**Sandia National Laboratories**

**Exceptional service in the national interest**

# COMPUTING-AS-A-SERVICE INFRASTRUCTURE FOR ACCELERATING DIGITAL ENGINEERING
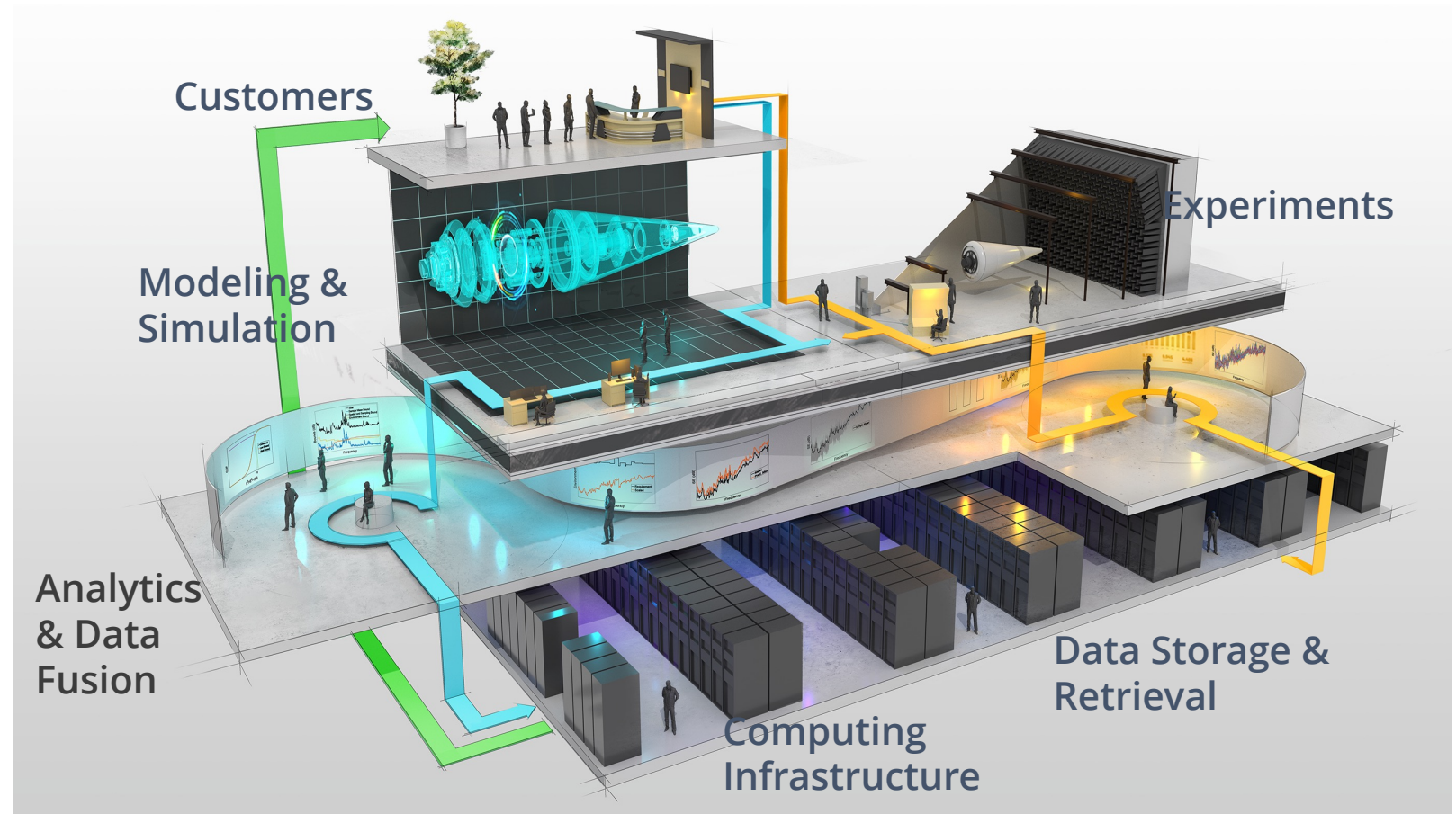
Eric Ho, Kevin Pedretti
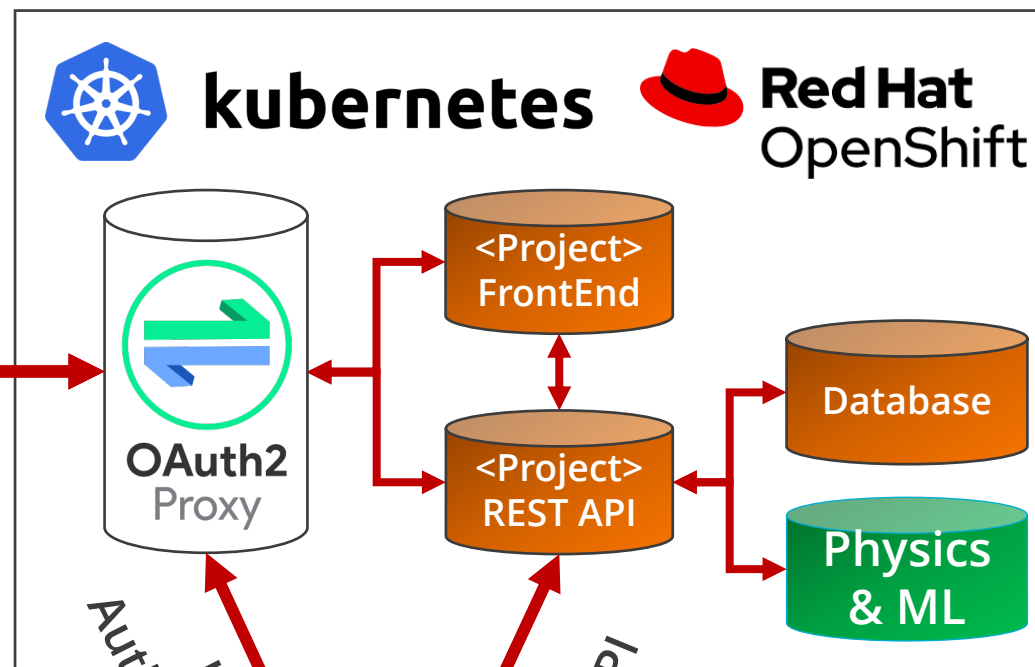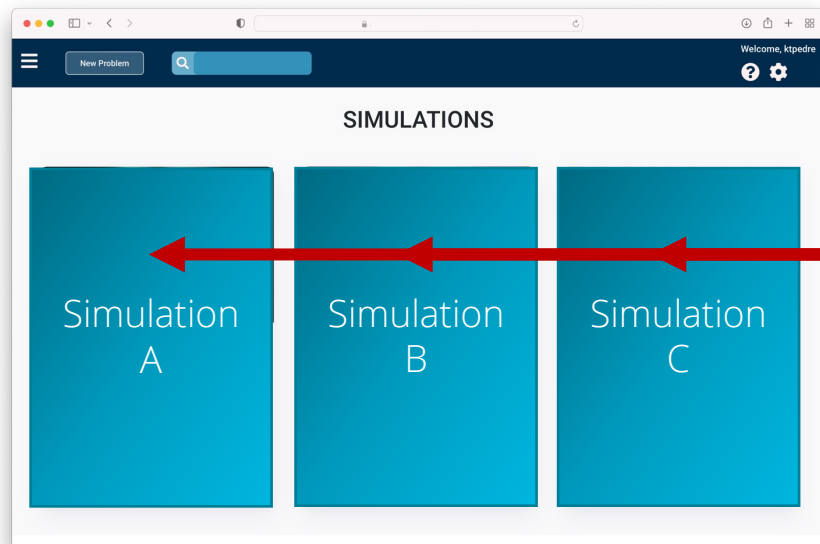
U.S. DEPARTMENT OF ENERGY

NNSA

# WHAT IS THE COMPUTING-AS-A-SERVICE PROJECT?

- Provide HPC as a cloud based service to teams with little to none HPC experience

- Customers interface with a GUI accessible through a web browser

- Jobs are intelligently routed to available HPC resources

- Customers are unaware of where their job is actually run

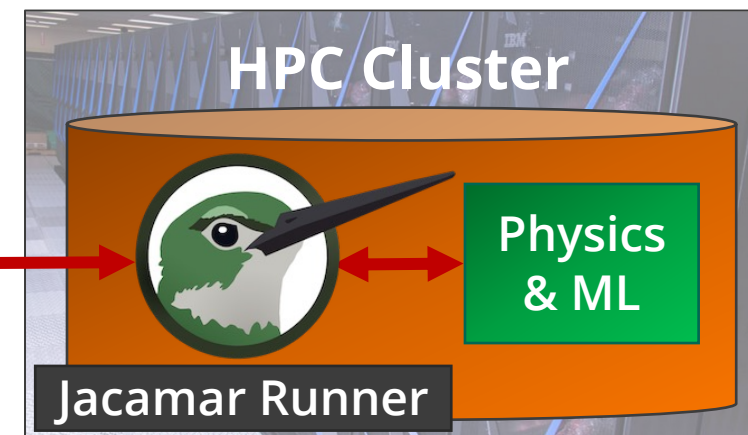# COMPUTING-AS-A-SERVICE ARCHITECTURE

**Customer's Front-end User Interface**



- Containerize all components (UI, Cloud, HPC)
- Deploy frontend via Kubernetes
- Job is sent to most available/fastest HPC cluster

Jacamar @ Sandia push from Scot Swan & Allen Robinson

# THE NEED FOR AUTOMATED DEPLOYMENT

## Versioned Code & Containers

## Kubernetes Clusters @ Sandia

**DEV**
MyApp-dev.sandia.gov
Sandia Enhanced
Azure Kubernetes Cluster

**PROD**
MyApp-prod.sandia.gov
Sandia Enhanced
Azure Kubernetes Cluster

**PROD2**
MyApp-prod2.sandia.gov
Sandia Common Eng. Env.
RedHat OpenShift Cluster

Application Deployment Code Repo

Application Container Registry

HELM

```
git clone <project>.git
# for each cluster
helm install <project> .
```

# CURRENT SUCCESSES AND CHALLENGES

- Deployed OpenShift Kubernetes testbed system to accelerate development

- 2 Production "Clusters".
  - OpenShift Kubernetes System
  - A100 DGX Station with Slurm

- Added GPU resources to the production OpenShift Kubernetes system so CaaS jobs have more nodes to run parallel codes

- Working with 4 different teams to provide CaaS to their projects. Each team requires unique frontend, backend, and hpc containers.

- Continue to develop, test, and deploy containerized apps for projects using CaaS