

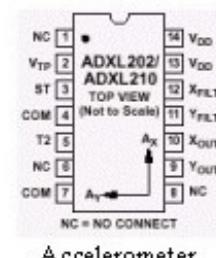
# ENM 540: Data-driven modeling and probabilistic scientific computing

## *Lecture #1: Introduction to data-driven modeling*

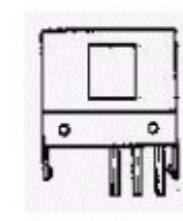
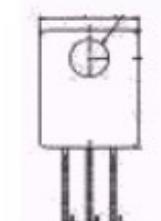
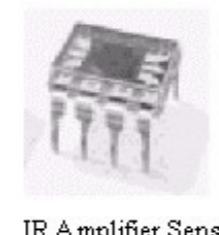
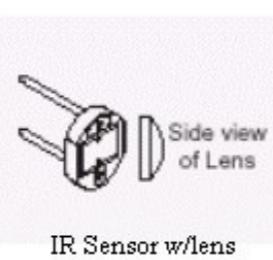
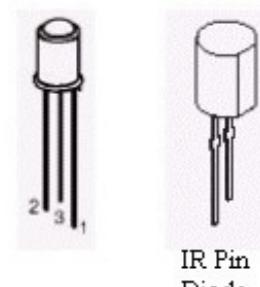
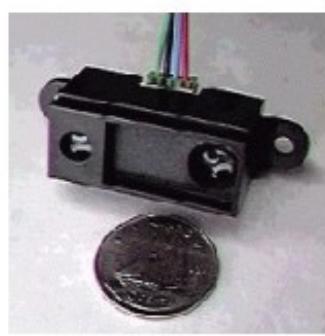
Paris Perdikaris  
January 11<sup>th</sup>, 2018



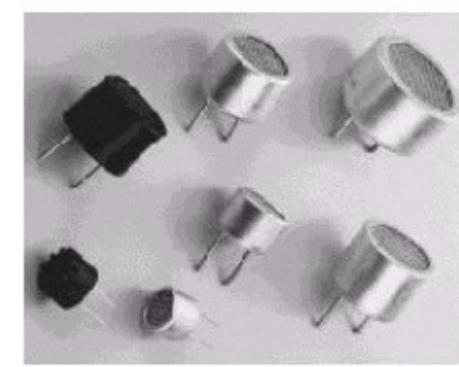
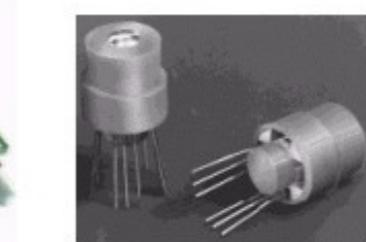
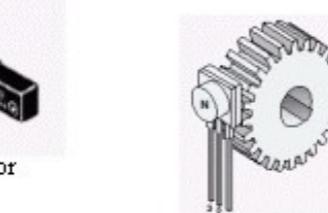
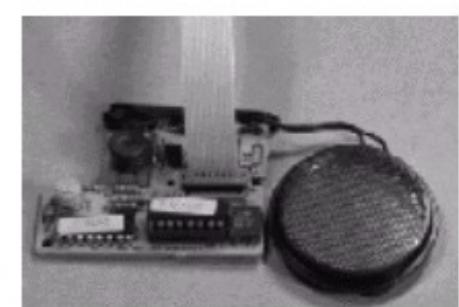
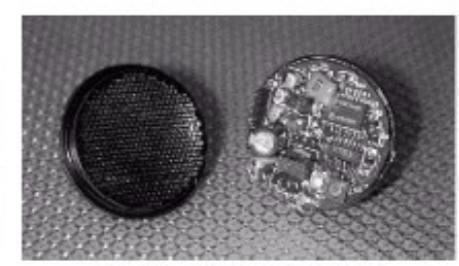
# Roadmap to Trillion Sensors



Accelerometer



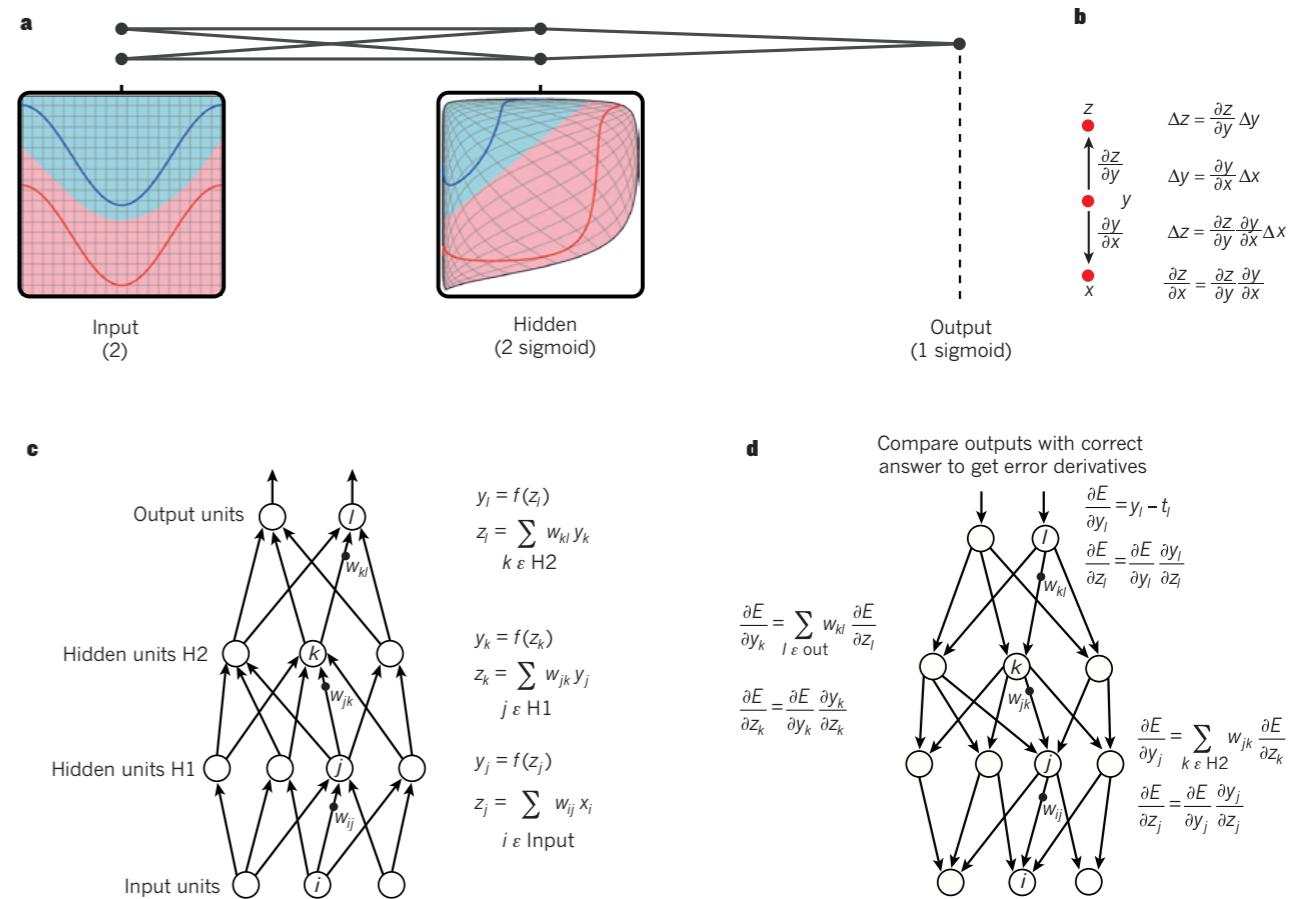
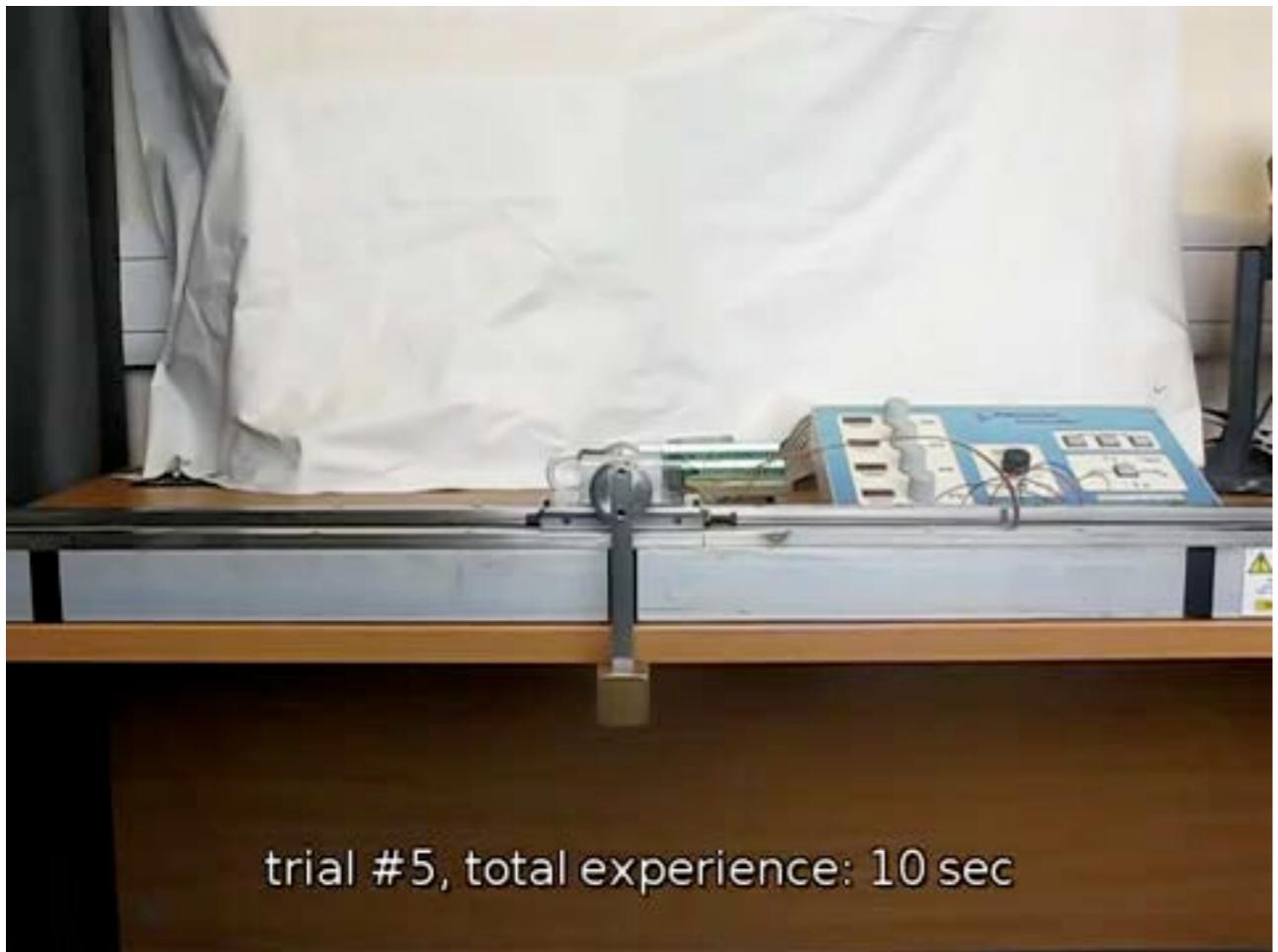
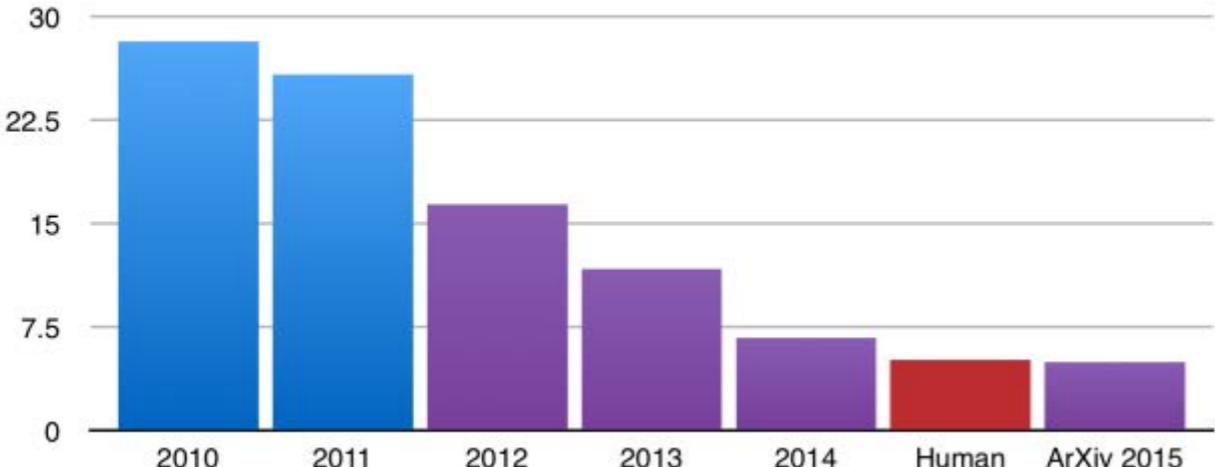
## Sensors Overview



Sensors are expected to generate a bronto-bytes (1000 trillion trillion) of data by 2025!

# Recent success of machine learning

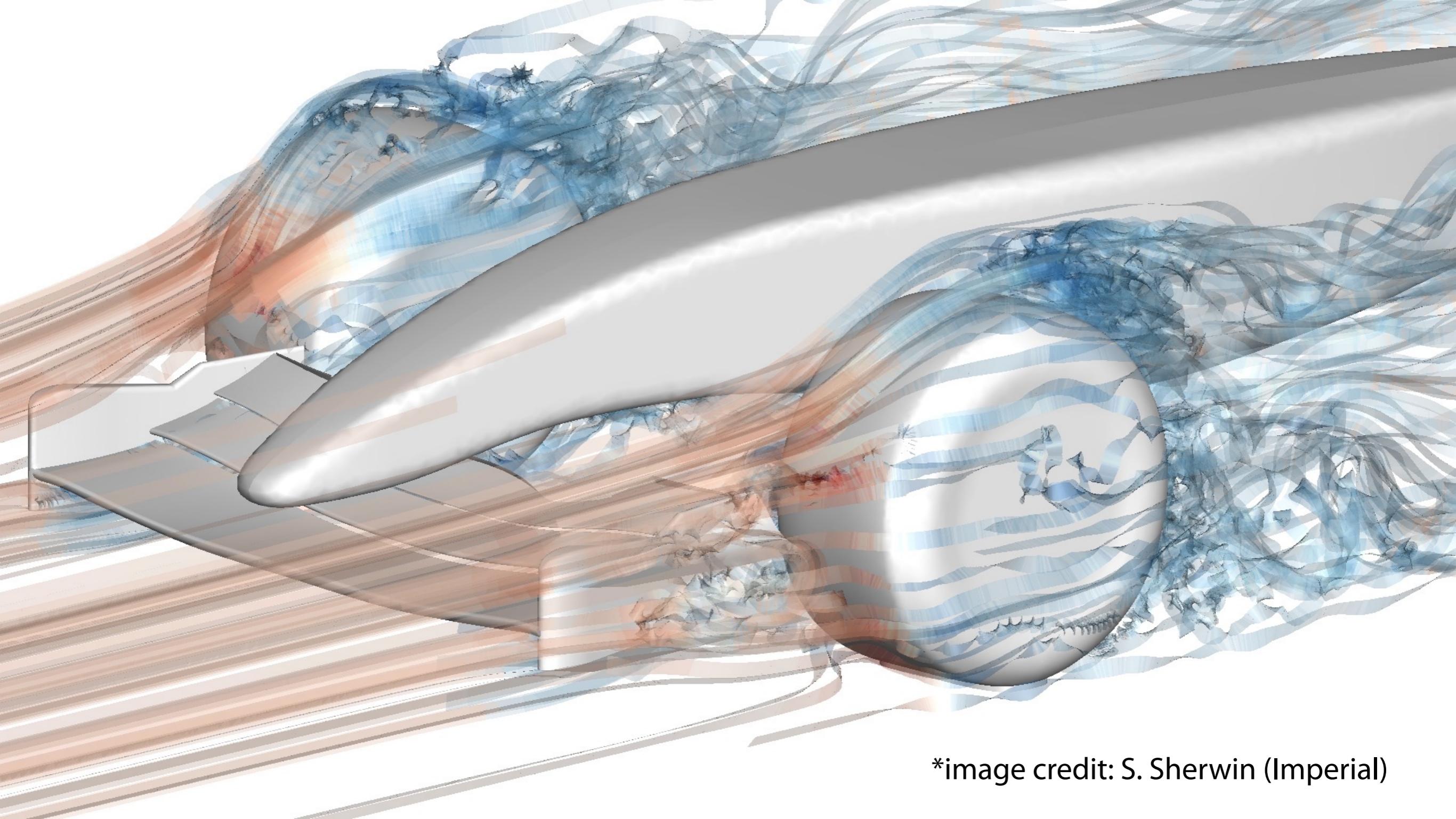
ILSVRC top-5 error on ImageNet



*All models are wrong  
but some are useful*



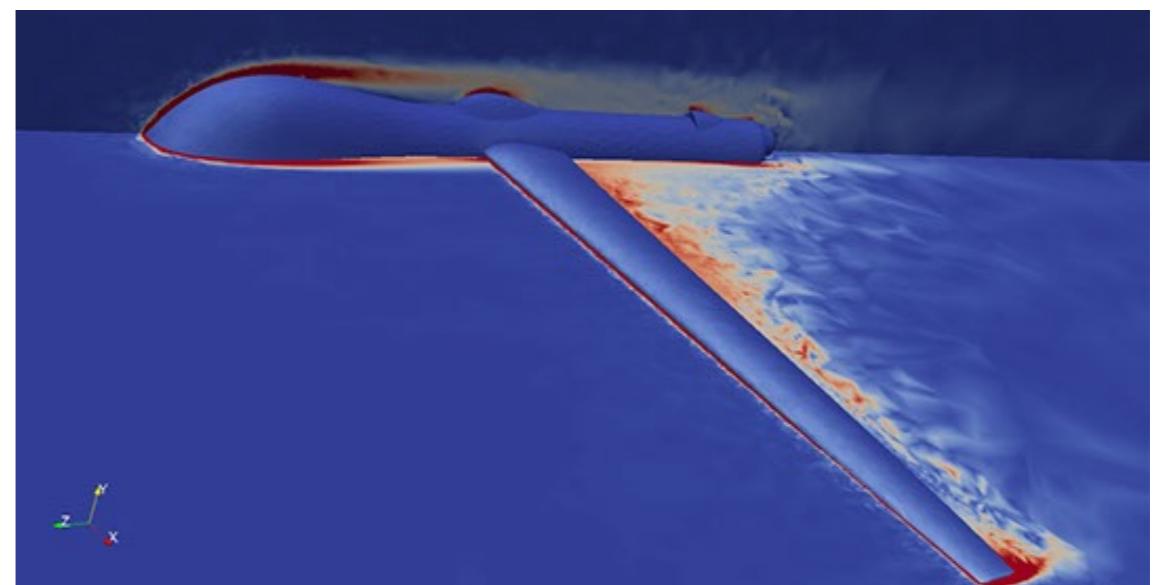
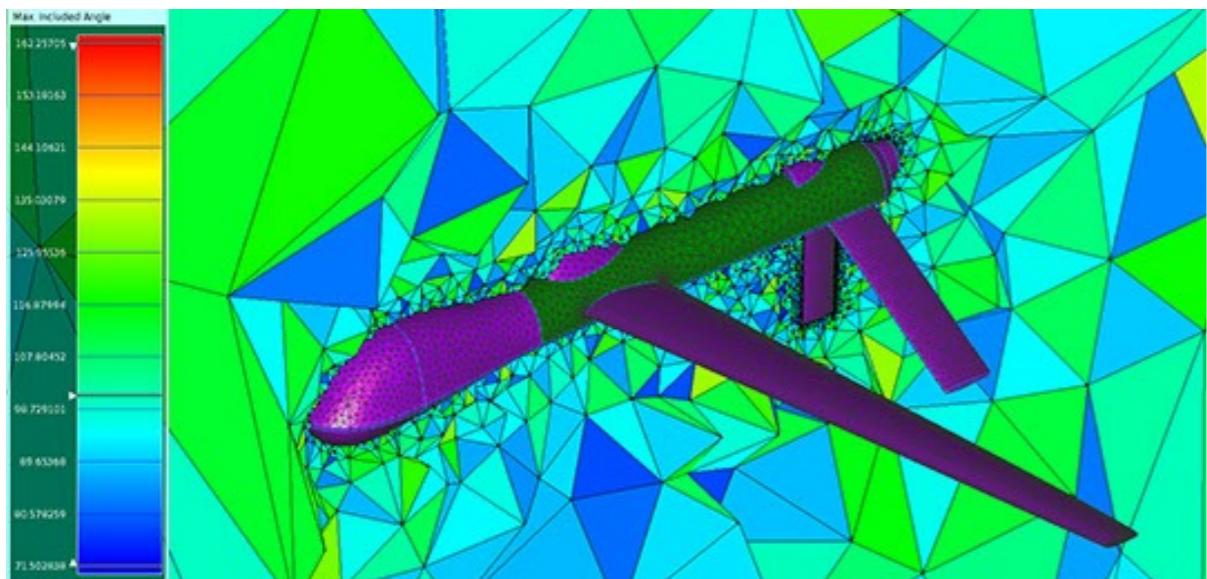
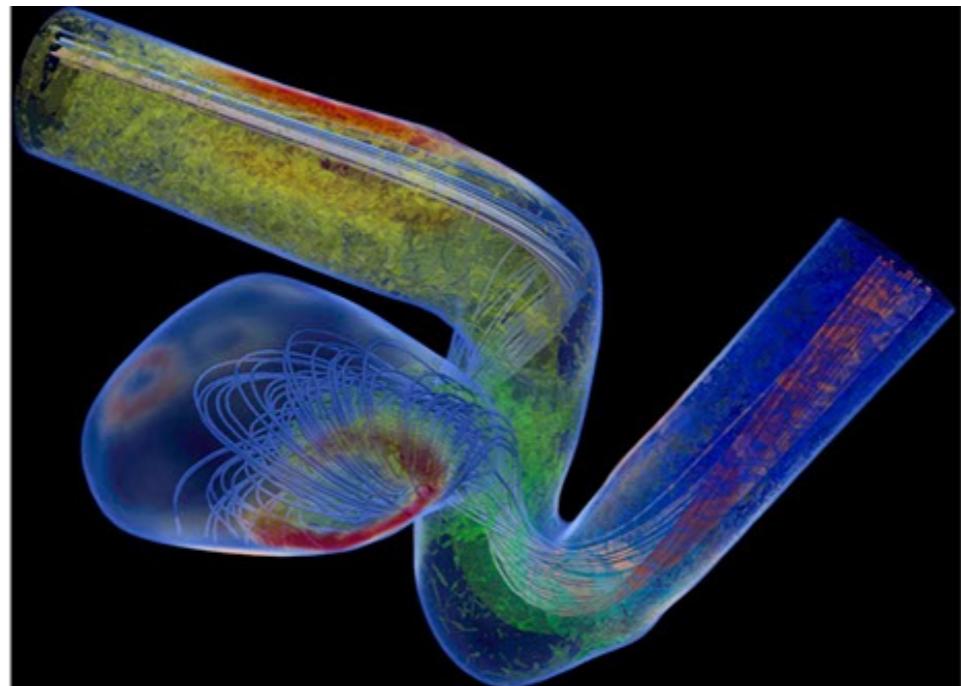
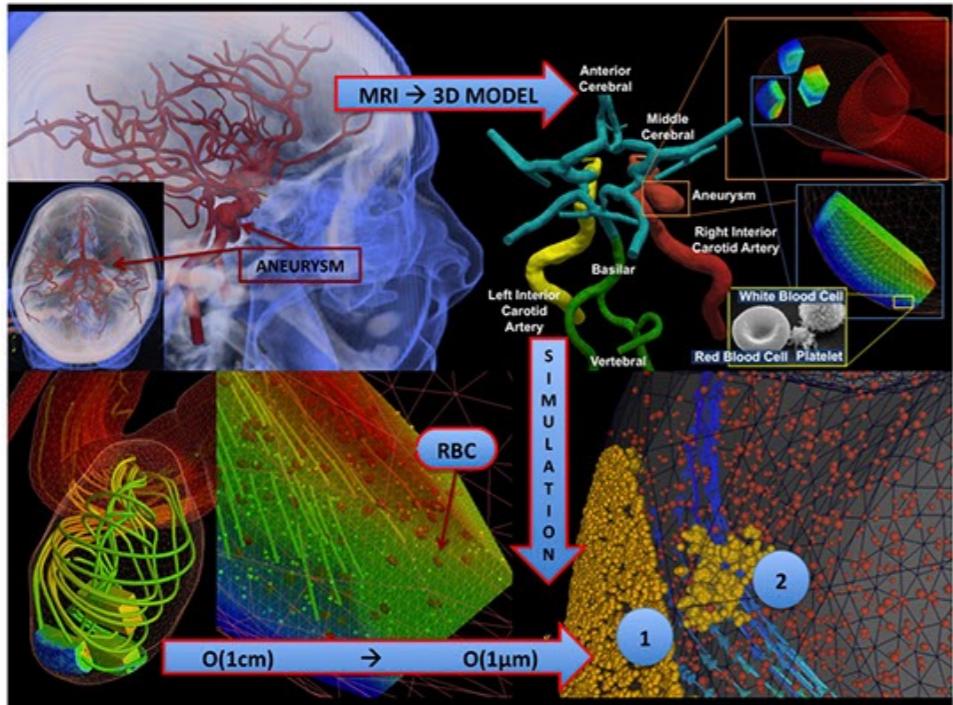
George E.P. Box



\*image credit: S. Sherwin (Imperial)

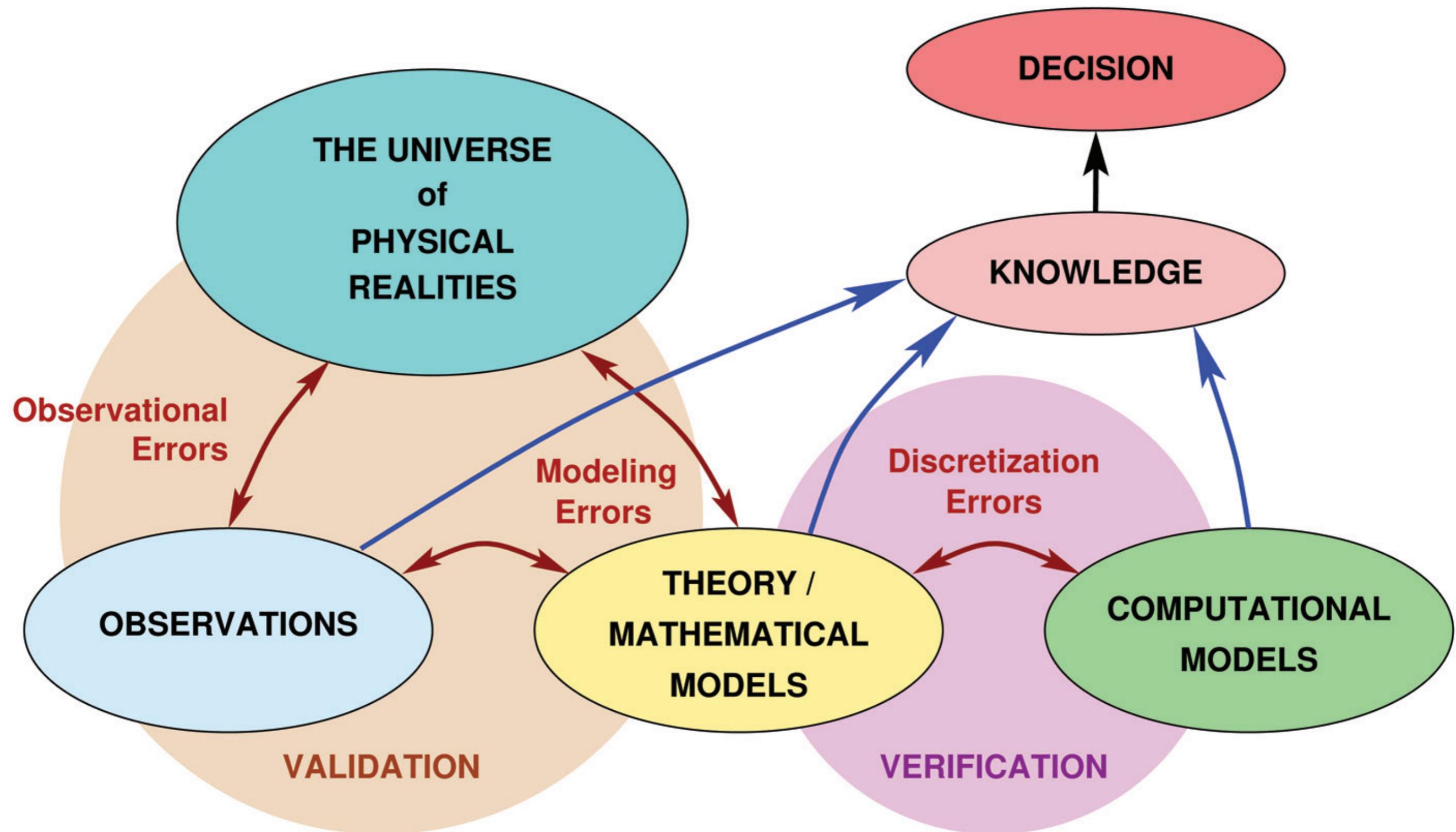
$$\frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} = -\frac{1}{\rho} \nabla p + \nu \nabla^2 \mathbf{u}$$
$$\nabla \cdot \mathbf{u} = 0$$

# Computational science & engineering



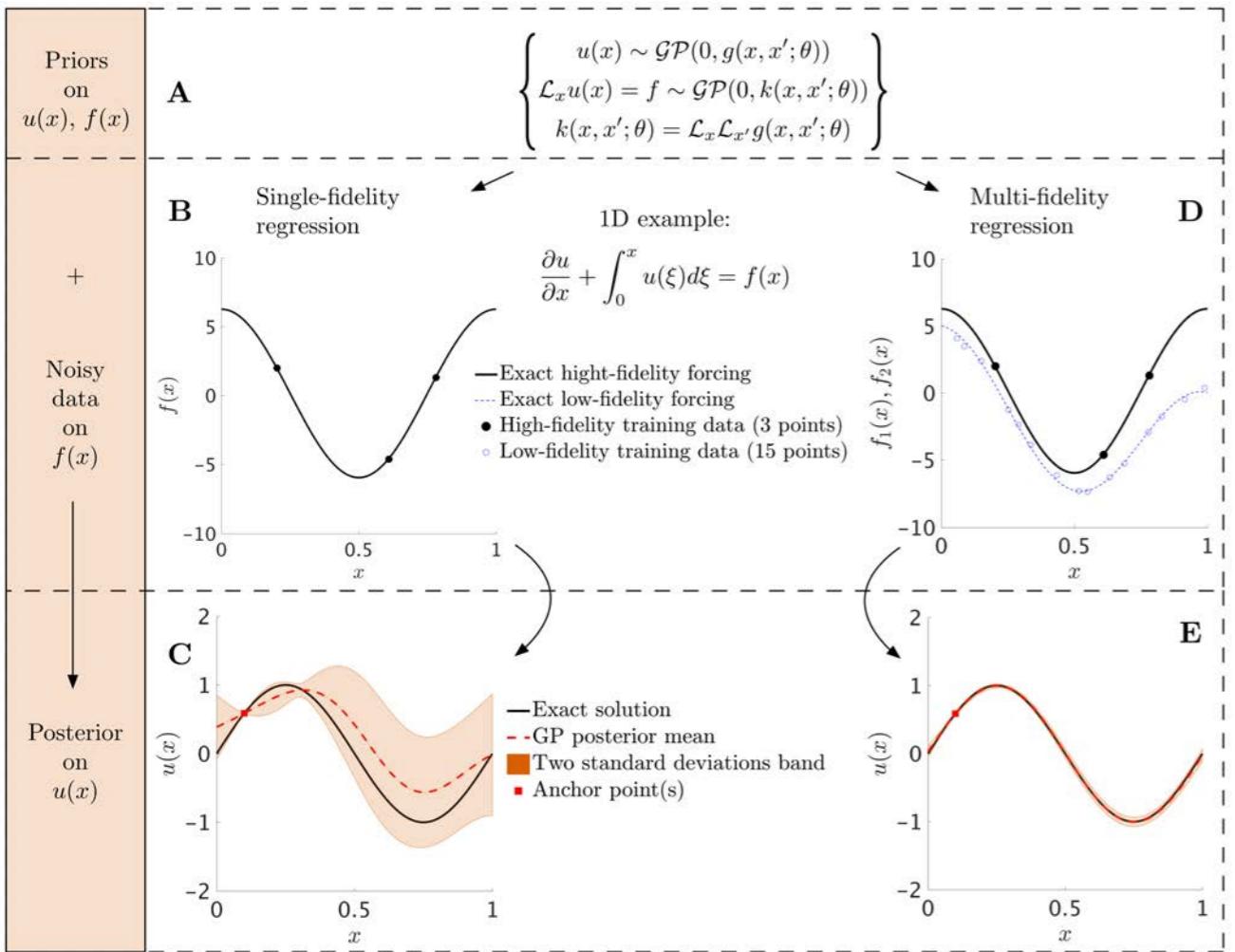
One exaflop is a thousand petaflops or a quintillion ( $10^{18}$ ) floating point operations per second!

# Predictive science



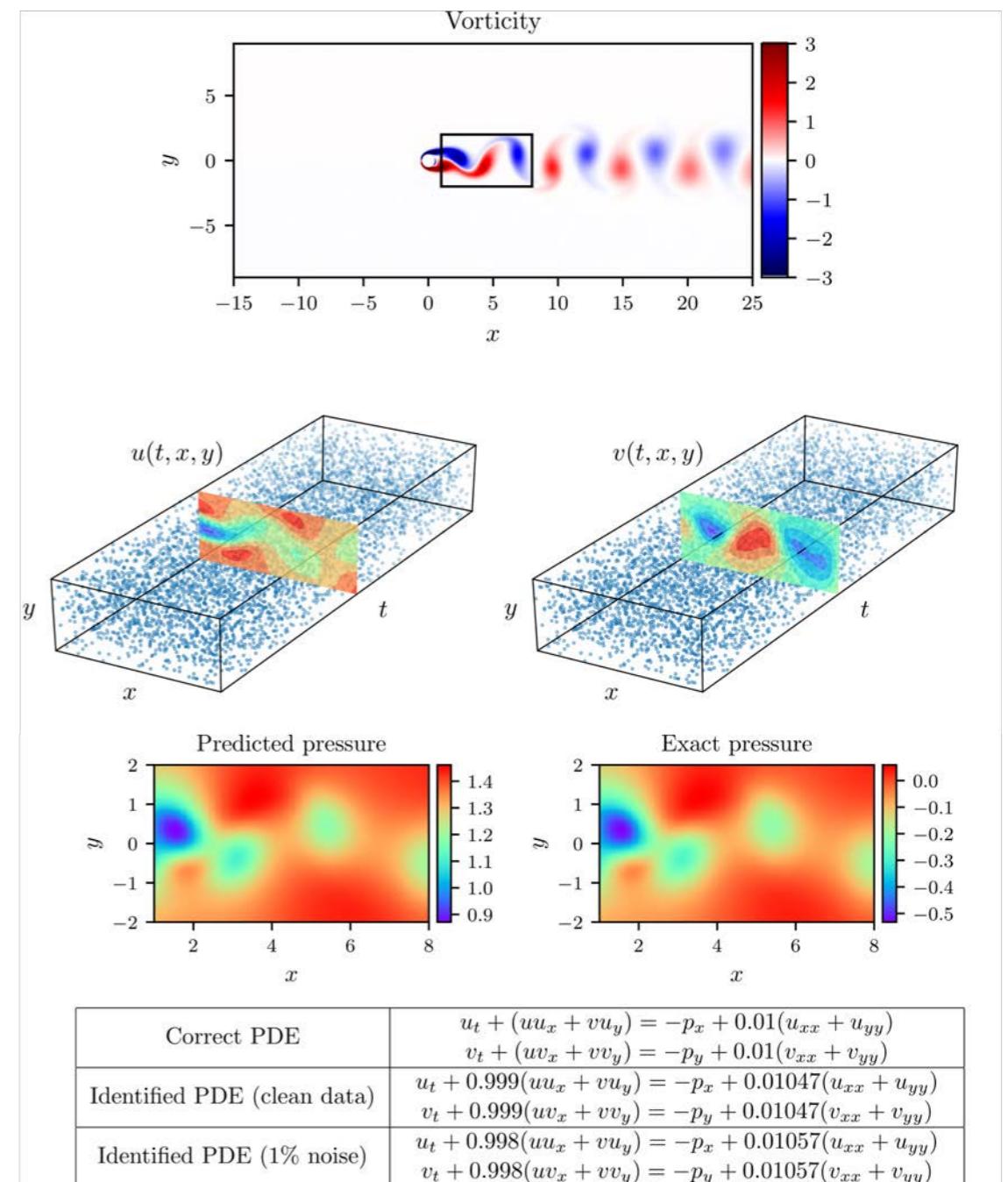
# CSE for ML

Leverage domain-specific prior information to construct more structured, interpretable, and accurate and data-efficient learning algorithms.



# ML for CSE

Leverage recent developments in machine learning to construct robust data-driven solvers, discover new models, and address longstanding challenges in scientific computing (e.g. high-dimensional problems, UQ, etc.)



$$f:\mathcal{X}\rightarrow\mathcal{Y}$$

# Supervised learning

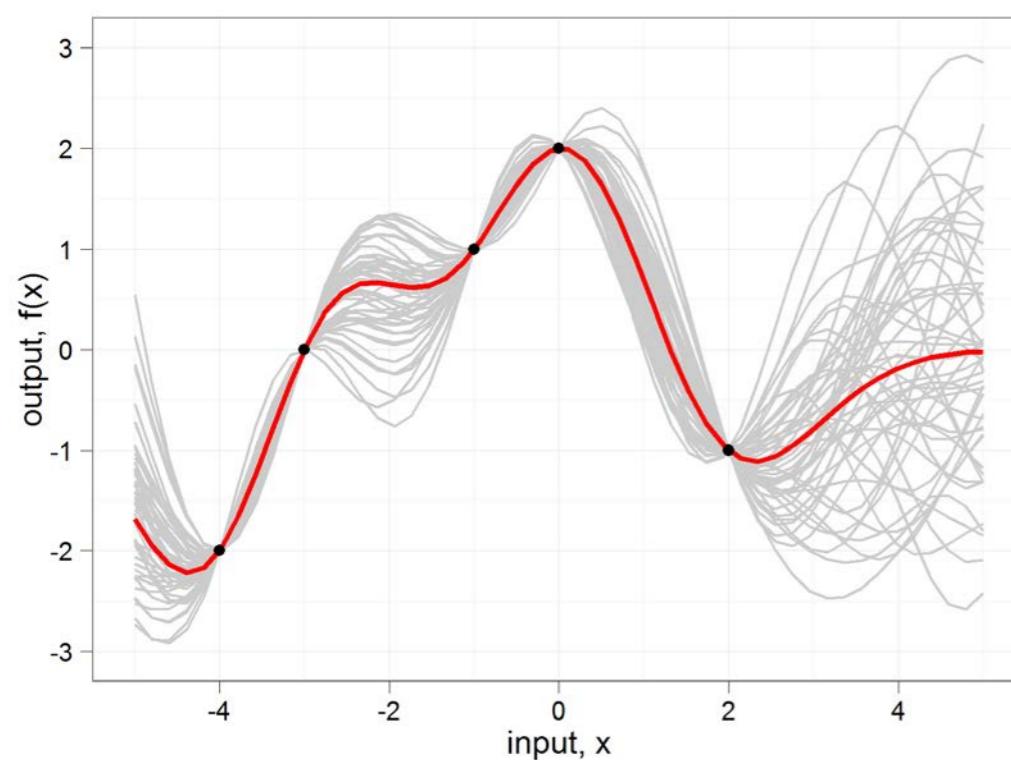
$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

$$\mathcal{D} = \{\mathbf{x}, \mathbf{y}\}, \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}$$

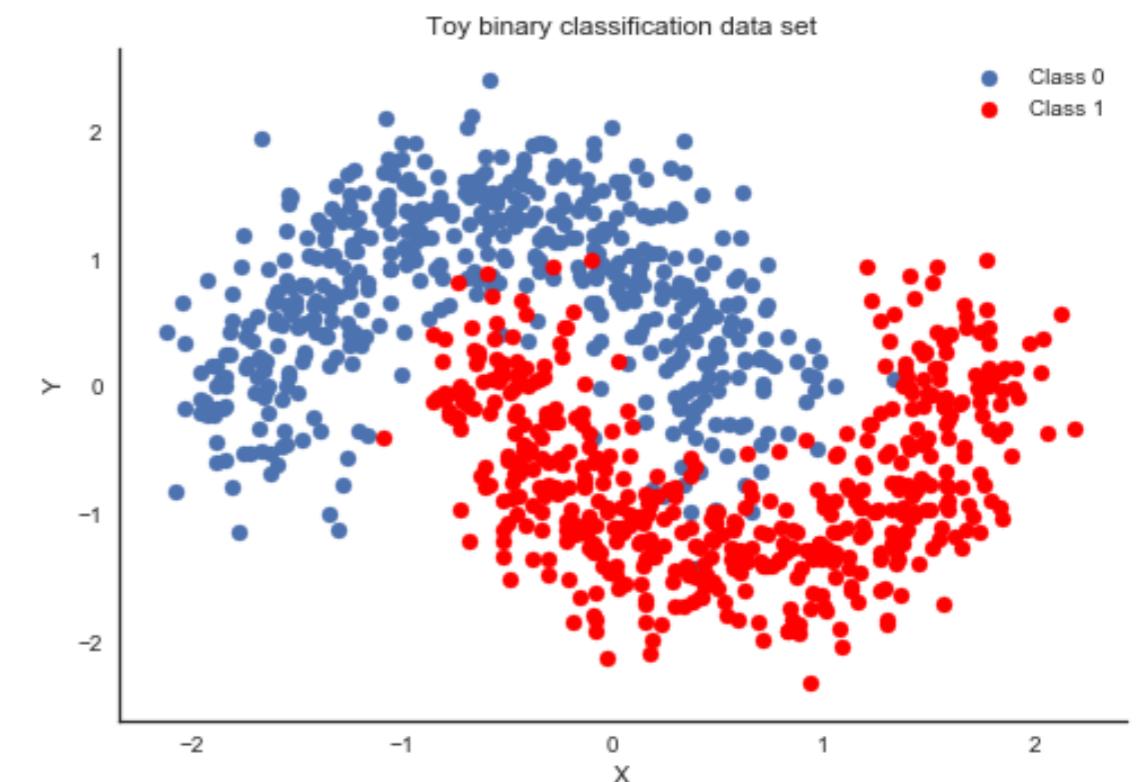
$$\mathbf{y} = f(\mathbf{x}) + \epsilon$$

$$p(f(\mathbf{x}^*) | \mathbf{x}^*, \mathcal{D})$$

Regression



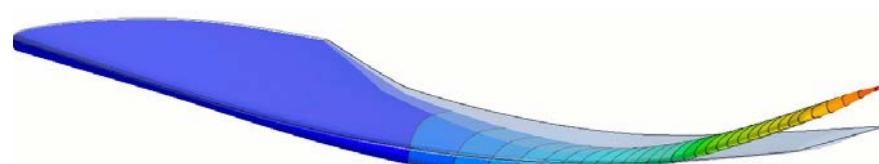
Classification



# Regression

## A motivating example:

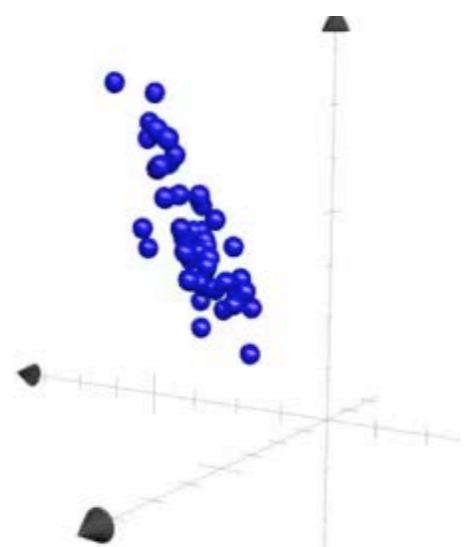
Estimation of wing deflection  
subject to fluid flow



$$\frac{\partial \sigma_{rr}}{\partial r} + \frac{1}{r} \frac{\partial \sigma_{r\theta}}{\partial \theta} + \frac{\partial \sigma_{rz}}{\partial z} + \frac{1}{r} (\sigma_{rr} - \sigma_{\theta\theta}) + F_r = \rho \frac{\partial^2 u_r}{\partial t^2}$$

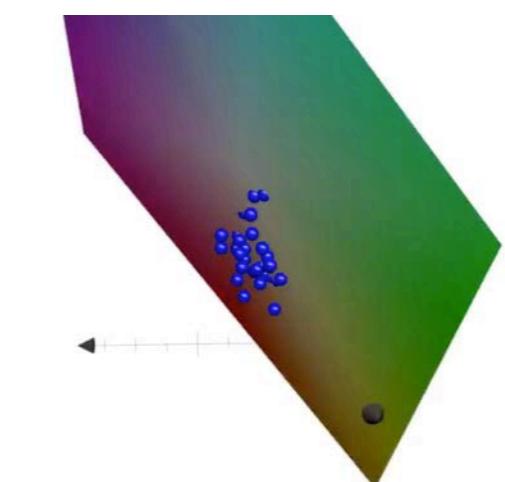
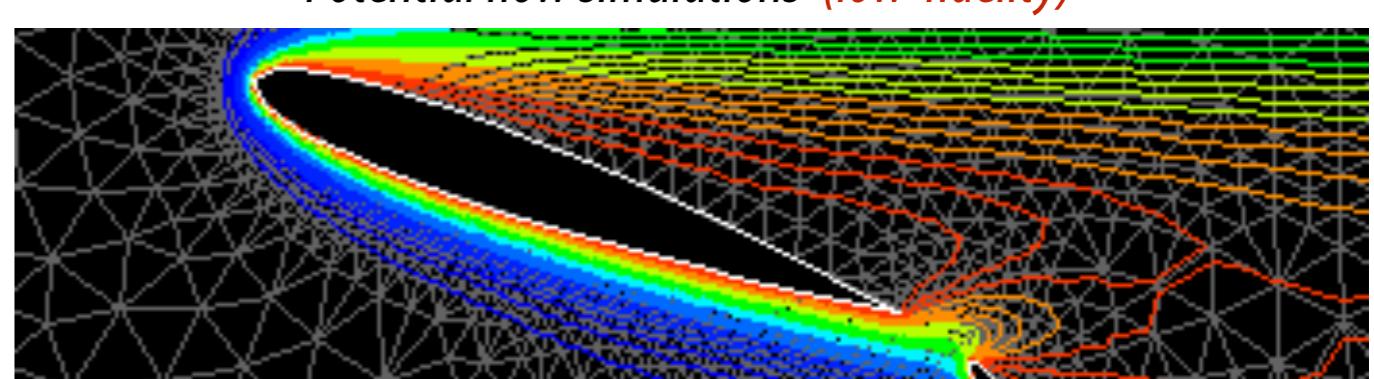
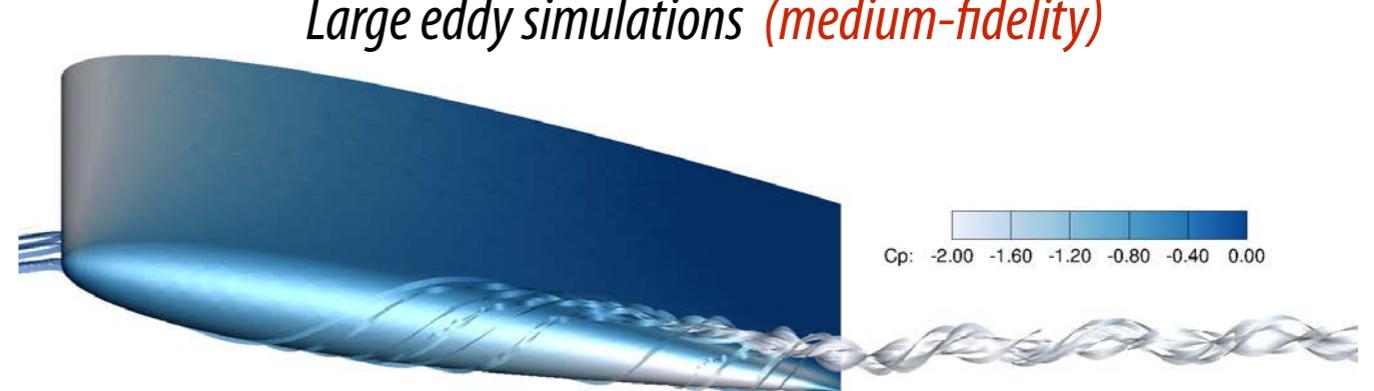
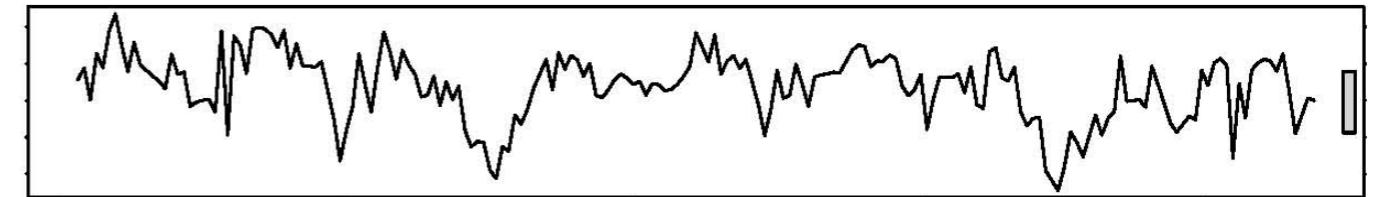
$$\frac{\partial \sigma_{r\theta}}{\partial r} + \frac{1}{r} \frac{\partial \sigma_{\theta\theta}}{\partial \theta} + \frac{\partial \sigma_{\theta z}}{\partial z} + \frac{2}{r} \sigma_{r\theta} + F_\theta = \rho \frac{\partial^2 u_\theta}{\partial t^2}$$

$$\frac{\partial \sigma_{rz}}{\partial r} + \frac{1}{r} \frac{\partial \sigma_{\theta z}}{\partial \theta} + \frac{\partial \sigma_{zz}}{\partial z} + \frac{1}{r} \sigma_{rz} + F_z = \rho \frac{\partial^2 u_z}{\partial t^2}$$



$$\mathcal{D} = \{x, y\}, \quad x \in \mathcal{X}, \quad y \in \mathcal{Y}$$

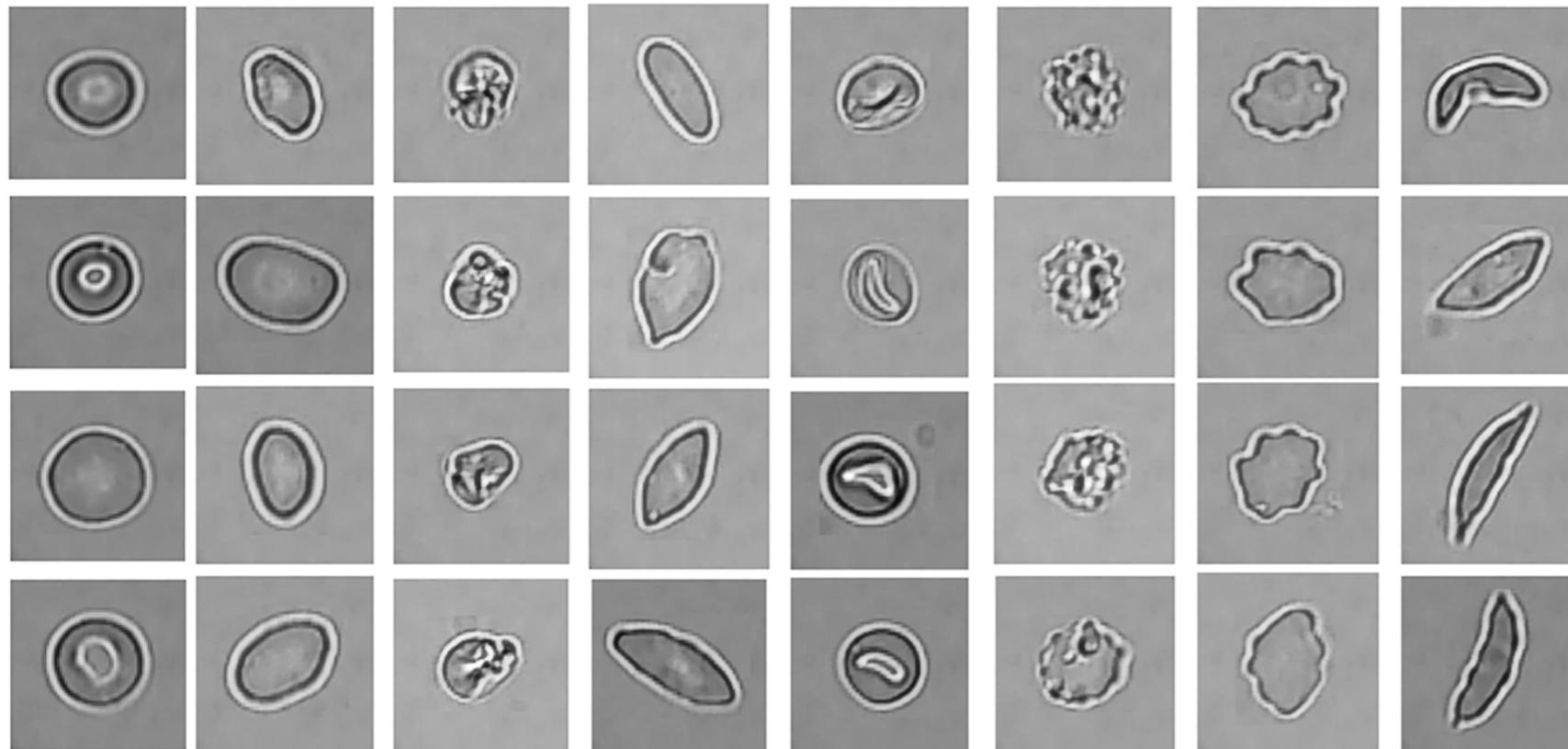
Multi-fidelity data



$$y = f(x) + \epsilon$$

# Classification

Discocytes   Oval   Reticulocytes   Elongated   Stomatocyte   Echinocytes   Granular   Sickle



$$\mathcal{D} = \{x, y\}, \quad x \in \mathcal{X}, \quad y \in \mathcal{Y}$$



$$y = f(x) + \epsilon$$

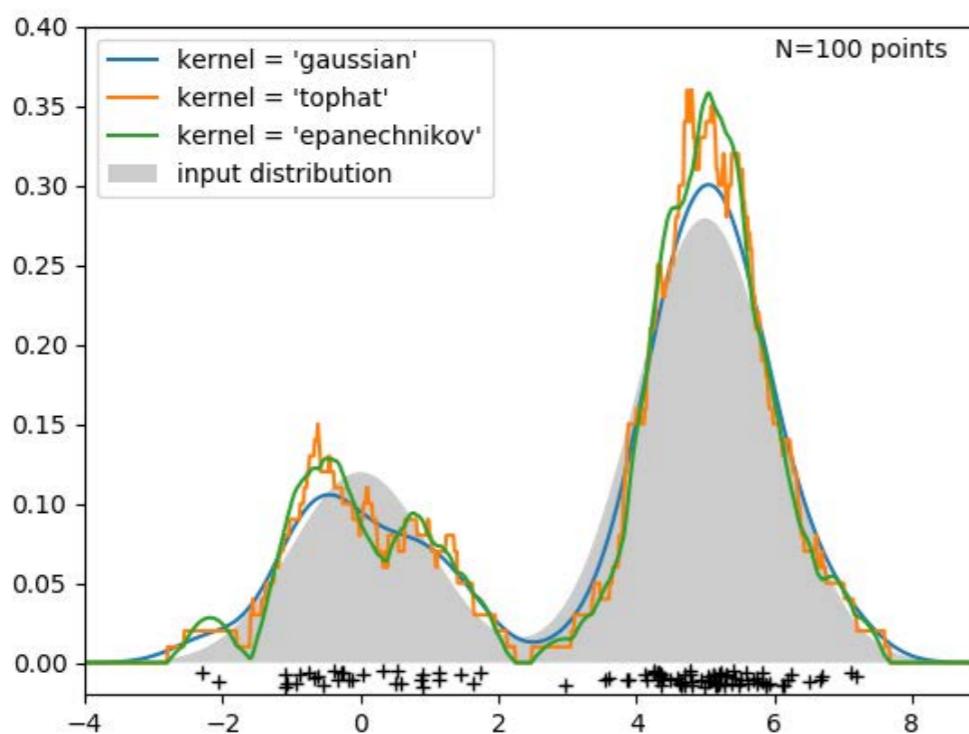
# Unsupervised learning

$$\{\mathbf{y}\}, \mathbf{y} \in \mathcal{Y}$$

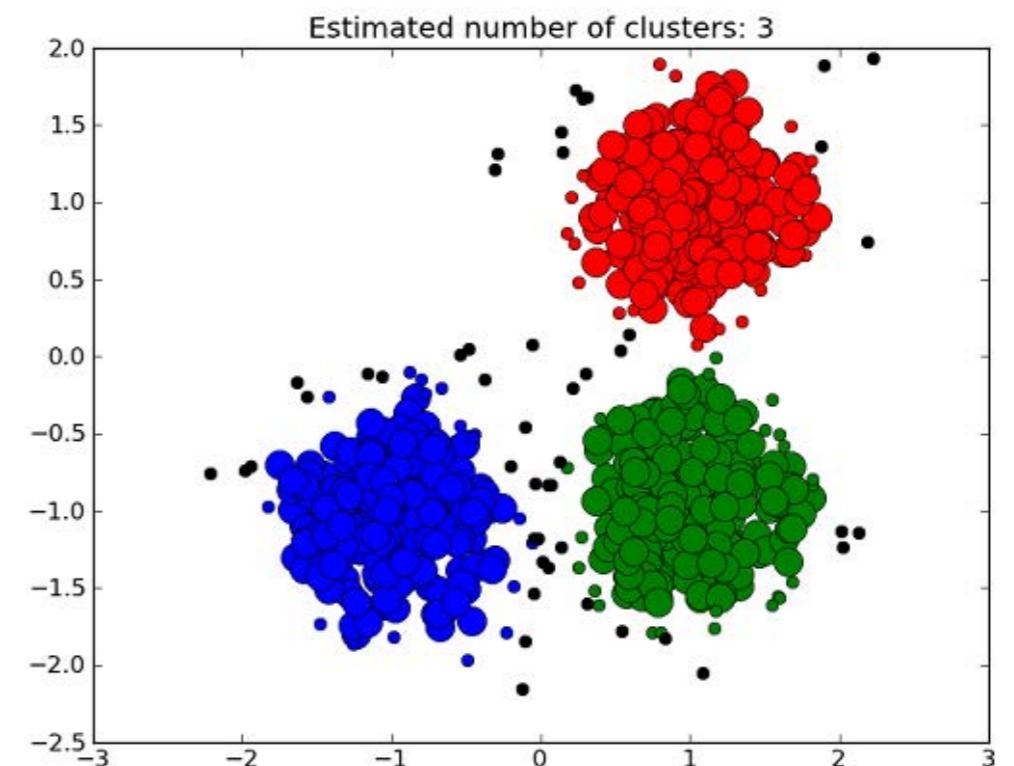
$$\mathbf{y} = f(\mathbf{z}) + \epsilon$$

$$p(\mathbf{y}), \mathbf{z} \sim p(\mathbf{z})$$

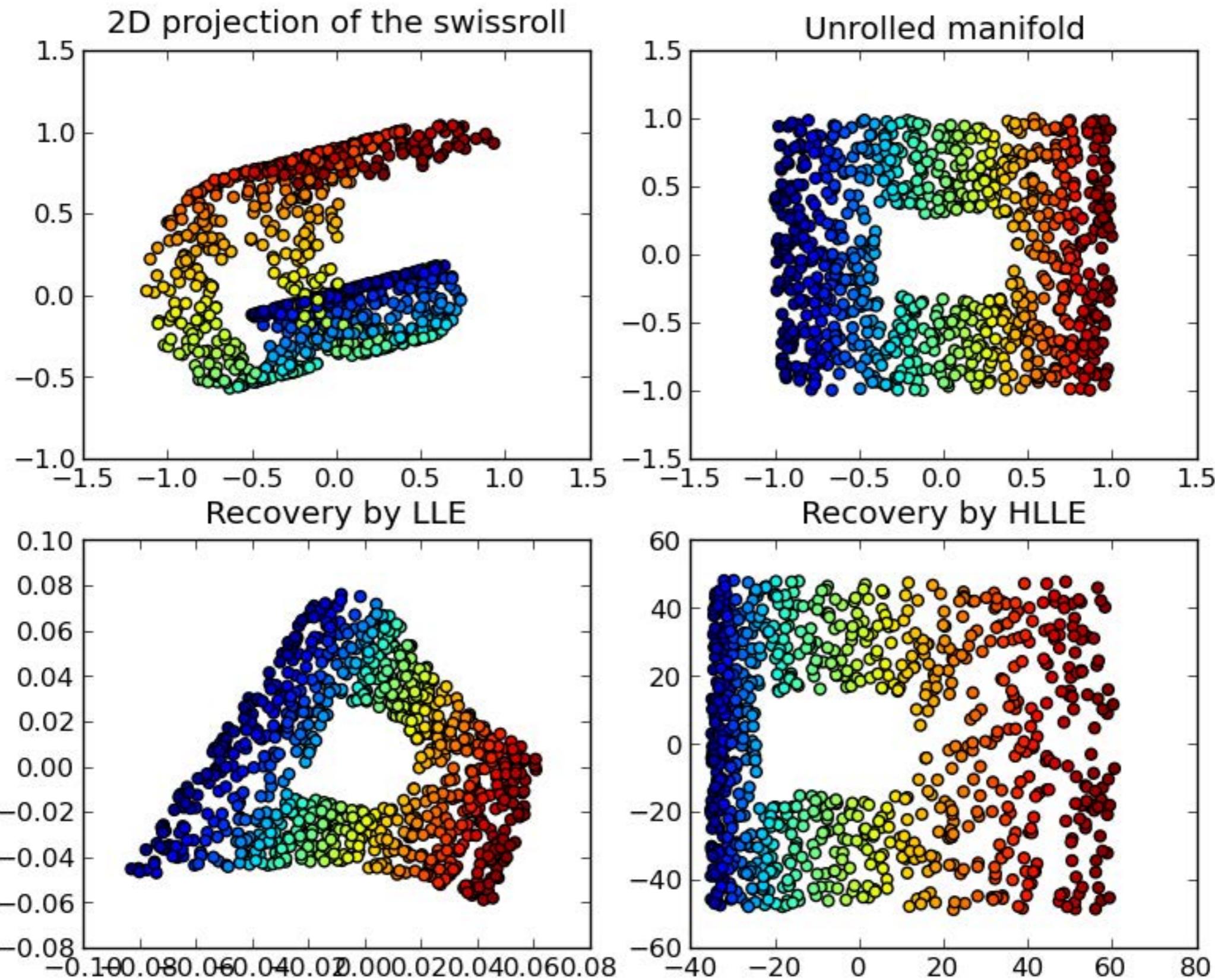
Density estimation



Clustering



# Dimensionality reduction



# Generative modeling



# Generative modeling

Gómez-Bombarelli, R., el. al. (2016).  
*Automatic chemical design using a data-driven continuous representation of molecules.* arXiv preprint arXiv:1610.02415.

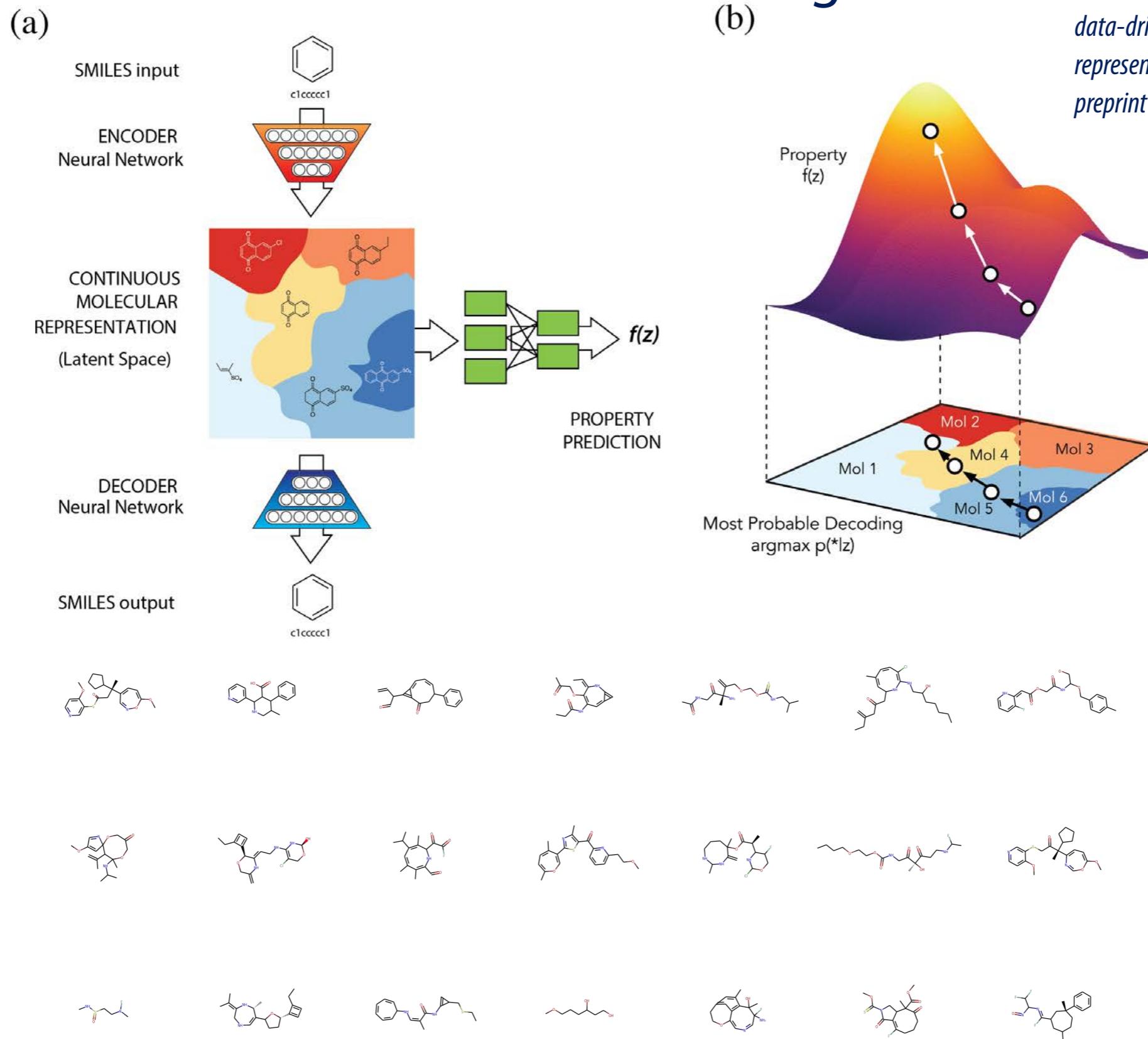


Figure 8: Molecules decoded from randomly-sampled points in the latent space of the ZINC VAE.

SPINGER BRIEFS IN STATISTICS

Jordi Vallverdú

# Bayesians Versus Frequentists A Philosophical Debate on Statistical Reasoning

Springer

# Uncertainty quantification



# Course goals

- Introduce a set of enabling machine learning tools for computational science and engineering applications.
- Utilize prior knowledge and modeling expertise to construct more structured and efficient predictive algorithms.
- Introduce a new type of scientific computing in which machine learning tools are used to develop new solvers, discover new models, perform sensitivity analysis/uncertainty quantification, guide the acquisition of new information, and solve design and optimization problems.

## Course logistics

- Duration: 14 Weeks
- Time: Tuesdays and Thursdays, 4:30 pm to 6:00 pm
- Dates: January 10 to April 25, 2018. No classes on Tuesday March 6, and Thursday March 8.
- Location: Towne Building, Room 309
- Office Hours: Thursdays 11:00 am to 1:00 pm at Office 359C (3401 Walnut Street).

For more details visit the course website:

<https://www.seas.upenn.edu/~enm540/>

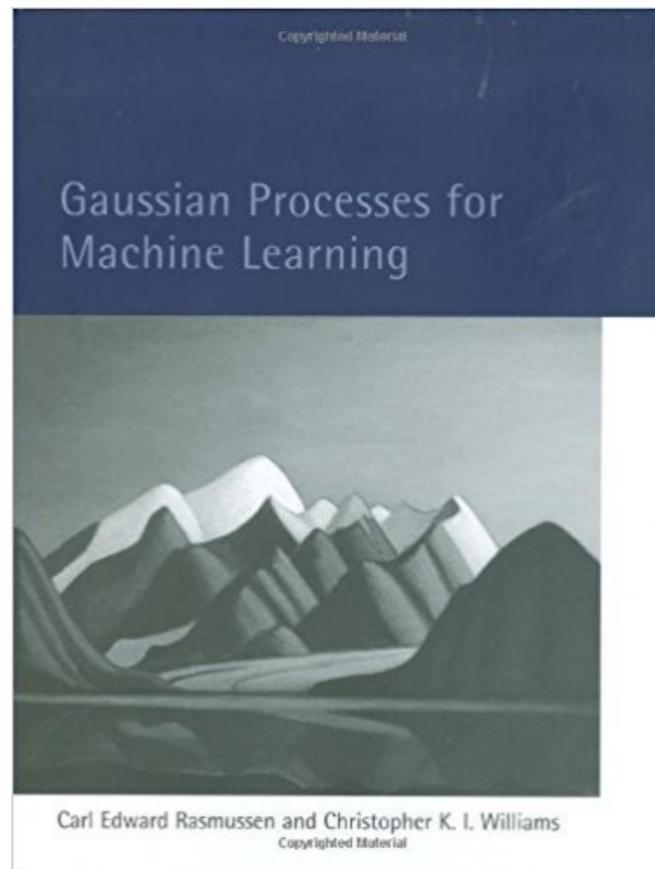
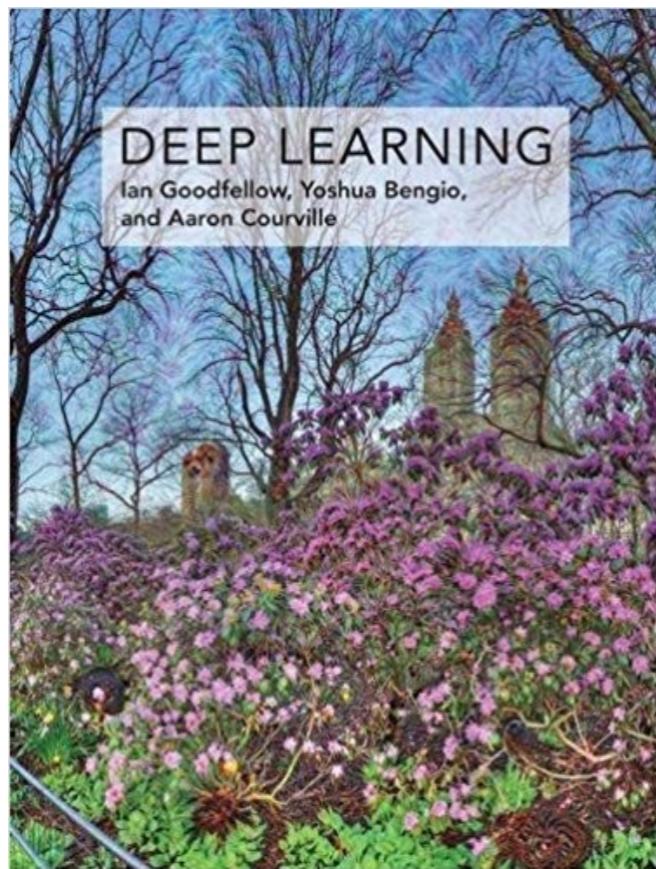
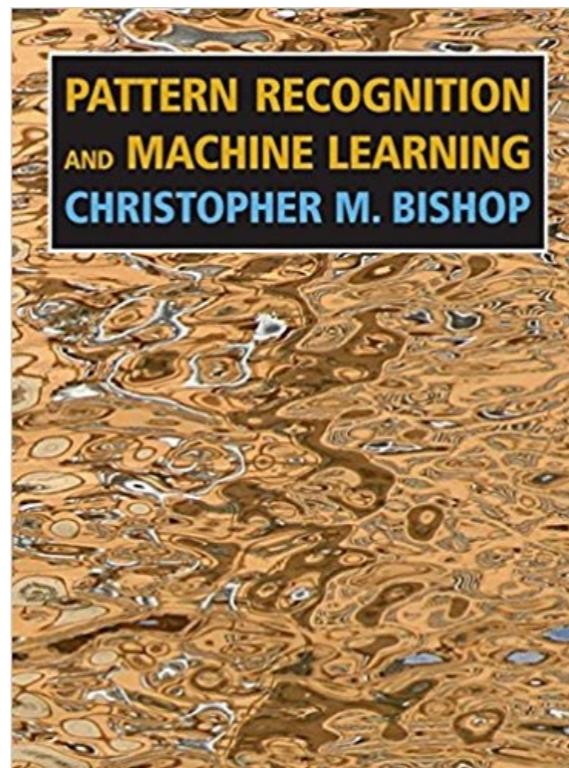
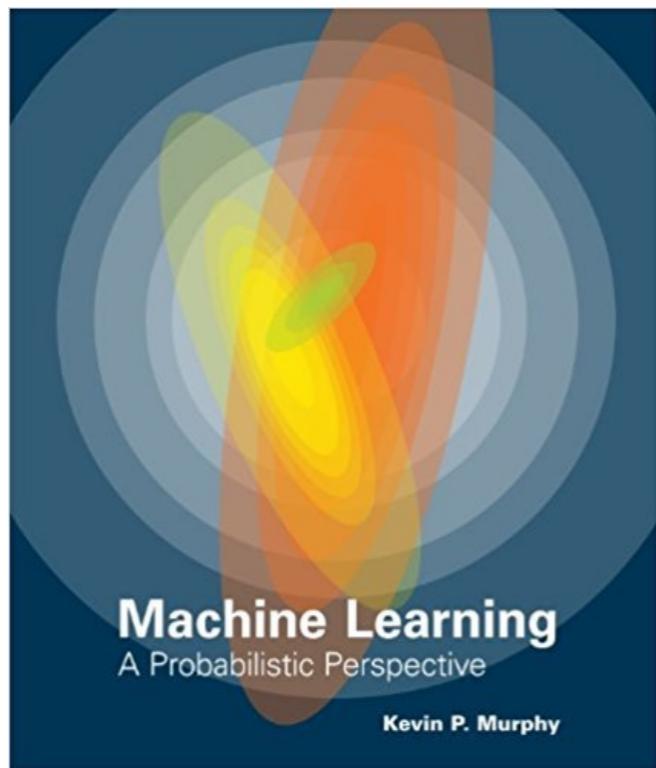
Slides and code can be found on Github:

<https://github.com/PredictiveIntelligenceLab/ENM-540>

# Course outline

- Week 1: Introduction to data-driven modeling, motivating examples and course outline.
- Week 2: Probability theory primer through the lens of Bayesian linear regression.
- Week 3: Optimization: gradients and Hessians, gradient descent, Newton's algorithm, stochastic gradient descent.
- Week 4: Deep neural networks: back-propagation, over-fitting and regularization, automatic differentiation.
- Week 5: Image classification with convolutional neural networks.
- Week 6: Recurrent neural networks and LSTMs.
- Week 7: Sampling and quantification of uncertainty.
- Week 8: Gaussian processes, multi-output Gaussian processes and multi-fidelity modeling.
- Week 9: Bayesian optimization and active learning.
- Week 10: Probabilistic scientific computing: Gaussian process regression meets differential equations.
- Week 11: Unsupervised learning: principal component analysis, Gaussian process latent variable models.
- Week 12: Variational auto-encoders and conditional variational auto-encoders.
- Week 13: Generative adversarial networks.
- Week 14: Class presentations of final papers.

# Textbooks



# Software

## Week 1 - Setting up your computing environment

Consult the Web for instructions and make sure you have the following software properly installed and working on your laptop:

- MATLAB: Free academic license and installation instructions available through [CETS](#).
- Python 2.7 or higher, set up for scientific computing. The simplest way to set this up is by installing the [Anaconda](#) distribution.
- [Tensorflow](#). If you have access to GPU hardware, make sure you install a version with GPU support.
- [PyTorch](#). If you have access to GPU hardware, make sure you install a version with GPU support.
- [Autograd](#). Efficiently computes derivatives of numpy code.
- [Jupyter notebook](#). You will need this in order to follow some of the in-class tutorials.
- [Git](#). You will need this in order to download and stay in sync with the latest code we will develop in class.

Make sure you can successfully run tutorials provided in the course Github page:

- [Linear regression in Numpy](#)
- [Linear regression in Numpy with Autograd](#)
- [Linear regression in PyTorch](#)
- [Linear regression in Tensorflow](#)
- [Interactive Jupyter notebook](#)