# OLCF Container Orchestration for HPC Middleware

Jason Kincl, HPC Systems Administrator

OAK RIDGE
National Laboratory

# Multiple Container Strategies at OLCF

- **Container orchestration**: Automate deploying and operating service containers with Kubernetes/OpenShift
  - Focused on framework for providing resources (cpu, memory, network, …) for running services
  - Uses own scheduler and is a separate resource from HPC
- **HPC container runtimes** with Singularity
  - Focused on using containers to run applications in a batch job
  - Uses scheduler from batch job submission system
  - Allows users to provide a portable environment to run jobs on HPC resources

**OAK RIDGE**
National Laboratory

# What is HPC Middleware?

- The collection of applications and services that helps a project achieve it's scientific goals

- Some examples:
  - Workflows
    - Data movement and job submission
    - Continuous integration and testing
    - Automation
  - Collaboration web portals for viewing data stored in OLCF
    - Jupyter notebooks
  - Streaming data
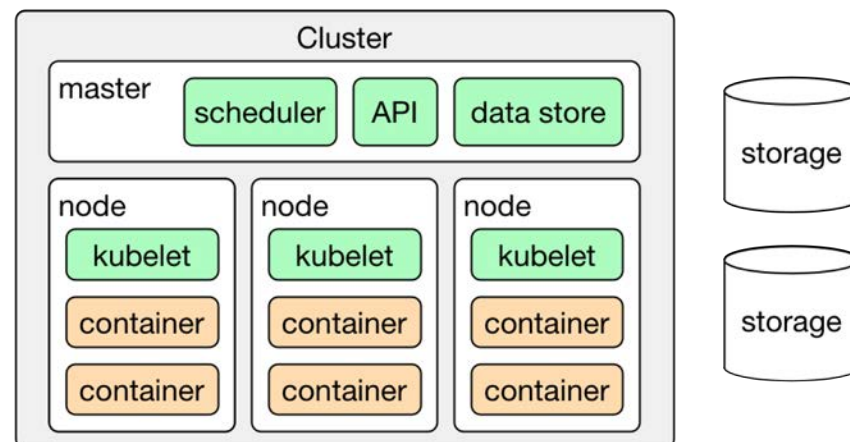
OAK RIDGE
National Laboratory

# Basic HPC Workflow Requirements

- Need ways for users to manage their workflow system
  - Diverse ecosystem of workflow systems makes it difficult for NCCS Operations to support every one

- Why not just use SSH keys?
  - Our moderate security controls require remote actions to be authenticated with RSA two-factor credentials
  - Instead we are working on providing ability for running workflow services locally

- Upon surveying existing workflow systems we came up with the following requirements:
  - Run a persistent service locally as a "daemon" that stays up
  - Talk to batch submission system for current queue information and job submission
  - Interact with files on GPFS/Lustre/NFS

**OAK RIDGE**
National Laboratory

# OpenShift

- Distribution of Kubernetes developed by Red Hat

- Kubernetes manages containerized applications across nodes and provides mechanisms for deployment, maintenance, and application-scaling.

- Analogous to but separate from a batch scheduler on a compute cluster
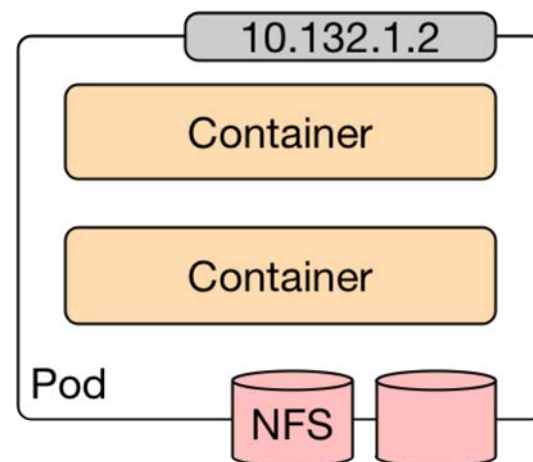
# Kubernetes Architecture

- Configuration: YAML or JSON data that describes the application being deployed

- Configuration can define:
  - Containers to run
  - HTTP routes and network ports to expose outside of the cluster
  - Mounting data volumes

```
apiVersion: v1
kind: ReplicationController
metadata:
  name: nginx
spec:
  replicas: 3
  template:
    metadata:
      labels:
        app: nginx
    spec:
      containers:
      - name: nginx
        image: nginx:1.10
```
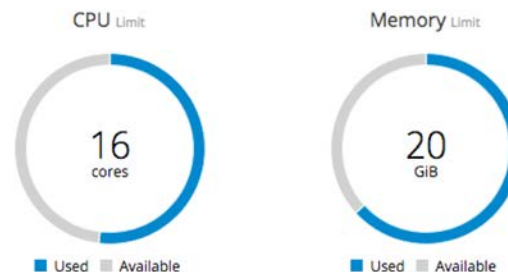
OAK RIDGE
National Laboratory

# Kubernetes – Pods

- Atomic unit of Kubernetes

- Made up of one or more containers deployed together on one host

- Pod lifecycle is defined, pod is assigned to run on a node and runs until the container(s) exit or it is removed for some other reason

- Volumes can be attached that do not share pod lifecycle for persistent data

- Each pod gets its own IP address that is accessible in the cluster

**OAK RIDGE**
National Laboratory

# Kubernetes – Replication Controllers

- Pod will not recreate itself if killed or deleted for some reason such as cluster maintenance or quota limit exceeded

- A ReplicationController ensures desired number of pods is running in the cluster

- For example: "I want to have three pods running nginx:1.10 image"

```
apiVersion: v1
kind: ReplicationController
metadata:
  name: nginx
spec:
  replicas: 3
  template:
    metadata:
      labels:
        app: nginx
    spec:
      containers:
      - name: nginx
        image: nginx:1.10
```

OAK RIDGE
National Laboratory

# Cluster Resources

- Resource allocation is different from the traditional core hours or node hours on Titan/Rhea

- Quota system based on CPU and memory limits

- User defines what CPU and memory are required for each container, if container exceeds limits it is killed

Quota  Learn More ⧉

compute-resources

Limits resource usage within this project.

CPU Limit

16
cores

■ Used  ▒ Available

Memory Limit

20
GiB

■ Used  ▒ Available

| Resource Type | Used | Max |
|---|---|---|
| CPU (Limit) | 8300 millicores | 16 cores |
| Memory (Limit) | 12956 MiB | 20 GiB |

OAK RIDGE
National Laboratory

# Exposing services

- OpenShift gives users the ability to expose services outside of the cluster
  - For HTTP-based services, NCCS will handle initial authentication to ensure service is accessed only by members of that project

# Accessing NCCS resources

- All containers run as an automation user that is tied to a project and has access to the project's allocation and files like a regular user

- Batch job submission from container
  - Users can base their container image off our NCCS golden image which comes with the tools to schedule batch jobs or get queue status

- Accessing shared filesystems (GPFS/Lustre/NFS)
  - Shared filesystems can be mounted in the container by Kubernetes allowing access just like a login or compute node

**OAK RIDGE**
National Laboratory

# Accessing OpenShift

- **OpenShift provides a command-line client as well as a web user interface**

# Current Clusters

- NCCS Moderate cluster – Marble
  - Available now
  - Same zone as Summit/Titan/Rhea/Atlas

- NCCS Open cluster – Onyx
  - Coming soon
  - Supports ORNL XCAMS user/password authentication instead of two-factor tokens (but no access to moderate resources like Atlas filesystem)
  - Access to Wolf filesystem

OAK RIDGE
National Laboratory

# Some Pilot Projects

- Continuous Integration with Jenkins and testing on Titan

- Project dashboard for viewing job output

- Running Fireworks workflow DB backend

- BigPanDA workflow services

- Database-based data portal for accessing project data

OAK RIDGE
National Laboratory

# Interested? Get in touch with us!

- Currently in a pilot phase, if you have a use case that may benefit please reach out to us

    Email: help@olcf.ornl.gov

    Web: https://www.olcf.ornl.gov/for-users/user-assistance/

OAK RIDGE
National Laboratory