# Kubernetes Event-Driven Autoscaling (KEDA)

Making application autoscaling on Kubernetes dead-simple

*Zbynek Roubalik - Principal Software Engineer at Red Hat*
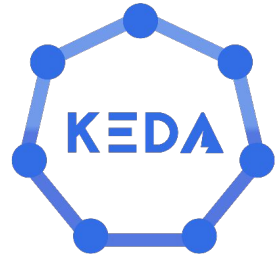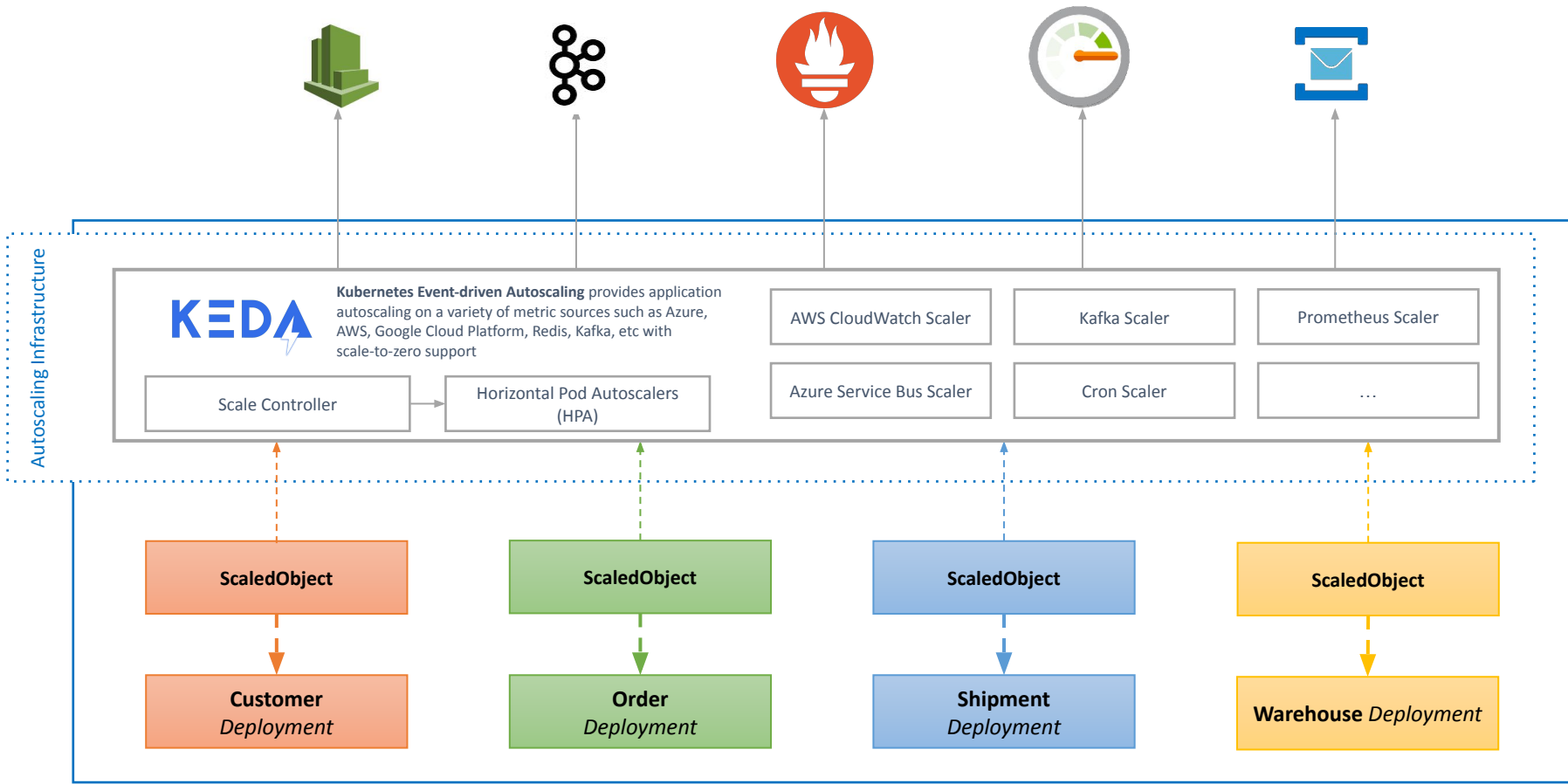*Tom Kerkhove - Senior Software Engineer at Microsoft*

# What is KEDA?

# Application autoscaling made simple with Kubernetes Event-driven Autoscaling (KEDA)

- Automatically scale Deployments, Jobs, /scale subresources

- Provides 55+ built-in scalers, but you can build your own
  - Support for external scaler, external push or Metrics API

- Production-grade authentication

- Save resources with scale to 0 or pause autoscaling

- Runs on Linux or ARM

**Focus on scaling your app, not the scaling internals**

KEDA

**Kubernetes Event-driven Autoscaling** provides application autoscaling on a variety of metric sources such as Azure, AWS, Google Cloud Platform, Redis, Kafka, etc with scale-to-zero support

Scale Controller → Horizontal Pod Autoscalers (HPA)

AWS CloudWatch Scaler

Kafka Scaler

Prometheus Scaler

Azure Service Bus Scaler

Cron Scaler

…

ScaledObject

ScaledObject

ScaledObject

ScaledObject

**Customer** *Deployment*

**Order** *Deployment*

**Shipment** *Deployment*

**Warehouse** *Deployment*
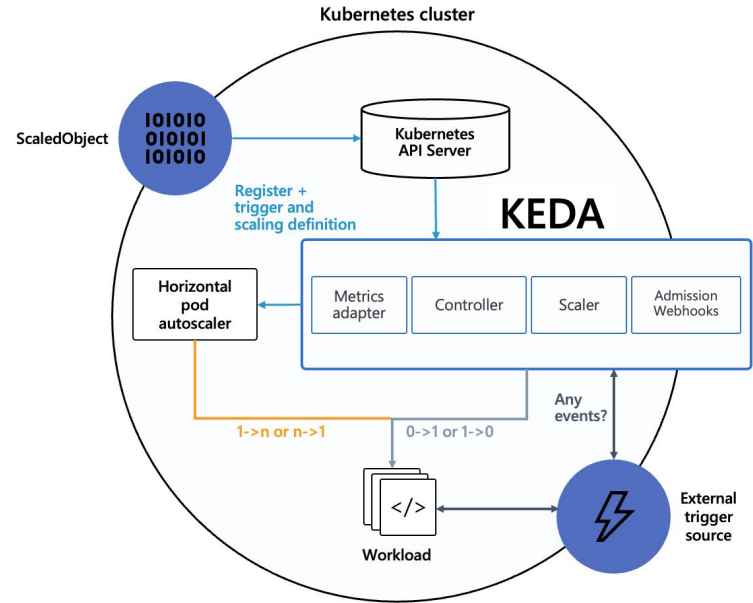
kubernetes

# How does it work?

## KEDA is built on top of Kubernetes

- Manages workloads to provide scale to 0
- Registers itself as a metric adapter
- Provides metrics for HPA to scale on

## Out-of-the-box & external scalers

## Easy to install

- Helm
- Operator Hub

# Production-grade authentication
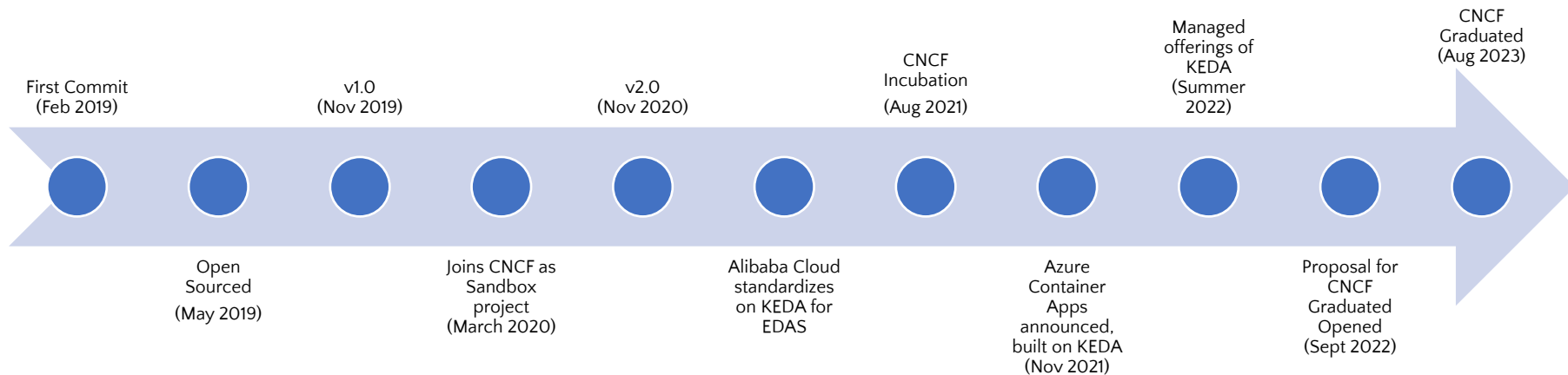
## Typical security concerns:

- Re-use secrets from scaled target – No separate identities

- Duplication of secrets – Harder to manage & rotate

## Re-use trigger authentication across `ScaledObject/ScaledJobs` with `TriggerAuthentication` (namespaced) or `ClusterTriggerAuthentication`

## Provides out-of-the-box integration with sources such as:

- Environment variables (on scale target)

- Kubernetes secrets

- Pod Identity ("No secret authentication" – Azure / AWS / EKS)
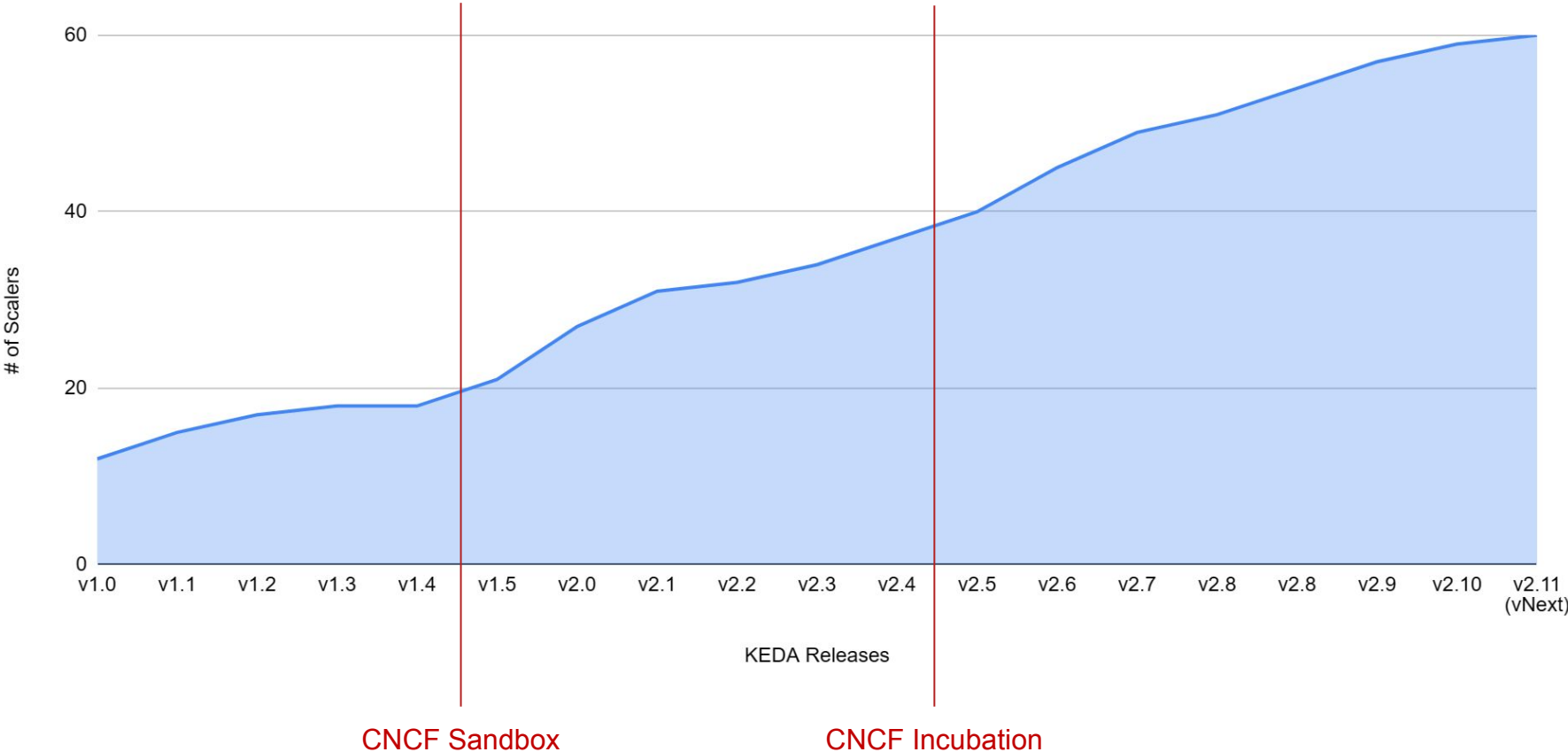
- HashiCorp Vault

- Azure Key Vault

# History of KEDA

First Commit
(Feb 2019)

Open
Sourced
(May 2019)

v1.0
(Nov 2019)

Joins CNCF as
Sandbox
project
(March 2020)

v2.0
(Nov 2020)

Alibaba Cloud
standardizes
on KEDA for
EDAS

CNCF
Incubation
(Aug 2021)

Azure
Container
Apps
announced,
built on KEDA
(Nov 2021)

Managed
offerings of
KEDA
(Summer
2022)

Proposal for
CNCF
Graduated
Opened
(Sept 2022)

CNCF
Graduated
(Aug 2023)

# What changed since CNCF Incubation?

# KEDA's Scaler Catalog Growth

# What did KEDA ship?

- **Autoscaling**
    - Support for 33 new scalers & 3 new authentication providers
    - Support for pausing autoscaling
    - Introduce admission webhooks to enforce autoscaling best-practices
    - POC with TAG Environmental Sustainability to reduce impact on environment
- **Artifacts & Deployment Scenarios**
    - Sign container images & migrate to GitHub Container Registry
    - Produce reproducible builds
    - Support for running on ARM machines
    - Support for non-public Azure clouds
- **Security**
    - Secure-by-default and runs as non-root
    - Support for identity segregation when using pod identities
    - Support for using custom CA for TLS
    - Extend security scanning suite for code & container images ([link](#))
    - Security Audit

# What did KEDA ship? (cont)

- **Operability & Production-readiness**
    - Provide operational metrics in Prometheus
    - Provide off-the-shelf Grafana dashboard for application autoscaling
- **Quality**
    - Provide chatops for running e2e tests in PR
    - Provide automation to manage test infrastructure
- **Governance**
    - Introduction deprecation & breaking change policy ([link](#))
    - Introduce scaler governance policy ([link](#))
    - Introduce Kubernetes Compatibility overview
    - Introduce roadmap and release cycle for release predictability ([link](#))

# KEDA on Artifact Hub

- Artifact Hub is the central place for cloud-native artifacts

- Build an ecosystem around external scalers

- Provide a better way to discover external scalers

- https://github.com/kedacore/external-scalers

- https://bit.ly/keda-artifact-hub

# KEDA 💘 Community

- 6,1k stars on GitHub

- ~260 contributors, incl.
  - Microsoft
  - Red Hat
  - Lidl
  - Reddit
  - IBM

- Bi-weekly community standups

# KEDA's Adoption Growth

- 42 listed end-users (+280% growth)
    - This is compared to when we opened our proposal for CNCF Incubation
    - And there are more exciting ones, which we cannot mention unfortunately

- 11.8% of Kubernetes users run KEDA (+151% growth)
    - Was 4.7% last year, based on CNCF Survey - Source

- Azure Container Apps, a cloud service built on top of KEDA, has become generally available

- Azure Kubernetes Service (AKS) and Red Hat OpenShift are offering a managed version of KEDA (preview)
    - Both managed offerings will offer support when they are generally available this summer.

# What's on the horizon

# Roadmap

**Introduce new scalers & secret sources**

- Azure IoT Hub, Kubernetes CRD, NATS Jetstream, Apache Pulsar, ...

**Provide CloudEvents to extend KEDA**

**Ship first-class support for HTTP-based autoscaling to GA**

**Allow multiple KEDA instances in one cluster** (depends on upstream)

**Allow run KEDA highly available** (depends on upstream)

**Expand our batteries-included approach with anti-pattern prevention**

- Currently supporting to only allow one ScaledObject/ScaledJob per scale target, don't allow arbitrary HPAs, etc

# Roadmap (cont)

**Embrace OpenTelemetry with scaler & runtime metrics**

**Capability to define autoscaling rules** (in progress)

**Historical analysis & predictive scaling** (considering)

- Considering to make it part of KEDA, but PredictKube is offering this as a service

**Bring KEDA's autoscaling engine outside of Kubernetes** (considering)

**Expand collaboration with CNCF's Environmental Sustainability TAG**

- Expand our "don't scale if it impacts our carbon neutrality too much, this is not high-prio" PoC

# The Autoscaling Sweetspot

# kubernetes

## Node #1

Customer App

Order App

Customer App

Order App

Customer App

## Node #2

Customer App

Customer App

Order App

Order App

Order App

## Node #3

Customer App

Customer App

Order App

### Infrastructure

**KEDA**

**Kubernetes Event-driven Autoscaling** provides application autoscaling on a variety of metric sources such as Service Mesh Interface, Service Bus, Monitor, etc.

| Scale Controller | → | Horizontal Pod Autoscalers (HPA) |

**Cluster Autoscaler** allows you to automatically add/remove nodes to the Kubernetes cluster.

**Virtual Kubelet & Virtual Nodes** allow you to run workloads outside of the Kubernetes cluster.

Order App

Customer App

Customer App

Order App

Order App

Customer App

Customer App

Order App