SCALITY METALK8S

AN OPINIONATED KUBERNETES DISTRIBUTION WITH A FOCUS ON LONG-TERM ON-PREM DEPLOYMENTS

Nicolas Trangez - Technical Architect nicolas.trangez@scality.com @eikke | @Scality | @Zenko



ABOUT SCALITY



ONE PURPOSE GIVING FREEDOM & CONTROL PEOPLEWHO



GLOBAL PRESENCE



GLOBAL CLIENT BASE









EndemolShine

MediaHub

Group

foxtel

OUR JOURNEY TO KUBERNETES

Scality RING, S3 Connector & Zenko



Scality RING

On-premise
Distributed Object & File Storage

- Physical servers, some VMs
- Only the OS available (incl. 'Legacy' like CentOS 6)
- Static resource pools
- Static server roles / configurations
- Solution distributed as RPM packages, deployed using SaltStack
- De-facto taking ownership of host, difficult to run multiple instances
- Fairly static post-install



Scality S3 Connector

On-premise S3-compatible Object Storage

- Physical servers, sometimes VMs
- Static resource pools
- "Microservices" architecture
- Solution distributed as Docker container images, deployed using Ansible playbooks
- No runtime orchestration
- Log management, monitoring,...
 comes with solution



Scality Zenko

Multi-Cloud Data Controller

- Deployed on-prem or 'in the Cloud': major paradigm shift
- New challenges, new opportunities
- Multi-Cloud Data Controller, must run on multiple Cloud platforms



Scality Zenko

Deployment Model

- Embraced Docker as distribution mechanism
 - Some shared with Scality S3 Connector
- For Cloud deployments, started with Docker Swarm
 - Ran into scaling, reliability and other technical issues
- Decided to move to Kubernetes
 - Managed platforms for Cloud deployments, where available (GKE, AKS, EKS one day)
 - On-prem clusters



Scality Zenko

Kubernetes Benefits

- Homogenous deployment between in-cloud and on-prem
- Various services provided by cluster:
 - Networking & policies
 - Service restart, rolling upgrades
 - Service log capturing & storage
 - Service monitoring & metering
 - Load-balancing
 - TLS termination
- Flexible resource management
 - If needed, easily add resources to cluster by adding some (VM) nodes
 - HorizontalPodAutoscaler



OUR JOURNEY TO KUBERNETES

MetalK8s



On-prem Kubernetes

- Can't expect a Kubernetes cluster to be available, provided by Scality customer
- Looked into various existing offerings, but in the ends needs to be supported by/through Scality (single offering)
 - Most existing solutions don't cover enterprise datacenter requirements
- Decided to roll our own



SCALITY METALK8S

AN OPINIONATED KUBERNETES DISTRIBUTION WITH A FOCUS ON LONG-TERM ON-PREM DEPLOYMENTS



OPINIONATED

We offer an out-of-the-box experience, no non-trivial choices to be made by users



LONG-TERM

Zenko solution is mission-critical, can't spawn a new cluster to upgrade and use ELB (or similar) in front



ON-PREM

Can't expect anything to be available but (physical) servers with a base OS



Scality MetalK8s

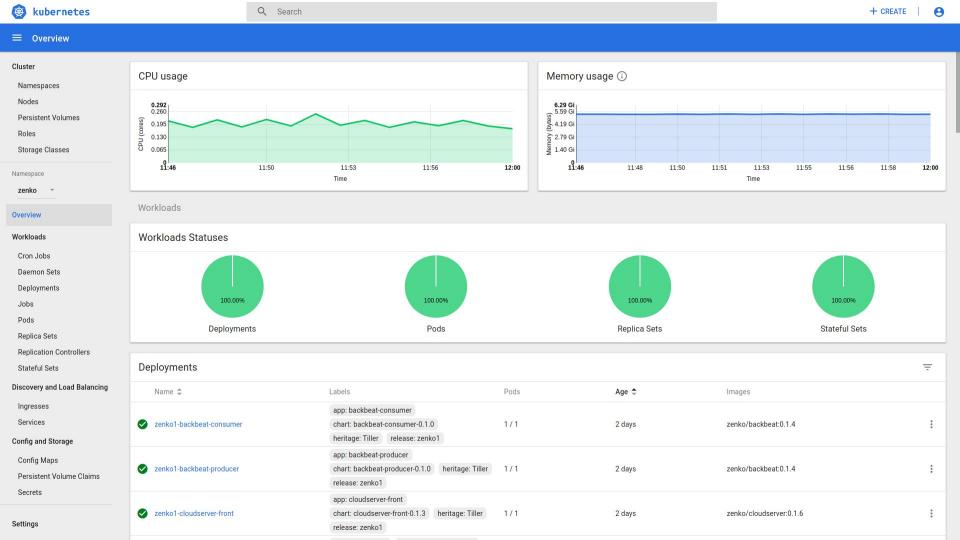
- Scope: 3-20 physical machine, pre-provisioned by customer or partner
- Built on top of the Kubespray Ansible playbook
- Use Kubespray to lay out a base Kubernetes cluster
 - Also: etcd, CNI
- Add static & dynamic inventory validation pre-checks, OS tuning, OS security
 - Based on experience from large-scale Scality RING deployments
- Augment with various services, deployed using Helm
 - Operations
 - Ingress
 - Cluster services
- Take care of on-prem specific storage architecture

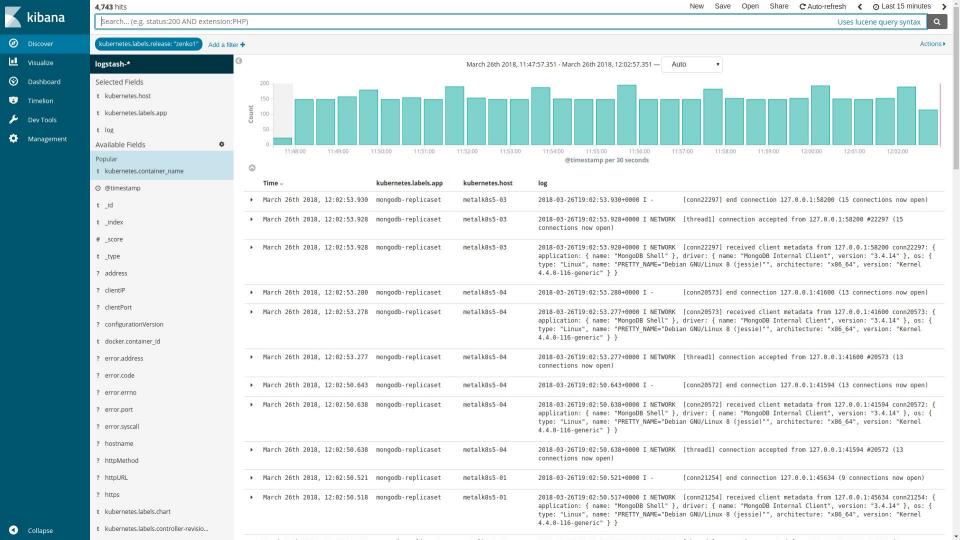


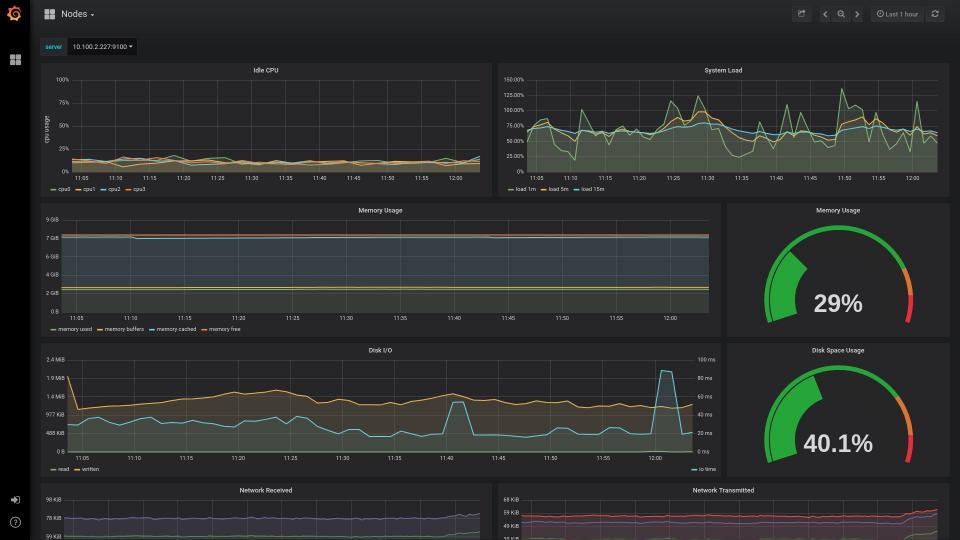
Scality MetalK8s: Cluster Services

- "Stand on the shoulders of giants"
- Heapster for dashboard graphs, `kubectl top`,...
- metrics-server for HorizontalPodAutoscaler
 - Looking into k8s-prometheus-adapter
- Ingress & TLS termination: nginx-ingress-controller
- Cluster monitoring & alerting: Prometheus, prometheus-operator, Alertmanager, kube-prometheus, Grafana
 - Host-based node_exporter on all servers comprising the cluster, including etcd
- Host & container logs: ElasticSearch, Curator, fluentd, fluent-bit, Kibana
- All of the above gives a great out-of-the-box experience for operators









Scality MetalK8s: Storage

- On-prem: no EBS, no GCP Persistent Disks, no Azure Storage Disk,...
- Also: can't rely on NAS (e.g. through OpenStack Cinder) to be available
- Lowest common denominator: local disks in a node
- PVs bound to a node, hence PVCs bound, hence Pods bound
 - Thanks PersistentLocalVolumes & VolumeScheduling!
- Decided not to use LocalVolumeProvisioner, but static approach (for now)
 - Based on LVM2 Logical Volumes for flexibility
 - PV, VG, LVs defined in inventory, created/formatted/mounted by playbook
 - K8s PV objects created by playbook
 - May support whole partitions/drives depending on application need
- Dynamic local volume provisioning (using LVM) is getting there...
 - Future: volume encryption?



Scality MetalK8s: Deployment

- Based on years of years of experience deploying Scality RING at enterprise customers, service providers,...
- Constraints in datacentra often very different from 'VMs on EC2'
 - No direct internet access: everything through HTTP(S) proxy, no non-HTTP traffic
 - Dynamic server IP assignment
 - Security rules requiring services to bind to specific IPs only, different subnets for control & workload,...
 - Fully air gapped systems: requires 100% offline installation
 - Non-standard OS/kernel
 - Integration with corporate authn/authz systems
- Not all of the above supported yet, tackling one by one
 - Relevant patches to be upstreamed to Kubespray
- Only support RHEL/CentOS family of Linux distributions
 - Support for Ubuntu and others can be community-driven, Kubespray supports them
 - RHEL/CentOS sometimes difficult targets for containers/Docker/Kubernetes



Scality MetalK8s: Ease of Deployment

```
$ # Requirements: a Linux or OSX machine with Python and coreutils-like
$ # Create inventory
$ vim inventory/...
$ make shell # Launches a 'virtualenv' with Ansible & deps, 'kubectl',
'helm'
 # Demo @ https://asciinema.org/a/9kNIpBWg4KiwjT5mNSrH0tmj9
$ ansible-playbook -i inventory -b playbooks/deploy.yml
 # Grab a coffee, and done
```



Scality MetalK8s: Non-technical goodies

- Documentation
 - Various guides: Installation, Operations, Reference
 - https://metal-k8s.readthedocs.io
- Extensive testing
 - Installation
 - Upgrade
 - Services
 - Failure testing



Scality MetalK8s: Shifting focus

- Today: general-purpose deployment tool, fulfil K8s cluster pre-req of \$product

- Future: use-case specific component a vendor (you!) can embed in on-prem solution/product running on K8s without being a K8s product
 - More configurable to match solution requirements
 - Tighten out-of-the-box security



Scality MetalK8s: The road forward

- Increase documentation coverage
- Considering removing Kubespray
 - Too 'big' for our purposes
 - kubeadm brought kubelet bootstrapping and other goodies
 - Non-trivial to implement certain requirements/features
- Considering removing Docker
 - Plain CRI with containerd or cri-o
- Looking into cluster federation (multi-site solutions), built-in over-the-wire encryption (Wireguard?), 'active' cluster controller (refresh short-TTL TLS certs, provision new nodes,...), netboot (like CoreOS/Matchbox/Tectonic, but plain CentOS/RHEL), other CNIs, integration of failover (VIP) and load-balancing service,...



SCALITY A METALK8S

AN OPINIONATED KUBERNETES DISTRIBUTION WITH A FOCUS ON LONG-TERM ON-PREM DEPLOYMENTS

https://zenko.io

https://github.com/scality/metalk8s

@Scality | @Zenko

