SAND2023-02999C

**Sandia National Laboratories**

# Exceptional service in the national interest

# Computing-as-a-Service Infrastructure for Accelerating Digital Engineering

Salishan Conference on High Speed Computing
Gleneden Beach, Oregon
April 24 - 27, 2023

Kevin Pedretti
Principal Member of Technical Staff
Scalable System Software, Org. 1423
ktpedre@sandia.gov

# Collaborators

**Computing-as-Service Team:**

- **Sylvain Bernard**
- **Ron Brightwell**
- **Wesley Coomber**
- **Mike Glass**
- **Eric Ho**
- **Todd Kordenbrock**
- **Cory Lueninghoener**
- **Aaron Moreno**
- **Kevin Pedretti**
- **Elliott Ridgway**
- **Gary Templet**
- **Andrew Younge**

**Guidance and Slide Material:**

- **Matthew Curry**
- **Ernest Friedman-Hill**
- **Chris Garasi**
- **Brenna Hautzenroeder**
- **Martin Heinstein**
- **Rob Hoekstra**
- **Jim Laros**
- **Scott Roberts**
- **Scot Swan**

# Outline

- Introduction

- Computing-as-a-Service Architecture

- R&D Directions / What's Missing

- Conclusion

# Computing-as-a-Service

*Computing as a service is a computing job that someone $ you to do*

11 year old's definition

- Cloud industry built around delivering things as a service

- Huge Business & Talent

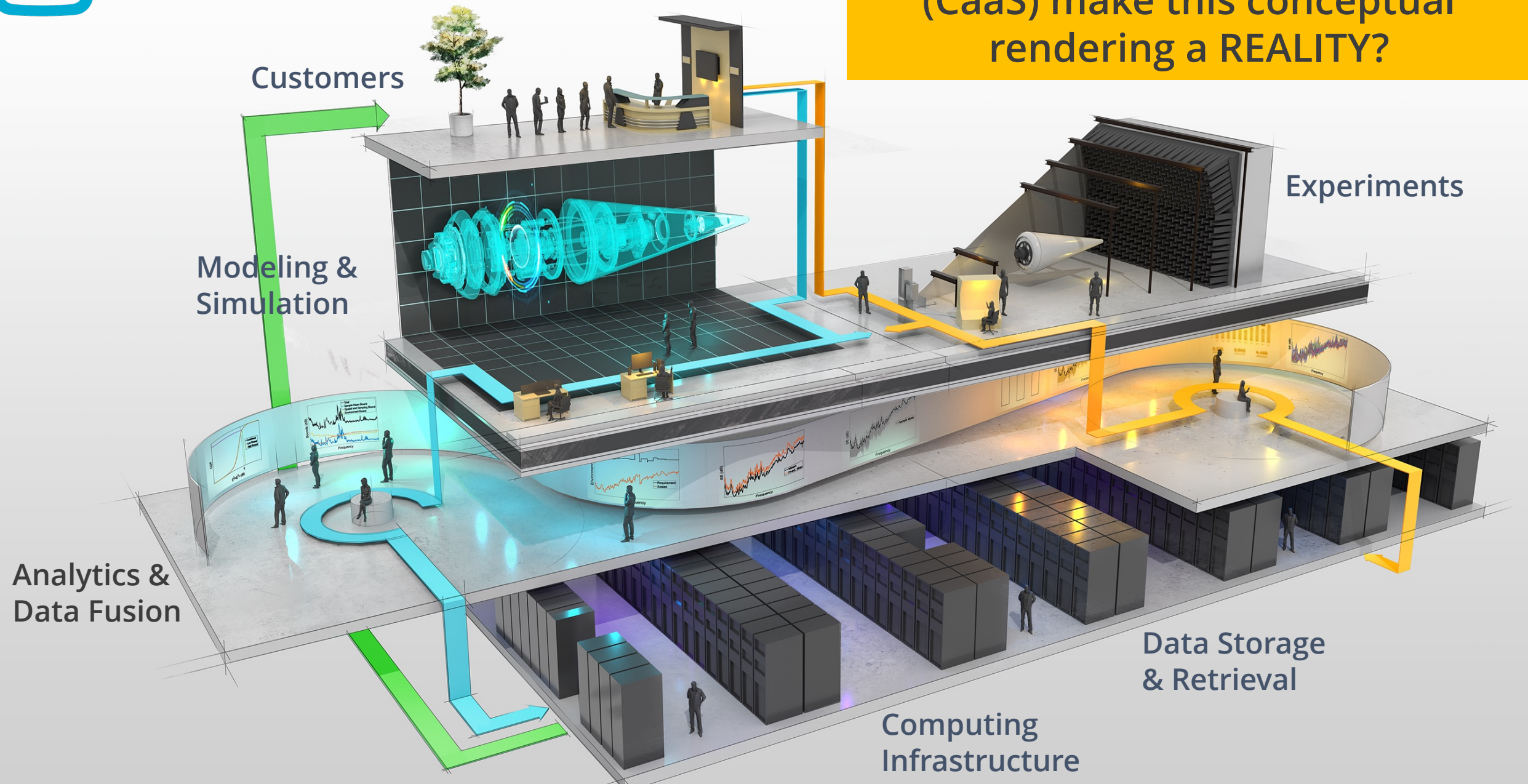- Software ecosystem for deploying turnkey services

# Cloud vs. HPC – Different Usage Models, Customs, and Practices

1. They use the same underlying technology – servers, storage, and networks

2. Cloud has 100's of services, HPC has ~ 1 **(HPC is the service)**

3. Cloud has APIs for managing all infrastructure and services, **(HPC APIs are ad hoc)**

4. Cloud uses token-based authentication, **(HPC uses passwords)**

5. Cloud runs the customer's software stack, **(HPC runs the facility's SW stack)**

6. Cloud charges by the hour (encouraging paranoia), **(HPC cycles are free)**

## Cross-Pollination of Cloud & HPC Mutually Beneficial

# Cloud vs. HPC – Different Usage Models, Customs, and Practices

1. They use the same underlying technology – servers, storage, and networks

2. Cloud has 100's of services, HPC has ~ 1 **(HPC is the service)**

3. Cloud has APIs for managing all infrastructure and services, **(HPC APIs are ad hoc)**

4. Cloud uses token-based authentication, **(HPC uses passwords)**

5. Cloud runs the customer's software stack, **(HPC runs the facility's SW stack)**

6. Cloud charges by the hour (encouraging paranoia), **(HPC cycles are free)**

**Cross-Pollination of Cloud & HPC Mutually Beneficial**

How can "Computing-as-a-Service" (CaaS) make this conceptual rendering a REALITY?

Customers

Modeling & Simulation

Experiments

Analytics & Data Fusion
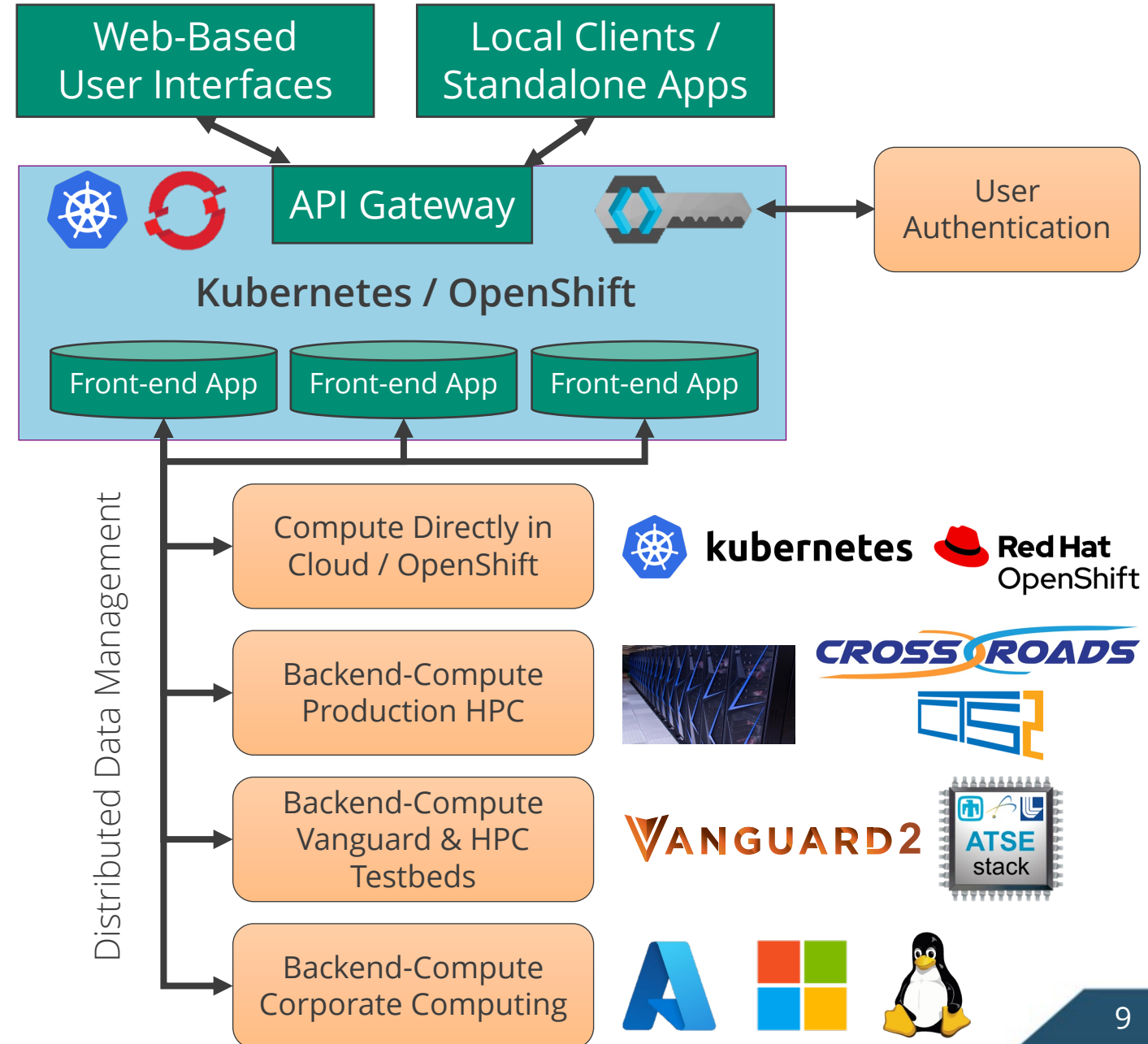
Data Storage & Retrieval

Computing Infrastructure

# Outline

- Introduction

- Computing-as-a-Service Architecture

- R&D Directions / What's Missing

- Conclusion

# Connecting the Cloud Services to Backend Compute

- Kubernetes is a distributed operating system for managing **_containerized workloads & services_**
  - Google open-sourced 2014
  - Now industry standard "Cloud OS"

- Sandia deploying production Kubernetes / OpenShift clusters

- Kubernetes not well suited for HPC
  - Bridging to HPC "on your own"
  - R&D efforts to improve Kubernetes support for HPC ( RedHat Partnership)

# DetNet Takes Leap of Faith and Teams with CaaS

DetNet PI: Chris Garasi

Goal:
- ***Provide simulation-as-a-service (SaaS) to the detonator community***
- *Speed development, reduce cost, reduce risk*
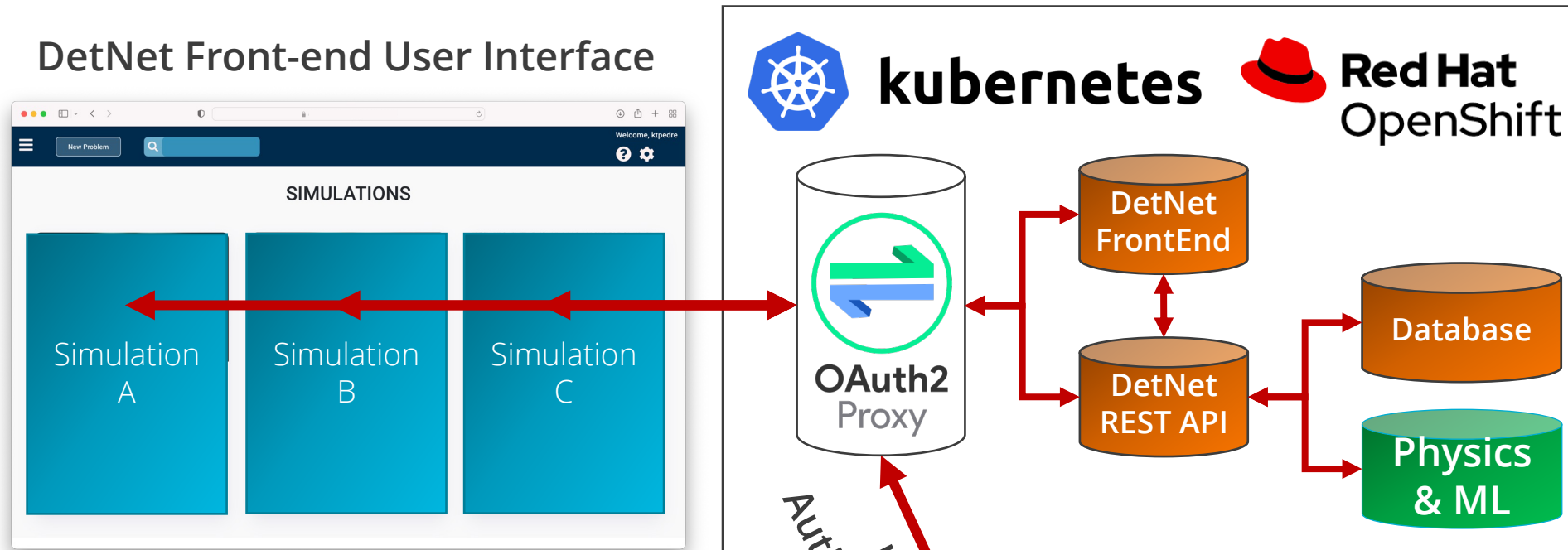
Challenges Experienced over 10+ years:
- Deploying software with antiquated input deck UI
- Software installation & upgrades
- User training – ***end-users were not HPC experts and didn't want to be***

Why now?
- ***Massively faster compute***
- ***Cloud infrastructure & containerization***

> ## Put the tools directly in the hands of the engineers
> ## (not the analysts)

# Challenge 1: Containerize and Demonstrate End-to-End Prototype



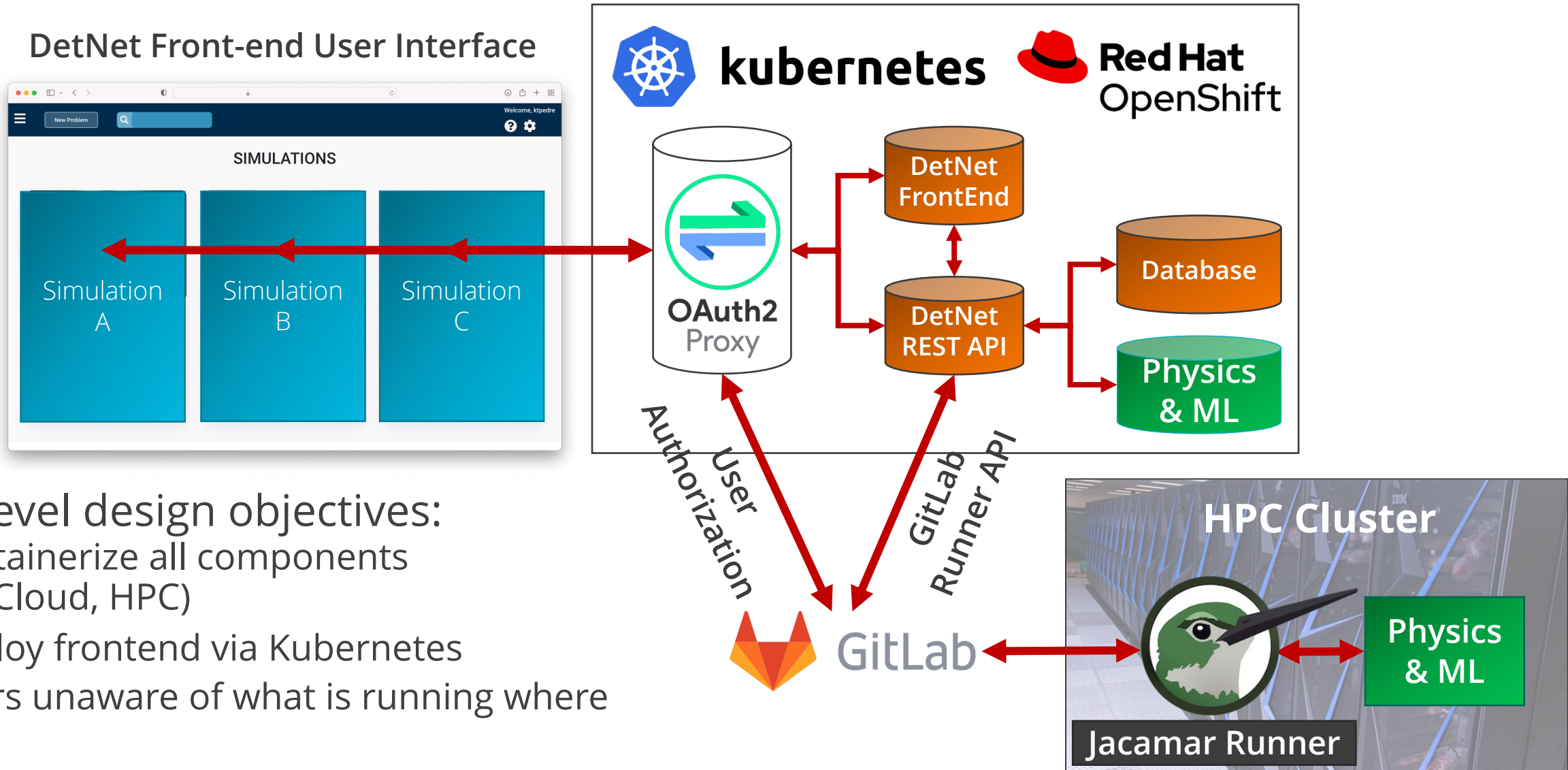**DetNet Front-end User Interface**

High-level design objectives:
- Containerize all components (UI, Cloud, HPC)
- Deploy frontend via Kubernetes
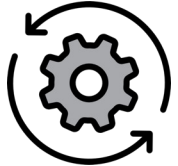- Users unaware of what is running where

# Challenge 2: Bridge to HPC with Jacamar Runners

**DetNet Front-end User Interface**



**kubernetes** | **Red Hat OpenShift**

SIMULATIONS

Simulation A | Simulation B | Simulation C

OAuth2 Proxy

DetNet FrontEnd

DetNet REST API

Database

Physics & ML

User Authorization

GitLab Runner API

GitLab

**HPC Cluster**

Physics & ML

Jacamar Runner

**High-level design objectives:**
- Containerize all components (UI, Cloud, HPC)
- Deploy frontend via Kubernetes
- Users unaware of what is running where

Jacamar @ Sandia push from Scot Swan & Allen Robinson

# DetNet is Operational, Demonstrating Key Pieces

***Designers / Engineers*** navigate web browser to DetNet front end
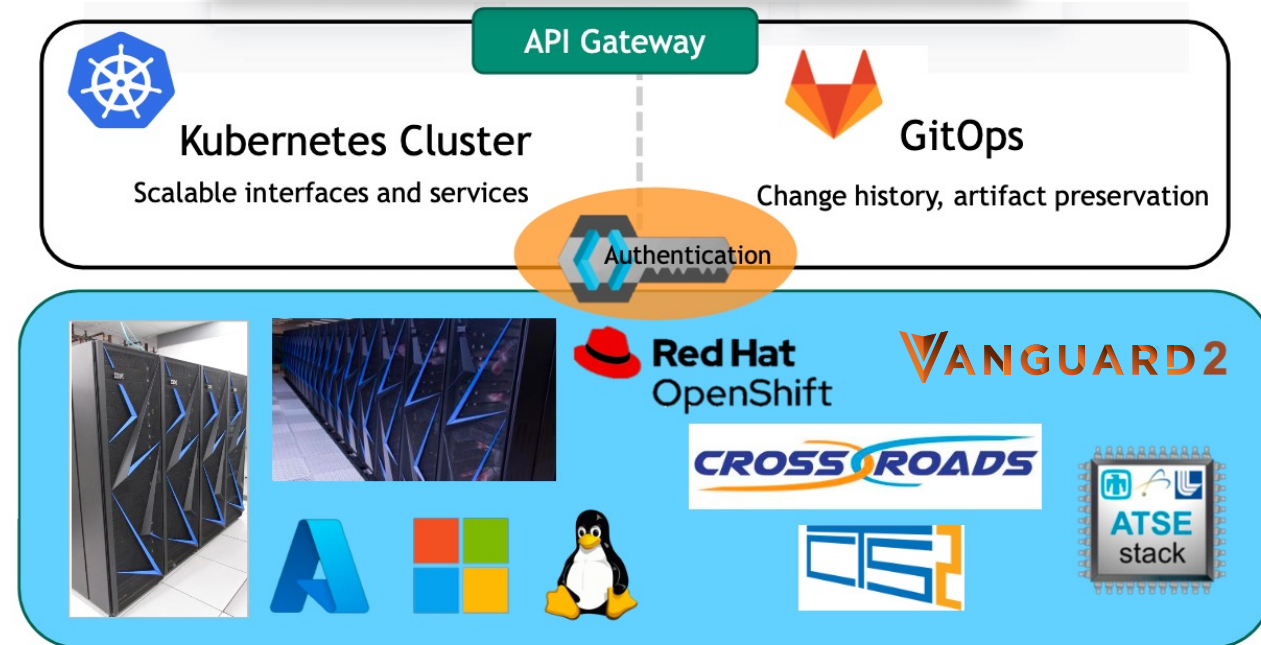
Presented with ***menu of simulations***, able to customize as needed

Computing-as-a-Service layer executes ***containerized simulation*** on appropriate computing resources

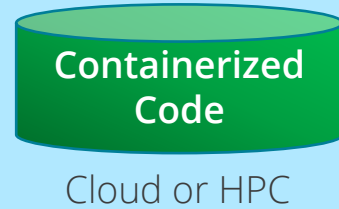Results presented ***interactively*** and stored for later retrieval & analysis

# What DetNet is Putting in the Hands of the Engineers
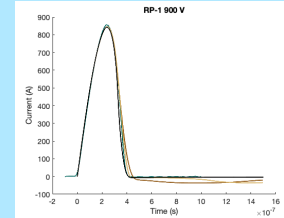
## Simple Example: 2-D Detonator Simulation

**1)** User Enters Desired Parameters

| Param A | 100 |
| Param B | 200 |
| Param C | 300 |
| Simulate | |

**2)** 2D Simulation, ~1 min turnaround
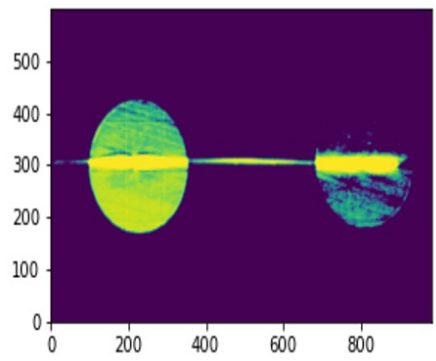
Containerized Code

Cloud or HPC

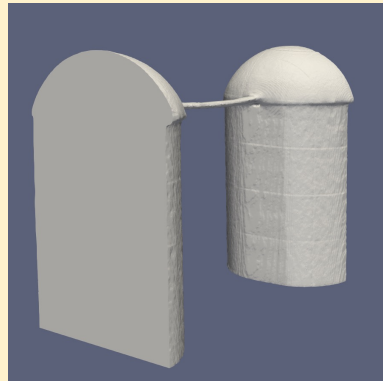**3)** Interactive Exploration of Result



**Ensemble Runs & Uncertainty Quantification**

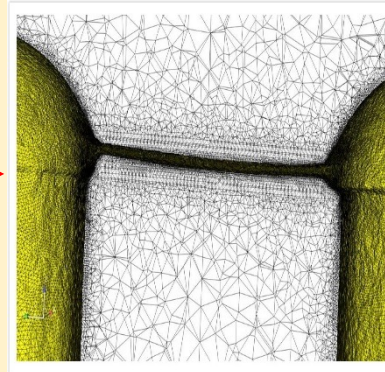## More Complex Example: Credible Automated Meshing of Images (CAMI), Surveillance Pathfinder
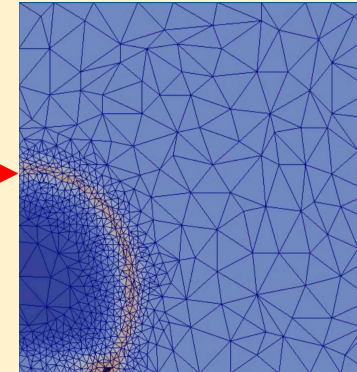
**1)** CT Scan Segmentation (ML, Python)
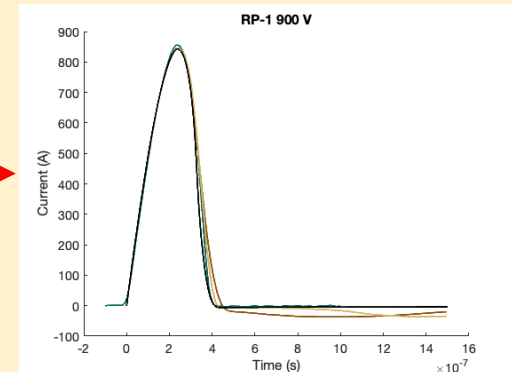
**2)** Surface Extraction & Smoothing (Python)

**3)** Surface & Volumetric Meshing (Krino & CUBIT)

**4)** Shock physics simulation (LGR)
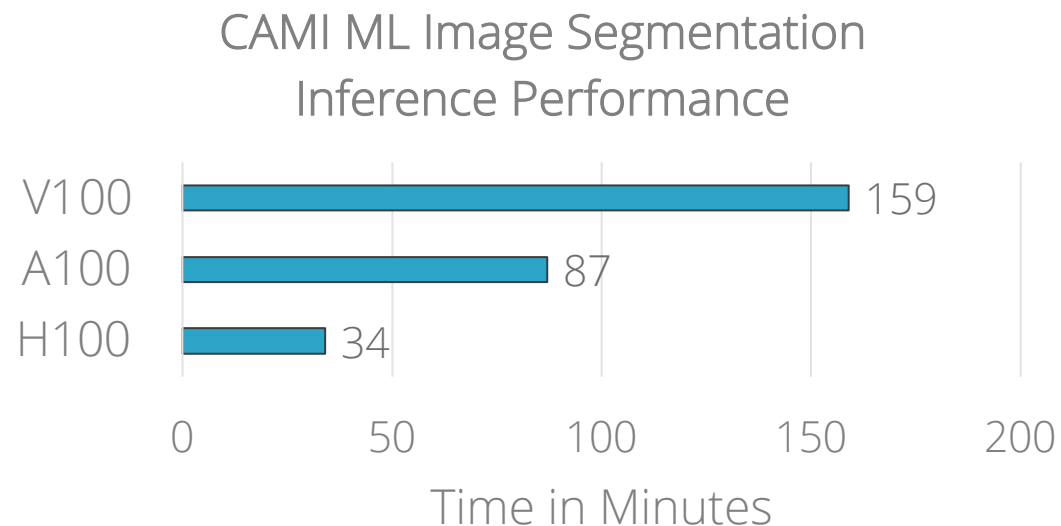
**5)** Simulated Response With UQ

# DetNet Success Driving CaaS Expansion

- Hands-on engagement with DetNet team built relationships & translated "Leap of Faith" healthy skepticism to **"This is Working!"** ☺

    - Key aspect was cross-disciplinary teaming (HPC ModSim, Web Apps, Infrastructure)

- Rapid progress & demos have attracted attention from other teams

- Adding GPU hardware to Sandia OpenShift clusters

Example:
CAMI ML-based
Image Segmentation
Requires GPUs For
"Coffee Break"
turnaround speed

### CAMI ML Image Segmentation Inference Performance

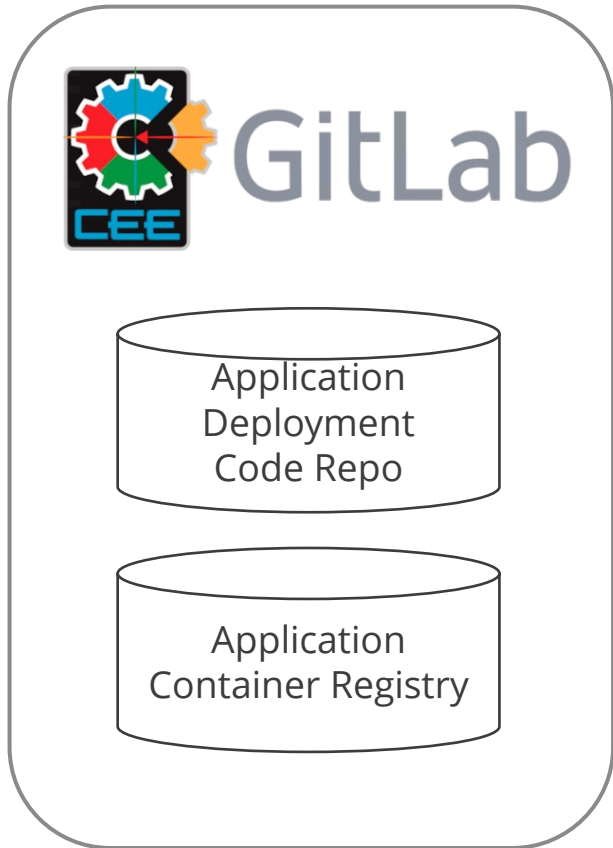| GPU | Time in Minutes |
|-----|-----------------|
| V100 | 159 |
| A100 | 87 |
| H100 | 34 |

Time in Minutes

# Outline

- Introduction

- Computing-as-a-Service Architecture

- R&D Directions / What's Missing

- Conclusion

# The Need for Automated Deployment

## Versioned Code & Containers



GitLab CEE

Application Deployment Code Repo

Application Container Registry

## 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

**HELM**

```
git clone detnet.git
# for each cluster
helm install detnet .
```

# The Need for an Intelligent Job Routing Layer

## Where's the best place to run this job?

**ADE Frontend Clients**

**?**

### Traditional HPC Platforms

SLURM & Flux APIs

K8s API

Vendor APIs

Fuzzball API

### Emerging "HPC 2.0"

FUZZBALL  https://ciq.co/

Workflows

Cloud

On-Prem

Exploring in VANGUARD

### HPC & AI Directly In Kubernetes

KUBERNETES BATCH + HPC DAY NORTH AMERICA

kubeflow/mpi-operator
multi-nic-cni-operator
Fluence (KubeFlux)

### Cloud Services & HW

AWS Lambda & Batch

Azure Functions & Batch

Google Cloud Dataflow

# The Need for Integrated Data Management Solutions

Spawned "informal" data management working group to find solutions

Sandia has deployed S3 (Simple Storage Service)

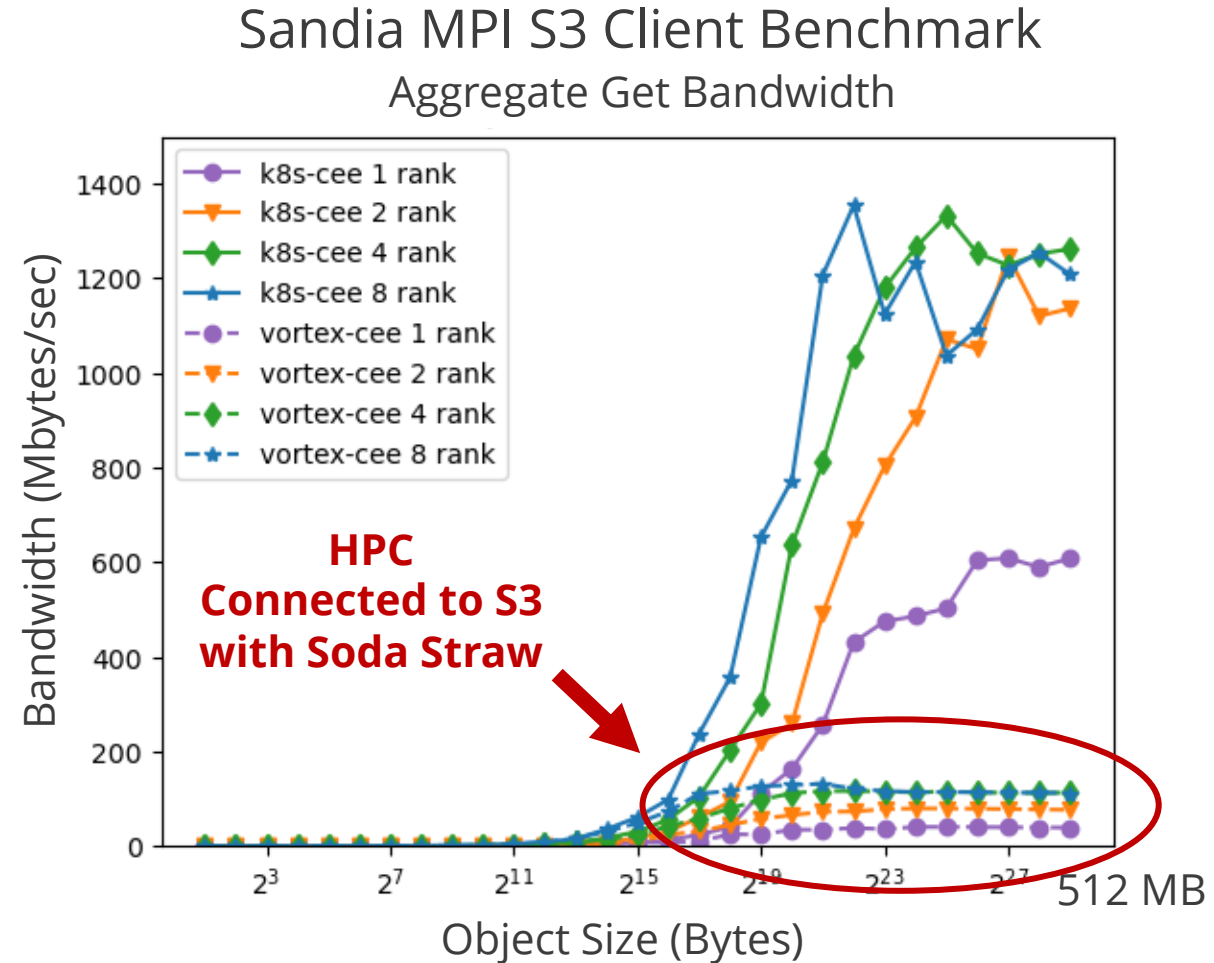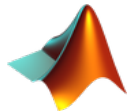Sandia data management services:

**DataSEA**
Data for Systems Engineering Applications

**SAW SDM**
Simulation Data Management

### Sandia MPI S3 Client Benchmark
Aggregate Get Bandwidth



Legend:
- k8s-cee 1 rank
- k8s-cee 2 rank
- k8s-cee 4 rank
- k8s-cee 8 rank
- vortex-cee 1 rank
- vortex-cee 2 rank
- vortex-cee 4 rank
- vortex-cee 8 rank

**HPC Connected to S3 with Soda Straw**

Y-axis: Bandwidth (Mbytes/sec)
X-axis: Object Size (Bytes), $2^3$, $2^7$, $2^{11}$, $2^{15}$, $2^{19}$, $2^{23}$, $2^{27}$, 512 MB
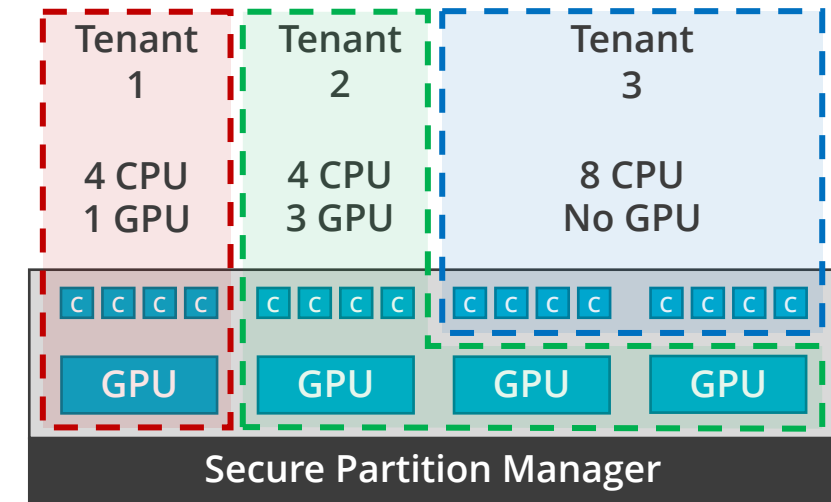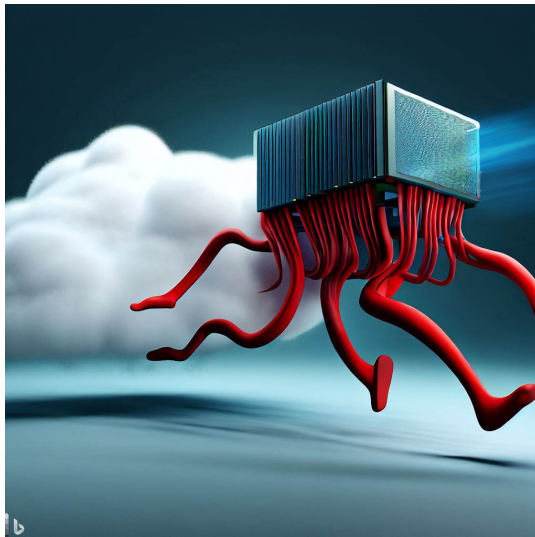
# Future Directions

- Cloud-style multi-tenancy

- Digital twins linking sensors & simulation

- ChatGPT style interfaces to ModSim tools



| Tenant 1 | Tenant 2 | Tenant 3 |
|----------|----------|----------|
| 4 CPU 1 GPU | 4 CPU 3 GPU | 8 CPU No GPU |

Secure Partition Manager

Split-up complex nodes to reduce waste & improve security



**Supercomputer sprouting legs and running from a cloud**



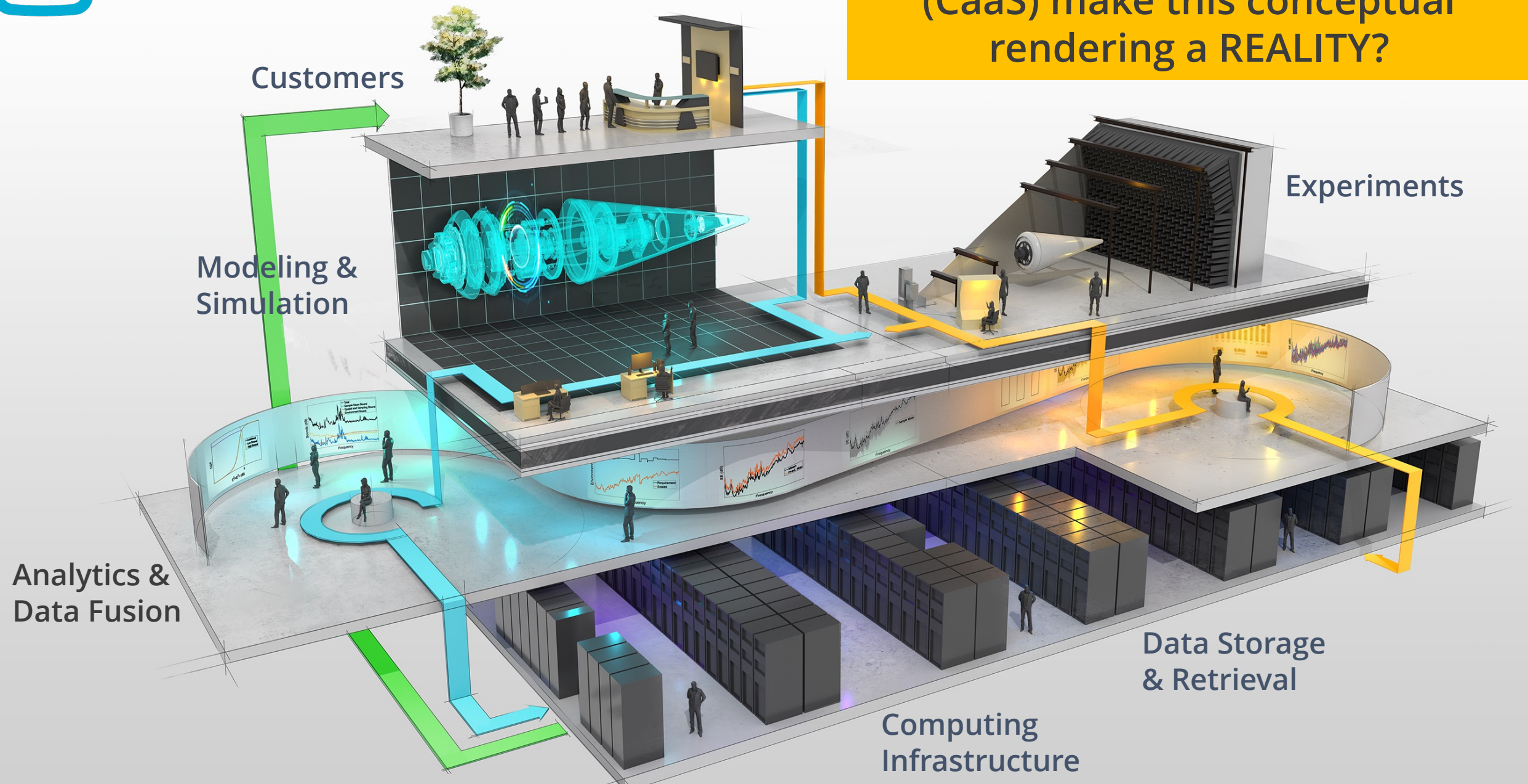**Supercomputer sprouting wings and flying in the clouds**

# Outline

- Introduction

- Computing-as-a-Service Architecture

- R&D Directions / What's Missing

- Conclusion

How can "Computing-as-a-Service" (CaaS) make this conceptual rendering a REALITY?

Customers

Modeling & Simulation

Experiments

Analytics & Data Fusion

Data Storage & Retrieval

Computing Infrastructure

# CaaS is Enabling Digital Engineering

- Delivering simulation capabilities as turnkey services

- Forming cross-disciplinary teams key to success

- Pioneered Approach with DetNet

  BEFORE: HPC Specialist Required,
  Turnaround time **Days**

  AFTER: Engineers directly access web-based detonator performance assessment tools, Containerized backend HPC & ML pipelines, Results within **1 hour**



**Thank you to the Sandia CEE, Testbeds, CapViz, ASC DevOps, Azure Stack, CSSE, and ADE teams for supporting this work.**

**Reshaping How We Deploy Codes to Users and Designing our Computing Infrastructure to Match**