| ConfigMap Watcher | Tags Watcher | Events Watcher | Nodes Watcher |
| ASGPool | | EventPool | NodePool |

Docplanner

# Possible solutions: Calculations overview

ASGPool

EventPool

NodePool

Calculations

Send AWS Request

Docplanner

# Possible solutions: Drain problem

**10 Ready Nodes**

**17 Ready Nodes**

{

**10 Ready Nodes**

Docplanner

# Possible solutions: Batch Drain Solution

```
NodePool ──┐         ┌─────────────────┐      ┌─────────────────┐      ┌───────┐      ┌──────────────────────┐
           ├────────▶│ Recently Ready  │─────▶│ Get node under  │─────▶│ Drain │─────▶│        AWS:          │
EventPool ─┘         │      Node       │      │ risk from event │      │       │      │ Destroy old instance │
                     └─────────────────┘      └─────────────────┘      └───────┘      └──────────────────────┘
                                  Batch
                     ┌─────────────────┐      ┌─────────────────┐      ┌───────┐      ┌──────────────────────┐
                     │ Recently Ready  │─────▶│ Get node under  │─────▶│ Drain │─────▶│        AWS:          │
                     │      Node       │      │ risk from event │      │       │      │ Destroy old instance │
                     └─────────────────┘      └─────────────────┘      └───────┘      └──────────────────────┘
```

# Possible solutions: Drain Result

17
Ready
Nodes

10
Ready
Nodes

10
Ready
Nodes

Docplanner

# Will it work?

Case 1: Stable situation

# Will Work?: real case



Spots nodes over the time
Feb 23, 5:39am – Feb 23, 11:39am

Docplanner

# WILL WORK?: 1. INITIAL SPIKE OF CAPACITY



**Spots nodes over the time**
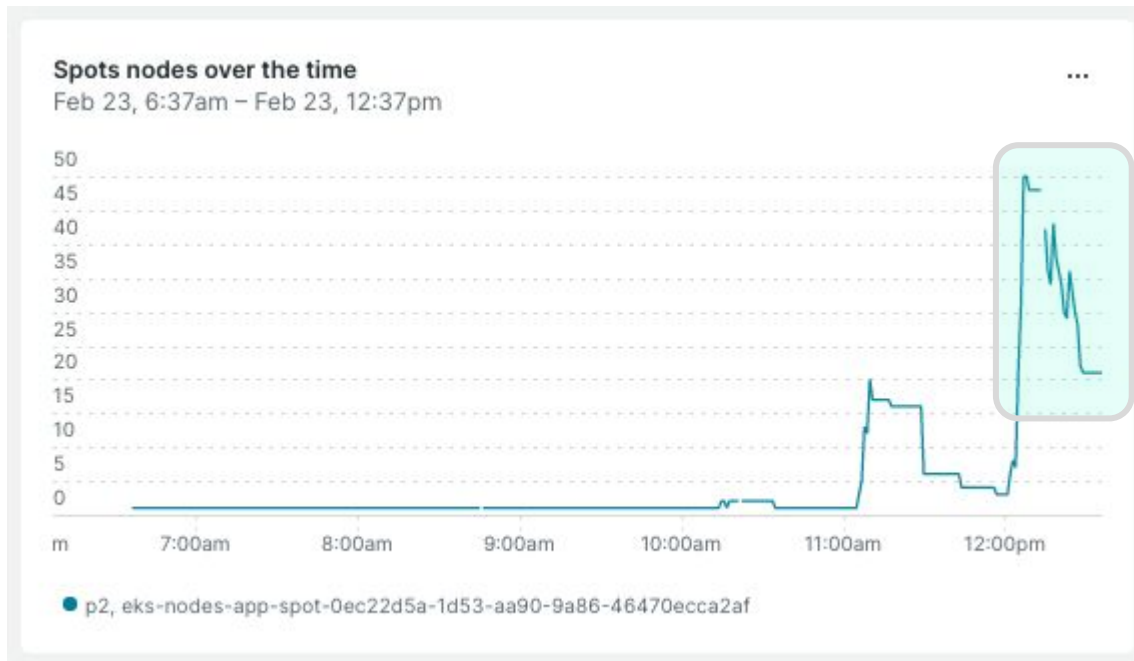Feb 23, 5:39am – Feb 23, 11:39am

Docplanner

# Will Work?: 2. Drain ladder

# Will Work?: 3. Cluster Autoscaler adjustments

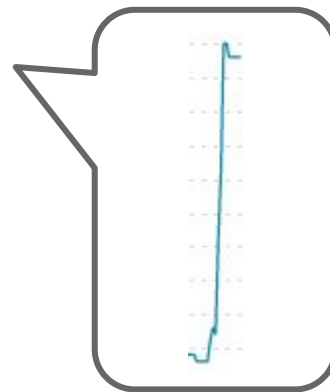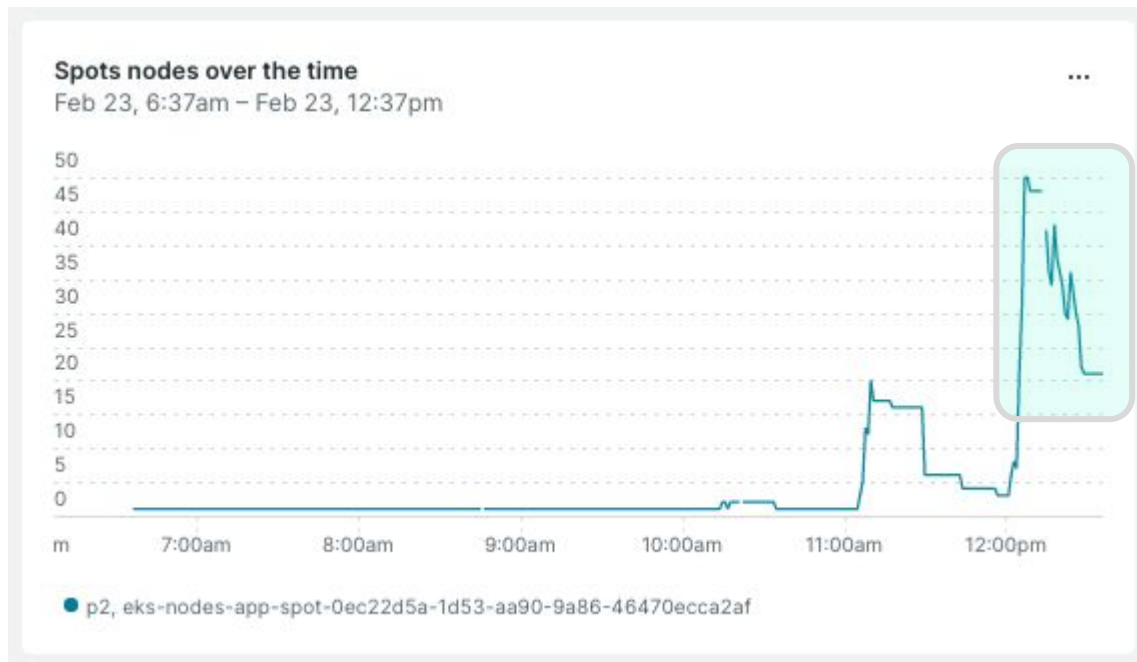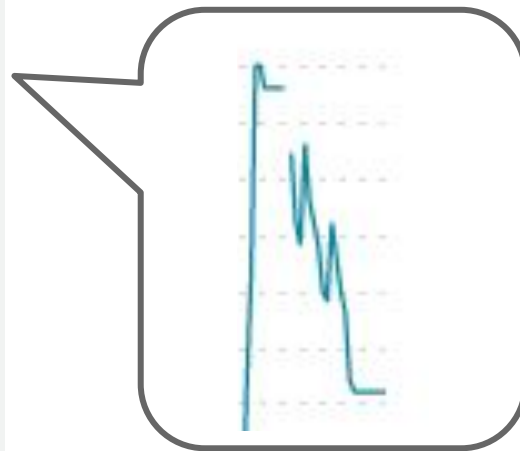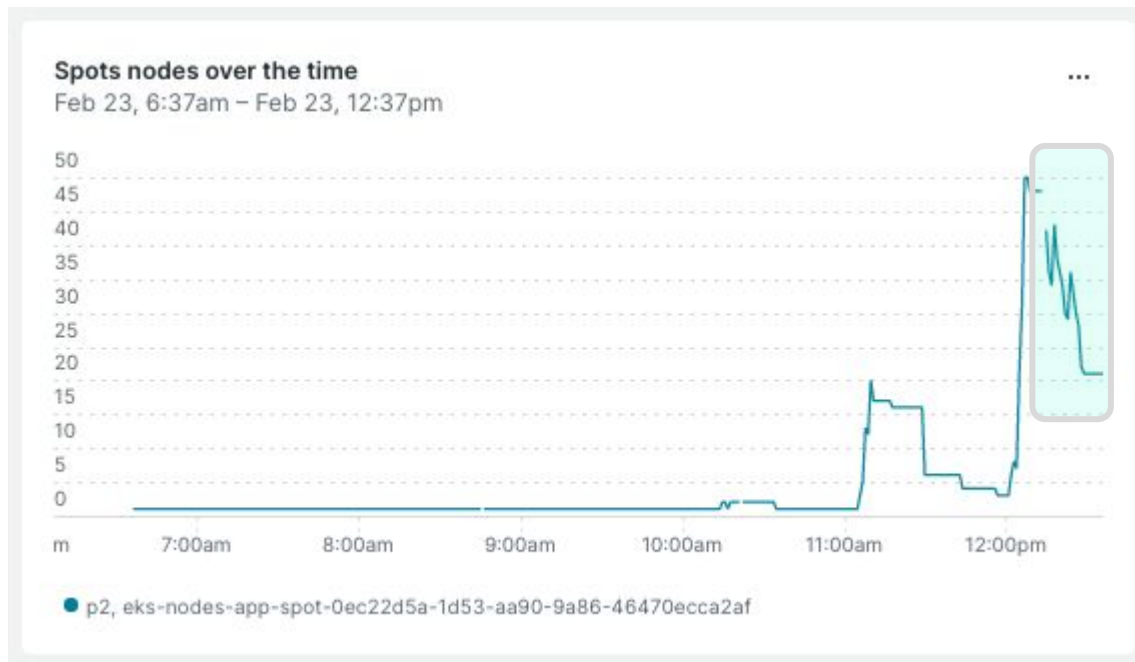# Will it work?

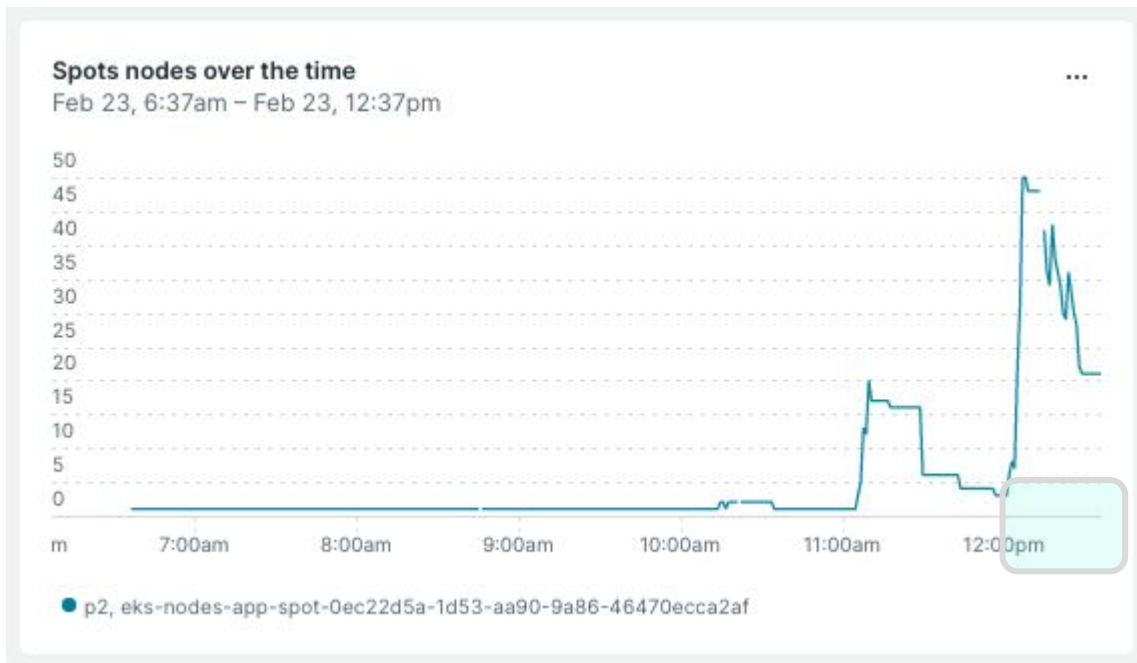Case 2: AWS being greedy

# Will Work?: Another Real case



Spots nodes over the time
Feb 23, 6:37am – Feb 23, 12:37pm

• p2, eks-nodes-app-spot-0ec22d5a-1d53-aa90-9a86-46470ecca2af

# Will Work?: 1. Initial Spike of capacity

# Will Work?: 2. Drain... ladder?

# WILL WORK?: 2. Drain... ladder?



**Spots nodes over the time**
Feb 23, 6:37am – Feb 23, 12:37pm

- p2, eks-nodes-app-spot-0ec22d5a-1d53-aa90-9a86-46470ecca2af

– Bad config:
Nodes are considered
new since 10 minutes
ago

– Real case:
30 minutes ago a lot
of nodes were joined

Docplanner