# Thoughts on Building Machine Learning Models with Scientific Data

Jay Thiagarajan

# Data-Driven Methods are Rapidly Transforming Workflows in Science and Engineering

A U.S. Department of Energy initiative could refurbish existing supercomputers, turning them into high-performance artificial intelligence machines. U.S. DEPARTMENT OF ENERGY
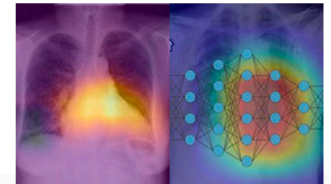
## Department of Energy plans major AI push to speed scientific discoveries

**National Institute of Biomedical Imaging and Bioengineering**
Creating Biomedical Technologies to Improve Health

## National experts chart roadmap for AI in medical imaging

Implications and opportunities for AI implementation in diagnostic medical imaging formulated in workshop report published in the journal, Radiology

Radiologists train for years to attain the skills to interpret subtle and not-so-subtle distinctions in medical images. Artificial intelligence (AI) is poised to make profound impact on their efforts, assisting human experts with computer-powered algorithms to recognize anatomical anomalies, enhance interpretation, and improve classification of medical imaging results.
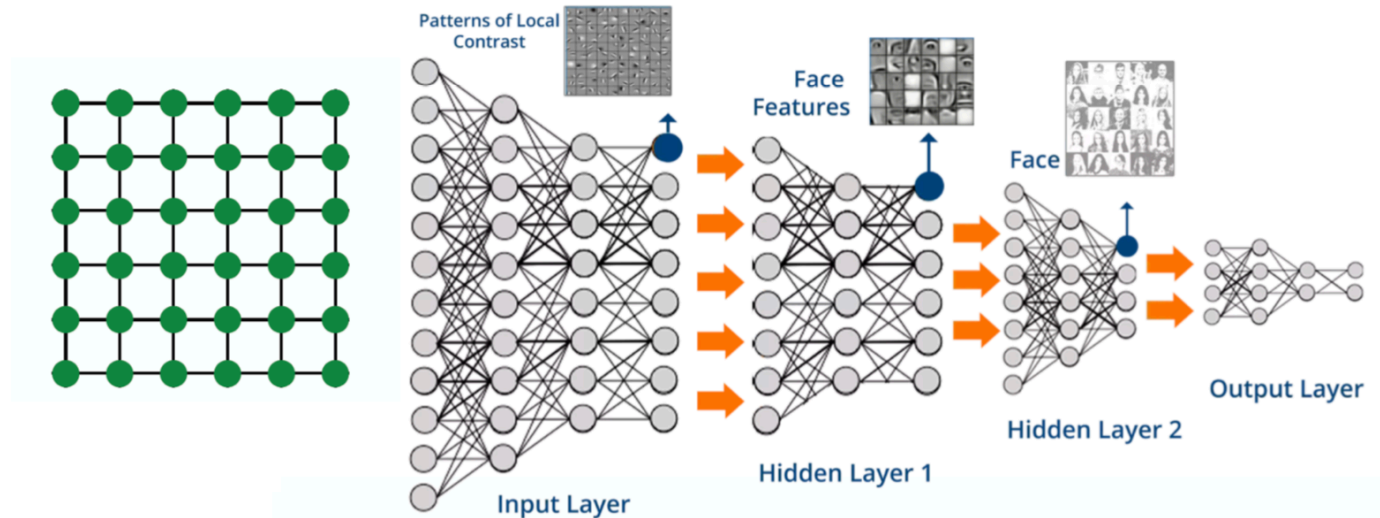
## Artificial Intelligence and Technology Office

Home » Science & Innovation » Artificial Intelligence and Technology Office

**AI** at **ENERGY.GOV**

## Vision:

Transform DOE into a world-leading AI enterprise by accelerating the research, development, delivery, and adoption of AI.

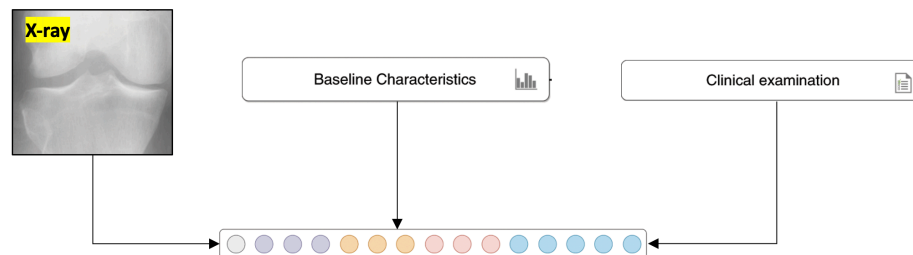# Modern Machine Learning Techniques are Highly Effective in Modeling Complex Data

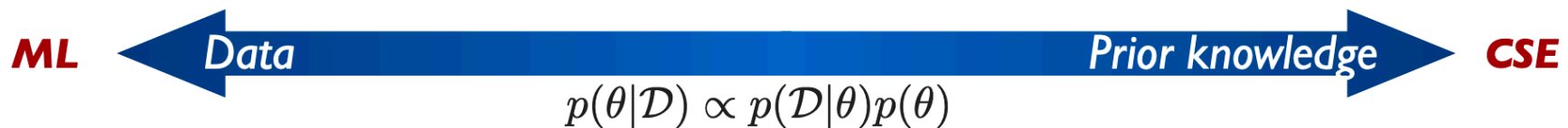*Data defined on spatial grids*

*Time-varying data*

*Multi-modal data*

# What is Different About Applying Machine Learning Tools to Scientific Data?

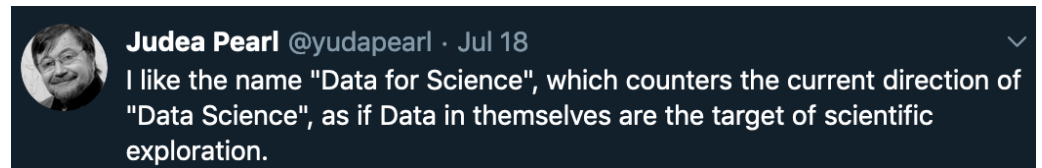Scientific problems present unique challenges, often requiring custom solutions:

- High cost of data acquisition

- High-dimensional observations and large parameter spaces

- Complex noise processes

- Uncertainty quantification

- Robust design/control

ML ⟷ Data          Prior knowledge ⟷ CSE

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$$

*A key step forward is to bridge conventional domain-informed modeling and machine learning to tackle these challenges*

# This Talk...

Insights from our experience in building predictive models for scientific data.



*"Though model choices, learning methods and constraints are highly specific to the application and data characteristics, there are useful ideas to share!*

*Use representation learning to explore and study complex relationships in scientific data!*

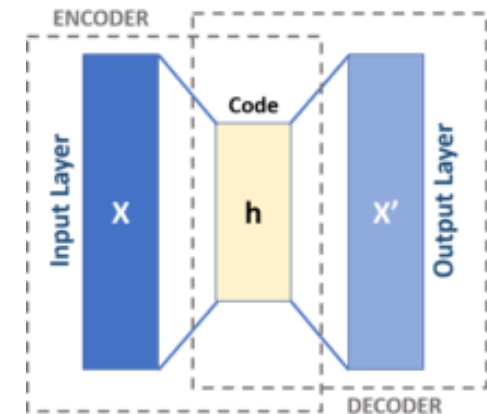## Identifying the True Generative Factors of Data Can Help Build Reliable Predictive Models

Representation learning allows us to infer latent features that succinctly describe the governing physical process.

Desired properties: (i) task-agnostic; (ii) low-dimensional; (iii) robust to data noise; (iv) preserves key relationships; (v) disentangled factors; (vi) known generative process.

A long-standing problem in machine learning – since the advent of principal component analysis.

# Representation Learning Methods can be Grouped into Three Categories

(i) <u>Generative</u> – Ability to sample data using a parametric/non-parametric generative model with low number of latent factors.



(ii) <u>Context prediction</u> – Predict missing information using its context.
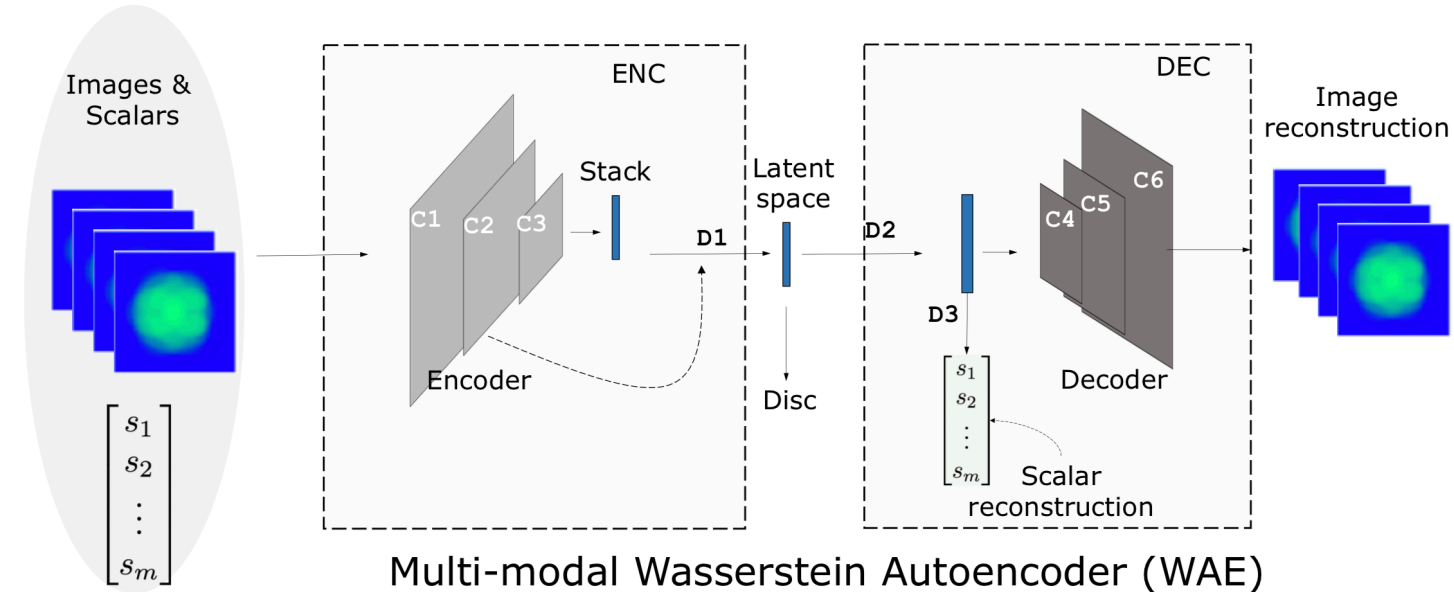


(iii) <u>Contrastive</u> – Build representations by learning to contrast.
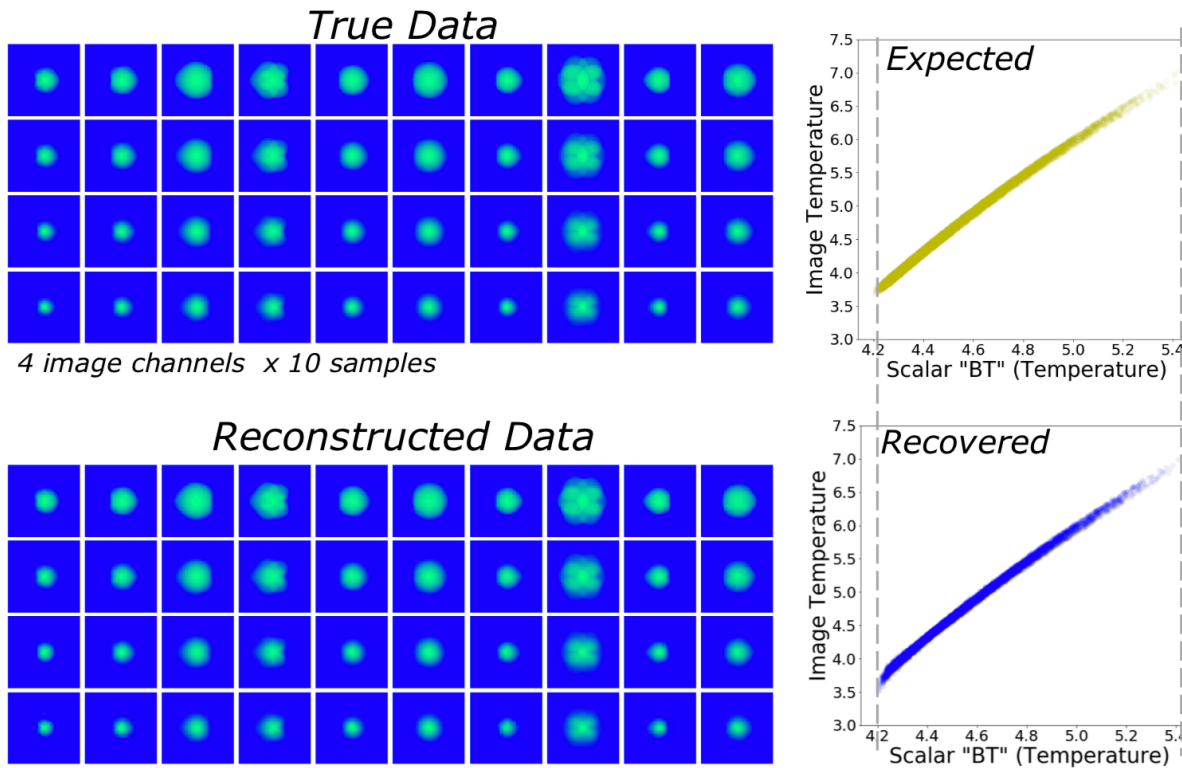
# Explore Learned Representations using Scientific Priors

**Example 1**: Multi-modal measurements from an inertial confinement fusion simulator



Autoencoding with a reconstruction loss

# Explore Learned Representations using Scientific Priors

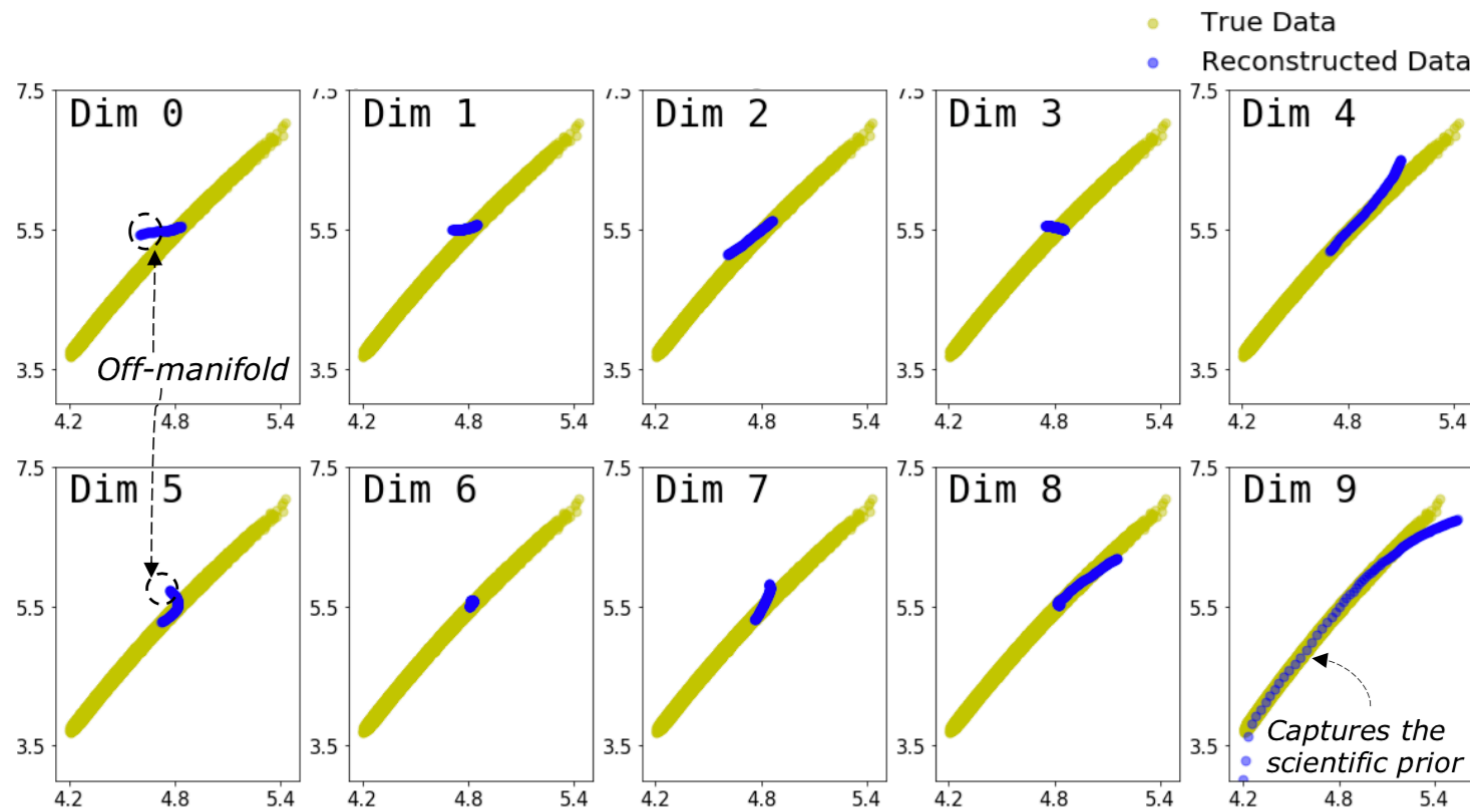**Example 1**: Multi-modal measurements from an inertial confinement fusion simulator



True Data

4 image channels  x 10 samples

Reconstructed Data

Evaluating the scientific prior on reconstructed and true samples

We use notions of thermal equilibrium to relate predicted ion temperatures to estimates of electron temperature formed by ratios of x-ray image brightness.

# Explore Learned Representations using Scientific Priors

**Example 1**: Multi-modal measurements from an inertial confinement fusion simulator



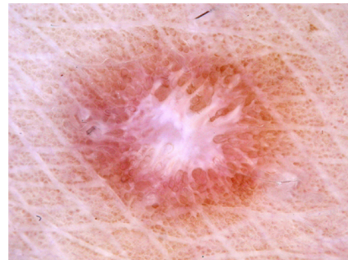Walking along different directions on the physics manifold → evaluating the scientific prior

# Explore Learned Representations using Scientific Priors

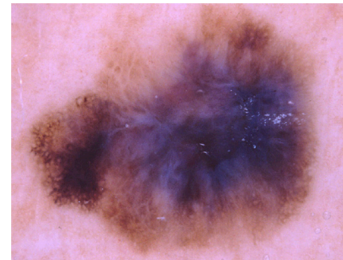**Example 2**: Dermoscopy images from subjects diagnosed with different types of skin lesions
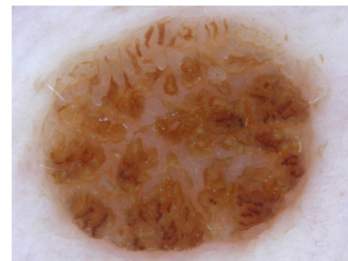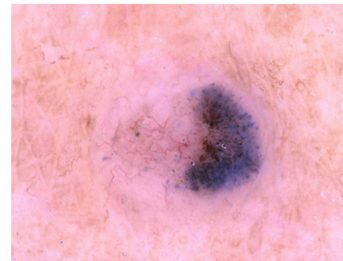
Nevus
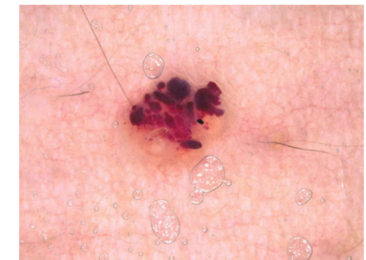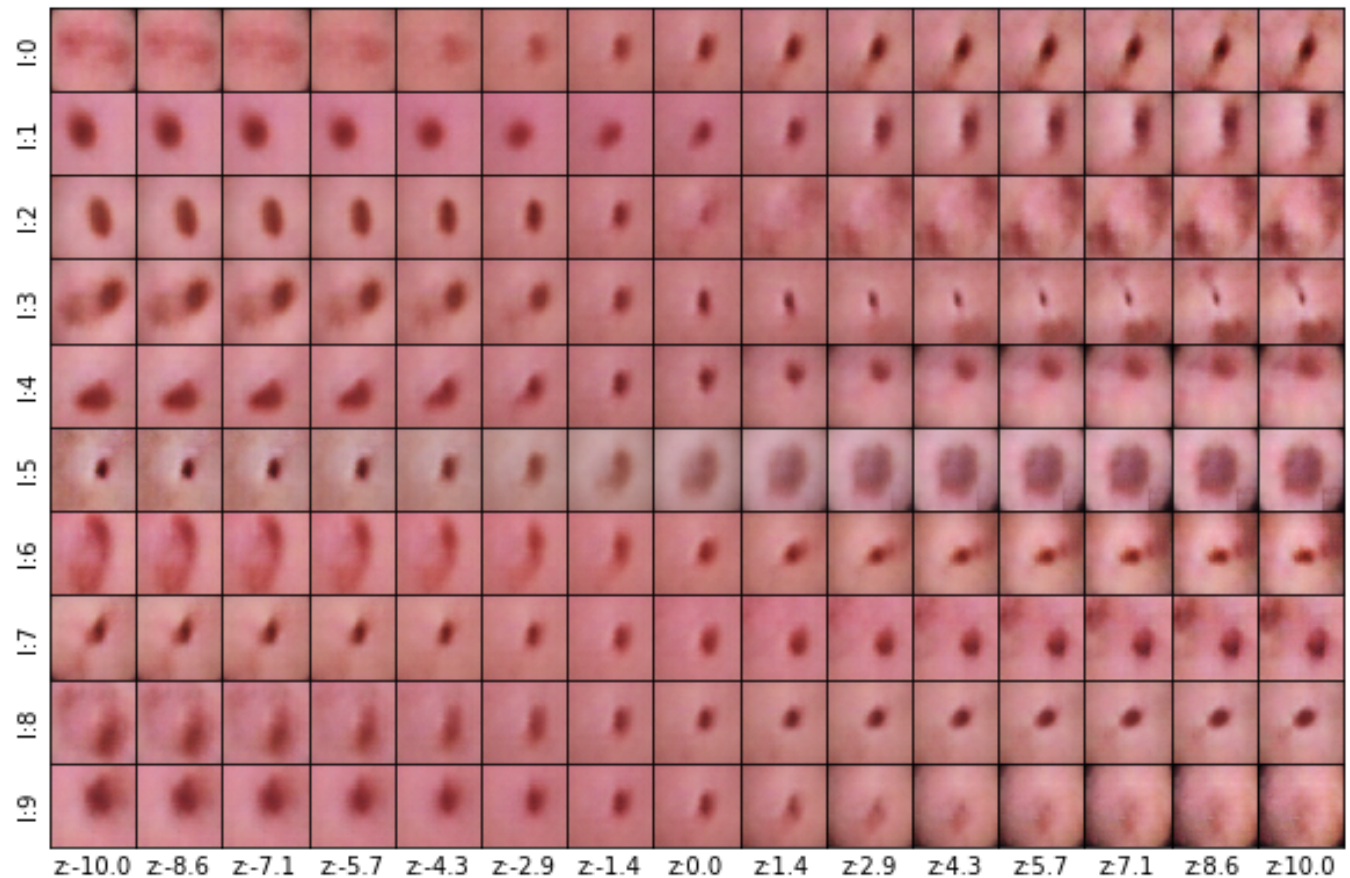
Dermatofibroma

Melanoma

Vascular

Pigmented Bowen's

Pigmented Benign Keratoses

Basal Cell Carcinoma

# Explore Learned Representations using Scientific Priors

**Example 2**: Dermoscopy images from subjects diagnosed with different types of skin lesions

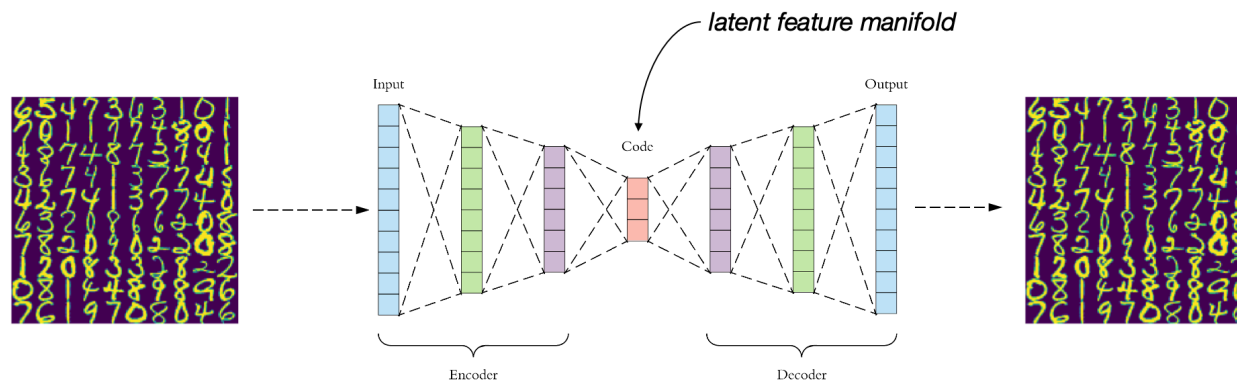Latent space traversal from a disentangled VAE
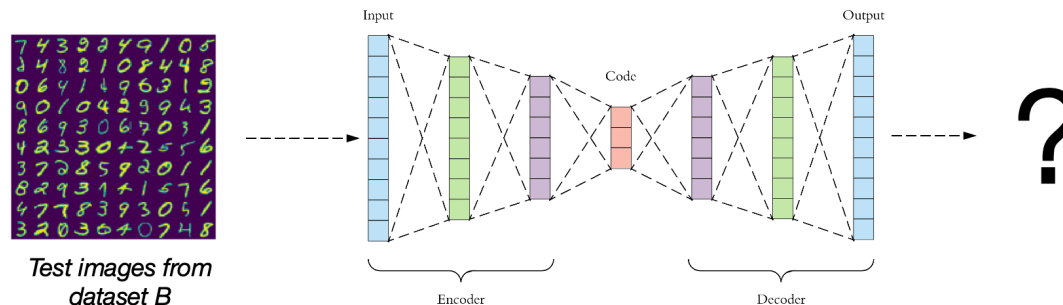
A: Asymmetry
B: Border
C: Color
D: Diameter

*Distribution shifts are features, not bugs!*

# Understanding Inductive Biases of ML Models is Critical to Characterizing their Behavior
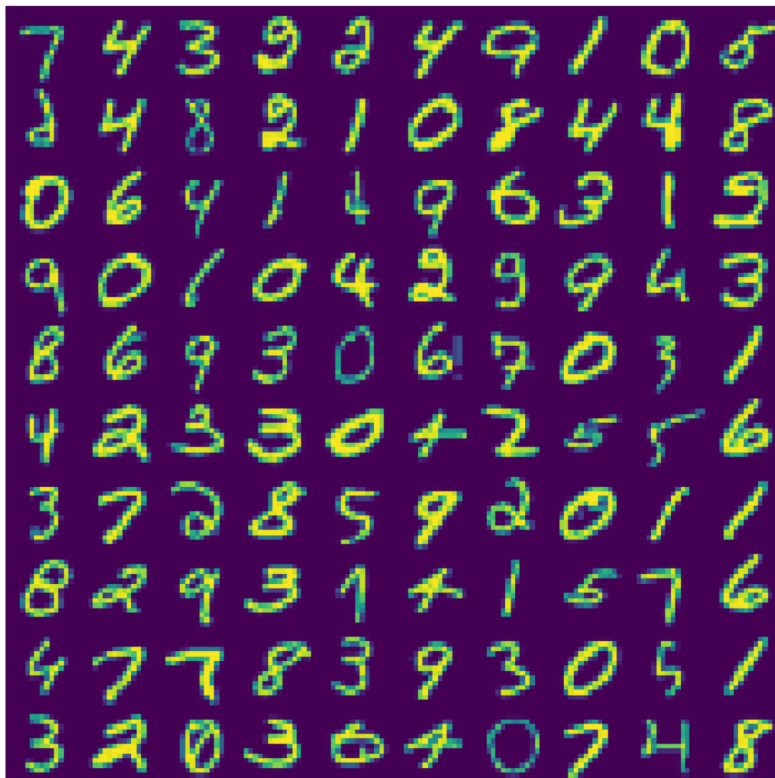
Data you expect vs. Data you get



latent feature manifold

Digits Autoencoder
*(trained on dataset A)*

Test images from dataset B

Digits Autoencoder

?

# Understanding Inductive Biases of ML Models is Critical to Characterizing their Behavior
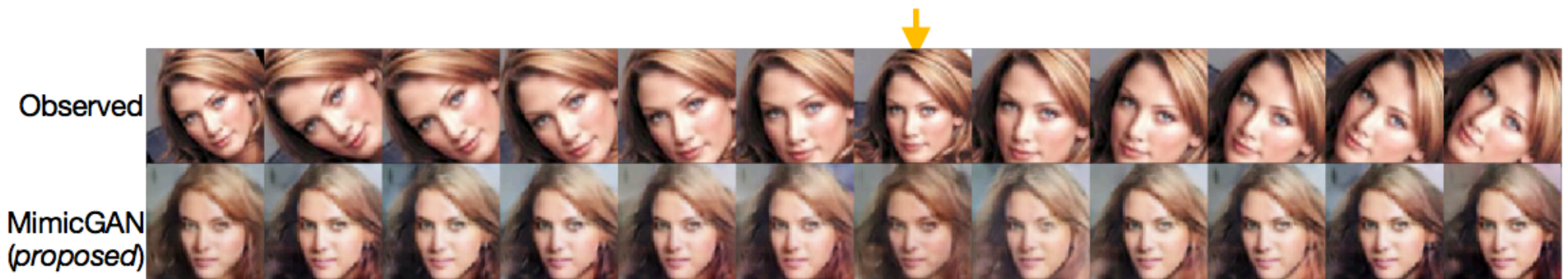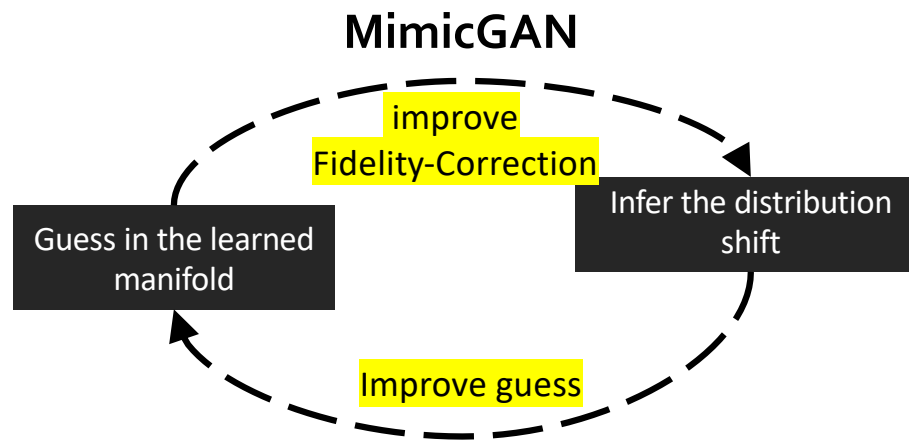
What we put in



What we got out

# Inferring the Unknown" Distribution Shifts is Essential to Effectively Utilize the Representations in Practice
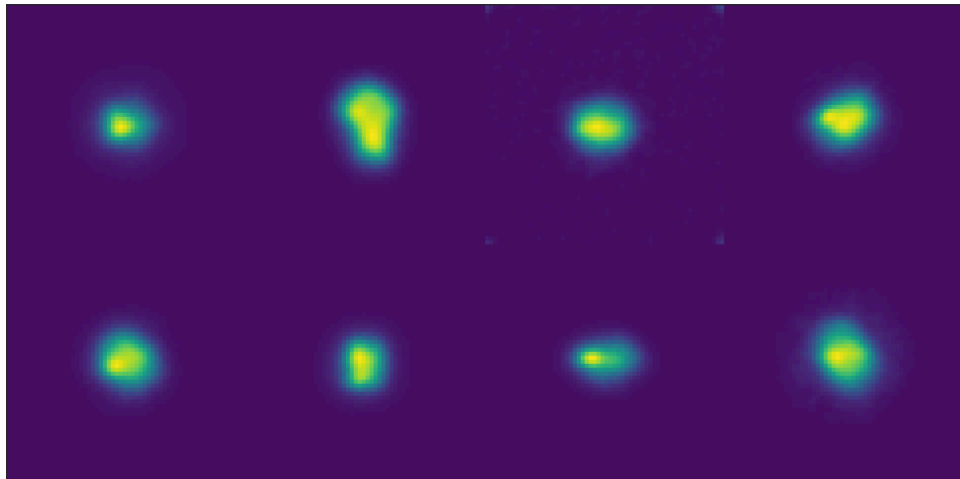
How should we transform "unseen" observations to look like they were produced by the "known" data generation process?
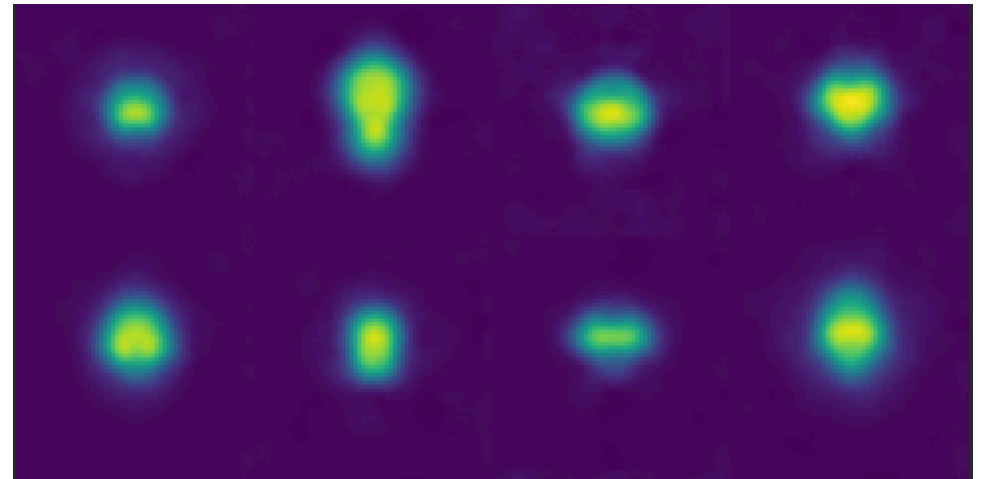
# Understanding Inductive Biases of ML Models is Critical to Characterizing their Behavior

**Example**: Matching real-experiments (X-ray images) to simulations from an inertial confinement fusion simulator (hydra).

Experiments – Asymmetric

Equivalent Simulations – Symmetric

*Your model is as good as your loss function!*

# Machine Learning Eventually Boils Down to Optimizing Parameters Based on a Fidelity Metric

Does patient A have cancer? → Does the patient's scan look different enough from the database of healthy subjects? → Minimize the cross-entropy loss

In many predictive modeling problems in sciences, we deal with continuous-valued targets and we use mean-squared loss as the objective.

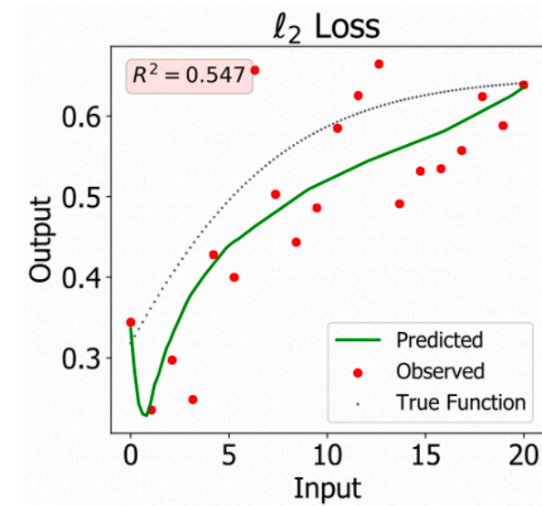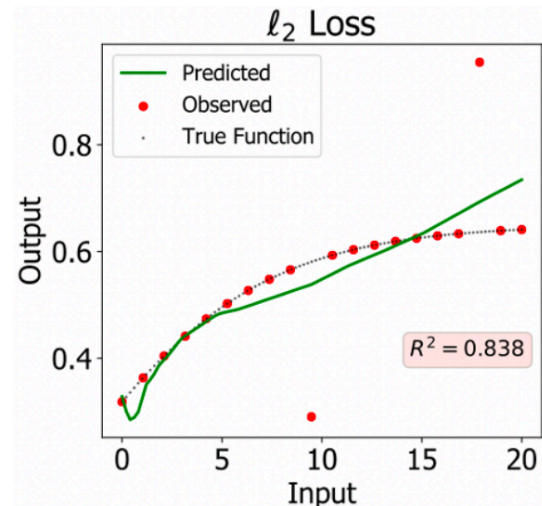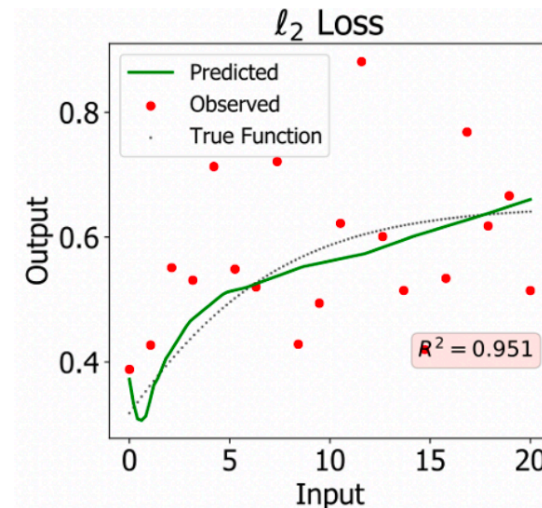Choice of loss function places a prior on the distribution of residuals

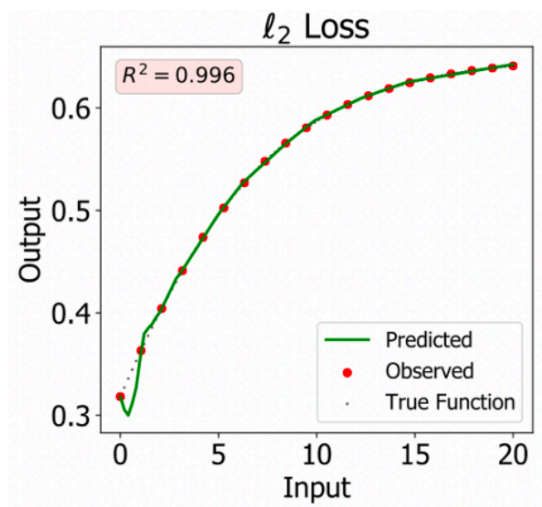$$\mathbf{y} - f(\mathbf{x})$$

L2 error optimal when the distribution is symmetric

Not robust when there are outliers

# Machine Learning Eventually Boils Down to Optimizing Parameters Based on a Fidelity Metric

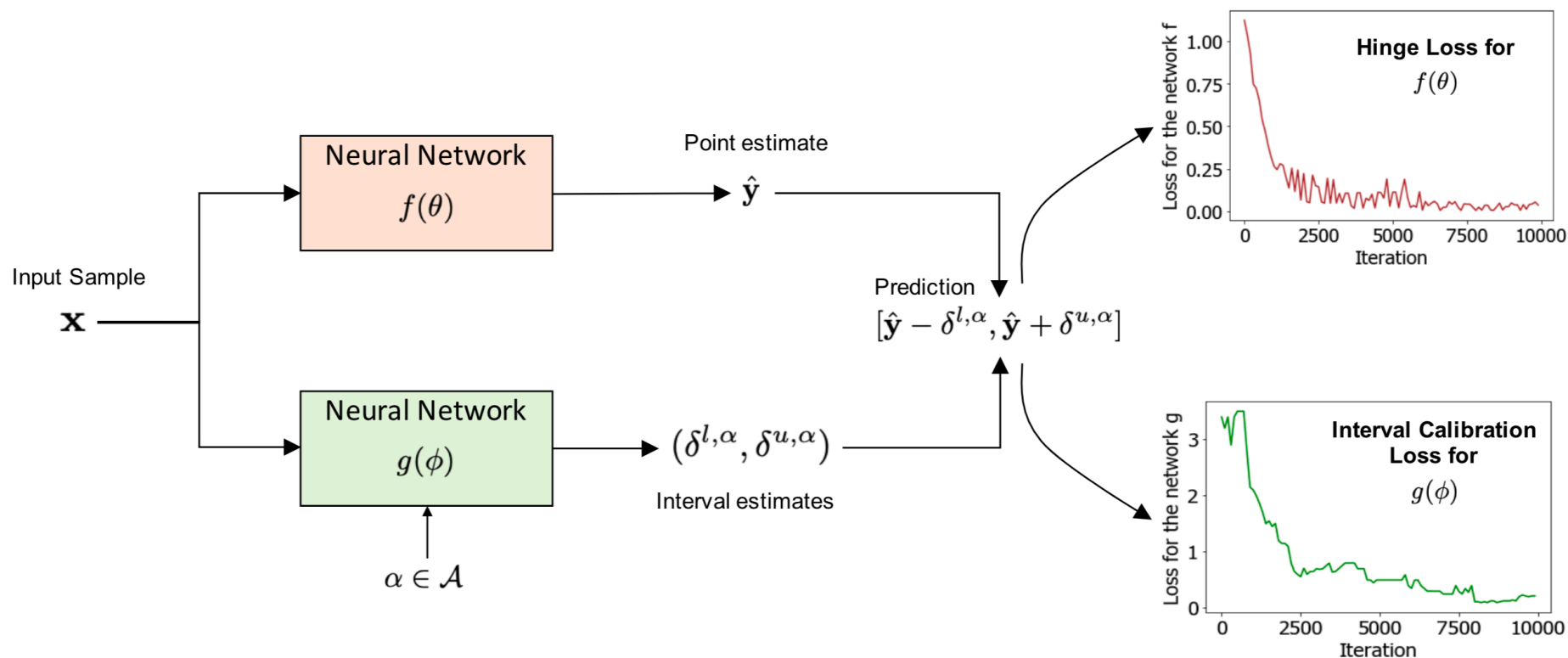**Example**: Approximating a synthetic function using a 1-layer MLP

# Learn-by-Calibrating: Prior-free Loss Function Design Based on Interval Calibration

If a model produces a prediction interval $[\hat{\mathbf{y}} - \delta^l, \hat{\mathbf{y}} + \delta^u]$

$$p\big((\hat{\mathbf{y}} - \delta^l) \leq \mathbf{y} \leq (\hat{\mathbf{y}} + \delta^u)\big) = \alpha \quad \longrightarrow \text{Confidence level}$$

# Machine Learning Eventually Boils Down to Optimizing Parameters Based on a Fidelity Metric

**Example**: Approximating a synthetic function using a 1-layer MLP
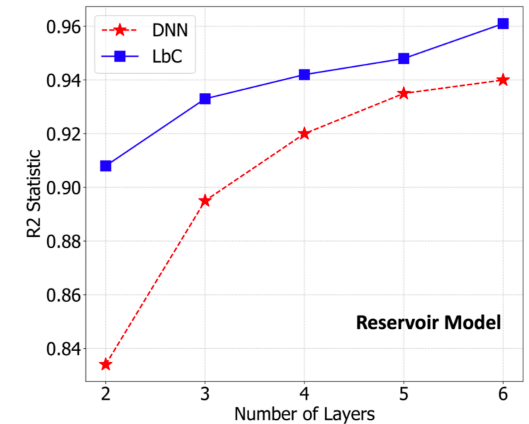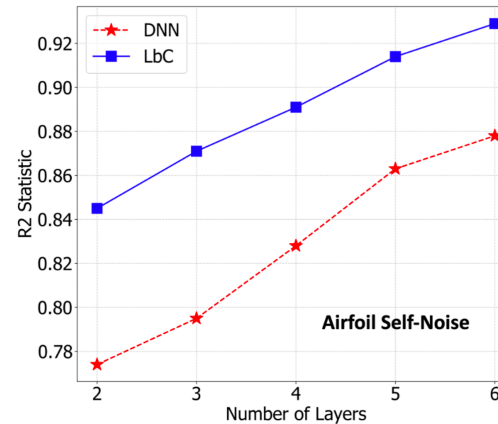
# Machine Learning Eventually Boils Down to Optimizing Parameters Based on a Fidelity Metric

LbC consistently outperforms existing symmetric losses when the residual distribution is skewed

LbC requires lesser number of model parameters to match the performance of a standard model

*Sometimes, we need to build our own architectures…*

# In Deep Learning, Model Architectures Act as a Prior in Defining the Hypothesis Space

Convolutional networks that enforce local smoothness act as a strong prior for images – Even untrained networks can provide useful feature representations.
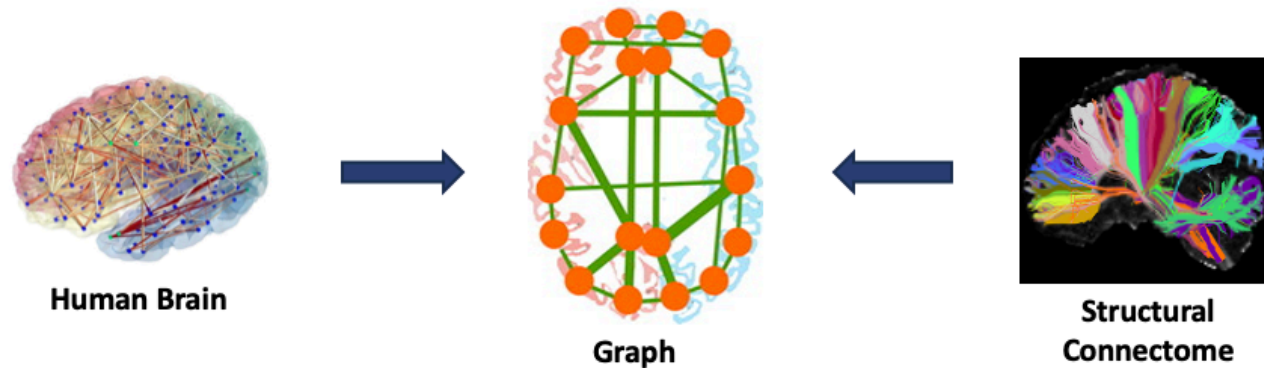
In scientific problems, sometimes we need to design custom architectures to provide strong inductive biases.

These solutions rarely generalize, however, can provide significant performance gains for the problem at hand.

# Designing Custom Model Architectures to Provide Strong Inductive Biases

**Example**: Modeling human brain connectomes for predictive modeling

Connectome is a comprehensive map of structural and functional connectivity of neural pathways in the brain.
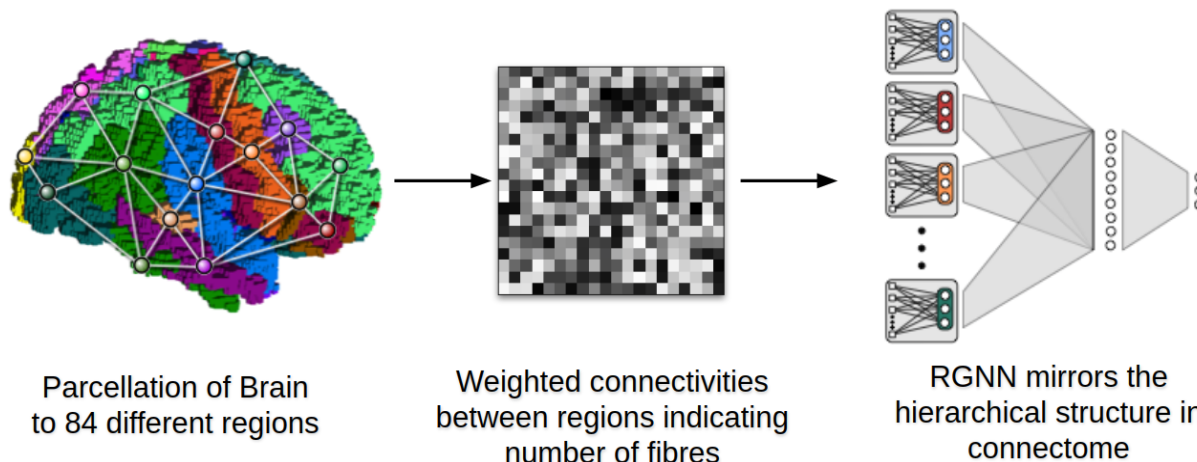


**Human Brain** → **Graph** ← **Structural Connectome**

*Connectome Graph generation from T1 images and probabilistic Tractography*

# Designing Custom Model Architectures to Provide Strong Inductive Biases

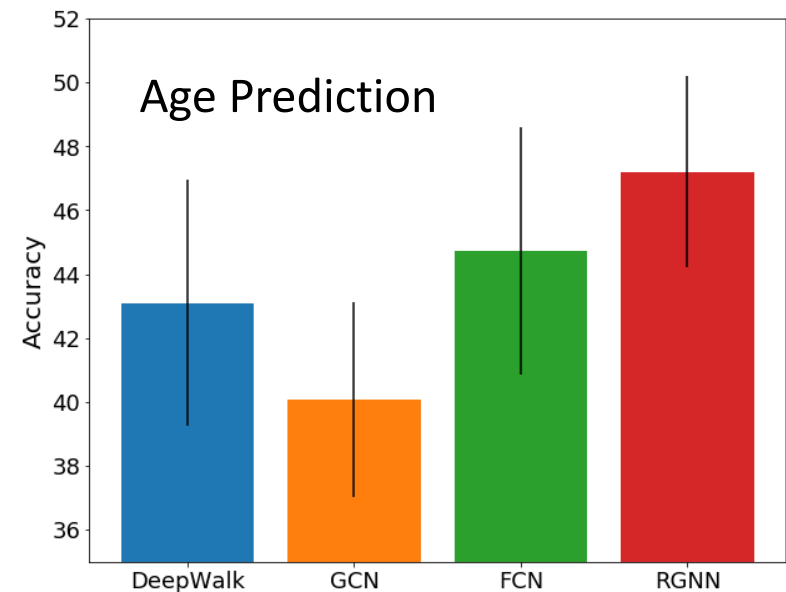**Example**: Modeling human brain connectomes for predictive modeling

Since connectomes can be viewed as a graph, one might be tempted to use graph neural networks (e. g. molecular chemistry) – Information diffusion property does not hold.



Parcellation of Brain to 84 different regions

Weighted connectivities between regions indicating number of fibres

RGNN mirrors the hierarchical structure in connectome

We proposed a novel message passing framework that is aptly suited to process relational structure without the need for diffusion.

# Designing Custom Model Architectures to Provide Strong Inductive Biases

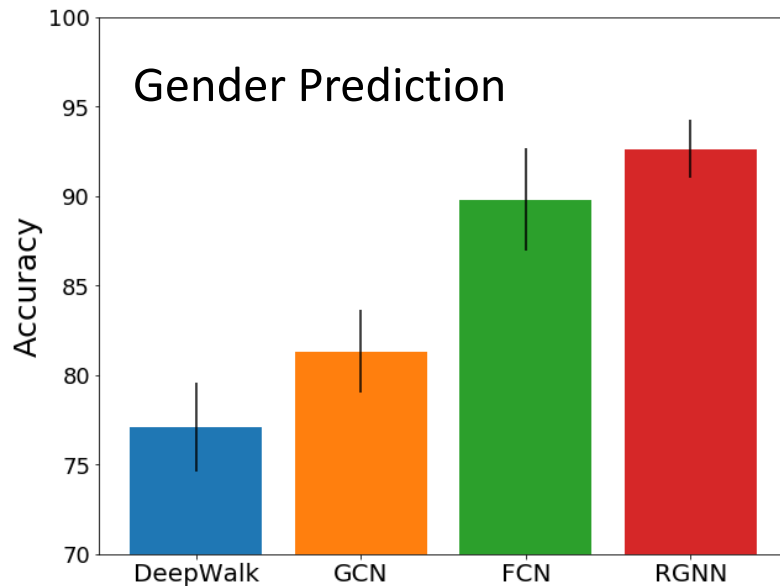**Example**: Modeling human brain connectomes for predictive modeling

*Do not just predict, but also introspect!*

**Using Machine Learning Models in Critical Applications Requires a Rigorous Characterization of its Behavior**
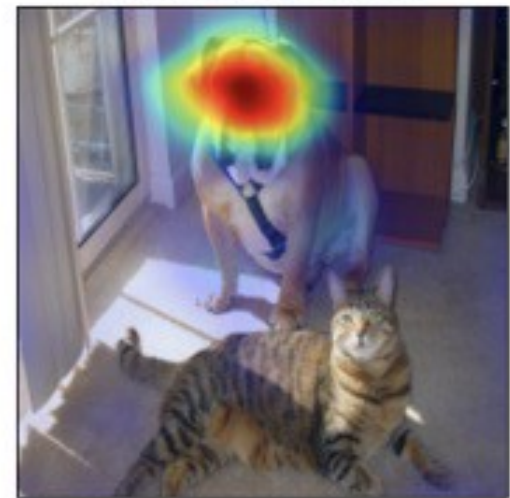
Gaining scientific insights from learned models needs introspection driven by *user-specified hypotheses*.

A broad class of interpretability techniques exist to uncover the "most plausible" explanations for a decision.
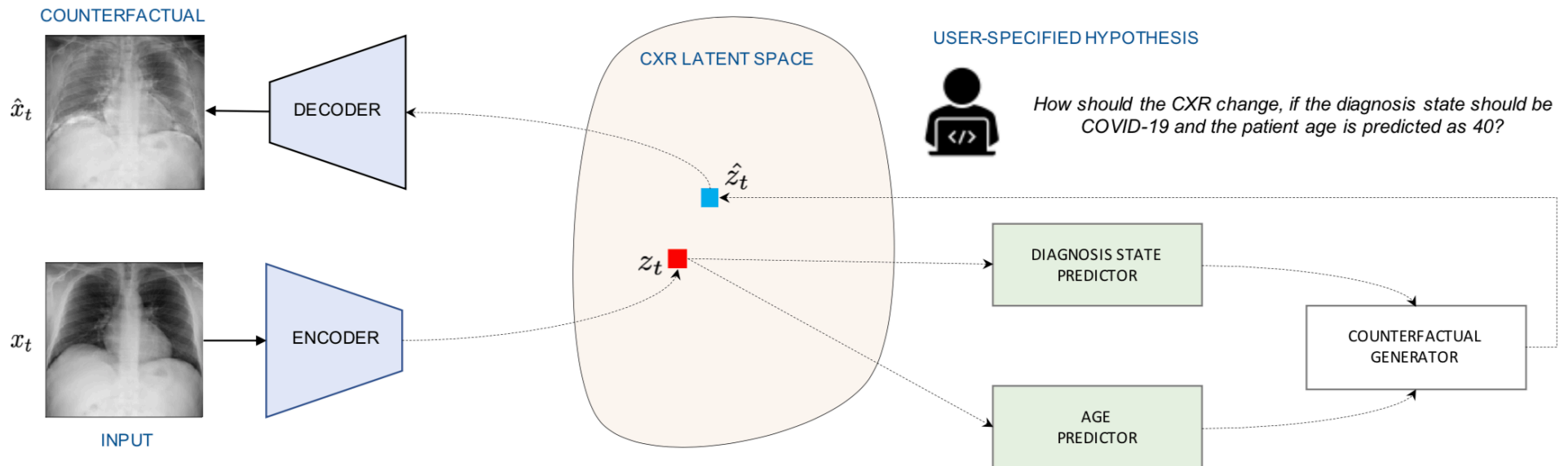


Grad-CAM for "Cat"

Grad-CAM for "Dog"

# Using Machine Learning Models in Critical Applications Requires a Rigorous Characterization of its Behavior
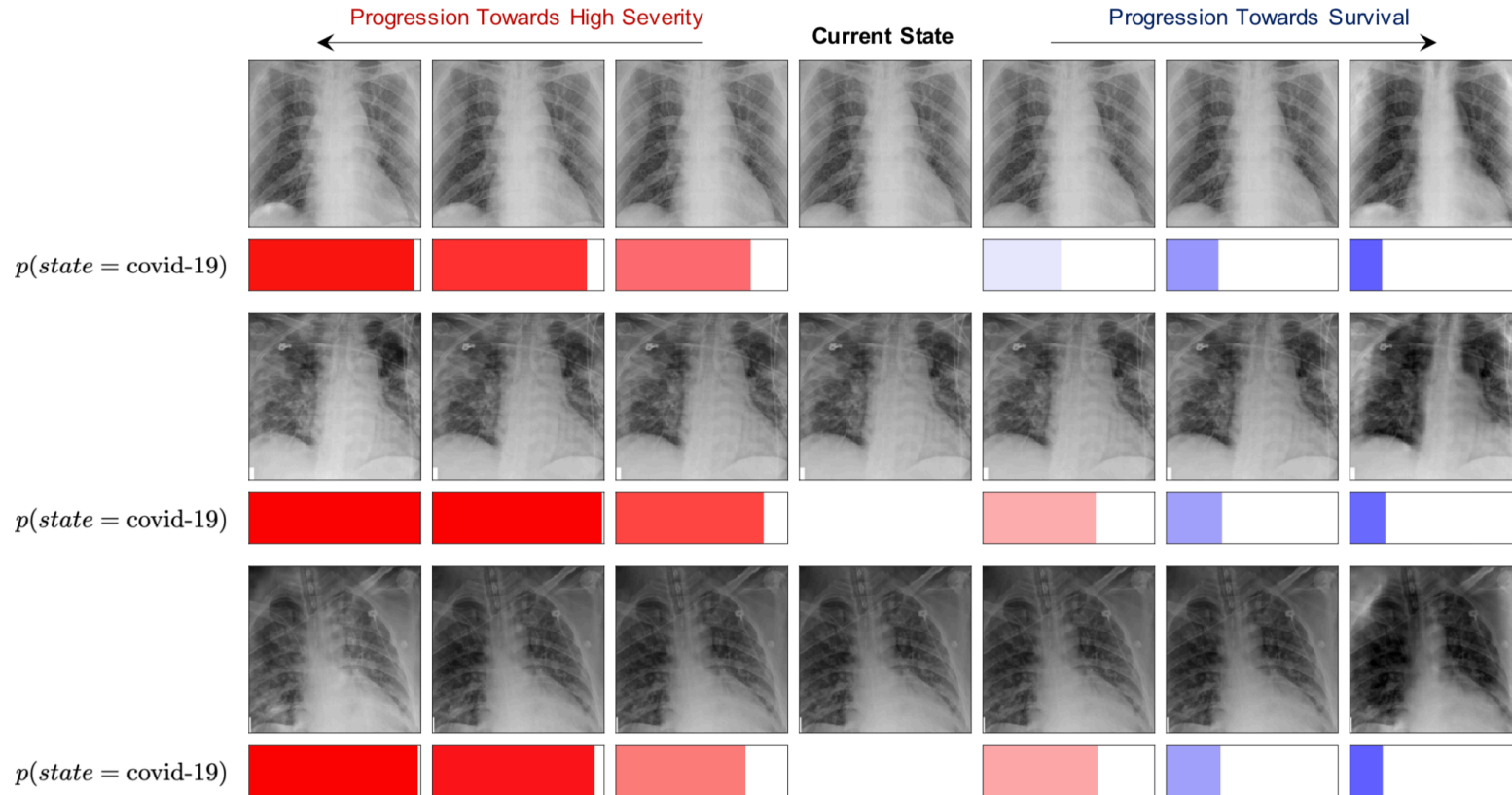
However, in science problems, we need to deal with more complex hypotheses.

**Example:** Analysis of COVID-19 infections from CXR data using counterfactual reasoning
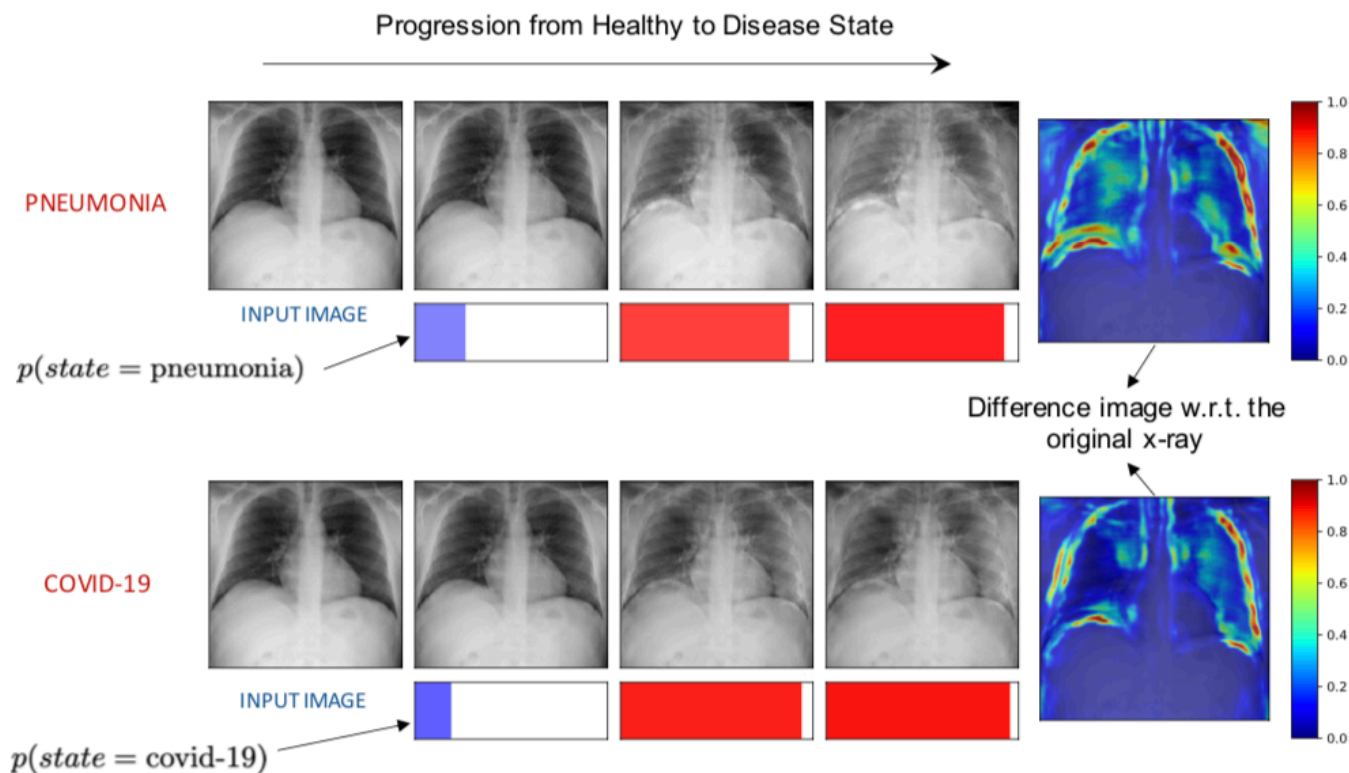
# Using Machine Learning Models in Critical Applications Requires a Rigorous Characterization of its Behavior

**Hypothesis:** The severity of the disease infection increases or decreases

# Using Machine Learning Models in Critical Applications Requires a Rigorous Characterization of its Behavior

**Hypothesis:** Manifestation of COVID-19 is different from other types of pneumonia known before



Progression from Healthy to Disease State

# Summary

Scientific data provide unique challenges and opportunities towards advancing machine learning

Using suitable machine learning methods will help accelerate scientific analysis and discovery

# Team



Rushil Anirudh

Bhavya Kailkhura

Timo Bremer

Vivek N

Bindya Venkatesh

Uday Shankar S

# Relevant Papers

## Improved surrogates in inertial confinement fusion with manifold and cycle consistencies

## Designing Accurate Emulators for Scientific Processes using Calibration-Driven Deep Models

Jayaraman J. Thiagarajan, Bindya Venkatesh, Rushil Anirudh, Peer-Timo Bremer, Jim Gaffney, Gemma Anderson, Brian Spears

## Exploring Generative Physics Models with Scientific Priors in Inertial Confinement Fusion

Rushil Anirudh*, Jayaraman J. Thiagarajan, Shusen Liu,
Peer-Timo Bremer, Brian K. Spears

Lawrence Livermore National Laboratory, Livermore, California.

**Abstract**

There is significant interest in using modern neural networks for scientific applications due to their effectiveness in modeling highly complex, non-linear problems in a data-driven fashion. However, a common challenge is to verify the scientific plausibility or validity of outputs predicted

## Designing Deep Inverse Models for History Matching in Reservoir Simulations

**Vivek Sivaraman Narayanaswamy***, **Jayaraman J. Thiagarajan**[†], **Rushil Anirudh**[†],
**Fahim Forouzanfar**[‡], **Peer-Timo Bremer**[†], **Xiao-Hui Wu**[‡]
*Arizona State University, [†]Lawrence Livermore National Laboratory,
[‡]ExxonMobil Upstream Research Company

**Abstract**

## Modeling Human Brain Connectomes using Structured Neural Networks

Uday Shankar Shanthamallu · Qunwei Li · Jayaraman J. Thiagarajan ·
Show all 6 authors · Peer-Timo Bremer

## MimicGAN: Robust Projection onto Image Manifolds with Corruption Mimicking

Rushil Anirudh ✉, Jayaraman J. Thiagarajan, Bhavya Kailkhura & Peer-Timo Bremer

# Questions?