2024

# Scaling/Managing RTC Workloads in K8s

**Joe Mordica**

Founder/CTO @ VOXO

**Things we've learned operating a multi-region communications platform 100% in k8s**

voxo.co

# Agenda

# GitHub Repo



**github.com/voxoco/k8s**

# Visual Diagram



**voxo.co/technology**
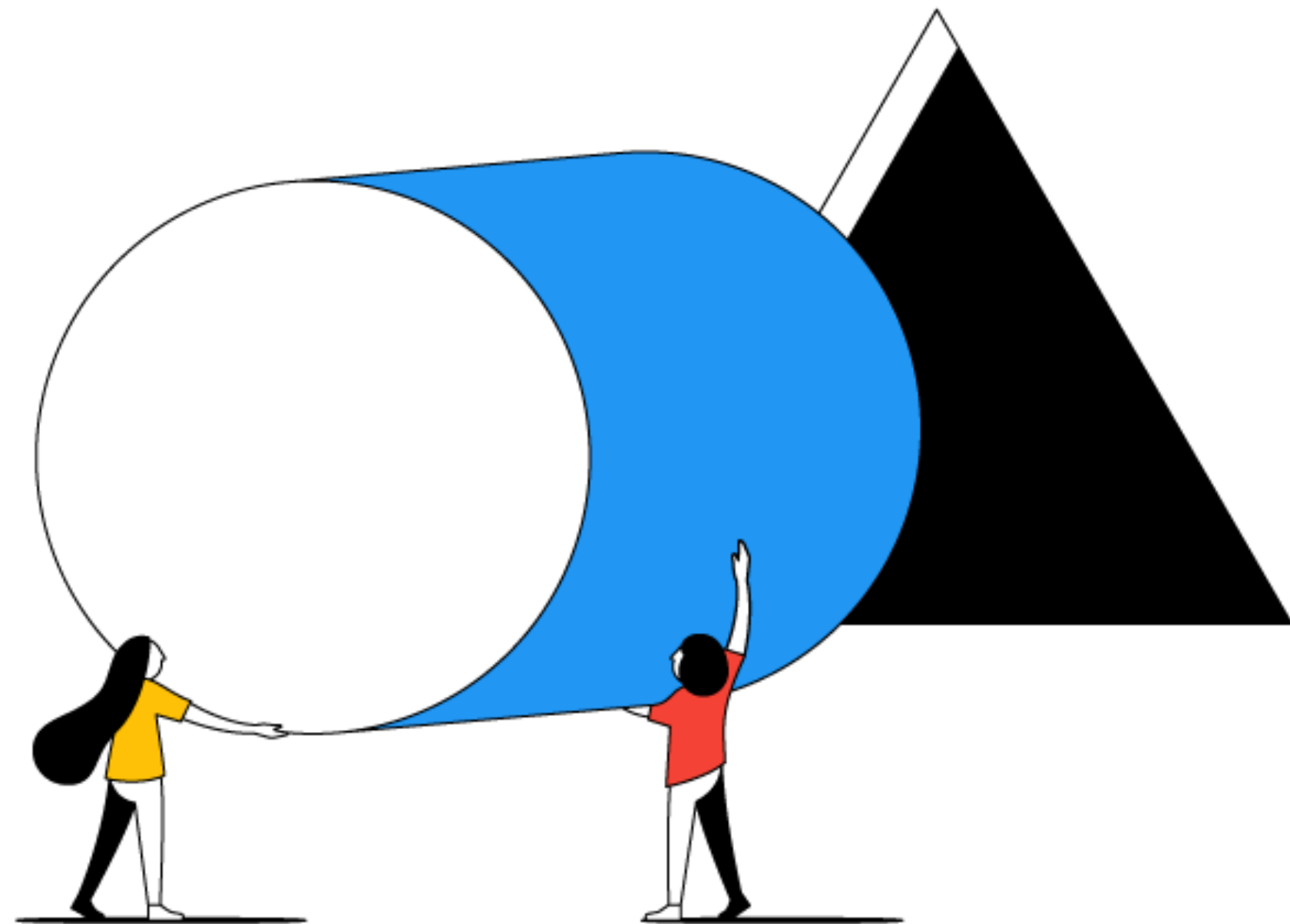
| **3+** | **0** | **0** |
|---|---|---|
| Regions | Single point of failure | Maintenance windows |

# Managing <u>stateful</u> workloads globally in the very <u>ephemeral</u> world of Kubernetes.

## Problems

**1** Plan for Failure (mindset shift)

**2** Dropped calls on shutdown

**3** Scaling

**4** Maintenance downtime


ONE DOES NOT SIMPLY DEPLOY KUBERNETES
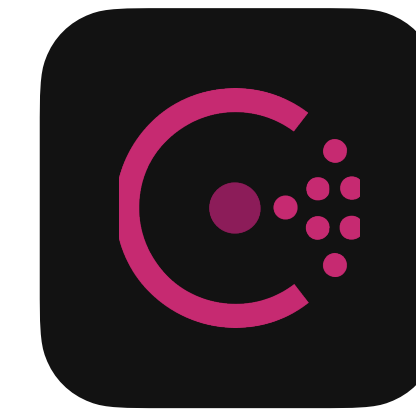
# Failure (Mindset shift)

### Kubernetes

Let Kubernetes do what it does best. Scale your apps and move stuff around based on CPU/Node health

### Ephemeral

- Make apps as ephemeral as possible. Expect them to fail and/or move between nodes anytime
- No PVC's.
- Orchestrator Example
- Use sidecar containers

### Global KV

We use Consul KV by HashiCorp as a global service discovery. Consul coupled with consul-template is awesome

# Dropped calls on shutdown

**2**

- ✓ **Graceful shutdowns**

- ✓ **Use sidecar containers**

- ✓ **Update KV store**

- ✓ **terminationGracePeriodSeconds**

- ✓ **Ship logs**

# Scaling

**3**

✓ **Let the HPA take care of it**

✓ **Tweak HPA scale up/scale down policy**

✓ **Set proper CPU requests on pods**

✓ **Application lifecycle needs to be dialed in**

✓ **Make small iterations and test**

# Maintenance downtime

**4**

- ✓ 0 downtime maintenance

- ✓ Create the ability to delete anything anytime

- ✓ Use rolling restarts for resources

- ✓ Use Google Global Load Balancer

- ✓ QA well

# Questions

joe@voxo.co

OUR WEBSITE

**voxo.co**