# Motivation

Self-adaptive system, deployed on K8



DevOps Engineer

## Motivation



A violates response time SLO

Still acceptable workload for B as A cannot process all incoming requests in time
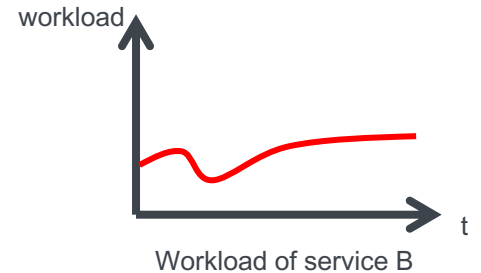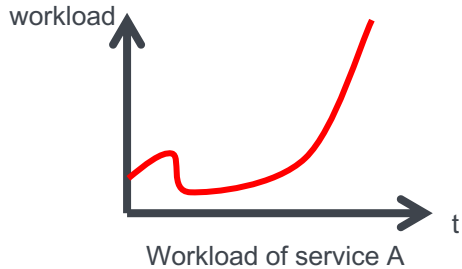
A

B

A violates CPU load SLO

C

workload

t

Workload of service A

workload

t

Workload of service B

## Motivation



A handles more requests

A

B

C

Each instance of A has lower CPU load

workload

Workload of service A

t

workload

Workload of service B

t

# Motivation



B violates response time SLO

B violates CPU load SLO

A

B

C

workload

Workload of service A

t

workload

Workload of service B

t

## Motivation

# Motivation



B handles more requests

Each instance of B has lower CPU load

A

B

C

workload

Workload of service A

t

workload

Workload of service B

t

Towards a Visual Explanation System for Self-Adaptation Events on Kubernetes

# Motivation

# Current Solutions

## Ereignisse

| Nachricht | Quelle | Sub-Objekt | Anzahl | Zuerst gesehen | Zuletzt gesehen |
|---|---|---|---|---|---|
| Scaled up replica set orchestrator-855c5c64df to 3 | deployment-controller | - | 2 | 2 hours ago | 24 minutes ago |
| New size: 3; reason: external metric tomatotraffictomato(&LabelSelector{MatchLabels:map[string prometheus,},MatchExpressions:[]LabelSelectorRequirement{ target | horizontal-pod-autoscaler | - | 1 | 24 minutes ago | 24 minutes ago |

Disorganized

Information is only saved for 1h (with default settings)

Only metric name is displayed and no other information about metric visible

Depending on the tool large sum of logs at different places to look at

# Current Solutions

```
Metrics:                                     ( current / target )
  "tomatocputomato" (target average value):    813m / 150m
  "tomatolatencytomato" (target average value): 17m / 150m
  "tomatoerrortomato" (target average value):   0 / 5
  "tomatotraffictomato" (target average value): 25m / 5
Min replicas:                                1
Max replicas:                                4
Deployment pods:                             4 current / 4 desired

  Type            Status   Reason            Message
  ----            ------   ------            -------
  AbleToScale     True     ReadyForNewScale  recommended size matches current size
  ScalingActive   True     ValidMetricFound  the HPA was able to successfully calculate a replica count from external metric tomatocputomato(&LabelSe
ector{MatchLabels:map[string]string{type: prometheus,},MatchExpressions:[]LabelSelectorRequirement{},})
  ScalingLimited  True     TooManyReplicas   the desired replica count is more than the maximum replica count

  Type     Reason            Age      From                          Message
  ----     ------            ----     ----                          -------
  Normal   SuccessfulRescale 2m6s     horizontal-pod-autoscaler     New size: 4; reason: external metric tomatotraffictomato(&LabelSelector{MatchLabels:map
[string]string{type: prometheus,},MatchExpressions:[]LabelSelectorRequirement{},}) above target
```

Metric value and conditions only accessible through CLI or APM tools, e.g., Prometheus

Does not show the effect of scaling decision

# Problem

There is a lack of clear explanations for adaptation decisions, which increases the time spent understanding/verifying the scaling behavior of self-adaptive systems

# Explainability Questions
Elicited using an expert survey

How much did it scale?

When did it scale?

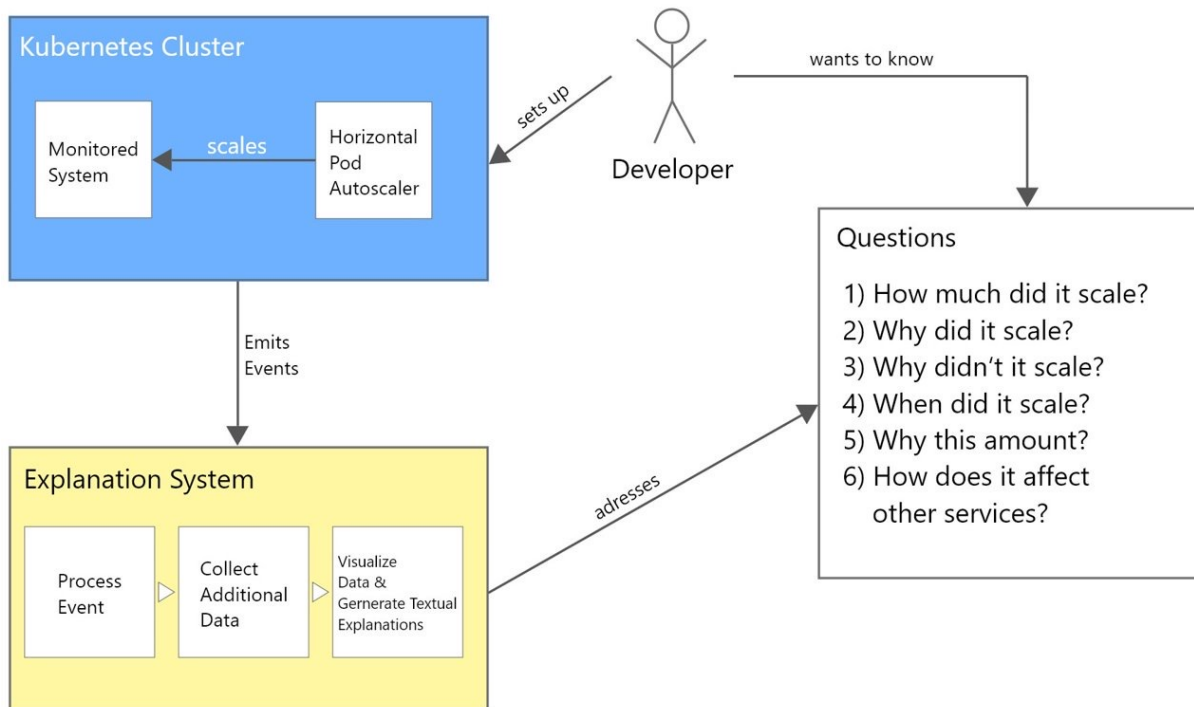Why did it scale?

Why this amount?

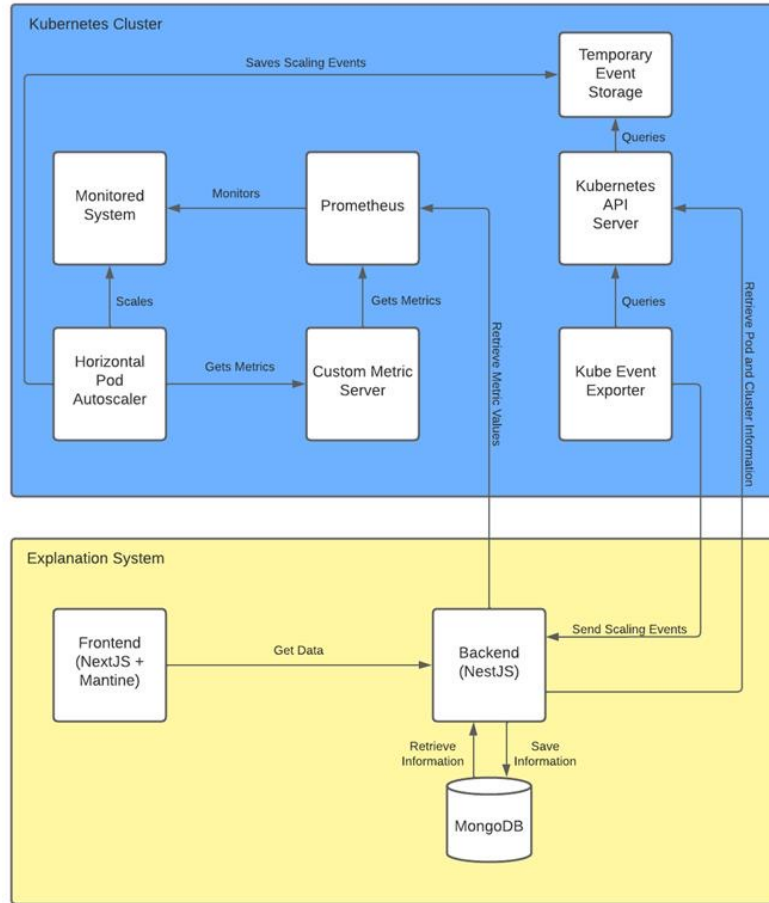Why didn't it scale?

How does it affect other services?

# Objective

# Architecture

# Dashboard

# Metrics

**EXPA**utoscaling

- ▦ Dashboard
- ▭ Analytics
- ⊙ Metrics

## Metrics

| Metric Name | Targeted Deployment | Metric Query | Target Value | Type | Min Replicas | Max Replicas | Created At |
|---|---|---|---|---|---|---|---|
| cpu | default/cart | sum(rate(container_cpu_usag | 150m | AverageValue | 1 | 4 | Mon, 06 Feb 2023 15:02:09 GMT |
| error | default/cart | rate(http_server_requests_s | 5 | AverageValue | 1 | 4 | Mon, 06 Feb 2023 15:02:09 GMT |
| latency | default/cart | sum(rate(http_server_reques | 150m | AverageValue | 1 | 4 | Mon, 06 Feb 2023 15:02:09 GMT |
| traffic | default/cart | sum(rate(http_server_reques | 5 | AverageValue | 1 | 4 | Mon, 06 Feb 2023 15:02:09 GMT |
| cpu | default/creditinstitute | sum(rate(container_cpu_usag | 150m | AverageValue | 1 | 4 | Mon, 06 Feb 2023 15:02:09 GMT |
| error | default/creditinstitute | rate(http_server_requests_s | 5 | AverageValue | 1 | 4 | Mon, 06 Feb 2023 15:02:09 GMT |
| latency | default/creditinstitute | sum(rate(http_server_reques | 150m | AverageValue | 1 | 4 | Mon, 06 Feb 2023 15:02:09 GMT |
| | | | | | | | Mon, 06 Feb |

# Event Summary



**EXPA**utoscaling

- Dashboard
- Analytics
- Metrics

## Events

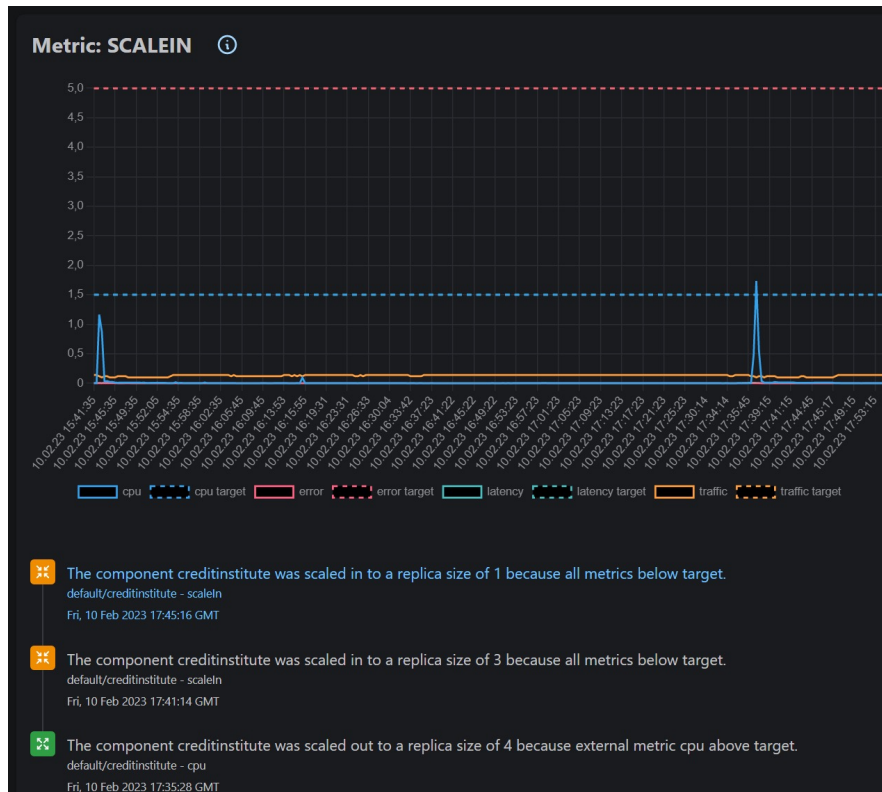| Timestamp | Deployment | Scaling type | Metric | Replica size | Status |
|---|---|---|---|---|---|
| Fri, 24 Feb 2023 16:44:51 GMT | default/orchestrator | scaleIn | scaleIn | 1 ⟼ 2 | SuccessfulRescale |
| Fri, 24 Feb 2023 16:43:54 GMT | default/creditinstitute | scaleIn | scaleIn | 1 ⟼ 3 | SuccessfulRescale |
| Fri, 24 Feb 2023 16:43:51 GMT | default/cart | scaleIn | scaleIn | 1 ⟼ 3 | SuccessfulRescale |
| Fri, 24 Feb 2023 16:43:50 GMT | default/uibackend | scaleIn | scaleIn | 1 ⟼ 3 | SuccessfulRescale |
| Fri, 24 Feb 2023 16:43:50 GMT | default/payment | scaleIn | scaleIn | 1 ⟼ 3 | SuccessfulRescale |
| Fri, 24 Feb 2023 16:43:50 GMT | default/order | scaleIn | scaleIn | 1 ⟼ 3 | SuccessfulRescale |
| Wed, 15 Feb 2023 15:14:55 GMT | default/creditinstitute | scaleIn | scaleIn | 3 ⟼ 1 | SuccessfulRescale |
| Wed, 15 Feb 2023 15:14:54 GMT | default/cart | scaleIn | scaleIn | 1 ⟼ 1 | SuccessfulRescale |
| Wed, 15 Feb 2023 15:14:53 GMT | default/orchestrator | scaleIn | scaleIn | 3 ⟼ 1 | SuccessfulRescale |
| Wed, 15 Feb 2023 15:14:53 GMT | default/order | scaleIn | scaleIn | 3 ⟼ 1 | SuccessfulRescale |
| Wed, 15 Feb 2023 15:14:53 GMT | default/inventory | scaleIn | scaleIn | 2 ⟼ 1 | SuccessfulRescale |
| Wed, 15 Feb 2023 15:14:52 GMT | default/uibackend | scaleIn | scaleIn | 1 ⟼ 1 | SuccessfulRescale |
| Wed, 15 Feb 2023 15:14:52 GMT | default/payment | scaleIn | scaleIn | 1 ⟼ 1 | SuccessfulRescale |
| Wed, 15 Feb 2023 15:07:06 GMT | default/order | scaleOut | cpu | 3 ⟼ 3 | SuccessfulRescale |

# 📏 How much did it scale?

# Why did it scale? / ✖ Why didn't it scale?

# Why did it scale?  /  ✖ Why didn't it scale?

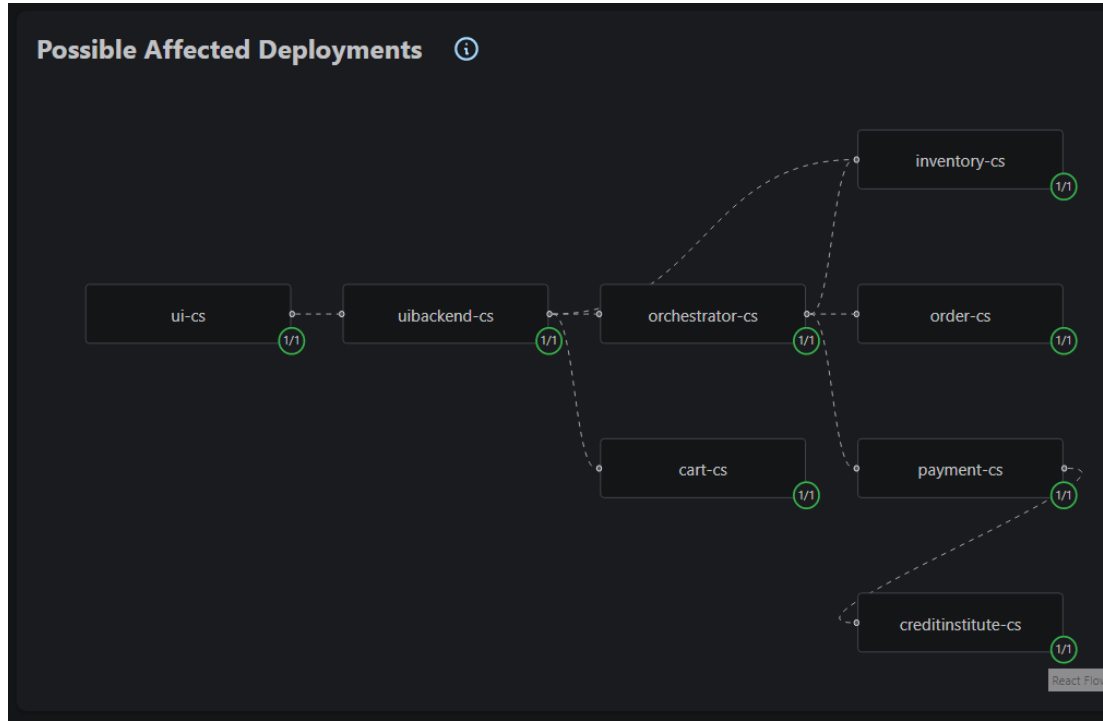# Why did it scale?  /  ✕ Why didn't it scale?

# ? Why this amount?

**Desired Replicas** ⓘ

desiredReplicas = ceil[currentReplicas * (currentMetricValue/desiredMetricValue)]

= ceil[3 * (0.01/1.5)]

= 1

The component was scaled to a size of 1 because the currentMetricValue was smaller than the desiredMetricValue, so only 0.67% of the 3 replicas were needed.
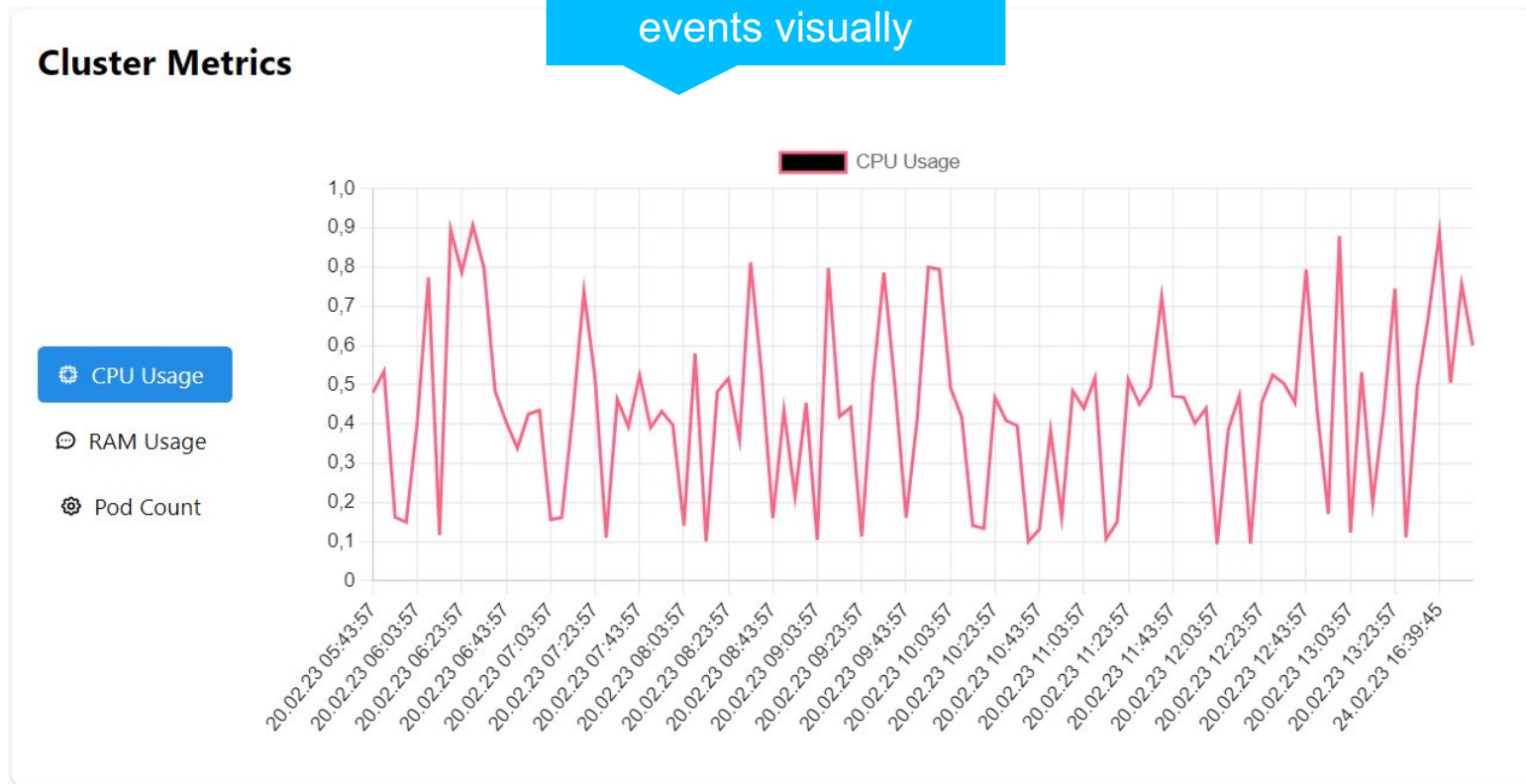
# How does it affect other services?

# Cluster Metrics

# Preliminary Evaluation based on Experiment & Explainability Questions
Using an expert survey

# Preliminary Evaluation based on Experiment & Explainability Questions
Using an expert survey

EXPA answers questions satisfactorily

How much did it scale?

Why did it scale?

Why didn't it scale?

When did it scale?

Why this amount?

EXPA answers questions partially

How does it affect other services?

# Limitations and Threats to Validity

Only addresses reactive autoscaling

Experts in interview were selected based on availability and convenience

Survey does not cover all the scenarios and challenges regarding scaling in Kubernetes clusters

EXPA cannot identify and show causal chains of events right now

Does not provide interactive explanations (yet)

# Discussion

What would you like to see?

Which information do you require?

Any ideas for a better visual explainability?

**University of Stuttgart**
Institute of Software Engineering (ISTE)
Software Quality and Architecture Group (SQA)

X @spethso

/in/sandro-Speth

@SandroSpeth

# Thank you!

**Sandro Speth**

e-mail   sandro.speth@iste.uni-stuttgart.de

phone   +49 (0) 711 685-61693

www.    iste.uni-stuttgart.de/sqa/team/Speth

University of Stuttgart
Institute of Software Engineering,
Software Quality and Architecture Group

Universitätsstraße 38,
70569 Stuttgart
Room 1.336

# EXPAutoscaling

Backend

https://github.com/lMaxTl/explaining-autoscaling-backend

Frontend

https://github.com/lMaxTl/explaining-autoscaling-frontend

Open Source available under MIT License