



# Highly scalable SaaS Apps on Oracle Kubernetes Engine (OKE) : Real Life Case Studies

Practical Experiences with OKE in Production at Scale

---

**Ram Kailasanathan**

Senior Director  
Product, Oracle  
Cloud Native

**Richard Bair**

Senior Director  
Engineering, Oracle IoT  
Apps and BlockChain



# Oracle Cloud Infrastructure: Complete set of services

**Governance**  
IAM, Tagging, Cost Analysis

**Security**  
IAM, Audit, KMS, CASB

**Management**  
Monitoring, Notifications, Alarms

**Automation**  
Resource Manager, Ansible

**Analytics / Integration / SOA Suite / Identity / Management / Content / API Platform / Developer / Visual Builder / Digital Assistant/ DataFlow / Data Science / Data Safe**

**Containers**  
**Containers and Kubernetes**



Fully managed, certified Kubernetes service with Docker containers

**Data Movement**  
**Storage appliance, Data Transfer**



Software NAS gateway, data ingest service with full chain of custody (HDD or appliance)

**Autonomous Database**  
**Transactions, Data Warehouse**



Fast provisioning. Automatic tuning, patching, securing. 99.995% availability.

**Cloud Native**  
**Events, Streaming, Functions**



Fully-managed FaaS, event-triggered functions, high-volume data ingest, notifications

**Compute**  
**Bare metal/VM, CPUs/GPUs**



Up to 64 CPU cores, 8 GPUs, 768 GB RAM, 51 TB local NVMe SSD, 5M IOPS, AMD and Intel processors

**Storage**  
**NVMe, Block, File, Object, Archive**



Predictable IOPS Block Storage for up to 98% less, storage for whole lifecycle

**Database**  
**Bare metal, VMs, Exadata**



Millions of TPS; Full RAC and Active Data Guard support

**Networking**  
**VCN, LBaaS, FastConnect, VPN**



Isolated networks with reserved IPs, security lists, firewalls, lowest cost private connectivity

**Public regions**

**Government regions**



# Oracle Strategy for Cloud Native

## Complete Cloud Native Stack

Deliver tools and services that are complete, integrated and open

- Continuous Integration & Deployment, Container Registry, Orchestration /Scheduling, Management /Operations, Analytics /Introspection
- With an application development platform for serverless and microservices

## Open Source

Actively participate in community driven open source container technologies

- Investing in Kubernetes, Docker, Fn, & CNCF, DevTools, and DevOps, with engineering resources, code contributions & sponsorships
- Active support from Oracle's portfolio of open source assets (Java, etc.)

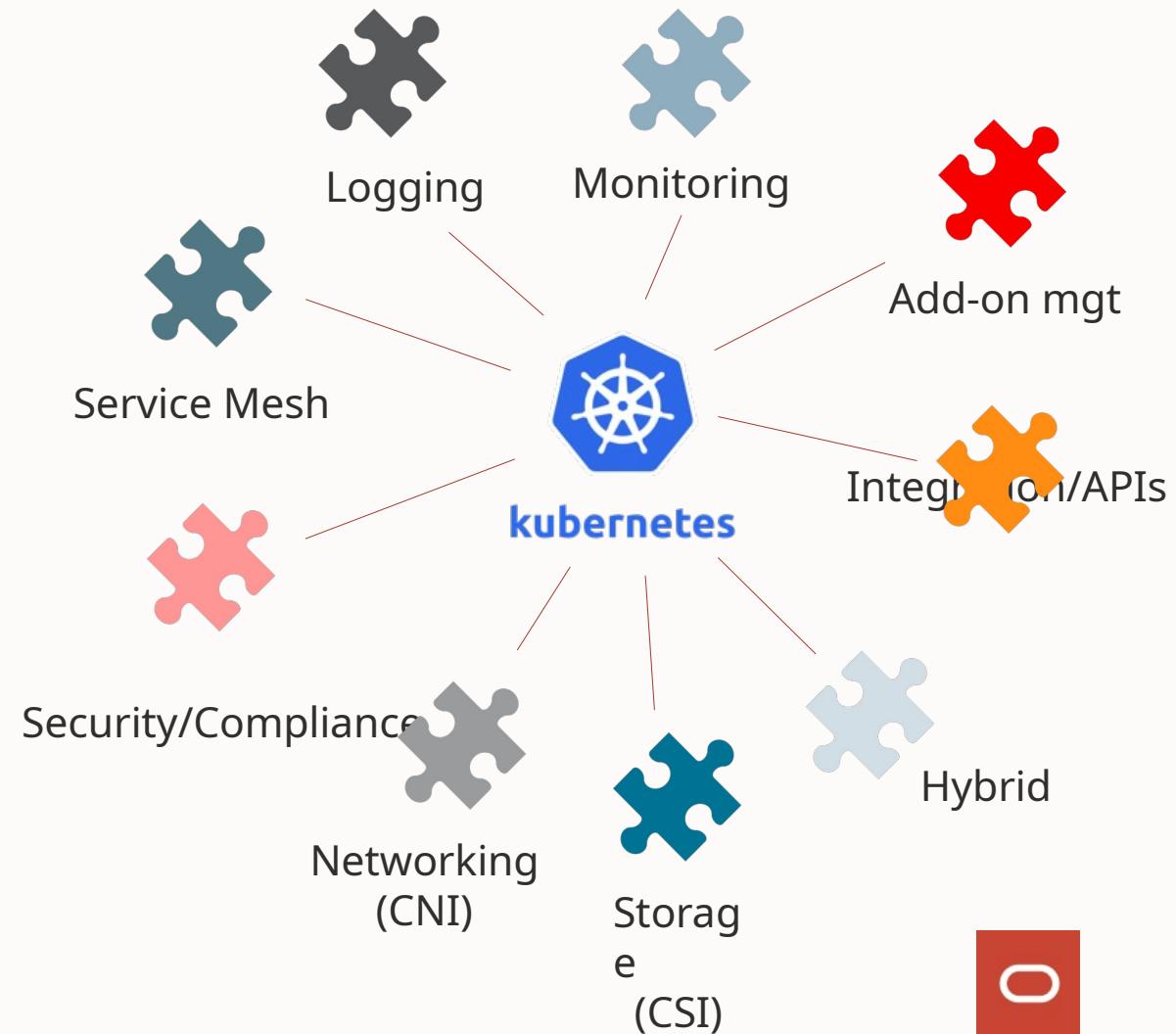
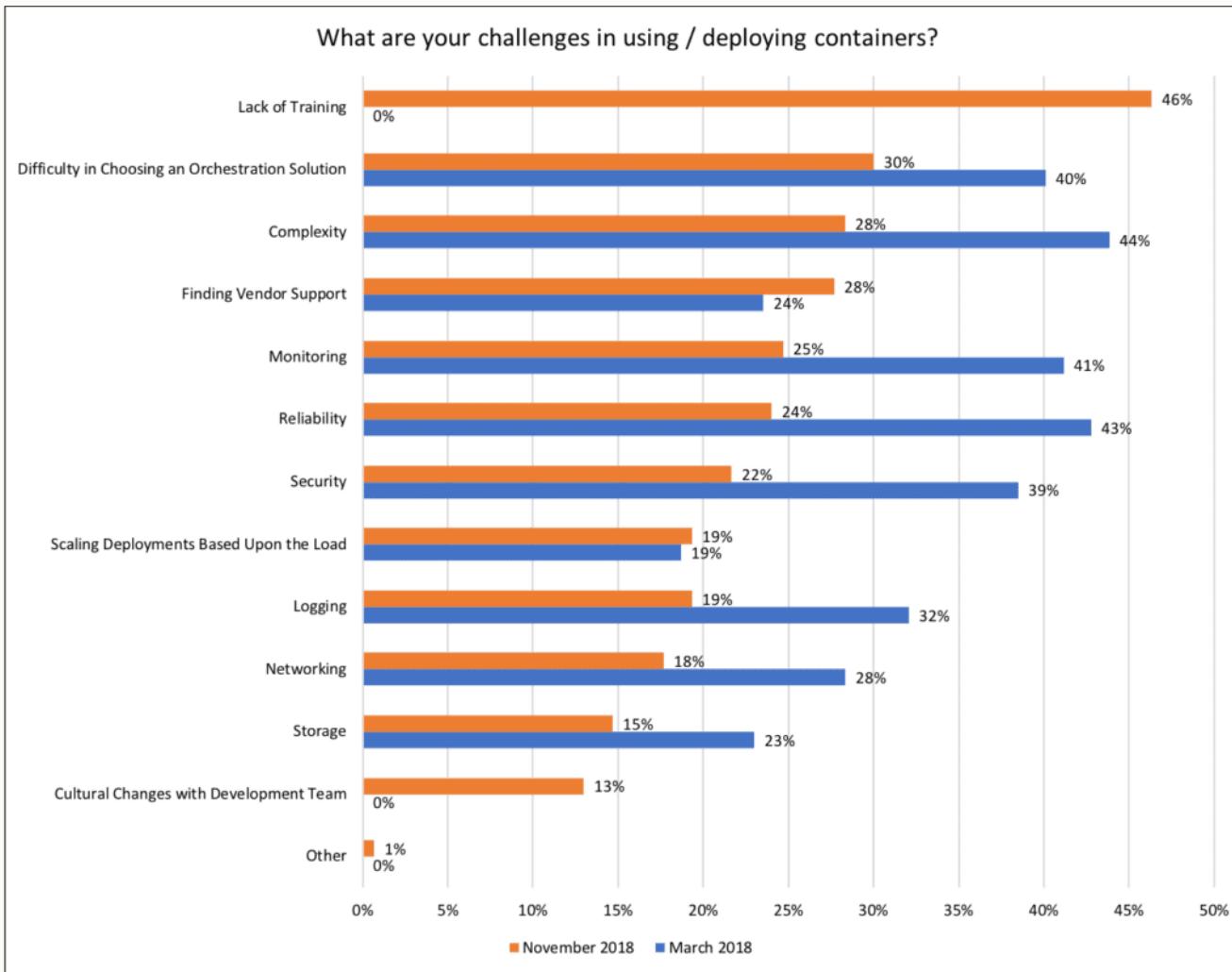
## Managed Services

Differentiate on quality of implementation, service and operational excellence

- Full, transparent management
- Open Source Compatible
- Standards **compliant**
- Deployed to Oracle Cloud Infrastructure
- Enterprise grade performance, **security**, HA, and governance



# Kubernetes - Necessary but not sufficient



Source: [CNCF Survey](#)

# Comprehensive Developer Services Portfolio



Easy to Start  
(Adopt)



App Development  
(Develop)



Deploy



Operate

Console

API Design\*

API Gateway

fluentd

Monitoring

QuickStart

SDK/CLI

Cloud Shell

Resource Manager

CI/CD\*

Streaming

Logging

Marketplace

Source Control\*

Container Registry

Container Engine  
(OKE)

Functions

Events

Always Free Tier

Cloud Native Java  
 helidon GraalVM ORACLE Coherence

CLOUD NATIVE COMPUTING FOUNDATION

cloudevents

Notification

Email

OCI Core Infrastructure Services

\*Under Development

# OKE @ Scale

---

Big Clusters, and a lot of 'em



# 17,000+

—  
Cores

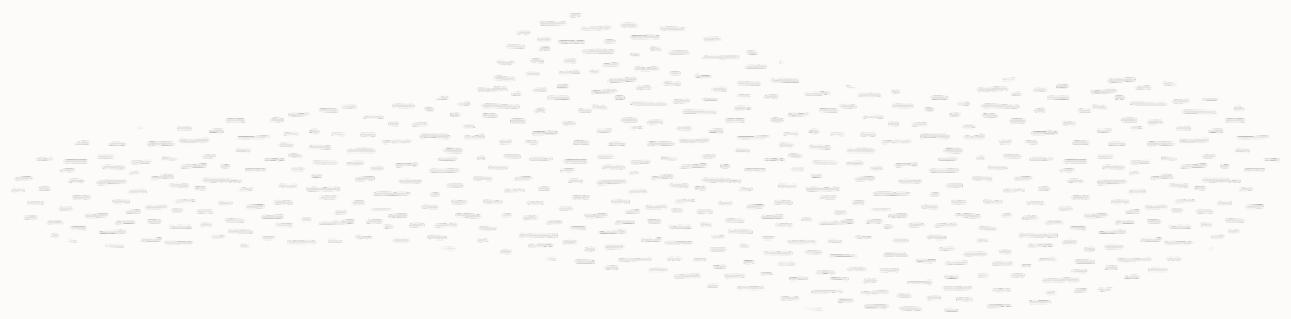
# 50+

---

Clusters

# OKE all the way

1. Started using OKE ~ late 2018
2. We use OKE on OCI for all dev, test, staging, and production workloads
3. The OKE team has been a very good partner in allowing us to run at scale
4. OKE makes Kubernetes a **lot** easier than running it ourselves



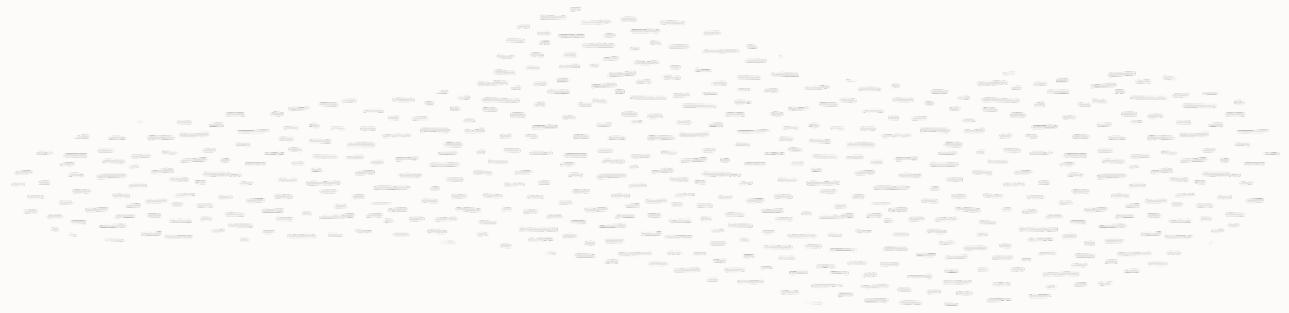
# Lesson #1

---

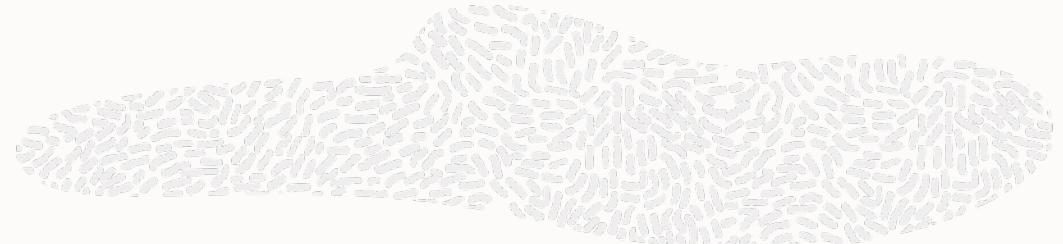
Things die. Deal with it.

# Kubernetes provides all the tools for high availability

1. Deployments and Stateful Sets try to maintain a specified number of *replicas*
2. Kubernetes monitors each pod's *readiness* and *health*. Unhealthy pods are automatically replaced
3. *Anti-affinity* rules can be defined to cause Kubernetes to place pods into different availability / fault domains
4. *Pod disruption budgets* will prevent Kubernetes from bringing down more pods than your application can handle



# Things die. Deal with it.



A pod may die if:

- The pod itself becomes unhealthy (stops responding, or runs out of memory, etc)
- Somebody kills the pod manually (*kubectl delete pod*)
- The cluster was scaled out and you want to rebalance the workload
- **The node that the pod is on dies**

Why would a node die?

- Upgrading the VM image requires recreating the node
- Upgrading the kubelet requires recreating the node
- It just stops working.\*

\* This happens to us a lot with X6 shapes.

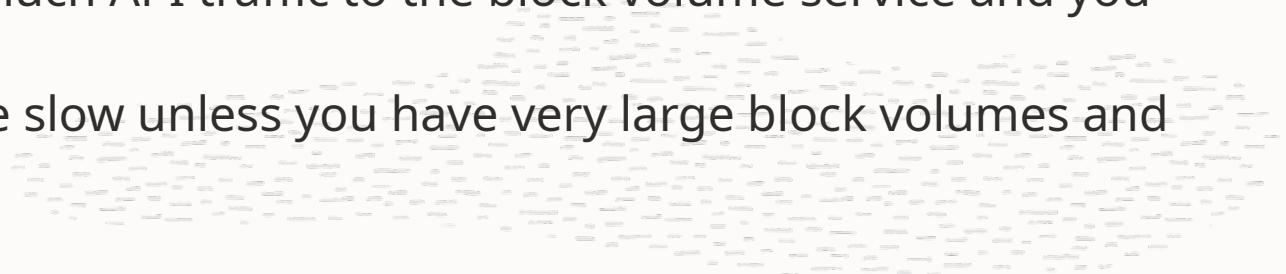
# Lesson #2

---

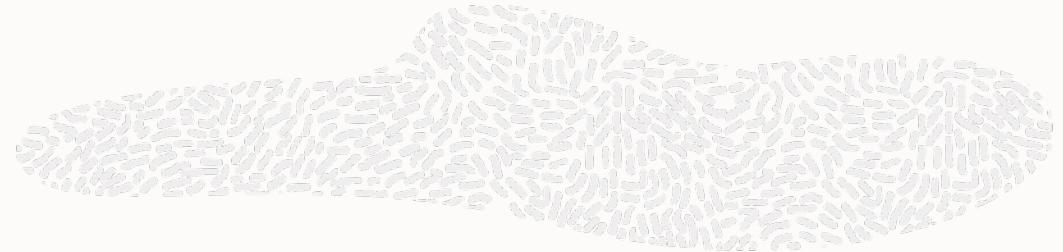
Persistent pains aren't just in old age.

# Persistent volumes caused us persistent problems

1. Kubernetes pods run on Virtual Machine nodes. Each OCI VM has a limit on the number of attached volumes they support. The more persistent volumes, the fewer pods per node.
2. Sometimes we found persistent volumes would get mounted read-only
3. A persistent volume is per-AD. If a pod moves from one AD to another, then it will fail to start because it was unable to attach to its persistent volume.
  - CSI (replaces flex-volume) is designed to help resolve this
4. Persistent volumes have a minimum size of 50GB in OCI
  - This is far too expensive for development, so we setup our own Ceph cluster on some Dense IO bare metal machines in dev so that we could hand out much smaller volumes
5. Too many persistent volumes cause too much API traffic to the block volume service and you will have attach errors
6. I/O intensive tasks with block volumes are slow unless you have very large block volumes and VMs.



# Try using FSS and ephemeral storage



File System Service is NFS for OCI

- If you need a distributed filesystem, you can use an FSS-backed PV and attach it to many pods
- Store heap dumps on an FSS-backed PV

Ephemeral storage for logs and tmp

- Logs are always shipped to ElasticSearch / Kibana ASAP
  - We only need enough ephemeral storage as a buffer
- Every pod in Kubernetes can have ephemeral storage
  - Very useful for very small disk usage requirements
  - Not to be used for storing customer data
  - Ephemeral storage **dies with the pod** and is stored locally on the OKE Node

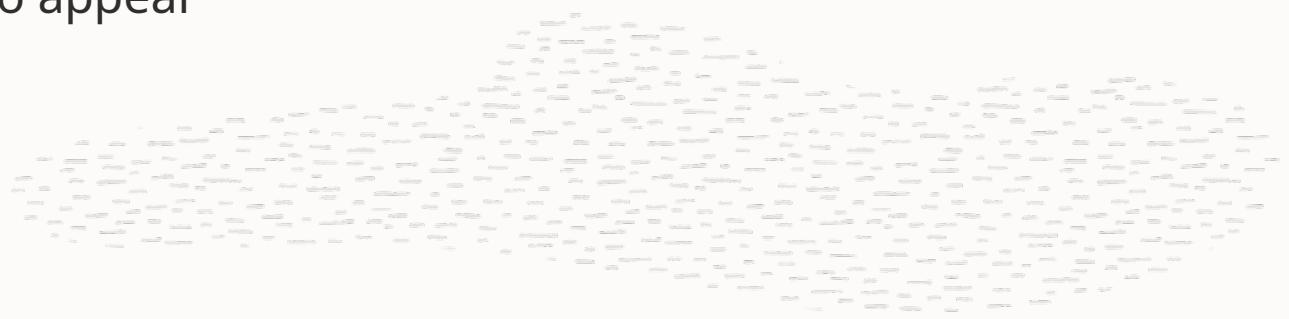
# Lesson #3

---

Know Thy Limits

# Every pod should have requests and limits specified

1. A *request* tells Kubernetes what a pod's *minimum* resource needs are. When Kubernetes schedules your pod, it will make sure the node it schedules the pod on has sufficient resources.
  - Requests ensure a pod has the resources it needs
2. A *limit* tells Kubernetes what a pod's *maximum* resource needs are. The pod's CPU will be limited at the Linux Kernel level. If the pod exceeds its memory *limit*, then the pod will die.
  - Limits prevent a pod from stealing resources from other pods
3. If the node begins to encounter pressure, *request* and *limit* are taken into account to decide what to evict first.
4. Problems due to missing limits and requests may not be visible when cluster load is low, but when the load increases, problems start to appear



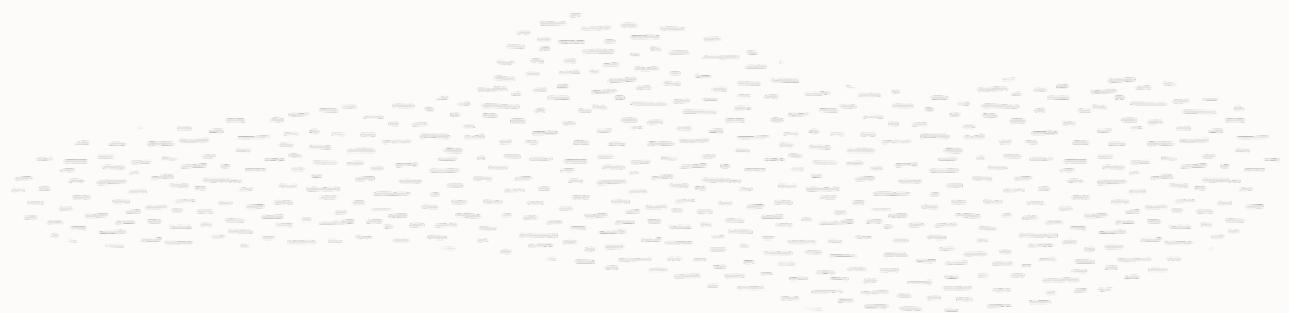
# Lesson #4

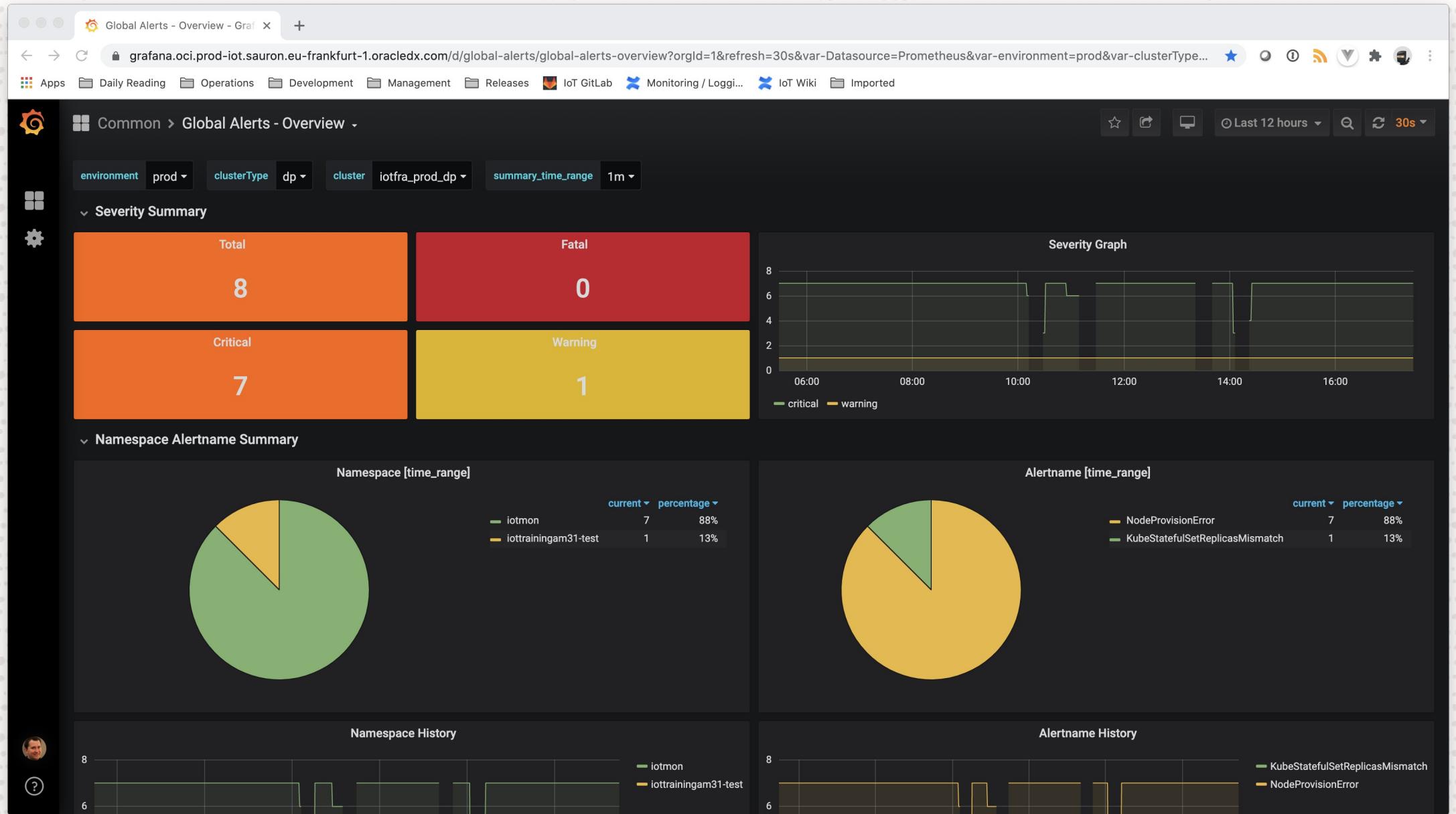
---

Monitoring is essential

# Without good monitoring tools we'd have no chance

1. Every OKE cluster we create runs Prometheus to gather metrics
2. Thanos provides distributed Prometheus queries giving us a single pane of glass for monitoring everything!
3. Single Kibana / ES cluster for collecting all our logs
4. Alert Manager raises PagerDuty alerts based on Kibana / Prometheus data
5. Considering providing some “health check” REST endpoint to give infrastructure **and** application level health data and make it available to monitoring.
  - A pod may be alive, but the underlying application still unhealthy. You need to know this
    - We add health check status to our Prometheus metrics to expose them to Graphana





Kubernetes / Compute Resource

grafana.oci.prod-iot.sauron.eu-frankfurt-1.oracledx.com/d/k8s-compute-resources-cluster-iot-dp/kubernetes-compute-resources-cluster-iot-dp?orgId=1&refresh=5m&var-dataSource=Pro...

Apps Daily Reading Operations Development Management Releases IoT GitLab Monitoring / Loggi... IoT Wiki Imported

Kubernetes > Kubernetes / Compute Resources / Cluster / IOT DP

Last 1 hour 5m

environment prod cluster iotiad\_prod\_dp availabilityDomain All vmType All threshold 0

Summary

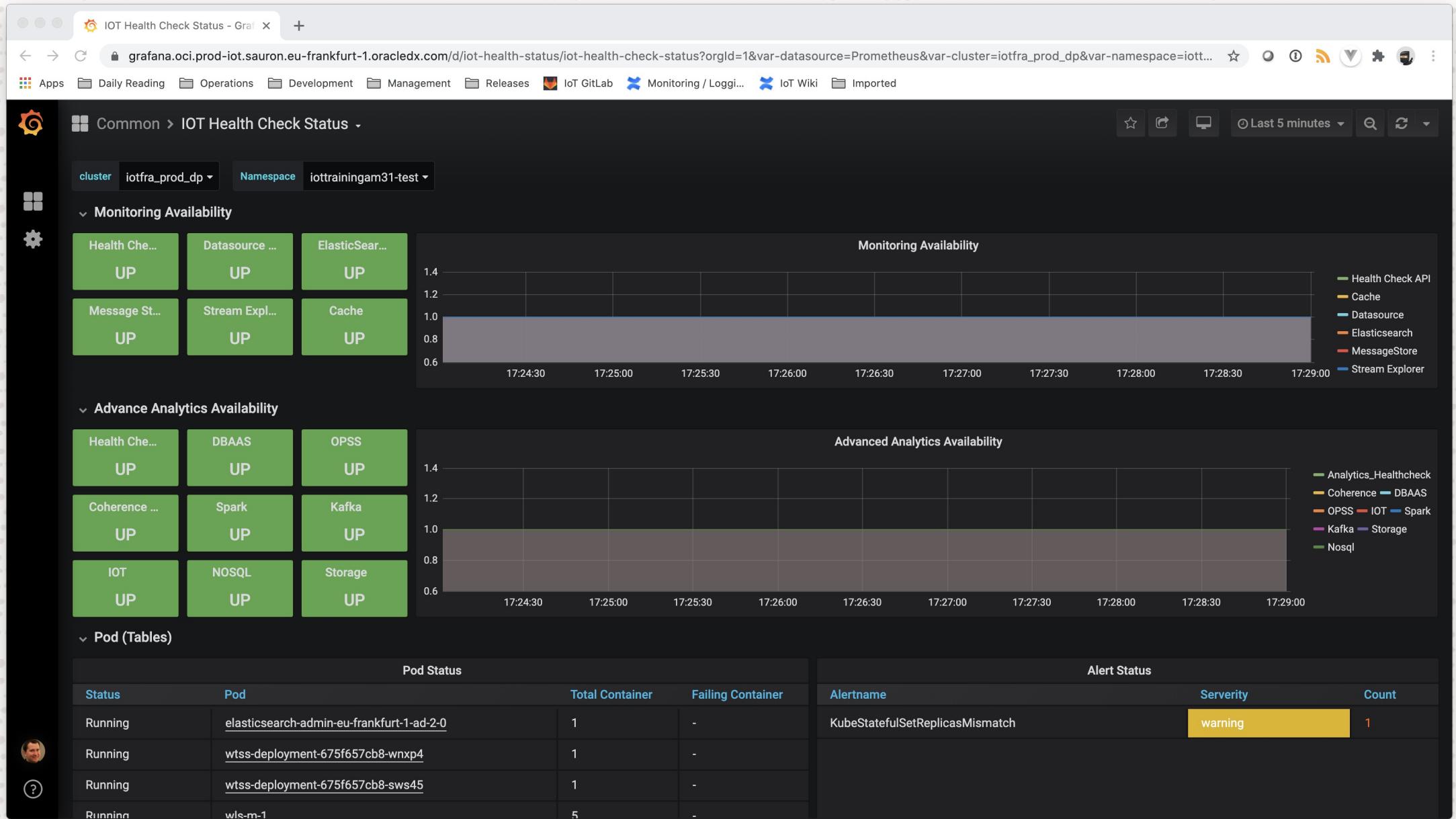
Nodes	Failed Nodes	CPU Nodes	SPARK Nodes	DATA Nodes
282	N/A	Availability Domain Count US-ASHBURN-AD-1 38 US-ASHBURN-AD-3 38 US-ASHBURN-AD-2 38	Availability Domain Count US-ASHBURN-AD-1 24 US-ASHBURN-AD-3 24 US-ASHBURN-AD-2 24	Availability Domain Count US-ASHBURN-AD-1 32 US-ASHBURN-AD-3 32 US-ASHBURN-AD-2 32

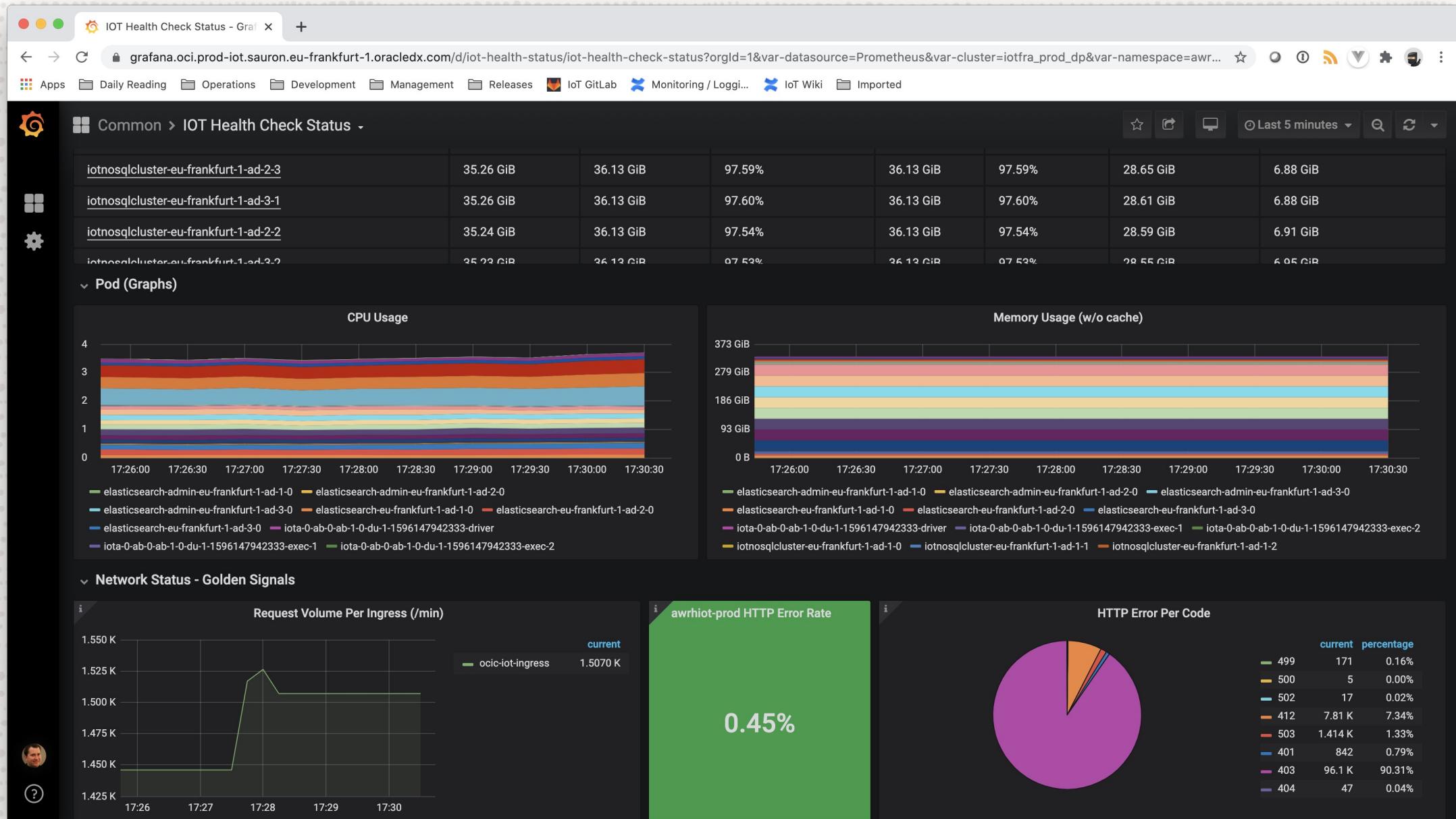
CPU Requests Total(CPU/DATA)		CPU Requests Total(SPARK)		Memory Requests Total(CPU/DATA)		Memory Requests Total(SPARK)	
Availability Domain	Usage	Availability Domain	Usage	Availability Domain	Usage	Availability Domain	Usage
US-ASHBURN-AD-3	38.15%	US-ASHBURN-AD-3	23.06%	US-ASHBURN-AD-3	28.42%	US-ASHBURN-AD-3	12.48%
US-ASHBURN-AD-2	42.74%	US-ASHBURN-AD-2	22.12%	US-ASHBURN-AD-2	32.26%	US-ASHBURN-AD-2	12.07%
US-ASHBURN-AD-1	45.69%	US-ASHBURN-AD-1	23.20%	US-ASHBURN-AD-1	33.89%	US-ASHBURN-AD-1	12.70%

Failed Nodes

No data to show

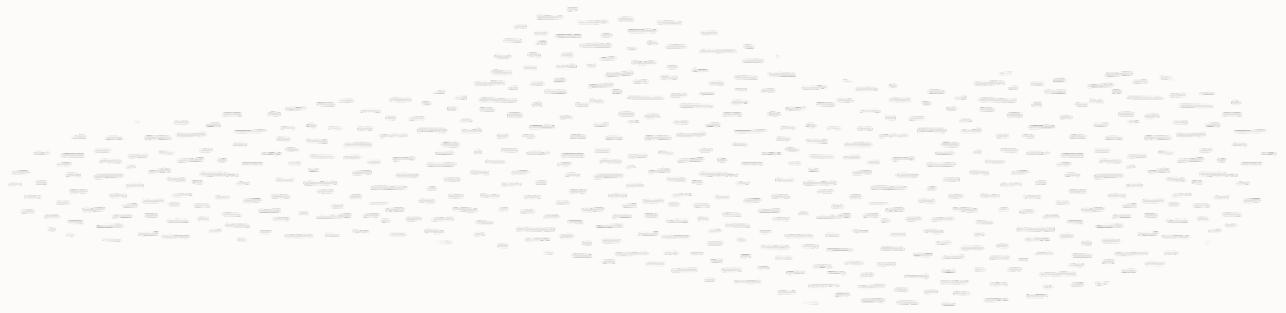
?



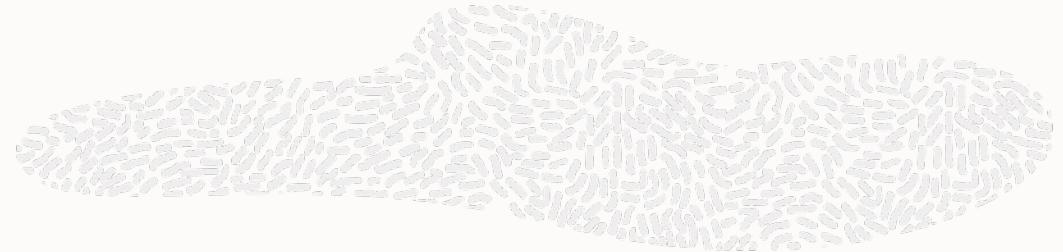


# Lessons

1. Things die. Design for HA from the beginning
2. Use persistent volumes wisely. They can help, but they come at a cost.
3. Limits and Requests help make things stable.
4. Monitoring is essential



# Do I recommend OKE? Yes!



What does OKE add over Kubernetes?

- Fully managed Kubernetes Control Plane
- Future feature to also have a fully managed Kubernetes Data Plane

What does Kubernetes add over Virtual Machines?

- Managed deployment of Docker containers
- All tools for liveness, readiness, capacity limits, scaling, HA, etc.

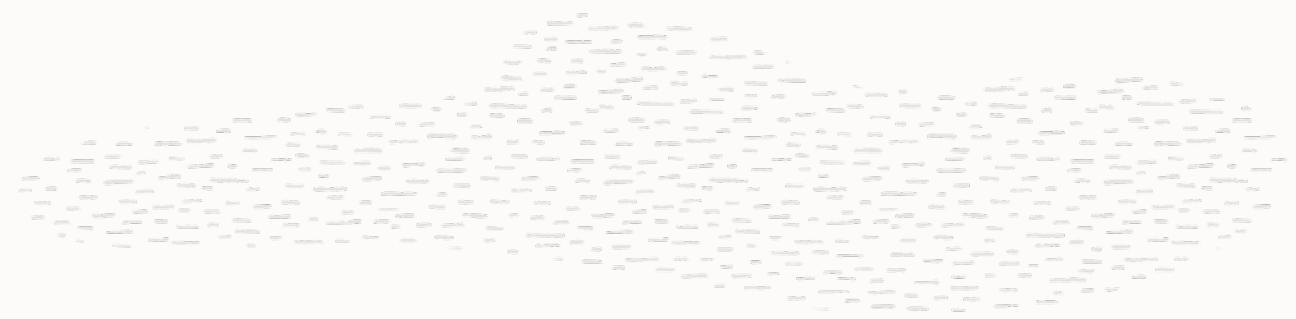
What do Virtual Machines add over Bare Metal?

- More variety in VM shapes and sizes
- Better support for HA / DR (e.g. network attached block volumes)

**The higher up the stack, the better!**

# Resources

1. <https://developer.oracle.com>
2. <https://www.oracle.com/cloud/free/>
3. <https://docs.cloud.oracle.com/en-us/iaas/Content/ContEng/Concepts/contengoverview.htm>
4. <https://www.oracle.com/cloud-native/>
5. <https://docs.cloud.oracle.com/en-us/iaas/Content/Resources/Assets/whitepapers/oci-security.pdf>
6. <https://www.oracle.com/a/ocom/docs/cloud/oracle-cloud-infrastructure-platform-overview-wp.pdf>



Thank you

---

