# Demo sprint 1

## HPC scheduling with Kubernetes

Nidhi Shah, Yilin Xu, Siyuan Chen, Juhi Paliwal, Soufiane Jounaid

# Agenda

- Background

- Problem and end goal

- Environment setup

- Use cases

- Burndown chart

- Plan for next sprint

# Background

Kubernetes is a framework that runs distributed systems

- Pods: Containers containing application
- Nodes: Worker Machine that run Pods
- Control Plane: Manages Nodes

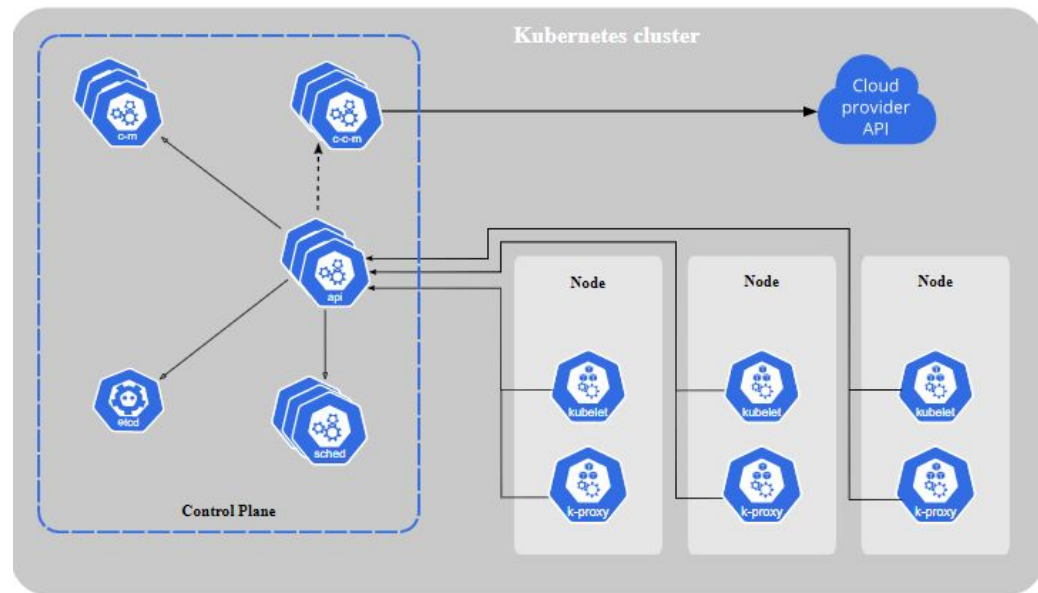**Kube-Scheduler** selects the nodes that will run the new pods

**Kubelet** starts a pod

**APIserver** provides communication between components.

Limitations of HPC in Kubernetes:

- No Batch, Co-Scheduling
- No proper queuing system
- Kube-scheduler not adapted to HPC.

Kubernetes allows you to add your own scheduler.



Components of Kubernetes

# Kube-flux

Flux is a framework that uses a graph based scheduler called Fluxion.

Kube-Flux is an HPC scheduler that employs the Fluxion library to take scheduling decisions in Kubernetes.

Scheduling Cycle:

- Transform Kubernetes Pod specification to Flux Job Specification
- Submit placement request to Fluxion Library.
- The Fluxion library finds the node on the cluster and returns its name.
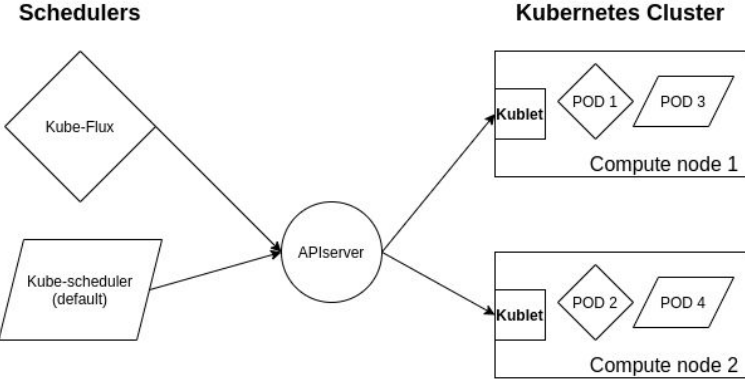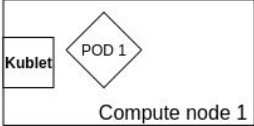- Pass the selected node name to the API to bind with a pod.

Binding Cycle

- Pod are bound to a node

# Problem

- Flux is not in charge of starting the actual pod, flux simply submits the pod to the Kubernetes APIserver which in turn works with Kublet to start the pod.
- Since Flux is not directly responsible for starting pods, it is not directly aware of their state on the cluster.
- If a pod scheduled by kube-flux finishes or is cancelled, kubernetes will remove said pod from the cluster and free the resources; this information however does not get propagated back to Kube-flux.
- Kube-Flux now sees a resource that has been freed as still occupied, which affects the correctness of future scheduling decisions.
- Consequently, since kubernetes does not receive state updates from the cluster, it is also unaware of the effects of scheduling decisions taken by other schedulers.
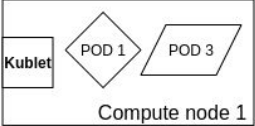
# Illustration of the problem

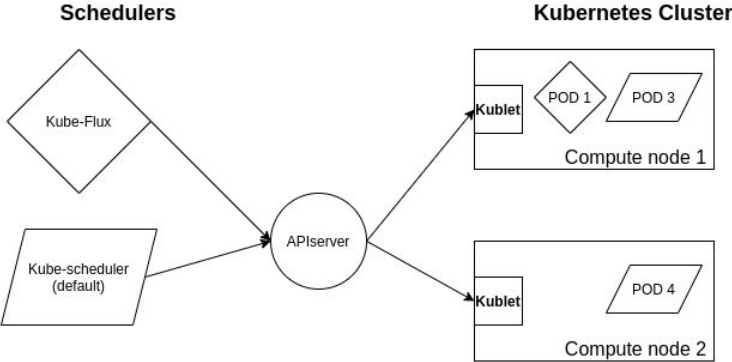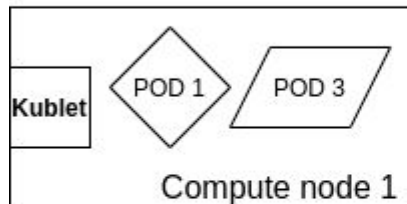# Goal of the project

- Full goal: Continuously synchronize the internal state of Flux with the state of the Kubernetes cluster of compute nodes.

    - Subgoal 1: Synchronize the Kube-flux scheduler to reflect changes in the state of its <u>own scheduled pods.</u>
    - Subgoal 2: Synchronize the Kube-flux scheduler to reflect changes in the state of <u>pods scheduled by other schedulers.</u>
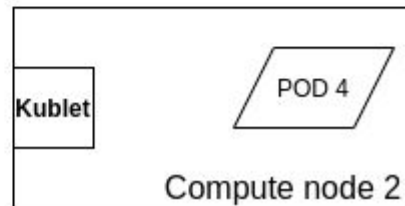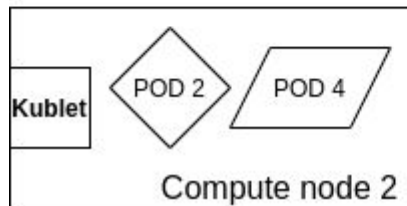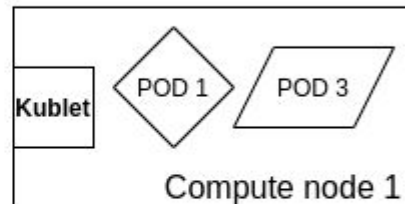
# Illustration of our goals (simplified)



What kube-Flux should start seeing by the end of this project

# Potential solutions

- Obvious solution
  - Before each kube-flux scheduling round, we re-collect information about the compute nodes current utilization data and rebuild Kube-flux's internal graph using that information.
- Slightly less obvious solution
  - Collect the changes to about the state of the compute nodes and relay this information to Kube-flux which can then use the info to update its internal state without having to rebuild the graph.

# Environment

- OS: Linux
- Dependencies:
  - Go: the language used in the development
  - Minikube: used to setup a local Kubernetes cluster
  - Helm: manage Kubernetes application
  - Docker: build images
  - CodeReady Containers (TBD): run local OpenShift clusters on a laptop
- Git Repository
  - Work together with mentors on a GitHub repo
- Local Repository
  - Clone the git repository
  - Test `make local-image`, the process takes about 8 mins

# Environment

# Use cases

- Allow HPC applications that can run on a cloud cluster as well as they can run on an HPC cluster to be scheduled on a cloud cluster.
- Run both HPC and cloud workloads on a unified cluster.

# Burndown Chart



20 Points
8 Tasks

# Sprint 2

- Setting up OpenShift resources
- Get familiar with Flux's internal representation of the state of the Kubernetes cluster.
- Devise a testing plan by composing expected correct Flux states to test actual states against.

# Thank You