

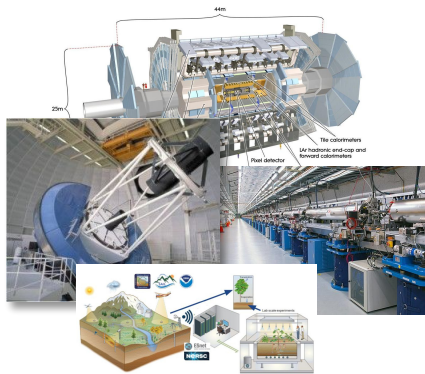
# Directions in Data at NERSC

Wahid Bhimji  
*Data Day 2024*



# Why are we here? Science

Observed data



Secrets of the

$\Omega_c$  universe

$\Omega_m$

$\sigma_8$

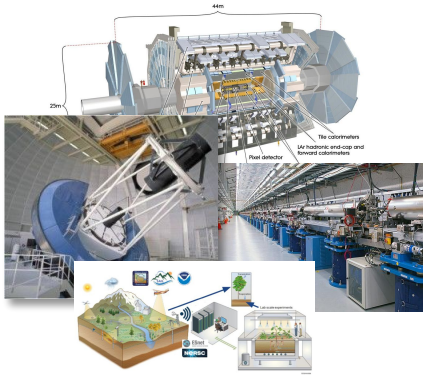
$\theta_{12}$

$m_H \dots$

Data Acquisition; Availability; Access; Analysis; AI; Archiving

# Science data presents challenges

Observed data



Secrets of the universe

$\Omega_c$

$\Omega_m$

$\sigma_8$

$\theta_{12}$

$m_H \dots$

Data Acquisition; Availability; Access; Analysis; AI; Archiving

SCIENCE  
LOST



Inefficient Filtering  
Storage capacity  
limitations  
Bad data quality  
Human steering

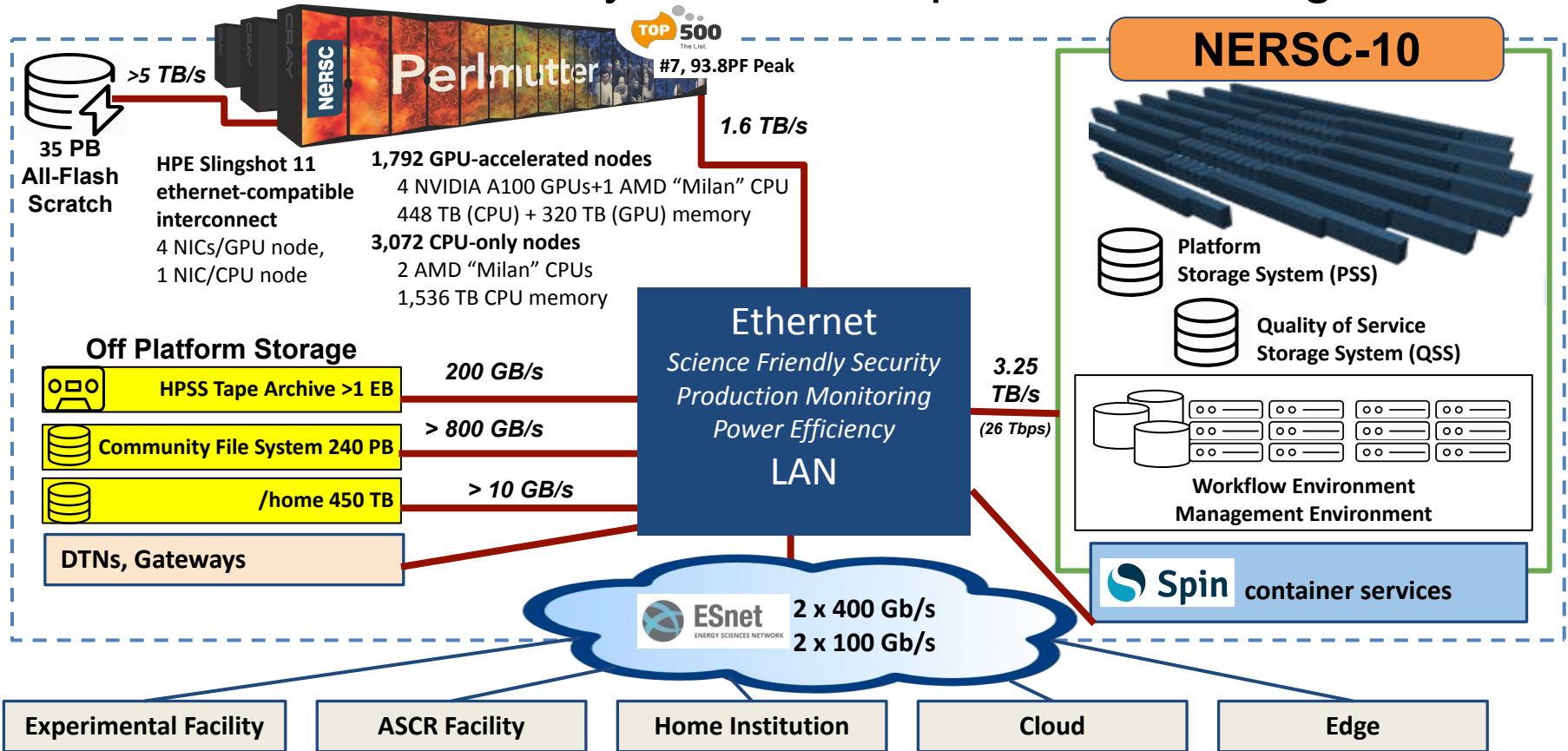
Inadequate  
data transfer  
Poor data  
management  
tools

Limited user  
interfaces  
Unproductive  
software  
I/O Bottlenecks

Insensitive  
analysis  
techniques

Irreproducible  
science  
Poor data  
preservation

# NERSC has rich data systems to help these challenges



# NERSC also has rich data software



globus online



jupyter



Data Transfer and Access



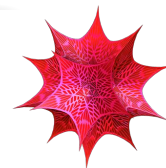
mongoDB®

MySQL™

Data Management



julia



Data Analytics



PyTorch



AI / Machine Learning



Visualization



SHIFTER



Spin

Containers/Cloud-tech



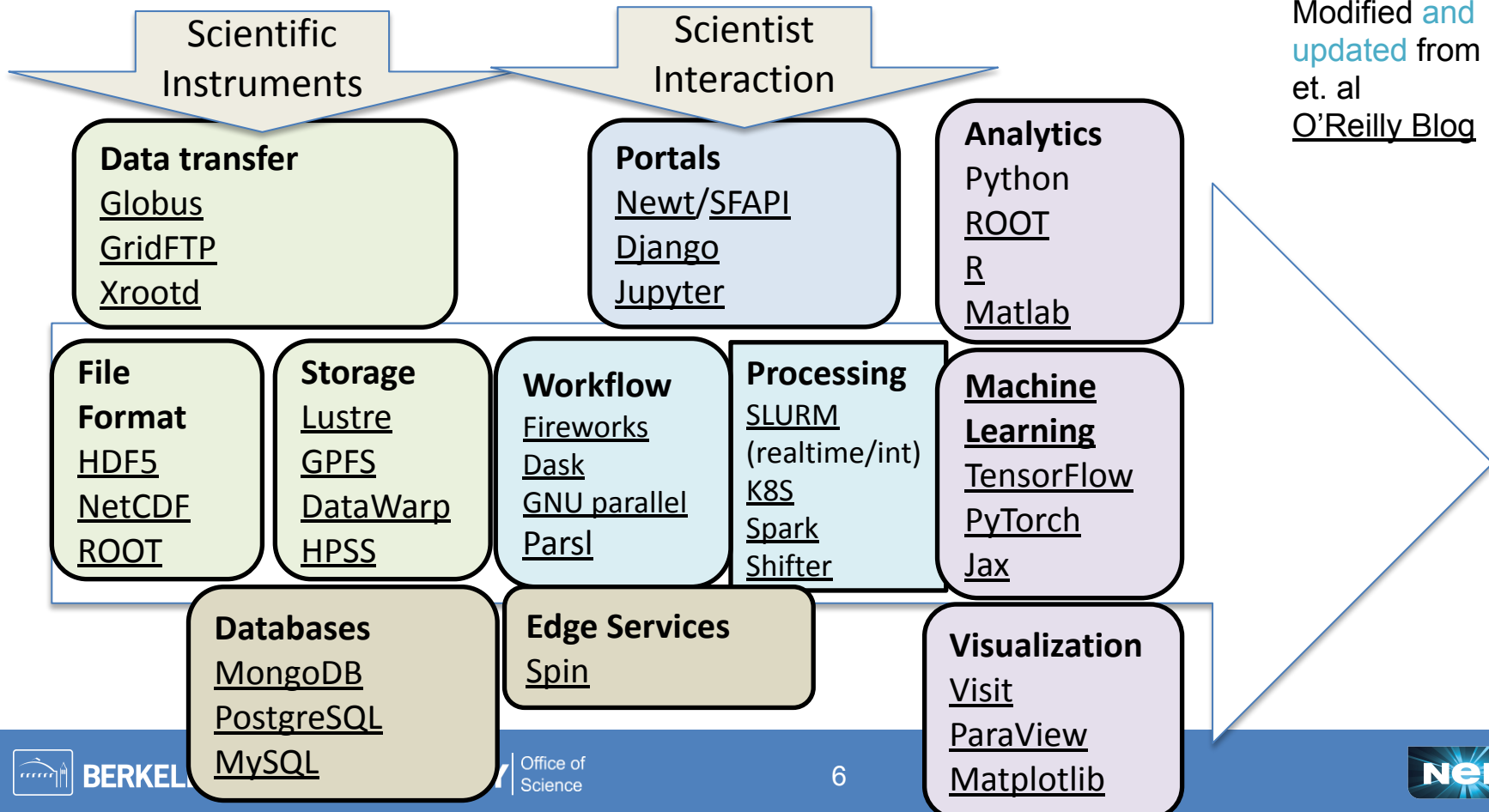
GNUparallel



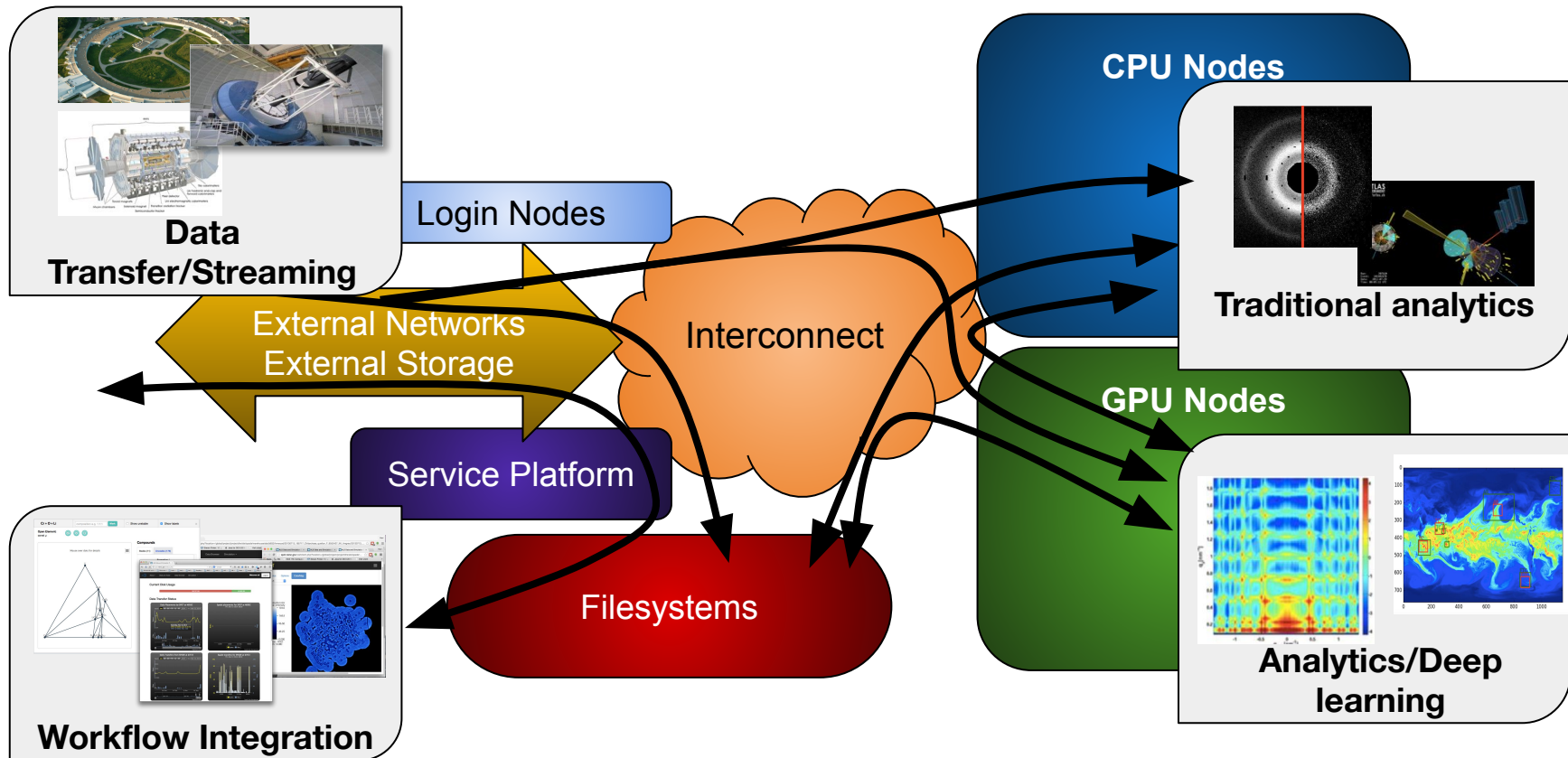
Workflows



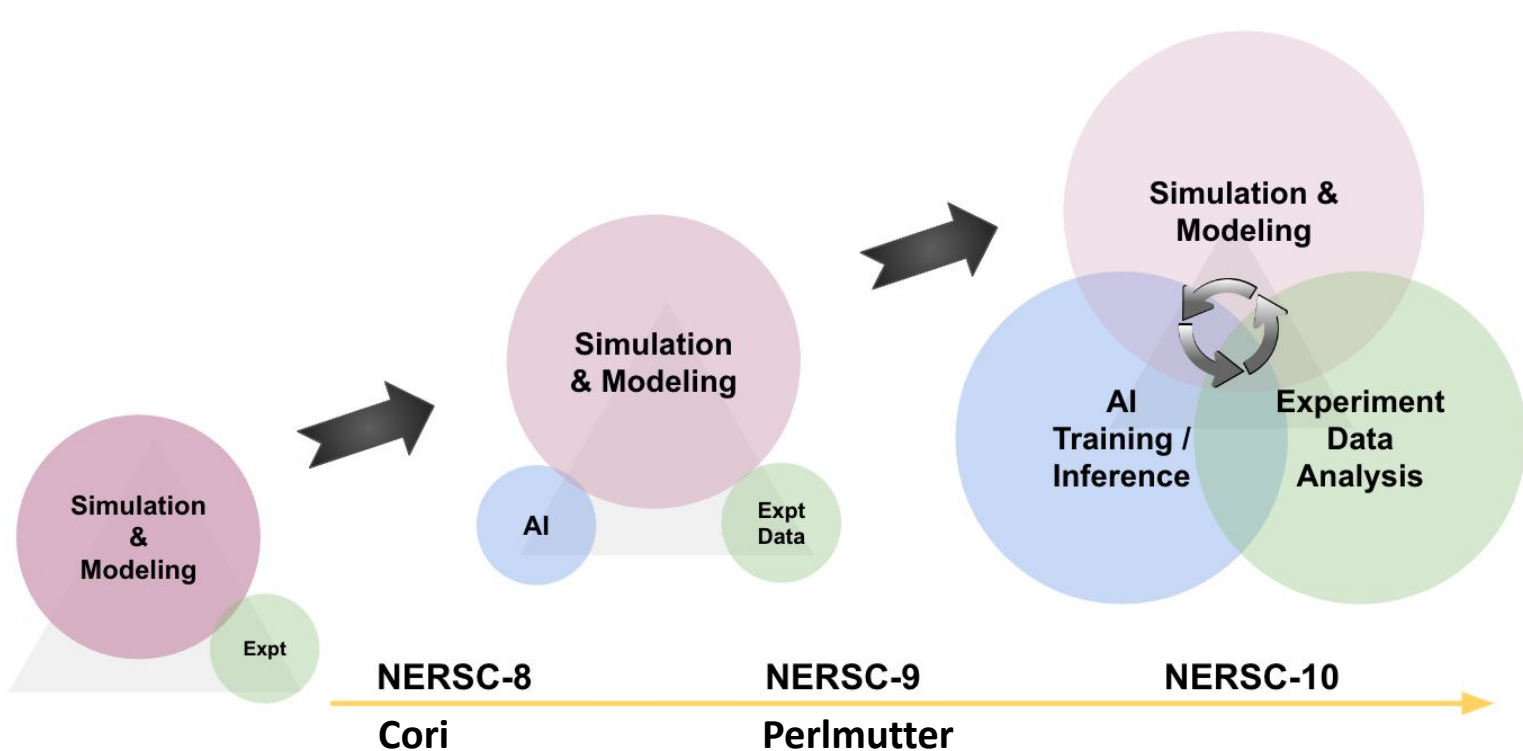
# NERSC provides a rich data ecosystem



# Data ecosystem impacts entire workflow



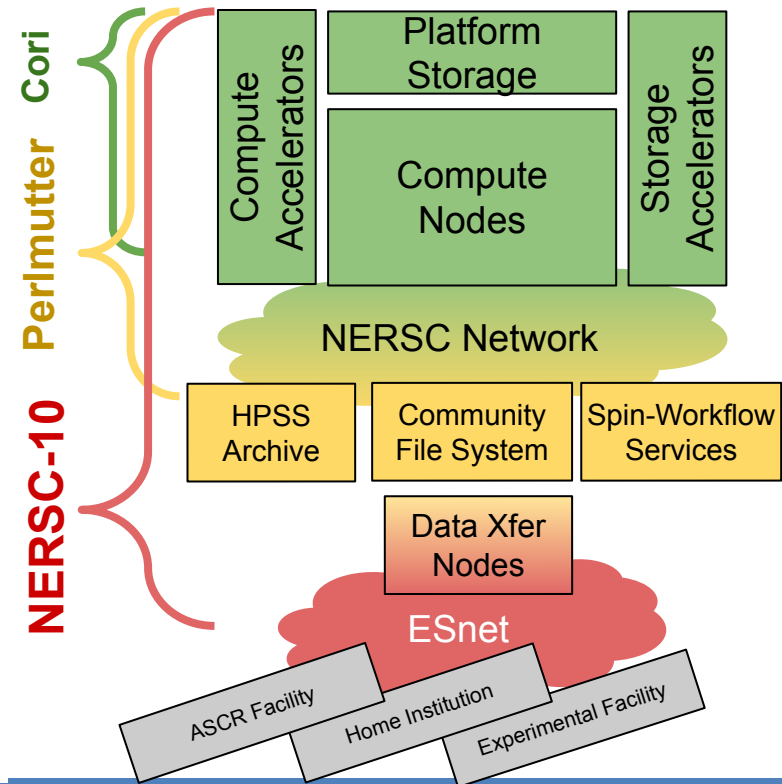
# Data Directions: NERSC-10 - integrated scientific workflows





# Data Directions: NERSC-10 workflows extend beyond system

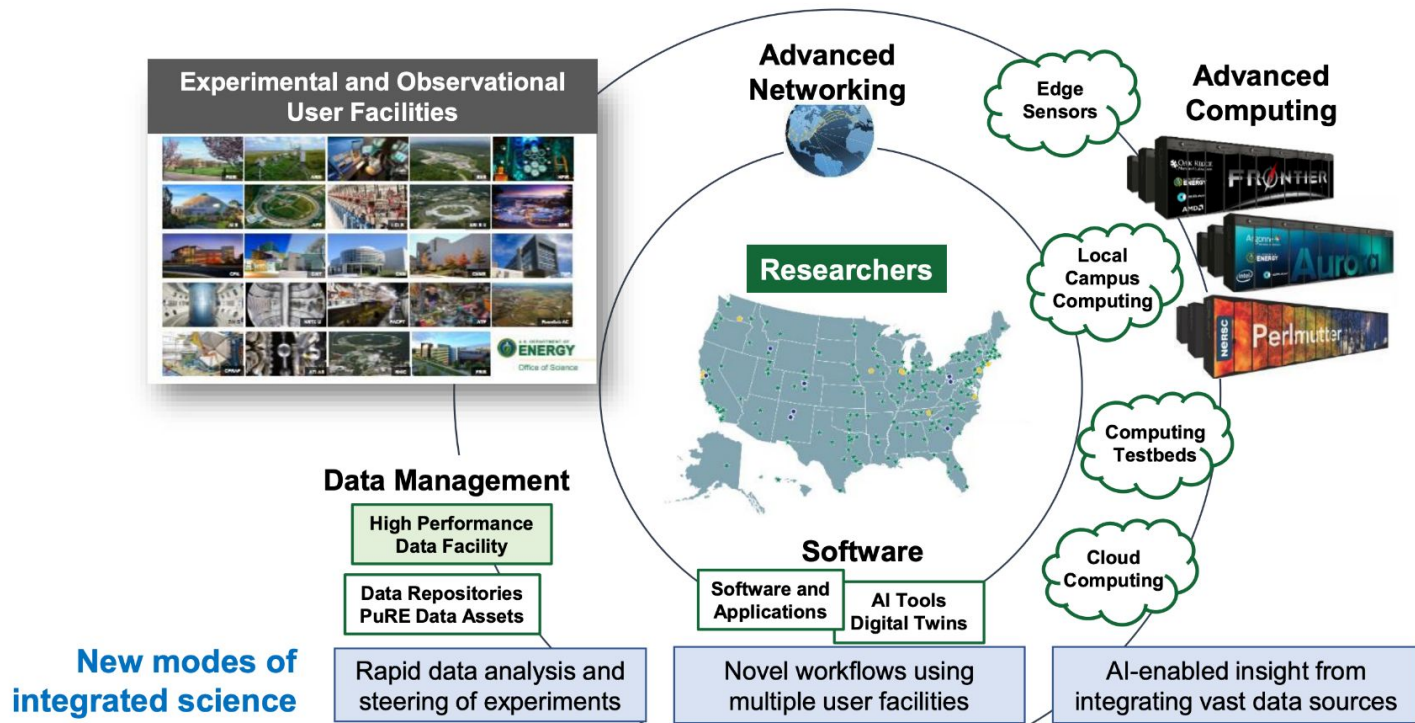
- **Quality of Service** – computation, storage and networking designed to emphasize response-time plus throughput/utilization.
- **Seamlessness** – tight integration of system components to enable high performance across workflow steps.
- **Portability** – Modular workflow execution across heterogeneous HPC, edge and cloud.
- **Programmability** – APIs to manage data, execute distributed code, and interact with system resources.
- **Orchestration** – coordinate resource management across different resource domains.
- **Security** – authentication, authorization and auditing (e.g., identify proofing, access/privacy control, records of transactions).



# Data Directions: IRI workflows extend beyond system

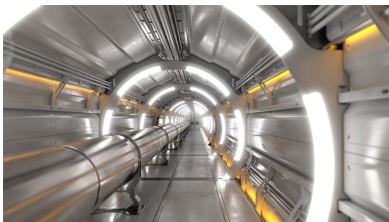
DOE's Integrated Research Infrastructure (IRI) Vision:

*To empower researchers to meld DOE's world-class research tools, infrastructure, and user facilities seamlessly and securely in novel ways to radically accelerate discovery and innovation*

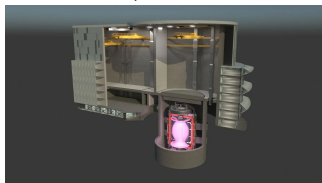


# Data Directions: Saving the science

Future circular collider, CERN

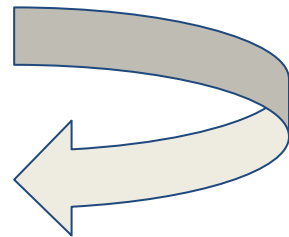


Commercial fusion reactor, UKAEA



Secrets of the universe

42+/- 6



## Data Acquisition; Availability; Access; Analysis; AI; Archiving

SCIENCE  
ENABLED ↑

Extremely low latency data transfer and access and/or compute workflow across edge and HPC

Compose reusable, resilient automated services and workflows

Granular, automatically discoverable, data

Experiment with seamless, smart, human computing interfaces

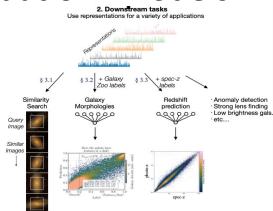
Leverage reusable, fast AI-inference with rich science models and uncertainties

Curate with re-interpretation and smart archiving

# Data Directions at NERSC - example projects

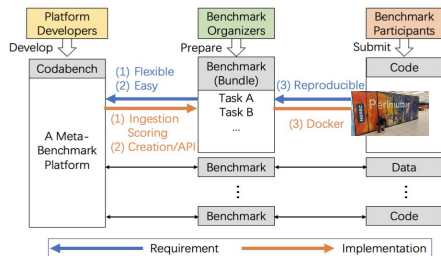
## Pervasive AI Science

Large multi-purpose, "foundation" models



[Sky survey self-supervised learning](#)

Platforms for scientific AI



["Fair Universe"](#)

## Converged HPC+Cloud

Slurm and k8s orchestration

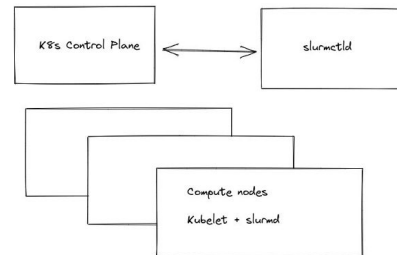


"Containers Everywhere"

Towards a "Containers Everywhere" HPC Platform

4:30pm - 5pm

Shane Canon, Laurie Stephey, Daniel Fulton, and Brandon Cook (L)



["Adjacent" perspective from SC22 SchedMD talk](#)

## FAIR Data Services

[PI Toolbox](#)

## Scalable analytics

[Project dragon](#)

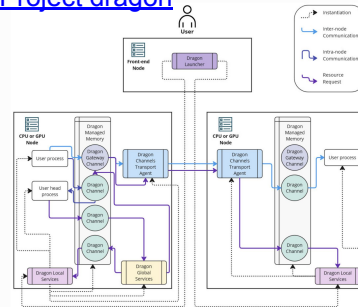


Fig. 1 Dragon Infrastructure Architecture

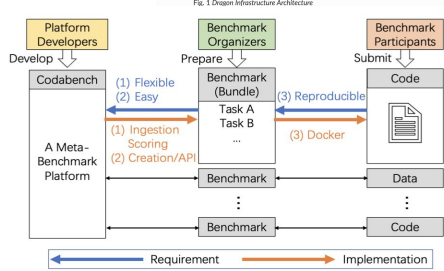
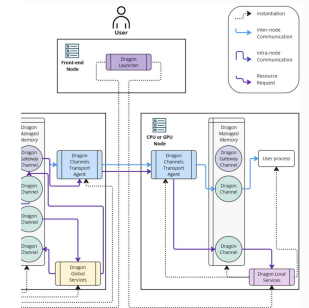
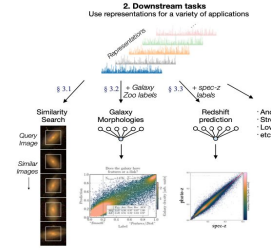
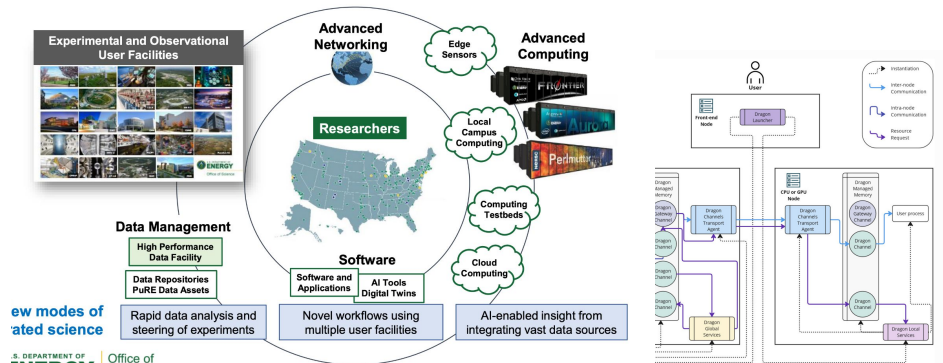
# Directions in data services

**Compose** services and compute seamlessly across distributed facilities

**Experiment** with and apply performant, productive analytics at scale

**Leverage** large AI models, fine-tune to new problems, apply to new data pipelines

**Curate** and re-use data through FAIR management services



# Directions in data services

**Compose** services and compute seamlessly across distributed facilities

**Experiment** with and apply performant, productive analytics at scale

**Leverage** large AI models, fine-tune to new problems, apply to new data pipelines

**Curate** and re-use data through FAIR management services

## Day 1

Time (PDT)	Topic	Presenters
9:00 - 9:10 a.m.	Introduction/Welcome	
9:10 - 9:30 a.m.	Directions in Data at NERSC	<a href="#">Wahid Bhimji</a>
9:30 - 10:00 a.m.	HPC Containers, Podman-HPC, Shifter	<a href="#">Adam Lavelly</a>
10:00 - 10:20 a.m.	Science Applications: Checkpoint/Restart in Containers	<a href="#">Madan Timalisina</a>
10:20 - 10:30 a.m.	Break	
10:30 - 11:00 a.m.	Choosing the right storage for your data	<a href="#">Steve Leak</a> , <a href="#">Ravi Cheema</a>
11:00 - 11:30 a.m.	Distributed Python at NERSC	<a href="#">Daniel Margala</a>
11:30 - 12:00 a.m.	Julia	<a href="#">Johannes Blaschke</a>
12:00 - 1:30 p.m.	Lunch	
1:30 - 2:00 p.m.	Deep Learning at Scale	<a href="#">Shashank Subramanian</a>
2:00 - 2:30 p.m.	Science Applications in AI	<a href="#">Jared Willard</a>
2:30 - 3:00 p.m.	Nvidia Triton Demo: Incorporating AI inference into workflows	<a href="#">Andrew Naylor</a>

## Day 2

Time (PDT)	Topic	Presenters
9:00 - 9:30 a.m.	Integrated Research Infrastructure and N10	<a href="#">Debbie Bard</a>
9:30 - 10:00 a.m.	Superfacility API	<a href="#">Bjoern Enders</a>
10:00 - 10:30 a.m.	Science Applications: Data Streaming to NCEM	<a href="#">Sam Welborn</a> , <a href="#">Peter Ercius</a>
10:30 - 11:00 a.m.	Science Applications: EJFAT/JIRIAF	Vardan Gyuriyan, Jeng-Yuan Tsai
11:00 - 11:30 a.m.	Science Applications: JAWS	<a href="#">Daniela Cassol</a>
11:30 - 12 p.m.	Data Transfers and the Globus	<a href="#">Nick Tyler</a>





# Questions? Collaboration?

Wahid Bhimji

wbhimji@lbl.gov

<https://docs.nersc.gov/analytics/analytics/>

<https://docs.nersc.gov/machinelearning/>

<https://docs.nersc.gov/services/>

And questions for you:

What are you hoping to get out of this data day?

What would you have liked to have on the agenda that you don't see?



# Questions

What are you hoping to get out of this data day?

What would you have liked to have on the agenda that you don't see?