



U.S. DEPARTMENT OF
ENERGY

Office of
Science

DOE Scientific Machine Learning & AI Overview: Machine Learning Seminar @ Jefferson Lab

**Presented by Steven Lee (DOE Program Manager)
Office of Advanced Scientific Computing Research
November 6, 2018**

Based on Basic Research Needs (BRN) Workshop

Held January 30 - February 1, 2018

Workshop Chair: Nathan Baker, PNNL

<https://www.ornl.gov/ScientificML2018/>



Summary of Charge Letter for Scientific Machine Learning Workshop

Greater machine learning-based prediction & decision-support capabilities are needed to address & anticipate DOE mission challenges:

- **DOE scientific user facilities drive rapid growth in data** from experiments, observations, and simulations
- Increasingly powerful science technologies are driving the need for **algorithms & automation to facilitate the use of advanced technologies** for science breakthroughs

The charge for the workshop is:

- First consider the status, recent trends, and broad use of machine learning for scientific computing
- Examine the opportunities, barriers, & potential for high scientific impact through fundamental advances in the underlying research foundations
- ASCR grand challenges & resulting priority research directions should span several major machine learning categories & state-of-the-art modeling & algorithms research
- ***Identify the basic research needs & opportunities that can potentially enable machine learning-based approaches to transform the future of science and energy research.***

Working Definitions of Machine Learning

Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.

- Arthur Samuel, 1959

Machine Learning: A set of rules that allows systems to learn directly from examples, data and experience.

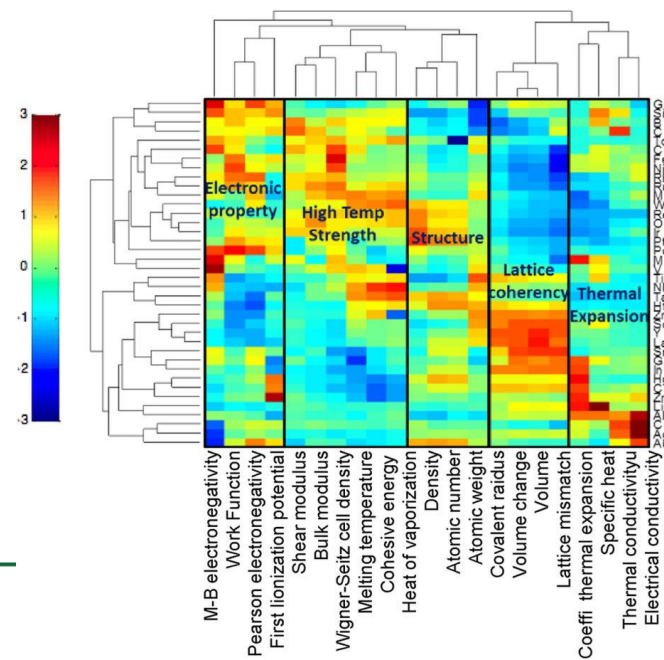
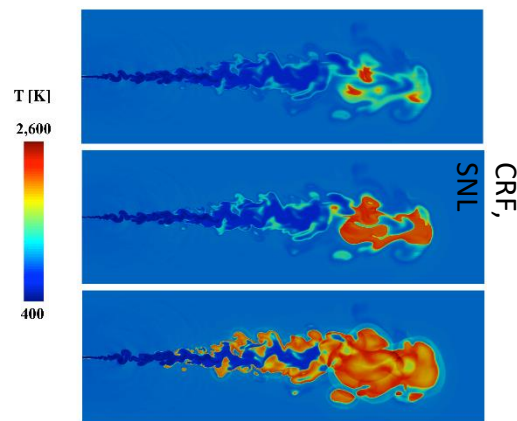
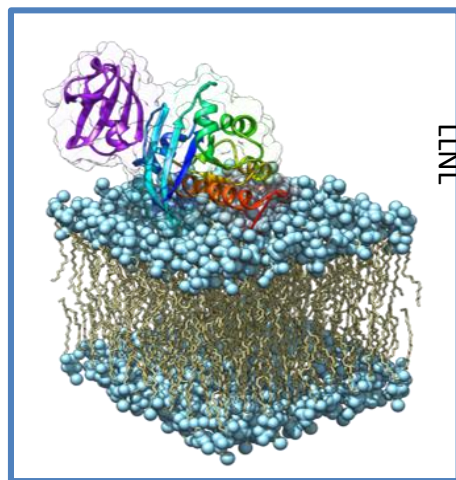
- Royal Society, 2017

“Learning” is the process of transforming information into expertise or knowledge; “Machine learning” is automated learning.

- Paraphrased from Jordan et al., 2015

Why: Define research challenges and directions for Scientific Machine Learning

- Machine learning use is on the rise throughout science domains
- However, many popular ML methods lack mathematical approaches to understand robustness, reliability, etc.
- ASCR Applied Mathematics has a long track record for building mathematical foundations to critical computational tools
- **Workshop to help ASCR define the grand challenges and priority research directions for scientific machine learning**



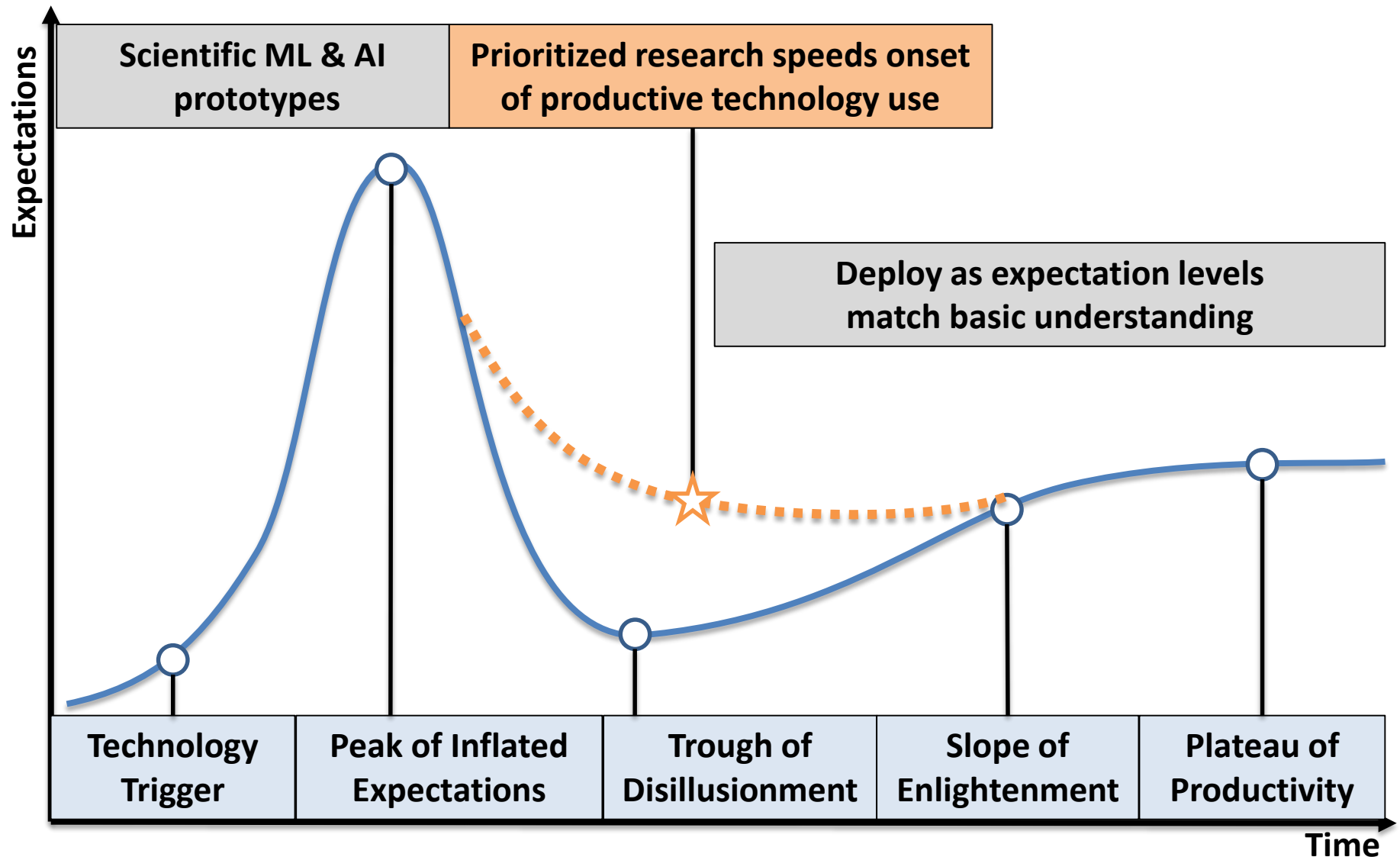
Krishna Rajan, Buffalo



U.S. DEPARTMENT OF
ENERGY

Office of
Science

Foundational Research will increase our basic understanding of Scientific Machine Learning & AI technologies



What: Deliverables and Products from this Workshop

- A BES Basic Research Needs-inspired process
- Pre-workshop report (January 2018)
 - Factual status document describing the Scientific Machine Learning landscape as it relates to ASCR
- Workshop deliverables (April 2018)
 - Articulation & refinement of grand challenges for Scientific ML
 - Priority research directions for Scientific ML
- Post-workshop report (~November 2018)
 - Incorporate updated factual status document
 - Incorporate workshop deliverables

What is a “Priority Research Direction”?

- High-priority area of research for scientific machine learning
- Has following components (ala Heilmeyer):
 - Clear statement of key challenge
 - Context in the current scientific landscape to establish timeliness and competition
 - Plausible research pathway(s)
 - Clear scientific impact
- *It is not*
 - A proposal for a specific project
 - Your favorite area of research without connection to Scientific ML themes

Please provide your title here

Key challenge	State of the art
Please provide a brief overview of the underlying science challenge	Please answer the following questions: <ul style="list-style-type: none">• Why is this a timely challenge?• Who else is doing this?
New research direction	Potential scientific impact
Please answer the following questions: <ul style="list-style-type: none">• What will you do to address the challenge?	Please answer the following questions: <ul style="list-style-type: none">• What new scientific capabilities will follow?• What new methods and techniques will be developed?

Scientific Machine Learning 2018 Workshop Please provide list of authors and affiliations here.

How: Workshop Components

- **Plenary talks**
 - Highlight status of machine learning, challenges, open questions
- **Panel discussions**
 - Summarize pre-workshop report
 - Provide perspectives across DOE ASCR facilities, Exascale Computing Project, and other organizations & programs
- **Breakout sessions**
 - Organized ~140 submitted Position Papers presented as flash talks
 - The “work” in workshop: Crucible for new Priority Research Directions
 - Need high levels of interaction and input (long days...)
 - Brainstorming (Day 1), Refining (Day 2) & Presenting (Day 3) Priority Research Directions



Scientific Machine Learning: Priority Research Directions (PRDs)

Foundational Themes

PRD1. Domain-Aware Scientific ML
Leveraging scientific domain knowledge

PRD2. Interpretable Scientific ML
Explainable & understandable results

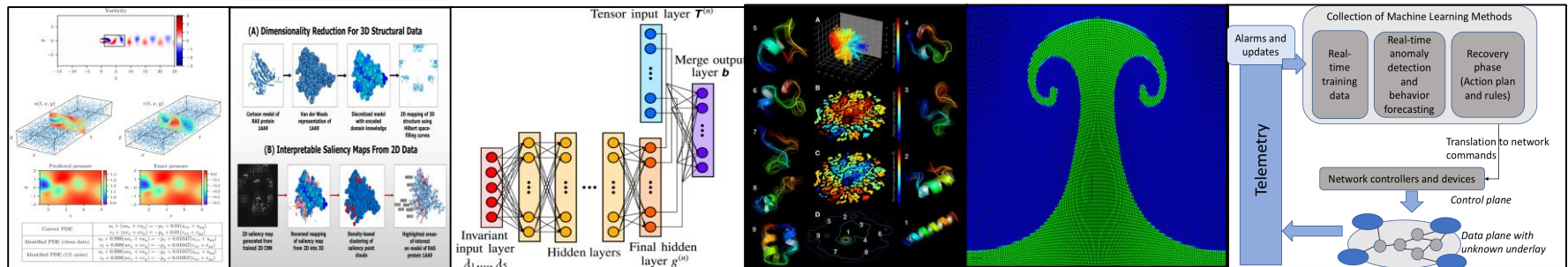
PRD3. Robust Scientific ML
Stable, well-posed & efficient formulations

Capabilities Research

PRD4. Data-Intensive Scientific ML
Automated scientific inference & data analysis

PRD5. ML-Enhanced Modeling & Simulations
Machine learning-enhanced models & algorithms for better scientific computing tools

PRD6. ML-Enhanced Decision Support
Automated decision-support, optimization, resilience, & control for complex systems & processes

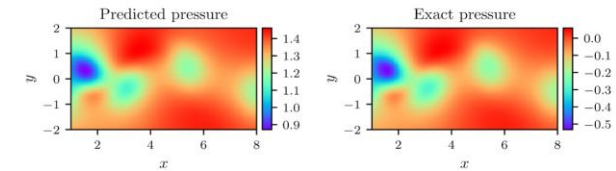
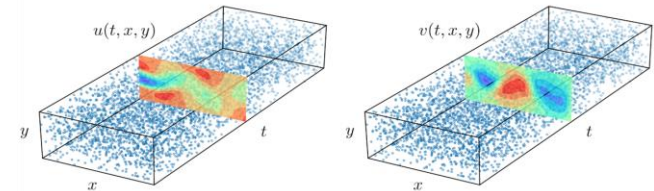
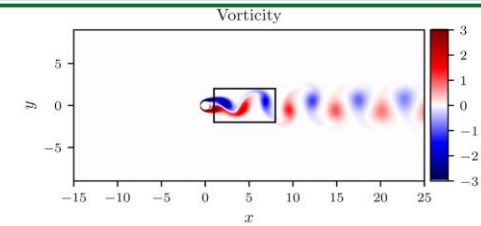


PRD1: Domain-Aware Scientific Machine Learning

Leveraging Scientific Domain Knowledge

Key Points: How can domain knowledge be effectively incorporated into Scientific ML methods?

- Established domain models based on physical mechanism & scientific knowledge
- Scientific ML offers significant opportunity to complement traditional domain models
- Domain knowledge: Physical principles, symmetries, constraints, computational predictions, uncertainties, etc
- Potential to improve accuracy, interpretability, & defensibility while reducing data requirements & accelerating training process



Correct PDE	$u_t + (uu_x + vv_y) = -p_x + 0.01(u_{xx} + v_{yy})$ $v_t + (uv_x + vv_y) = -p_y + 0.01(v_{xx} + v_{yy})$
Identified PDE (clean data)	$u_t + 0.999(uu_x + vv_y) = -p_x + 0.01047(u_{xx} + v_{yy})$ $v_t + 0.999(uv_x + vv_y) = -p_y + 0.01047(v_{xx} + v_{yy})$
Identified PDE (1% noise)	$u_t + 0.998(uu_x + vv_y) = -p_x + 0.01057(u_{xx} + v_{yy})$ $v_t + 0.998(uv_x + vv_y) = -p_y + 0.01057(v_{xx} + v_{yy})$

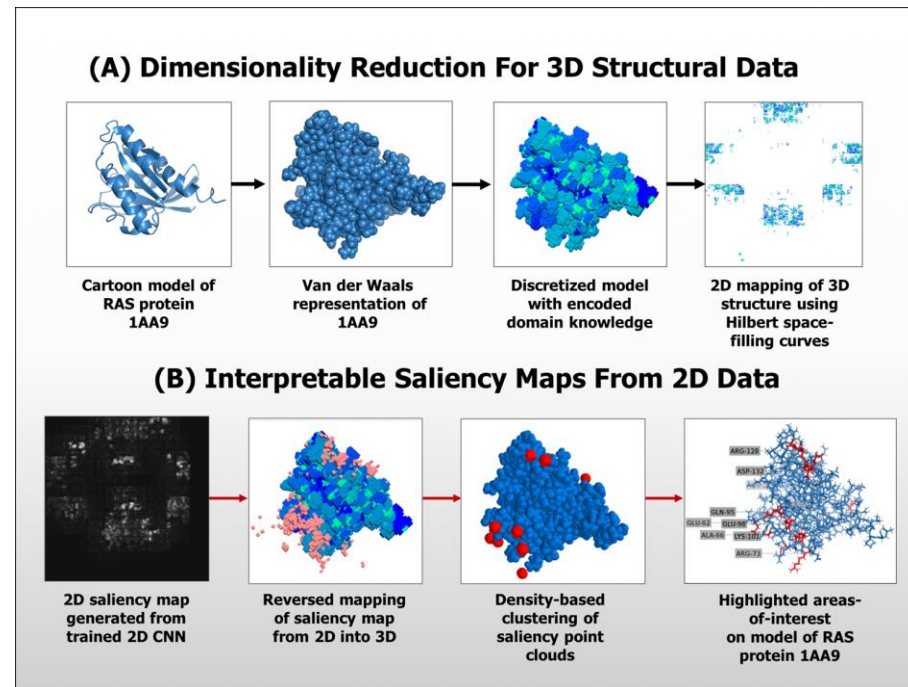
This example illustrates the capabilities obtained by incorporating domain knowledge into a deep neural network. Given scattered and noisy data components of an incompressible fluid flow in the wake of a cylinder, we can employ a physics-informed neural network that is constrained by the Navier-Stokes equation in order to identify unknown parameters, reconstruct a velocity field that is guaranteed to be incompressible and satisfy any boundary conditions, as well as recover the entire pressure field. Figure from: Raissi et al.

PRD2: Interpretable Scientific Machine Learning

Explainable and Understandable Results

Key Points: How to balance the use of increasingly complex ML models with the need for users to understand conclusions & derive insights?

- Physical understanding has been the bedrock of modeling
- User confidence linked to the conviction that model accounts for domain knowledge (variables, parameters, physical laws, etc.)
- Need exploration & visualization approaches for “debugging” complex machine learning models
- Need metrics to quantify model differences



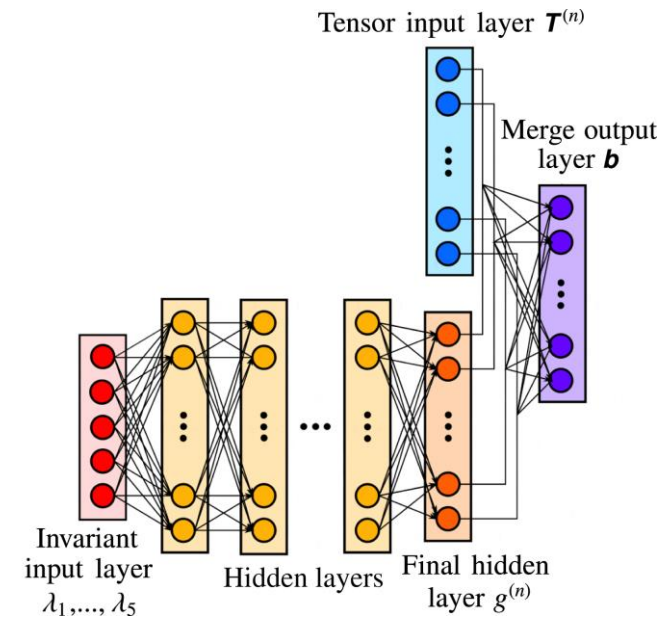
High-level data pipeline overview for dimensionality reduction of 3D protein structures (A) and interpretation of saliency maps from trained CNN model (B). Saliency maps generated from CNN models can then be clustered to identify areas along the 3D structure that are regions that highly influence the output of the CNN model. From these salient regions, specific residues can be identified that fall in close proximity to the salient regions. Image credit: Rafael Zamora-Resendiz and Silvia Crivelli, LBNL.

PRD3: Robust Scientific Machine Learning

Stable, Well-Posed, and Efficient Formulations

Key Points: How can computationally efficient Scientific ML methods be developed and implemented to ensure outcomes are not unduly sensitive to perturbations in training data and model selection?

- Scientific ML methods need to establish the properties of robustness & reliability
- Integration of protocols for verification & validation are in their infancy
- Progress will require research proving that developed methods and implementations are stable and well-posed



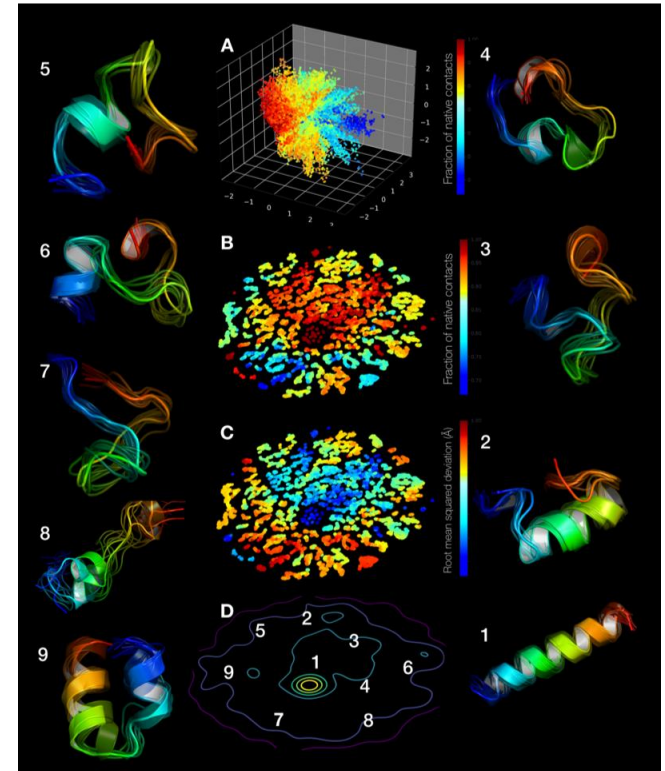
In the context of Reynolds averaged incompressible turbulence modeling, a neural network has been used in an eddy viscosity turbulence closure model. From physical arguments, the model needs to satisfy rotational invariance, ensuring that the physics of the flow is independent of the orientation of the coordinate frame of the observer. A special network architecture, a tensor basis neural network (TBNN), embeds rotational invariance by construction. Without this guarantee, the NN model evaluated on identical flows with the axes defined in different directions could yield different predictions.
Image credit: SNL.

PRD4: Data-Intensive Scientific Machine Learning

Automated Scientific Inference & Data Analysis

Key Points: What novel approaches can be developed for reliably finding signals, patterns or structure within high-dimensional, noisy, uncertain input data?

- Scientific ML methods require the development of improved methods for statistical learning in high-dimensional Scientific ML systems with noisy and complex data
- Need approaches required to identify structure in complex high-dimensional data
- Scientific ML requires efficient sampling in high-dimensional parametric and model spaces



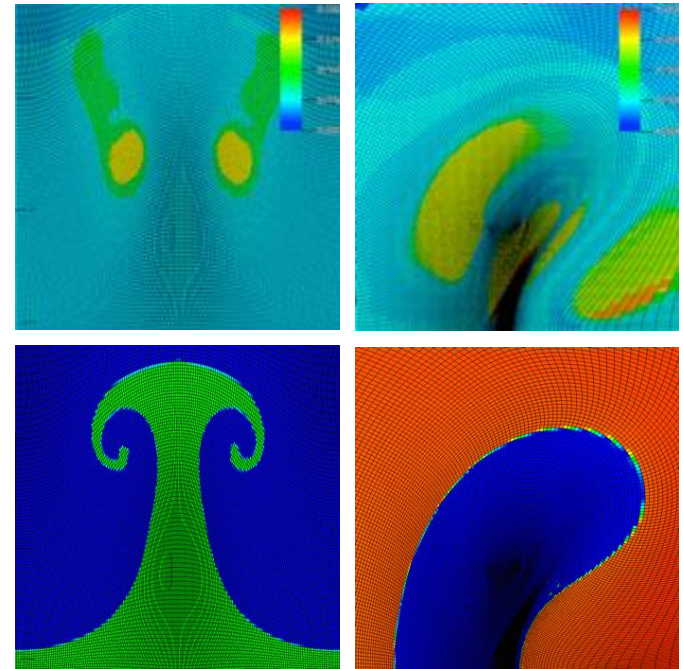
ML techniques reveal Fs-peptide folding events from long time-scale molecular dynamics simulations. A low dimensional embedding of the simulation events reveal transitions from fully unfolded states (blue) to fully folded states (red). A two dimensional embedding using t-test stochastic neighborhood embedding shows the presence of near native states (labeled state 1) versus partially unfolded (2-7) and fully unfolded states (8-9) in the picture. Image Credit: Arvind Ramanathan, ORNL.

PRD5: ML-Enhanced Modeling & Simulations

Hybrid Machine Learning, Models, & Algorithms

Key Points: What is the role and potential advantages of ML-enhanced approaches in computational model and algorithm development?

- Combination of scientific computing with learned adaptivity for more efficient simulations
- ML for in-situ parameter tuning
- ML for sub-grid physics models
- Progress will require the development of new methods to quantify tradeoffs and optimally manage the interplay between traditional and ML models and implementations



The arbitrary Lagrangian-Eulerian (ALE) method is used in a variety of engineering and scientific applications for enabling multi-physics simulations. Unfortunately, the ALE method can suffer from simulation failures, such as mesh tangling, that require users to adjust parameters throughout a simulation just to reach completion. A supervised ML framework for predicting conditions leading to ALE simulation failures was developed and integrated into a production ALE code for modeling high energy density physics.

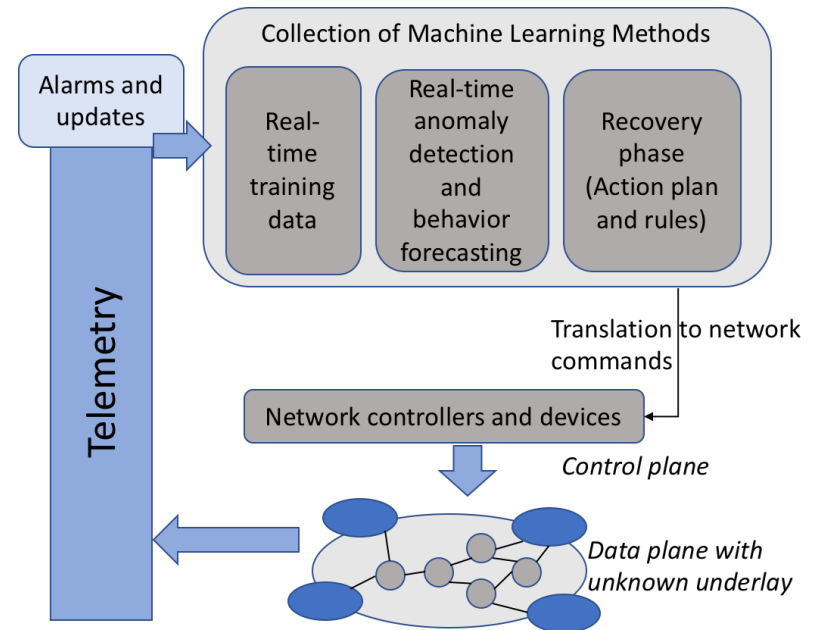
Image credit: M. Jiang, LLNL.

PRD6: ML-Enhanced Decision Support & Automation

Automated Decision Support, Optimization, Resilience, & Control

Key Points: What are the challenges in managing the interplay between automation & human decision-making?

- Outer-Loop applications include optimization, uncertainty quantification, inverse problems, data assimilation, & control.
- New mathematically & scientifically justified methods to guide data acquisition and ensure data quality and adequacy.
- Scientific ML methods for improving system resilience or responsiveness.



Exascale applications are exponentially raising demands from underlying DOE networks such as traffic management, operation scale and reliability constraints. Networks are the backbone to complex science workflows ensuring data is delivered securely and on-time for important compute to happen. In order to intelligently manage multiple network paths, various tasks such as pre-computation and prediction are needed to be done in near-real-time. ML provides a collection of algorithms that can add autonomy and assist in decision making to support key facility goals, without increased device costs and inefficiency. In particular, ML can be used to predict potential anomalies in current traffic patterns and raise alerts before network faults develop. Image credit: Prabhat, LBNL.

Compelling Scientific Machine Learning Examples

Scientific Machine Learning has widespread Science & Energy uses

Three Capability PRDs (and combinations) cover nearly all examples

- Data-Intensive Scientific ML
- ML-Enhanced Modeling & Simulations
- ML-Enhanced Decision Support & Automation

Compelling Big Science use cases include:

- Improved operational capabilities of scientific user facilities
- Better computational models from data-compute convergence
- Automation & adaptivity within scientific method (systems, processes)
- Many more ...

Core Agenda for Scientific Machine Learning Research

Scientific Machine Learning: Priority Research Directions (PRDs)

Foundational Themes	Capabilities Research
PRD1. Domain-Aware Scientific ML Leveraging scientific domain knowledge	PRD4. Data-Intensive Scientific ML Automated scientific inference & data analysis
PRD2. Interpretable Scientific ML Explainable & understandable results	PRD5. ML-Enhanced Modeling & Simulations Machine learning-enhanced models & algorithms for better scientific computing tools
PRD3. Robust Scientific ML Stable, well-posed & efficient formulations	PRD6. ML-Enhanced Decision Support Automated decision-support, optimization, resilience, & control for complex systems & processes

- **Lens of Scientific Computing & Applied Mathematics → Scientific ML**
- **Capabilities Research: “Taxonomy” & PRDs for major use cases**
- **Foundational Themes: Basic research is essential**
- **Core technologies for use of Artificial Intelligence in scientific context**



History: DOE Applied Math Base Program & Research Initiatives are Key Foundations for Scientific Machine Learning

DOE Applied Math Base Program ⇒ **Foundational Themes in Scientific ML**

Fundamental research in robust & stable formulations, Data-intensive analysis, Multi-physics & multi-scale models, Scalable linear algebra & solvers, Optimization under uncertainty, UQ, etc

DOE Applied Math Research Initiatives

Scientific Inference & Data Analysis ⇒ **Data-Intensive Scientific ML**

- **2009** - Mathematics for Analysis of Petascale Data
- **2013** - DOE Data-Centric Science at Scale

Multiscale Models & Algorithms ⇒ **ML-Enhanced Modeling & Simulations**

- **2005** - Multiscale Mathematics Research and Education
- **2008** - Multiscale Mathematics for Complex Systems (also MMICCs in **2012, 2017, 2018**)

Integrated Capabilities for Complex Systems ⇒ **ML-Enhanced Decision Support**

- **2009** - Mathematics for Complex, Distributed, Interconnected Systems
- **2010 & 2013** - Uncertainty Quantification for Complex Systems; UQ for Extreme-Scale Science
- **2012, 2017, 2018** - Mathematical Multifaceted Integrated Capability Centers

Scientific Machine Learning will leverage basic research investments, widespread Science & Energy use cases, & DOE workforce expertise.



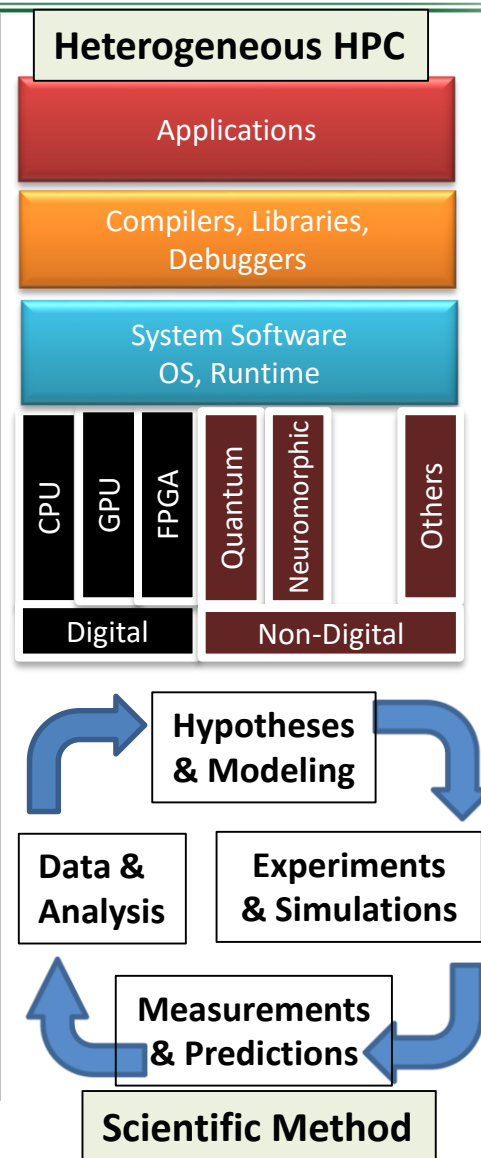
Strategic Vision for ASCR's Research Program

Emerging trends are pointing to a future that is increasingly

1. **Instrumented:** Massive data, high-tech sensors, detectors, satellites
2. **Interconnected:** Internet of Things, heterogeneous & composable resources
3. **Automated:** Machine learning for complex processes, real-time requirements
4. **Accelerated:** Faster & flexible research pathways for science & research insights

What is the role of ASCR's Research Program in transforming the way we carry out energy & science research?

1. **Post-Moore technologies:** Need basic research in new algorithms, software stacks, and programming tools for quantum and neuromorphic systems
2. **Extreme Heterogeneity:** Need new software stacks, programming models to support the heterogeneous systems of the future
3. **Adaptive Machine Learning, Modeling, & Simulation for Complex Systems:** Need algorithms and tools that supports automated decision making for intelligent operating systems, in situ workflow management, resilient infrastructure, & improved operational capabilities
4. **Uncertainty Quantification:** Need basic research in uncertainty quantification and artificial intelligence to enable statistically and mathematically rigorous foundations for advances in science domain-specific areas
5. **Data Tsunami:** Need to develop the software and coordinated infrastructure to accelerate scientific discovery by addressing challenges and opportunities associated with research data management, analysis, and reuse



DOE Applied Math develops mathematical foundations & computational building blocks for accelerated scientific insights

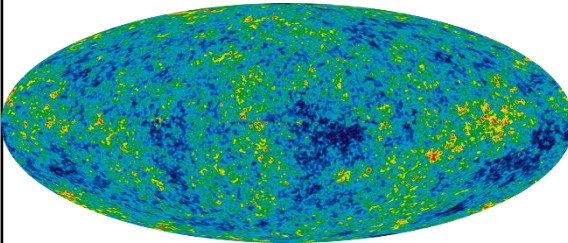
Scientific progress will be driven by

- Massive Data: sensors, simulations, networks
- Predictive Models & Adaptive Algorithms
- Heterogeneous High-Performance Computing
- Scientific Machine Learning & AI.

Trend: Human-AI collaborations will transform the way science is done.

Exemplars of Scientific Achievement

Cosmic Microwave Background

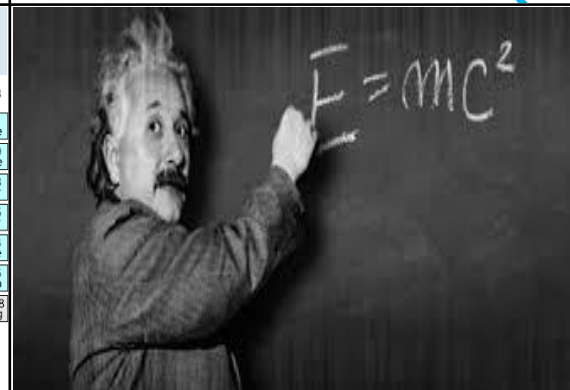


DNA Structure



Periodic Table of the Elements

Group	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
1	1 H																	2 He
2	3 Li	4 Be											5 B	6 C	7 N	8 O	9 F	10 Ne
3	11 Na	12 Mg											13 Al	14 Si	15 P	16 S	17 Cl	18 Ar
4	19 K	20 Ca	21 Sc	22 Ti	23 V	24 Cr	25 Mn	26 Fe	27 Co	28 Ni	29 Cu	30 Zn	31 Ga	32 Ge	33 As	34 Se	35 Br	36 Kr
5	37 Rb	38 Sr	39 Y	40 Zr	41 Nb	42 Mo	43 Tc	44 Ru	45 Rh	46 Pd	47 Ag	48 Cd	49 In	50 Sn	51 Sb	52 Te	53 I	54 Xe
6	55 Cs	56 Ba	57 La	72 Hf	73 Ta	74 W	75 Re	76 Os	77 Ir	78 Pt	79 Au	80 Hg	81 Tl	82 Pb	83 Bi	84 Po	85 At	86 Rn
7	87 Fr	88 Ra	89 Ac	104 Rf	105 Db	106 Sg	107 Bh	108 Hs	109 Mt	110 Ds	111 Rg	112 Cn	113 Nh	114 Fl	115 Mc	116 Lv	117 Ts	118 Og
	58 Ce	59 Pr	60 Nd	61 Pm	62 Sm	63 Eu	64 Gd	65 Tb	66 Dy	67 Ho	68 Er	69 Tm	70 Yb	71 Lu				
	90 Th	91 Pa	92 U	93 Np	94 Pu	95 Am	96 Cm	97 Bk	98 Cf	99 Es	100 Fm	101 Md	102 No	103 Lr				



Human-AI insights enabled via scientific method, experimentation, & AI reinforcement learning.



Long-Term Scientific Computing Research: DOE Applied Mathematics Themes

ASCR Criteria for Themes I - III

Basic research will develop & sustain

1. Computational Leadership
2. Discovery-Enabling Technologies
3. S/W Tools, Prototypes, Ecosystems
4. High-Tech DOE Workforce in advanced scientific computing.

I. Intelligent Automation & Decision Support

- Synergistic Human-AI collaboration.
Management of complex systems & processes.
Integration of foundational R&D.
- * AI-enhanced scientific method & discoveries
 - * Optimal experimental design
 - * Systems resilience, reliability, & control
 - * Automated performance optimization

II. Predictive Multifaceted Modeling & Simulations

- Predictive scientific computing.
Hierarchical, coupled, & hybrid simulations.
Novel formulations, meshing, & interfaces.
- * Smart machine learning-enhanced models
 - * Data-driven models, Surrogate sub-models
 - * Uncertainty & Error propagation, Validation
 - * Forward, Adjoint, & Parameter sensitivities

III. Adaptive High-Performance Solvers & Algorithms

- Adaptive HPC solvers & Scalable data analysis.
Forward, Inverse, & Optimization problems.
Post-Moore computational motifs/patterns.
- * Massively parallel, Asynchronous, Ensembles
 - * Preconditioners, Statistics, Learning from data
 - * Randomized, Graphs, Blackbox/Legacy codes
 - * Poly-algorithms, Adaptive-precision arithmetic

DOE Scientific Machine Learning and AI: Summary & Outlook

Machine Learning is a powerful scientific enabling technology

- More than Data. Also for Modeling, Complex Systems, & Science
- Scientific computing & mathematical foundations are essential
- Fast moving area → Need roadmap, blueprint, strategy
- Compelling: Re-visit ML, Re-think scientific computing uses

Pump is Primed for DOE leadership

- Roots from previous decade(s) of Applied Math basic research
- Ready: Researchers & expertise, Professional communities, etc

Future of Science & Energy Research

- Advanced technologies: More complex, more heterogeneous
- Automation & adaptivity needed for research breakthroughs & insights
- Scientific Machine Learning research is the basis for AI crosscutting initiative for accelerated scientific progress