# Level Up your Kubernetes Scaling with KEDA

Wolfgang Ofner

# Agenda

- Architecture in SW projects

- Introduction to KEDA

- Scaling with messages in Azure Service Bus Queue

- KEDA Conclusion

- Q&A

# About Me

Senior Software Architect, bbv Software, Zürich

Consultant and Speaker

Focus on Azure, Kubernetes, DevOps and .NET

https://programmingwithwolfgang.com
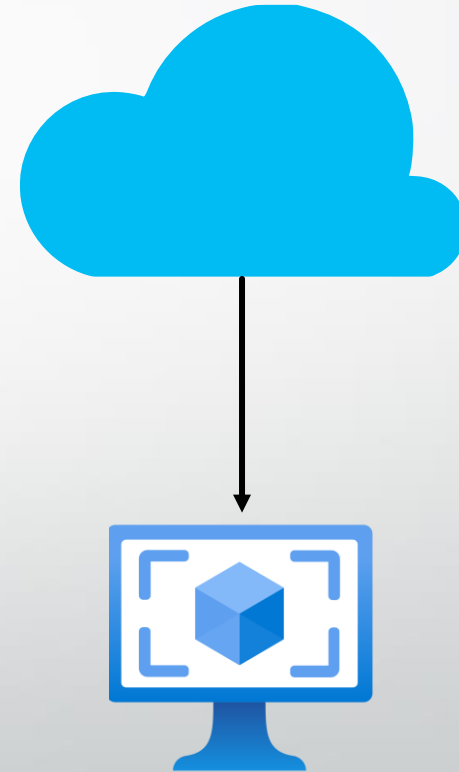
https://www.linkedin.com/in/wolfgangofner

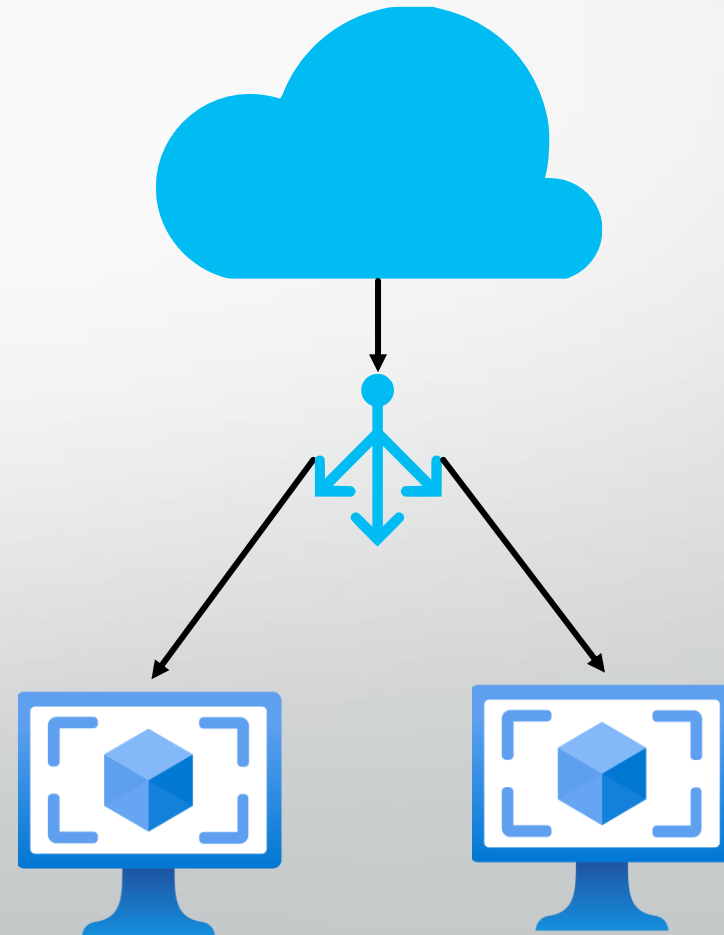https://twitter.com/wolfgang_ofner

# Simplified Architecture History

- Server – Client Architecture
- Only few clients
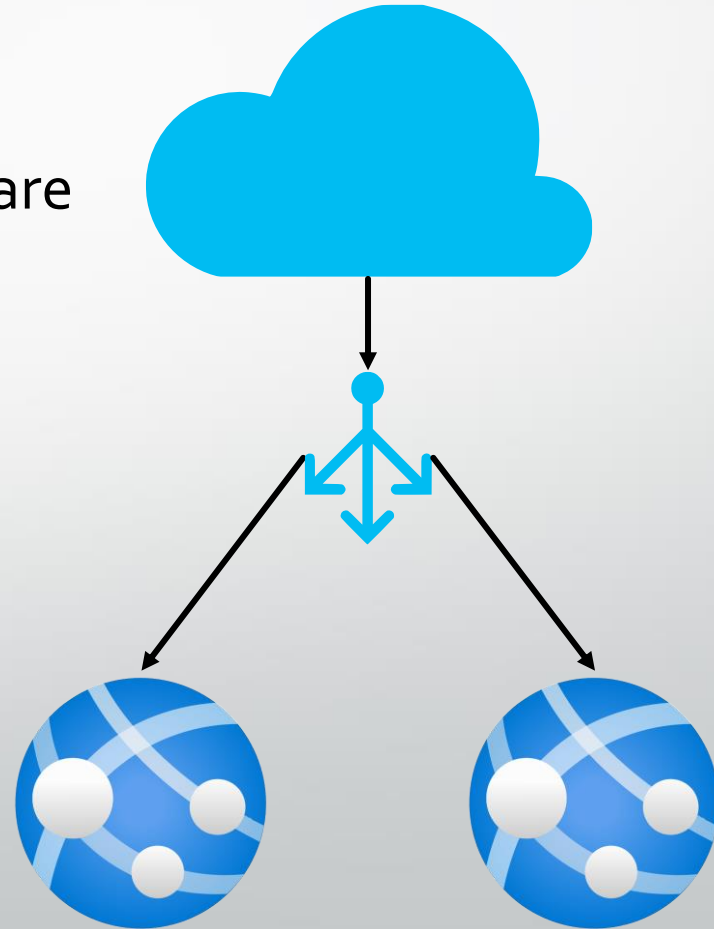- No redundancy
- No high availability

# Simplified Architecture History

- Static load balancing

- New VMs need to be added by hand
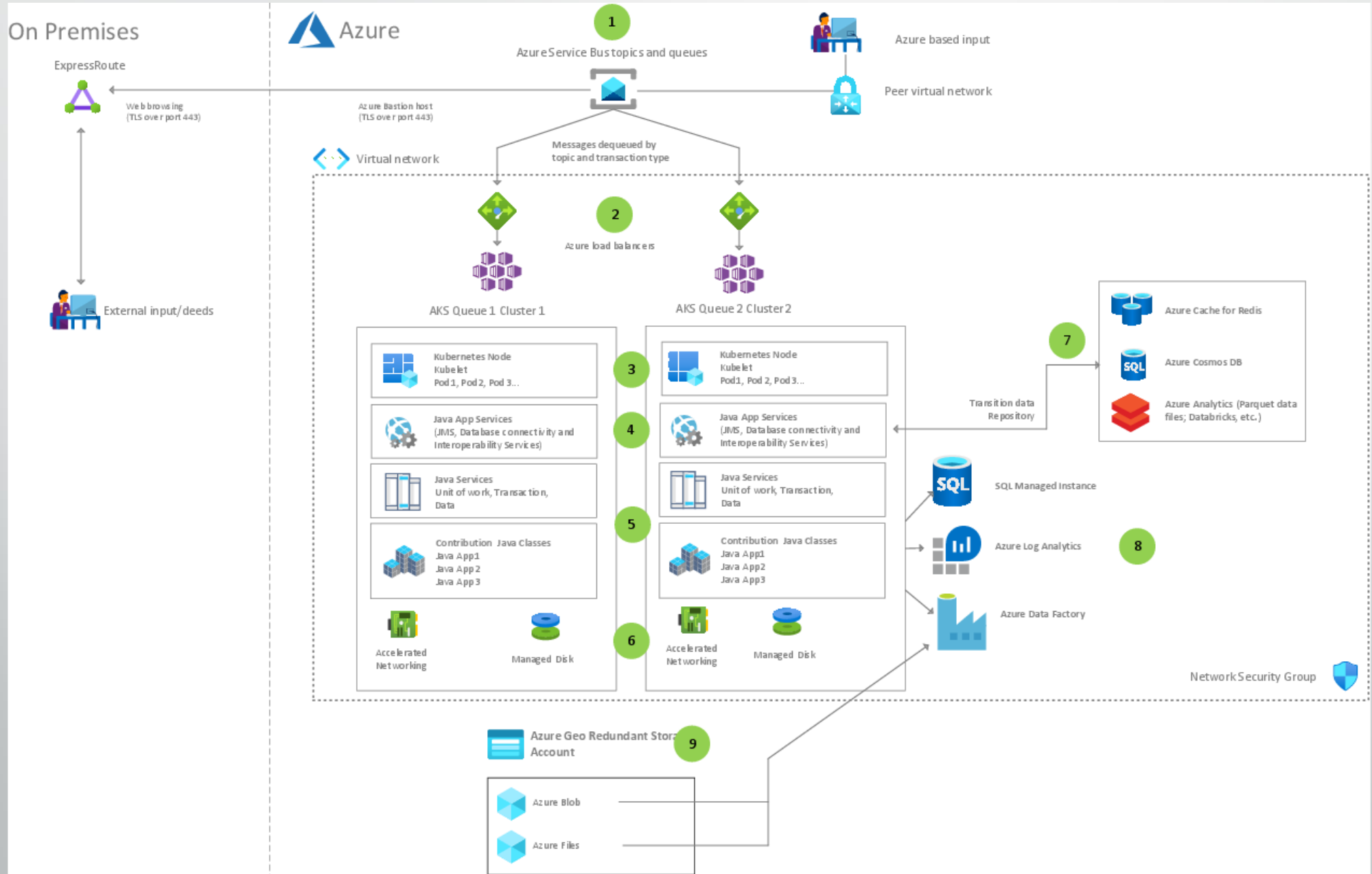
- Expensive on-premises hardware

# Simplified Architecture History

- Automatically adding additional hardware

- Pay only what you need

- Mostly CPU or RAM based scaling

# Modern Architecture

# Kubernetes

- Horizontal Pod Autoscaler (HPA)
  - Scaling according to CPU and/or RAM
- Architectures get more and more complex
- Dependencies on external components
- Applications have to react to events
  - Database
  - Service Bus
  - Streams

# Horizontal Pod Autoscaler

- Scales Deployments or StatefulSets

- Adds or removes pods

- Scaling based on CPU or RAM usage

- Scaling on custom metrics

  - Query custom metrics from Kubernetes API

  - Prometheus

  - requests per second

# Horizontal Pod Autoscaler Configuration

```yaml
apiVersion: autoscaling/v1
kind: HorizontalPodAutoscaler
metadata:
  name: customerapi
  namespace: customerapi-test
spec:
  maxReplicas: 10
  minReplicas: 1
  averageCpuUtilization: 50
  scaleTargetRef
    apiVersion: apps/v1
    kind: Deployment
    name: customerapi

behavior:
  scaleDown:
    policies:
    - type: Pods
      value: 4
      periodSeconds: 60
    - type: Percent
      value: 10
      periodSeconds: 60
    selectPolicy: Min

scaleUp:
  policies:
  - type: Pods
    value: 5
    periodSeconds: 60
  - type: Percent
    value: 12
    periodSeconds: 60
  selectPolicy: Max
```

# Limitation of the HPA

- Black Friday

- Thousands of orders are stored in a queue

- Scaling using CPU or RAM is not sufficient

- No option for scaling in this scenario

# KEDA – Kubernetes Event-driven Autoscaling

- Kubernetes Event-driven Autoscaling

- Open source

- CNCF Project

- Maintained by

  - Docplanner Tech

  - Microsoft

  - Red Hat

# KEDA

- ~53 built-in Scaler
  - Apache Kafka
  - Azure Blob Storage
  - Azure Monitor
  - Azure Service Bus
  - Elastic Search
  - MongoDB
  - Prometheus
  - Redis Streams

# KEDA Use Cases

- Scale according to external events

- Scale to Zero

  - Bring serverless to your datacenter

  - Recreate Azure Functions architecture

  - Better resource usage

# KEDA Installation

- Installation via Helm charts
- Namespace: keda

# KEDA Installation

kubectl create namespace keda

helm repo add kedacore https://kedacore.github.io/charts

helm repo update

helm install keda kedacore/keda --namespace keda

# KEDA Resources

# KEDA Architecture

- 2 components for KEDA
  - Agent or Operator
  - Metrics Server
- Uses HPA for scaling
- Seamless integration into existing architecture

# KEDA Architecture

- 2 components for KEDA
  - Agent
  - Metrics Server
- Uses HPA for scaling
- Seamless integration into existing architecture
- 2 custom K8s resources for scaler
  - ScaledObject
  - TriggerAuthentication

# ScaledObject

```yaml
apiVersion:
keda.sh/v1alpha1

kind: ScaledObject

metadata:

  name: kedademoapi-scaler
```

```yaml
spec:

  scaleTargetRef:

    name: kedademoapi

  minReplicaCount: 0

  maxReplicaCount: 10

  pollingInterval: 30

  cooldownPeriod: 30
```

```yaml
triggers:

  - type: azure-servicebus

  metadata:

    queueName: KedaDemo

    queueLength: '5'

  authenticationRef:

    name: trigger-
authentication-kedademoapi
```

# TriggerAuthentication

```yaml
apiVersion: keda.sh/v1alpha1
kind: TriggerAuthentication
metadata:
  name: trigger-authentication-kedademoapi
spec:
  secretTargetRef:
  - parameter: connection
    name  kedademoapi-connectionstrings
    key: AzureServiceBus__ConnectionString
```

# Kubernetes Secret

```
PS C:\Users\Wolfgang> kubectl get secrets
NAME                                             TYPE                                  DATA   AGE
default-token-88lzb                              kubernetes.io/service-account-token   3      26h
kedademoapi-connectionstrings                    Opaque                                1      26h
kedademoapi-tls                                  kubernetes.io/tls                     2      26h
sh.helm.release.v1.kedademoapi-kedademoapi-test.v1   helm.sh/release.v1                1      26h
sh.helm.release.v1.kedademoapi-kedademoapi-test.v2   helm.sh/release.v1                1      22h
PS C:\Users\Wolfgang> kubectl describe secret kedademoapi-connectionstrings
Name:         kedademoapi-connectionstrings
Namespace:    kedademoapi-test
Labels:       app.kubernetes.io/managed-by=Helm
Annotations:  meta.helm.sh/release-name: kedademoapi-kedademoapi-test
              meta.helm.sh/release-namespace: kedademoapi-test


Type:  Opaque


Data
====
AzureServiceBus__ConnectionString:   165 bytes
```

# Kubernetes Secret

# Demo

- Scale with messages in an Azure Service Bus Queue
- Scale to 0
- Scale to 1

# kedademo (wolfgangkedademo/kedademo) | Service Bus Explorer

Service Bus Queue

Search (Ctrl+/)                          «

⟳  Refresh

Authentication type  ⓘ

Access key    Active Directory

Send    **Receive**    Peek

Receive performs a destructive read (ReceiveAndDelete) from Queue *kedademo*.
from the Queue. Messages shown here are no longer stored.

| Active | Dead-Lettered | Scheduled |
|---|---|---|
| **1** MESSAGES | **0** MESSAGES | **0** MESSAGES |

Please Select Queue or DeadLetter

⦿ Queue    ◯ DeadLetter

## Settings

Overview

Access control (IAM)

Diagnose and solve problems

Shared access policies

Service Bus Explorer (preview)

Properties

Locks

**POST** `/v1/ServiceBusProcessing` Action to add new messages to the queue.

## Parameters

| Name | Description |
|---|---|
| numberOfQueueItems<br>**integer($int32)**<br>*(query)* | 270 |

**Execute**

## Responses

**Curl**

```
curl -X 'POST' \
  'https://test.kedademo.programmingwithwolfgang.com/v1/ServiceBusProcessing?numberOfQueueItems=270' \
  -H 'accept: */*' \
  -d ''
```

**Request URL**

```
https://test.kedademo.programmingwithwolfgang.com/v1/ServiceBusProcessing?numberOfQueueItems=270
```

**Server response**

| Code | Details |
|---|---|
| 200 | **Response headers**<br><pre>content-length: 0<br>date: Fri,18 Feb 2022 15:45:21 GMT<br>strict-transport-security: max-age=15724800; includeSubDomains</pre> |

| GET | **/v1/ServiceBusProcessing** | Action to start processing the queue items. |
| --- | --- | --- |

**Parameters**

No parameters

<div style="text-align:center">Execute</div>

**Responses**

**Curl**

```
curl -X 'GET' \
  'https://test.kedademo.programmingwithwolfgang.com/v1/ServiceBusProcessing' \
  -H 'accept: application/json'
```

**Request URL**

```
https://test.kedademo.programmingwithwolfgang.com/v1/ServiceBusProcessing
```

**Server response**

| Code | Details |
| --- | --- |
| 200 | **Response body** |

```
271
```

**Response headers**

```
content-type: application/json; charset=utf-8
date: Fri,18 Feb 2022 15:51:31 GMT
strict-transport-security: max-age=15724800; includeSubDomains
```

```
PS C:\Users\Wolfgang> kubectl get pods
No resources found in kedademoapi-test namespace.
```

503 Service Temporarily Unavailable

test.kedademo.programmingwithwolfgang.com/index.html

# 503 Service Temporarily Unavailable

nginx

# kedademo (wolfgangkedademo/kedademo) | Service Bus Explorer

Service Bus Queue

🔍 Search (Ctrl+/)          «

↻  Refresh

▭▭ Overview

👥 Access control (IAM)

🔧 Diagnose and solve problems

**Settings**

🔑 Shared access policies

▦ Service Bus Explorer (preview)

⚙ Properties

🔒 Locks

---

Authentication type ⓘ

( **Access key**   Active Directory )

Send      Receive      Peek

Send Message to Queue *kedademo*

Content Type *

Text/Plain

new message

# ▦ kedademo (wolfgangkedademo/kedademo) | Service Bus Explorer

Service Bus Queue

🔍 Search (Ctrl+/)    «

⟳ Refresh

📧 Overview

Authentication type ⓘ

👥 Access control (IAM)

( **Access key**   Active Directory )

🔧 Diagnose and solve problems

Send    **Receive**    Peek

## Settings

🔑 Shared access policies

Receive performs a destructive read (ReceiveAndDelete) from Queue **kedademo**.
from the Queue. Messages shown here are no longer stored.

▦ Service Bus Explorer (preview)

| Active | Dead-Lettered | Scheduled |
| **1** MESSAGES | **0** MESSAGES | **0** MESSAGES |

⚙ Properties

Please Select Queue or DeadLetter

🔒 Locks

◉ Queue    ○ DeadLetter

```
PS C:\Users\Wolfgang> kubectl get pods
NAME                          READY   STATUS    RESTARTS   AGE
kedademoapi-6f986c4b76-b8pgj  1/1     Running   0          40s
```

Swagger
Supported by SMARTBEAR

# KedaDemo Api v1 OAS3

/swagger/v1/swagger.json

A simple API to read items from an Azure Service Bus Queue

Wolfgang Ofner - Website
Send email to Wolfgang Ofner

## ServiceBusProcessing

**GET** /v1/ServiceBusProcessing   Action to start processing the queue items.

**POST** /v1/ServiceBusProcessing   Action to add new messages to the queue.

# KEDA Scaling Logs

- keda-operator pod writes logs during scaling events

{"scaledobject.Name": "kedademoapi-scaler", "scaledObject.Namespace": "kedademoapi-test", "scaleTarget.Name": "kedademoapi", "Original Replicas Count": 6, "New Replicas Count": 0}
er kind": "ScaledObject", "name": "kedademoapi-scaler", "namespace": "kedademoapi-test"}
"scaledObject.Namespace": "kedademoapi-test", "scaleTarget.Name": "kedademoapi", "Original Replicas Count": 0, "New Replicas Count": 1}

# Limitations

- Scaler not available for used technology
- Cluster runs out of resources

# Pods

| | Name | Labels | | | Ready | | Phase | | Restarts | | Node | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ⋮ | ⓘ kedademoapi-6f986c4b76-2zfxc | app:kedademoapi | draft:draft-app | 1+ | 0/1 | | Pending | | 0 | | <not scheduled> | |
| ⋮ | ⓘ kedademoapi-6f986c4b76-6w9tc | app:kedademoapi | draft:draft-app | 1+ | 0/1 | | Pending | | 0 | | <not scheduled> | |
| ⋮ | ⓘ kedademoapi-6f986c4b76-777r8 | app:kedademoapi | draft:draft-app | 1+ | 0/1 | | Pending | | 0 | | <not scheduled> | |
| ⋮ | ✓ kedademoapi-6f986c4b76-9vs76 | app:kedademoapi | draft:draft-app | 1+ | 1/1 | | Running | | 0 | | aks-nodepool1-35436033-vmss000000 | |
| ⋮ | ⓘ kedademoapi-6f986c4b76-jdd8x | app:kedademoapi | draft:draft-app | 1+ | 0/1 | | Pending | | 0 | | <not scheduled> | |
| ⋮ | ✓ kedademoapi-6f986c4b76-mdj62 | app:kedademoapi | draft:draft-app | 1+ | 1/1 | | Running | | 0 | | aks-nodepool1-35436033-vmss000000 | |
| ⋮ | ⓘ kedademoapi-6f986c4b76-qg298 | app:kedademoapi | draft:draft-app | 1+ | 0/1 | | Pending | | 0 | | <not scheduled> | |
| ⋮ | ⓘ kedademoapi-6f986c4b76-rzgfm | app:kedademoapi | draft:draft-app | 1+ | 0/1 | | Pending | | 0 | | <not scheduled> | |
| ⋮ | ⓘ kedademoapi-6f986c4b76-s56q6 | app:kedademoapi | draft:draft-app | 1+ | 0/1 | | Pending | | 0 | | <not scheduled> | |
| ⋮ | ⓘ kedademoapi-6f986c4b76-wb7rr | app:kedademoapi | draft:draft-app | 1+ | 0/1 | | Pending | | 0 | | <not scheduled> | |

# Pods

| Name | Labels | Ready | Phase | Restarts | Node |
|---|---|---|---|---|---|
| ⓘ kedademoapi-6f986c4b76-2zfxc | app:kedademoapi draft:draft-app 1+ | 0/1 | Pending | 0 | \<not scheduled\> |
| ⓘ kedademoapi-6f986c4b76-6w9tc | app:kedademoapi draft:draft-app 1+ | 0/1 | Pending | 0 | \<not scheduled\> |
| ⓘ kedademoapi-6f986c4b76-777r8 | app:kedademoapi draft:draft-app 1+ | 0/1 | Pending | 0 | \<not scheduled\> |
| ✓ kedademoapi-6f986c4b76-9vs76 | app:kedademoapi draft:draft-app 1+ | 1/1 | Running | 0 | aks-nodepool1-35436033-vmss000000 |
| ⓘ kedademoapi-6f986c4b76-jdd8x | app:kedademoapi draft:draft-app 1+ | 0/1 | Pending | 0 | \<not scheduled\> |
| ✓ kedademoapi-6f986c4b76-mdj62 | app:kedademoapi draft:draft-app 1+ | 1/1 | Running | 0 | aks-nodepool1-35436033-vmss000000 |
| ⓘ kedademoapi-6f986c4b76-qg298 | app:kedademoapi draft:draft-app 1+ | 0/1 | Pending | 0 | \<not scheduled\> |
| ⓘ kedademoapi-6f986c4b76-rzgfm | app:kedademoapi draft:draft-app 1+ | 0/1 | Pending | 0 | \<not scheduled\> |
| ⓘ kedademoapi-6f986c4b76-s56q6 | app:kedademoapi draft:draft-app 1+ | 0/1 | Pending | 0 | \<not scheduled\> |
| ⓘ kedademoapi-6f986c4b76-wb7rr | app:kedademoapi draft:draft-app 1+ | 0/1 | Pending | 0 | \<not scheduled\> |

## Events

| Message | | Reason |
| --- | --- | --- |
| 0/1 nodes are available: 1 Insufficient cpu. | | FailedScheduling |

# Limitations

- Scaler not available for used technology
- Cluster runs out of resources
  - Azure Cluster Autoscaler
  - Define replica limit
  - Monitor cluster usage

# KEDA in Production

- Azure Container Apps use KEDA for scaling
  - Serverless containers
- KEDA 1.0.0 → 17. Nov 2019
- Currently 2.7.1
- Over 5k GitHub stars

# Resources

- Demo Application
  - https://github.com/WolfgangOfner/MicroserviceDemo/tree/master/KedaDemoApi
- Slides
  - https://github.com/WolfgangOfner/Presentation
- KEDA
  - https://keda.sh
- KEDA GitHub
  - https://github.com/kedacore/keda
- KEDA Architecture Screenshot
  - https://keda.sh/docs/2.6/concepts/#architecture

# Q&A

Level Up your Kubernetes Scaling with KEDA

Wolfgang Ofner

https://programmingwithwolfgang.com

https://www.linkedin.com/in/wolfgangofner

https://twitter.com/wolfgang_ofner