



FORTH

INSTITUTE OF COMPUTER SCIENCE

**Cloud automation and
interfacing to HPC – What
happens with Data ?**

Fotis Nikolaidis
Antony Chazapis
Manolis Marazakis
Angelos Bilas

Per3S
2023

About this presentation

- Objective:
 - Run Kubernetes applications on HPC infrastructure.

- Agenda
 - Background on Cloud/HPC Convergence.
 - Rootless architecture for deploying Kubernetes as Slurm job.
 - Show-case.
 - Key takeaways.

- Slurm dominates HPC.
 - Specialized in launching MPI jobs at scale.
 - **Singularity containers is an option, but requires external management.**
- Kubernetes dominates Cloud.
 - Specialized in **container management** (i.e healing, scaling, replication).
 - Originally designed for DevOps, is now getting traction for **repeatable data science**.

“

Why should I care ?

”

1. **Cloud-User PoV: Scale-out workflows written for Kubernetes.**
 - a. Genomics and Bioinformatics
 - b. ML Training

2. **HPC-user PoV: Exploit Cloud-native Data Science Tools.**
 - a. Combine HPC codes with Cloud-native data analysis → Visualization, querying, ...
 - b. Interactive code execution → Jupyter
 - c. Workflow management → Argo Workflows, Apache Airflow, ...

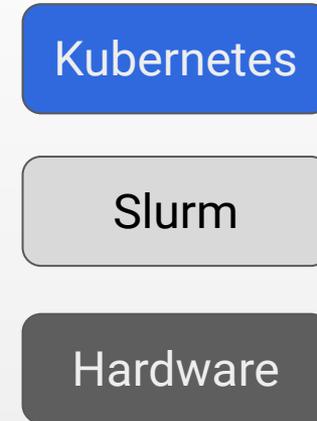
3. **HPC-center PoV: Increase utilization of data center.**
 - a. Attract Cloud users

Architectures for Convergence

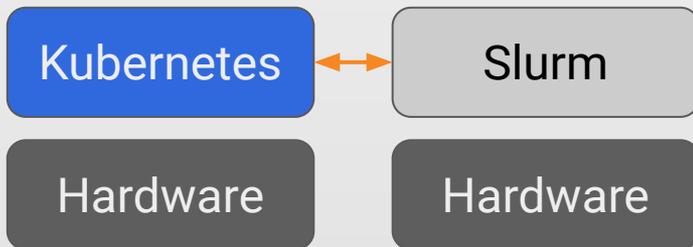


Case: Run MPI jobs on Cloud resources.

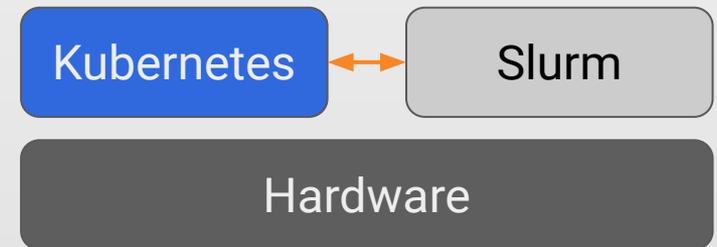
**Today's
Topic!!!**



Case: Run Kubernetes on HPC cluster.



Case: Bridge different clusters.



Case: Support both, avoid data transfers.

Kubernetes

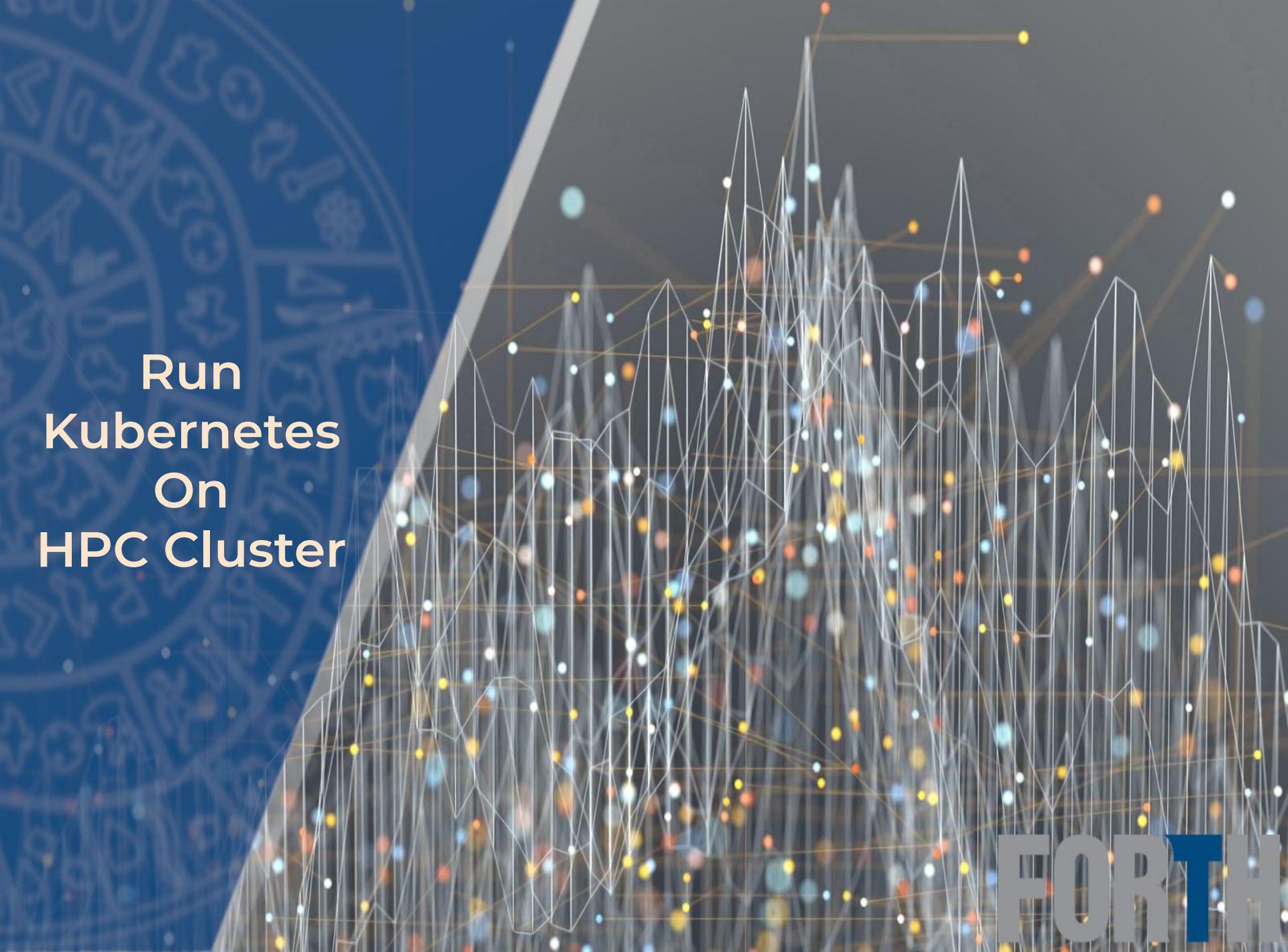
Slurm

Hardware

Objective: Run K8s workloads within **unmodified** Slurm environment.

Requirements:

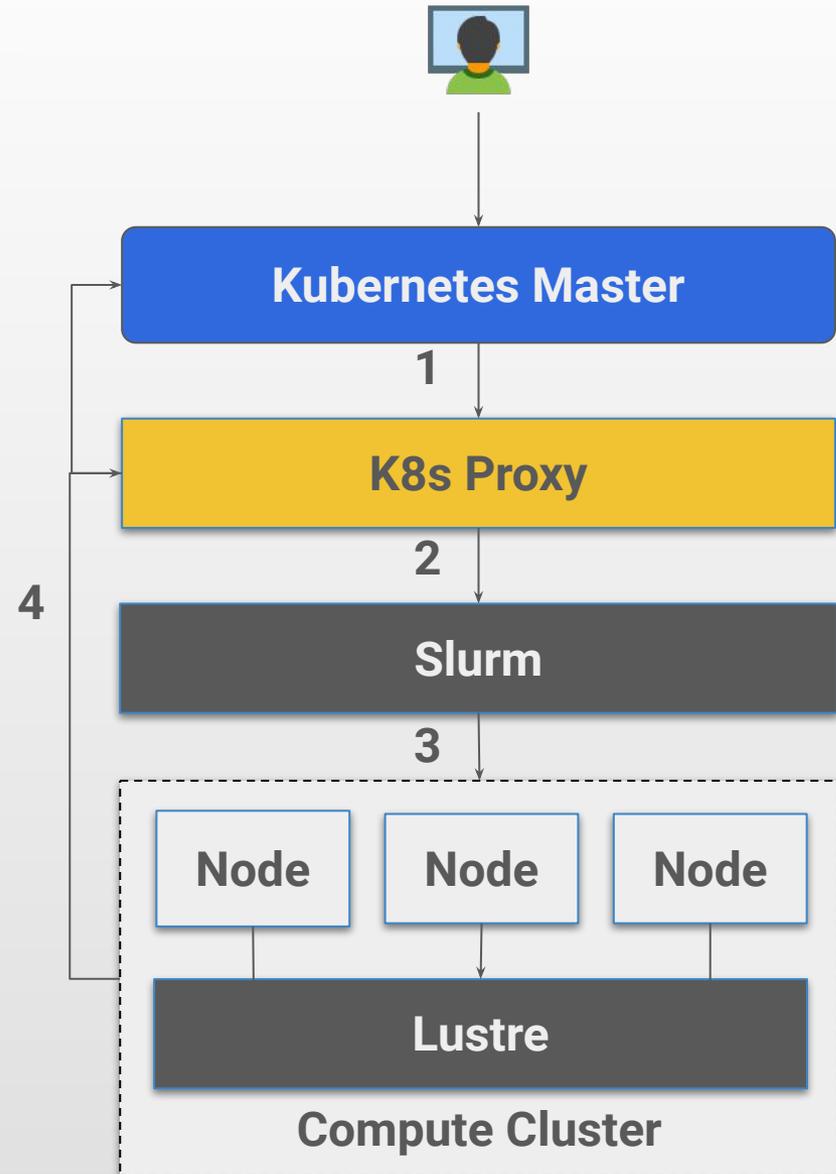
1. Create ephemeral k8s clusters as **user jobs**.
2. Support all Kubernetes abstractions, except privileged.
3. Scale across all nodes of the cluster.
4. Minimal pre-installed software.



Run
Kubernetes
On
HPC Cluster

Four-step process.

1. Proxy receives Job Request from Kubernetes.
2. Proxy translates K8s Job to Slurm Job.
3. Slurms runs the job (**as the host user**).
4. Proxy keeps Slurm and Master in sync.



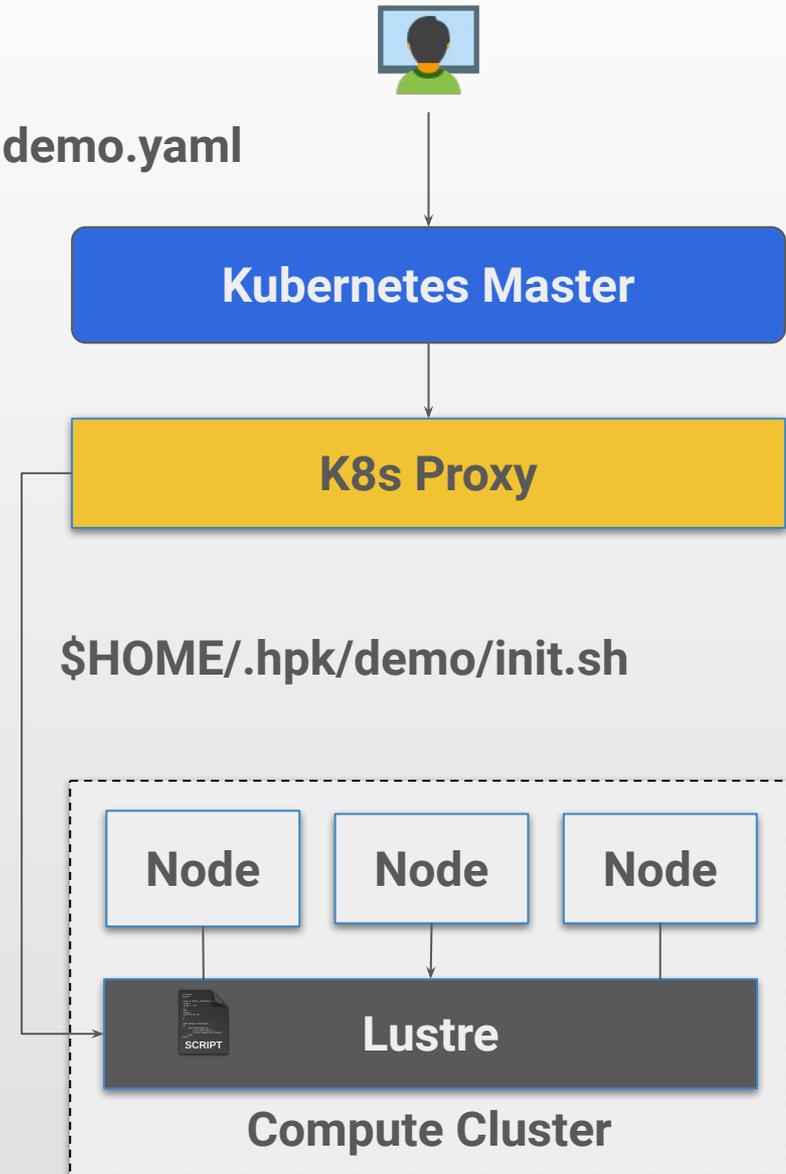
K8s Proxy:

1. Parses YAML fields from k8s requests.
2. Generates equivalent sbatch scripts.
3. Stores the sbatch script (*init.sh*) in Lustre.

Supported Fields:

- Container image to run
- Resource Requirements
- Volumes to be mounted
- ... and more ...

kubectl demo.yaml



Volume Preparation

K8s requires certain data to be present.

- Credentials, Configs, ...

Application needs scratch space for the runtime.

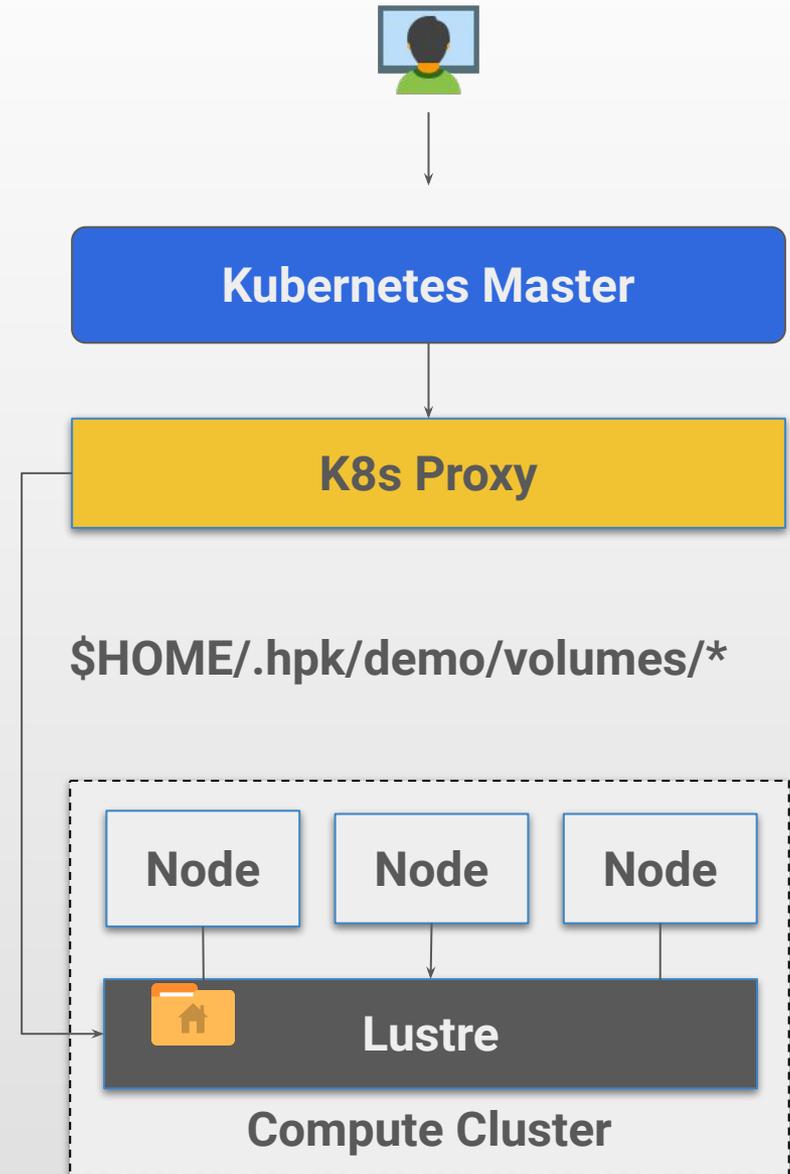
- Logs, temporary files, ...

User needs access to persistent storage.

- Load dataset, write results.

K8s Proxy:

1. Convert volumes into Lustre files/dir.
 - a. Downloads volumes data from master.
 - b. Create temp dir for scratch.
 - c. Create symlinks to other Lustre files.
2. Set the volume paths in **init.sh**.

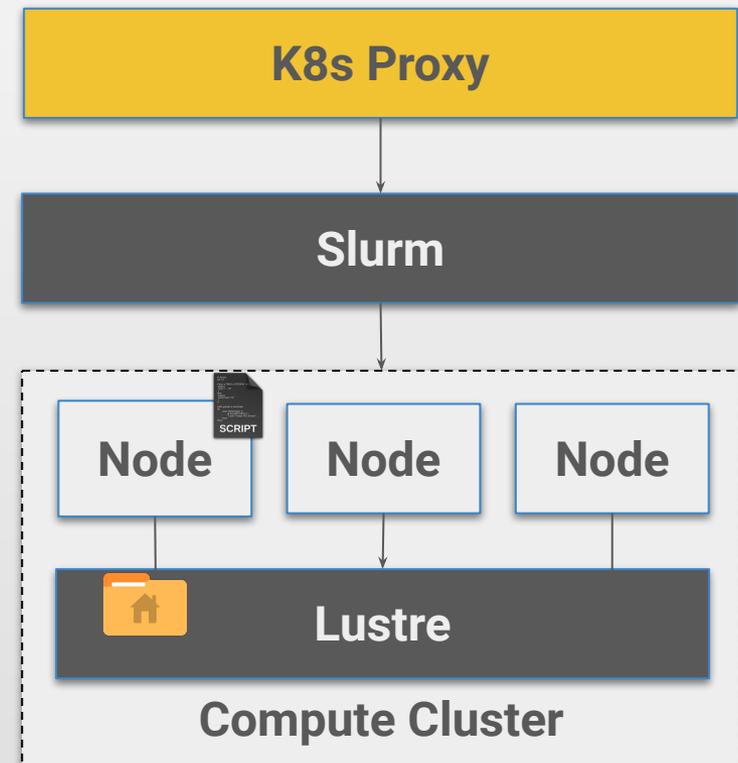


K8s Proxy: Submits */demo/init.sh* to Slurm.

Slurm: Schedules job to the compute nodes.

Compute node: Runs the *init.sh* locally.

Init.sh: Mounts */demo/volumes/** to the container.



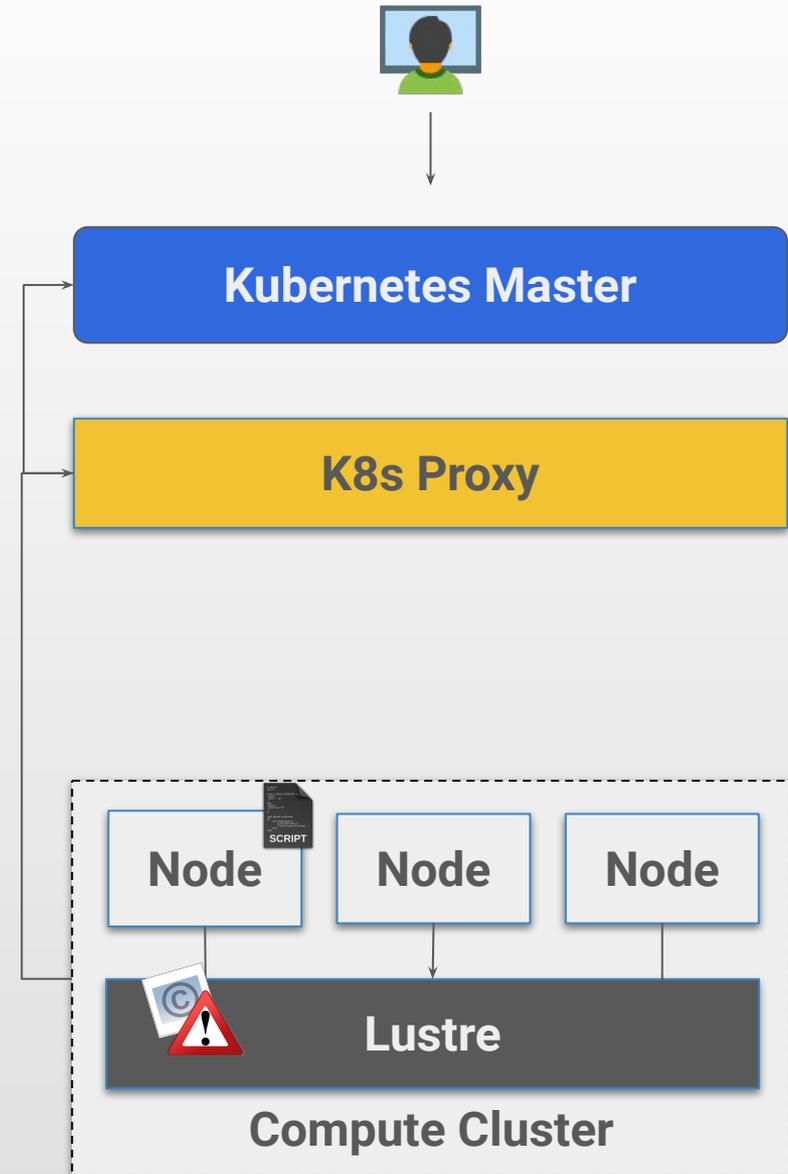
How does Kubernetes know the status of the job ?

Init.sh

Through its execution, it creates *Control Files* on Lustre to indicate the state of execution.

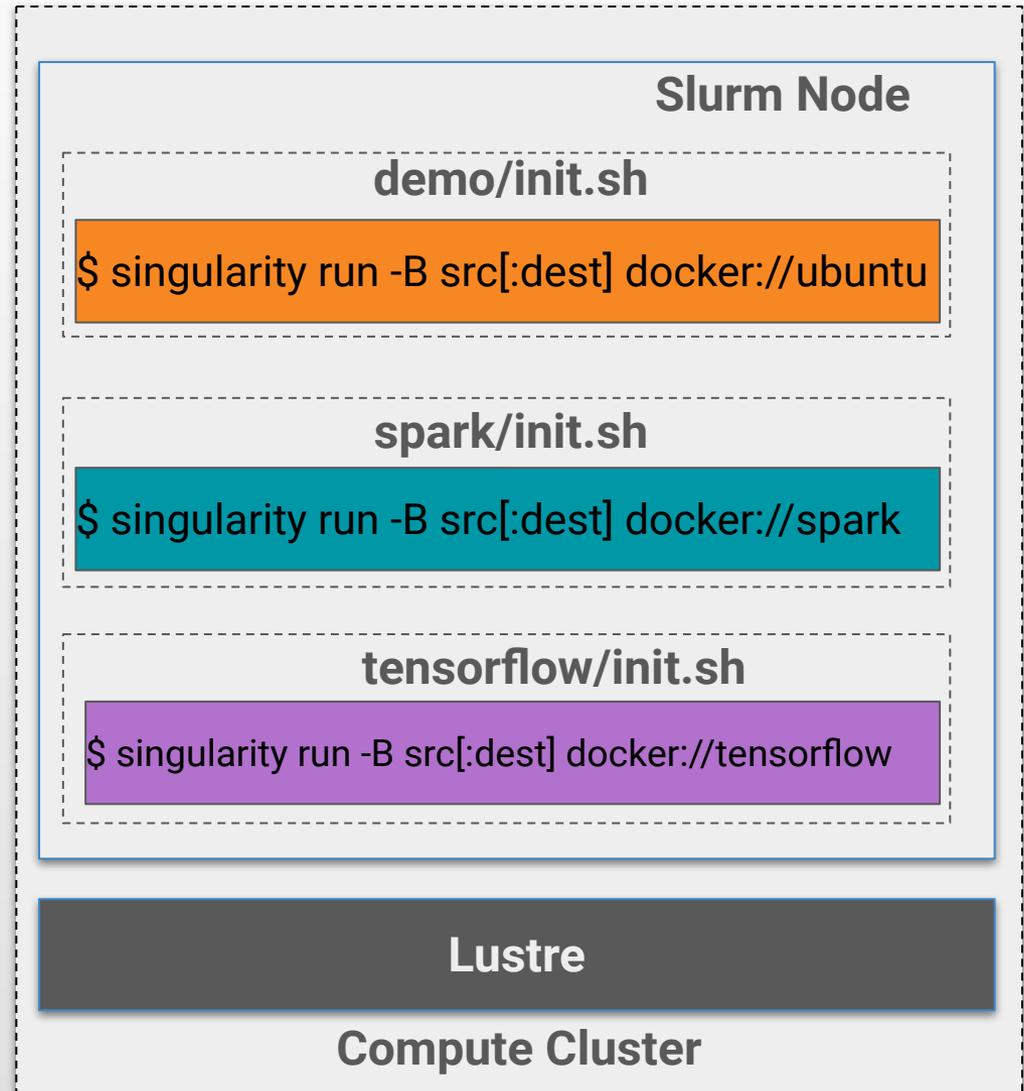
K8s Proxy:

- Tracks changes on *Control Files* by using an Inotify-like mechanism.
- Updates the job's status according to defined semantics.
- Updates Kubernetes master about the new status.



Singularity is the reference container technology for HPC.

- Daemonless: just a binary.
- Rootless: runs as a simple user.
- Integrated: inherits user's env and mountpoints.
- Backward-compatible with Docker images.

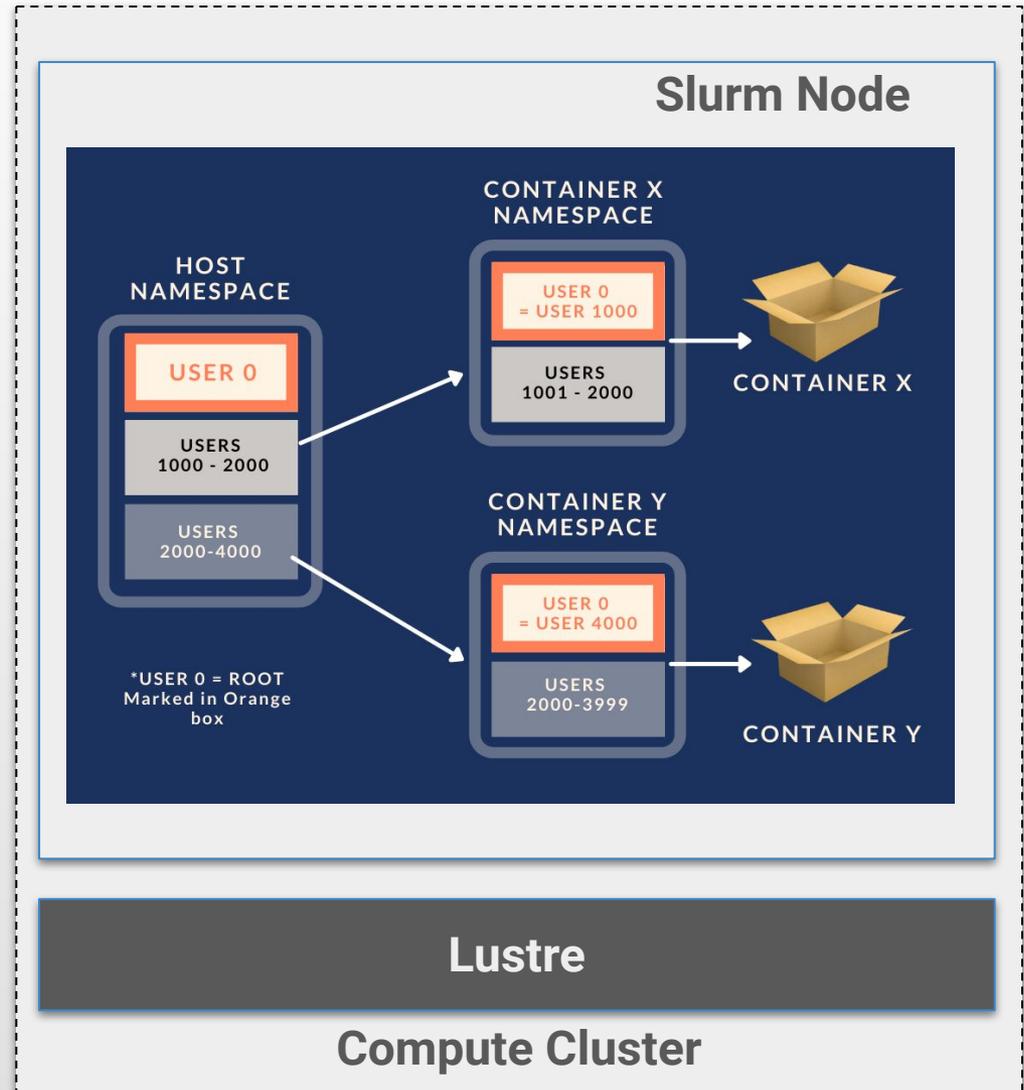


Execute container runtimes as an unprivileged user, by using Linux User Namespaces.

User Mapping: map UIDs/GIDs in the container namespace to unprivileged range in the host namespace.

- Within container: *root:root*
- Outside container: *user:group*

Fakeroot Capabilities: Full capabilities, except for inserting kernel modules, rebooting, ...

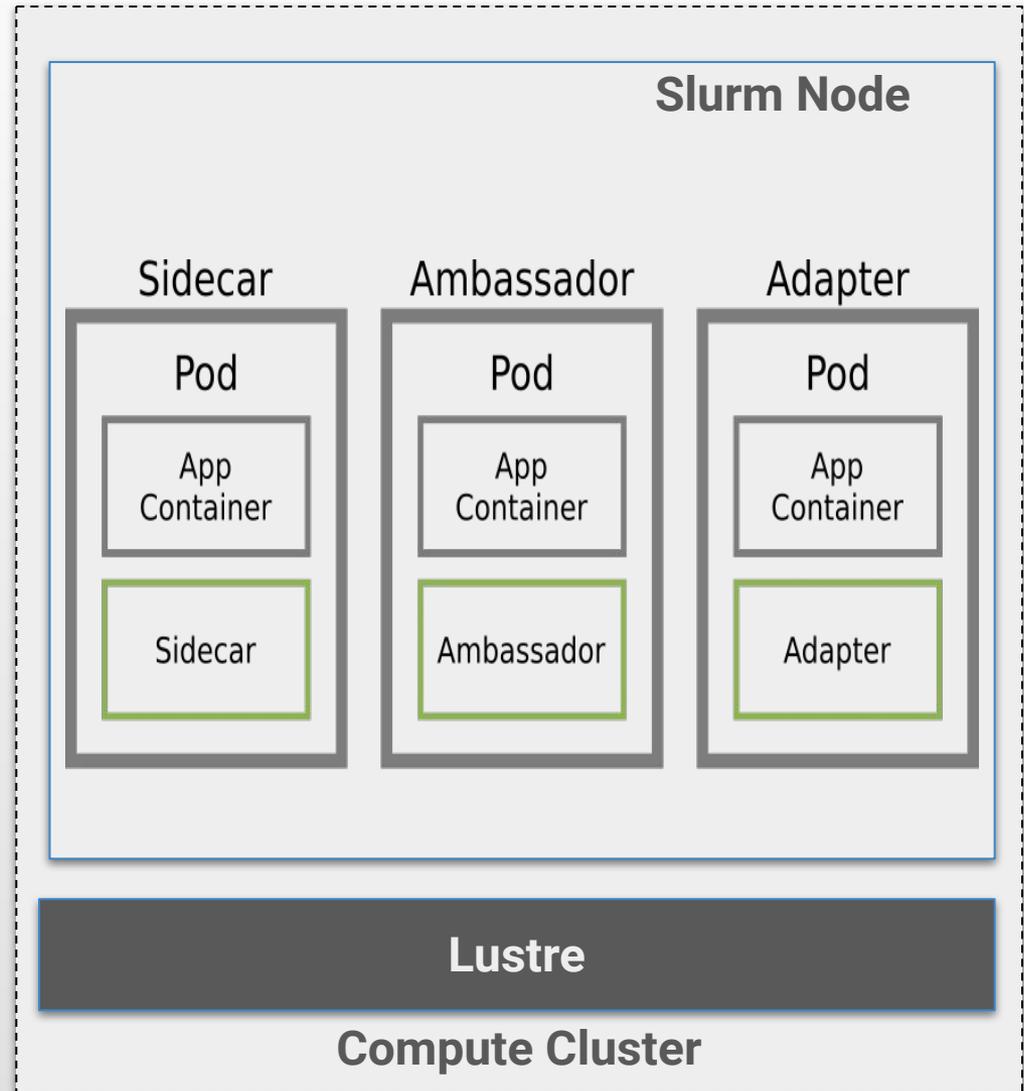


... But not exactly what we want

Kubernetes attributes its success to the concept of Multi-Container Pods.

The Pod is Logical group of containers with shared storage and network resources.

Issue: How can we support Pods with Singularity ?



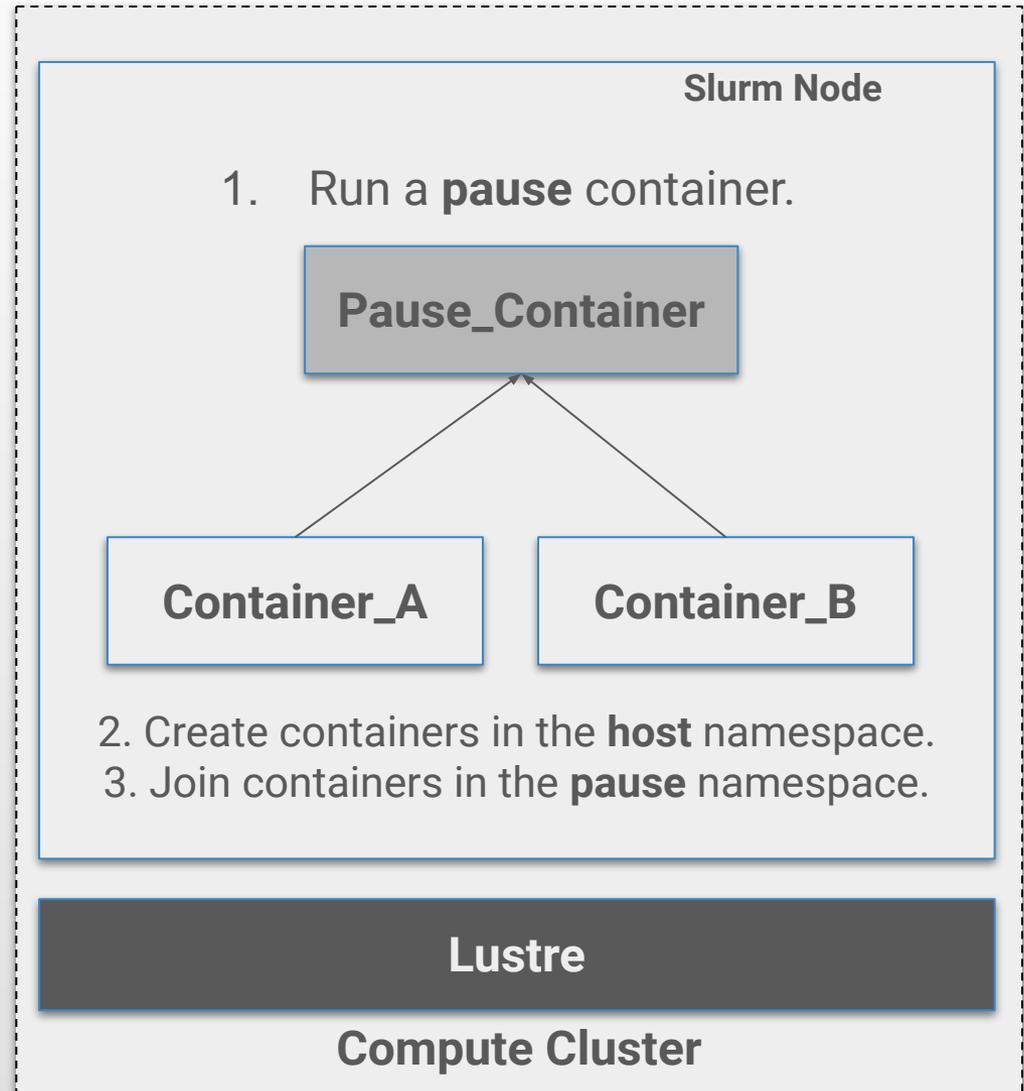
Pause container:

Empty container which establishes namespaces and reservations before individual containers are created.

Main Containers:

Containers are created on the host namespace, and then join the namespace of the pause container using **nsenter** command.

Issue: nsenter requires root.



Rootless Pods (nested containers)

Step 1:

Run a pause container.

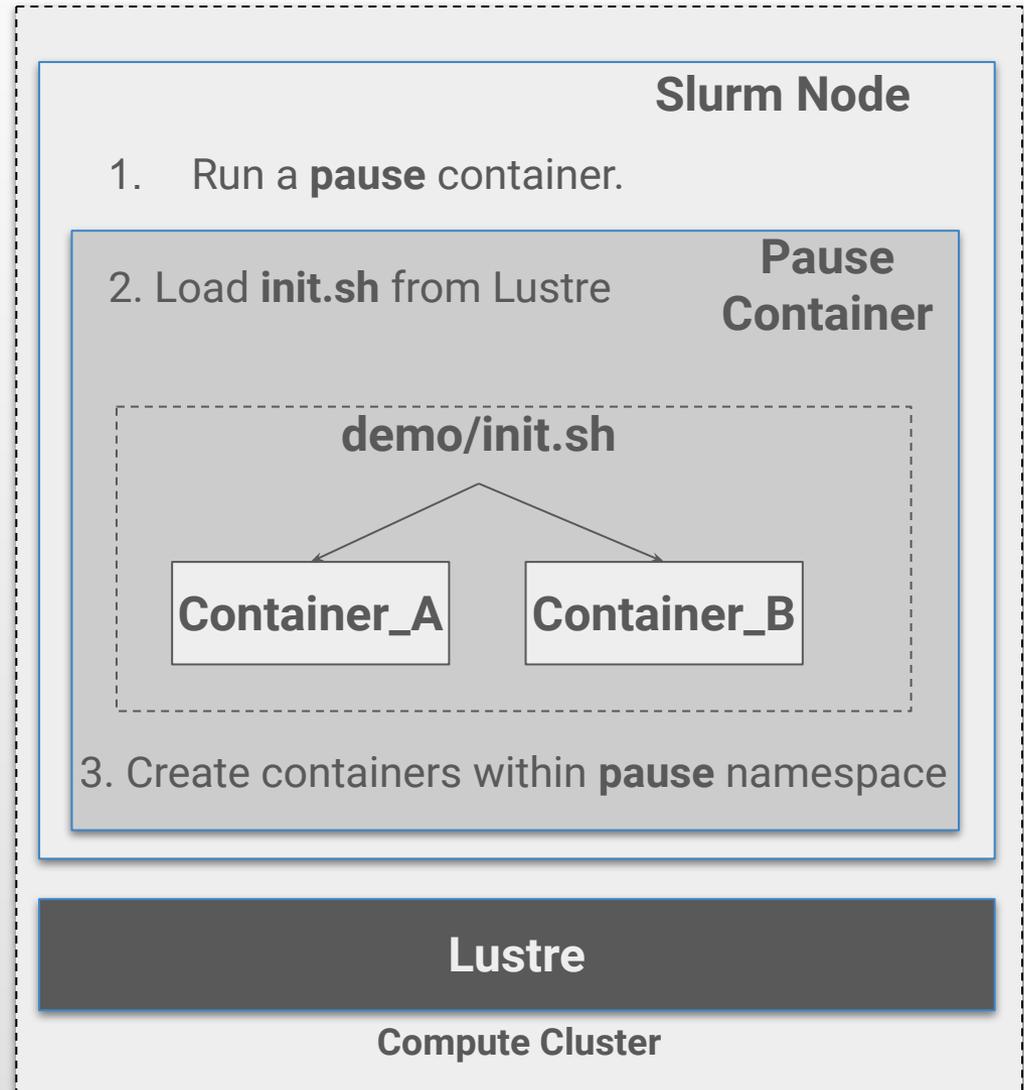
Step 2:

Within the pause container, load and run **init.sh**.

Step 3:

The script creates the containers.

Limitations: all containers must be known in advance.



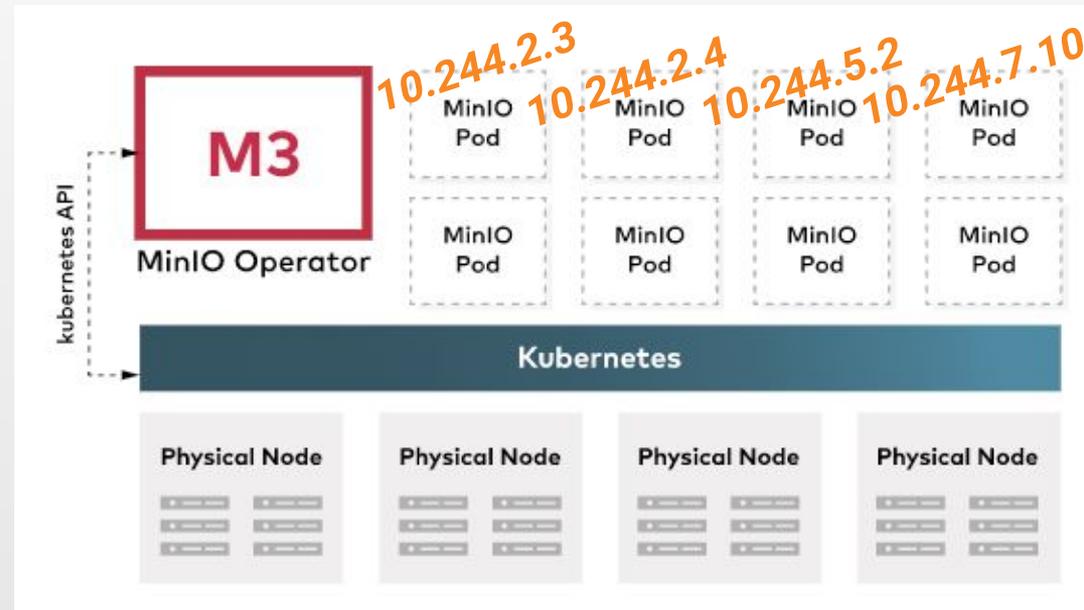
What about object storage ?

Minio is an S3-compatible implementation that is already Kubernetes-native.

Issue #1: Minio, like most other tools, is web-based and requires a routable IP.

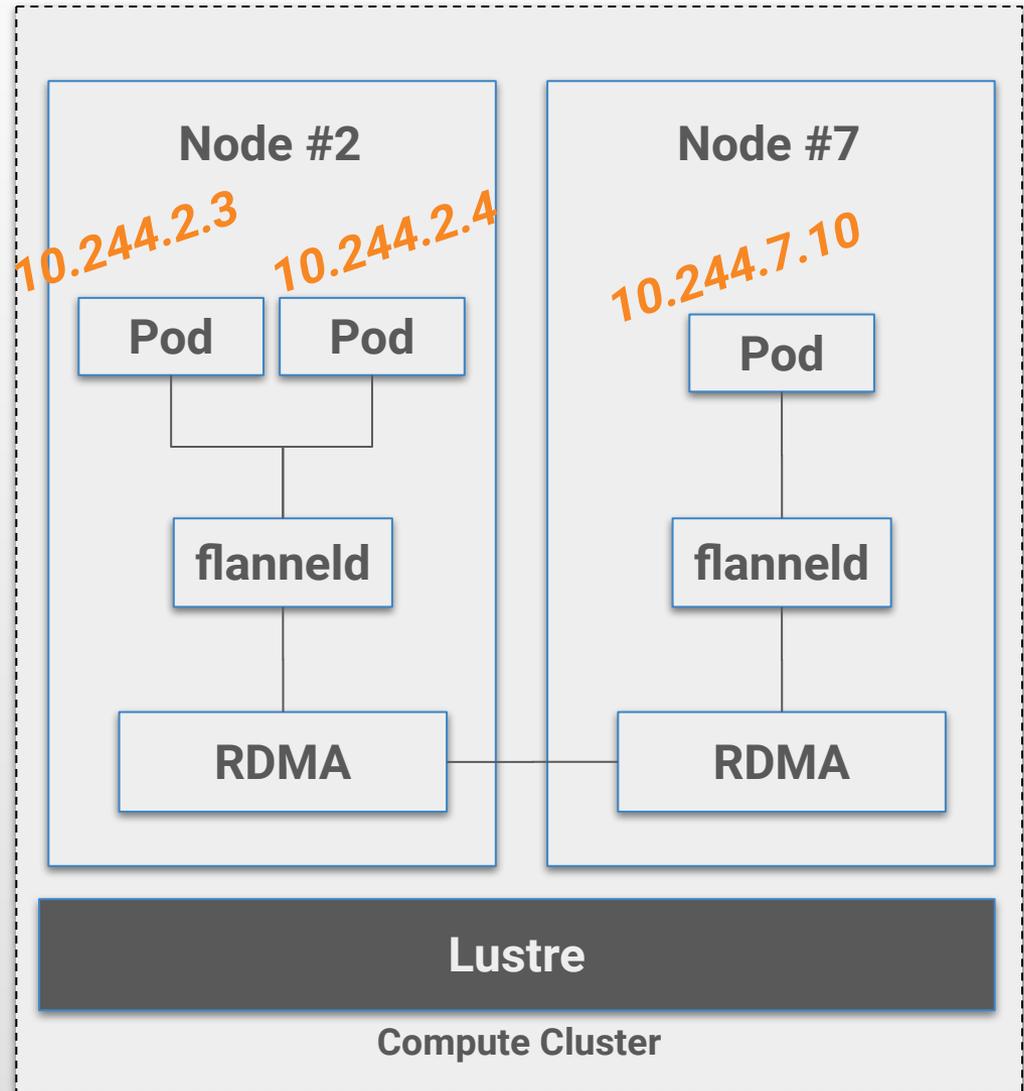
Issue #2: You can't pollute the host IP range.

Issue #3: How can Minio servers become discoverable?



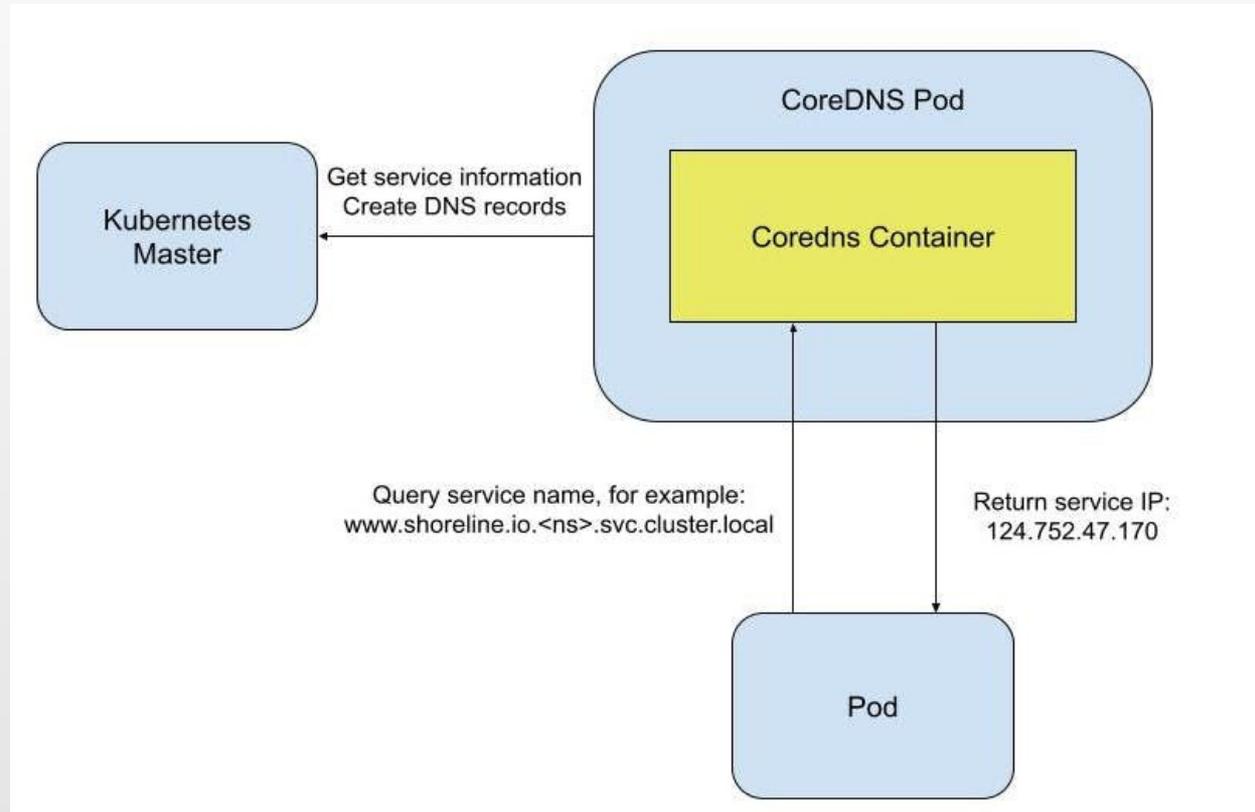
Container Networking

- Within the Pod, we create a network namespace.
- The namespace takes IPs from flanneld that runs on the host.
- Flannel implements container-to-host and host-to-host by modifying the routing tables of the host.

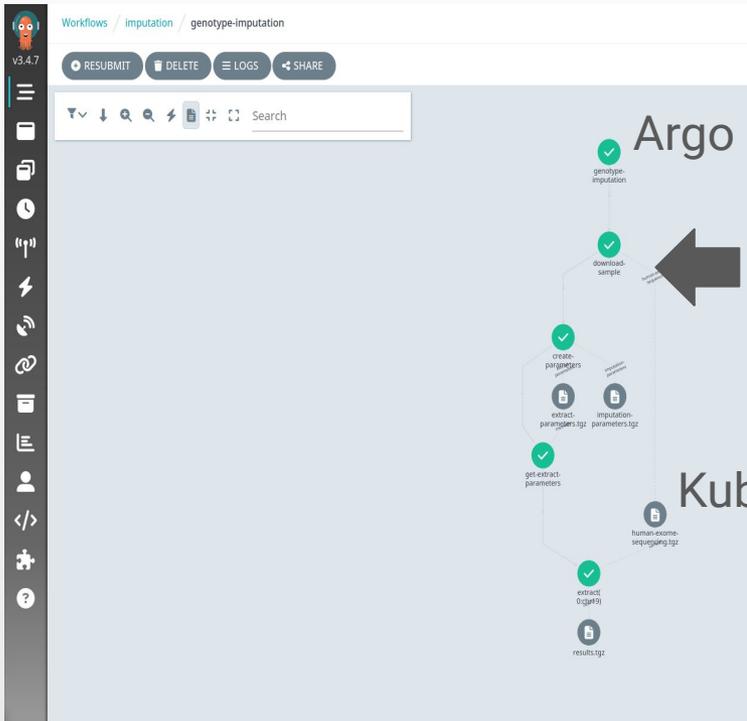


CoreDNS: DNS records for Kubernetes.

- Exclude load balancing as it requires changes on iptables.
- **Support direct Service-To-Pod mappings.**



Demo: Genotype Analysis Workflow



Argo Workflows + Minio

Object Browser

artifacts

Created on: Tue, May 16 2023 19:51:29 (GMT+3) Access: PRIVATE

artifacts / genotype-imputation

Name

genotype-imputation-create-parameters-1769642917

genotype-imputation-download-sample-2118888266

genotype-imputation-run-command-731355685

Kubernetes over Slurm

```

fnikol@jedi12:~$ kubectl get pods -A -o wide
NAMESPACE      NAME                                     READY   STATUS    RESTARTS   AGE
imputation     argo-artifacts-66d557d8b7-qd7xp        1/1     Running   0           26m
imputation     argo-workflows-server-f9fc657b6-2qb42  1/1     Running   0           26m
imputation     argo-workflows-workflow-controller-75c87844b7-9fx2t  1/1     Running   0           26m
imputation     genotype-imputation-create-parameters-1769642917  0/2     Success   0 (25m ago)
imputation     genotype-imputation-download-sample-2118888266  0/2     Success   0 (26m ago)
imputation     genotype-imputation-run-command-731355685  0/2     Success   0 (25m ago)
imputation     genotype-imputation-workflow-controller-75c87844b7-9fx2t  0/2     Success   0 (23m ago)
fnikol@jedi12:~$ queue --format="%.18i %.9P %.30j %.8T %.10M %.9L %.6D %R" --me
JOBID PARTITION NAME USER STATE TIME TIME_LIMI NODES NODELIST (REASON)
10853 jedi argo-artifacts-66d557d8b7-9fx2t fnikol RUNNING 27:10 365-00:00:00 1 jedi1
10855 jedi argo-workflows-workflow-controller-75c87844b7-9fx2t fnikol RUNNING 26:50 UNLIMITED 1 jedi1
10856 jedi argo-workflows-server-f9fc657b6-2qb42 fnikol RUNNING 26:50 UNLIMITED 1 jedi1
    
```

Overlay Networking

Virtual Node

User job

- **Objective:** Run Kubernetes on Slurm.
- **Challenges**
 - Easy Deployment
 - All key components (API server, etcd, CoreDNS, ...) are packaged in a container.
 - Rootless execution
 - Implemented using Singularity containers.
 - Pod Support
 - Implemented using nested containers.
 - Network Services
 - Implemented using flannel.
- **Available at Github:** <https://github.com/CARV-ICS-FORTH/HPK>
- **System requirements:**
 - Singularity should allow running as fakeroot.
 - Singularity configured with Flannel (or other CNI) for assigning cluster-wide IPs.

The background features a dark blue color with a complex network diagram of white lines and small colored dots (blue, yellow, red) scattered across it. On the left side, there is a faint, circular seal or logo with text and symbols, likely representing a university or research institution.

Thank you !

We thankfully acknowledge the support of the European Commission and the Greek General Secretariat for Research and Innovation under the EuroHPC Programme through project EUPEX (GA-101033975). National contributions from the involved state members (including the Greek General Secretariat for Research and Innovation) match the EuroHPC funding.

FORTH

- To allow multi-user mappings, shadow-utils provides `newuidmap` and `newgidmap` (packaged by most distributions).
 - SETUID binaries writing mappings configured in `/etc/sub[ug]id`

```
/etc/subuid:  
1000:420000:65536
```

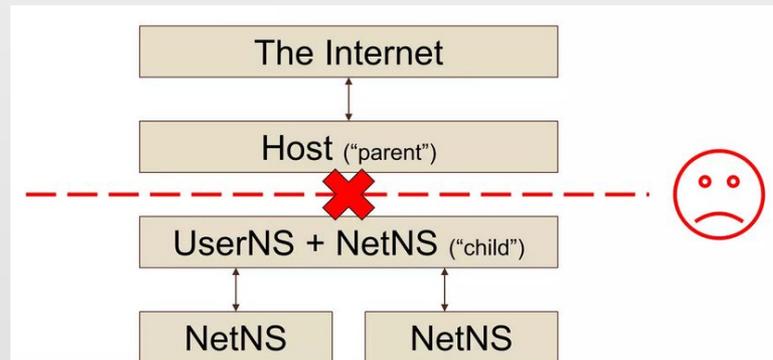
Provided by the admin (real root)

```
/proc/42/uid_map:  
0 1000 1  
1 420000 65536
```

User can configure map UIDs after unsharing a user namespace

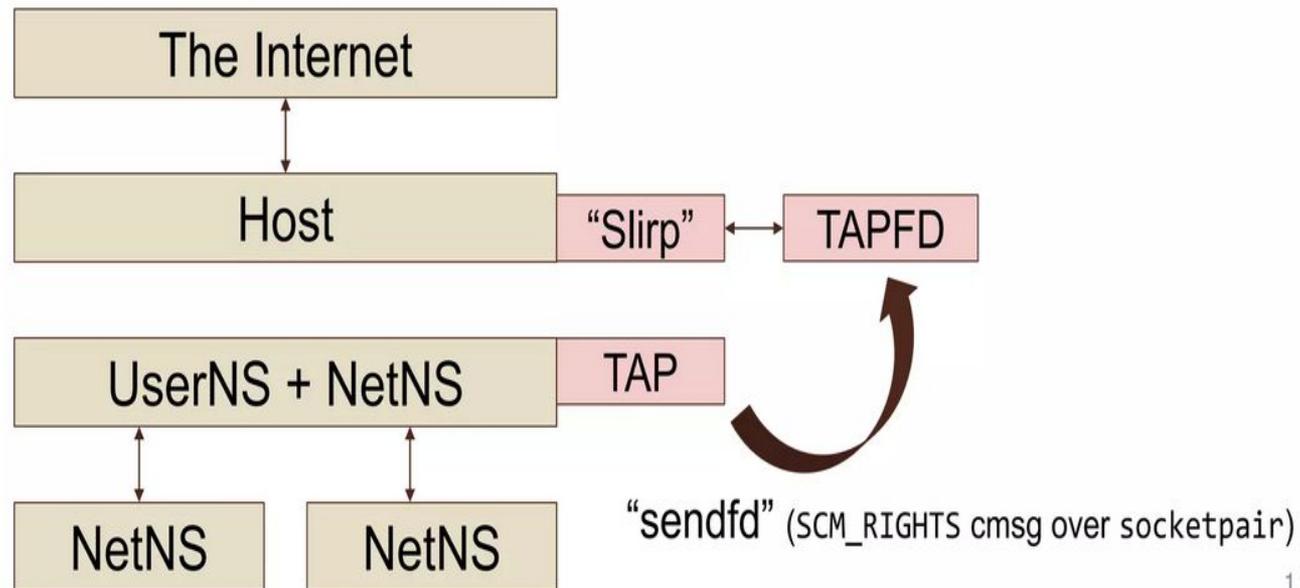
Challenge: Networking

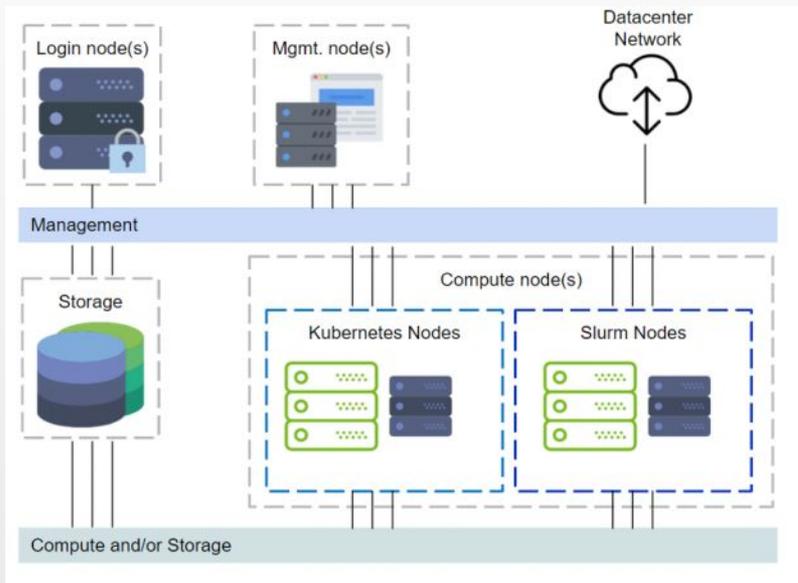
- An unprivileged user can create network namespaces along with user namespaces
 - For iptables, VXLAN, abstract socket isolation...
- But an unprivileged user cannot set up `veth` pairs across the host and namespaces, i.e. No internet connection
 - User-mode network stack (“Slirp”) can be used instead



- **Prior work: LXC uses SETUID binary (lxc-user-nic) for setting up the VETH pair across the parent and the child namespaces**
- **Problem: SETUID binary can be dangerous!**
 - CVE-2017-5985: netns privilege escalation
 - CVE-2018-6556: arbitrary file open(2)

- Our approach: use usermode network (“Slirp”) with a TAP device (<https://github.com/rootless-containers/slirp4netns>)
 - Similar to `qemu -netdev user`
 - Completely unprivileged
 - iperf3 benchmark on Travis: 9.21 Gbps

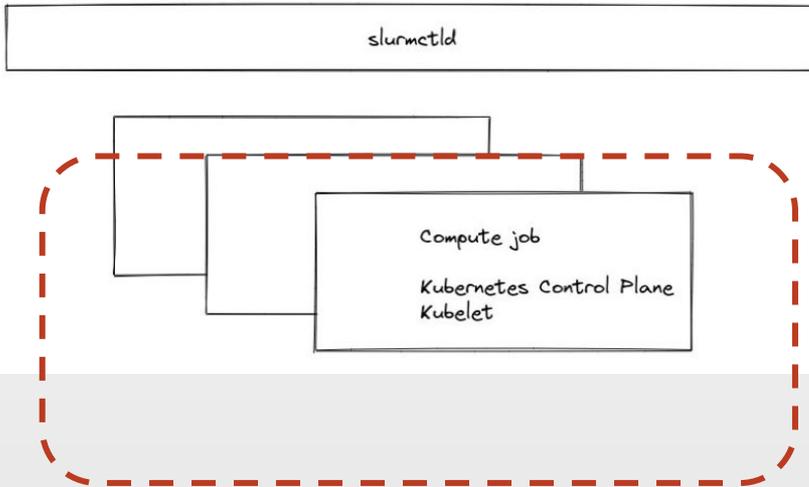




Setup: Partition K8s and Slurm nodes.¹

- Pros:
 - Full access to k8s capabilities
 - Full access to Slurm capabilities
- Cons:
 - Poor Interfacing between k8s and Slurm.
 - Data transfers from one partition to another.
 - Doubles the maintenance cost.

¹ KNoC is a Kubernetes node to manage container lifecycle on HPC clusters (InteractiveHPC 2022) <https://github.com/CARV-ICS-FORTH/knoc>

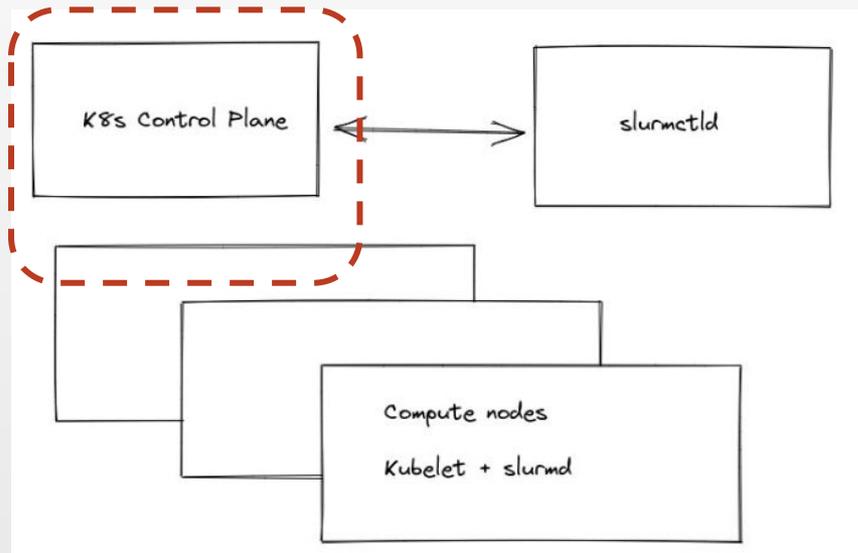


Setup: Run Slurm cluster(s) within a K8s environment.²

- Pros:
 - Elastic use of Cloud resources.
 - Portability of HPC solutions across Cloud.
 - Traditional experience for Slurm users.
- Cons:
 - Does not address the site underutilization issue.

² Genisys is a Kubernetes scheduler for running HPC jobs inside Virtual Clusters alongside other services (VHPC'22) <https://github.com/CARV-ICS-FORTH/genisys>

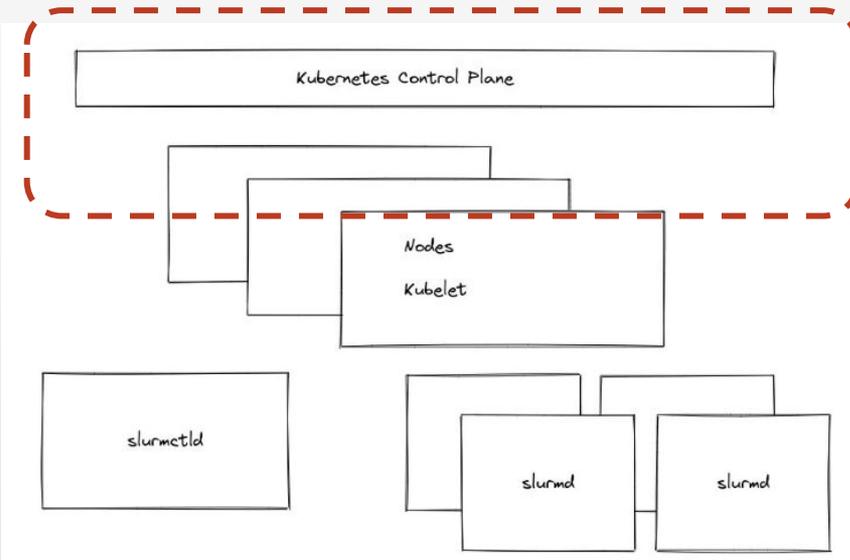
“Adjacent”



Setup: Run both Slurm and k8s on the same nodes.

- Pros:
 - Full access to Slurm capabilities
 - Full access to k8s capabilities
 - No data transfers required.
- Cons:
 - Security concerns (k8s runs as root)
 - Resource conflicts (nodes vs pods)
 - Increased Maintenance costs

Our vision - "Over"



Setup: Run K8s cluster(s) within unmodified Slurm environment.

Design Goals:

1. Create ephemeral k8s clusters as **user jobs**.
2. Support all Kubernetes abstractions, except privileged.
3. Scale across all nodes of the cluster.
4. Minimal pre-installed software.