



Aug 5, 2025

# Advancing Scalable and Trustworthy Foundational Models for Science at Oak Ridge National Laboratory

---

Prasanna Balaprakash



U.S. DEPARTMENT OF  
**ENERGY**

ORNL IS MANAGED BY UT-BATTELLE LLC  
FOR THE US DEPARTMENT OF ENERGY



# ORNL has a rich history leveraging AI for science



**1979**  
Oak Ridge  
Applied Artificial  
Intelligence  
Project



**1991**  
Automated  
machines



**Current**  
Frontier

- #2 HPL-MxP @10 exaflops for AI
- Scaled to 1T+ parameter AI model training

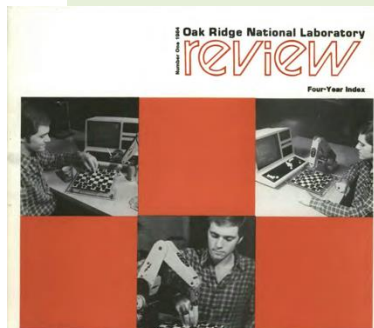
1980

1990

2000



**1981**  
AI infrastructure  
supports  
spectroscopy,  
environmental  
management,  
nuclear fuel  
reprocessing,  
and programming  
assistance



**1983**  
Robotics



**2012**  
Titan:  
First GPU-powered  
supercomputer



# AI transforming science and national security

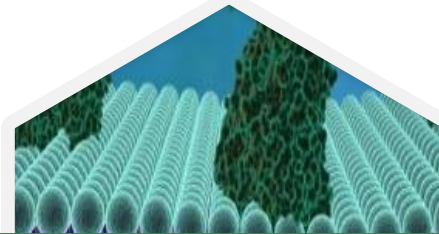
## ORNL facilities, expertise enable AI revolution



**Spallation  
Neutron Source**



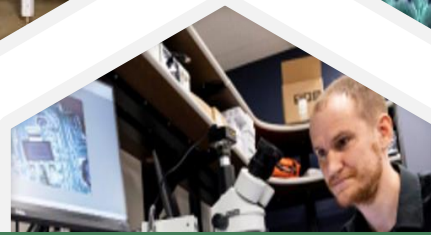
**Manufacturing  
Demonstration Facility**



**Center for Structural  
Molecular Biology**



**Oak Ridge Leadership  
Computing Facility**



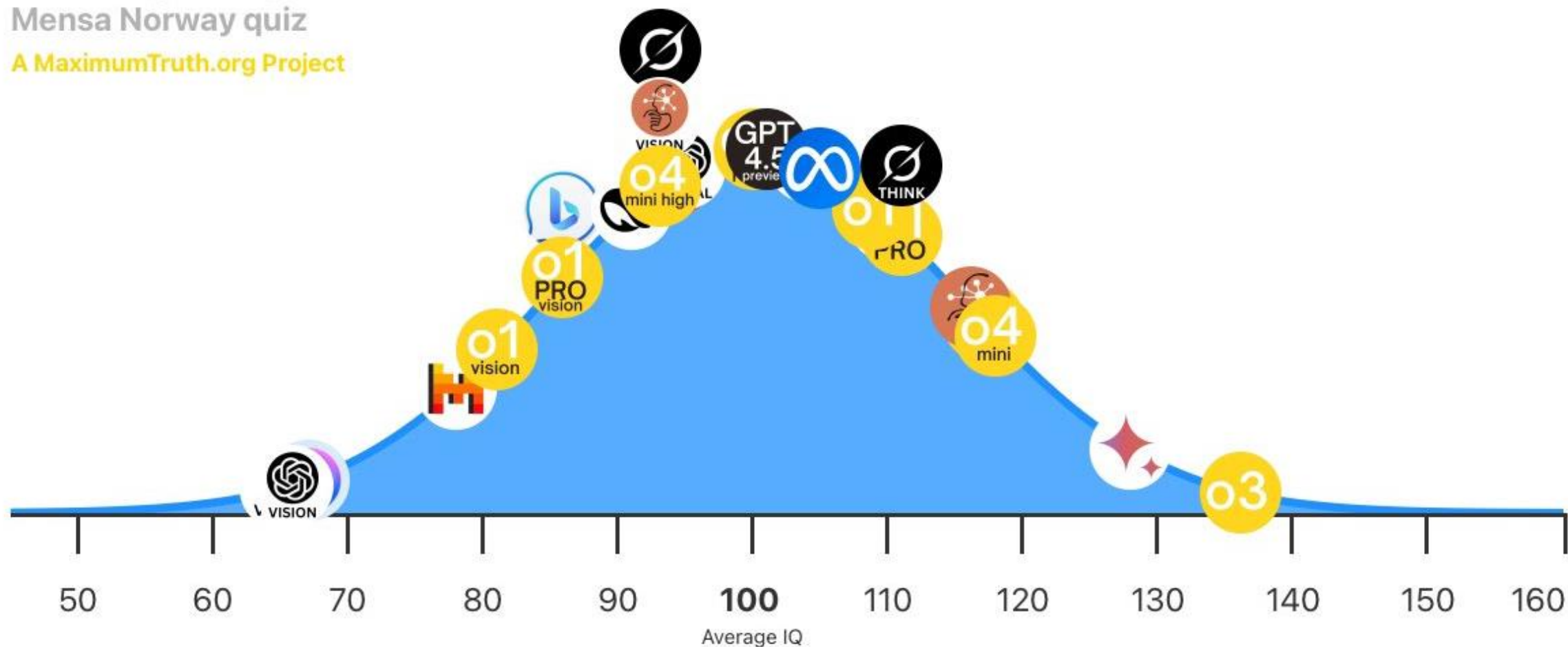
**Cyber Science  
Research Facility**



**High Flux  
Isotope Reactor**

# What is industry doing?

Mensa Norway quiz  
A MaximumTruth.org Project



# The Gentle Singularity — Sam Altman, 2025

- **The Takeoff Has Begun**
  - Humanity has crossed the AI event horizon.
  - AI systems (e.g. GPT-4, o3) already outperform humans in key areas.
  - Scientific breakthroughs behind them were the hardest part—momentum is now self-reinforcing.
- **AI's Transformative Impact**
  - 2025: Agents doing real cognitive work (coding, scientific discovery).
  - By 2027: Robots may enter the physical world.
  - Productivity, creativity, and research velocity are surging—scientists report 2–3X output boosts.
- **Self-Reinforcing Progress Loops**
  - Recursive AI research: AI used to build better AI.
  - Economic flywheel: AI → Value → Infrastructure → More AI.
  - Automation of datacenters, supply chains, and robot manufacturing looms ahead.
- **Toward Abundant Intelligence & Energy**
  - Intelligence cost converging to energy cost.
  - Exponential gains may redefine what we consider “real work” or “progress.”
  - Wonders become routine—and then table stakes.





Zuckerberg said the company plans on investing "hundreds of billions of dollars" to power AI



PRESS RELEASE • NUCLEAR POWER

## Tennessee Valley Authority submits application for construction of first BWRX-300 small modular reactor in the U.S.

• 5 min read

**KNOXVILLE, Tenn. (May 20, 2025)** – Tennessee Valley Authority (TVA) has submitted an application to the U.S. Nuclear Regulatory Commission to construct a GE Vernova Hitachi Nuclear Energy (GVH) BWRX-300 small modular reactor (SMR) at the Clinch River site in Oak Ridge, Tennessee. It is the first construction permit application for a BWRX-300 in the U.S.

[Home](#) > [Topics](#) > [Unleash American Energy Innovation](#) > DOE Announces Site Selection for AI Data Center and Energy Infrastructure Development on Federal Lands

## DOE Announces Site Selection for AI Data Center and Energy Infrastructure Development on Federal Lands

The forthcoming solicitations will drive innovation in reliable energy technologies, contribute to lower energy costs, and strengthen American leadership in artificial intelligence.

[Energy.gov](#)

July 24, 2025

# Scientific AI demand richer world models



## LLMs Today

- Trained on  $\sim 2.0E13$  text tokens ( $\sim 6.0E13$  bytes) on static data
- Reading equivalent would take a human 300,000 year
- Text input alone often lacks physical grounding

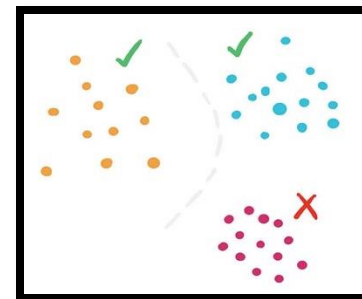
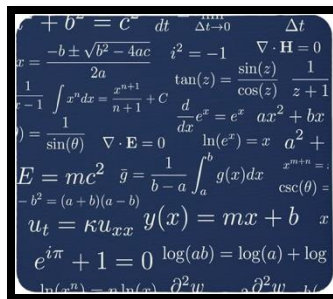
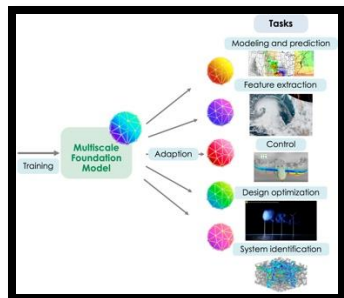
## 4 years old child

- 16K waking hours in 4 years
- $1.1E14$  bytes of multimodal, real-time data
- Vision, touch, language, causality

***Even a child's learning is vastly more multimodal than today's LLMs.  
Science demands even more!!***



# How can we capitalize?



Exascale and  
AI Integration

Multimodal  
AI Models

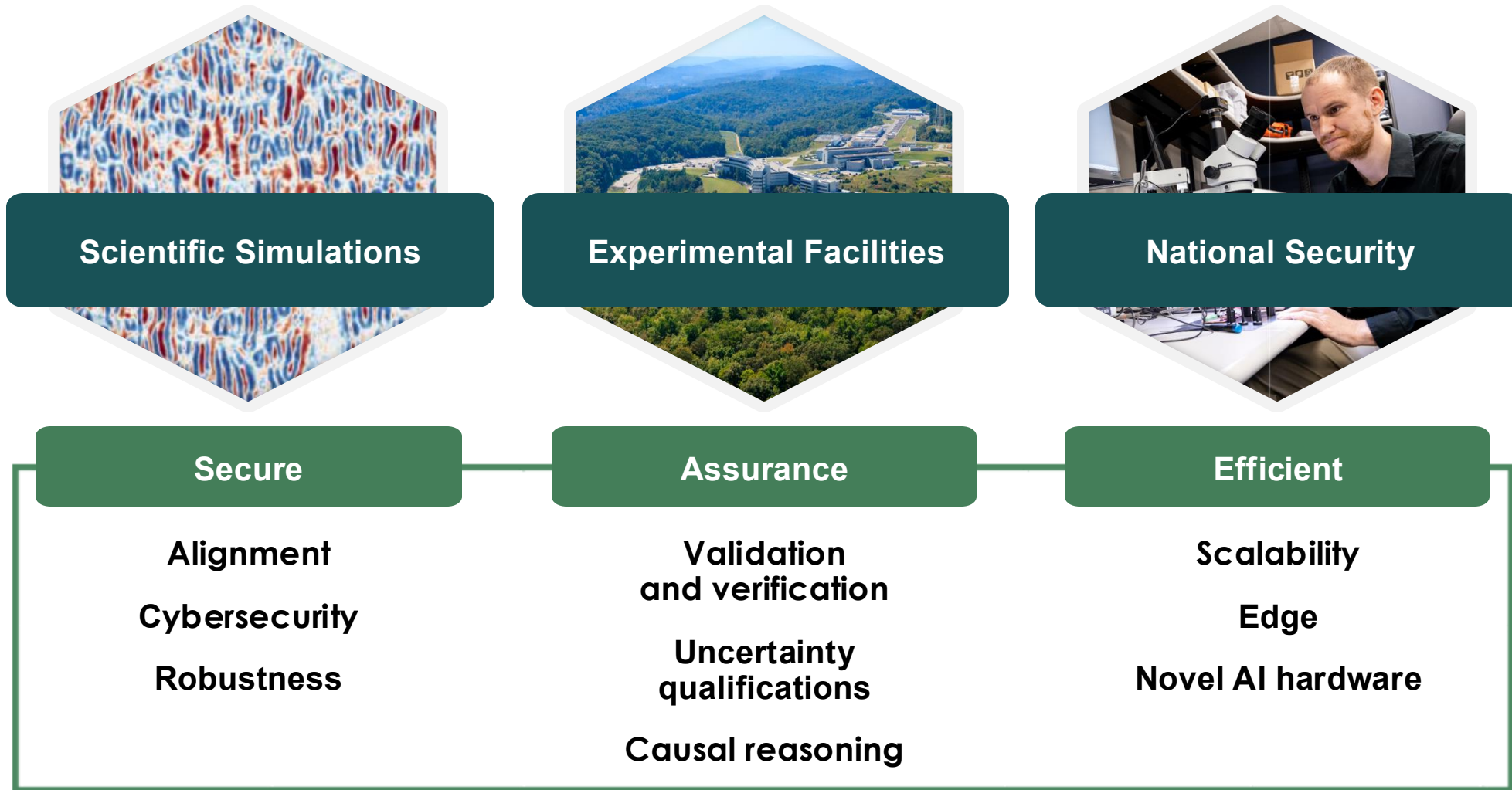
Domain-  
Knowledge

Robust  
Validation  
Frameworks

Secure,  
Usable AI  
Infrastructure

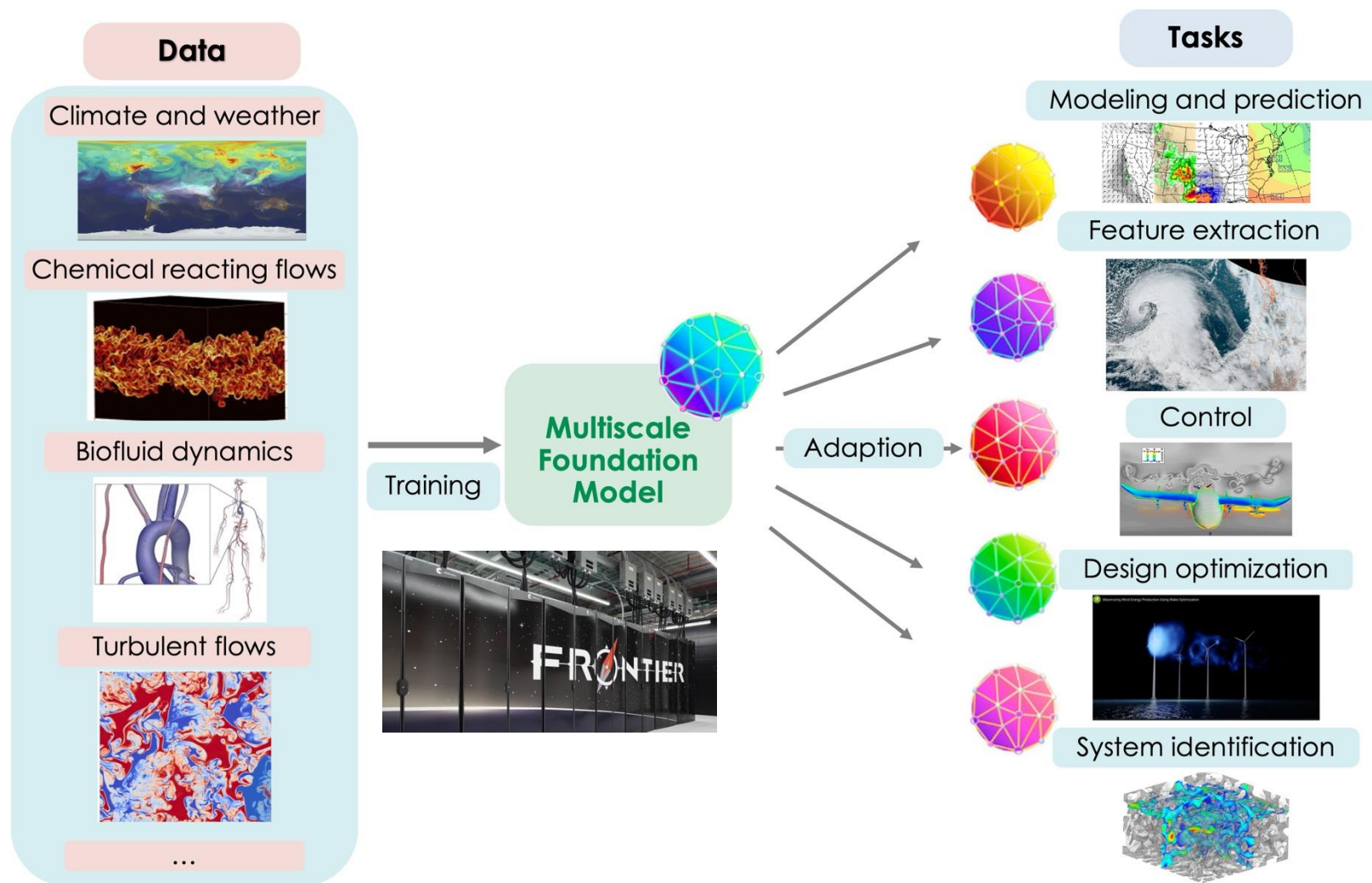
# ORNL's AI initiative

## Secure, assured, and efficient



*The initiative's portfolio comprises 15 advanced AI projects and involves over 50 researchers from 5 different directorates across the lab.*

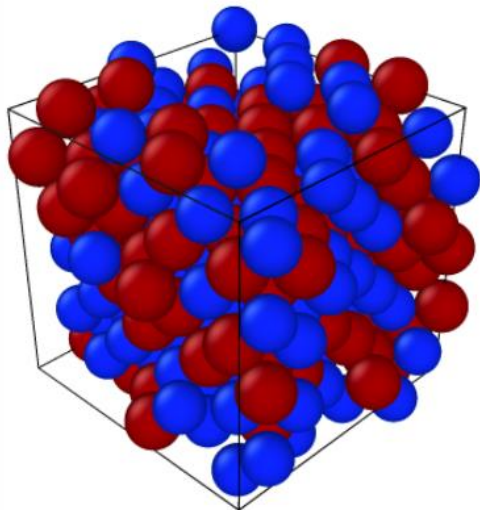
# Foundation AI model(s) for science





# Graph representation of materials at different scales

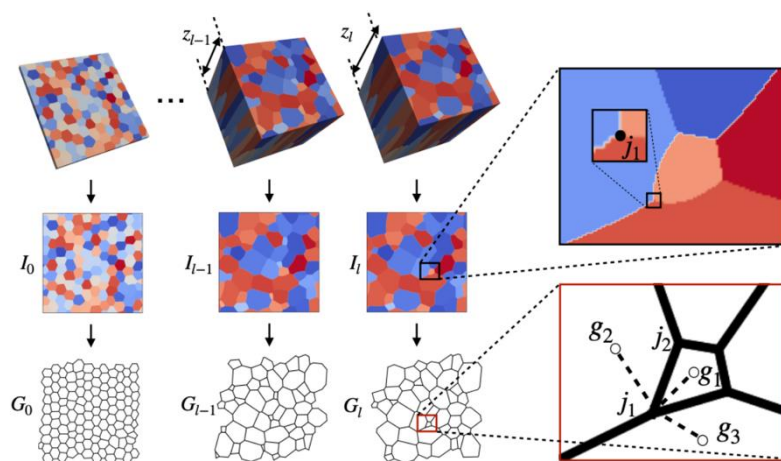
**Atomistic scale**



Nodes = atoms

Edges = interatomic bonds

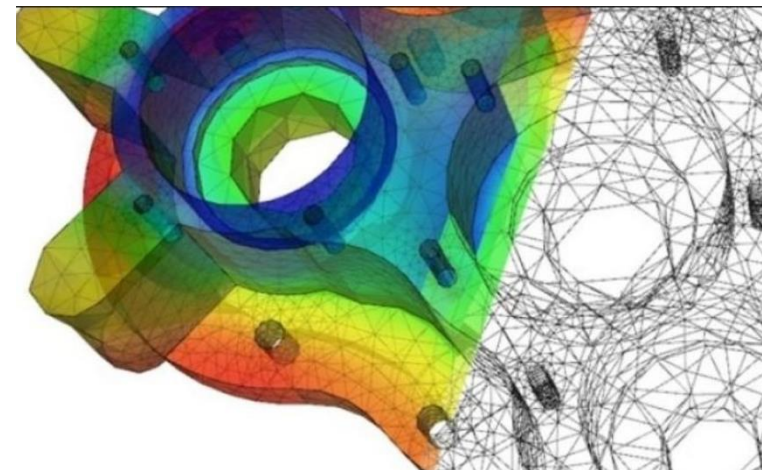
**Mesoscale**



Nodes = Voronoi centers

Edges = connection  
between Voronoi centers

**Continuum scale**

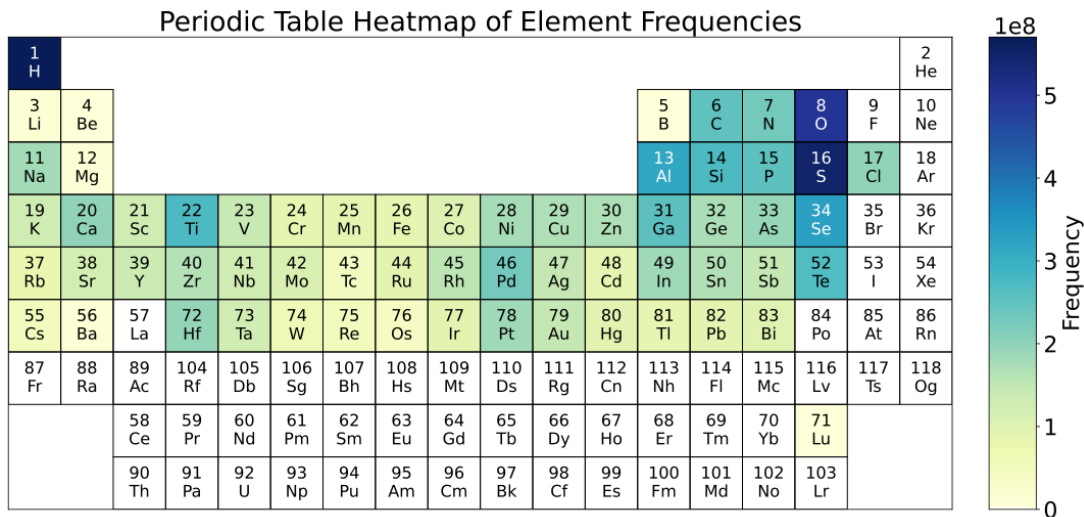


Nodes = vertices of the finite  
element mesh

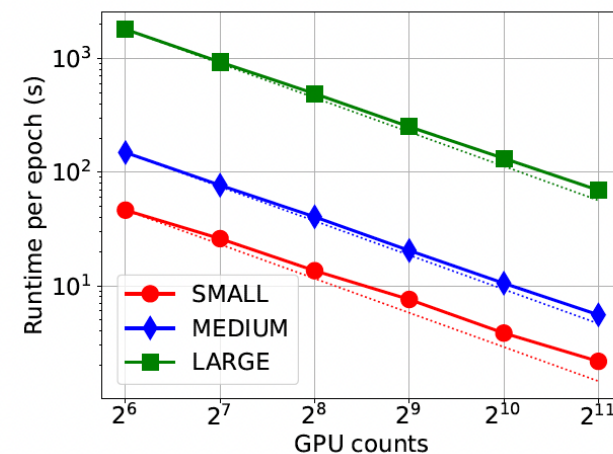
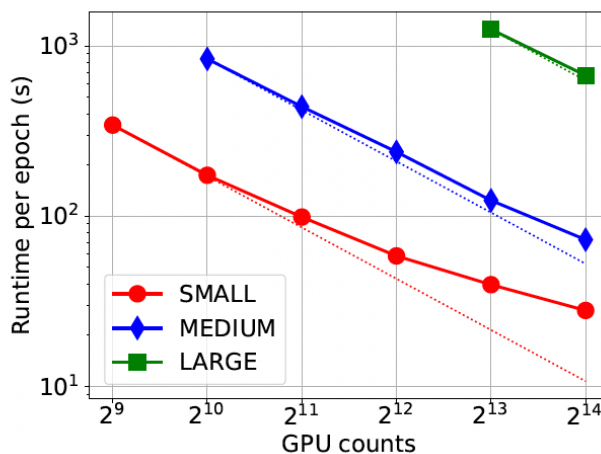
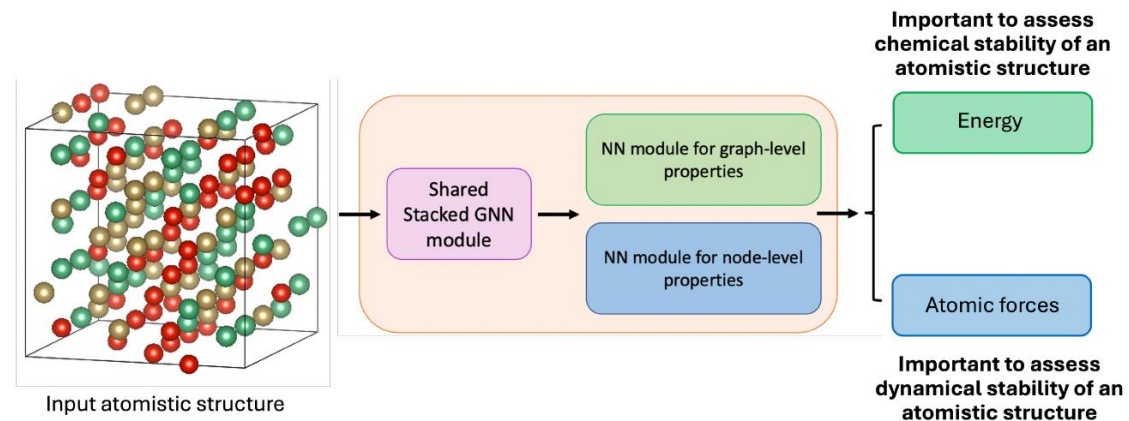
Edges = edges of the finite  
element mesh

Graph structured data maps naturally onto graph neural networks

# Scalable training of graph foundation model for materials science

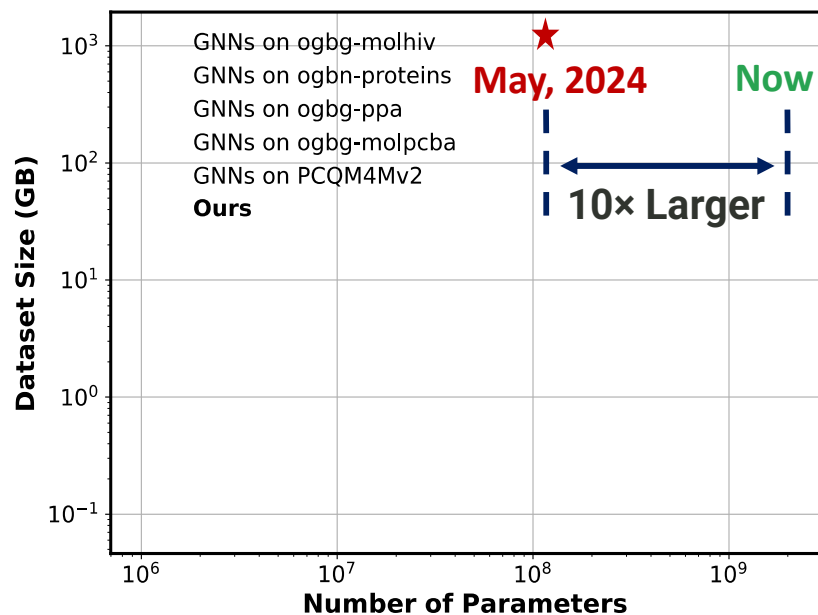


Dataset	Number of data samples	Size
ANI1x [63]	4,956,005	5.3 GB
QM7-X [64]	4,195,237	23 GB
OC2020 [39]	134,929,018	4.3 TB
OC2022 [40]	8,847,031	648 GB
MPTrj [37]	1,580,395	17 GB
Total	154,507,686	5.2 TB

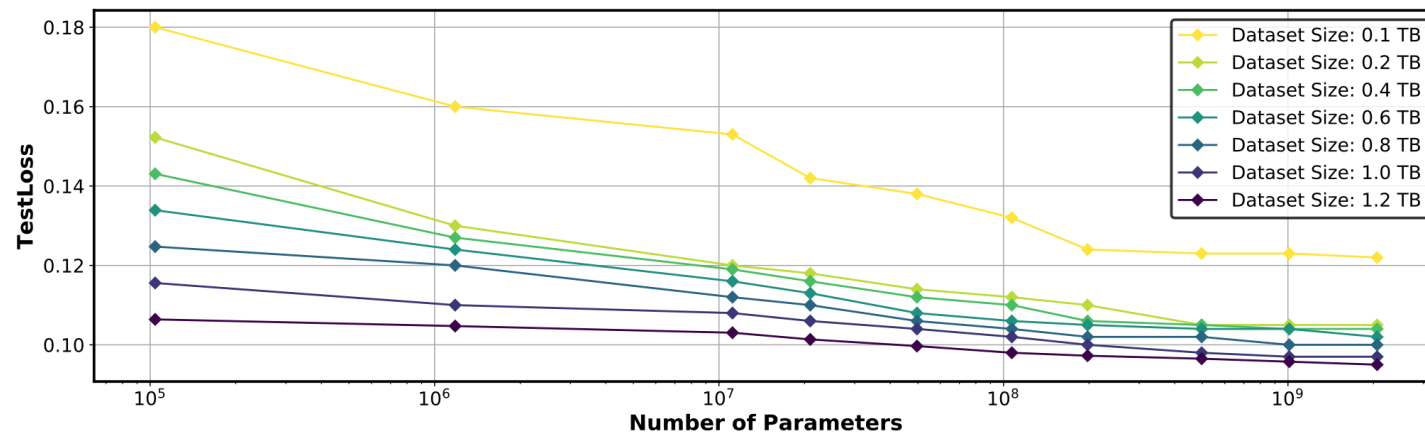


Strong scaling of HydraGNN multitasking pre-training on a problem of 120 million graphs on Frontier 16K GPUs and 2 million graphs on Perlmutter 2K GPUs with three GNN model sizes.

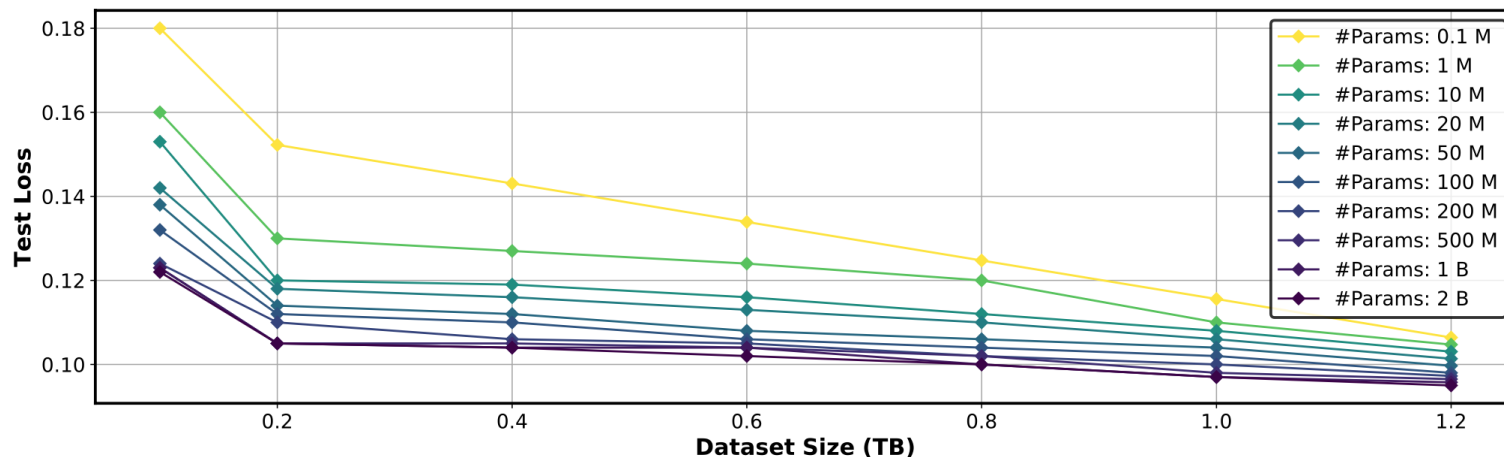
# Scaling laws for GNNs



Scaling law of GNN accuracy as a function of model size



Scaling law of GNN accuracy as a function of data size

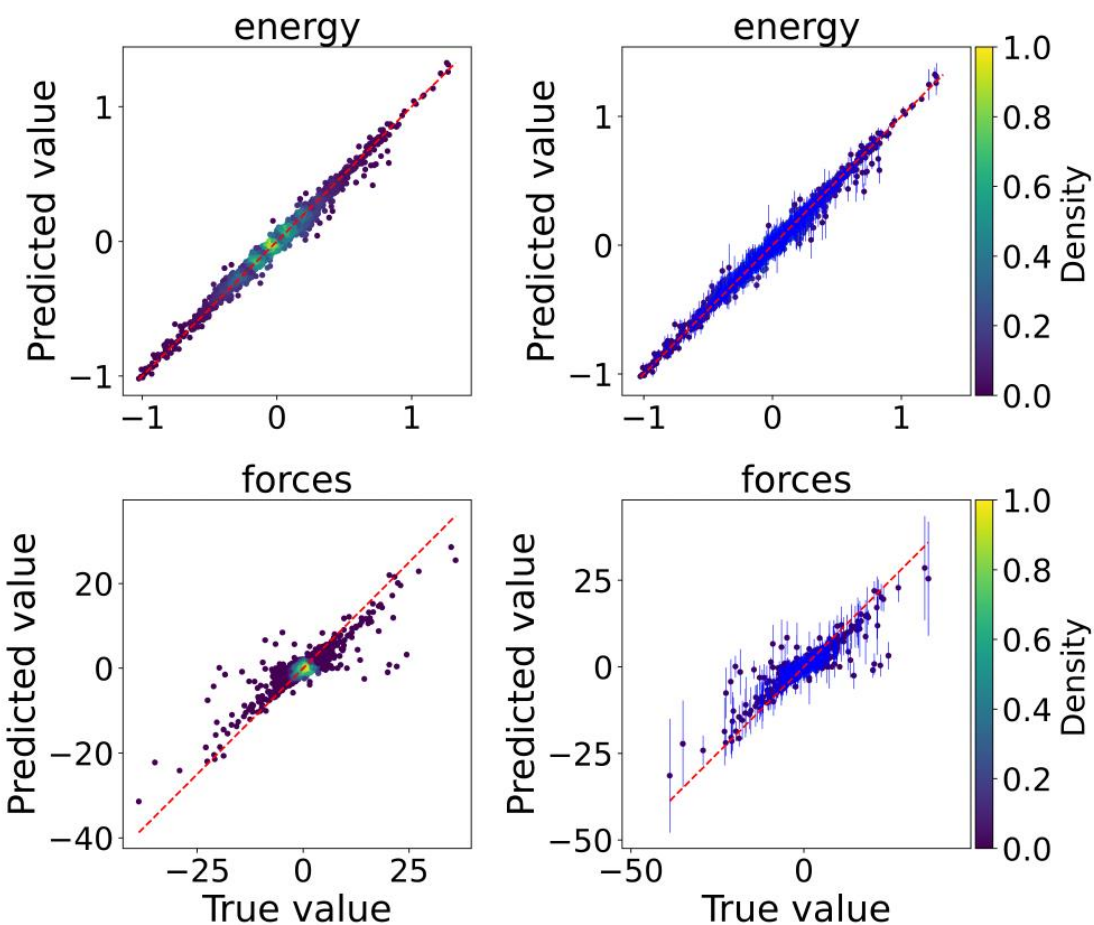
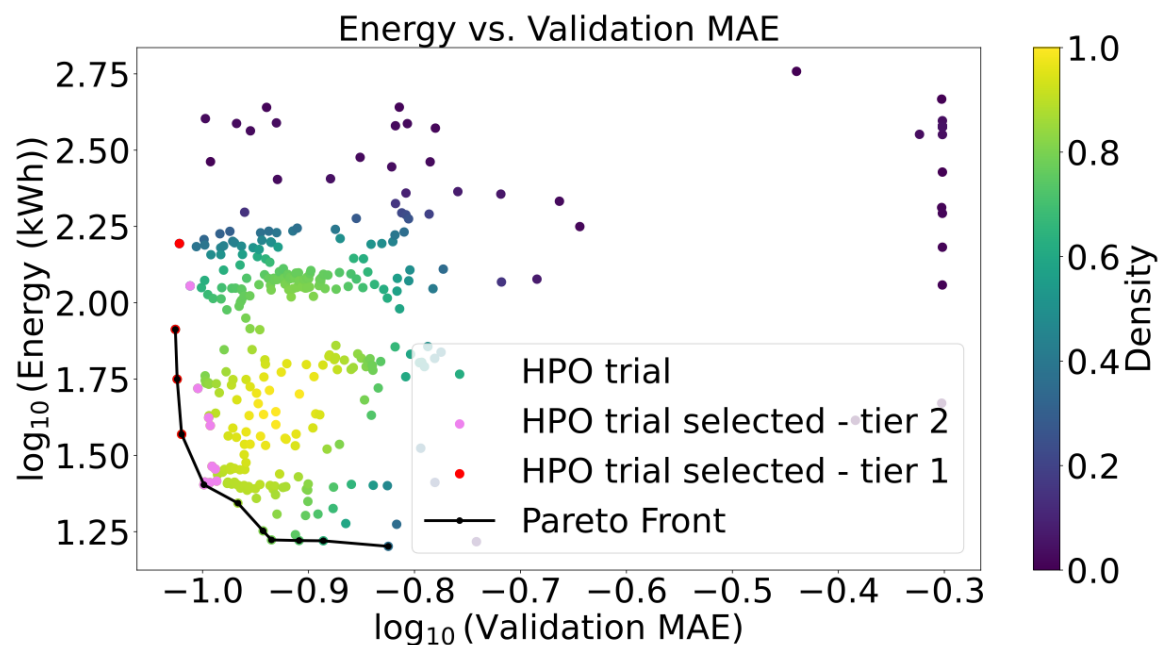


HydraGNN with 2 billion (10x larger than previous state-of-the-art) parameters on 1.2 TB of data.

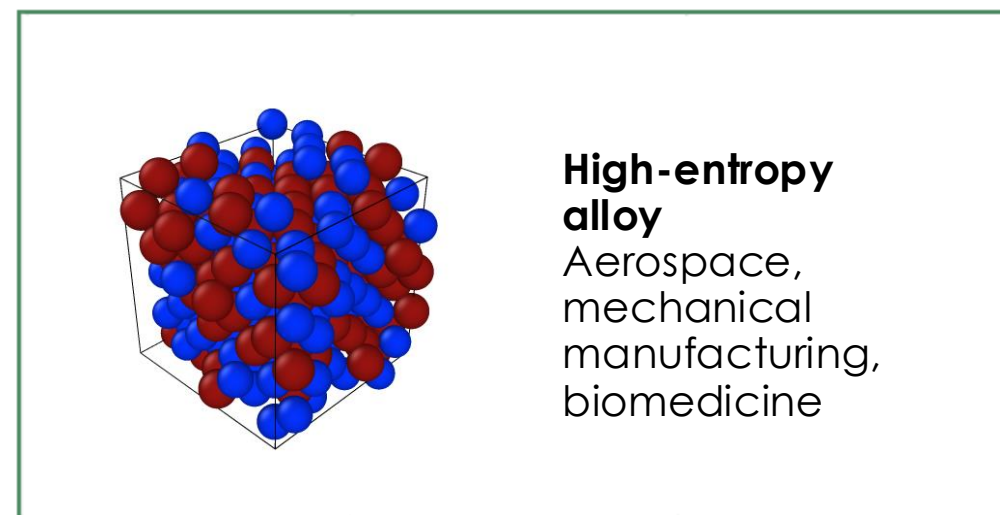
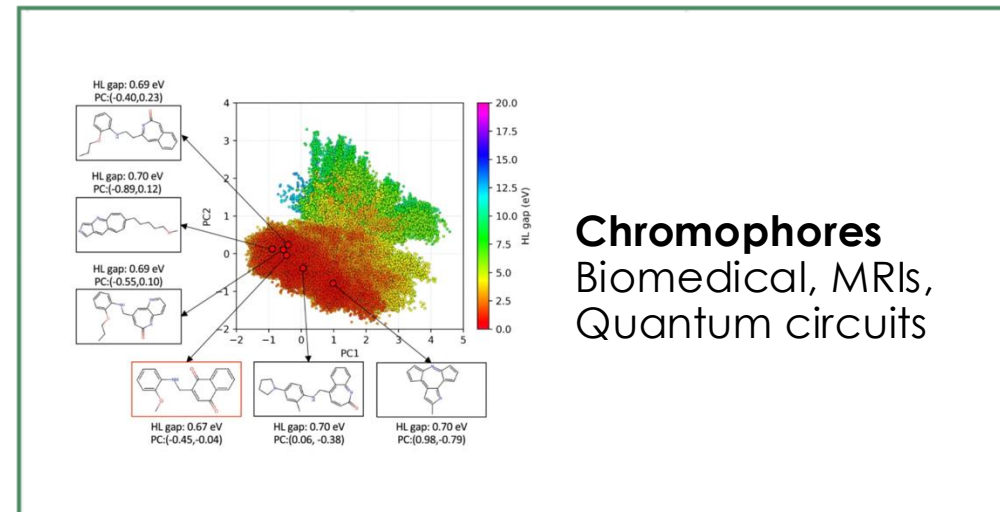
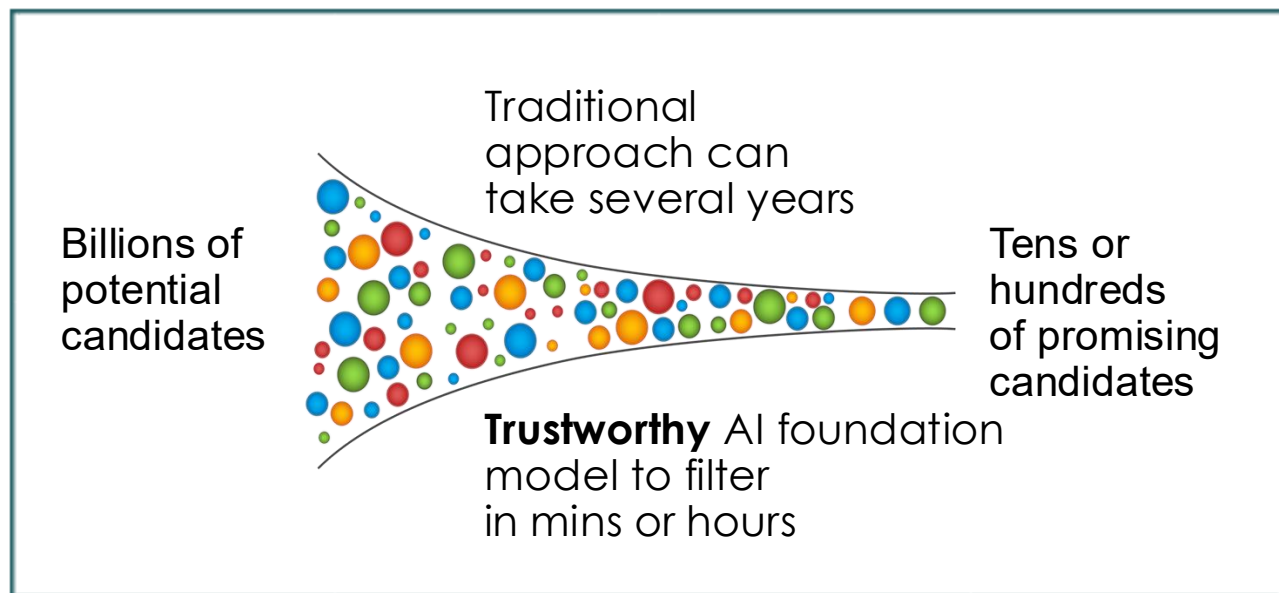


# Bi-objective optimization and ensemble uncertainty quantification

Hyperparameter	Type	Admissible values
Type of MPNN layer	Categorical	{PNA, EGNN, SchNet}
# MPNN layers	Integer	{1,...,6}
# neurons in MPNN layers	Integer	{100, ..., 2,000}
# FC layers	Integer	{2,3}
# neurons in FC layers	Integer	{300, ..., 1,000}
# batch size	Integer	{16, ..., 128}



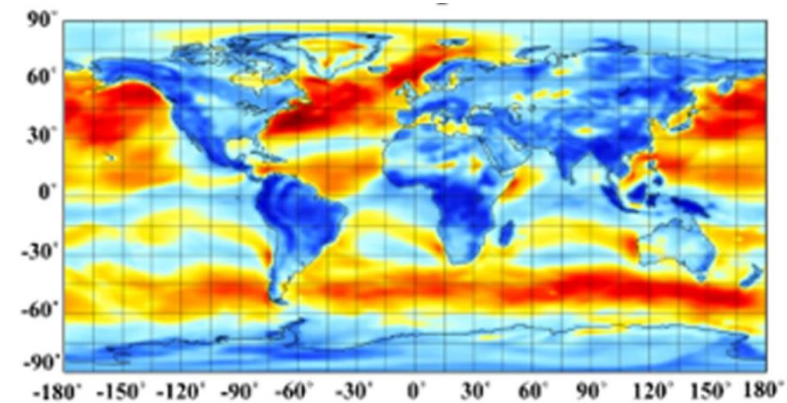
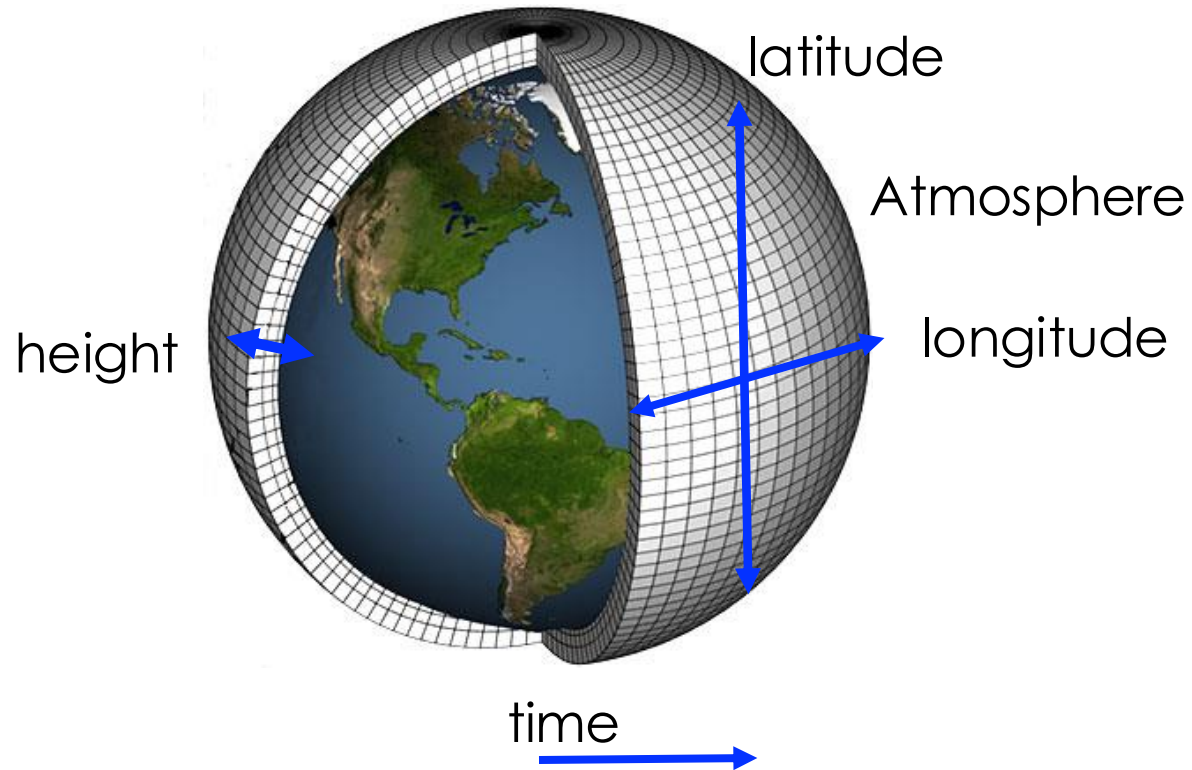
# Accelerated materials discovery via trustworthy AI models on Frontier



# Spatiotemporal data

4D+X

longitude, latitude, time, (#height, #weather/climate variables) (91)



2D visualization for a weather variable at a fixed height and time point



# Challenges for scaling up foundation models for spatiotemporal data

## Data Challenge:

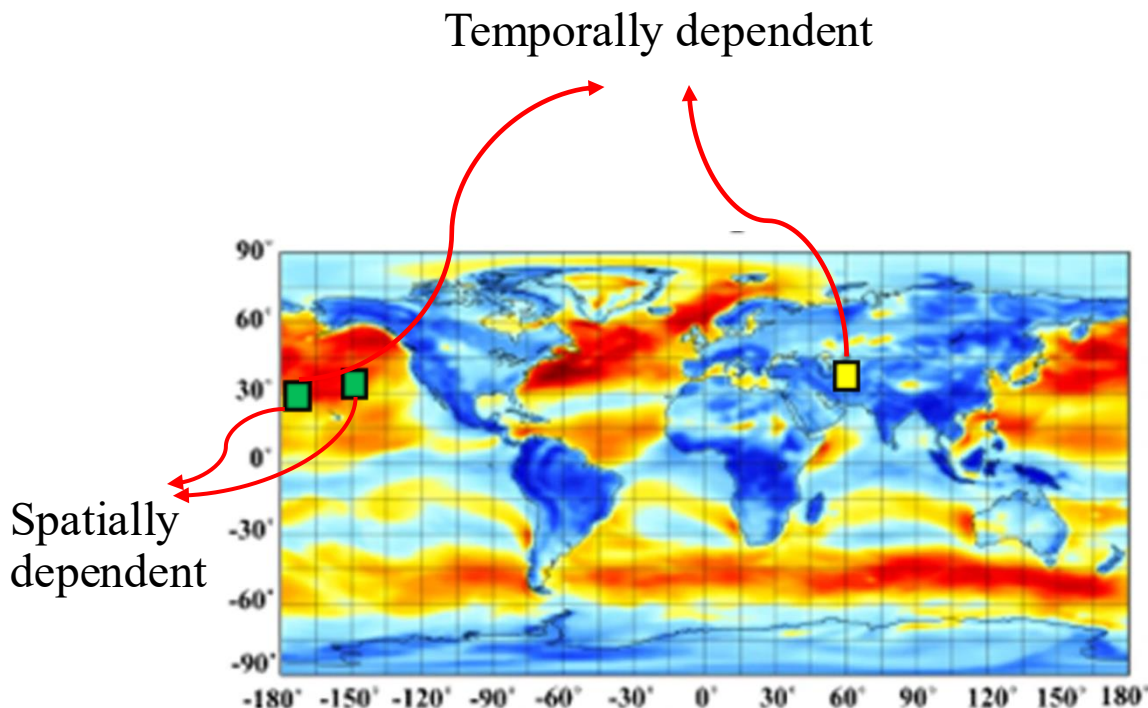
- 95D (4D+X) spatial-temporal data
- Non-linear memory and computing increase with resolution
- Complex spatial and temporal dependency

## Model Challenge:

- Larger activations, parameters, gradients, optimization states
- Non-linear memory and computing increase

Industry solutions: Pipeline, Tensor, FSDP

- Not optimized for this modality
- Limited scalability



## Implications:

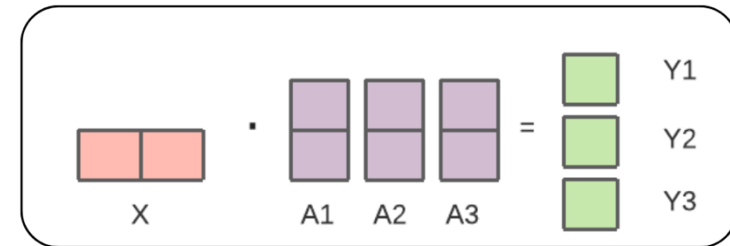
**(1) Much more expensive than LLM**

**(2) Small Vision Transformer Model Size**

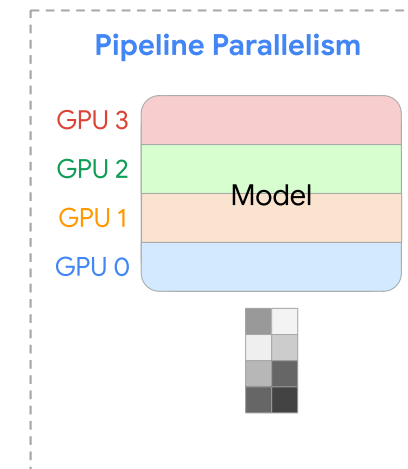
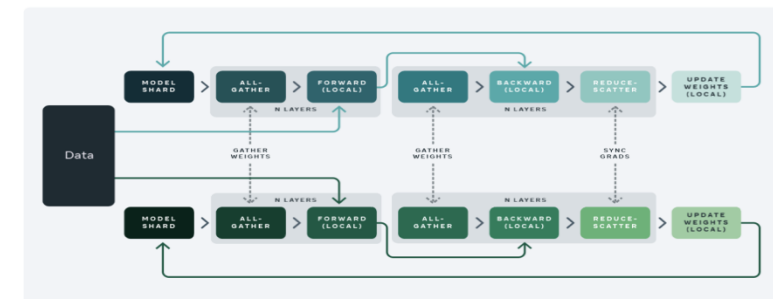
- largest dense ViT has 22 billion parameters
- largest climate ViT AI model has 115 to 500 million parameters

# Hybrid sharded tensor-orthogonal parallelism

- Tensor parallelism to ensure that compute-heavy layers don't bottleneck.
  - Dividing a single matrix multiplication across GPUs.
- Fully Sharded Data Parallelism to eliminate replicas
  - Every GPU holds just a slice of the model's weights
- Pipeline Parallelism to minimize GPUs idleness
  - the model is sliced into stages. While GPU 1 is working on Layer 1 of Batch A, GPU 2 is crunching Layer 2 of Batch B

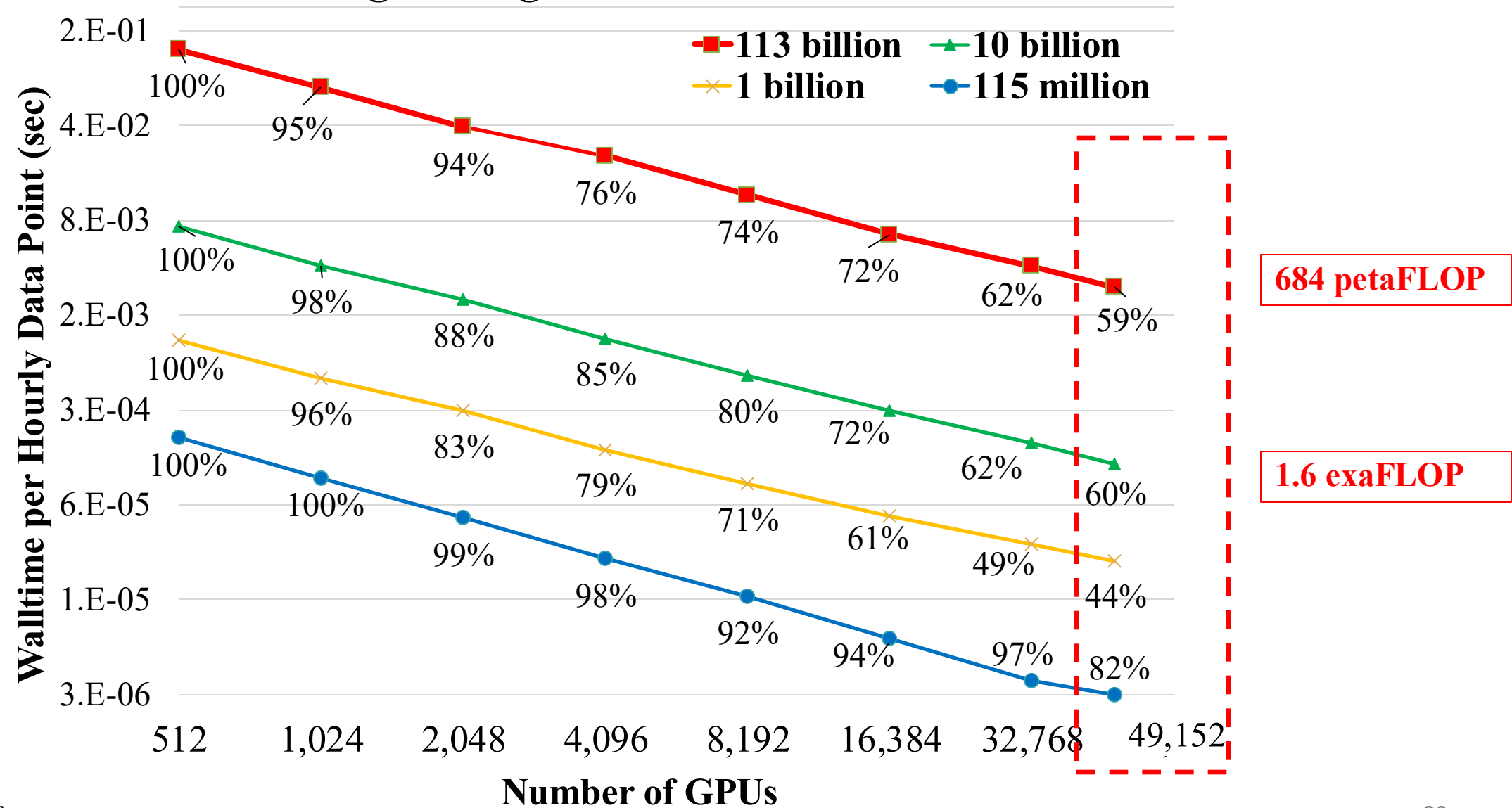


Fully sharded data parallel training



# Pretraining achieves ExaFLOP throughput with CMIP

## Strong Scaling at 48 Channel Variables





# ORBIT inferencing enables near real-time prediction

Model Size	115 million	1 billion	10 billion	113 billion
GPUs	1 GPU	1 GPU	8 GPUs	80 GPUs
Inference Speed (sec)	0.04 sec	0.24 sec	0.16 sec	0.5 sec

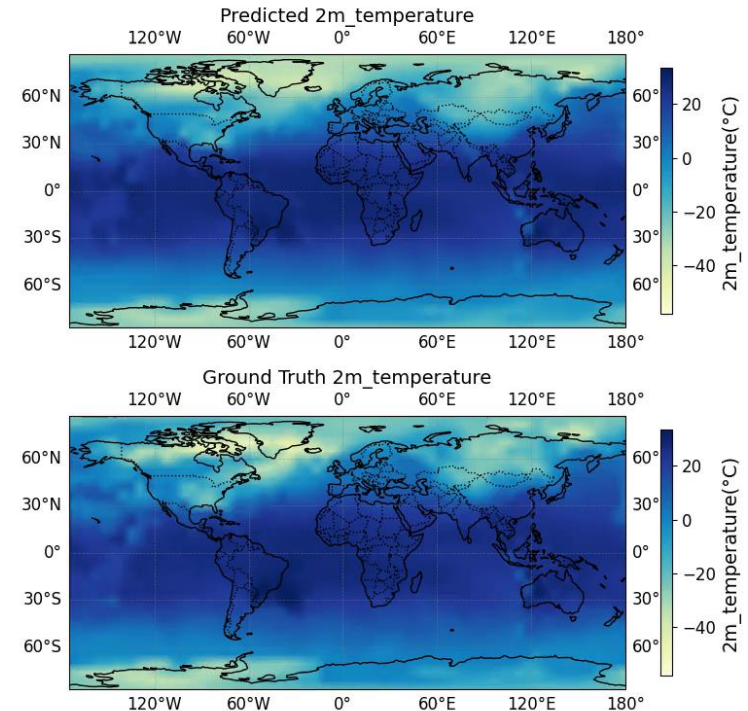
E3SM Atmosphere Model  
1.26 simulated year per day on Frontier  
supercomputer

***Potential for 500x  
speedup on limited  
resource***

# ORBIT: Oak Ridge Base Foundation Model for Earth System Predictability



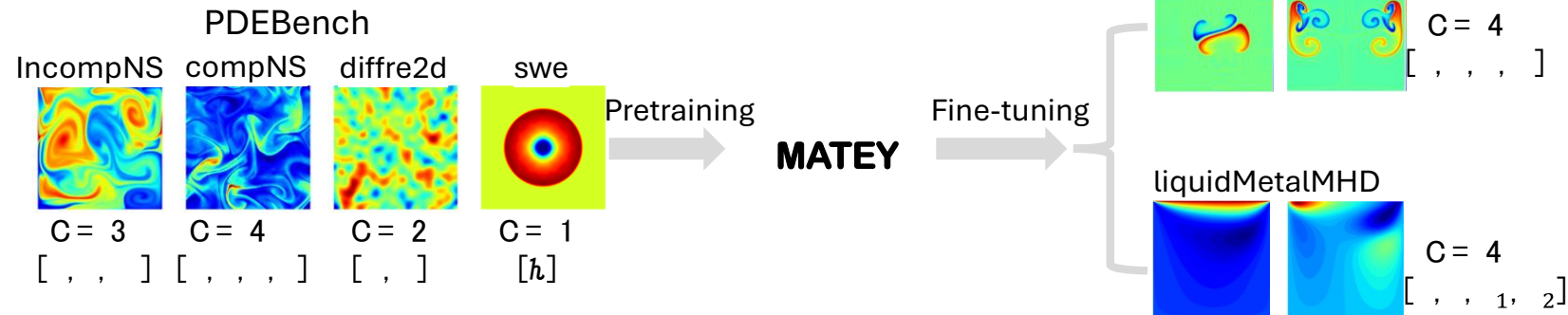
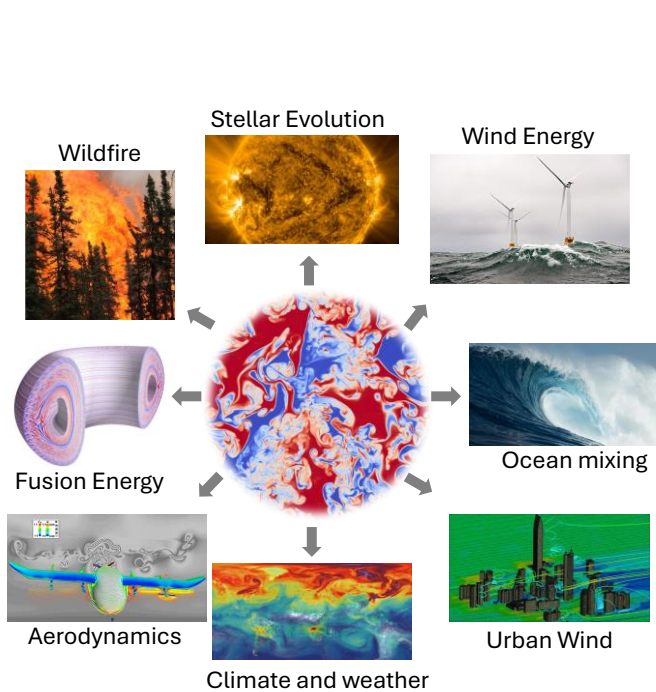
Variable 2m\_temperature, at time: 2017-01-04 02:00, lead time: 72 hrs



- Developed AI foundation model (FM), pretrained on CMIP6 model simulation data and adaptable to various Earth system modeling tasks.
- Using 49,152 GPUs on 6,144 Frontier nodes, ORBIT achieves 70% scaling efficiency with a computing throughput of 1.6 exaflops ([finalist for the 2024 and 2025 Gordon Bell Prize for Climate Modelling; 2025 SC best paper nomination](#)).
- ORBIT-2 achieves competitive or better accuracy for super resolution task

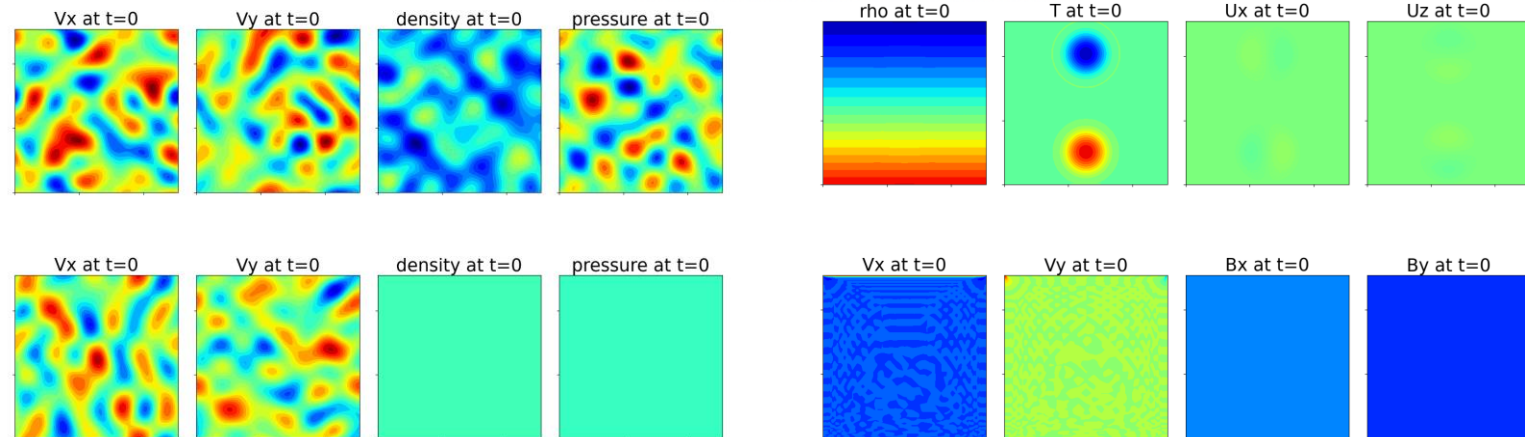


# MATEY: multiscale adaptive foundation models for spatiotemporal physical systems



$$\mathbf{u}_{t+t_{\text{lead}}} \approx \mathbf{f}_{\mathbf{w}}(\underbrace{\mathbf{u}_{t-T+1}, \dots, \mathbf{u}_t}_{\mathbf{U}_k}; t_{\text{lead}})$$

$$\mathbf{U}_k \in \mathbb{R}^{T \times H \times W \times C_k}$$



*Diverse applications characterized by the same core physics: turbulence*





# The DOE American Science Cloud

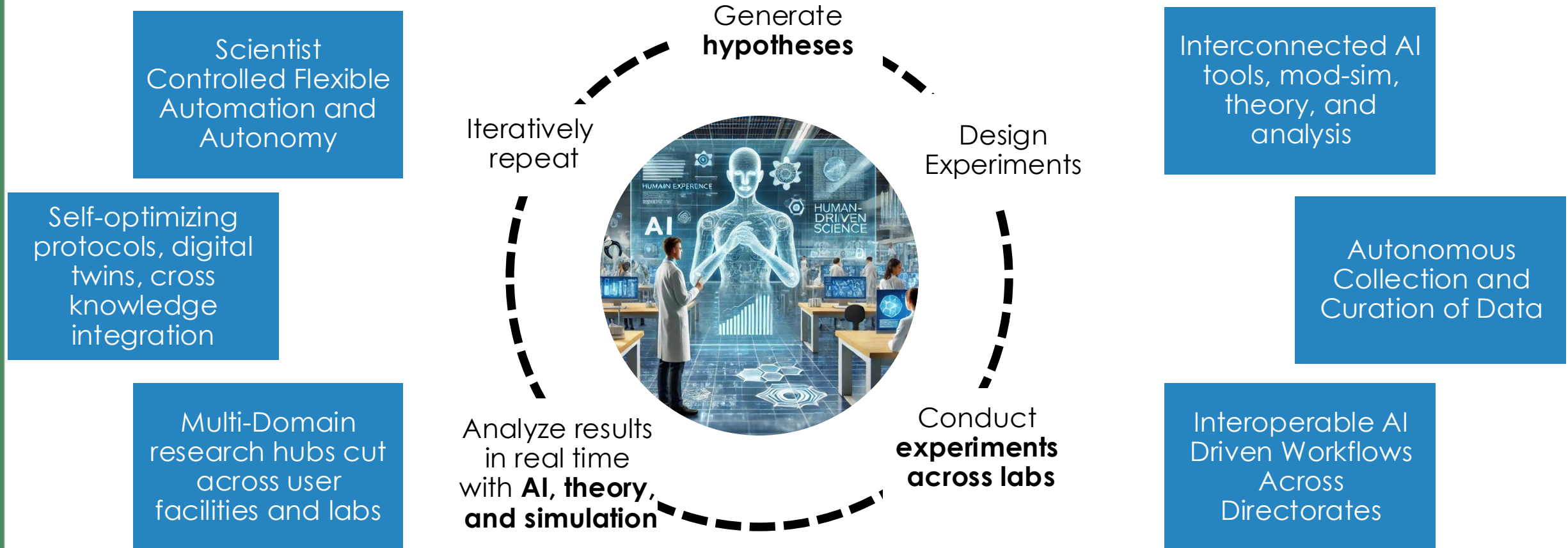


OBBS also includes forward-looking provisions that position government data as a strategic asset for American competitiveness:

- ***The American Science Cloud:*** The Department of Energy's \$150 million investment aims to mobilize National Laboratories to structure and preprocess scientific data for AI and machine learning applications. Importantly, the legislation requires these models to be made available to the broader scientific community through a cloud-based platform—embodying the principles of open data and collaborative innovation.

# AI-driven autonomy is reshaping the scientific workflow

Labs of the future: Interconnected network of multi-domain research hubs to drive new investigative approaches that **combine human creativity with evolving artificial intelligence (AI)**



**AI Agent** drive autonomous decisions  
**Robotic Platforms** automate experiments  
**Humans** in/out of the workflow loop