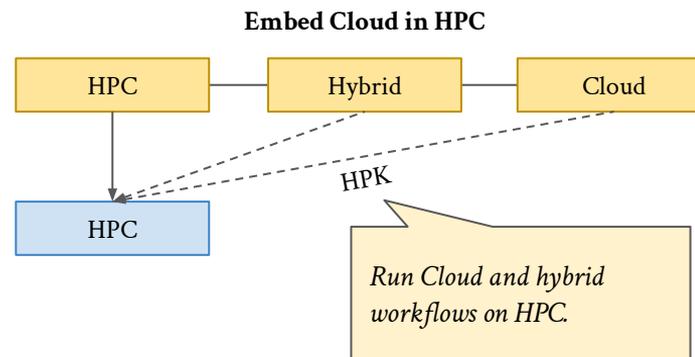# Challenges and Opportunities in Running Kubernetes Workloads on HPC

5/6/2024 • PASC24

A. Chazapis, E. Maliaroudakis, F. Nikolaidis, M. Marazakis, A. Bilas
FORTH-ICS/CARV

# Objectives

- HPC vs Cloud
  - Tight parallelization vs data distribution
  - Fixed setup vs unlimited resources
  - Low vs. high-level frameworks
  - Run binaries vs. webs of microservices
  - Console access vs interactive GUIs

- Combine best of both worlds using hybrid workflows → run on HPC
  - Provide a way to run Cloud software in HPC
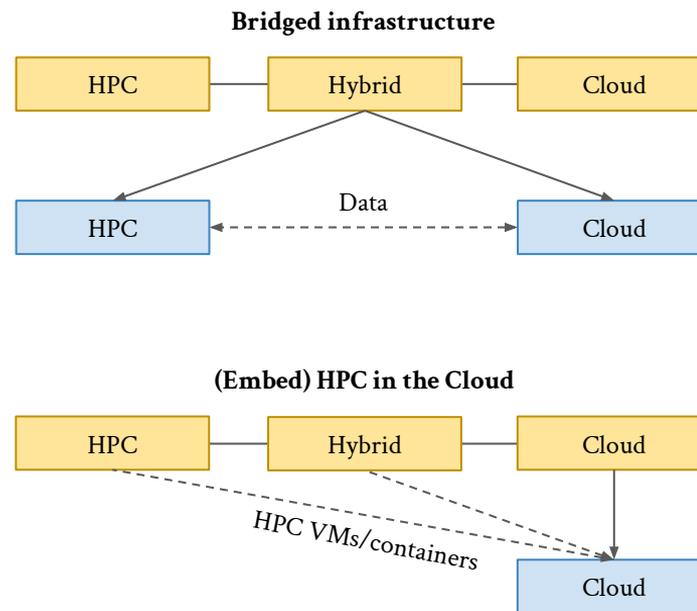  - Create ephemeral Kubernetes mini-Clouds in HPC

*We need to enable portable, reproducible, composite workflows*
*— Daniel Milroy, LLNL (WOCC'23)*

**Embed Cloud in HPC**

| HPC | — | Hybrid | — | Cloud |

HPC

HPK

*Run Cloud and hybrid workflows on HPC.*

# Running hybrid workloads

- **Bridge two separate environments**
  - Submit HPC jobs from the Cloud side or vice versa[1]
  - Deal with separate data and network contexts, unless running HPC in the Cloud (VM offerings available)
  - Hardware and maintenance costs
- **Run in Kubernetes**
  - Embed the HPC software stack in containers → Delegate scheduling to Kubernetes[2]
- **Install both on the same hardware**
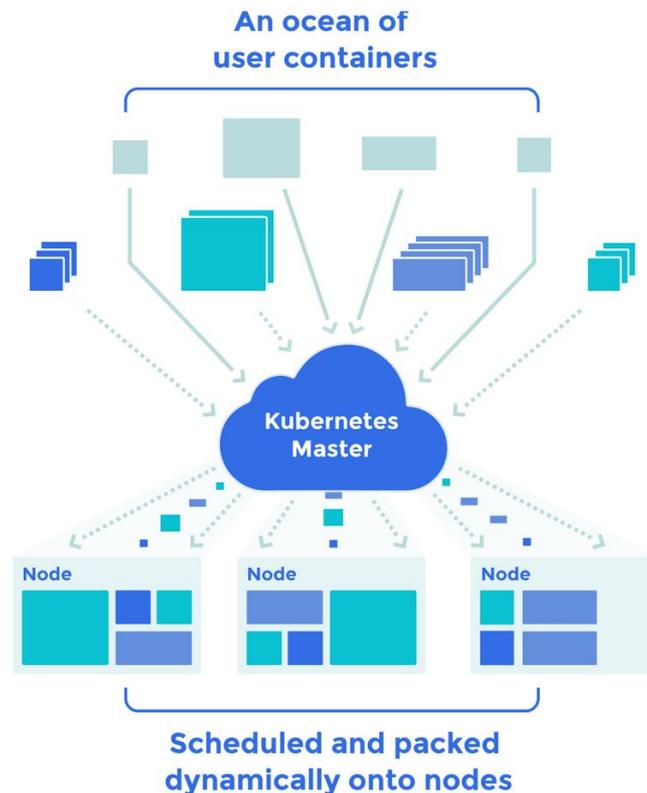  - Inpractical due to interference

**Bridged infrastructure**

| HPC | — | Hybrid | — | Cloud |

| HPC | ← Data → | Cloud |

**(Embed) HPC in the Cloud**

| HPC | — | Hybrid | — | Cloud |

HPC VMs/containers

| Cloud |

[1] KNoC is a Kubernetes node to manage container lifecycle on HPC clusters: https://github.com/CARV-ICS-FORTH/knoc (InteractiveHPC 2022)
[2] Genisys is a Kubernetes scheduler for running HPC jobs inside Virtual Clusters alongside other services: https://github.com/CARV-ICS-FORTH/genisys (VHPC'22)
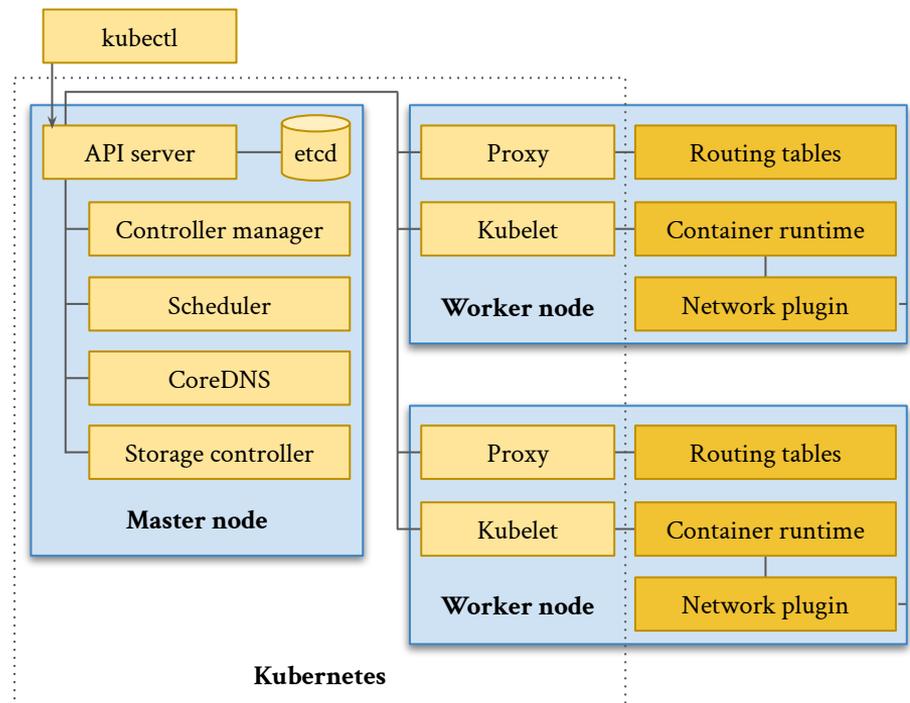
# So... Kubernetes?

Kubernetes is a *container orchestration runtime*

- It manages the container lifecycle
  - Containers are lightweight (vs. VMs) and portable
  - Interface to the container runtime → Docker/containerd

- It is not only a scheduler
  - Handles networking between containers
  - Provides service discovery and load balancing mechanisms
  - Reacts to load (scaling) and failures

- It runs almost everywhere
  - Any scale, most architectures → desktop to Cloud
  - Runs as a system service → Needs "elevated" permissions

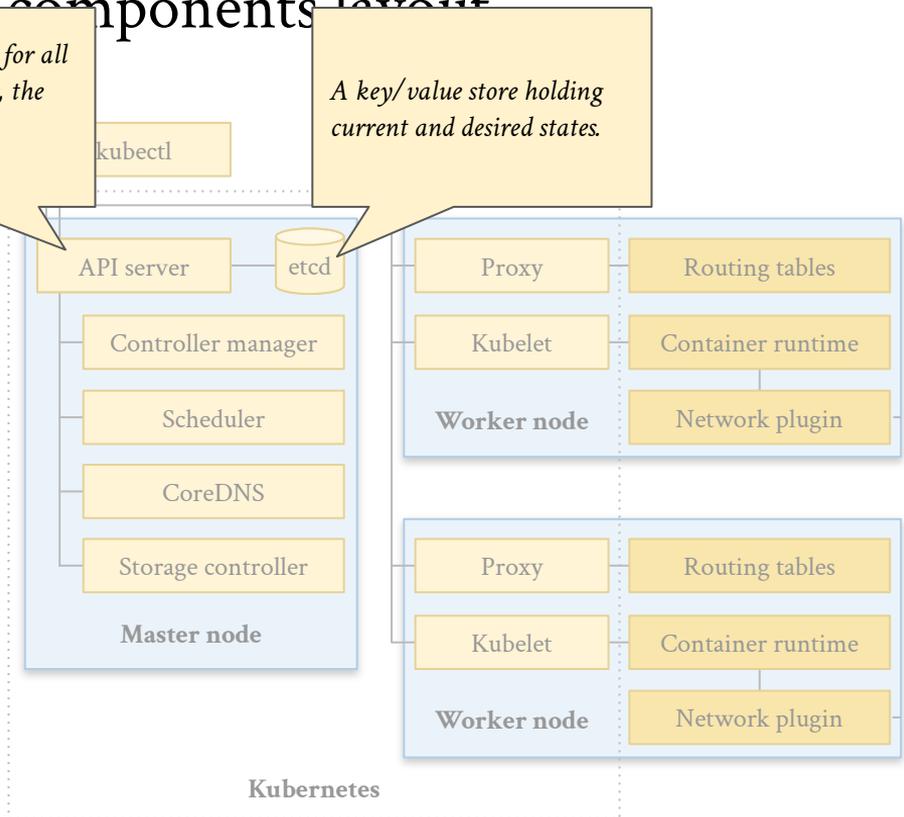**An ocean of user containers**

**Kubernetes Master**

**Node**   **Node**   **Node**

**Scheduled and packed dynamically onto nodes**

Image source: https://discuss.newrelic.com/t/relic-solution-what-you-need-to-know-about-new-relic-when-deploying-with-docker/52492

# Kubernetes concepts

- **Declarative vs imperative**
- **API endpoint & controllers**
- **Abstractions**
  - Pods → Collection of containers
  - Deployments → Replicated pod groups
  - Services → Microservice naming
  - Jobs → Pods that run to completion
  - Volumes → Mountable file collections
  - Labels → Queryable metadata
- **DevOps compliant**
  - Infrastructure as code
  - Version rollouts, CI/CD workflows
- **Typical distributed structure**
  - "Control plane" → Scheduling and placement
  - Node agents → Handle execution
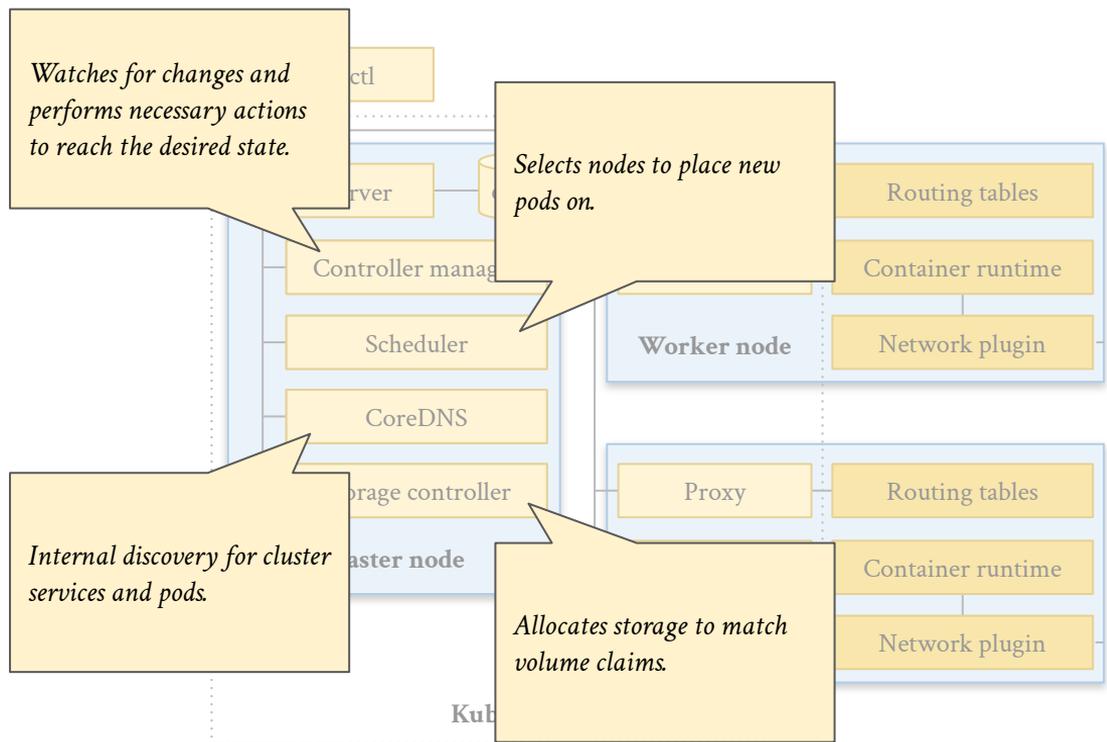  - Monitoring and accounting infrastructure

# Kubernetes core components layout



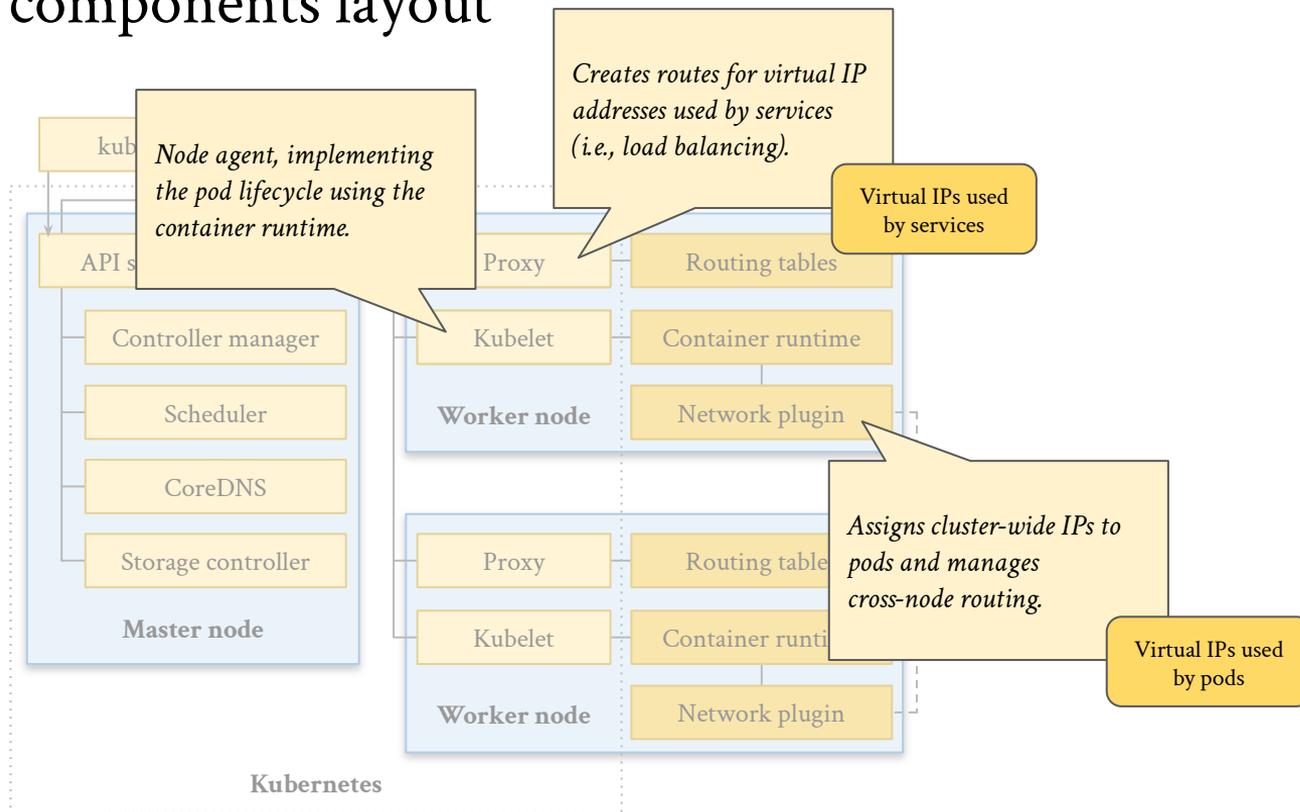*The main entrypoint for all requests to the cluster, the point of sync for all controllers.*

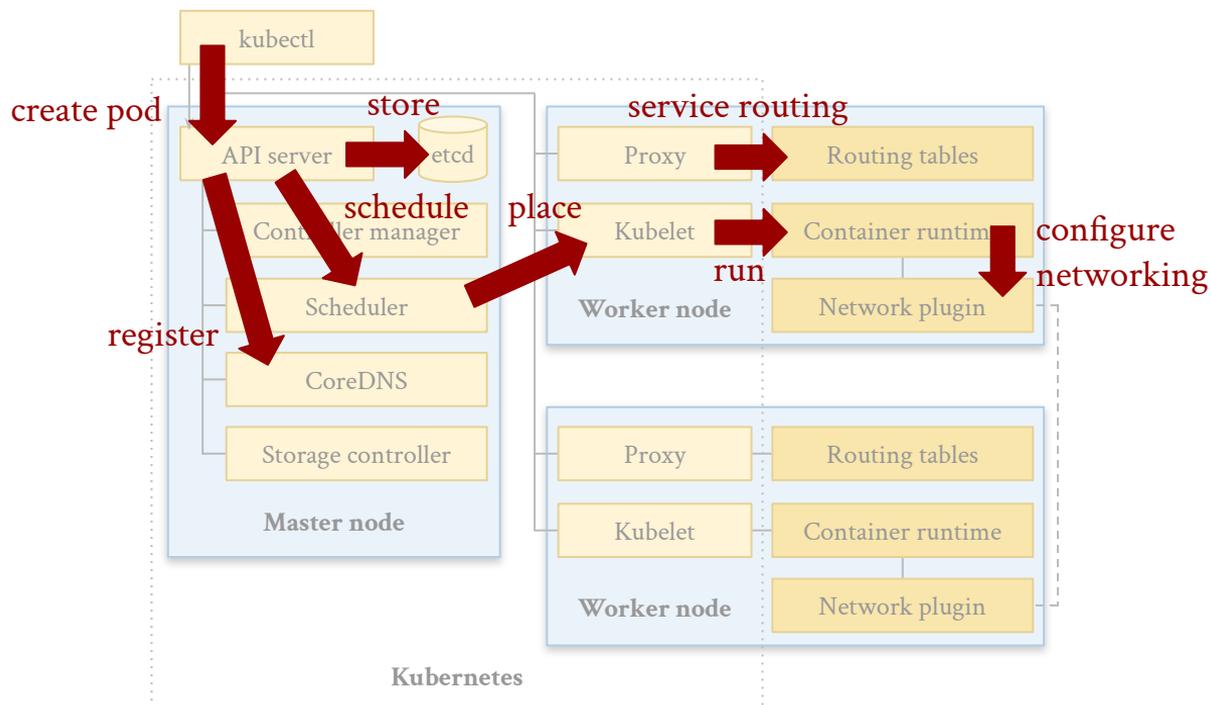*A key/value store holding current and desired states.*

kubectl

**Master node**
- API server
- etcd
- Controller manager
- Scheduler
- CoreDNS
- Storage controller

**Worker node**
- Proxy — Routing tables
- Kubelet — Container runtime — Network plugin

**Worker node**
- Proxy — Routing tables
- Kubelet — Container runtime — Network plugin

**Kubernetes**

# Kubernetes core components layout

# Kubernetes core components layout
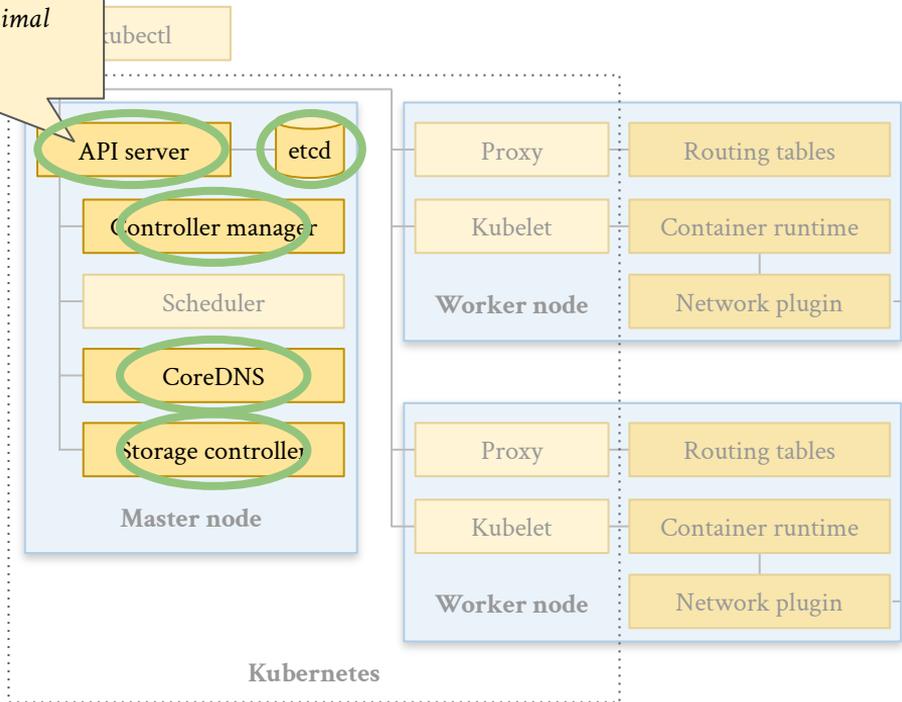
# Kubernetes core components layout

# Design goals

Run *Kubernetes in an HPC environment* as a user → High-Performance Kubernetes (HPK)
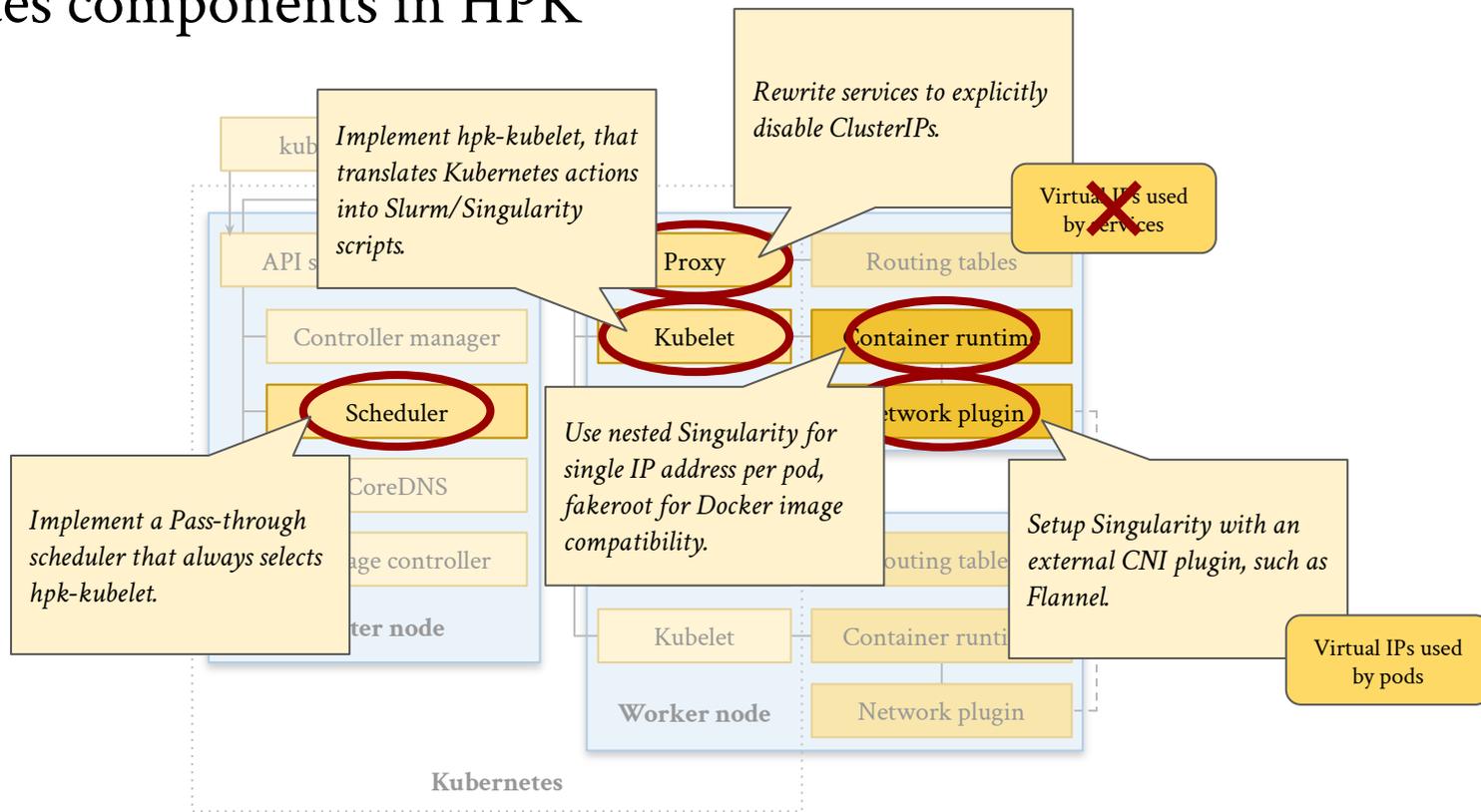
Requirements:

- **All Kubernetes abstractions should be available and fully functional**
  - Except those that affect physical hardware resources (like NodePort services)
  - Private, inter-container network and internal DNS should work as expected
- **Delegate resource management to Slurm**
  - Respect organization policies
  - Comply with established resource accounting mechanisms
  - Scale across all nodes of the cluster

- **Use Singularity as the container runtime**
  - Preinstalled in most HPC environments
- **Make it easy for HPC administrators**
  - No (or little) configuration changes should be required at the host level
  - No reliance on special libraries or binaries that execute with "elevated" permissions
- **Make it easy for users**
  - All neatly packaged up in one container
  - Simple, one-command deployment via Slurm
  - All relevant configuration and files should be in the user's home folder
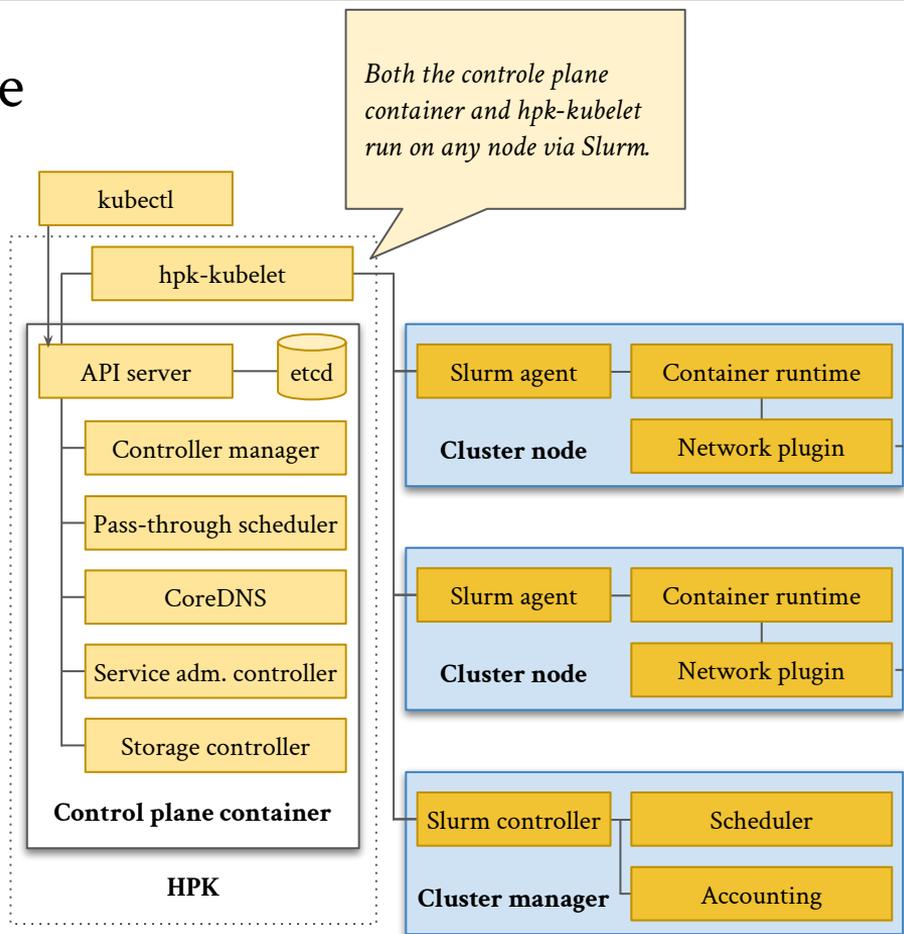
# Kubernetes components in HPK

*Services that run "out of the box", with no special dependencies, or minimal changes.*

kubectl

**Master node**

API server — etcd

Controller manager

Scheduler

CoreDNS

Storage controller

**Worker node**

Proxy — Routing tables

Kubelet — Container runtime

Network plugin

**Worker node**

Proxy — Routing tables

Kubelet — Container runtime

Network plugin

**Kubernetes**

# Kubernetes components in HPK



Implement hpk-kubelet, that translates Kubernetes actions into Slurm/Singularity scripts.

Rewrite services to explicitly disable ClusterIPs.

Virtual IPs used by services

Implement a Pass-through scheduler that always selects hpk-kubelet.

Use nested Singularity for single IP address per pod, fakeroot for Docker image compatibility.

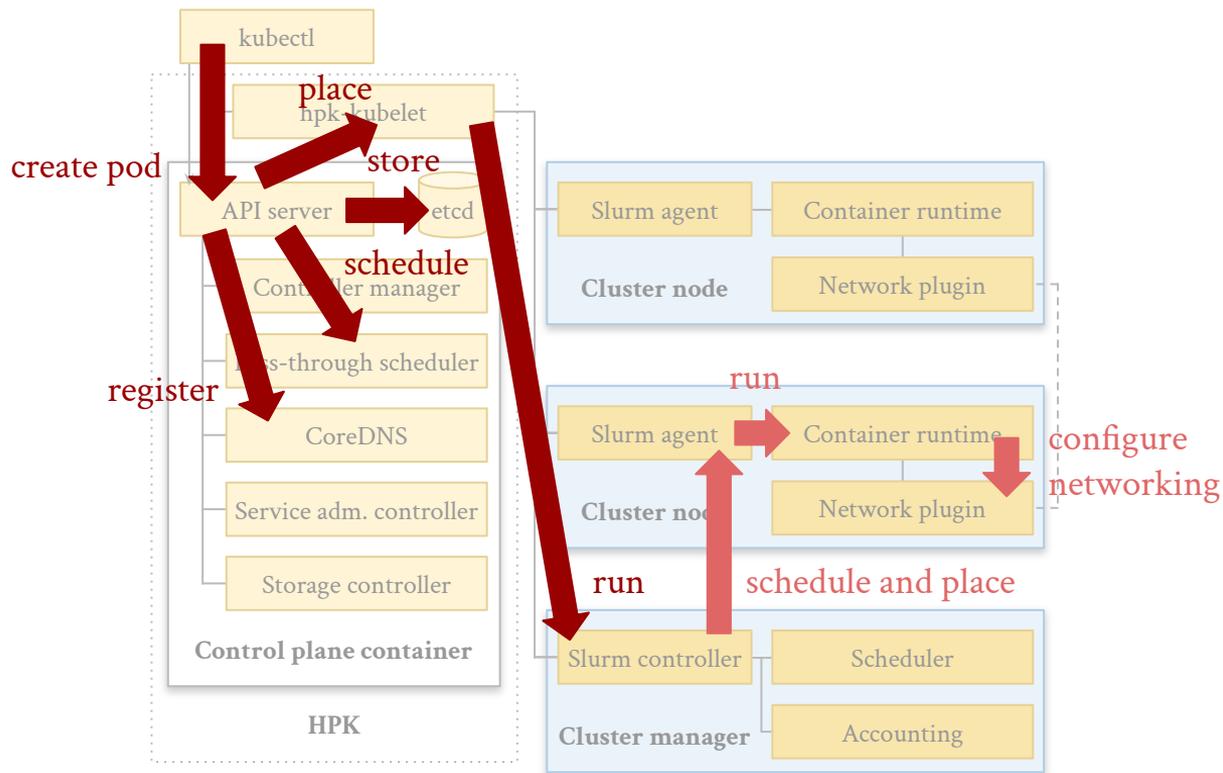Setup Singularity with an external CNI plugin, such as Flannel.

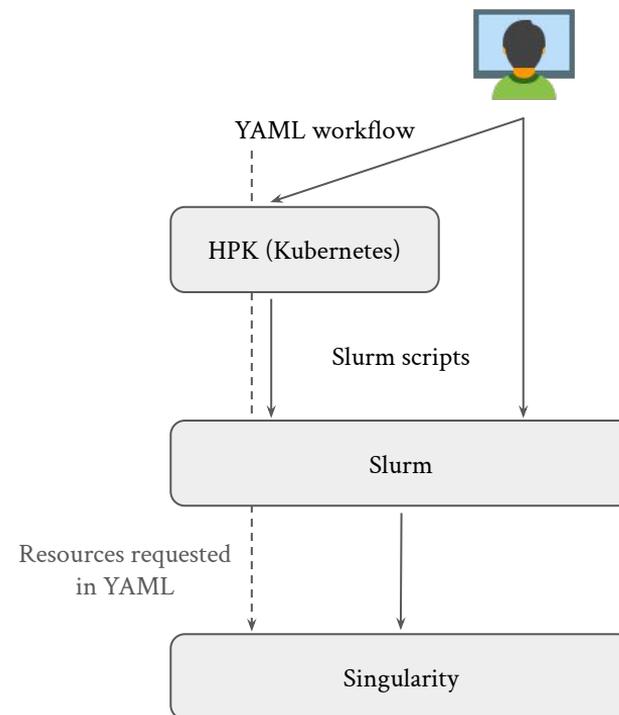Virtual IPs used by pods

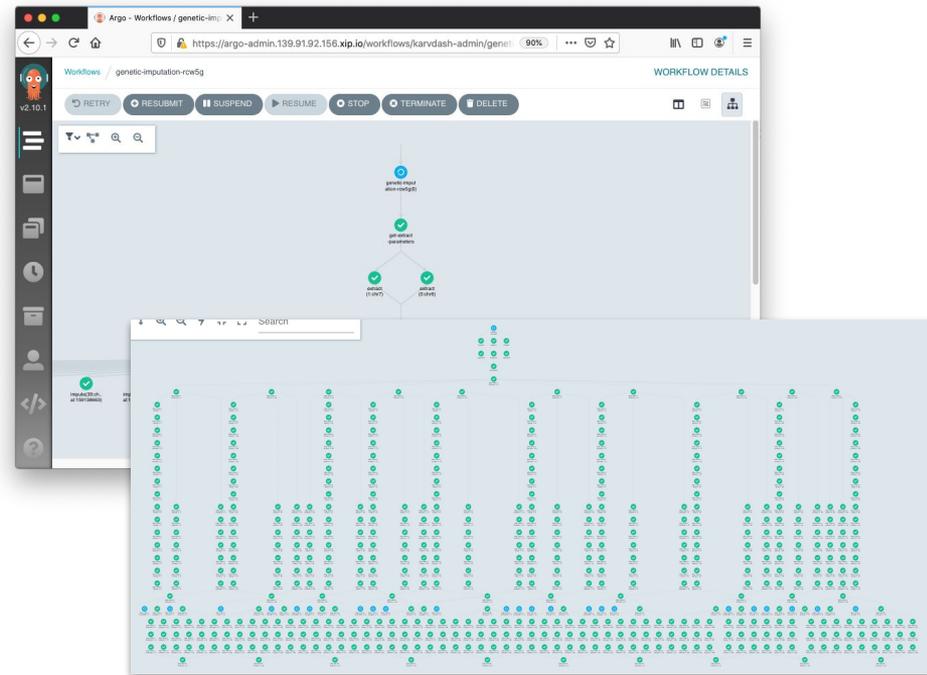# HPK architecture

# HPK architecture

# HPK implementation

- **HPK translates Kubernetes to Slurm scripts**
  - Pods/jobs show enter as YAML through the Kubernetes API, exit as Slurm scripts from hpk-kubelet → Pods show up as Slurm jobs
  - Kubernetes resource requirements end up in Slurm allocation requests → No changes to accounting
- **HPK runs as a user process via Slurm**
  - User can run both Kubernetes and Slurm workloads at the same time
  - No "special" allocation needed for HPK → 1 CPU, few GBs of RAM should be enough
  - Little support needed by the environment → No Slurm modifications, some Singularity configuration (inc. Flannel or other CNI)

# HPK for Cloud-native workflows

- **Using Argo Workflows**
  - All steps are containers → Containerized code is portable
  - Interactive UI → Monitor, control

- **Language features**
  - Loops, conditionals
  - High-level parameters (shared across all steps), workflow templates
  - Artifacts → Step I/O

- **HPK extensions → Integrate MPI steps**
  - Slurm passthrough flags via annotations → Control scalability, allocate GPUs
  - Still need to containerize code

# HPK for Cloud-native workflows

- Using Argo Workflows
  - All steps are containers → Containerized code is portable
  - Interactive UI → Monitor, control

- Language features
  - Loops, conditionals
  - High-level parameters (shared across all steps), workflow templates
  - Artifacts → Step I/O

- HPK extensions → Integrate MPI steps
  - Slurm passthrough flags via annotations → Control scalability, allocate GPUs
  - Still need to containerize code
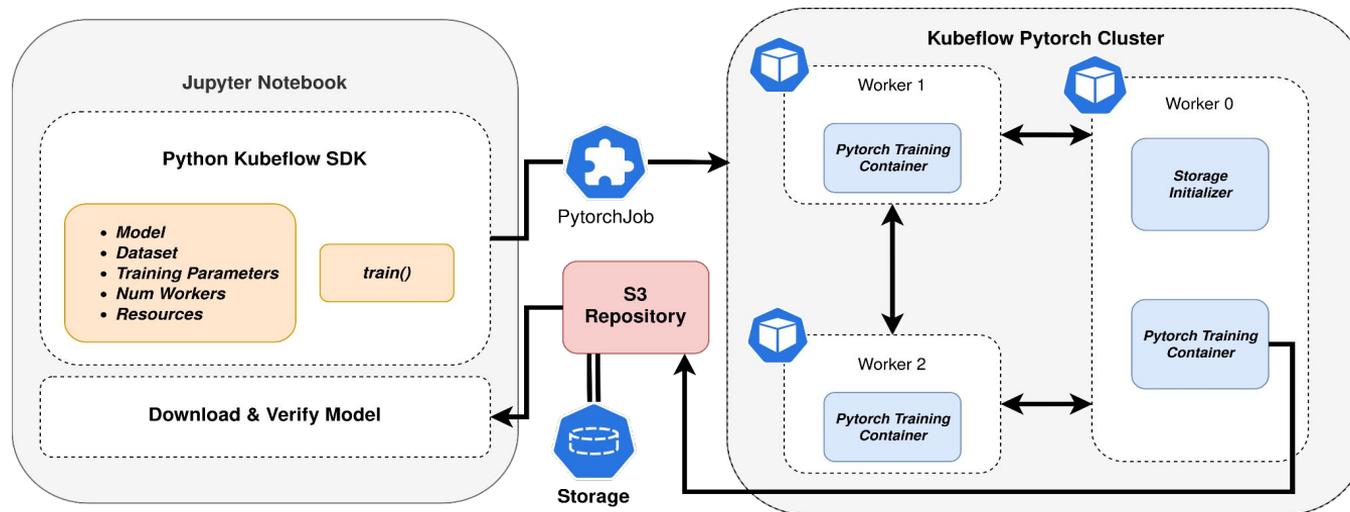
```
 1  kind: Workflow
 2  metadata:
 3    ...
 4  spec:
 5    entrypoint: npb-with-mpi
 6    templates:
 7    - name: npb-with-mpi
 8      dag:
 9        tasks:
10        - name: A
11          template: npb
12          arguments:
13            parameters:
14            - {name: cpus, value: "{{item}}"}
15          withItems:
16          - 2
17          - 4
18          - 8
19          - 16
20    - name: npb
21      metadata:
22        annotations:
23          slurm-job.hpk.io/flags: "--ntasks={{inputs.parameters.cpus}}"
24          slurm-job.hpk.io/mpi-flag : "..."
25      inputs: s
26        parameters:
27        - name: cpus
28      container:
29        image: mpi-npb:latest
30        command: ["ep.A.{{inputs.parameters.cpus}}"]
```

17

# HPK for ML Workflows: Fine-Tune BERT with PyTorch

- Diverse microservices
  - Jupyter notebook → Coordinates tasks (data preprocessing, model fine-tuning, testing)
  - Kubeflow controller → Manages Pytorch jobs
  - MinIO → S3-compatible object store
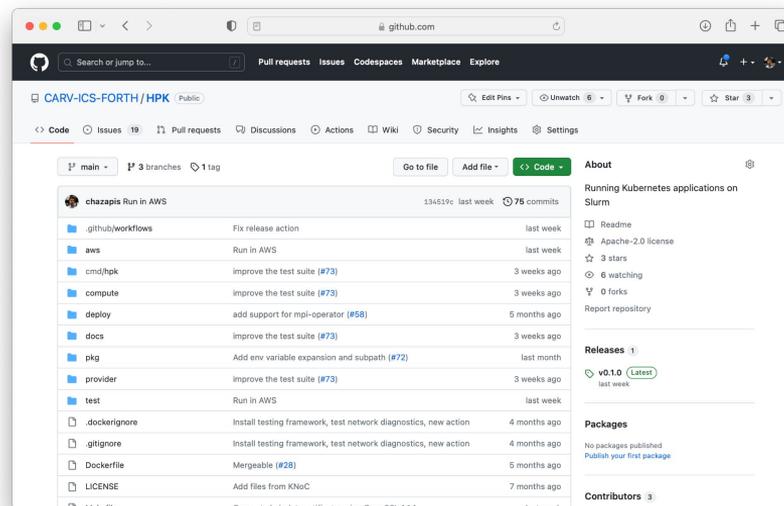
# HPK for Spark in HPC

- **Evaluated HPK + Spark at Jülich (JSC)**
  - TPC-DS benchmark
  - Spark Operator (EKS spark-benchmark) → Manages Spark applications on Kubernetes
  - MinIO → S3-compatible object store
- **Challenges**
  - Preparation of testbed → Apptainer (fakeroot support), Flannel (deploy on all nodes, configure subnets, requires etcd)
  - Integration with Slurm → Exclusive node policy does not align with container sizes

# HPK vision

- Cloud-user PoV → Run on HPC hardware

- HPC-user PoV → Exploit the Cloud software ecosystem
    - Combine HPC codes with backend services (database, queueing systems)
    - Interactive code execution → Jupyter
    - Workflow management → Argo Workflows, Apache Airflow, …
    - Monitoring utilities → Grafana
    - Frameworks for automatically optimizing and scaling code → Spark, DASK, …

- HPC centre PoV → Run Cloud workloads on the main HPC partition, attract users
    - The common practice is to have separate partitions for Cloud (analytics) and HPC

*HPK is available at* https://github.com/CARV-ICS-FORTH/HPK



*Acknowledgements*