# AI4Science – LANL's ArtIMis Investment

Nathan DeBardeleben, Ph.D.
HPC-DESign
ArtIMis Co-Lead
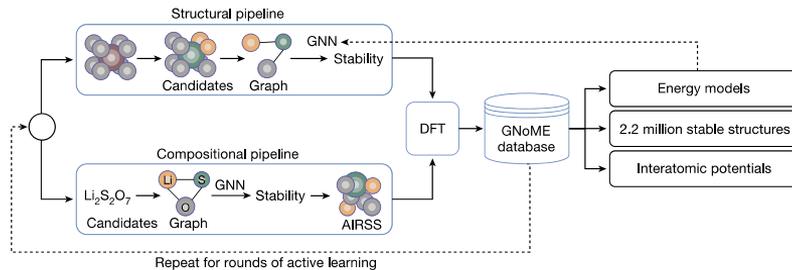~~ML Lead for Enabling Manufacturing~~

February, 2026

LA-UR-26-21430

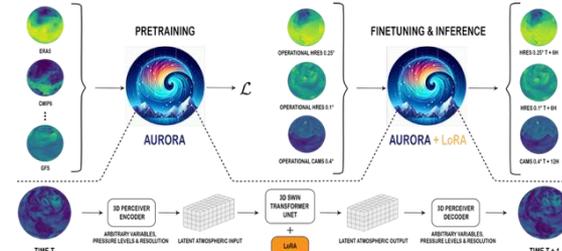# Industry has led the way on AI for science but is unlikely to lead a sustained effort.

- Need institutions with long horizons, open verification, and mission alignment
- Commercial incentives favor demo cycles, product pivots, and PR wins
- DeepMind's materials project walked back after analysis by LBNL
- Meta's Galactica taken offline days after release
- IBM sold off Watson Health – market alignment diverged from scientific utility
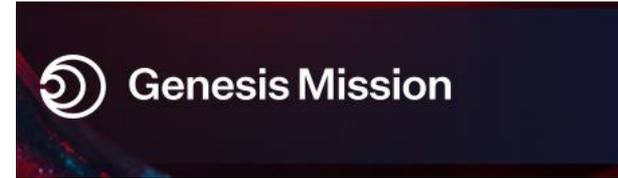
## Materials at DeepMind



Less successful than initial claims.

## Earth Systems at Microsoft



Effort abandoned.

# DOE has invested significant effort to develop an AI strategy.





Secretary Wright visited LANL to discuss Manhattan Project 2.0.

https://www.youtube.com/watch?v=Gyr_wOSy_xQ

Application Areas



**Our goal is to develop an ecosystem of agents and foundation models to address DOE mission needs.**

Enabling Technologies



ArtIMis
AI for Mission

Application Areas

Enabling Technologies

ArtIMis
AI for Mission

Emily discussed this yesterday

Moved to lab-wide effort

# We have built over a dozen special purpose science foundation models including . . .

Stochastic Latent Transformers and VAEs



GPT-like (discrete tokens)

UNet-based

Shape-agnostic PDE Foundation Model
https://huggingface.co/mahindrautela/MORPH



Transformer-based

Most models are:
- Trained on autoregressive roll-outs
- Multi-modal (input decks, equations, simulation data)
- Fine-tuned on new datasets (across domain)
- Fine-tuned for specific scalar outputs (time to fracture, physics numeric quantity, etc.)

**True foundation models** – not just surrogate models



Looking for a good name still ☺

Los Alamos
NATIONAL LABORATORY
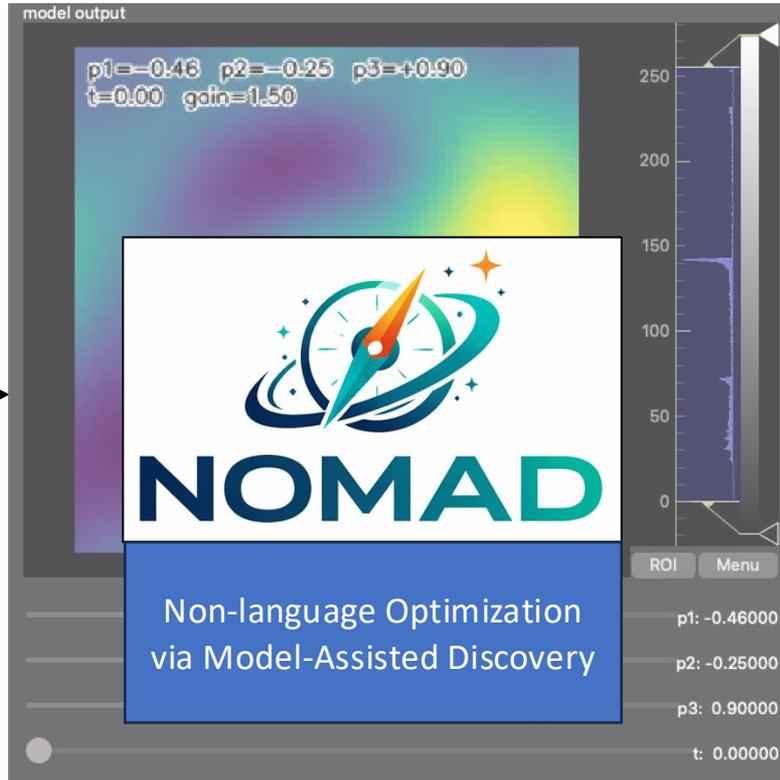
# LANL Staff Can Use Surrogate Models to Explore Designs

# But we can also drive those FMs with reasoning agents



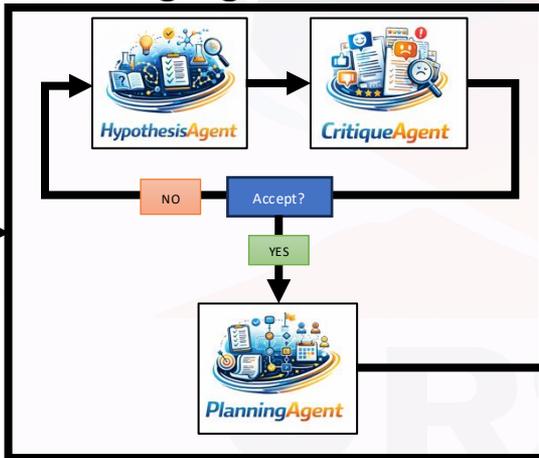Non-language Optimization via Model-Assisted Discovery

# We are building an agentic design and discovery system driven by industry reasoning models.

**Inputs**



*"Find a new alloy that resists fracture at low temperatures. Test it in a desired configuration and maximize time to failure."*

**Thinking Agents**



HypothesisAgent → CritiqueAgent

NO — Accept?

YES

PlanningAgent

**Task Agents**



SimulationAgent

CoderAgent

ControlAgent

ScienceSummarizationAgent

MORPH — PHYSICS FOUNDATION MODEL

TRANSPARENT EARTH

CHELATRON

NOMAD

FindingPhysics

AND MORE

SafetyAgent

Bomberman

MathOptAI.jl

WebSearchAgent

**Hypothesis Generation**

Nemotron-Super-49B

**Hypothesis Evaluation**

"Find a ligand that could become a cancer drug."

Group A

Group B

**Molecule Design**

GenMol

**Complex Construction**
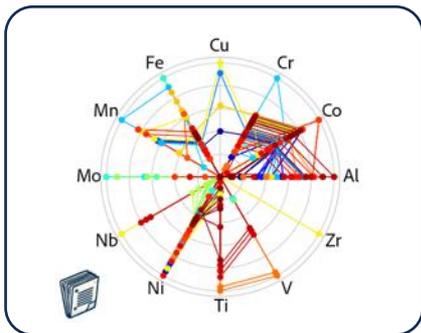
ARCHITECTOR
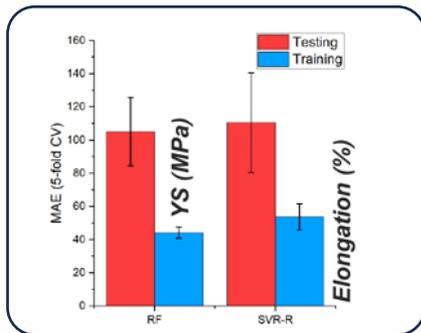
**Computational Modeling**
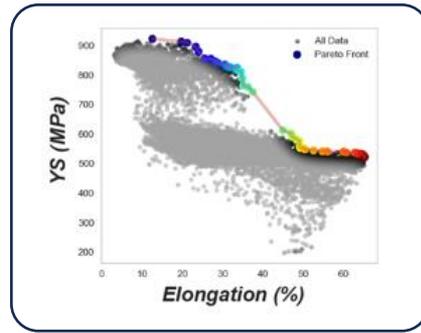
$$\widehat{H}\Psi = E\Psi$$

3/4/2026    19

# We are automating a materials design lab for alloy discovery.
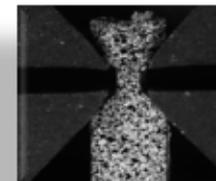


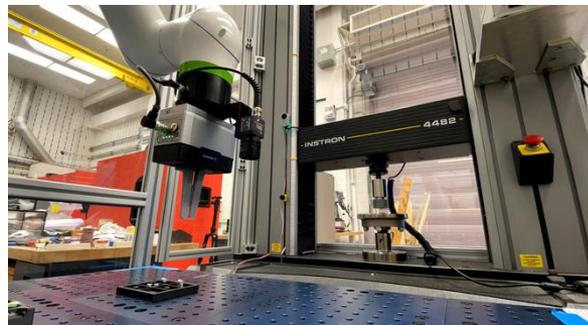**Data collection from literature. Automated.**



**Property prediction. Partially automated**



**Model Inversion. Partially automated.**
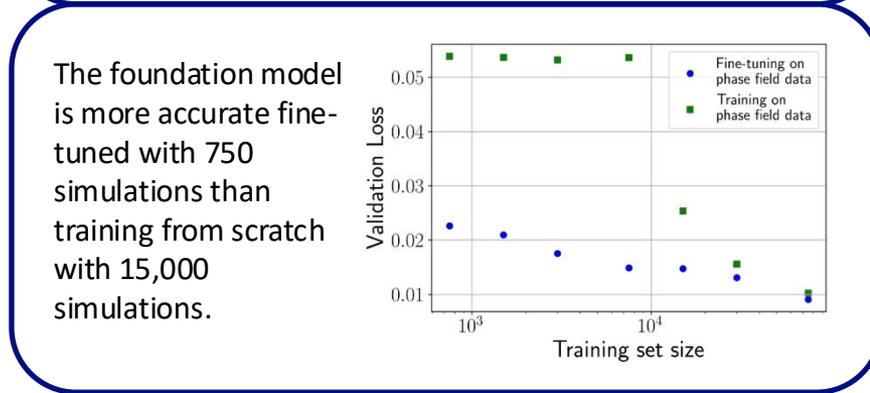


**Manufacturing and testing. Partially automated.**
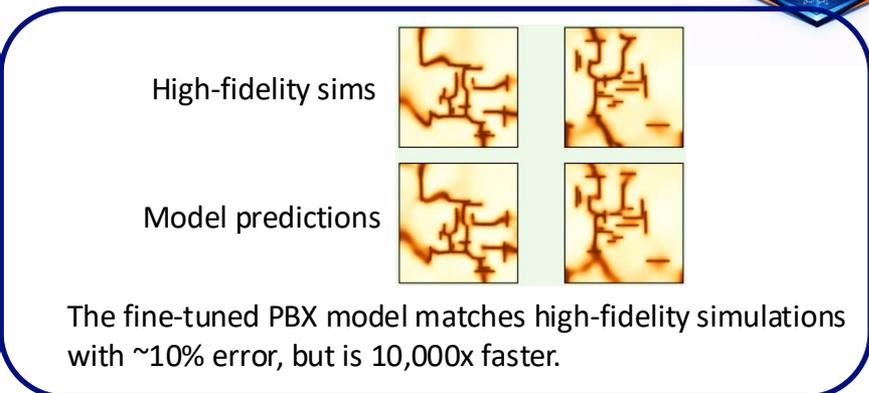






Team Leads: Saryu Fensin and Tim Germann

**Testing is automated. Automated manufacturing this summer.**

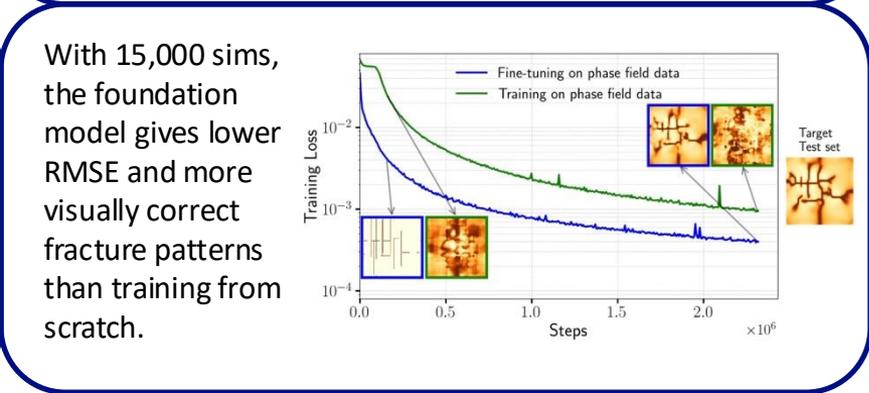# The fracture foundation model predicts PBX fractures that match physics sims, but 10,000x times faster.



The transformer-based model is pre-trained on simplistic fracture data and then fine-tuned with physics-based simulations.



High-fidelity sims

Model predictions

The fine-tuned PBX model matches high-fidelity simulations with ~10% error, but is 10,000x faster.

The foundation model is more accurate fine-tuned with 750 simulations than training from scratch with 15,000 simulations.



With 15,000 sims, the foundation model gives lower RMSE and more visually correct fracture patterns than training from scratch.

# Models include tuning we have found useful

- Curriculum learning – tasks get progressively harder
- Scheduled sampling – sample from either ground truth or rollout during training
- "Pushforward trick"
- Data augmentation (including noise injection)
- Domain-relevant loss functions
- And more

MESSAGE PASSING NEURAL PDE SOLVERS

**Johannes Brandstetter**[*]
University of Amsterdam
Johannes Kepler University Linz
brandstetter@ml.jku.at

**Daniel E. Worrall**[*]
Qualcomm AI Research[†]
dworrall@qti.qualcomm.com

**Max Welling**
University of Amsterdam
m.welling@uva.nl

# We have invested heavily in improving auto-regressive rollouts



BEFORE OPTIMIZATION

AFTER OPTIMIZATION

Ground Truth (simulation)

Auto-regressive Rollout

Difference (residual)

Ground Truth (simulation)

Auto-regressive Rollout

Difference (residual)

HEAT dataset: LA-UR-25-26490
https://oceans11.lanl.gov/

# Michigan's MIST model has found application in our bio and HE work.



Fine-tuned with LANL-curated data, MIST can predict properties for energetic molecules.



Fine-tuned MIST embeddings let the Bio team vastly improve token prediction related to small-molecule binding.

# We have an ICF design demo combining planning, research, simulation, and data processing agents.

# The planning agent does a careful breakdown with expected outputs and success criteria.

# URSA can use molecular science foundation models to explore promising compounds

# Available now: https://github.com/lanl/ursa

# University partnerships have expanded our ability to make progress on fundamentals.



Foundation models for materials and multiphysics.



Design, discovery & control with Georgia Tech.



Applications.

# Industry partnerships have provided access to tools and accelerated progress.



Training and support for model building.



Early access for test and evaluation.

# We have the only on-prem installation of an OpenAI model anywhere in the world.

## Lab partners with OpenAI to advance national security

**OpenAI's latest reasoning models will be used on Venado supercomputer**

Los Alamos National Laboratory has entered a partnership with OpenAI to install its latest o-series models — capable of expert reasoning for a broad span of complex scientific problems — on the Lab's Venado supercomputer, which uses NVIDIA GH200 Grace Hopper Superchips, to conduct national security research.

"As threats to the nation become more complex and more pressing, we need new approaches and advanced technologies to preserve America's security," said Laboratory Director **Thom Mason.** "Artificial intelligence models from OpenAI will allow us to do this more successfully, while also advancing our scientific missions to solve some of the nation's most important challenges."
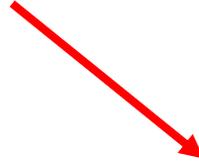
The Venado machine will be moved to a secure, classified network where it will be a shared resource for researchers from Los Alamos, Lawrence Livermore, and Sandia national labs.

Venado, the Lab's newest supercomputer, will use models from OpenAI.

**TODAY'S TOP NEWS**

# Utility vs. Spectacle

We don't need **THIS**,
for **THIS** to be useful

# 5 AI Myths (and Reality) at LANL

| MYTH | REALITY |
|------|---------|
| ❌ **AI replaces** mod/sim | ✅ **AI accelerates** mod/sim — physics stays truth |
| ❌ **GenAI** is too wild for science | ✅ GenAI is a **controlled interface** (read/write/code) |
| ❌ We need **AGI** | ✅ **Narrow models** deliver mission value now |
| ❌ **Agents** are for booking trips | ✅ **Agents** automate **workflows** (data → run → report) |
| ❌ We hand decisions to AI | ✅ AI advises; humans decide (bounded autonomy) |

## AI is a force multiplier—not a replacement.

# LANL's ArtIMis in Summary . . .

- Leverage frontier industry reasoning models
- AI ecosystem of reasoning agents powering non-language science foundation models
- Partnerships with industry and academia
- Dozens of papers – CS as well as in applied domains
- Deep application to LANL mission
- Feeding directly into Genesis Mission project

# URSA

https://github.com/lanl/ursa

## Thanks for the Invite!

Feel free to reach out

Nathan DeBardeleben, Ph.D.

High Performance Computing Design (HPC-DES)

ndebard@lanl.gov