


# High Performance Computing for AI/ML

Corey Adams  
Argonne National Laboratory

Authorship of this talk should be attributed to all of ALCF, but particularly to  
Mike Papka, Venkat Vishwanath, and Kyle Felker

November 01, 2022



MAN ACHIEVED HERE  
THE FIRST SELF-SUSTAINING CHAIN REACTION  
AND THEREBY INITIATED THE  
CONTROLLED RELEASE OF NUCLEAR ENERGY

## Argonne National Laboratory

The U.S. Department of Energy's Argonne National Laboratory delivers world-class research, technologies, and new knowledge that aim to make an impact — from the atomic to the human to the global scale.

# About Argonne

**Argonne is a multidisciplinary science and engineering research center located outside Chicago.**

- Born out of the University of Chicago's work on the Manhattan Project in the 1940s.
- **Managed by UChicago Argonne, LLC, for the U.S. Department of Energy's Office of Science.**
- Works with universities, industry, and other national labs on questions and experiments too large for any one institution to do by itself.

# Our one-of-a-kind facilities enable science from the nanoscale to the exascale

Argonne's five flagship facilities support one of the largest user communities in the U.S. Department of Energy complex.



**Advanced  
Photon Source**



**Argonne  
Tandem Linear  
Accelerator  
System**



**Argonne  
Leadership  
Computing  
Facility**



**Center for  
Nanoscale  
Materials**



**Atmospheric  
Radiation  
Measurement – The  
Southern Great  
Plains**

# DOE SC Advanced Scientific Computing Research User Facilities

The Advanced Scientific Computing Research (ASCR) program leads the nation and the world in supercomputing, high-end computational science, and advanced networking for science.

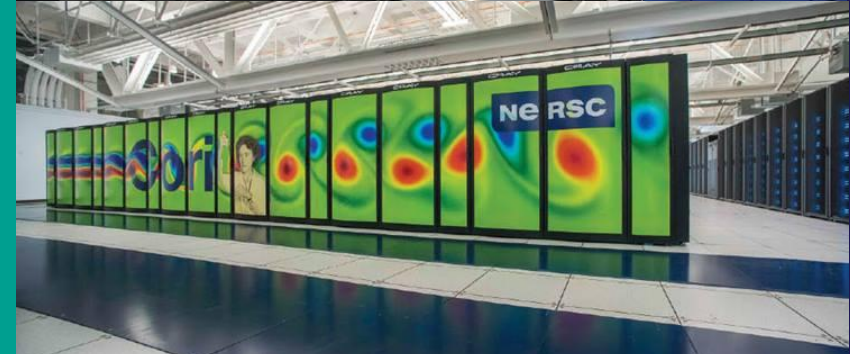
**ALCF and OLCF make up the  
DOE Leadership Computing Facility**

Argonne  
Leadership  
Computing  
Facility  
(ALCF)

Oak Ridge  
Leadership  
Computing  
Facility  
(OLCF)

National Energy  
Research Scientific  
Computing Center  
(NERSC)

Energy Sciences  
Network (ESnet)



# DOE Leadership Computing Facility

- Established in 2004 as a collaborative, multi-lab initiative funded by DOE's *Advanced Scientific Computing Research* program
- Operates as **one facility** with two centers, at Argonne and at Oak Ridge National Laboratory
- Deploys and operates at least two advanced architectures that are **10-100 times more powerful** than systems typically available for open scientific research
- **Fully dedicated** to open science to address the ever-growing needs of the scientific community

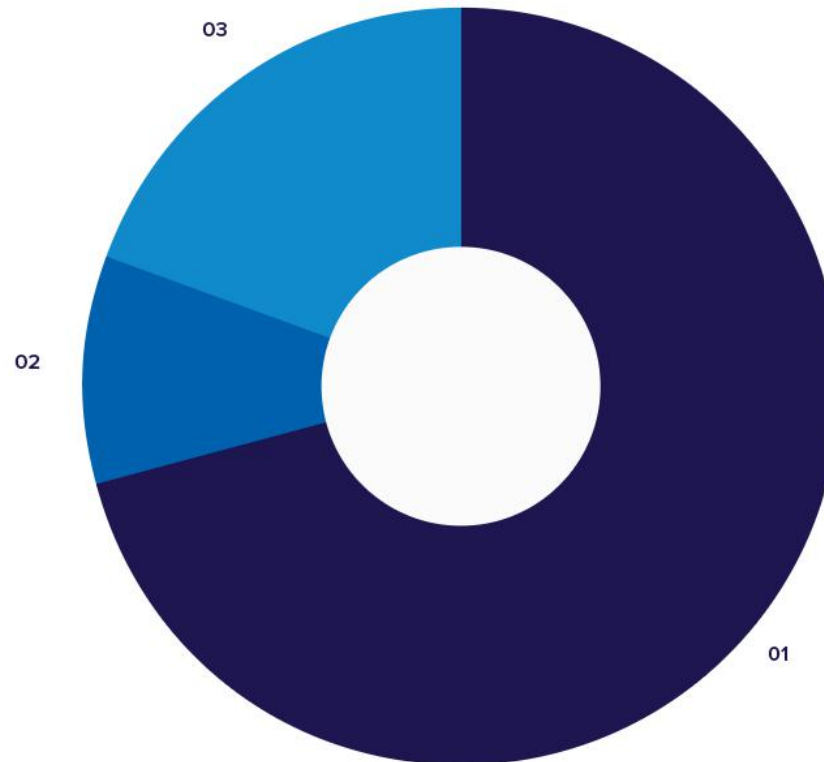


# ALCF Staff

To ensure facility users are able to get the most out of its supercomputers, the ALCF has assembled an exceptional team of:

- HPC system and network administrators
- computational scientists,
- computer scientists
- data scientists
- performance engineers
- visualization experts
- software developers
- user support staff

ALCF STAFF NUMBERS



01 Staff Members

113

02 Postdoctoral Researchers

16

03 Summer Students

33

# ALCF at a Glance in 2021

- Users pursue scientific challenges
- In-house experts to help maximize results
- Resources fully dedicated to open science

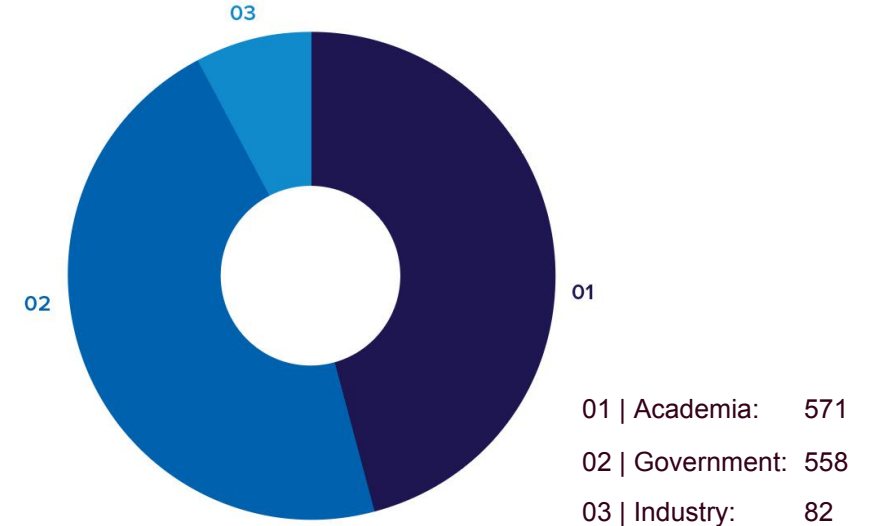
**33.5M** node-hours of compute time

**375** active projects

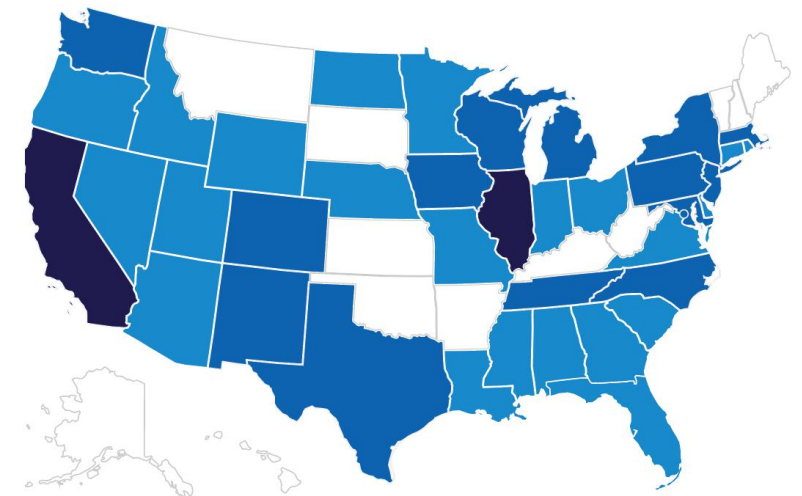
**1,168** facility users

**230+** publications

2021 ALCF Users by Affiliation

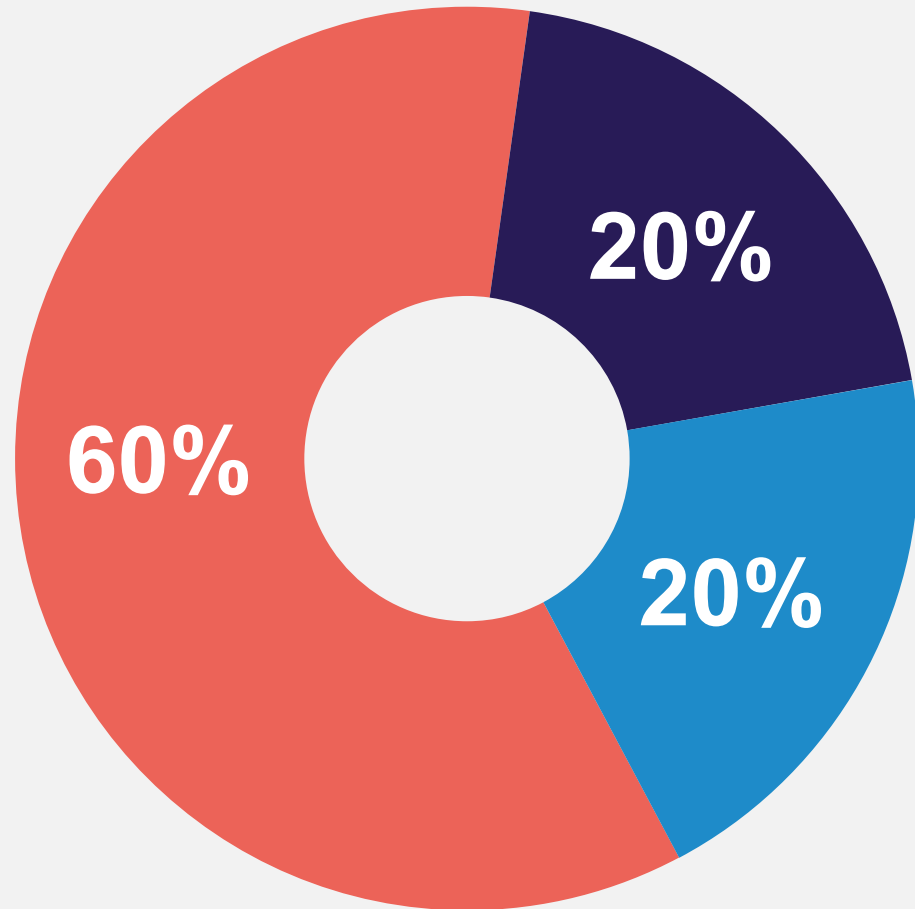


2021 U.S. ALCF Users by State





# ALCF Allocation Programs



## **INCITE: Innovative and Novel Computational Impact on Theory and Experiment**

- § Yearly call with computational readiness and peer reviews
- § Open to all domains and user communities

---

## **ALCC: ASCR Leadership Computing Challenge**

- § Yearly call with peer reviews
- § Focused on DOE priority

---

## **DD: Director's Discretionary Program**

- § Rapid allocations for project prep and immediate needs
  - Early Science Program (ESP)
  - Exascale Computing Project (ECP)
  - ALCF Data Science Program (ADSP)
  - Proprietary Projects

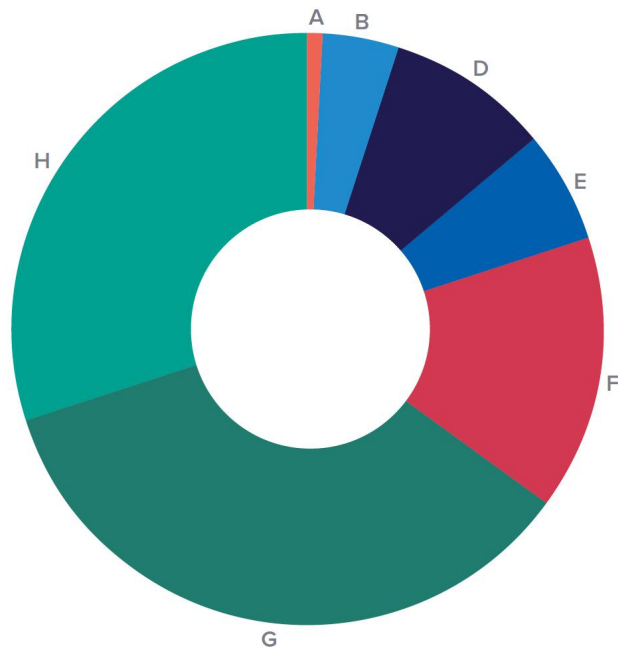
# Accessing ALCF Resources for Science

As a national user facility dedicated to open science, any researcher in the world with a large-scale computing problem can apply for time on ALCF computing resources.

## 2021 INCITE

**17.8M** NODE HOURS

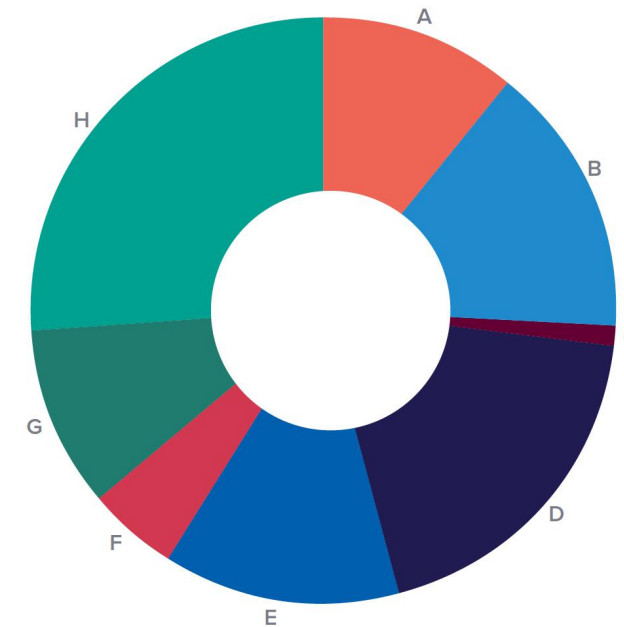
<b>A</b> Biological Sciences	1 %
<b>B</b> Chemistry	4
<b>C</b> Computer Science	–
<b>D</b> Earth Science	9
<b>E</b> Energy Technologies	6
<b>F</b> Engineering	15
<b>G</b> Materials Science	35
<b>H</b> Physics	30



## 2021 ALCC

**7.3M** NODE HOURS

<b>A</b> Biological Sciences	11 %
<b>B</b> Chemistry	15
<b>C</b> Computer Science	1
<b>D</b> Earth Science	19
<b>E</b> Energy Technologies	13
<b>F</b> Engineering	5
<b>G</b> Materials Science	10
<b>H</b> Physics	26



*ALCC data are from calendar year 2021.*



# Computing Resources

FEATURE	POLARIS	THETA: KNL NODES	THETA: GPU NODES	COOLEY
Purpose	Exascale Testbed	Production Supercomputer	Production Supercomputer	Data Analysis and Visualization Cluster
Architecture	HPE Apollo 6500 Gen10+	Intel-Cray XC40	NVIDIA DGX A100	Intel Haswell
Peak Performance	44 PF (double precision)	11.7 PF	3.9 PF	293 TF
Processors per Node	3rd Gen AMD EPYC	64-core, 1.3-GHz Intel Xeon Phi 7230	2 AMD EPYC 7742	2 6-core, 2.4-GHz Intel E5-2620
GPU per Node	4 NVIDIA A100 Tensor Core	—	8 NVIDIA A100 Tensor Core	NVIDIA Tesla K80
Nodes	560	4,392	24	126
Cores	560	281,088	576	1,512
Memory	280 TB (DDR4); 87.5 TB (HBM)	843 TB (DDR4); 70 TB (HBM)	24 TB (DDR4); 7.7 TB (HBM)	47 TB (DDR4); 3 TB (GDDR5)

## JLSE Experimental Testbeds

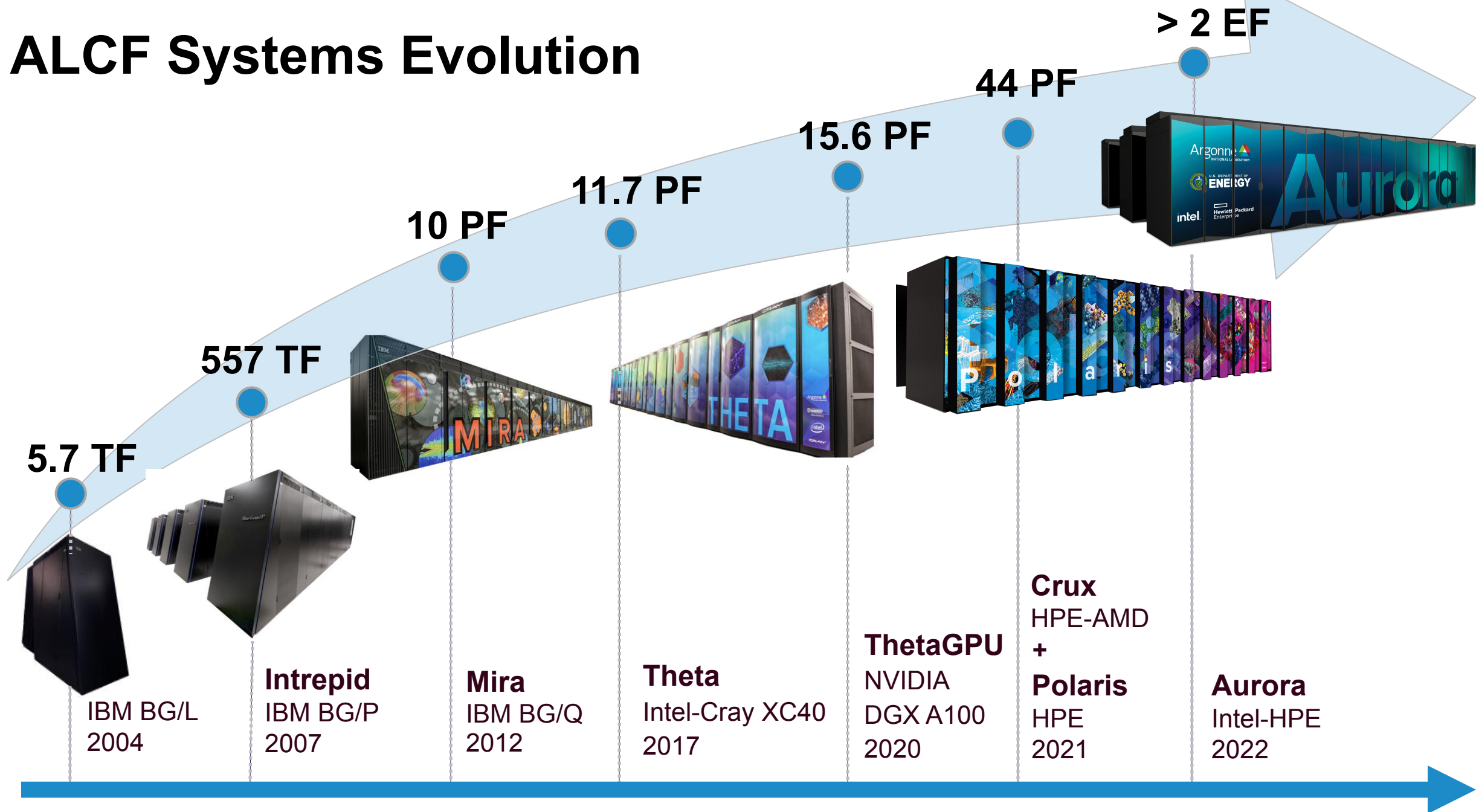
- 150 nodes
- Intel/AMD/IBM/Marvell/GPGPU
- EDR/100GbE/OPA
- Lustre/GPFS/DAOS

## Grand and Eagle (Storage)

Each system has:

- HPE ClusterStor E1000
- 100 petabytes of usable capacity
- 8,480 disk drives
- Lustre filesystem
  - § 160 Object Storage Targets
  - § 40 Metadata Targets
- HDR InfiniBand network
- 650 GB/s rate on data transfers

# ALCF Systems Evolution



# Aurora

Argonne's upcoming exascale supercomputer will leverage several technological innovations to support machine learning and data science workloads alongside traditional modeling and simulation runs.

SUSTAINED PERFORMANCE

**≥2 Exaflop DP**

X<sup>e</sup> ARCHITECTURE-BASED GPU

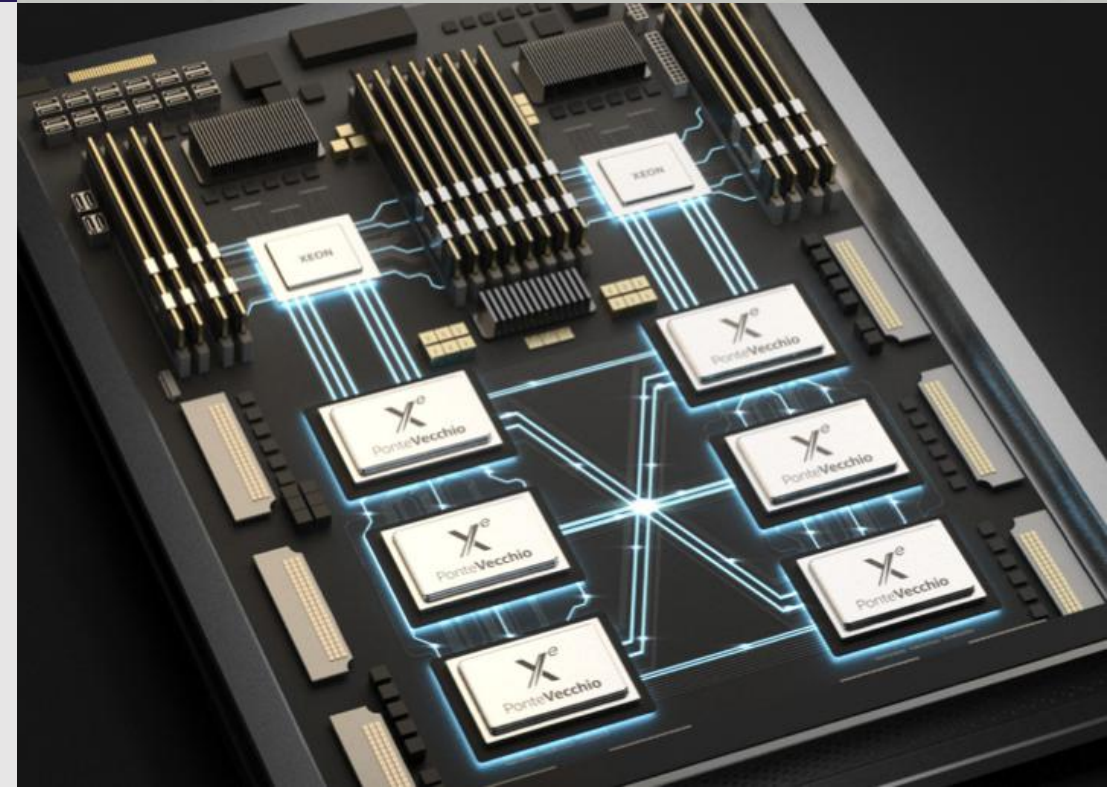
**Ponte Vecchio**

INTEL XEON SCALABLE PROCESSOR

**Sapphire Rapids**

PLATFORM

**HPE Cray EX**



## Compute Node

2 Intel Xeon scalable “Sapphire Rapids” processors; 6 X<sup>e</sup> arch-based GPUs; Unified Memory Architecture; 8 fabric endpoints; RAMBO

## GPU Architecture

X<sup>e</sup> arch-based “Ponte Vecchio” GPU; Tile-based chiplets, HBM stack, Foveros 3D integration, 7nm

## CPU-GPU Interconnect

CPU-GPU: PCIe  
GPU-GPU: X<sup>e</sup> Link

## System Interconnect

HPE Slingshot 11; Dragonfly topology with adaptive routing

## Network Switch

25.6 Tb/s per switch, from 64–200 Gbs ports (25 GB/s per direction)

## High-Performance Storage

≥230 PB, ≥25 TB/s (DAOS)

## Programming Models

Intel oneAPI, MPI, OpenMP, C/C++, Fortran, SYCL/DPC++

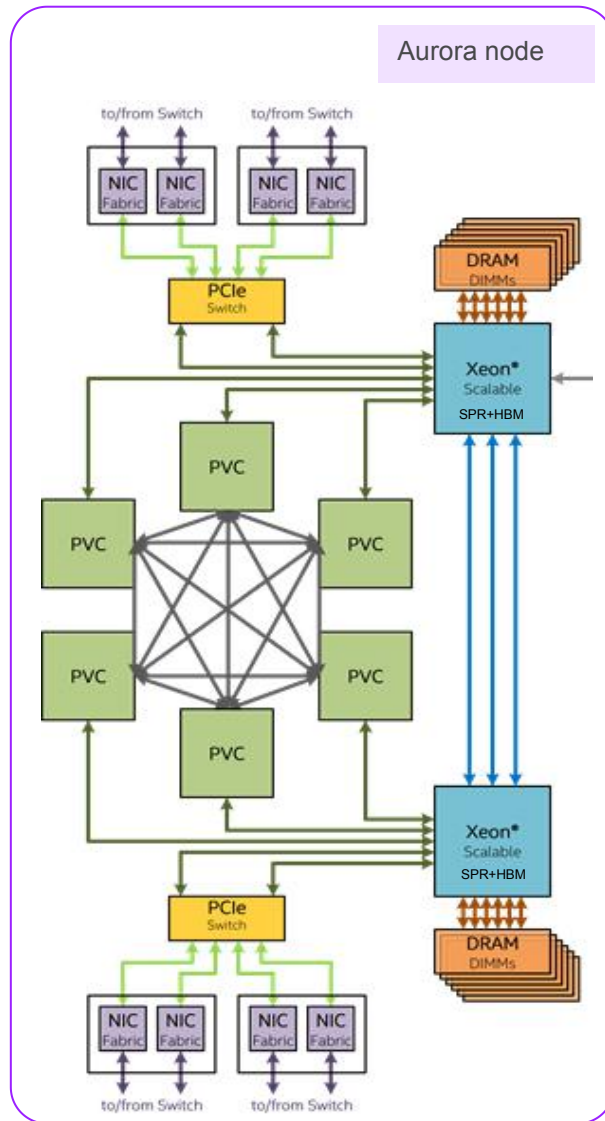
## Node Performance

>130 TF

## System Size

>9,000 nodes

# Aurora Compute Node



- 6 Xe Architecture based GPUs (Ponte Vecchio)
  - All to all connection
- 2 Intel Xeon (Sapphire Rapids) processors
- Unified Memory Architecture across CPUs and GPUs
- 8 Slingshot Fabric endpoints





Argonne's Aurora System > 60,000 Intel GPUs Science Starts in 2023

# Aurora Cabinets Installed at Argonne





# Data Science and Learning on Aurora

Aurora will provide for a familiar, productive and performant HPC and AI software stack

Intel AI Analytics Toolkit

## § Python Ecosystem

§ Numba, NumPy, etc.

## § Deep Learning Frameworks:

§ PyTorch, TensorFlow, Horovod, DDP,

## § Machine Learning

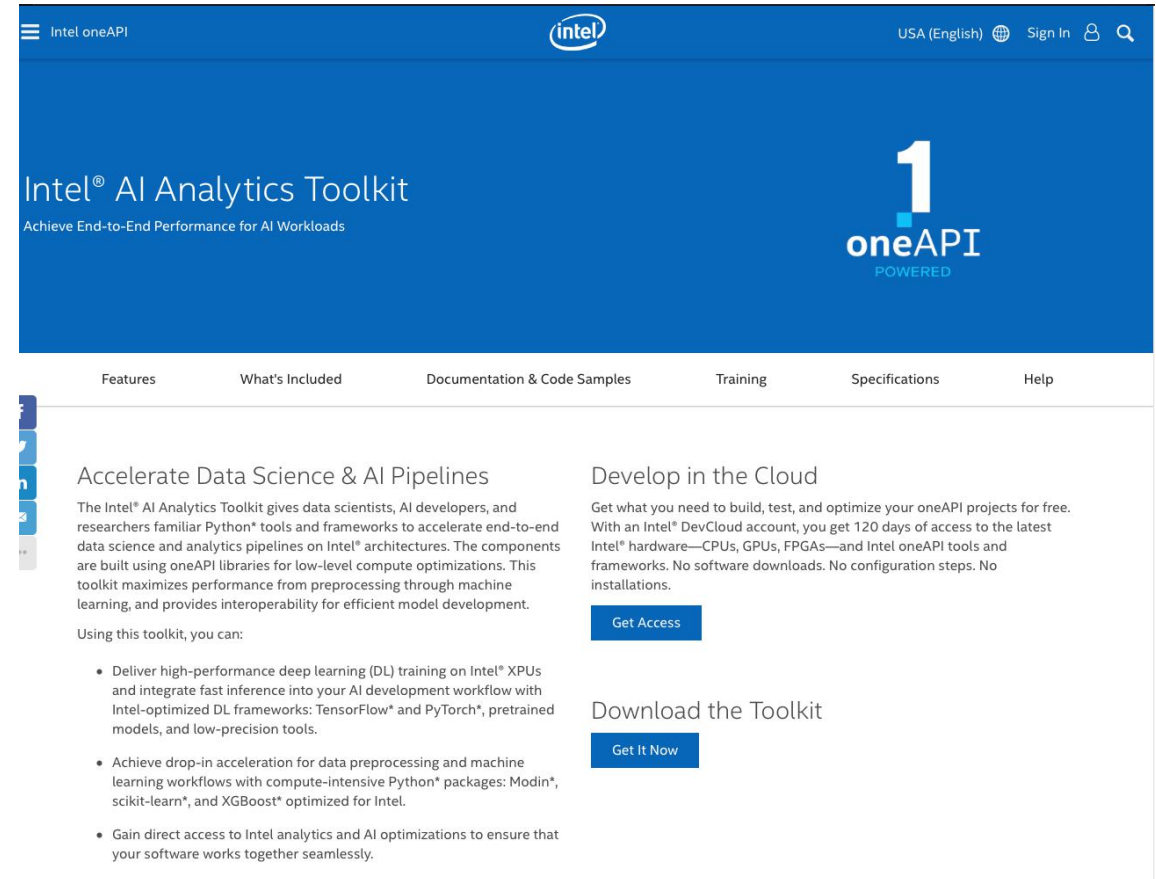
§ OneDAL, scikit-learn, XGBoost, etc.

## § Optimized and scalable communication using OneCCL

## § Spark BigData Analytics

## § DAOS Object storage for fast I/O and for workflows

## § Profiling and debugging tools



<https://software.intel.com/content/www/us/en/develop/tools/oneapi/ai-analytics-toolkit.html>

# Community Data Sharing with Eagle

- A global filesystem deployed to bring larger and more capable production-level file sharing to facility users
  - A space for broader distribution of reassembled data acquired from various experiments
    - Data originating at the ALCF
    - Greater scientific community
  - Science community can access uploaded data, and ALCF users are able to directly access the data for analysis
  - Designed to foster experimentation
    - Analysts are able to write new algorithms to attempt analyses that have never been performed
- **HPE ClusterStor E1000**
  - **100 petabytes of usable capacity**
  - **8,480 disk drives**
  - **Lustre filesystem**
  - **160 Object Storage Targets**
  - **40 Metadata Targets**
  - **HDR InfiniBand network**
  - **650 GB/s rate on data transfers**

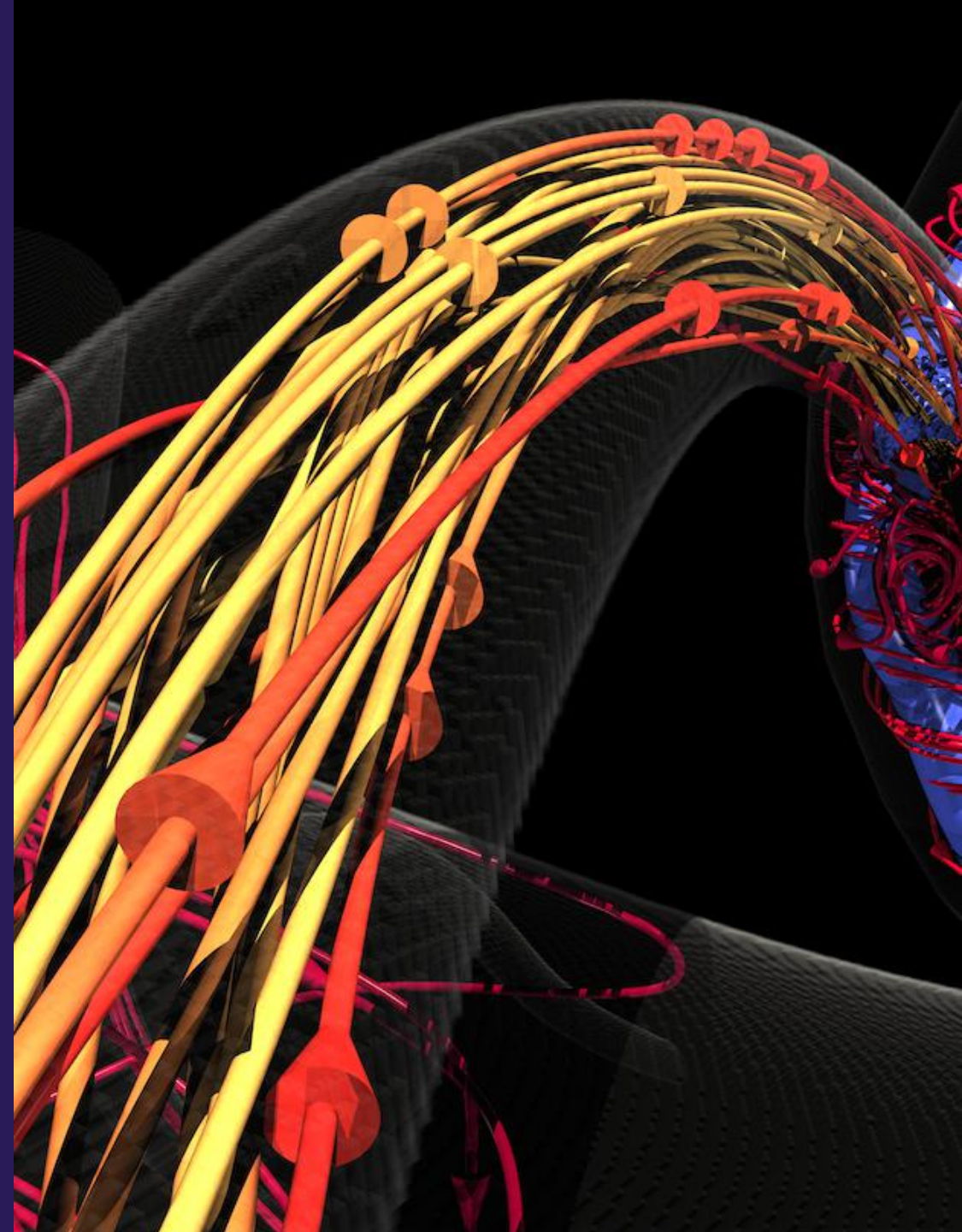
# ALCF AI-Testbed

## Advancing science with HPC

- ALCF AI pathfinding effort provides insights on cutting-edge AI technology and how it improves science outcomes
- Evaluates the usability and performance of machine learning-based applications running on these accelerators
  - a deep learning accelerator, reconfigurable dataflow units, intelligent processing unit- (IPU) based systems
- Ongoing work is guiding the facility toward a future marked by extreme heterogeneity in the compute: CPUs, GPUs, AI, and other accelerators

### AI testbeds include:

- SambaNova DataScale
- GraphCore MK1
- Groq
- Cerebras CS-2
- Habana Gaudi



# Science



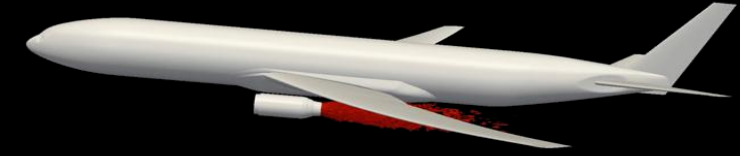
# Contribution to Science

To prepare for future exascale systems, the ALCF is driving a new paradigm for scientific computing.



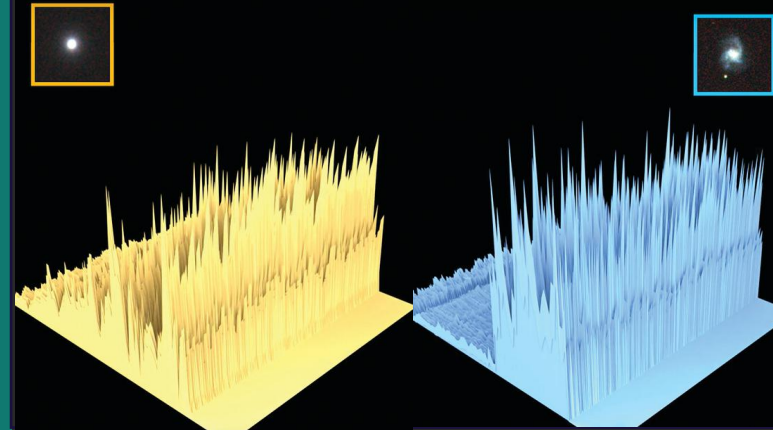
## Modeling & Simulation

Used to study things that are too big, too small, or too dangerous to study in a laboratory setting.



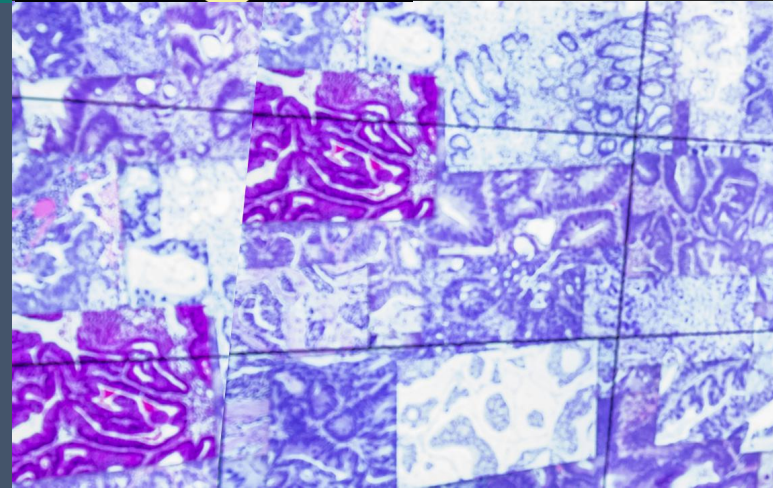
## Data Science

Researchers can glean insights from very large datasets produced by experimental, simulation, or observational methods.



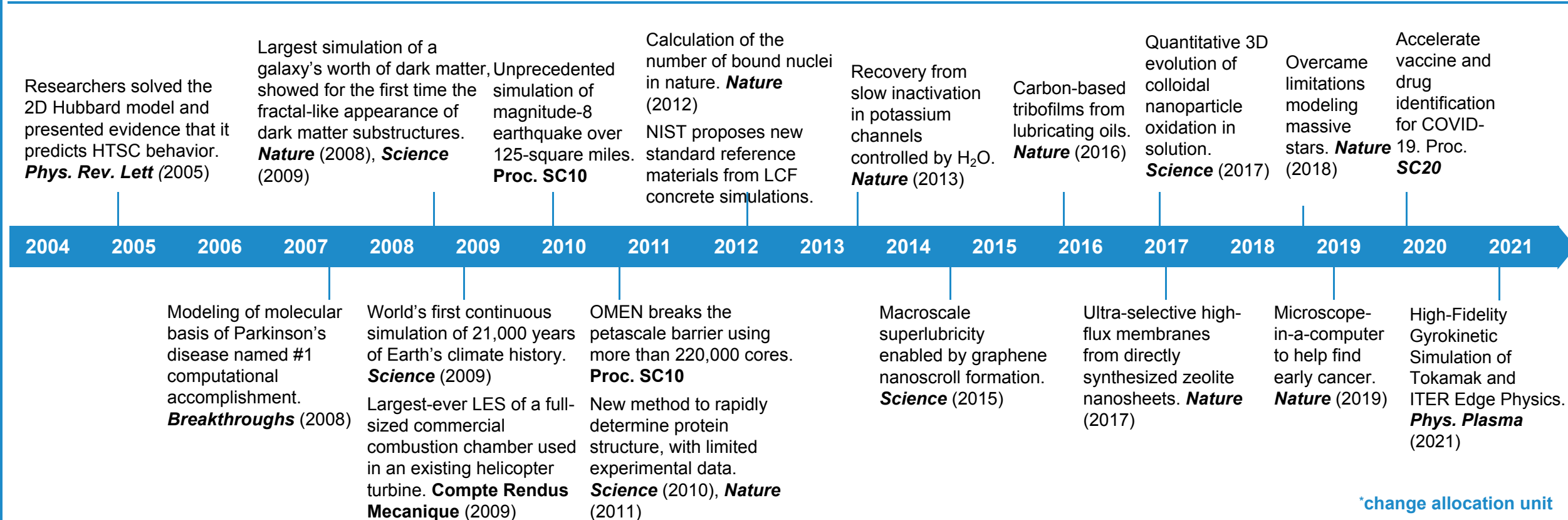
## Machine Learning

A type of artificial intelligence that trains computers to discover hidden patterns in data to make novel predictions without being explicitly programmed.



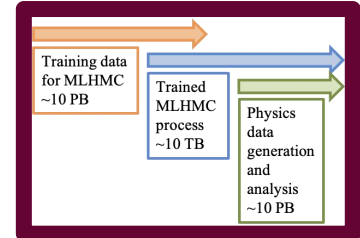
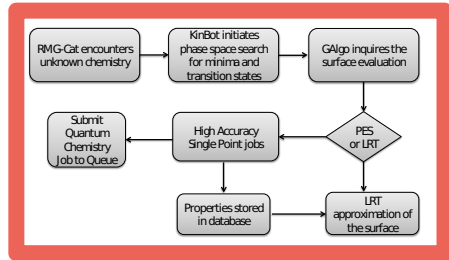
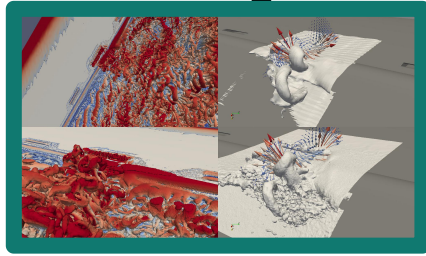
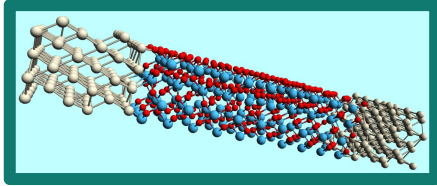
# LCF Growth and Impact of the INCITE Program

	2004	2005	2006	2007	2008	~2X per year					~3X per year					~4X per year		
	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019*	2020*	2021
Hours	4.9M	6.5M	18.2M	95M	268M	889M	1.6B	1.7B	1.7B	4.7B	5.8B	5.8B	5.8B	5.8B	5.9B	71M	37.6M	39.9M
Projects	3	3	15	45	55	66	69	57	60	61	59	56	56	55	55	62	47	47

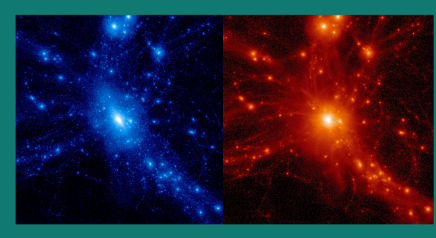
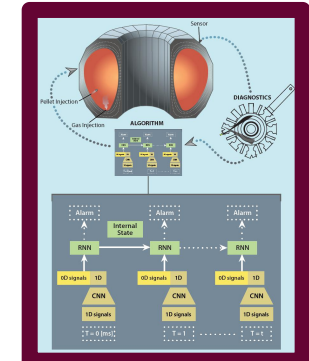
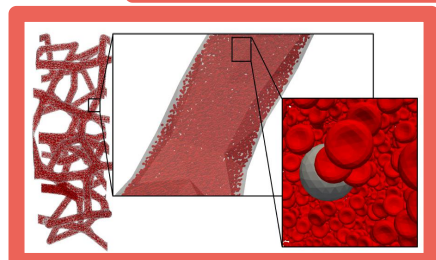
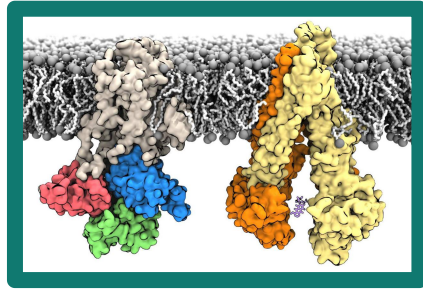
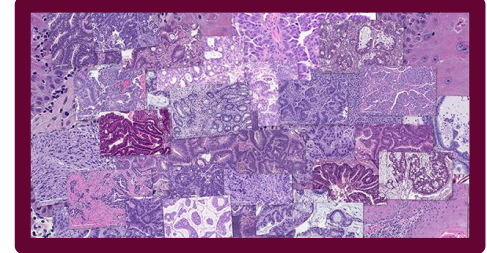
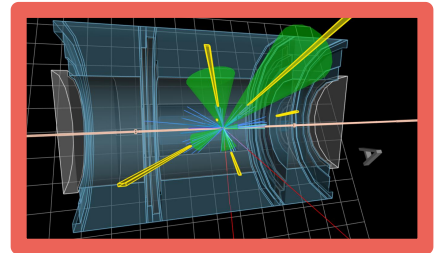
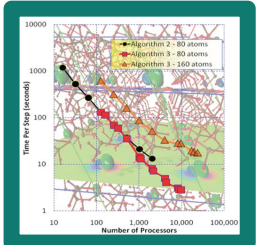
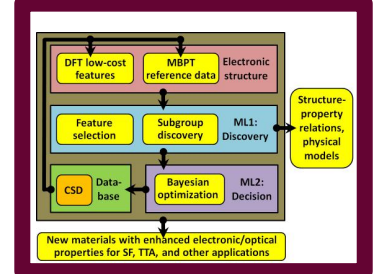
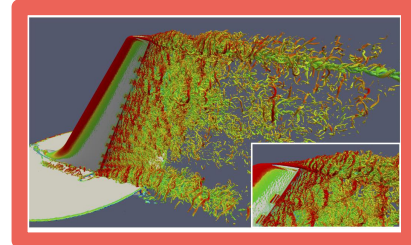
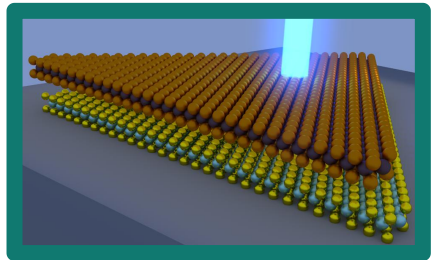
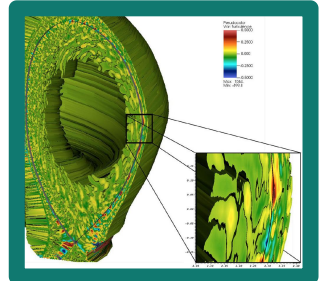
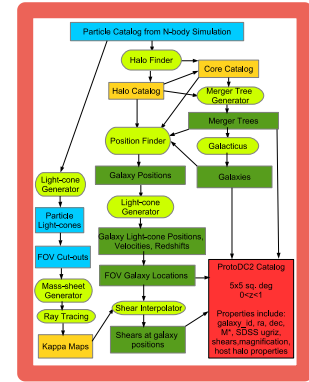
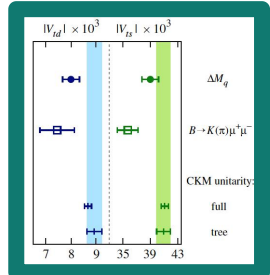
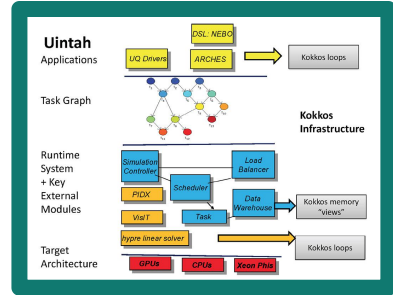


\*change allocation unit

# Aurora ESP Projects



S D L



Computing Facility

# ESP Projects

**Anouar Benali (ANL)**

Extending Moore's Law computing with Quantum Monte Carlo

**Martin Berzins (U. Utah)**

Design and evaluation of high-efficiency boilers for energy production using a hierarchical V/UQ approach

**CS Chang (PPPL)**

High fidelity simulation of fusion reactor boundary plasmas

**Theresa Windus (Ames)**

NWChemEx: Tackling Chemical, Materials & Biochemical Challenges in the Exascale Era

**Katrin Heitmann (ANL)**

Extreme-Scale Cosmological Hydrodynamics

**Ken Jansen (U. Colorado)**

Extreme Scale Unstructured Adaptive CFD: From Multiphase Flow to Aerodynamic Flow Control

**Norman Christ (Columbia)**

Lattice Quantum Chromodynamics Calculations for Particle and Nuclear Physics

**Aiichiro Nakano (USC)**

Metascalable Layered Materials Genome

**Benoit Roux (U. Chicago)**

Free Energy Landscapes of Membrane Transport Proteins

**David Bross (ANL)**

Exascale Computational Catalysis

**Salman Habib (ANL)**

Dark Sky Mining

**Ken Jansen (U. Colorado)**

Data Analytics and Machine Learning for Exascale CFD

**Walter Hopkins(ANL)**

Simulating and Learning in the ATLAS detector at the Exascale

**Amanda Randles (Duke U.)**

Extreme-scale In Situ Visualization and Analysis of Fluid-Structure-Interaction Simulations

**Will Detmold (MIT)**

Machine Learning for Lattice Quantum Chromodynamics

**Nicola Ferrier (ANL)**

Enabling Connectomics at Exascale to Facilitate Discoveries in Neuroscience

**Noa Marom (CMU)**

Many-Body Perturbation Theory Meets Machine Learning to Discover Singlet Fission Materials

**Rick Stevens (ANL)**

Virtual Drug Response Prediction

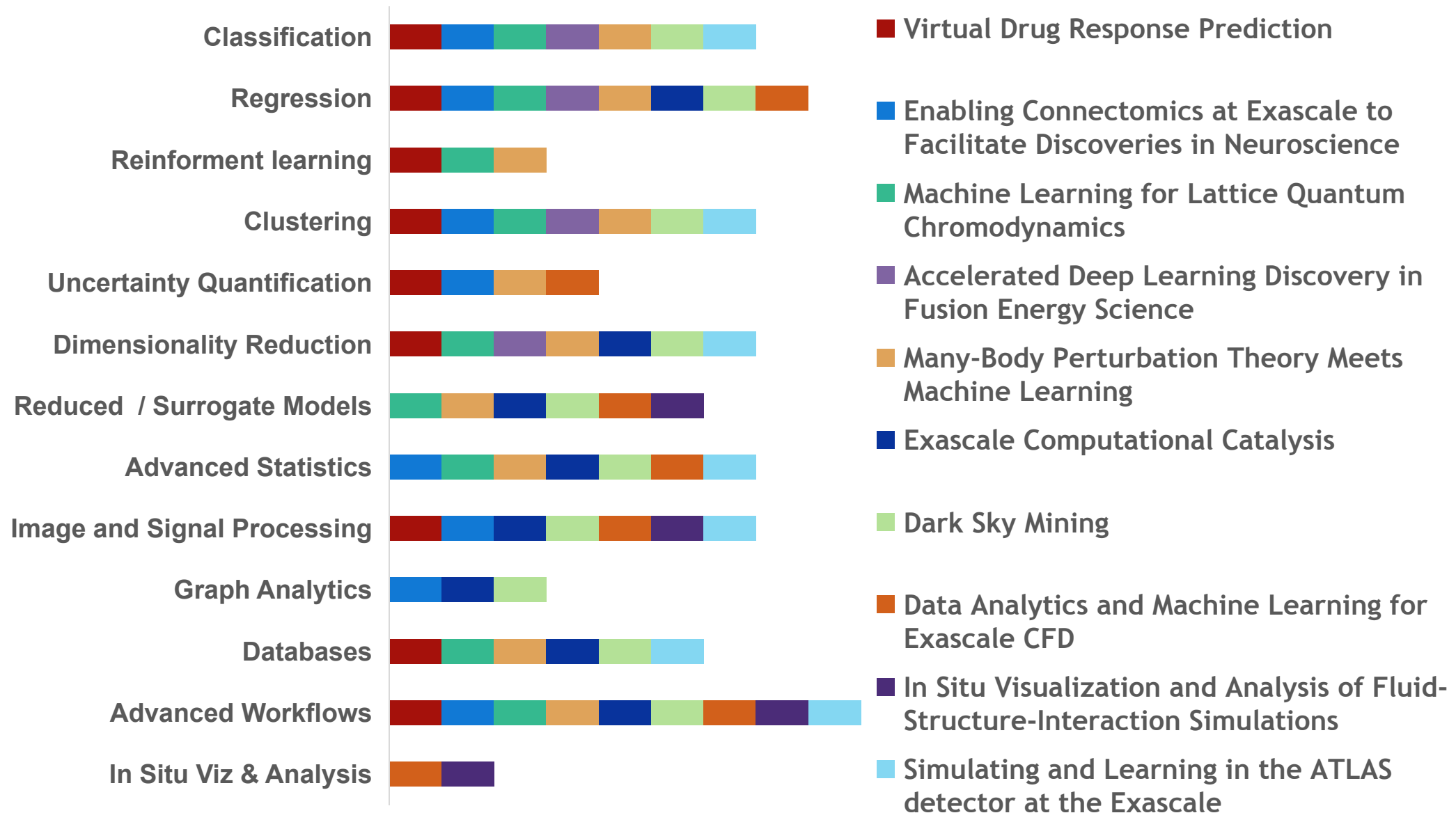
**Bill Tang (Princeton)**

Accelerated Deep Learning Discovery in Fusion Energy Science

S D L



# AURORA ESP Data and Learning Projects and Methods



Learning

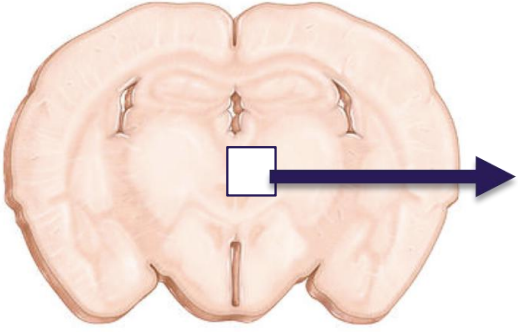
Data

# Connectomics Data-driven Models



Mouse

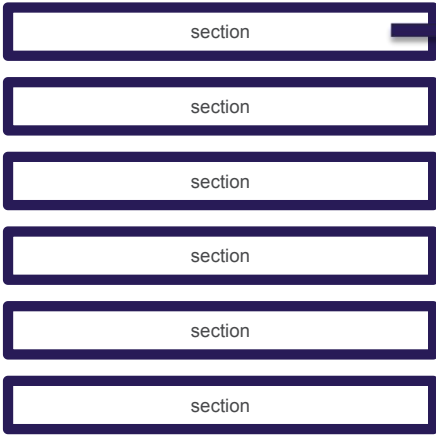
Mouse brain: 70M neurons



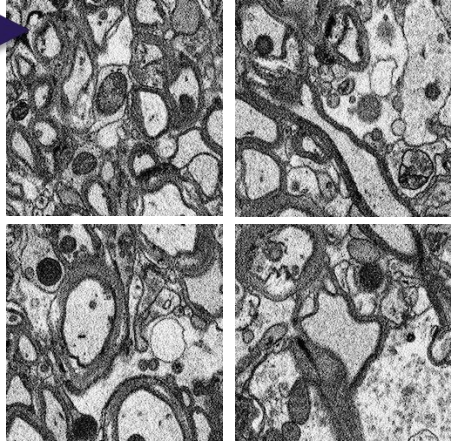
~1cm<sup>3</sup>



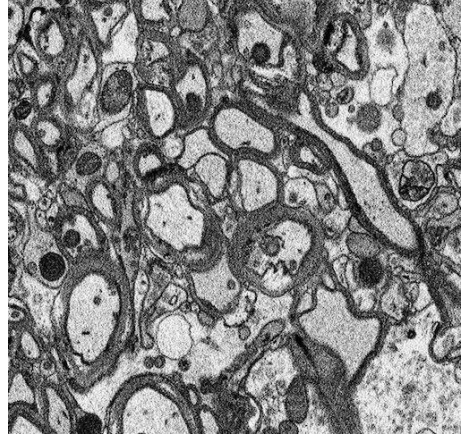
~1mm<sup>3</sup>



25000  
40nm sections  
1mm x 1mm  
(6nm resolution)



Each section  
imaged with EM as  
N tiles (8 bit)  
80K x 40K pixels



Sections  
stitched  
together

How much image data is 1mm<sup>3</sup> ? 1e15 voxels -> ~1 PB

# data challenges in connectomics



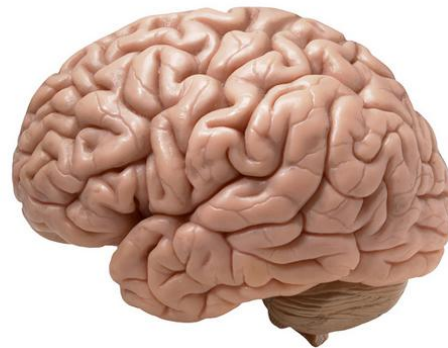
Mouse brain: 70M neurons



$\sim 1\text{cm}^3$

How much image data is  $1\text{cm}^3$  ?  **$\sim 1\text{EB}$**

Human brain: 80B neurons



$\sim 1000\text{cm}^3$

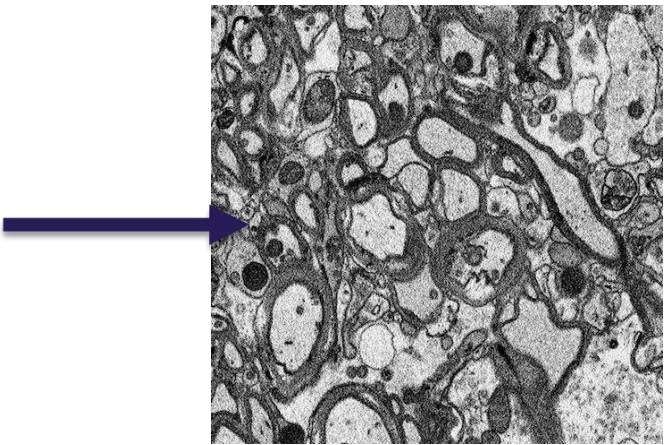
How much image data is  $1000\text{cm}^3$  ?  **$\sim 1000\text{ EB}$**   
(6nm x 6nm x 40nm)

Reconstructed data will be much larger:

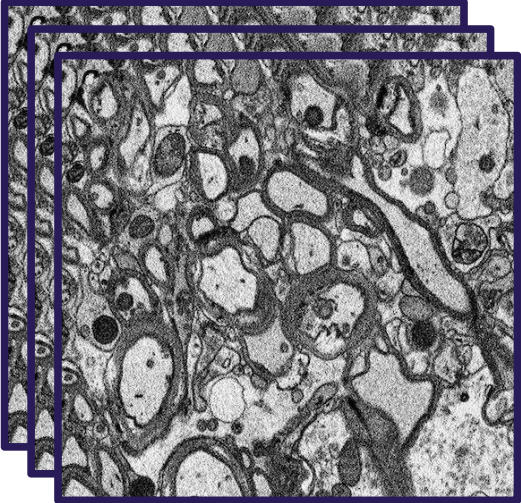
- Segmentation labels for each voxel (4x voxel data)
- 3D Mesh
- Skeleton

The structures are expected to be used to seed simulations to study flow in neuro transmitters, in better modeling the brain, among others.

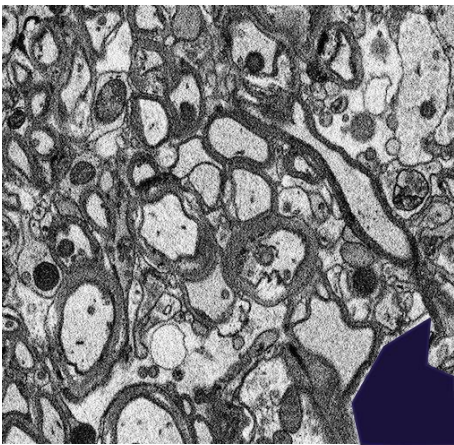
# Connectomics processing



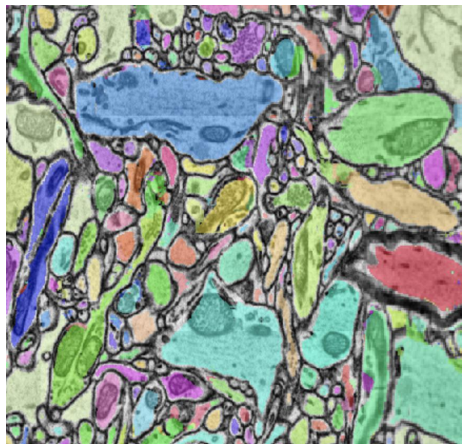
Sections stitched together



Align sections



Mask out non-target objects



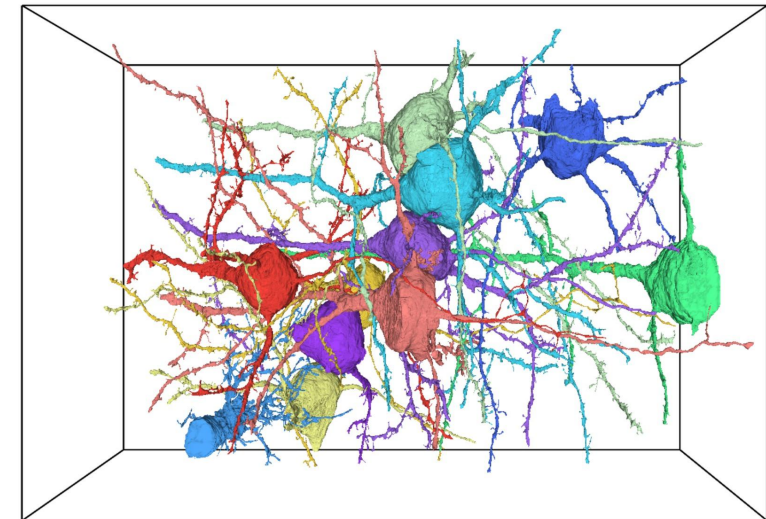
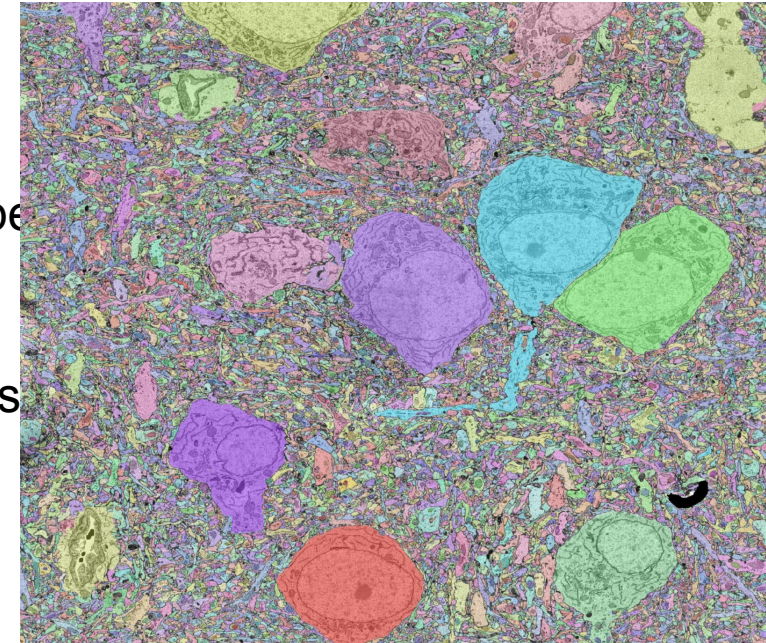
Segment target objects

# large-scale reconstruction

- **Inference (and training) has scaled** on CPU-based and GPU-based supercomputers (with high granularity: overlapping subvolumes)
  - Achieved million-way concurrency on Theta supercomputer
- Image stitching and alignment components are being scaled as well to ensure a robust pipeline

## Exascale Inference Problem:

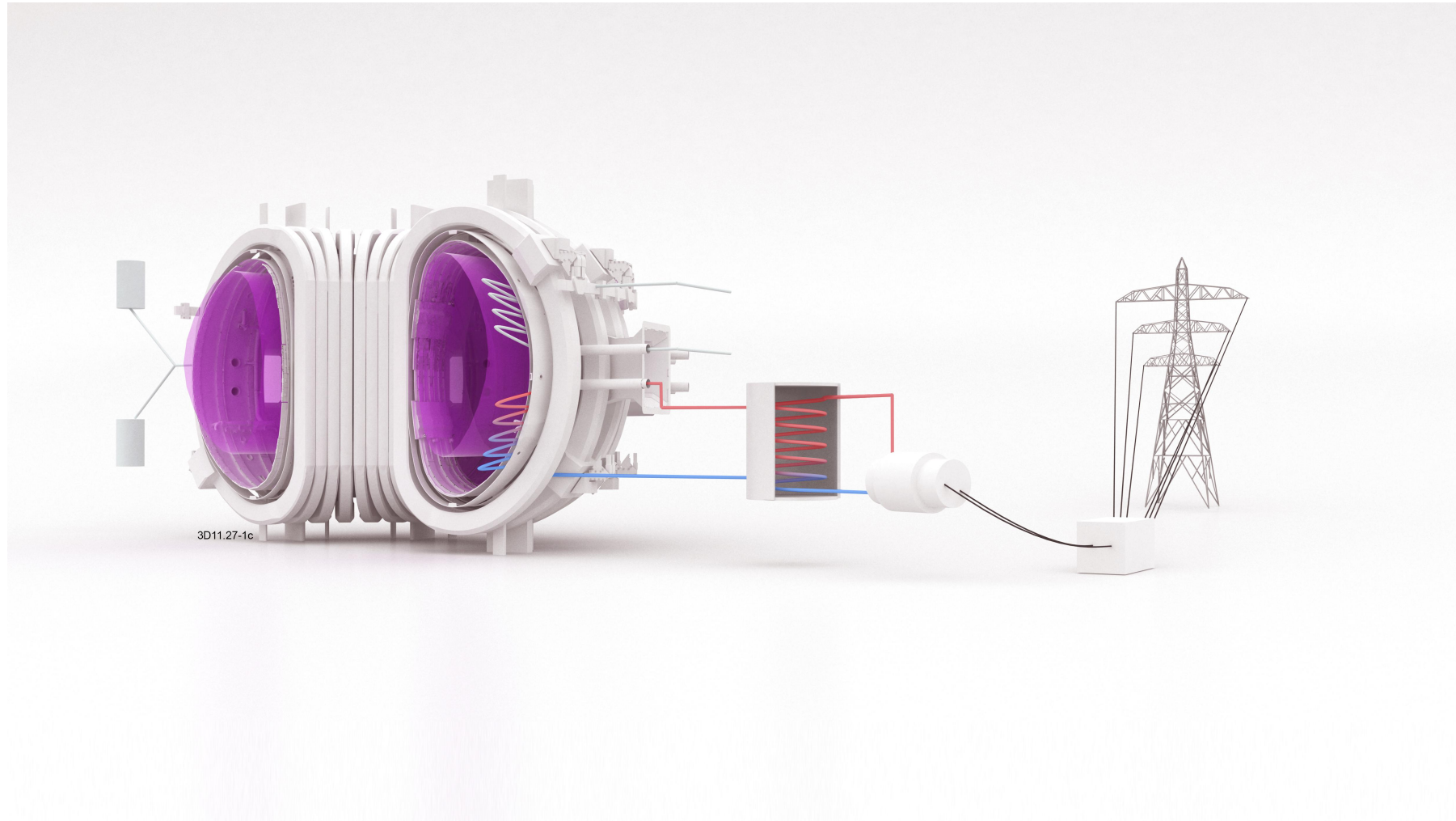
- On a single GPU (A100), we achieve ~80 MegaVoxels/hour using 32-bit (There is still room for improvement here)
- In reduced precision (8-16 bits), we expect ~1 GigaVoxel/hour per GPU
- 1 PetaVoxel (1mm<sup>3</sup>) will take ~1M GPU node hours
- Approximately, **24 hours on a system with 50K GPUs** (considering overlapping subvolumes)
- For a mouse brain (1cm<sup>3</sup>), 1 ExaVoxel, we would need **~3 years on an exascale system**



Dong, et al, "Scaling Distributed Training of Flood-Filling Networks on HPC Infrastructure for Brain Mapping", 2019 IEEE/ACM Third Workshop on Deep Learning on Supercomputers (DLS) at SC19

Vescovi, et al, "Toward an Automated HPC Pipeline for Processing Large Scale Electron Microscopy Data", 2020 IEEE/ACM 2nd Annual Workshop on Extreme-scale Experiment-in-the-Loop Computing (XLOOP) at SC19

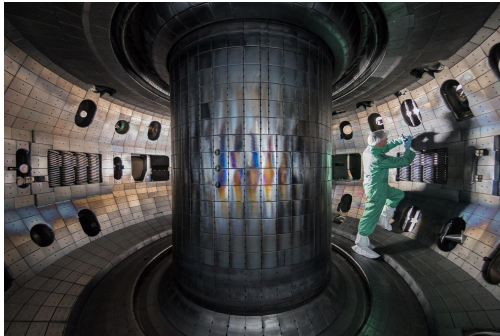
# Mission: delivery of fusion power on the grid



# Tokamaks

**Dead/broken tokamaks**  
(have access to these datasets, but not actively using them)

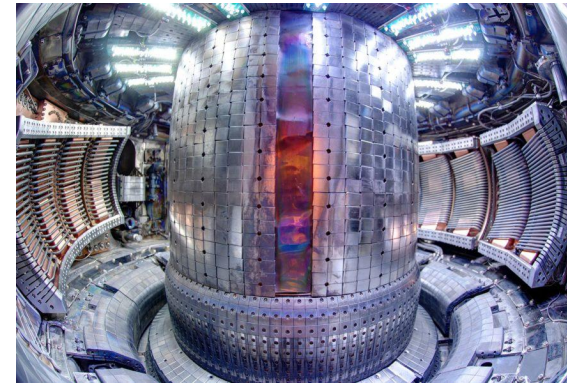
**Operational “traditional” tokamaks**  
(current main datasets)



General Atomics: DIII-D



CCFE: JET

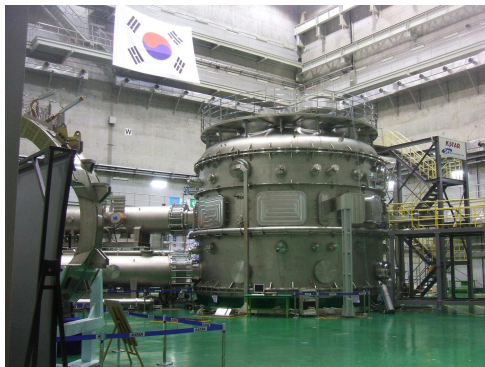


MIT: Alcator C-Mod (retired)

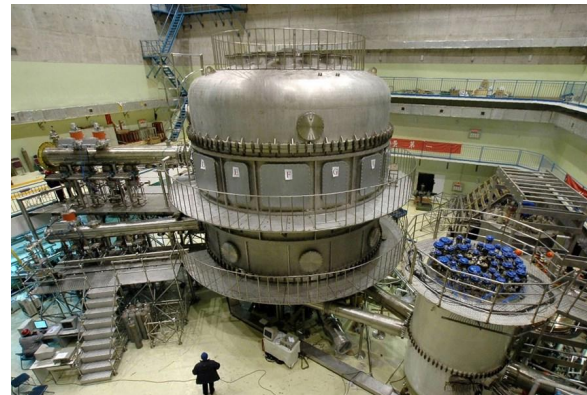


PPPL: NSTX-U (broken)  
Spherical Tokamak (ST)

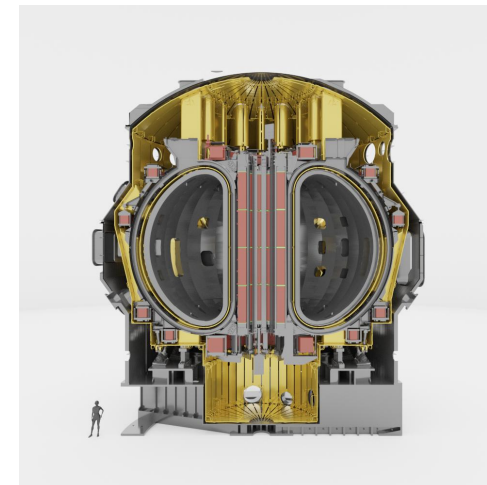
**Superconducting, long pulse, tokamaks**  
(want access to these datasets)



KFE: KSTAR



Hefei: EAST



JAEA: JT-60SA

# 2022: Fusion in the news

Science > Energy

## The World's Largest Tokamak Just Crushed the Record for Nuclear Fusion Energy

England's Joint European Torus (JET) produced 59 megajoules of energy for five seconds.



BY CAROLINE DELBERT PUBLISHED: FEB 16, 2022



### Ignition confirmed in a nuclear fusion experiment for the first time

A 2021 experiment achieved the landmark milestone of nuclear fusion ignition, which data analysis has now confirmed – but attempts to recreate it over the last year haven't been able to reach ignition again



PHYSICS 11 August 2022

By Karmela Padavic-Callaghan

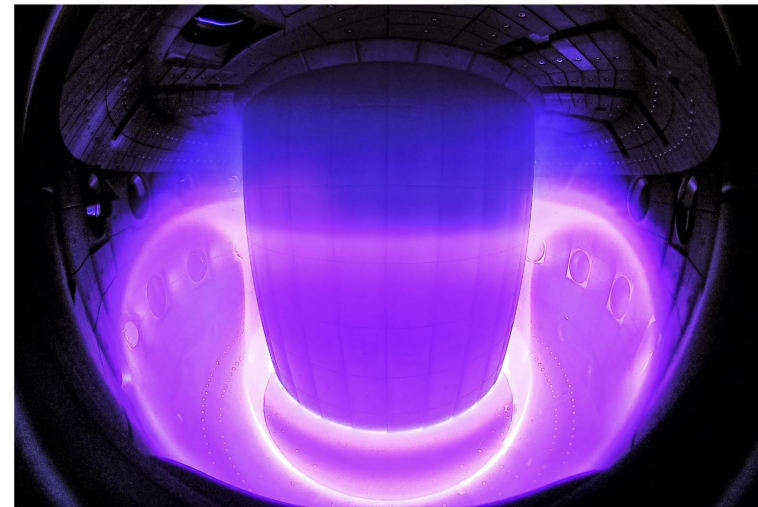


AMIT KATWALA

SCIENCE FEB 16, 2022 11:00 AM

### DeepMind Has Trained an AI to Control Nuclear Fusion

The Google-backed firm taught a reinforcement learning algorithm to control the fiery plasma inside a tokamak nuclear fusion reactor.



PHOTOGRAPH: CURDIN WÜTHRICH, SPC/EPFL



# 2022: Accelerated interest from government and industry

APRIL 19, 2022

## Readout of the White House Summit on Developing a Bold Decadal Vision for Commercial Fusion Energy

OSTP BRIEFING ROOM PRESS RELEASES

### Executive Summary

**The Biden-Harris Administration is developing a strategy to accelerate fusion energy— a clean energy technology that uses the sun and stars.** On March 17, 2022, the Department of Energy (DOE) and the U.S. Office of Science and Technology Policy (OSTP) and the U.S. Department of Energy hosted the first-ever White House Summit on Developing a Bold Decadal Vision for Commercial Fusion Energy. Over 1,200 viewers to witness fusion energy research, industry, academia, and other stakeholder groups and have inclusive conversations about an

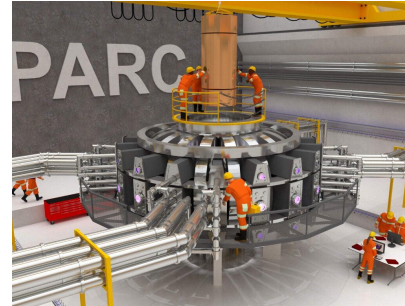
**DOE announced three new initiatives:**

### Congress Provides Record Funding for Fusion Energy and Initiates New Public Private Partnership

In the *2022 Consolidated Appropriations Act*, the \$1.5 trillion spending bill which funds the government for the 2022 Fiscal Year, the U.S. Congress provided record new funding for fusion energy research in the Department of Energy's Office of Fusion Energy Sciences and provides enough resources to initiate a new milestone-based public-private partnership program.

Long a priority for the Fusion Industry Association, this new milestone-based public private partnership program would allow the Department of Energy to partner with private companies to build new fusion energy devices, focused towards defined milestones, as agreed by a competitive application process. The program was created by Congress in the *Energy Act of 2020*, but the \$45 million directed to this program is the first funding that was directly appropriated for it. Although the funding will be enough to initiate the program and define the parameters of application, it will not be enough - on its own - to meet ambitious goals for fusion energy development.

The milestone program was authorized at \$325 million in funding over 5 years, a number that the House of Representatives has proposed to raise to \$800 million in the *Science for*



## Department of Energy Announces \$50 Million for a Milestone-Based Fusion Development Program

SEPTEMBER 22, 2022

Office of Science » Department of Energy Announces \$50 Million for a Milestone-Based Fusion Development Program

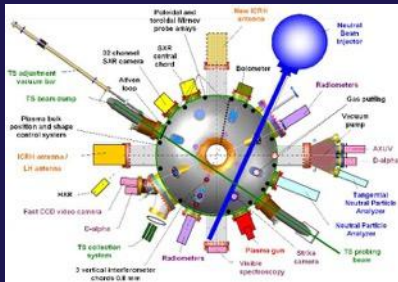
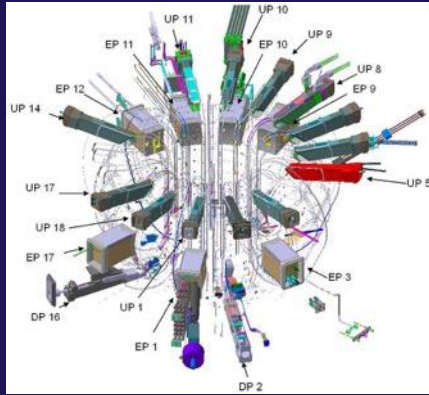
*This new public-private partnership program is the first step toward realizing the Administration's bold decadal vision for commercial fusion energy*

**WASHINGTON, D.C.** – Today, the **U.S. Department of Energy (DOE)** announced up to \$50 million to launch a new milestone-based fusion development program as authorized in the Energy Act of 2020. This program will support for-profit entities, who may team with national laboratories, universities, and others to meet major technical and commercialization milestones toward the successful design of a fusion pilot plant (FPP) that will help bring fusion toward technical and commercial viability. The program is informed by recent reports from the Fusion Energy Sciences Advisory Committee; the National Academies of Sciences, Engineering, and Medicine; community workshops; and input from private industry.

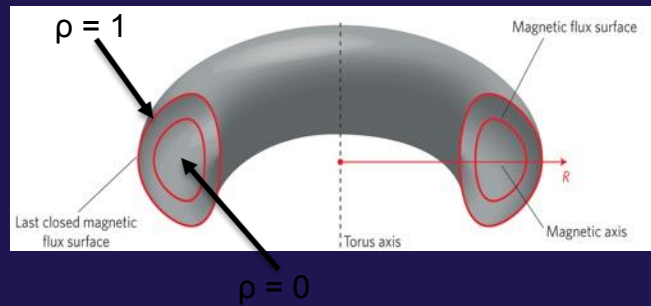
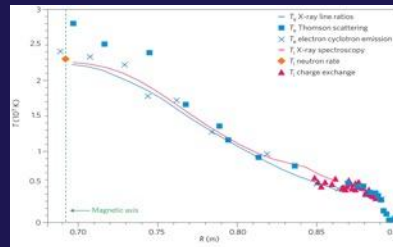
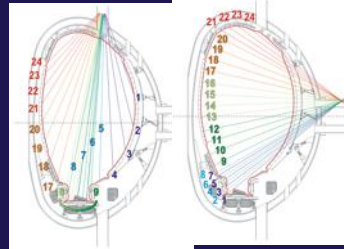
"Fusion holds the promise of being an on-demand, safe, abundant source of carbon-free primary energy and electricity, with the potential to transform the way we generate and use energy," said

# DIAGNOSTIC DATA SOURCES

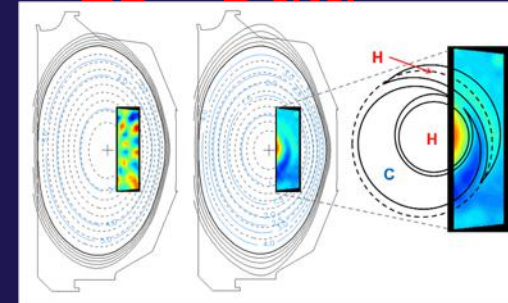
## 0D Scalar Data



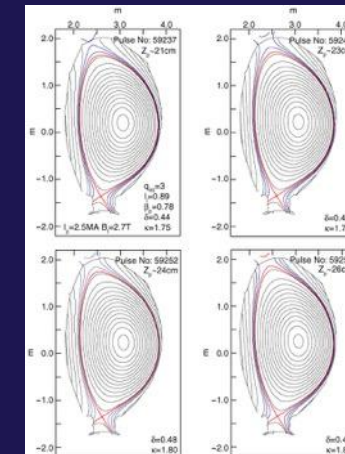
## 1D Profile Data



## 2D+ Data

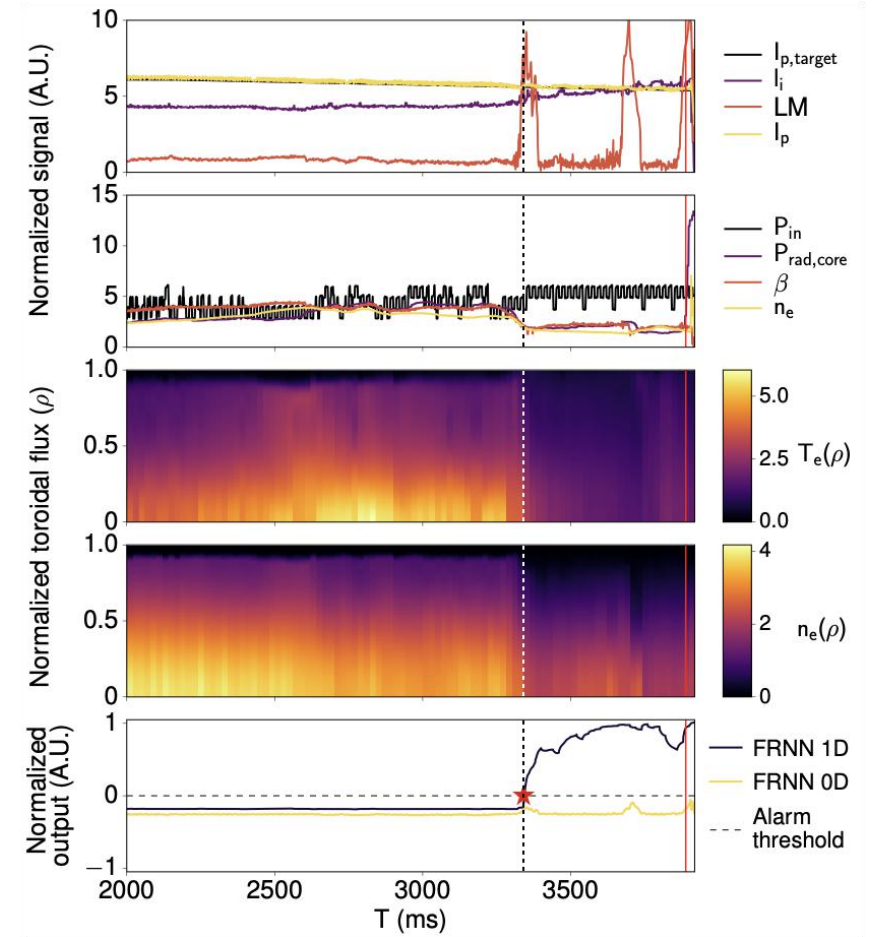
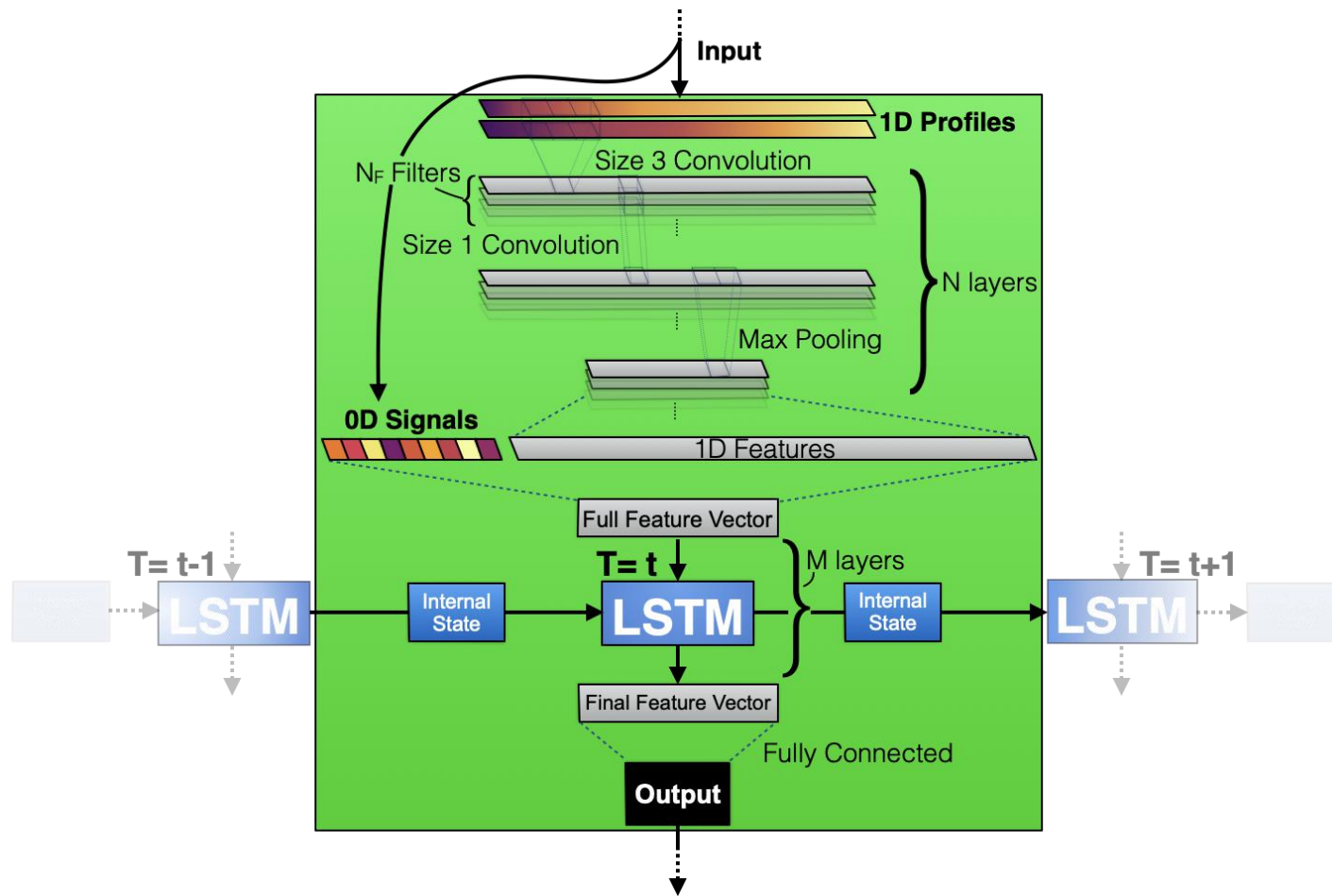


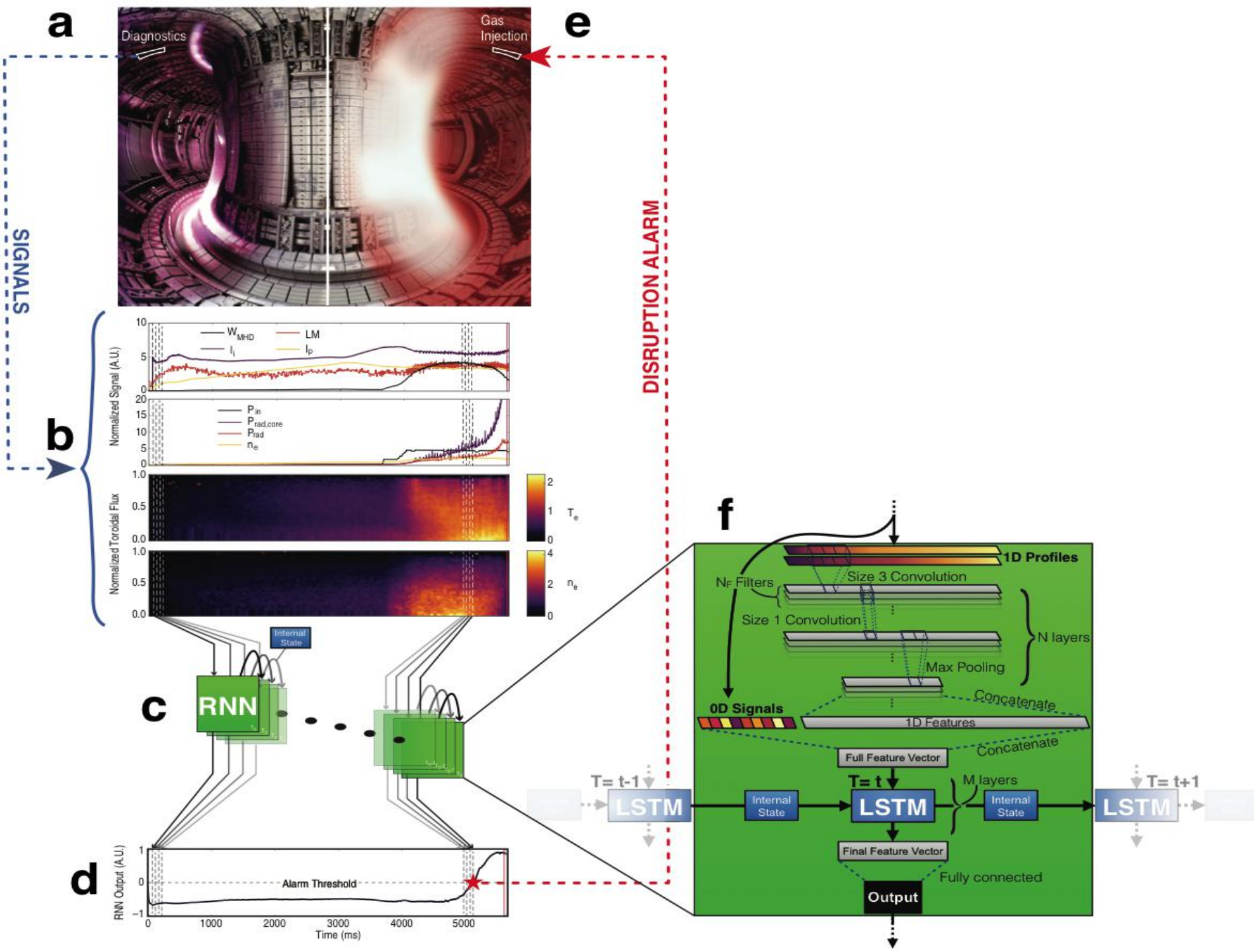
ECEI imaging



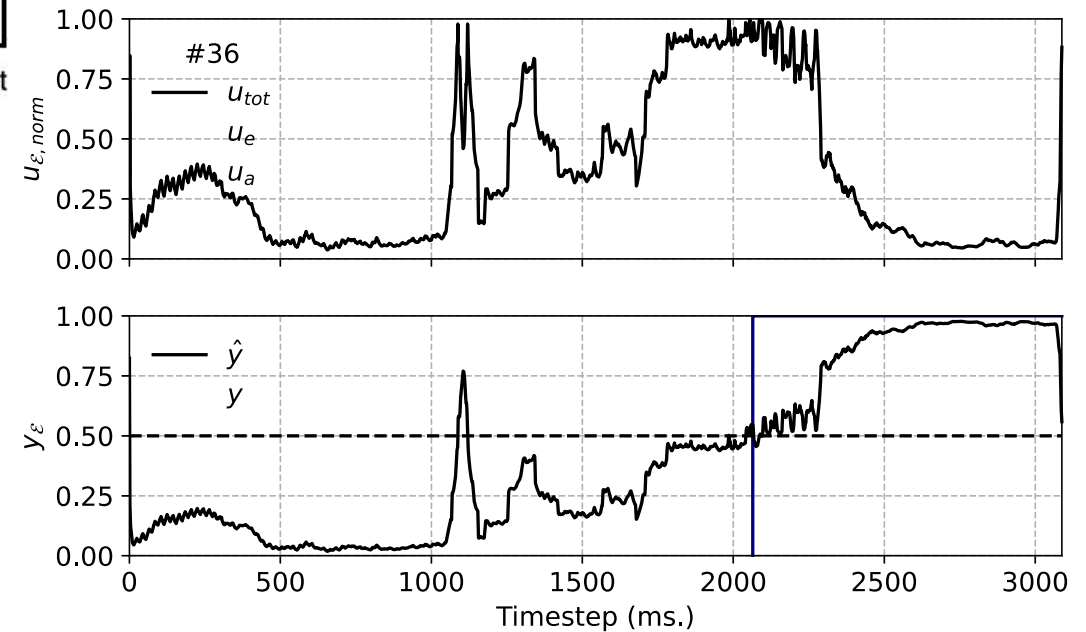
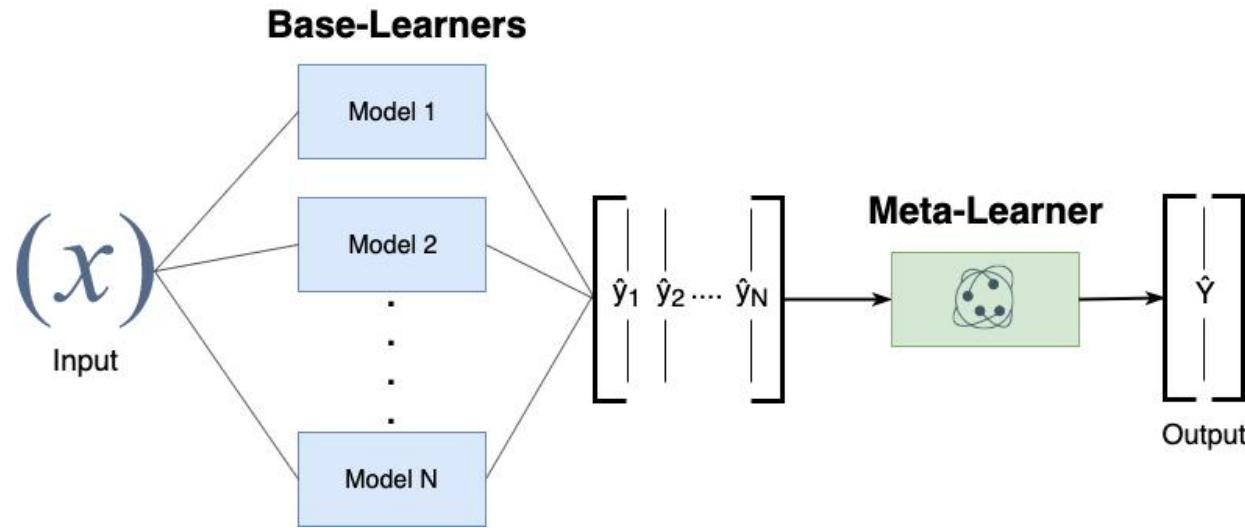
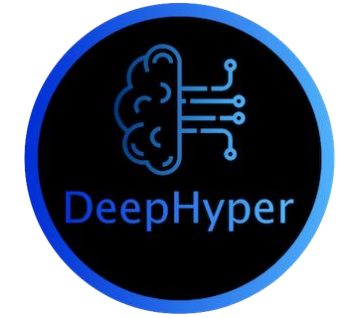
Magnetic equilibria

# LSTM-based architecture

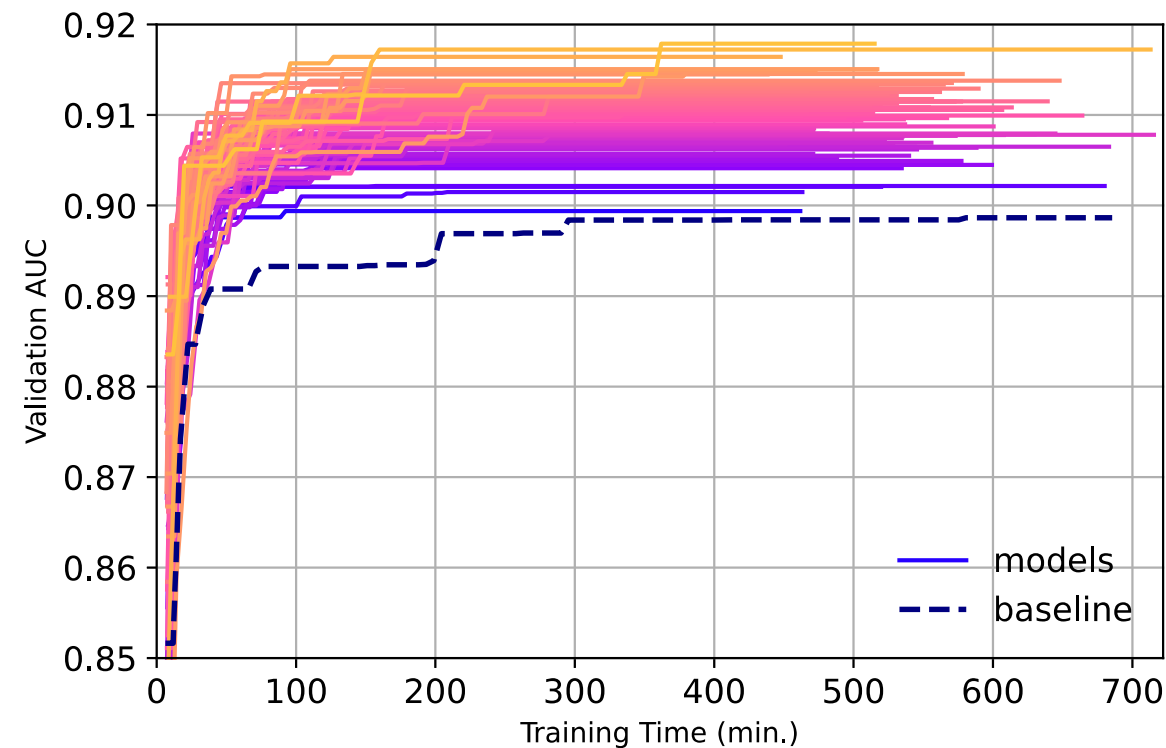
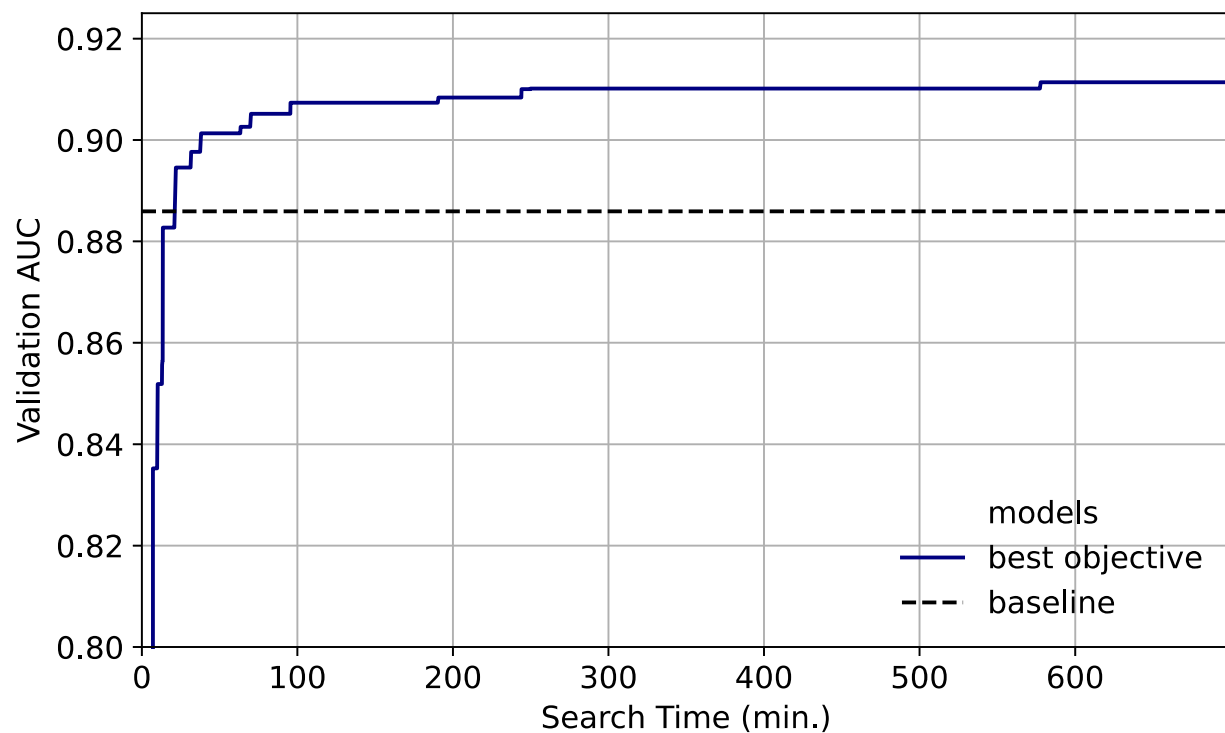




# Motivation for fast LSTM training at scale: gradient boosted ensembles



# Search trajectory and ensemble creation



# ALCF AI Testbeds

<https://www.alcf.anl.gov/alcf-ai-testbed>



Cerebras (CS-2)



SambaNova



Graphcore



Habana




Groq

- Infrastructure of next-generation machines with hardware accelerators customized for artificial intelligence (AI) applications.
- Provide a platform to evaluate usability and performance of machine learning based HPC applications running on these accelerators.
- The goal is to better understand how to integrate AI accelerators with ALCF's existing and upcoming supercomputers to accelerate science insights

	<b>Cerebras CS2</b>	<b>SambaNova Cardinal SN10</b>	<b>Groq GroqCard</b>	<b>GraphCore GC200 IPU</b>	<b>Habana Gaudi1</b>	<b>NVIDIA A100</b>
<b>Compute Units</b>	850,000 Cores	640 PCUs	5120 vector ALUs	1472 IPUs	8 TPC + GEMM engine	6912 Cuda Cores
<b>On-Chip Memory</b>	40 GB	>300MB	230MB	900MB	24 MB	192KB L1 40MB L2
<b>Process</b>	7nm	7nm	14nm	7nm	14nm	7nm
<b>System Size</b>	2 Nodes	2 nodes (8 cards per node)	4 nodes (8 cards per node)	4 nodes (16 cards per node)	2 nodes (8 cards per node)	Several systems
<b>Estimated Performance of a card (TFlops)</b>	>5780 (FP16)	>300 (BF16)	>188 (FP16)	>250 (FP16)	>150 (FP16)	312 (FP16), 156 (FP32)
<b>Software Stack Support</b>	Tensorflow, Pytorch	SambaFlow, Pytorch	GroqAPI, ONNX	Tensorflow, Pytorch, PopArt	Synapse AI, TensorFlow and PyTorch	Tensorflow, Pytorch, etc
<b>Interconnect</b>	Ethernet-based	Infiniband	RealScale™	IPU Link	Ethernet-based	NVLink





Director's Discretionary (DD) awards support various project objectives from scaling code to preparing for future computing competition to production scientific computing in support of strategic partnerships.

Getting Started on ALCF Systems including the AI Testbed:

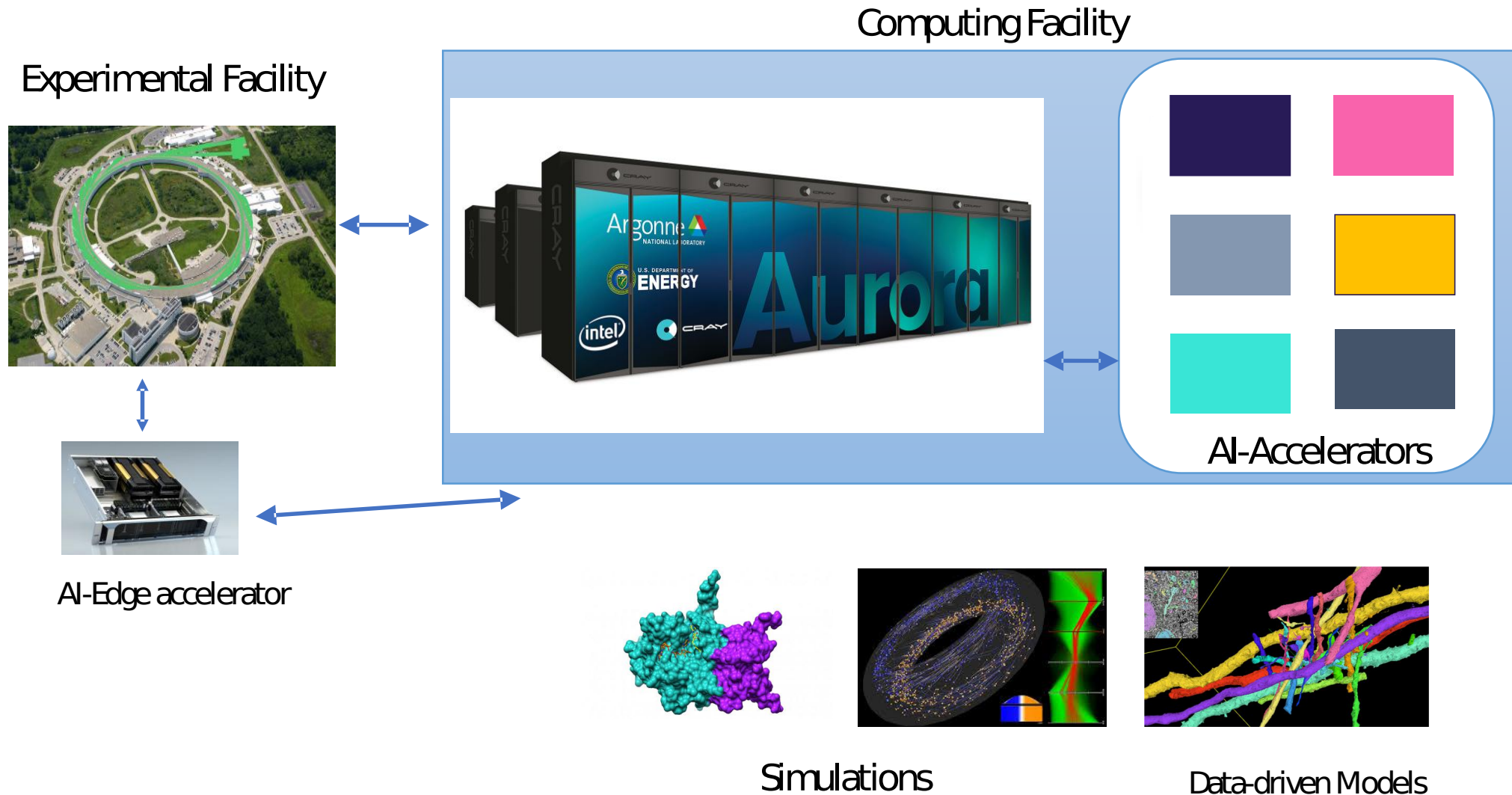
Apply for a Director's Discretionary (DD) Allocation Award

ALCF systems, including AI testbed systems - Cerebras CS-2 and SambaNova Datascale - are available for allocations.

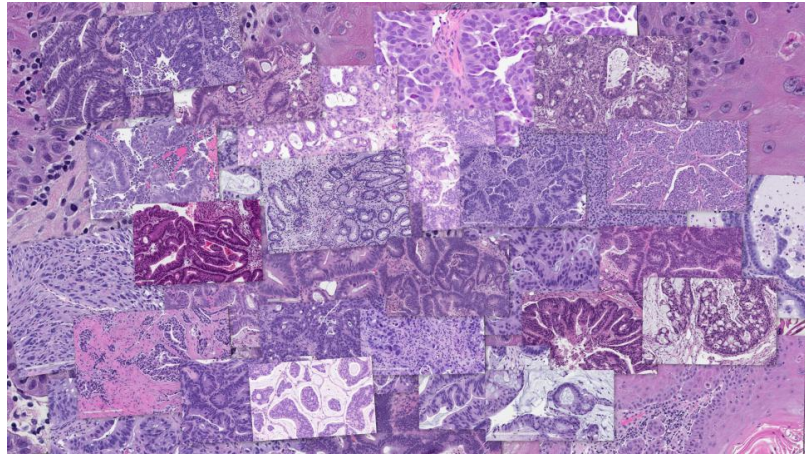
[Allocation Request Form](https://www.alcf.anl.gov/science/directors-discretionary-allocation-program)

<https://www.alcf.anl.gov/science/directors-discretionary-allocation-program>

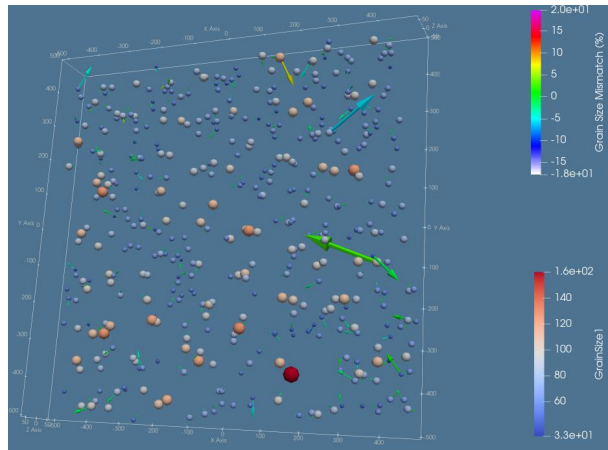
# Integrating AI systems in facilities



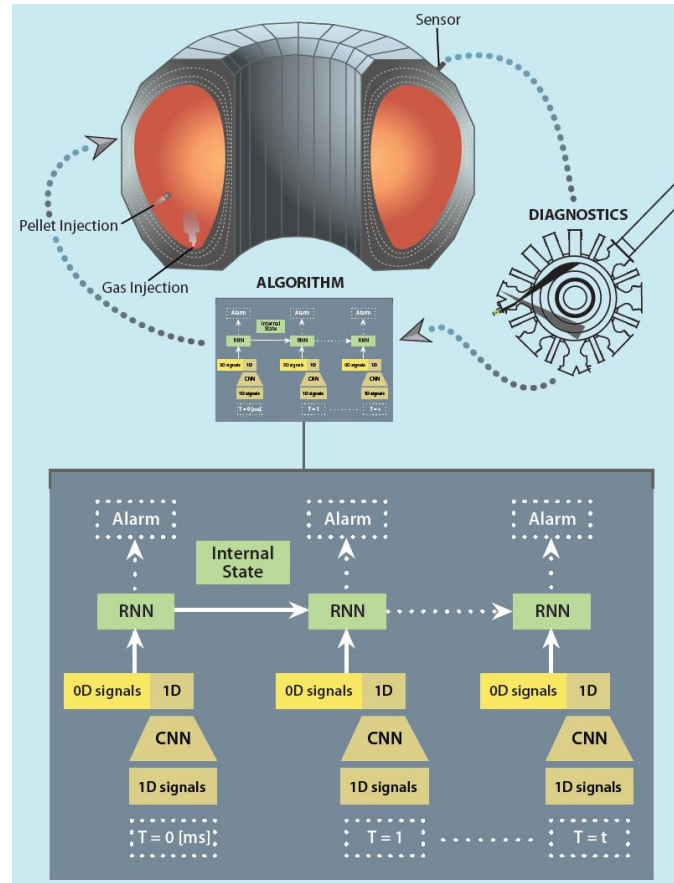
# AI FOR SCIENCE APPLICATIONS ON AI TESTBED



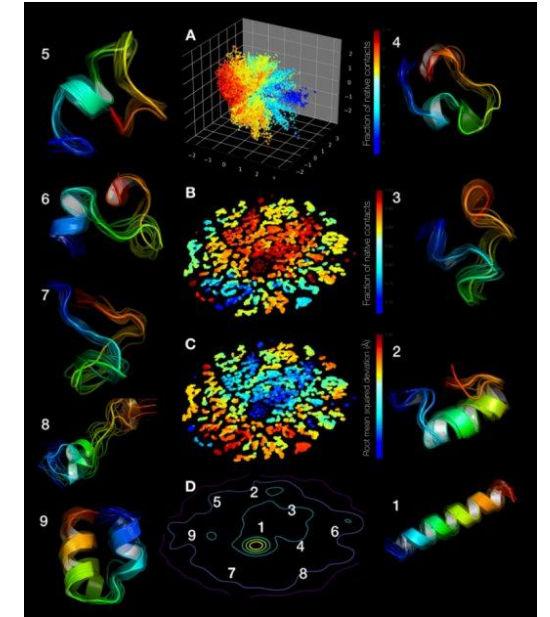
Cancer drug response prediction



Imaging Sciences-Braggs Peak



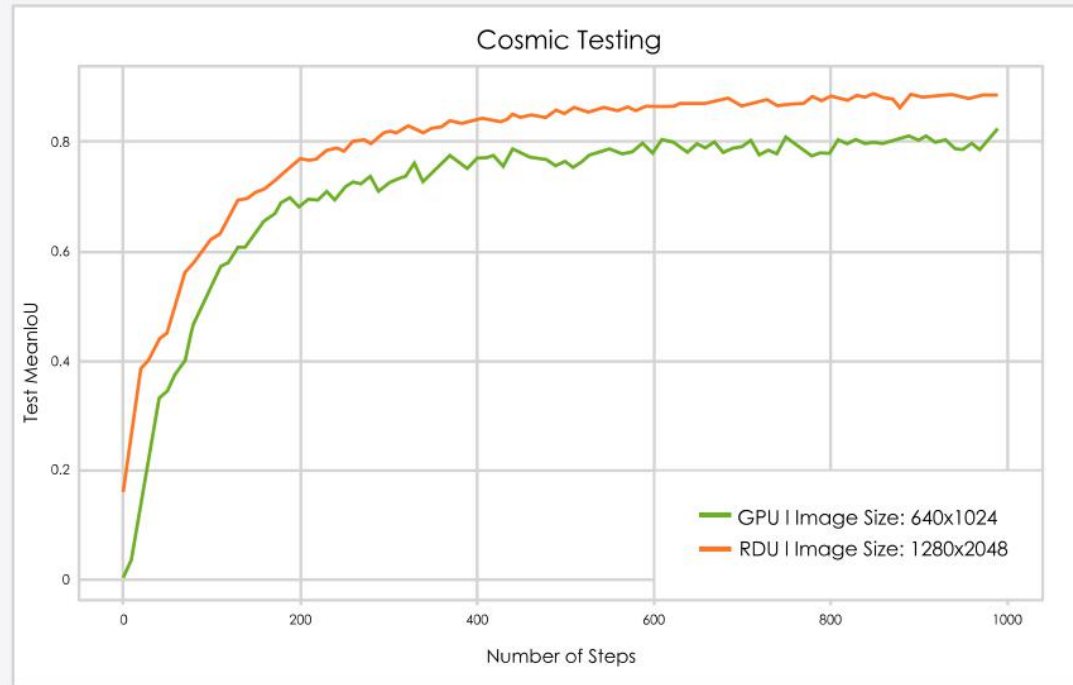
Tokamak Fusion Reactor operations



Protein-folding(Image: NCI)

and more..

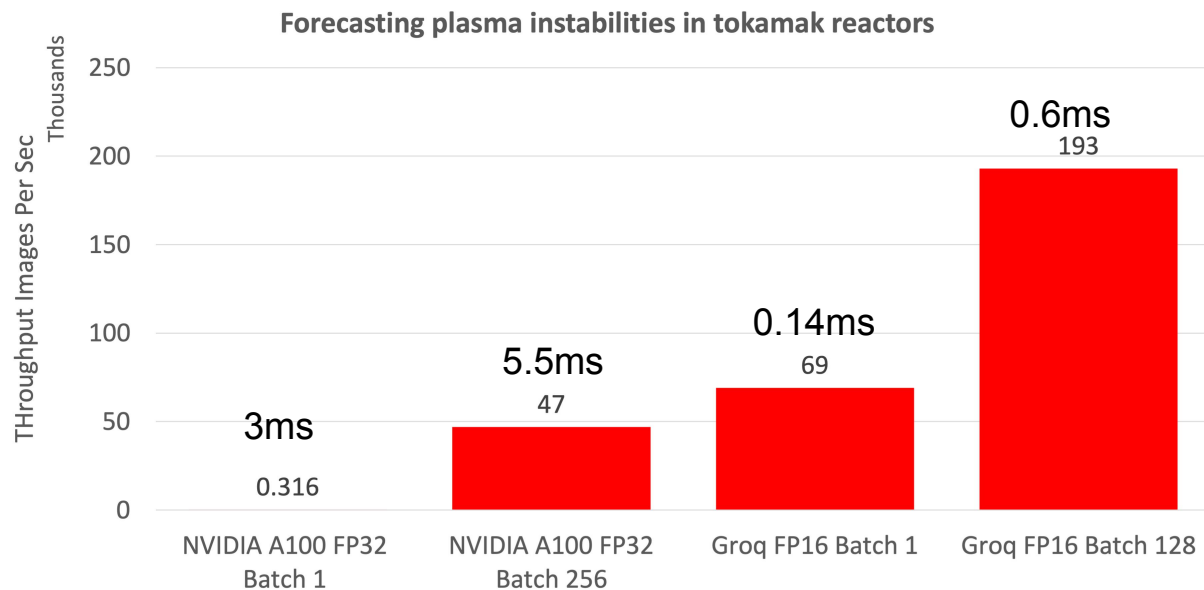
# Cosmic tagger on SambaNova datascale



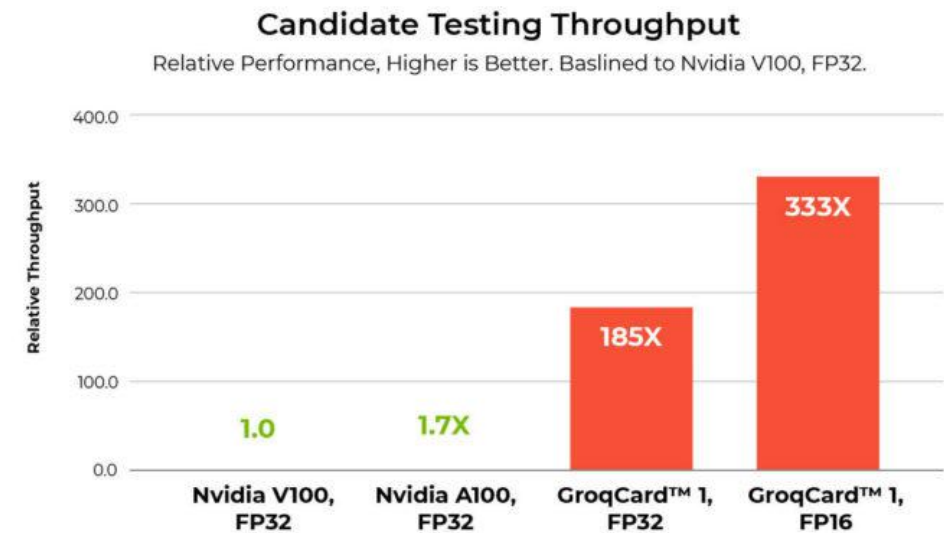
SambaNova RDUs able to accommodate larger image sizes and achieve higher accuracy

*M. Emani et al., "Accelerating Scientific Applications With SambaNova Reconfigurable Dataflow Architecture," in Computing in Science & Engineering, vol. 23, no. 2, pp. 114-119, 1 March-April 2021, doi: 10.1109/MCSE.2021.3057203.*

# Early Experience with Inference on Groq



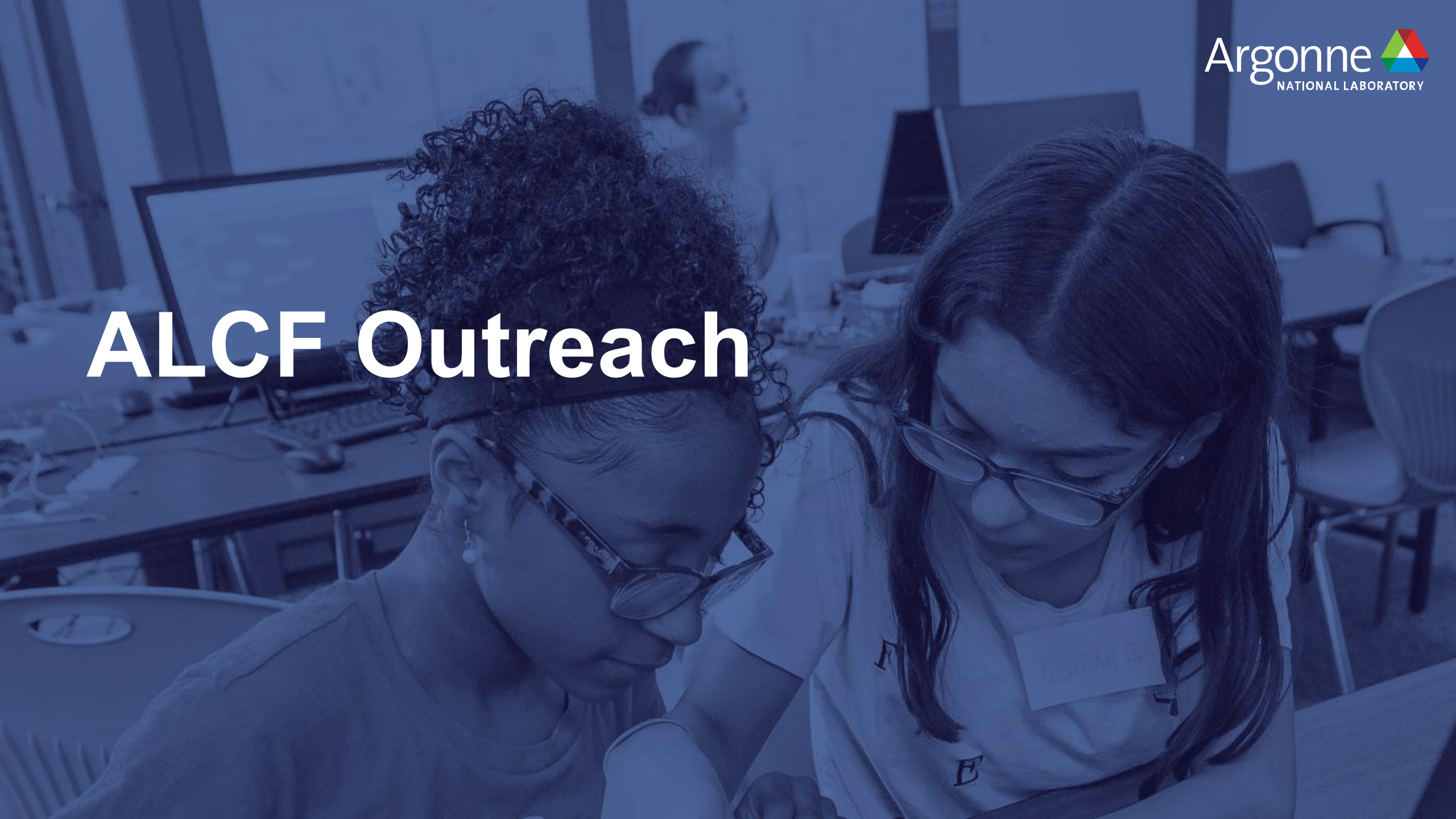
Forecasting Plasma Instability in Tokamak



COVID19 Candidate drug molecule screening

Promising results using GroqChip for science Inference use-cases with respect to latency and throughput in comparison to GPUs

# ALCF Outreach






# Summer Student Research Program

ALCF's internship program provides opportunities to work on real-world research projects.

- College students work side-by-side with staff mentors.
- Work utilizes some of the world's most powerful supercomputers.
- Opportunities in computational science, system administration, and data science.

# Opportunities at Argonne

## Undergraduate



### Undergraduate Programs

Here at Argonne, we work to make the world a better place through science and innovation, and we want to help empower undergraduates as they start their journey into the world of science and engineering.

We pursue discovery by pushing boundaries, challenging ourselves and each other, and stretching our abilities. This makes Argonne an excellent place for undergrads to explore and test their own ideas in science and technology.

Argonne can help undergrads grow, choose, and hone their areas of interest with programs that immerse them in cutting-edge research and discovery in all areas of the Lab. Over 100 students each year participate in the various internship opportunities Argonne offers. We also hire undergraduate students for part-time and temporary assignments to provide technical support to our scientists and engineers.

**NOTE: The laboratory is tentatively planning for student programs to be fully onsite during Summer 2022. However, a student's workplace may shift to a hybrid or virtual status contingent upon local COVID infection rates. Student candidates will be notified of changes in the lab's operational status at the time of offer and/or during their appointment.**

- Temporary Employment
- Internship Opportunities
- Visiting Students

**Communicating Science**  
From writing a report to delivering an oral presentation, our internship programs offer numerous opportunities for you to exercise your communication skills, and the Communicating Science resource is here to help.

[LEARN MORE](#)

**Educational Programs and Outreach**

- About Us
- STEM Outreach
- Learning Center
- Undergraduate Programs
- Graduate Programs
- Faculty Programs

SHARE

CONTACT US  
[Undergraduate Programs](#)  
[undergrad@anl.gov](#)  
[Subscribe to Education Email List](#)  
[Follow us on Instagram](#)


**Educational Programs and Outreach Mailing List**  
Educational Programs and Outreach uses immersive and engaging programs to create STEM pathways for students throughout their journeys. We offer many opportunities for students, families, education professionals, and others to help students learn and grow. Stay up-to-date on our latest science competitions, summer camps, undergraduate and graduate internships, and more by subscribing to our email list.

SUBSCRIBE TO EDUCATIONAL PROGRAMS AND OUTREACH UPDATES

email address  [SUBSCRIBE](#)

<https://www.anl.gov/education/undergraduate-programs>

## Graduate



### Graduate Programs

Argonne offers a variety of research appointments, cooperative education opportunities, and specialized schools to fit the needs – and schedules – of any graduate student. We encourage students pursuing a Master's degree as well as Ph.D. candidates to check out the many opportunities we offer. Come to Argonne, expand and create new knowledge – and change your world.

**NOTE: The laboratory is tentatively planning for student programs to be fully onsite during Summer 2022. However, a student's workplace may shift to a hybrid or virtual status contingent upon local COVID infection rates. Student candidates will be notified of changes in the lab's operational status at the time of offer and/or during their appointment.**

- Graduate Research
- Graduate Internships
- Graduate Temporary Employment
- Graduate Training Programs
- Visiting Student Program for Graduate Students

**Educational Programs and Outreach**

- About Us
- STEM Outreach
- Learning Center
- Undergraduate Programs
- Graduate Programs
- Faculty Programs


SHARE

CONTACT US  
[Graduate Programs](#)  
[graduate@anl.gov](#)  
[Subscribe to Education Email List](#)  
[Follow us on Instagram](#)

**FEATURED GRADUATE INTERNSHIP PROGRAM FLYERS**  
[W.J. Cody Associates Program](#)  
[Givens Associates Program](#)

<https://www.anl.gov/education/graduate-programs>

## Faculty



### Faculty Programs

Argonne offers a number of programs that provide faculty with opportunities to conduct research, network with Argonne experts and enhance their understanding of our mission and science through conferences and workshops.

Argonne is a place where scientists, engineers, and researchers immerse themselves in the grand challenges of our society. We strive to find solutions and make discoveries that provide an enduring value to our nation and to the future of the American scientific enterprise. We believe that bringing together the brightest minds and ideas from across disciplines helps answer the grand challenges of today. Through partnerships and collaboration, we push the boundaries of what seems possible. This is a great place for university faculty to connect, collaborate, and reinvigorate their passion for scientific research.

As a faculty researcher at Argonne, you can discover new ideas to take back to your students, shaping the next generation of scientists. Everyone benefits from the leading edge at Argonne.

- Faculty Research Participation
- Faculty Sabbatical
- Guest Faculty Research
- DOE Visiting Faculty Program

**Educational Programs and Outreach**

- About Us
- STEM Outreach
- Learning Center
- Undergraduate Programs
- Graduate Programs
- Faculty Programs

SHARE

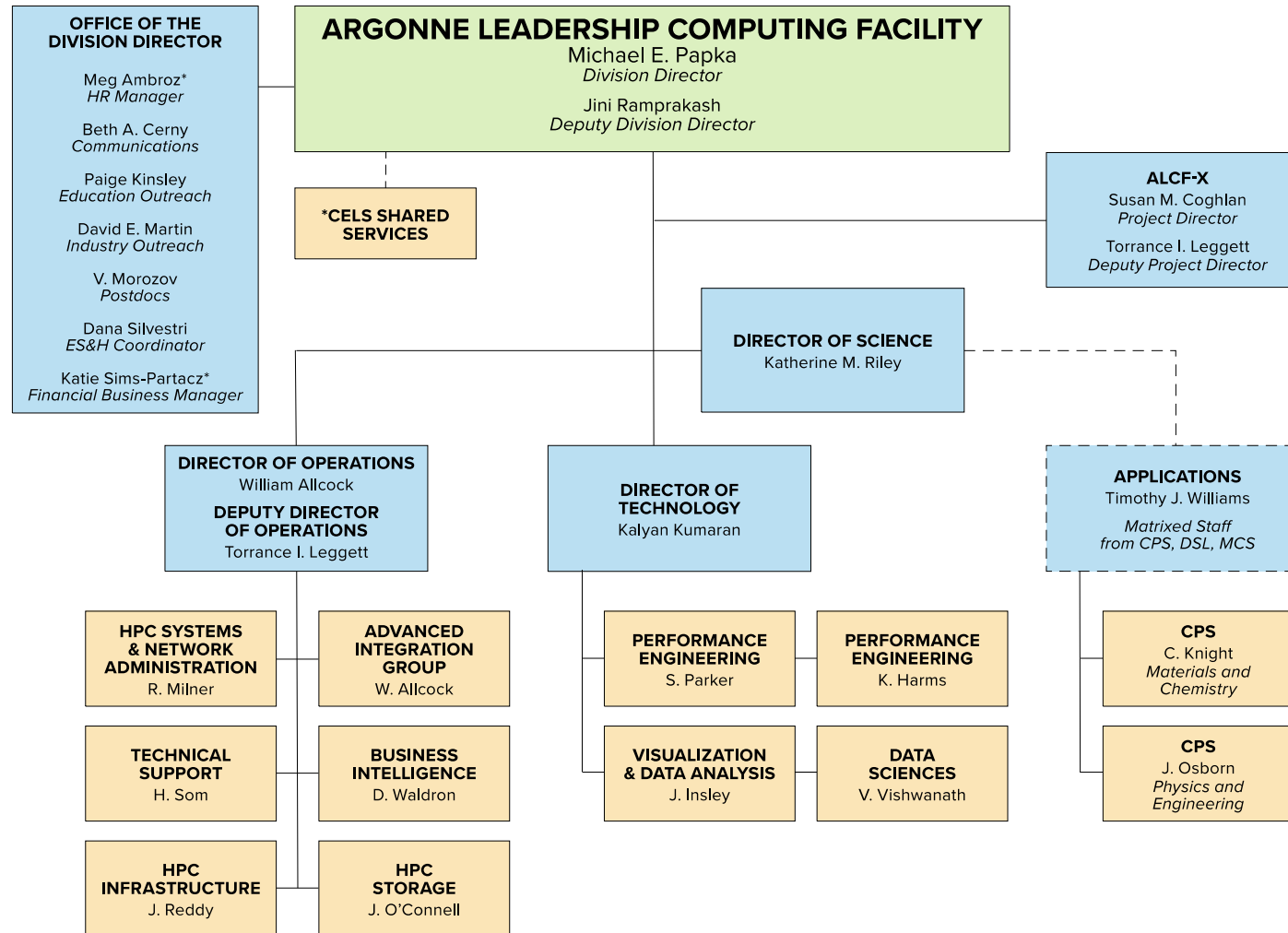
CONTACT US  
[Faculty Programs](#)  
[faculty@anl.gov](#)  
[Subscribe to Education Email List](#)

<https://www.anl.gov/education/faculty-programs>



Thank You

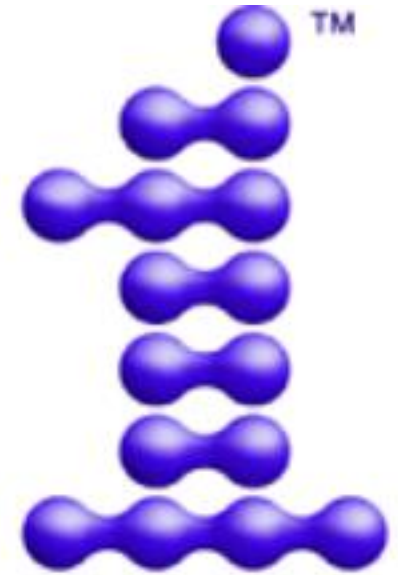
# ALCF Organizational Chart



April 11, 2022

# oneAPI

- Industry specification from Intel (<https://www.oneapi.com/spec/>)
  - Language and libraries to target programming across diverse architectures (DPC++, APIs, low level interface)
- Intel oneAPI products and toolkits (<https://software.intel.com/ONEAPI>)
  - Languages
    - Fortran (w/ OpenMP 5+)
    - C/C++ (w/ OpenMP 5+)
    - DPC++
    - Python
  - Libraries
    - oneAPI MKL (oneMKL)
    - oneAPI Deep Neural Network Library (oneDNN)
    - oneAPI Data Analytics Library (oneDAL)
    - MPI
  - Tools
    - Intel Advisor
    - Intel VTune
    - Intel Inspector



# oneAPI

<https://software.intel.com/oneapi>

# Available Aurora Programming Models

- Aurora applications may use:

- DPC++/SYCL
- OpenMP
- Kokkos
- Raja
- OpenCL



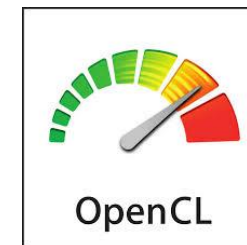
- Experimental

- HIP

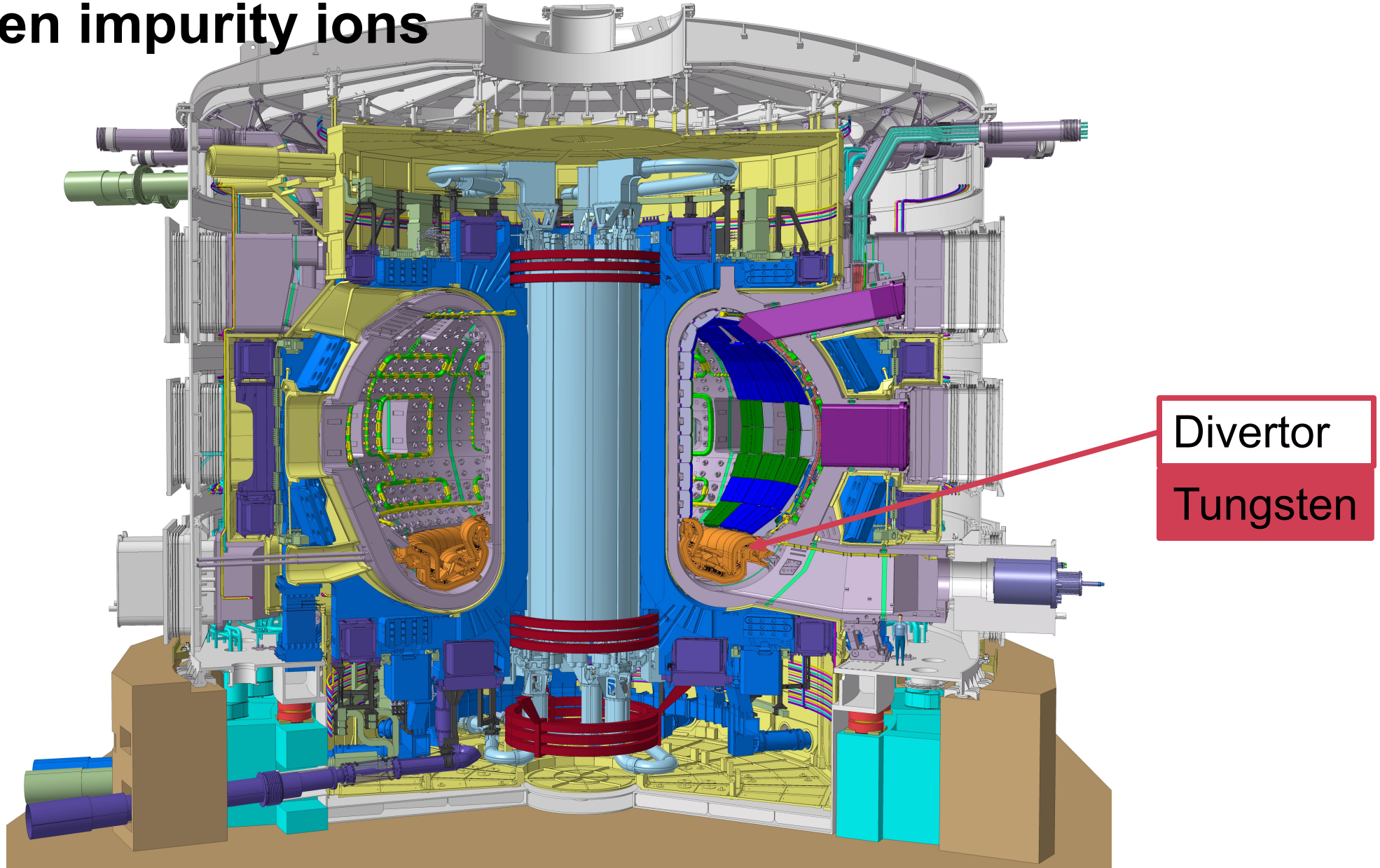


- Not available on Aurora:

- CUDA
- OpenACC



# ITER Tokamak - Predict ITER plasma behavior with Tungsten impurity ions



# Showcase

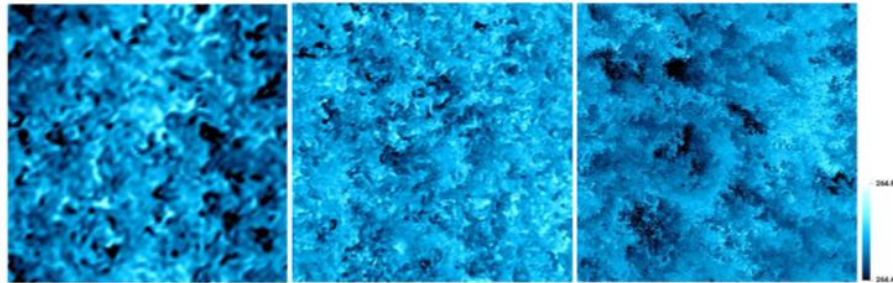
## ExaSMR: NekRS Performance on Ponte Vecchio

Ponte Vecchio with Intel OneAPI DPC++ implementation

1.5x performance lead

*ExaSMR: Small modular reactors (SMRs) and advanced reactor concepts (ARCs) will deliver clean, flexible, reliable, and affordable electricity while avoiding the traditional limitations of large nuclear reactor designs.*

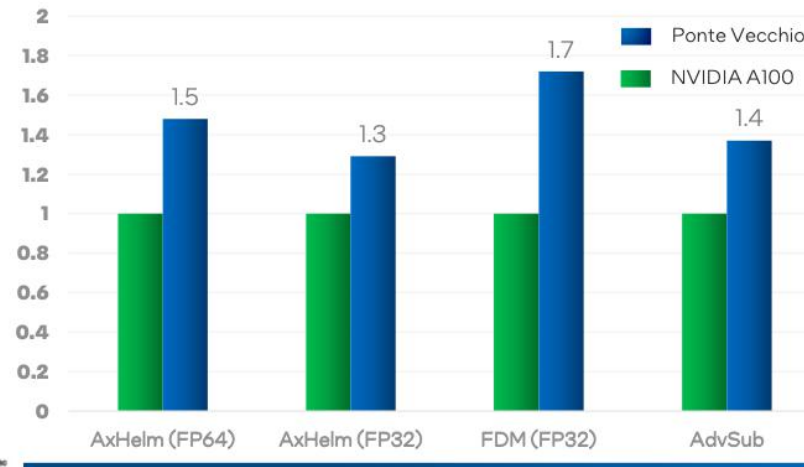
<https://www.exascaleproject.org/research-project/exasmr/>



**Figure 10:** NekRS: potential temperature distributions in [K] at time 6h and  $z=100\text{m}$  on different resolutions of  $\Delta x = 3.12\text{m}$  (left),  $1.56\text{m}$  (center), and  $0.78\text{m}$  (right) corresponding to the number of grid points,  $n=128^3$ ,  $256^3$ , and  $512^3$ , respectively.  $\Delta x$  represents the average grid-spacing for the spectral elements,  $E = 16^3$ ,  $32^3$  and  $64^3$  and the polynomial order  $N = 8$  on the domain  $400\text{m} \times 400\text{m} \times 400\text{m}$ .

<https://ceed.exascaleproject.org/docs/ceed-ms38-report.pdf>

Relative Performance of NekRS Benchmarks w/ problem size of 8196 (Averaged throughput, higher is better)



Application Summary:

NekRS is an open-source Navier Stokes solver based on the spectral element method targeting classical processors and accelerators like GPUs. The code started as a fork of libParanumal in 2019. For API portable programming OCCA is used.

<https://github.com/argonne-lcf/nekRS/>

OCCA is an open-source library which aims to make it easy to program different types of devices (e.g. CPU, GPU, FPGA). It provides a unified API for interacting with backend device APIs (e.g. OpenMP, CUDA, OpenCL), uses just-in-time compilation to build backend kernel, and provide a kernel language, a minor extension to C, to abstract programming for each backend.

<https://libocca.org>

H. Jiang, "Intel's Ponte Vecchio GPU Architecture, Systems & Software," 2022 IEEE Hot Chips 34 Symposium (HCS), 2022, pp. 1-29, doi: 10.1109/HCS55958.2022.9895631.

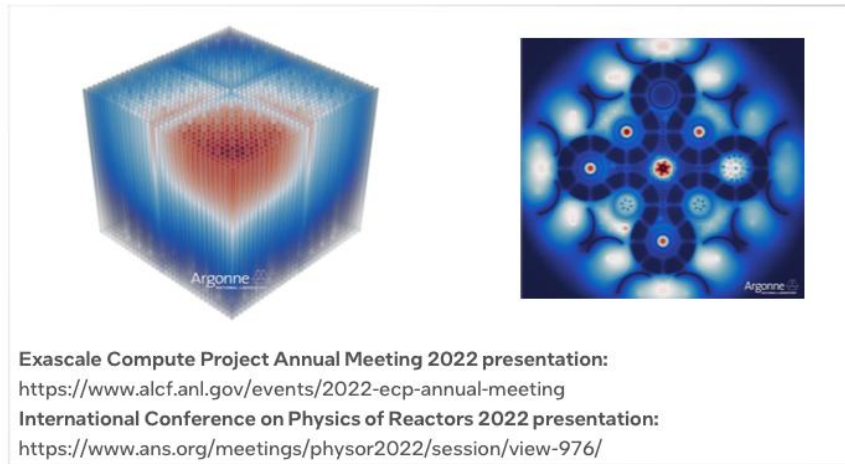
# Showcase

## ExaSMR: OpenMC Performance on Ponte Vecchio

Monte Carlo particle transport code for exascale computations

Ponte Vecchio with OpenMP Target offload

2x performance lead



Exascale Compute Project Annual Meeting 2022 presentation:

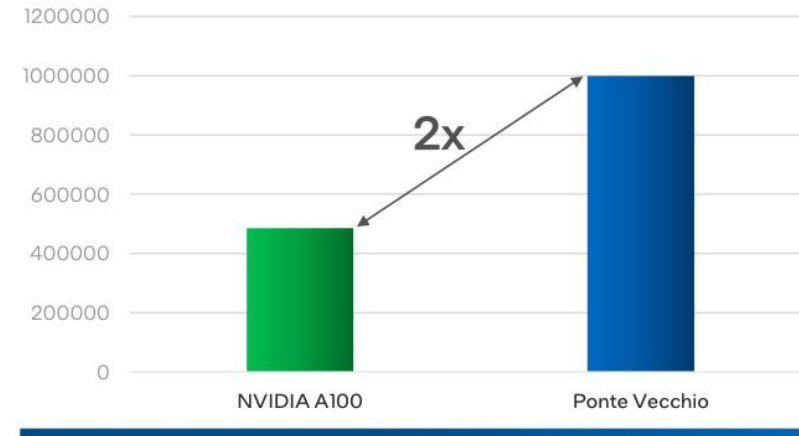
<https://www.alcf.anl.gov/events/2022-ecp-annual-meeting>

International Conference on Physics of Reactors 2022 presentation:

<https://www.ans.org/meetings/physor2022/session/view-976/>

<https://docs.openmc.org>

OpenMC Depleted Fuel Inactive Batch Performance on HM-Large Reactor with 40M particles (particles/second, Higher is better)



**Application Summary:** OpenMC is a Monte Carlo particle transport application that has recently been ported to the OpenMP target offloading programming model for use on GPU-based systems. The Monte Carlo method employed by OpenMC is considered the "gold standard" for high-fidelity simulation while also having the advantage of being a general-purpose method able to simulate nearly any geometry or material without the need for domain-specific assumptions. However, despite the extreme advantages in ease of use and accuracy, Monte Carlo methods like those in OpenMC often suffer from a very high computational cost. The extreme performance gains OpenMC has achieved on GPUs, as compared to traditional CPU architectures, is finally bringing within reach a much larger class of problems that historically were deemed too expensive to simulate using Monte Carlo methods. The leap in performance that GPUs are now offering carries with it the potential to disrupt a number of engineering technology stacks that have traditionally been dominated by non-general deterministic methods. For instance, faster MC applications may greatly expand the design space and simplify the regulation process for new nuclear reactor designs – potentially improving the economics of nuclear energy and therefore helping to solve the world's climate crisis.



- See backup for workloads and configurations. Results may vary.
- Intel does not endorse or warrant the use of its products for any specific application.

intel.

23

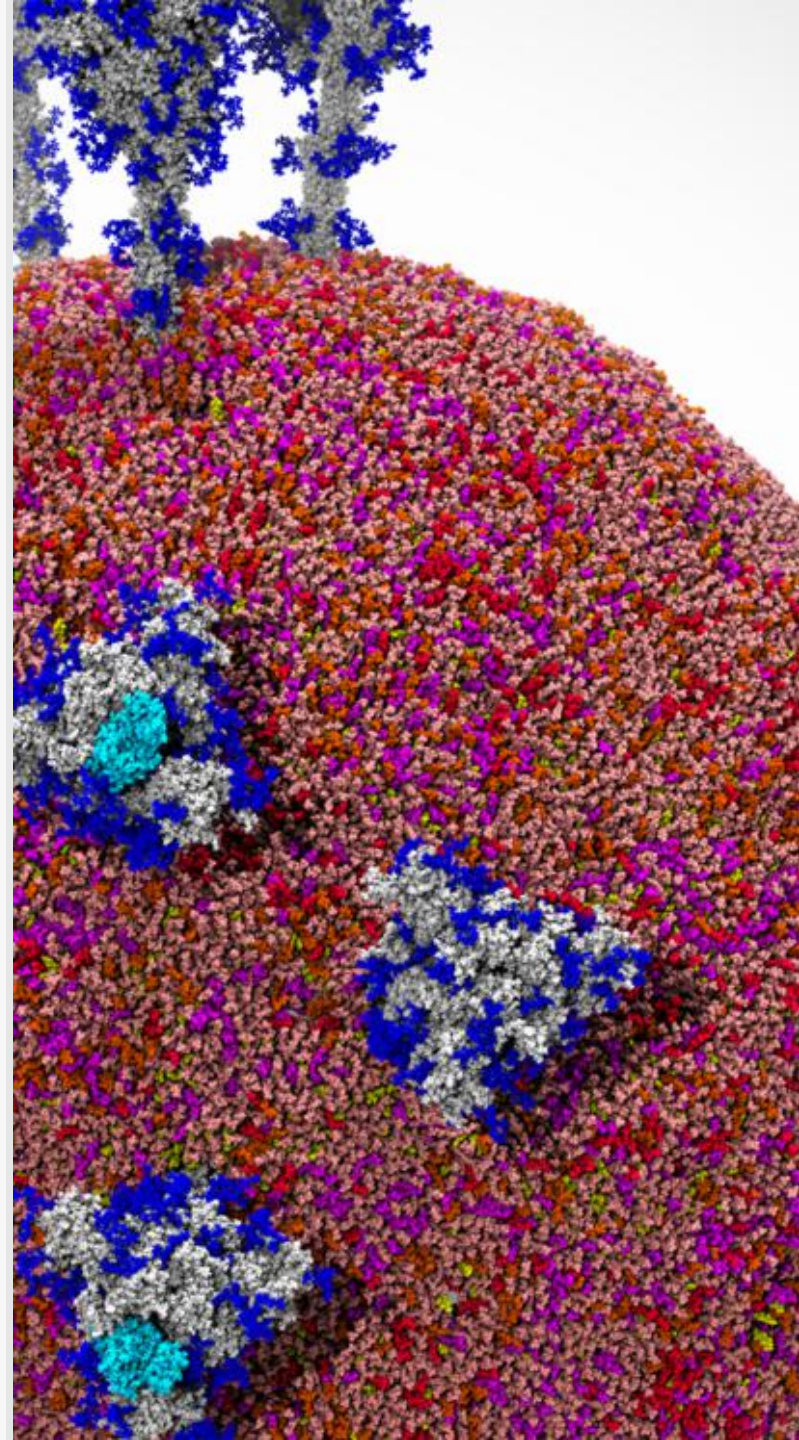
H. Jiang, "Intel® OneAPI Architecture and System Software, 2022 Intel® HCS Symposium (HCS), 2022, pp. 1-29, doi: 10.1109/HCS55958.2022.9895631.

# AI-Driven Drug Discovery for SARS-CoV-2 Proteome

- **PI: Arvind Ramanathan, Argonne National Laboratory**

---

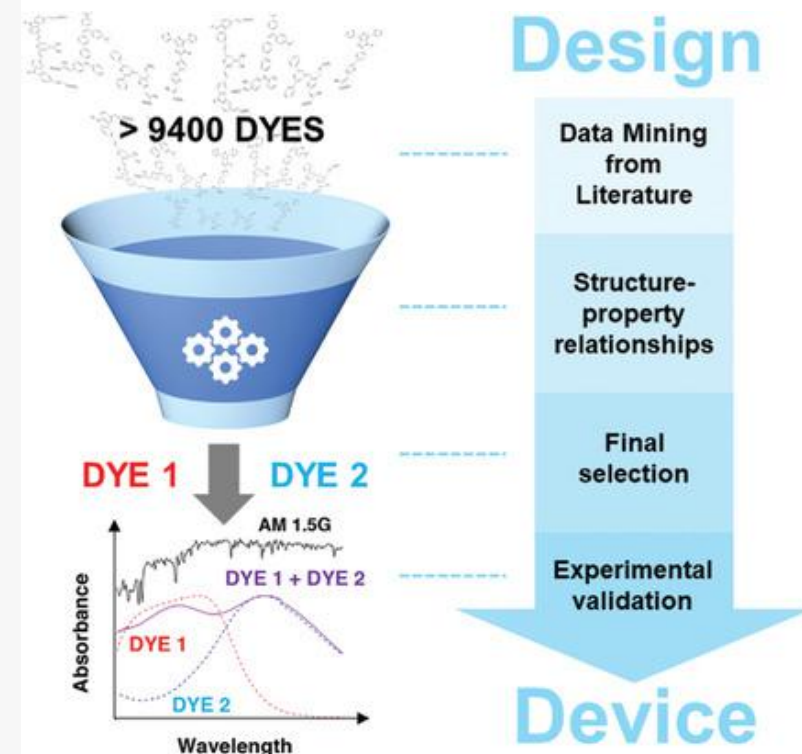
- **Science Summary:**  
Using AI techniques to screen over 6 million small molecules, researchers identified at least 20 partially active molecules that can potentially inhibit viral function in wet lab experiments.
- The 20 candidates are being validated in labs for activity against the virus. The work also generated new models for any small molecules and antibodies. Data are publicly available.
- **Impact:** This research will aid in the design of antibodies for the virus.





# Data-Driven Design of Solar Cells

- **PI:** Jacqueline Cole, University of Cambridge
- **Science Summary:** Light-absorbing dyes are promising, low-cost materials for organic dye-sensitized solar cells that can passively supply energy through tinted windows. Using data mining, machine learning, and computational modeling techniques, researchers identified two high-performing dyes, and then produced a solar cell competitive with common industry materials.
- **Impact:** The team's use of data mining, in conjunction with large-scale simulations and experiments, offers a novel approach to advance the design and discovery of new functional materials. In addition, the project's development of open-source databases and data-extraction software tools will help accelerate materials discoveries by removing the hurdle of manual database creation.

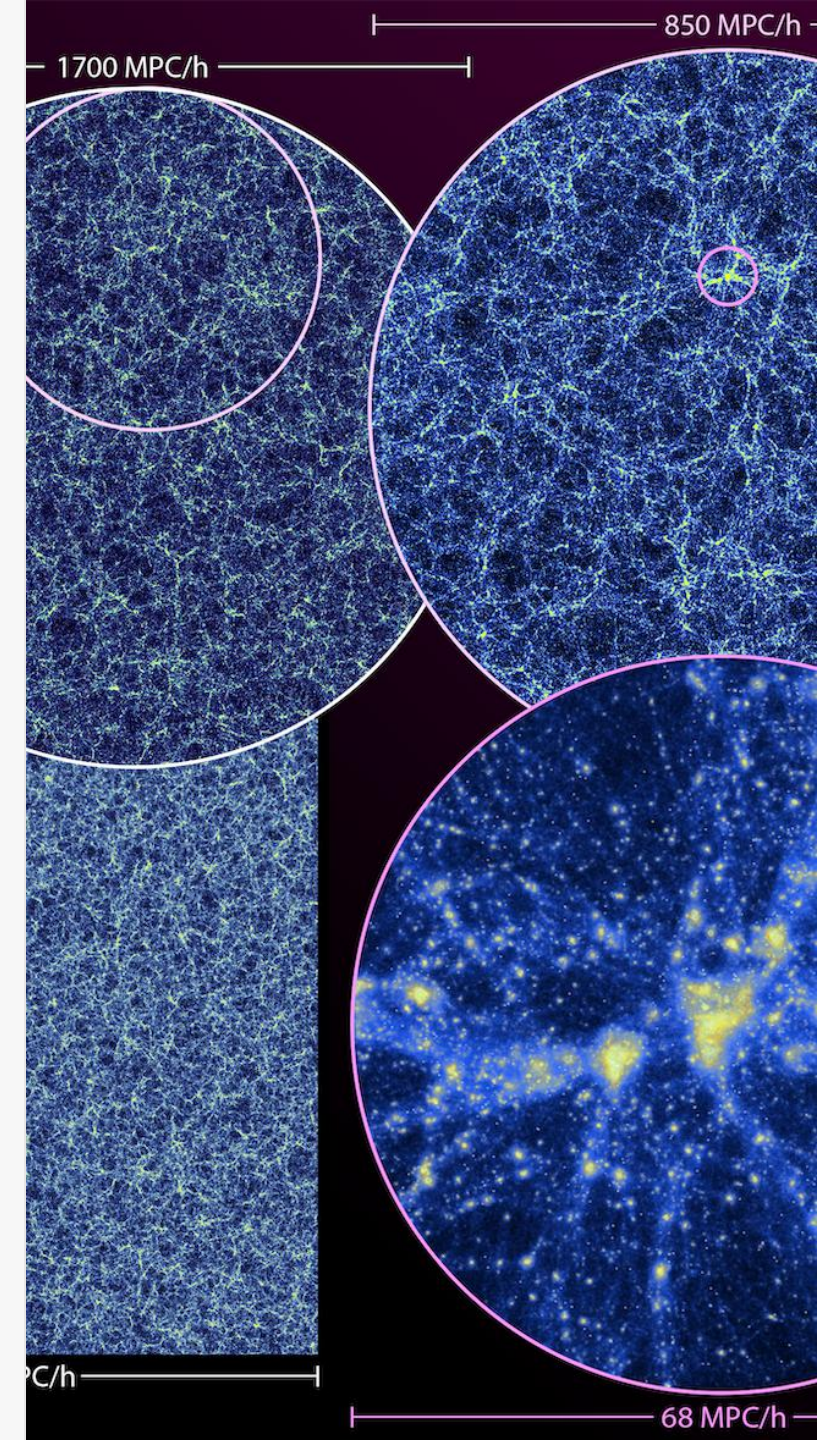


# The Last Journey

- **PI:** Katrin Heitmann, Argonne National Laboratory

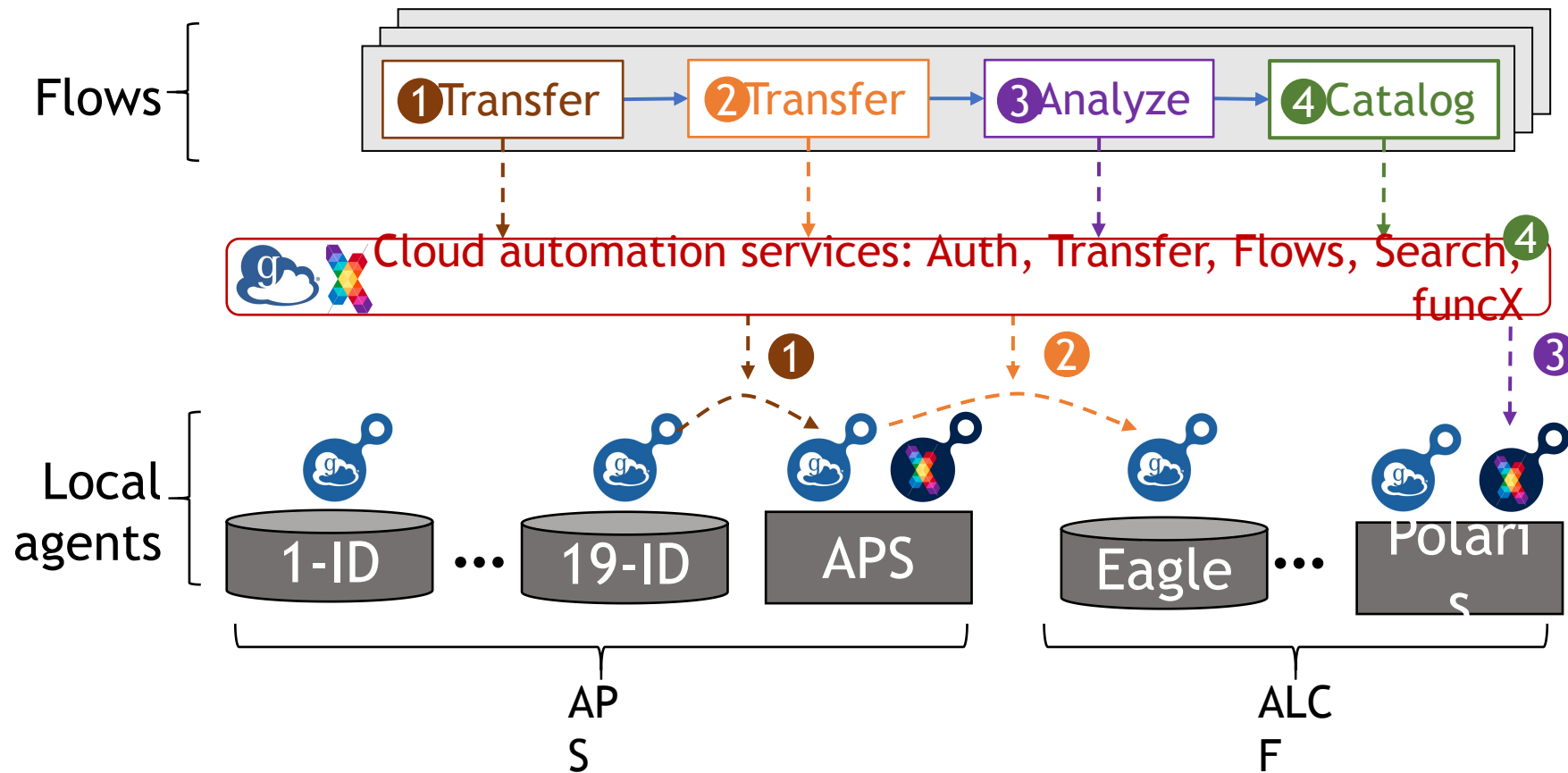
---

- **Science Summary:** In the final months of Mira's operation, researchers ran one of the largest cosmological simulations using cutting-edge observational advances from satellites and telescopes. Evolving a massive number of particles, the simulation was designed to help resolve mysteries of dark energy and dark matter. Results will form the basis for sky maps used by numerous surveys.
- **Impact:** The team's simulation was designed to address numerous fundamental questions in cosmology; the data produced are essential for enabling the refinement of existing predictive tools and aid the development of new models. Their research will impact both ongoing and upcoming cosmological surveys, including the Dark Energy Spectroscopic Instrument (DESI), the LSST, SPHEREx, and the "Stage-4" ground-based cosmic microwave background experiment (CMB-S4).



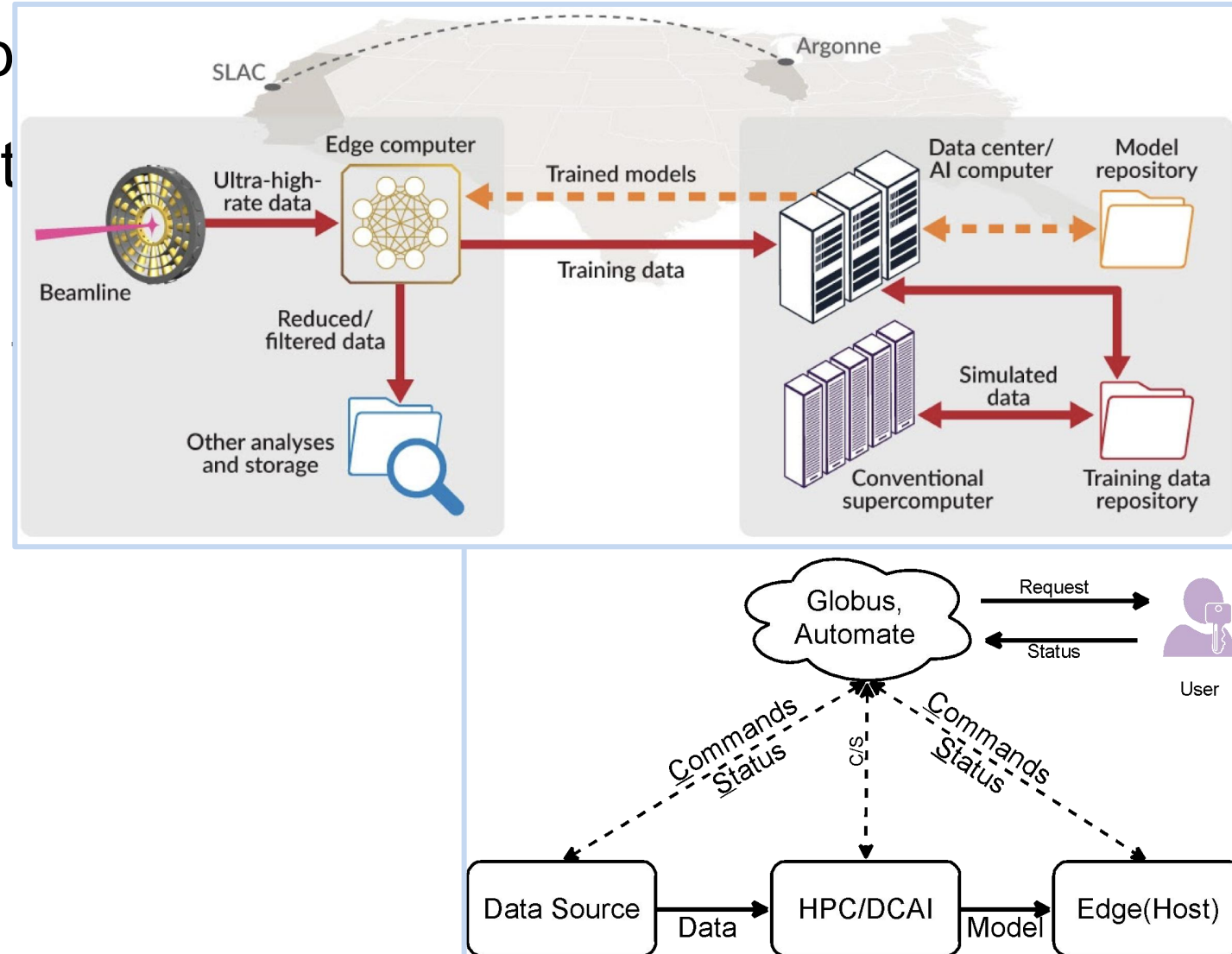
# Toward Coupling ALCF with Experimental Sciences

- On-demand computing facilitated by pre-emption and flexible resource provisioning, data services capabilities, together with 200Gbps+ external network connectivity, among others, enable real-time coupling at ALCF with experimental science



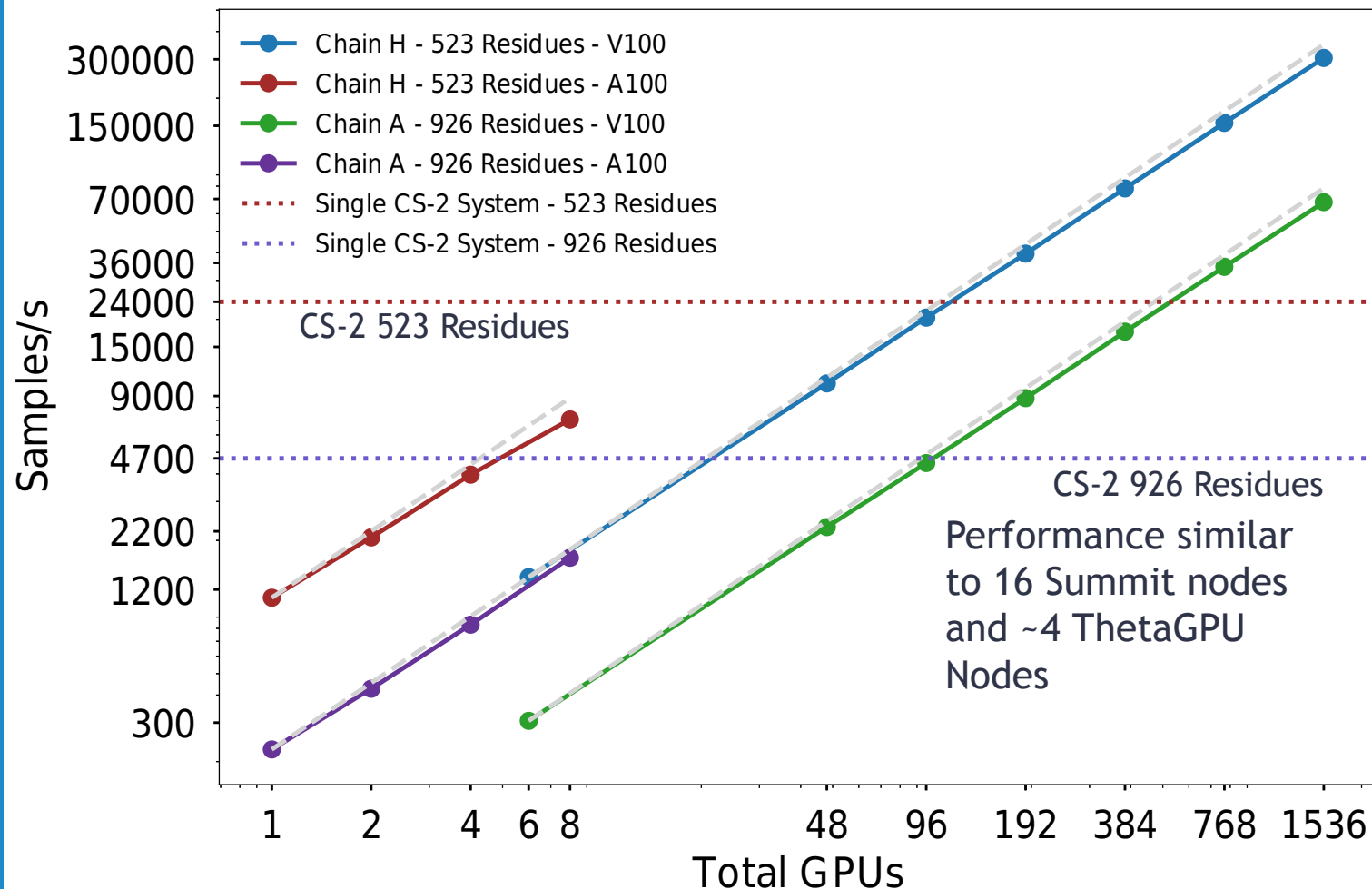
# Example: Rapid Training of Deep Neural Networks using Remote Resources

- DNN at the edge for fast p
- Requires tight coupling with data
- Near real-time steering of interest



See Zhengchun's talk next

# COVID-19 CVAE Training on Summit and Cerebras CS-2



- Single CS+0 delivers performance of over / . . GPUs on CVAE
- Results are for **out-of-the-box performance** based on model config not optimized for CS+0 ,

Performance	523 X 523	926 X 926
<b>Throughput (samples/sec)</b>		
/x CS+0 System	02* . . .	25 . .
/x V/ . . GPU	006	3/
/x A/ . . GPU	~// . .	~/3 .
<b>Speedup (CS2 vs. GPU )</b>		
/ x V/ . . GPU	//1x	/ . /x
/ x A/ . . GPU	~00X	~10X

*Intelligent Resolution: Integrating Cryo-EM with AI-driven Multi-resolution Simulations to Observe the SARS-CoV-2 Replication-Transcription Machinery in Action, SC21 COVID19 Gordon Bell Finalist, In IJHPCA 2022*

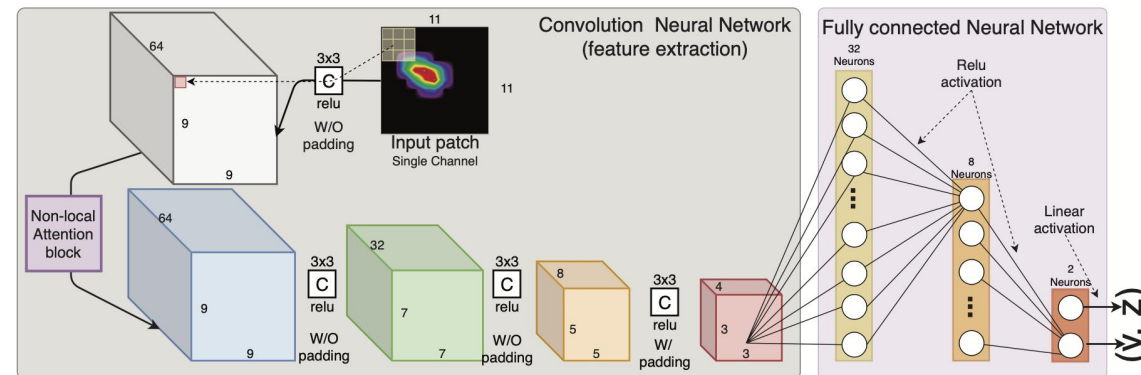
<https://www.biorxiv.org/content/10.1101/2021.10.09.463779v1.full.pdf>

# Fast X-Ray Bragg Peak Analysis

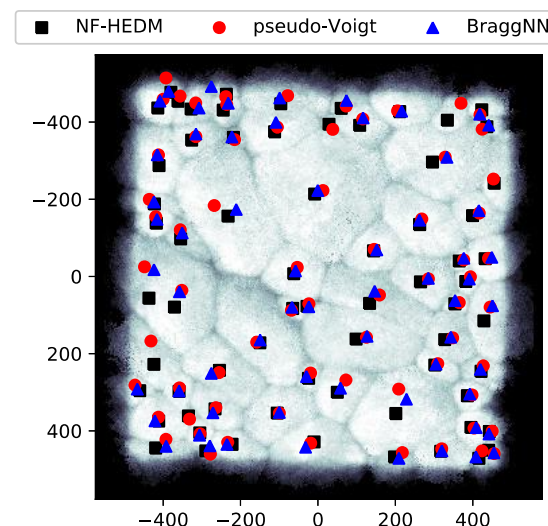
**Goal:** Enable rapid analysis and real-time feedback during an in-situ experiment with complex detector technologies

**Proposed Approach:** Deep learning-based method, BraggNN, for massive extraction of precise Bragg peak locations from far-field high energy diffraction microscopy data. BraggNN has achieved 200X improvement over conventional pseudo-Voigt profiling

**Challenges:** Model training capability is limited by the hardware



Application of the BraggNN deep neural network to an input patch yields a peak center position  $(y, z)$ . All convolutions are 2D of size  $3 \times 3$ , with rectifier as activation function. Each fully connected layer, except for the output layer, also has a rectifier activation function.

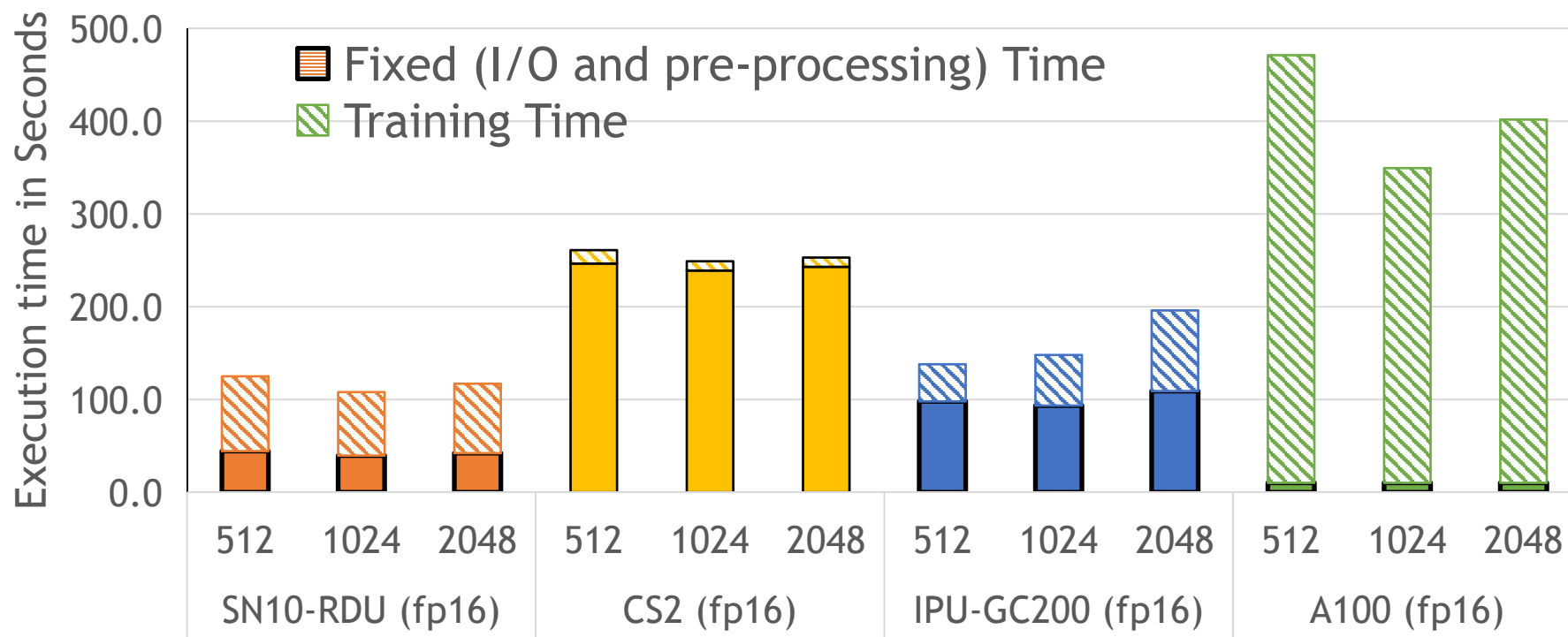


A comparison of BraggNN, pseudo-Voigt FF-HEDM and NF-HEDM. (a) Grain positions from NF-HEDM (black squares), pseudo-Voigt FF-HEDM (red circles) and BraggNN FF-HEDM (blue triangles) overlaid on NF-HEDM confidence map

Courtesy: Z. Liu et al. BraggNN: Fast X-ray Bragg Peak Analysis Using Deep Learning. International Union of Crystallography (IUCrJ), Vol. 9, No. 1, 2022

# Fast X-Ray Bragg Peak Analysis

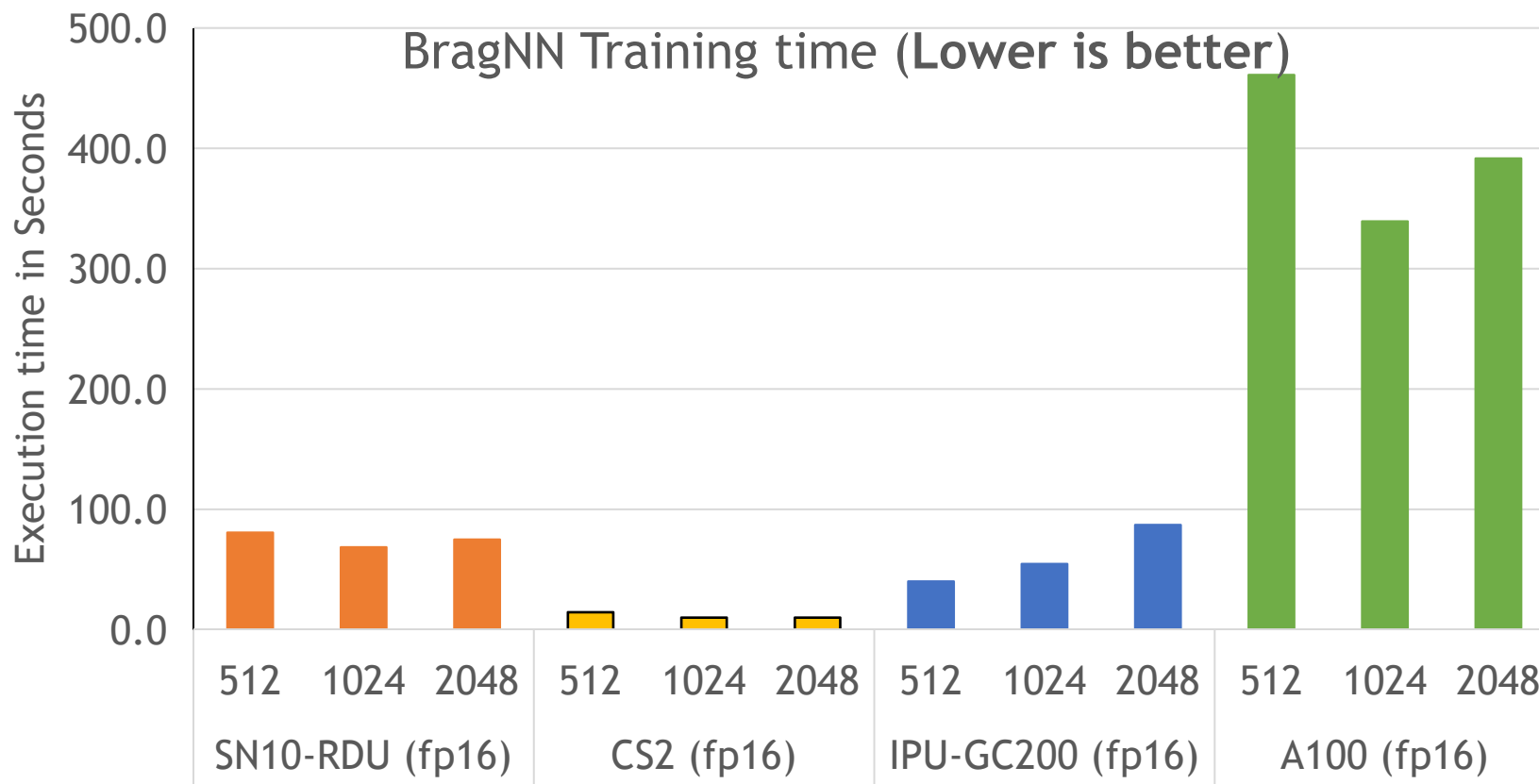
BragNN End-to-End execution time (Lower is better)



SambaNova and Graphcore achieve lowest time to solution and achieve up to 3.7X to 3.4X speedup in comparison to Nvidia A100 respectively. Cerebras achieves up to 80% improvement over A100

“A comprehensive evaluation of Novel AI accelerators for Deep Learning Workloads”, M. Emani et. al, To appear at PMBS workshop SC'22

# Fast X-Ray Bragg Peak Analysis



For training time, we ignore the data loading and pre-processing time (Fixed cost time). Cerebras CS2 achieves up to 33X improvement over A100 while SN and Graphcore achieve up to 6-11X improvement over A100 respectively for training. Note: Cerebras performance includes use of multi-replica optimization.