

CSCI 699 - ProbGen

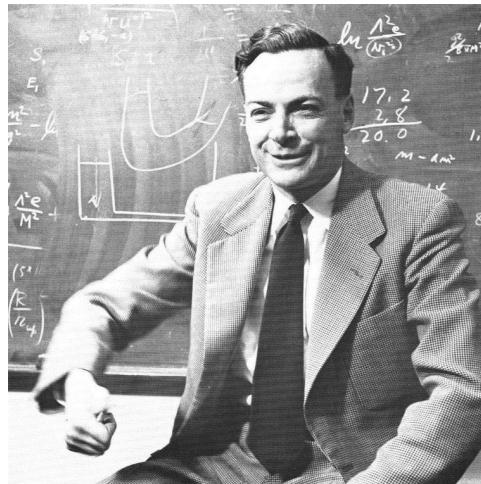
Probabilistic and Generative Models

Willie Neiswanger

Lecture 1 - Introduction

Generative Models

Richard Feynman

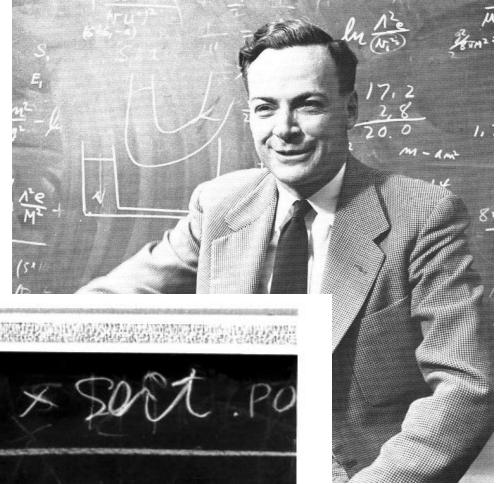
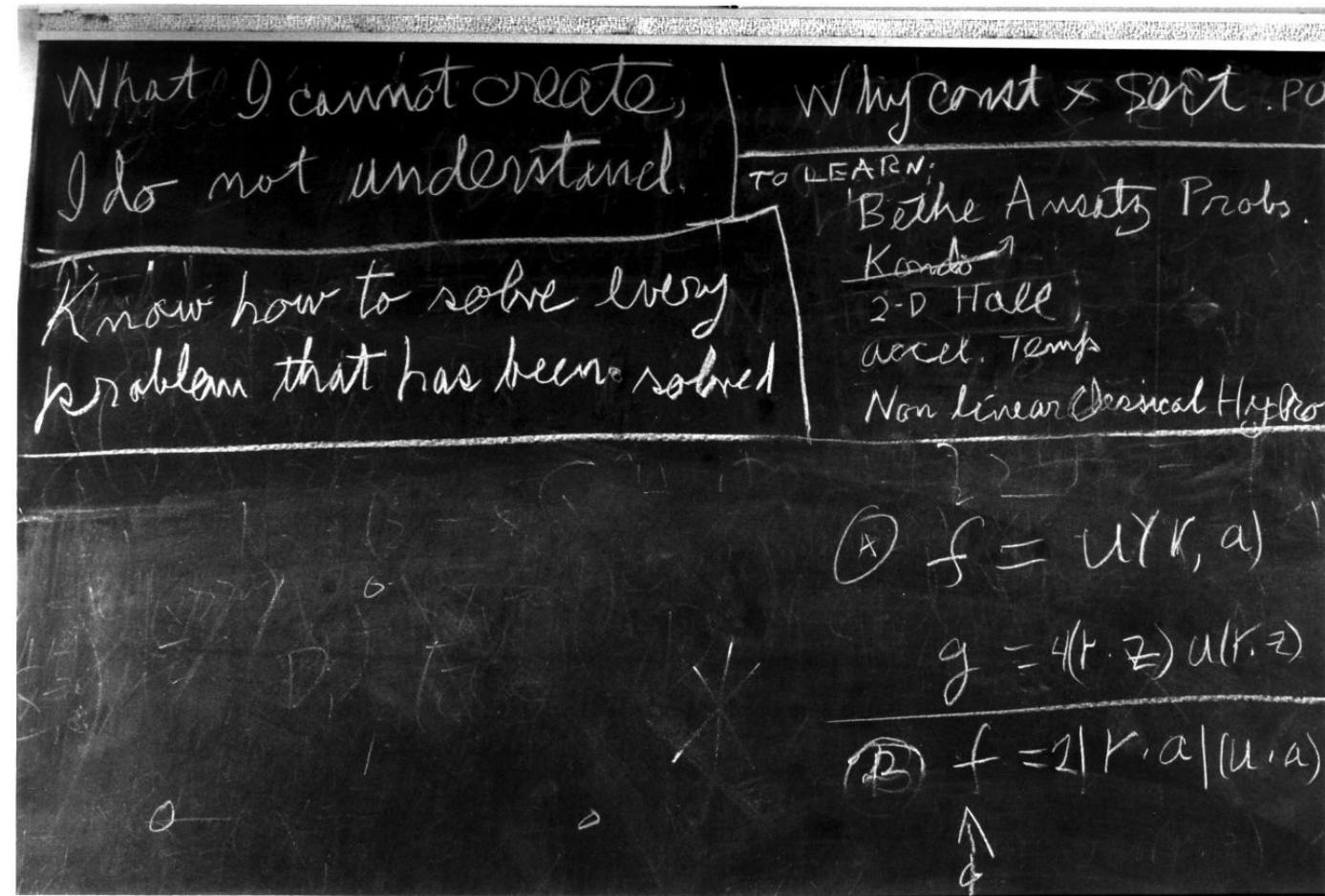


Source: Stefano Ermon, Deep Generative
Models Course

(famous physicist)

Generative Models

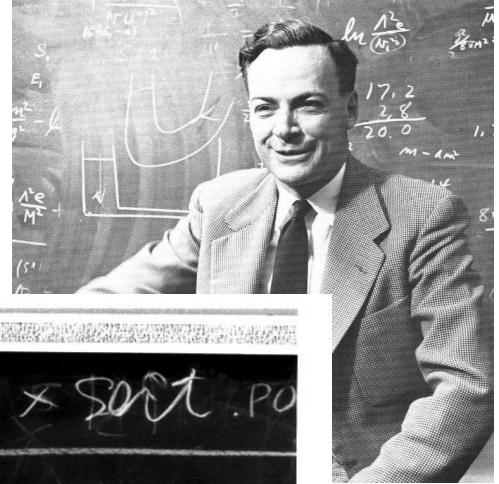
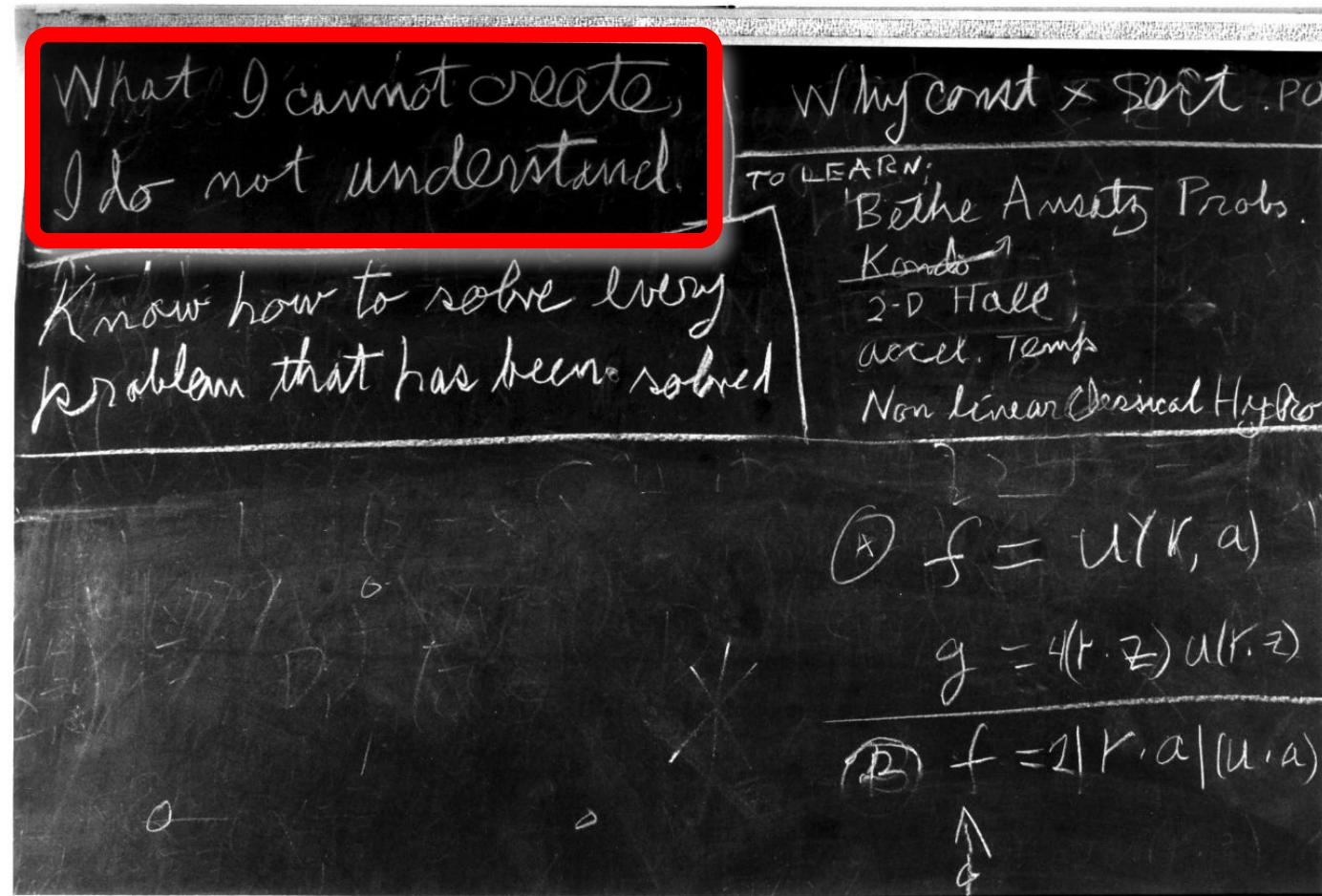
Richard Feynman's
blackboard at the time
of his death.



n, Deep
course

Generative Models

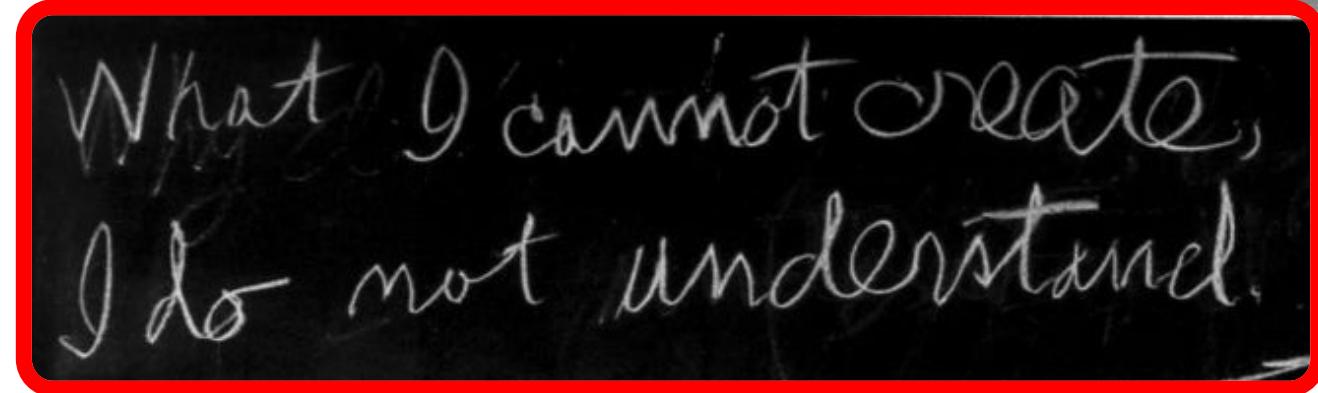
Richard Feynman's
blackboard at the time
of his death.



n, Deep
course

Generative Models

Richard Feynman's
blackboard at the time
of his death.



Stefano Ermon, Deep
Generative Models Course

Richard Feynman: “*What I cannot create, I do not understand*”

⇒ Meaning: to understand a thing, is to be able to create that thing.

Generative Models

And this is very much the mantra of **Generative Modeling in AI**:

⇒ Producing intelligence (*understanding*) via the ability to generate (*create*) things.

Generative Models

And this is very much the mantra of **Generative Modeling in AI**:

⇒ Producing intelligence (*understanding*) via the ability to generate (*create*) things.

Intuitively, producing intelligence via:

- Learning representations that provide the ability to generate realistic data from the world around us!
- Self-supervised learning.
- Next-token generation.
- Reconstruction of data.



Generative Models

Generative models are taking the world by storm!



ElevenLabs



Generative Models → Probabilistic Models

These models arise from a long history of probabilistic models in machine learning.

Journal of Machine Learning Research 1(2000) 1107-1115

Submitted 4/02; Published 2/03

A Neural Probabilistic Language Model

Yoshua Bengio

Réjean Ducharme

Pascal Vincent

Christian Jauvin

Département d'Informatique et de Recherche Opérationnelle,

Centre de Recherche Mathématiques

Université de Montréal, Montréal, Québec, Canada

BENGIO@CS.UMONTREAL.CA

DUCHARME@CS.UMONTREAL.CA

VINCENT@CS.UMONTREAL.CA

JAUVIN@CS.UMONTREAL.CA

Editor: Jim Kyndt, Thomas Hofmann, Tomasz Kipf and John Shawe-Taylor

Abstract

A goal of statistical language modeling is to learn the joint probability function of sequences of words in a language. This is a immensely difficult because of the curse of dimensionality: a word sequence to which the model will be tested is likely to be different from all the word sequences seen during training. We propose a solution to this problem by learning a distributed representation for words which allows each training example to contribute to the joint probability of every word sequence containing that word. The proposed approach is called a neural language model. The model learns simultaneously (1) a distributed representation for each word along with (2) a joint probability function for sequences of words based on the learned representations. Generalization is obtained because a sequence of words that has never been seen yet gets high probability. The proposed approach is able to learn large amounts of words forming an already known sentence. Training such large models (with millions of parameters) without a significant time is itself a significant challenge. We report on experiments using neural networks with a distributed representation of words. The proposed approach is able to learn words significantly improves on state-of-the-art n-gram models, and the proposed approach allows to take advantage of parallel processing.

Keywords: Statistical language modeling, artificial neural networks, distributed representation, curse of dimensionality

1. Introduction

A fundamental problem that makes language modeling and other learning problems difficult is the *curse of dimensionality*. It is particularly obvious in the case when one wants to model the joint distribution between many discrete random variables (such as words in a sentence, or discrete attributes in a database). For example, if we consider a vocabulary of size $V = 10^4$, there are potentially 10^{40} conceivable words in a natural language with a vocabulary V of size 10^4 , whereas there are potentially $10^{1000000} \approx 10^{300}$ parameter. When modeling continuous variables, we obtain problems that are similarly difficult to solve. For example, if we consider a Gaussian mixture model (or a Gaussian mixture function) the function to be learned can be expected to have some local smoothness properties. For discrete species, the generalization structure is not as obvious: any change of these few discrete variables may have a drastic impact on the value of the function to be estimated.

©2000 Yoshua Bengio, Réjean Ducharme, Pascal Vincent, Christian Jauvin.

Autoregressive LLMs, 2003

Journal of Machine Learning Research 6 (2005) 695-709
Submitted 11/04; Revised 3/05; Published 4/05

Estimation of Non-Normalized Statistical Models by Score Matching

Aapo Hyvärinen
Baltic Institute for Information Technology (BRU)
Department of Computer Science
FIN-00014 University of Helsinki, Finland

AAPO.HYVARNEN@HUT.HELPP.FI

Editor: Peter Dayan

Abstract

One often wants to estimate statistical models where the probability density function is known only up to a multiplicative normalization constant. Typically, one then has to resort to Markov Chain Monte Carlo methods, or approximations of the normalization constant. Here we show that such problems can be solved by estimating the so-called normalized distance between the gradient of the log-density given by the model and the gradient of the log-density given by the true underlying distribution. This normalized distance of the log-density function, in principle, a very difficult non-convex problem, we prove a surprising result: that it can be approximated by a sum of squares of derivatives of the log-density function. This provides a simple formula for this objective function. The density function of the observed data is estimated by a function which simplifies to a Gaussian if the data is a sum of some derivatives of the log-density given by the model. The validity of the method is demonstrated on several examples, including a non-Gaussian mixture, a nonparametric model, and by estimating an overcomplete filter set for natural image data.

Keywords: statistical estimation, non-normalized densities, pseudo-likelihood, Markov chain Monte Carlo, contrastive divergence

1. Introduction

In many cases, probabilistic models in machine learning, statistics, or signal processing are given only up to a factor of proportionality. This is because the normalization constant is an unknown normalization constant whose computation is too difficult for practical purposes.

Assume we observe a random vector $x \in \mathbb{R}^n$ which has a probability density function (pdf) denoted by $p(x|\theta)$, where θ is a parameter vector. We want to estimate the parameter θ from x , i.e. we want to approximate $p(x|\theta)$ by $\hat{p}(x|\hat{\theta})$ for the estimated parameter $\hat{\theta}$. (We shall here consider the case $n > 1$.)

The problem we consider here is why we only are able to compute the pdf given by the model up to a multiplicative constant $Z(\theta)$:

$$p(x|\theta) = \frac{1}{Z(\theta)} p_\theta(x|\theta)$$

That is, we do know the functional form of p as an analytical expression (or any form that can be easily computed), but we do not know how to easily compute Z which is given by

Score Matching, 2005

Auto-Encoding Variational Bayes

<p>Diederik P. Kingma Machine Learning Group University of Amsterdam dkpkingma@gmail.com</p>	<p>Max Welling Machine Learning Group University of Amsterdam wellling.mw@gmail.com</p>
--	---

Abstract

How can we perform efficient inference and learning in directed probabilistic models, in the presence of continuous latent variables with intricate posterior distributions, and large datasets? We introduce a stochastic variational inference and message passing framework that, under mild conditions, provides a lower bound on the marginal likelihood that is both a tight approximation to the true posterior, and even works in the intractable case. Our contributions are two-fold. First, we propose that a representation of the posterior distribution as a lower bound on the log-marginal likelihood can be straightforwardly optimized using standard stochastic gradient methods. Second, we show that for i.i.d. datasets with continuous latent variables, the variational lower bound can be computed very efficiently by fitting an approximate inference model (also called a recognition model) to the data. This allows us to obtain a lower bound estimator.

Theoretical advantages are reflected in experimental results.

1 Introduction

How can we perform efficient approximate inference and learning with directed probabilistic models whose continuous latent variables either have intractable posterior distributions? The variational Bayesian approach [1] provides a general solution to this problem by approximating the posterior. Unfortunately, the common mean-field approach requires analytical solutions of expectation-maximization (EM) steps, which are slow and often do not converge. We propose a new representation of the variational lower bound yields a simple differentiable unbiased estimator of the log-marginal likelihood. This unbiased estimator can be used to perform exact and efficient approximate posterior inference in almost any model with continuous latent variables, and straightforwardly optimise using standard stochastic gradient descent techniques.

In this paper, we propose a new variational lower bound estimator and a new algorithm, auto-Encoding VI (AEVI), the algorithm in the AVEI algorithm we make inference and learning especially efficient. The AVEI algorithm is based on a variational lower bound estimator that is able to perform very efficient approximate posterior inference using simple ancestral sampling, which in turn allows us to efficiently learn approximate posterior inference models. The proposed approximate posterior inference (such as MCMC) per datapoint. The learned approximate posterior inference model can also be used for a host of tasks such as recognition, denoising, representation and visualisation purposes. When a neural network is used for the recognition model, we arrive at the variational autoencoder.

2 Method

The strategy in this section is to seek to derive a lower bound estimator (a stochastic objective function) for a variety of directed graphical models with continuous latent variables. We will restrict ourselves here to the common case where we have an i.i.d. dataset with latent variables per datapoint, and hence the log-marginal likelihood is the sum of the log-marginal likelihoods over all datapoints, one on the (global) parameters, and variational inference on the latent variables. It is, for example,

Variational Autoencoder, 2013

Deep Unsupervised Learning using Non equilibrium Thermodynamics

Jascha Sohl-Dickstein
Stanford University
Eric A. Wang
University of California, Berkeley
Nirav Makharevichan
Stanford University
Surya Ganguli
Stanford University

JASCHA@STANFORD.EDU
EAWIESS@BERKELEY.EDU
NIRUM@STANFORD.EDU
SGANGULI@STANFORD.EDU

Abstract
A causal graphical model learning framework modeling complex data sets using highly flexible families of probability distributions in which learning is based on a thermodynamic principle that are still analytically or computationally tractable. This framework is able to learn models that easily achieve both flexibility and tractability.
The framework is inspired by non-equilibrium statistical physics, it is able to learn both the complex density structure in a data distribution through an iterative process of inference and learning, and learn a reverse diffusion process that restores more tractable probability distributions from the intractable generative model of the data. This approach allows us to rapidly learn deep generative models with thousands of layers or time steps, as well as learn a generative model for the joint probabilities under the learned model. We also introduce a release an open source reference implementation of the algorithm.

1. Introduction
Inherently probabilistic models suffer from a race between learning efficiency, tractability and flexibility. Models that are tractable can be analytically evaluated and easily fit to data (e.g., a Gaussian or Laplace). However, *Proceedings of the 2015 Conference on Machine Learning (ICML 2015), 2015. JMLR: W&G volume 31. Copyright 2015 by the author(s).*

these models are typically descriptive in rich datasets. On the other hand, models that are flexible do not fit in structure in arbitrary data. For example, we can learn a model that is a Gaussian mixture of \sqrt{Z} (yielding the flexible distribution $p(x) \propto \frac{1}{\sqrt{Z}}$), where Z is a normalization constant. However, computing this model's joint probability distribution requires fitting training, or drawing samples from such flexible model types is typically intractable.

A variety of analytical approximations exist which are tractable, but do not remove, but rather for instance mean-field theory (Hinton et al., 2006; Hinton & Zemel, 1989); variational Bayes (Jordan et al., 1999); approximate divergence (Welling & Hanson, 2002; Hinton, 2002); minimum KL construction (Lin, 2011), principal point scoring (Gillis, 2011); variational message passing (Wainwright & Jordan, 2009); variational expectation propagation (Häggström, 2003); parallelized belief (Beal, 1997), loopy belief propagation (Koller & Friedman, 2009), and many, many more. Non-parametric methods (Gershman & Blei, 2012) can also be very effective.

2.1. Diffusion probabilistic models

We present a novel way to learn probabilistic models that

1. extreme flexibility in model structure,
2. exact sampling.

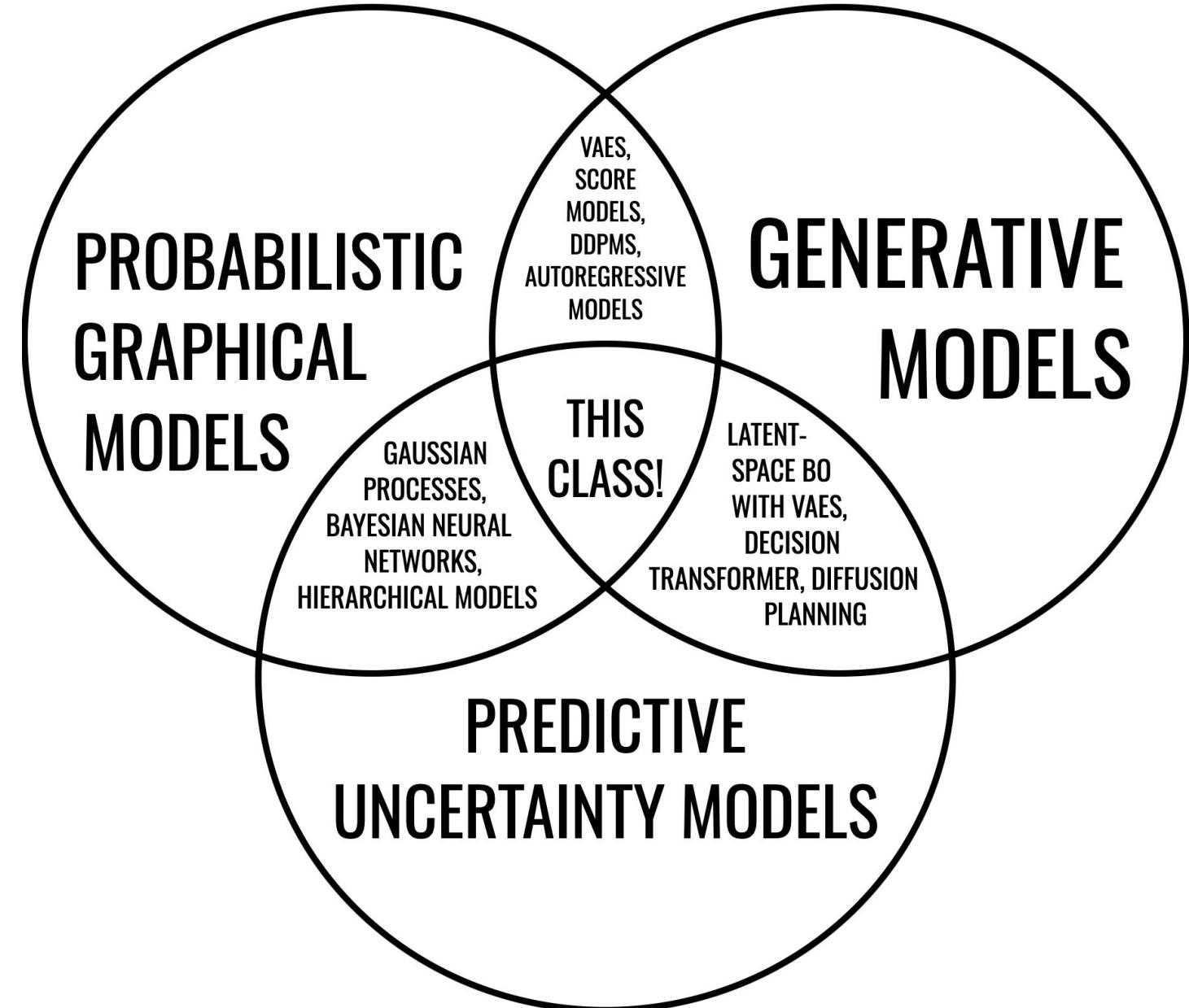
Non-parametric methods can be seen as maintaining several probability distributions over the same variable. The non-parametric Gaussian mixture model will represent a small number of components, while the non-parametric infinite data as a mixture of an infinite number of Gaussians.

Diffusion Models, 2015

Normalizing Flows, 2015

This Class

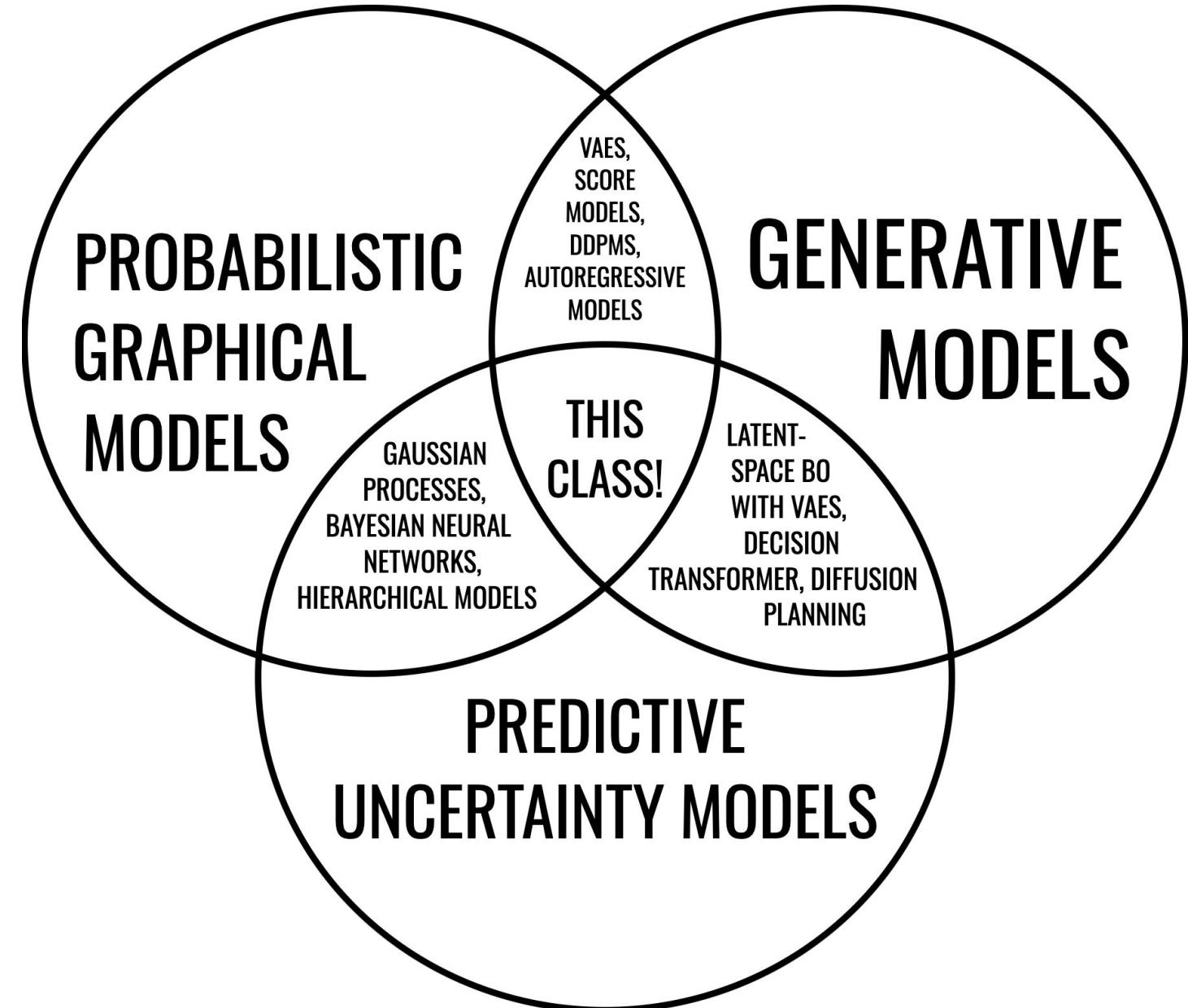
This course focuses on probabilistic models and their central role within modern machine learning and generative modeling.



This Class

This course focuses on probabilistic models and their central role within modern machine learning and generative modeling.

This course is inspired by some previous classes...



This Class – Inspiration

Probabilistic Graphical Models (PGM) class → e.g., at CMU & Stanford

10-708

[Home](#)
[Lectures](#)
[Homework](#)
[Project](#)
[Logistics](#)
[Previous](#)

10-708 – Probabilistic Graphical Models
2020 Spring

Many of the problems in artificial intelligence, statistics, computer systems, computer vision, natural language processing, and computational biology, among many other fields, can be viewed as the search for a coherent global conclusion from local information. The probabilistic graphical models framework provides an unified view for this wide range of problems, enables efficient inference, decision-making and learning in problems with a very large number of attributes and huge datasets. This graduate-level course will provide you with a strong foundation for both applying graphical models to complex problems and for addressing core research topics in graphical models.

- Instructor: [Eric P. Xing](#) (epxing@cs)
- Time: MW 12:00-1:20pm
- Location: Wean 7500
- Office Hours: Mon 1:30-2:30pm GHC 8101
- Piazza: <https://www.piazza.com/cmu/spring2020/10708>
- Gradescope: <https://www.gradescope.com/courses/80181>
- TAs (email, office hours):
 - Xun Zheng (xzheng1@andrew, Fri 4-5pm GHC 8013)
 - Ben Lengerich (blengeri@andrew, Thu 10-11am GHC 9005)
 - Haohan Wang (haohanw@andrew, Fri 5-6pm, GHC 5507)
 - Yiwen Yuan (yiweny@andrew, Tue 1:50-2:50pm, outside GHC 8011)
 - Xiang Si (xsi@andrew, Wed 2-3pm, GHC Citadel Commons)
 - Junxian He (junxian1@andrew, Mon 4-5pm GHC 6603)

Page generated 2020-04-29 16:58:34 EDT, by [jemdoc](#).

CS 228 - Probabilistic Graphical Models
Winter 2023-24

[Ed](#) [Calendar](#) [Course Notes](#)

[Logistics](#) | [Course Info](#) | [Syllabus](#) | [Other Resources](#)

Logistics

- Lectures: Tue, Thu, 10:30am-11:50am, Gates 81
- Office Hours and Sections: [Google Calendar](#)

For SCPD students, please email scpdsupport@stanford.edu or call 650-741-1542.

Instructor



Stefano Ermon
[ermon \[at\] cs.stanford.edu](mailto:ermon@cs.stanford.edu)

[Website](#)

Course Assistants



Charlie Marx [cmarx \[at\] stanford.edu](mailto:cmarx@stanford.edu)
Sofian Zalouk [szalouk \[at\] stanford.edu](mailto:szalouk@stanford.edu)
Garrett Thomas [gthomas \[at\] stanford.edu](mailto:gthomas@stanford.edu)
Chaitanya Patel [chpatel \[at\] stanford.edu](mailto:chpatel@stanford.edu)
Devansh Sharma [devansh \[at\] stanford.edu](mailto:devansh@stanford.edu)

Course Information

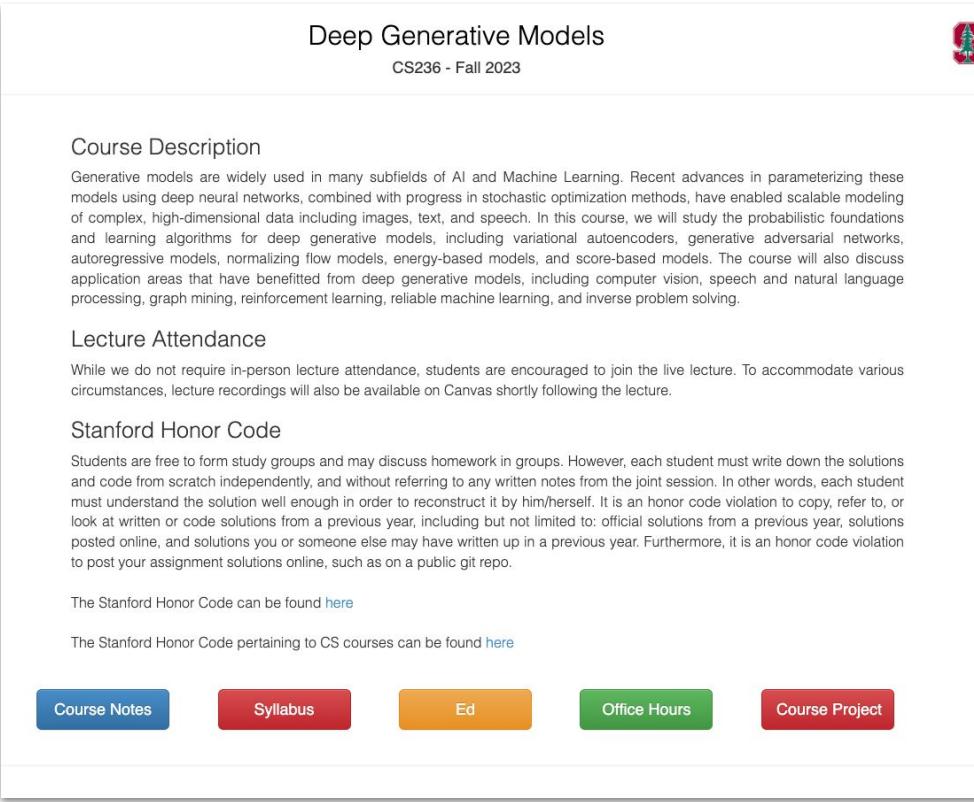
Course Description

Probabilistic graphical models are a powerful framework for representing complex domains using probability distributions, with numerous applications in machine learning, computer vision, natural language processing and computational biology. Graphical models bring together graph theory and probability theory, and provide a flexible framework for modeling large collections of random variables with complex interactions. This course will provide a comprehensive survey of the topic, introducing the key formalisms and main techniques used to construct them, make predictions, and support decision-making under uncertainty.

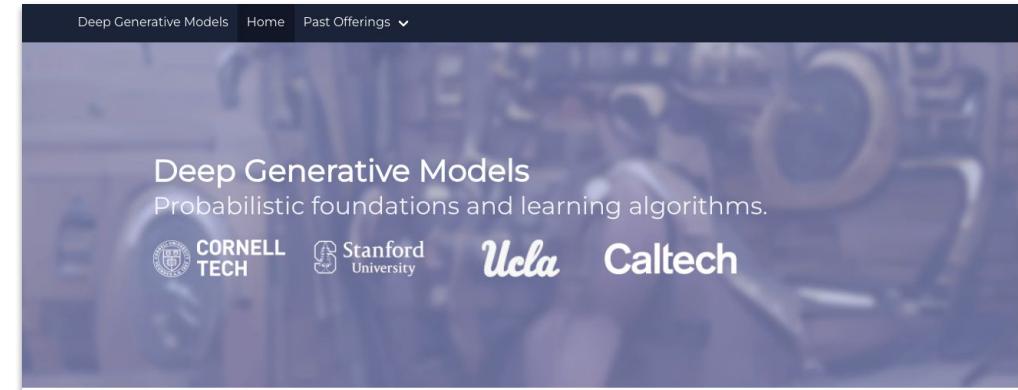
This Class – Inspiration

Deep Generative Models (DGM) class

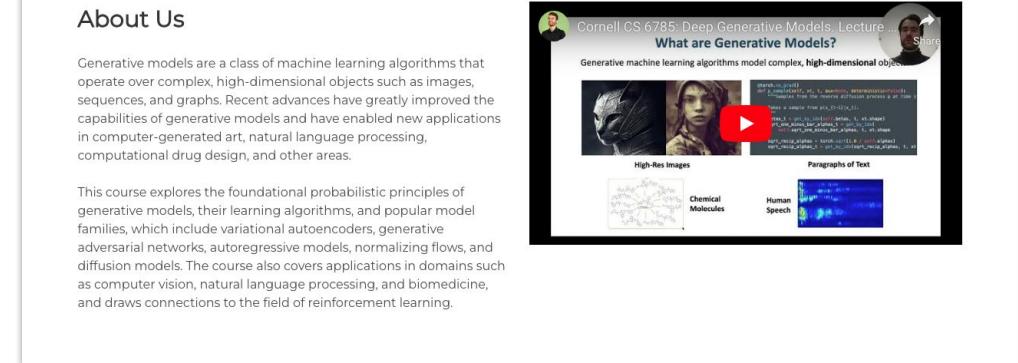
→ e.g., at Stanford, Cornell, UCLA, etc.



The screenshot shows the Stanford Deep Generative Models (CS236 - Fall 2023) course website. It features a header with the course name and a small red logo. Below the header is a 'Course Description' section with a detailed paragraph about generative models. Underneath is a 'Lecture Attendance' section with a note about lecture recordings. A 'Stanford Honor Code' section follows, with links to the general code and specific course details. At the bottom are five navigation buttons: 'Course Notes' (blue), 'Syllabus' (red), 'Ed' (orange), 'Office Hours' (green), and 'Course Project' (red).



The screenshot shows the Stanford Deep Generative Models course website. It has a dark header with the course name and navigation links for 'Home' and 'Past Offerings'. The main content area features a large banner with the text 'Deep Generative Models' and 'Probabilistic foundations and learning algorithms.' Below the banner are logos for Cornell Tech, Stanford University, UCLA, and Caltech. The background of the page is a blurred image of a person in a lab coat.



The screenshot shows the Cornell Deep Generative Models (CS 6785) course website. It includes an 'About Us' section with a detailed paragraph about generative models. Below this is a 'What are Generative Models?' section featuring a video player and three examples: 'High-Res Images', 'Chemical Molecules', and 'Human Speech'. The background of the page is a blurred image of a person in a lab coat.

Probabilistic Modeling – What is a probabilistic model?

Probabilistic Modeling – What is a probabilistic model?

It's common to define a mathematical model of the world in the form of an equation.

Probabilistic Modeling – What is a probabilistic model?

It's common to define a mathematical model of the world in the form of an equation.

Example: a simple model is a linear equation of the form:

$$y = \beta^T x$$

Where y is an outcome variable we want to predict, and x are known input variables that affect the outcome.

Probabilistic Modeling – What is a probabilistic model?

It's common to define a mathematical model of the world in the form of an equation.

Example: a simple model is a linear equation of the form:

$$y = \beta^T x$$

Where y is an outcome variable we want to predict, and x are known input variables that affect the outcome.

y – the price of a house.

x – features that affect the price (*location, # bedrooms, age of house, etc.*)

Assume that y is a linear function of these inputs, parameterized by β .

Probabilistic Modeling – Why be probabilistic?

Probabilistic Modeling – Why be probabilistic?

Often real-world has a significant amount of uncertainty
(e.g., the price y has a certain chance of going up if a new restaurant opens nearby that is popular...etc).

Probabilistic Modeling – Why be probabilistic?

Often real-world has a significant amount of uncertainty
(e.g., the price y has a certain chance of going up if a new restaurant opens nearby that is popular...etc).

It is natural to deal with this uncertainty by modeling the world in the form of a probability distribution: $p(x, y)$

Probabilistic Modeling – Why be probabilistic?

Often real-world has a significant amount of uncertainty
(e.g., the price y has a certain chance of going up if a new restaurant opens nearby that is popular...etc).

It is natural to deal with this uncertainty by modeling the world in the form of a probability distribution: $p(x, y)$

Given such a model, we could answer questions like:

Probabilistic Modeling – Why be probabilistic?

Often real-world has a significant amount of uncertainty

(e.g., the price y has a certain chance of going up if a new restaurant opens nearby that is popular...etc).

It is natural to deal with this uncertainty by modeling the world in the form of a probability distribution: $p(x, y)$

Given such a model, we could answer questions like:

- “What is the probability that housing price will be $> \$50,000$?”

Marginal Probability

Probabilistic Modeling – Why be probabilistic?

Often real-world has a significant amount of uncertainty

(e.g., the price y has a certain chance of going up if a new restaurant opens nearby that is popular...etc).

It is natural to deal with this uncertainty by modeling the world in the form of a probability distribution: $p(x, y)$

Given such a model, we could answer questions like:

- “What is the probability that housing price will be $> \$50,000$?”
- “Given that the house costs $\$100,000$, what is the probability that it has three bedrooms?”

Conditional Probability

Probabilistic Modeling – Why be probabilistic?

Often real-world has a significant amount of uncertainty

(e.g., the price y has a certain chance of going up if a new restaurant opens nearby that is popular...etc).

It is natural to deal with this uncertainty by modeling the world in the form of a probability distribution: $p(x, y)$

Given such a model, we could answer questions like:

- “What is the probability that housing price will be $> \$50,000$?”
- “Given that the house costs \$100,000, what is the probability that it has three bedrooms?”
- “Given that the house costs \$1,000,000, what is an example of a plausible house?”

Sampling

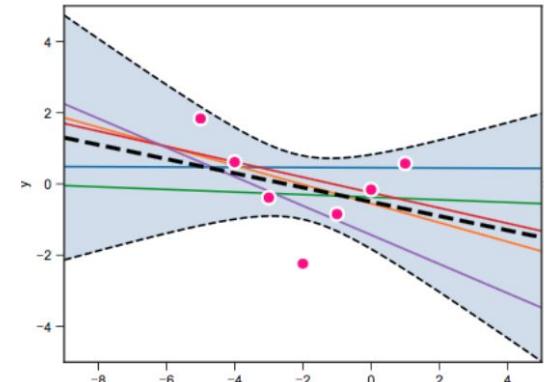
Probabilistic Modeling – Why be probabilistic?

Probability is important here! For a few key reasons:

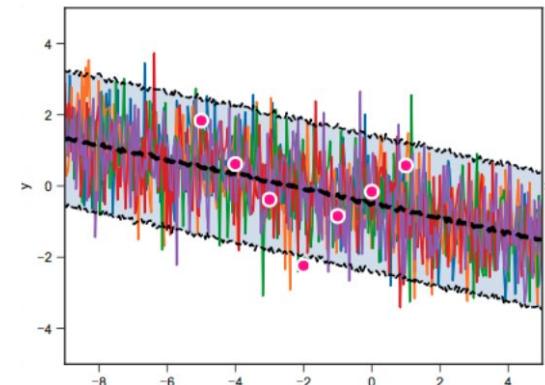
Probabilistic Modeling – Why be probabilistic?

Probability is important here! For a few key reasons:

1. Need to handle uncertainty about the environment.
We cannot make perfect predictions, and there is uncertainty due to:
 - lack of knowledge about the world (*epistemic uncertainty*).
 - inherent stochasticity on outcomes (*aleatoric uncertainty*).



Epistemic uncertainty over linear model parameters

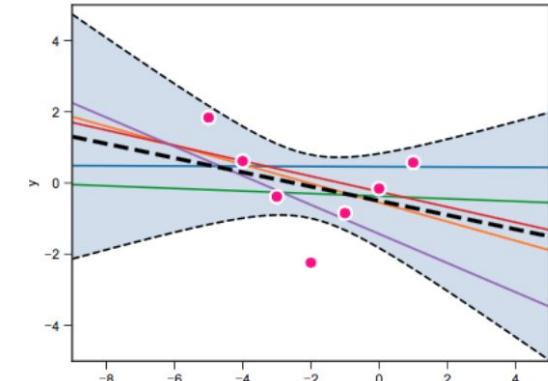


Aleatoric uncertainty shown for a single parameter

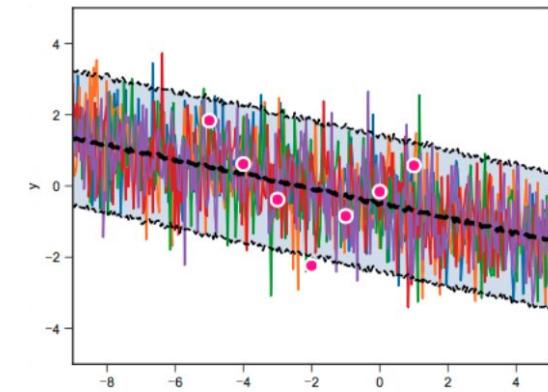
Probabilistic Modeling – Why be probabilistic?

Probability is important here! For a few key reasons:

1. Need to handle uncertainty about the environment.
We cannot make perfect predictions, and there is uncertainty due to:
 - lack of knowledge about the world (*epistemic uncertainty*).
 - inherent stochasticity on outcomes (*aleatoric uncertainty*).
2. Want to use uncertainty over our predictions.
Predicting a single value is often not enough.
We want a system that outputs its beliefs about what is going on in the world.



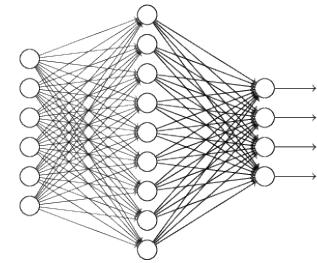
Epistemic uncertainty over linear model parameters



Aleatoric uncertainty shown for a single parameter

Probabilistic Modeling – Uncertainty

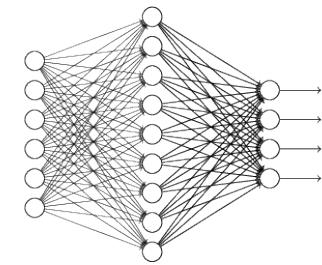
Machine learning is based on predictions – typically “point predictions”.



Dog

Probabilistic Modeling – Uncertainty

Machine learning is based on predictions – typically “point predictions”.



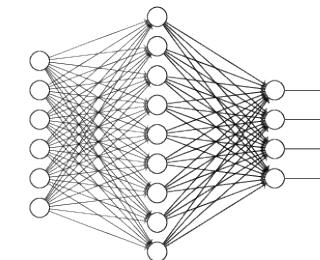
Dog?
Cat?

Rabid Cat?
Cotton Ball?
Trap?

Predictions **with reliable uncertainties** (e.g., with a measure of confidence) are useful for enabling trustworthy machine learning.

Probabilistic Modeling – Uncertainty

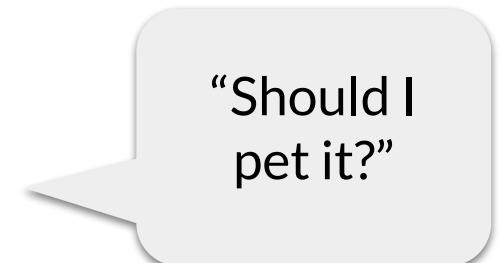
Machine learning is based on predictions – typically “point predictions”.



Predictions **with reliable uncertainties** (e.g., with a measure of confidence) are useful for enabling trustworthy machine learning.

This predictive uncertainty quantification is useful for *decision-making under uncertainty*.

(Applies to both classical probabilistic models and deep learning models.)



Deep Uncertainty Models

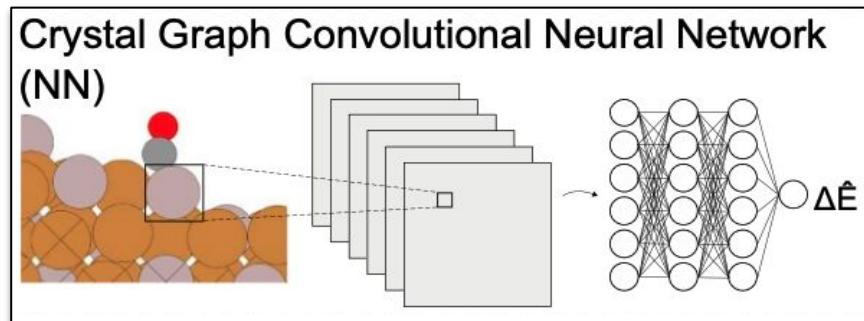
For example, we'll study methods for incorporating predictive uncertainty into neural networks:

Deep Uncertainty Models – Example: Computational Catalyst Design

For example, we'll study methods for incorporating predictive uncertainty into neural networks:

Deep Uncertainty Models – Example: Computational Catalyst Design

For example, we'll study methods for incorporating predictive uncertainty into neural networks:



Xie, Tian, and Jeffrey C. Grossman. "Crystal graph convolutional neural networks for an accurate and interpretable prediction of material properties." *Physical review letters*, 2018.

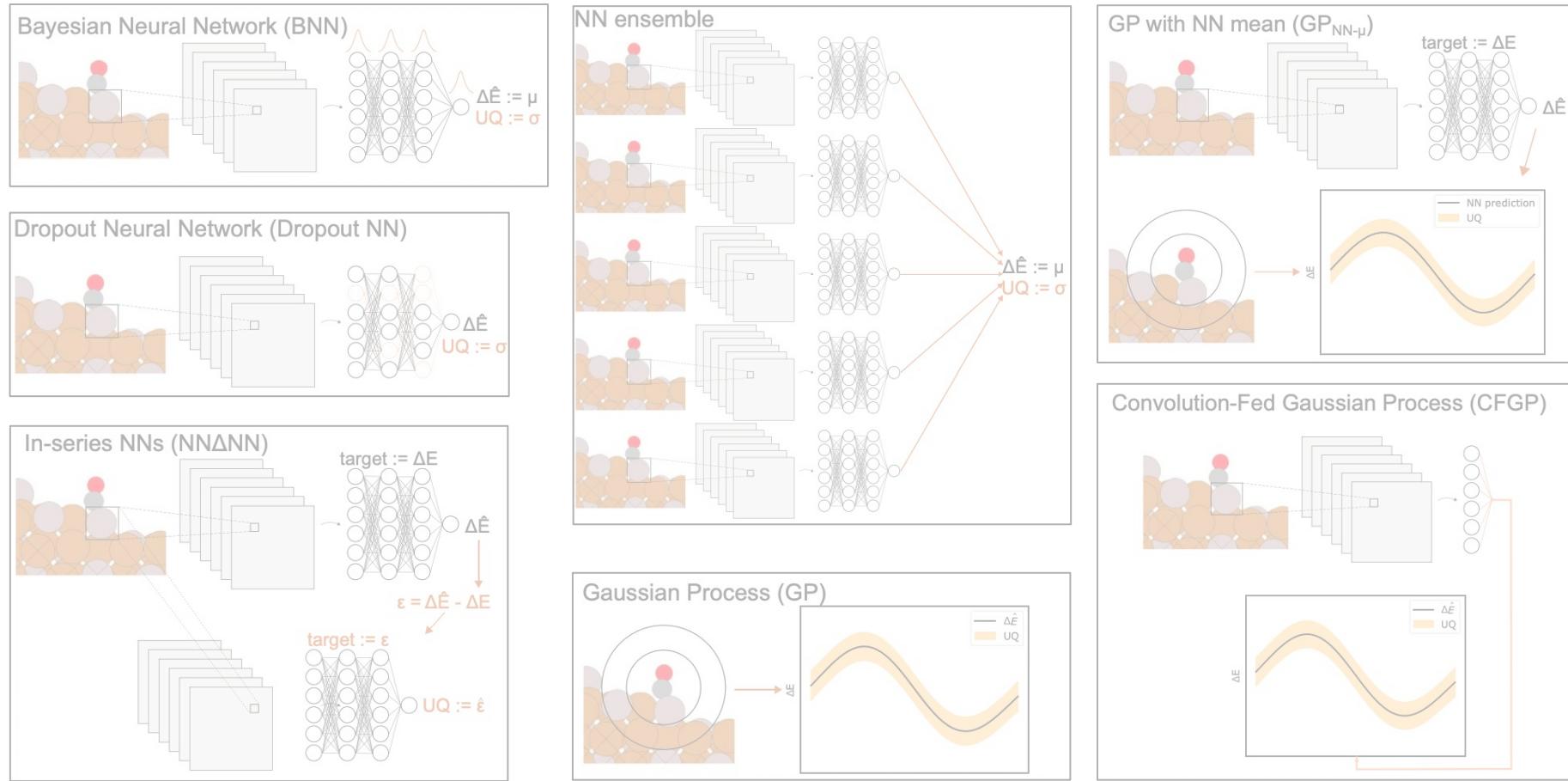
CGCNN Model (a graph convolutional neural network model)

- *Inputs:* three-dimensional atomic structure (a graph).
- *Outputs:* DFT-calculated site adsorption energies ΔE (*regression*).

"Methods for comparing uncertainty quantifications for material property predictions", *Tran, *Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

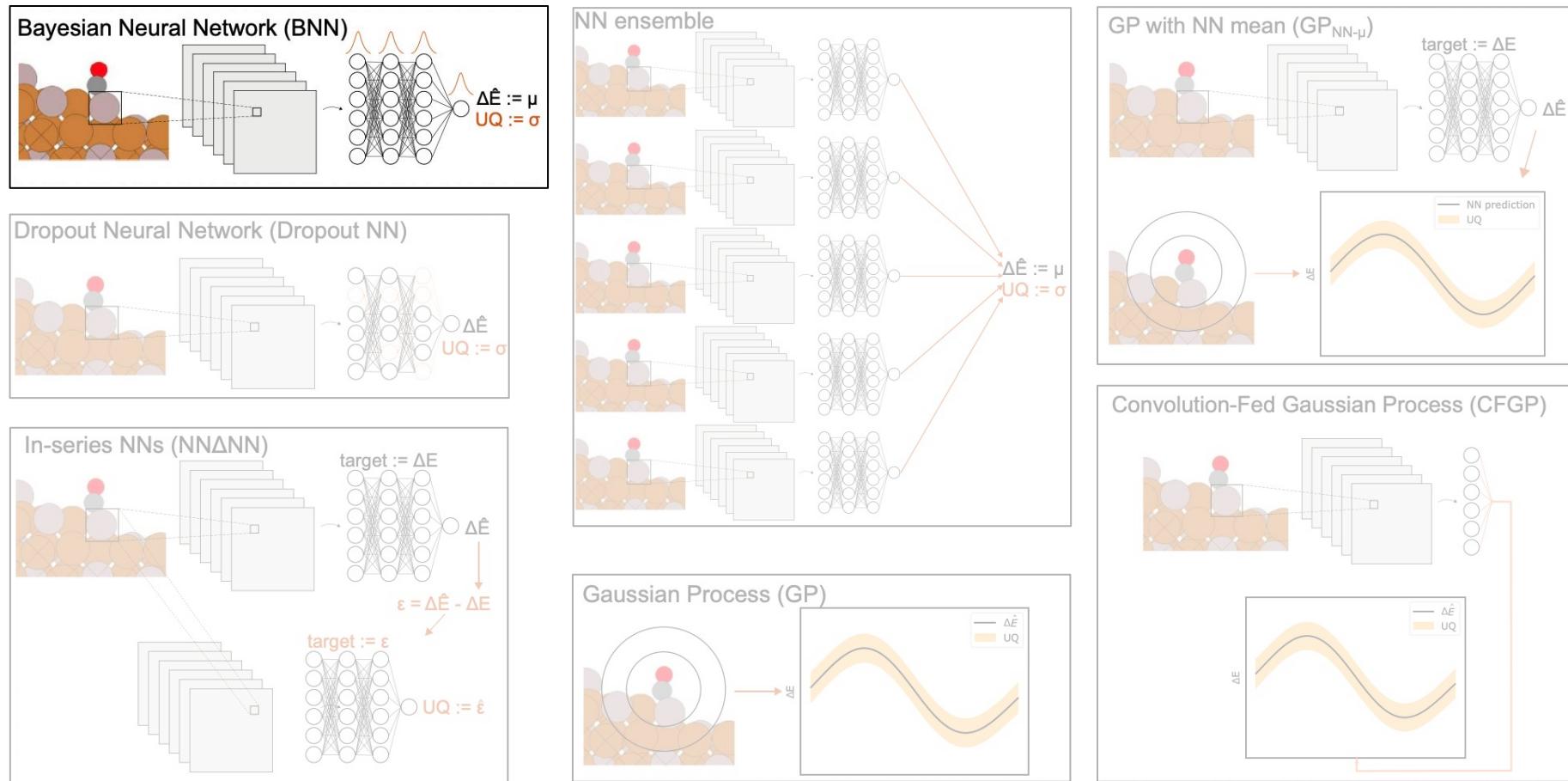
Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", *Tran, *Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

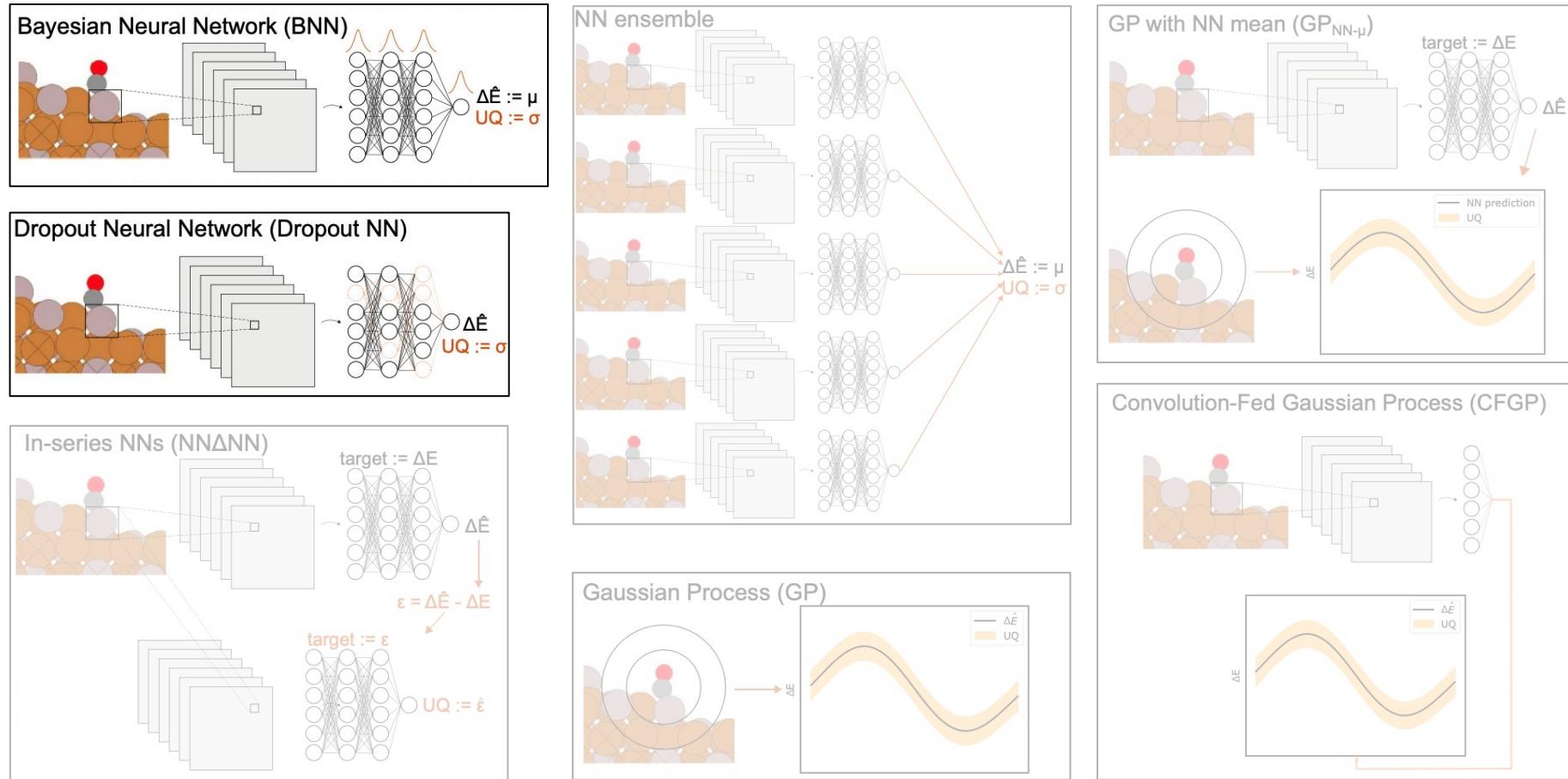
Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", *Tran, *Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

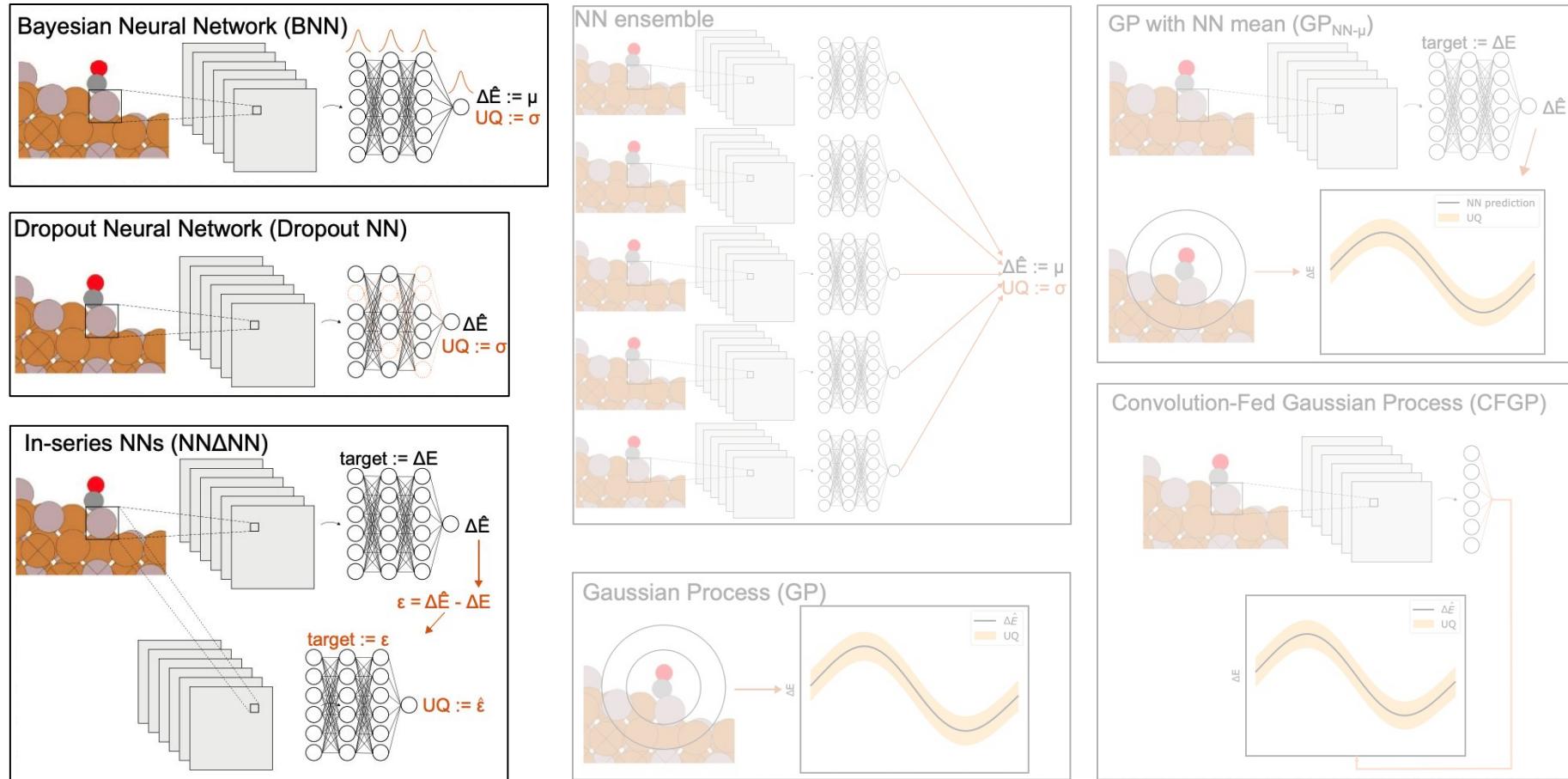
Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", *Tran, *Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

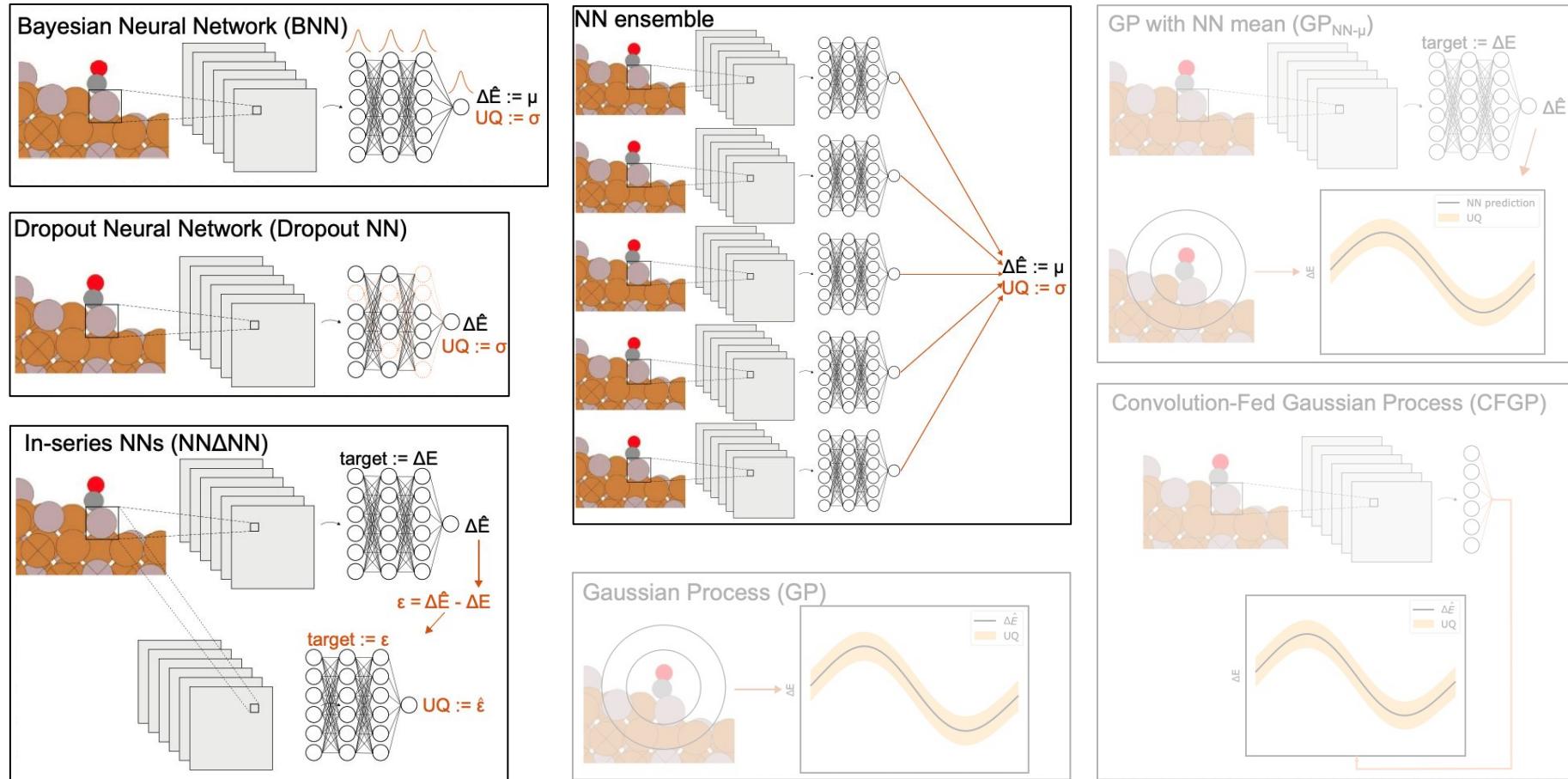
Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", *Tran, *Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

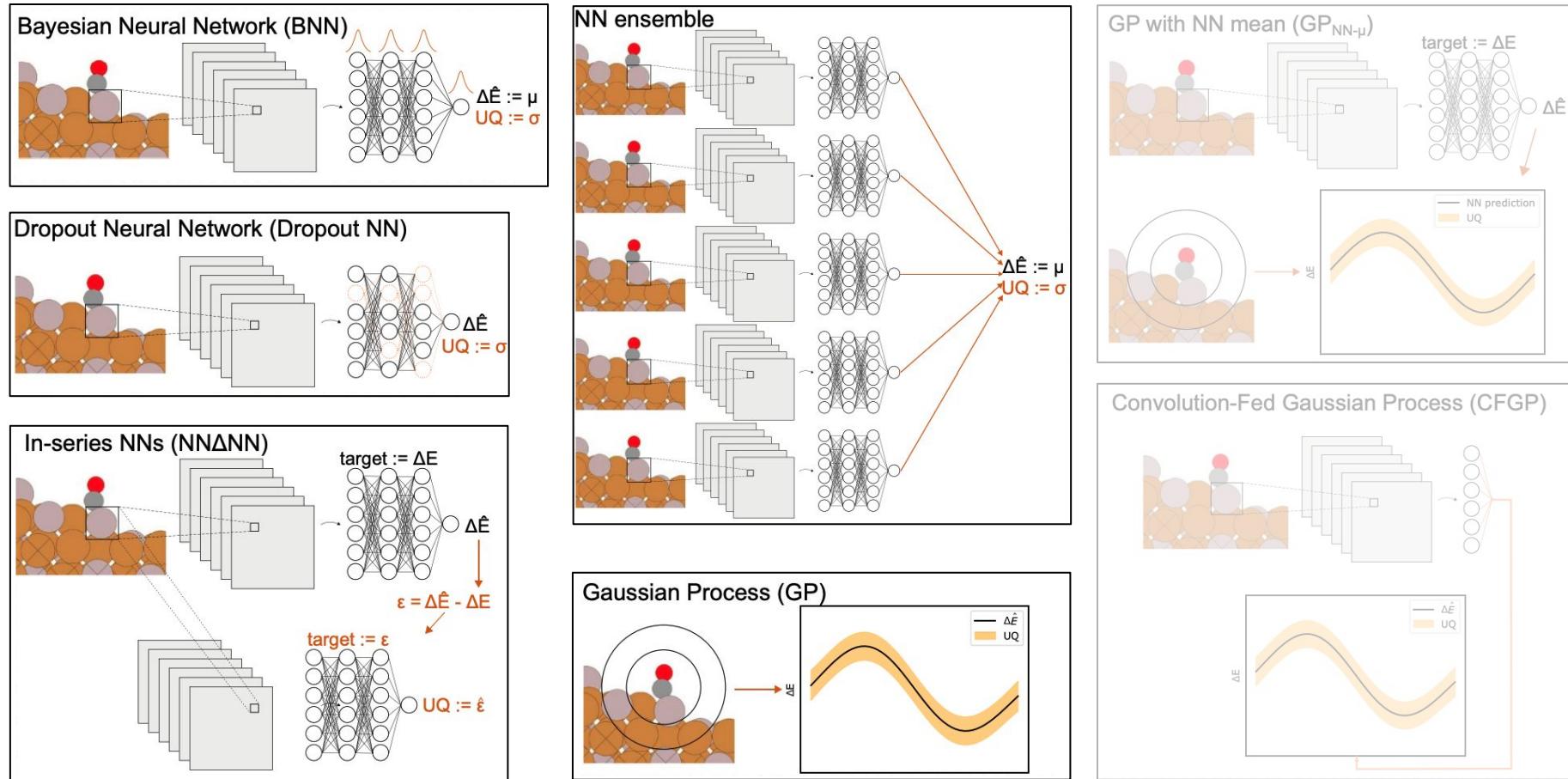
Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", *Tran, *Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

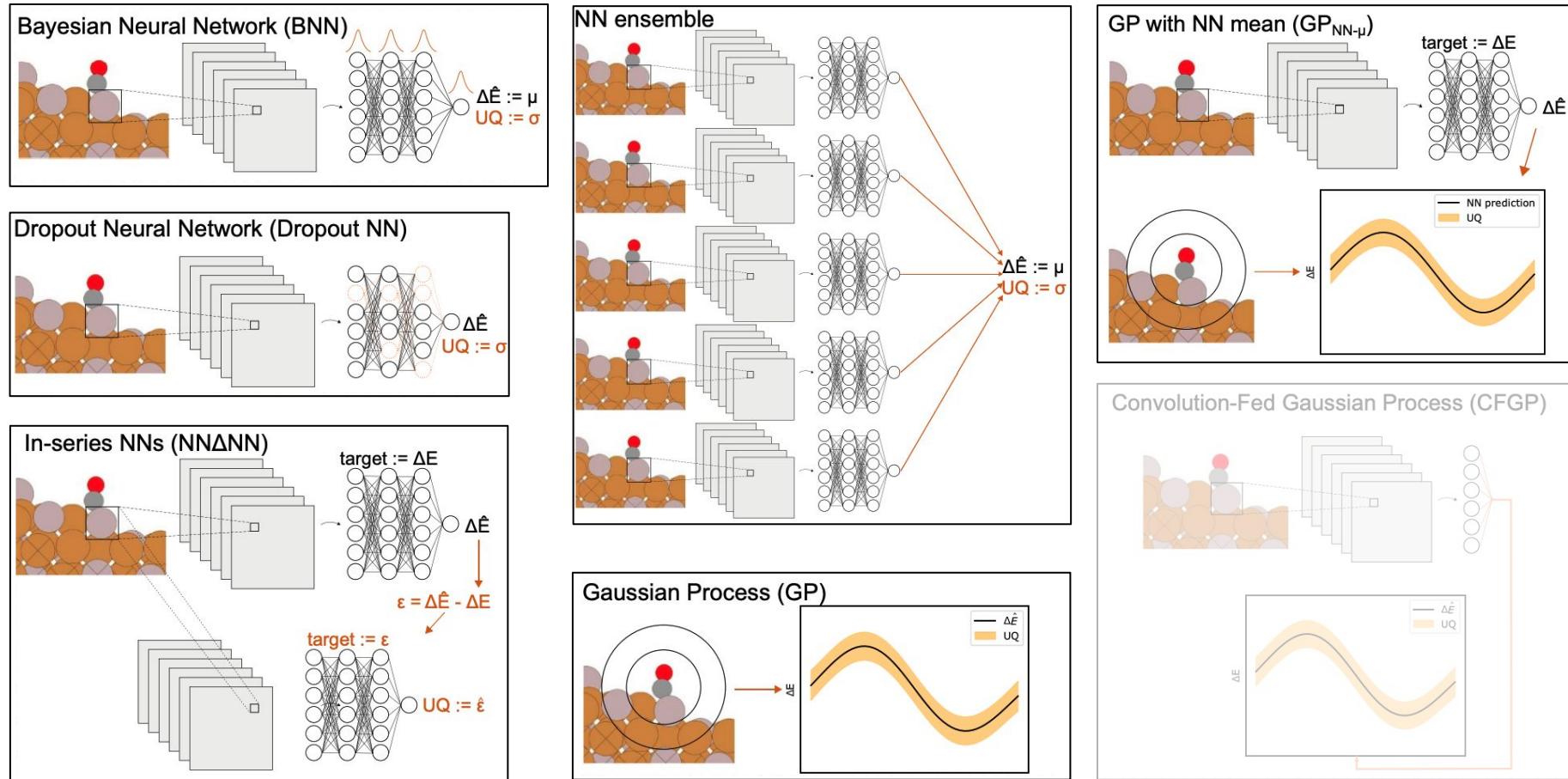
Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", *Tran, *Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

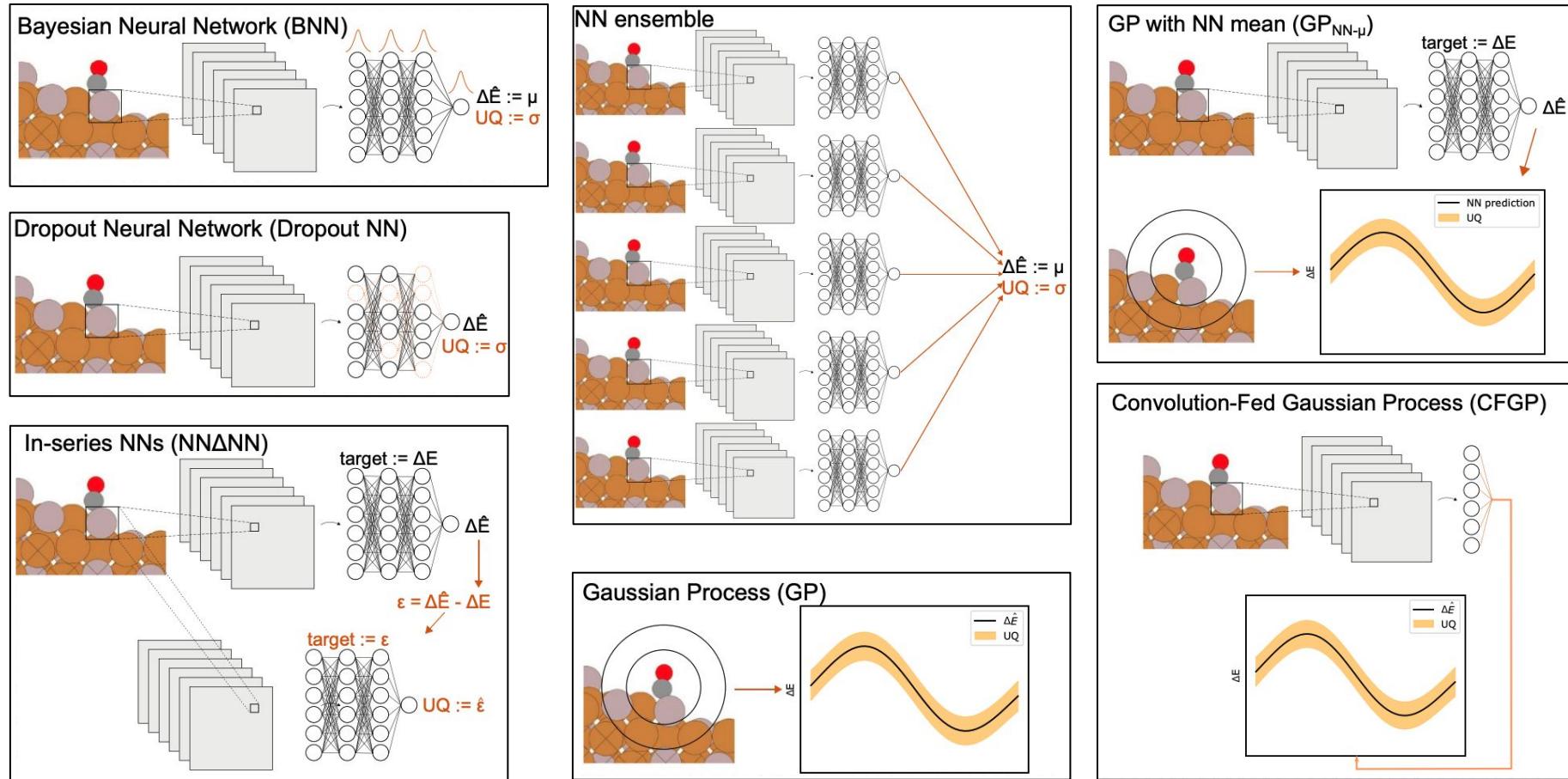
Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", *Tran, *Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", *Tran, *Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

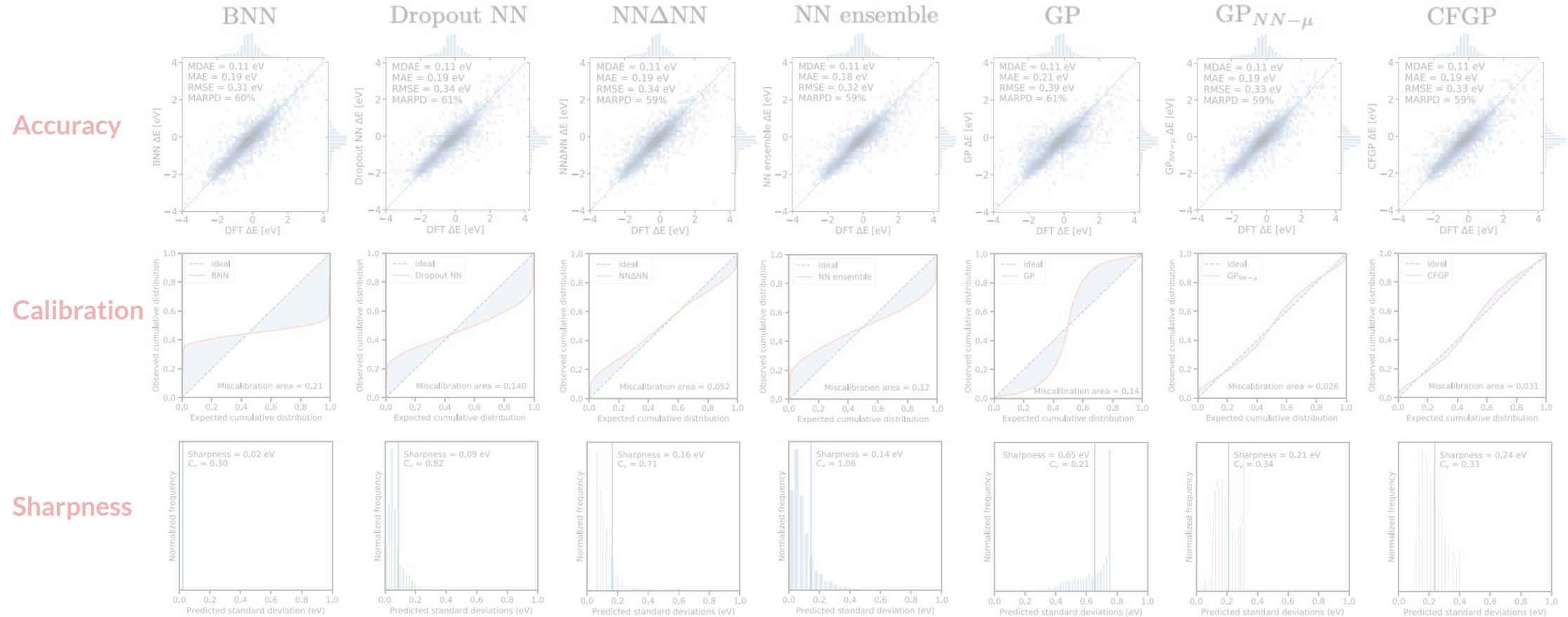
Deep Uncertainty Models – Example: Computational Catalyst Design

Metrics for deep uncertainty models.

“Methods for comparing uncertainty quantifications for material property predictions”, *Tran, *Neiswanger, et al., MLST. 2020

“Computational catalyst discovery: Active classification through myopic multiscale sampling”, *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

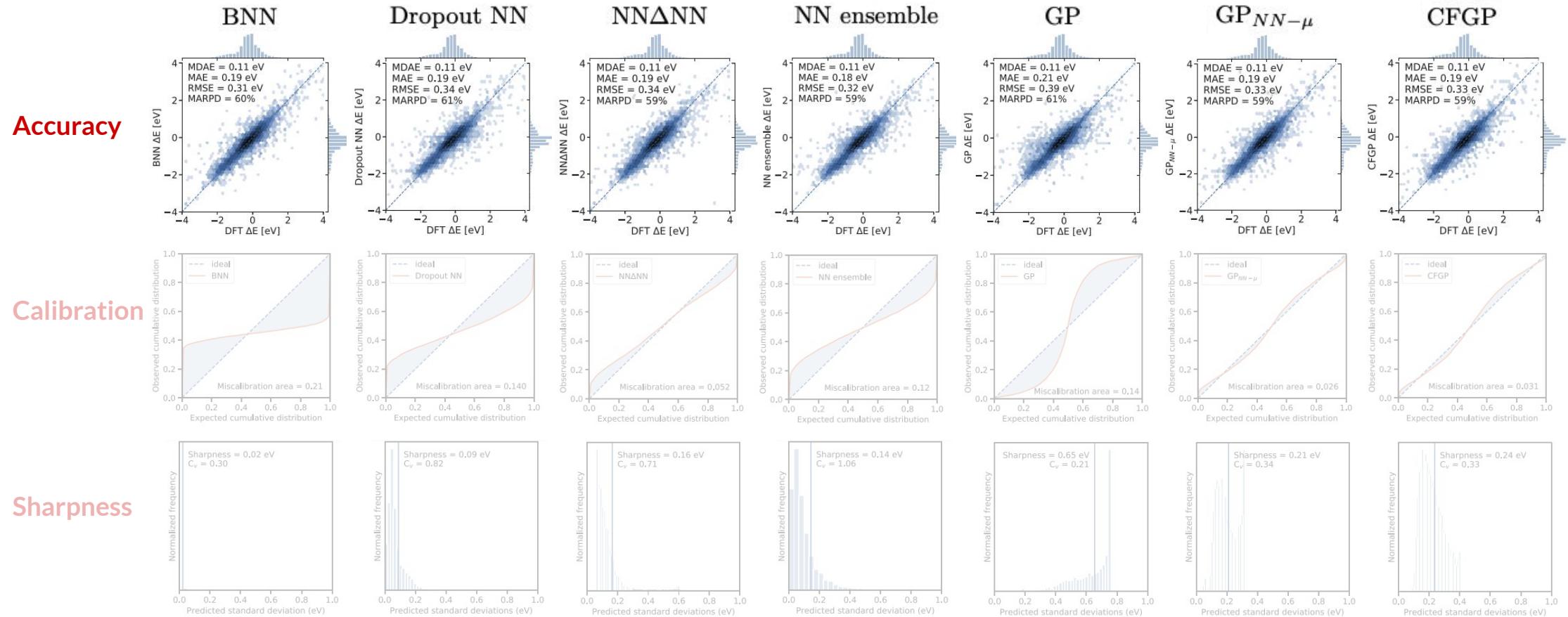
Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", *Tran, *Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

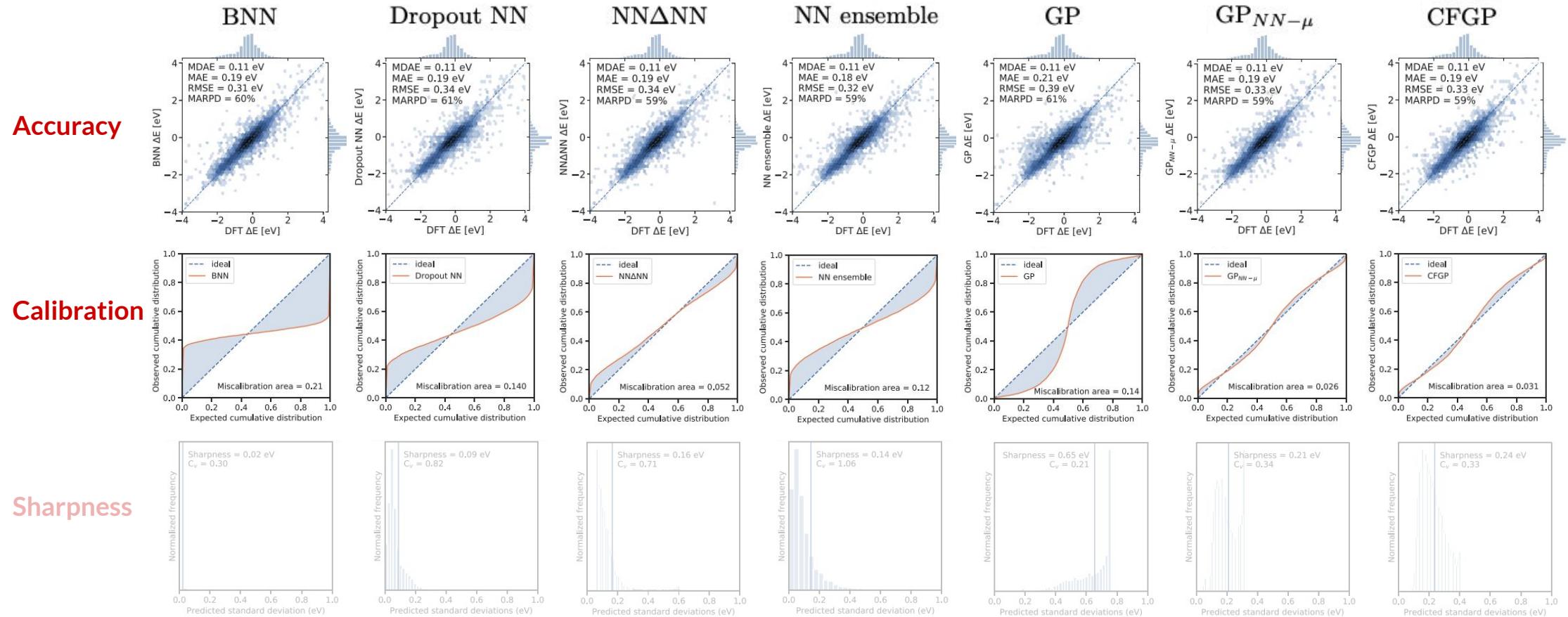
Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", *Tran, *Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

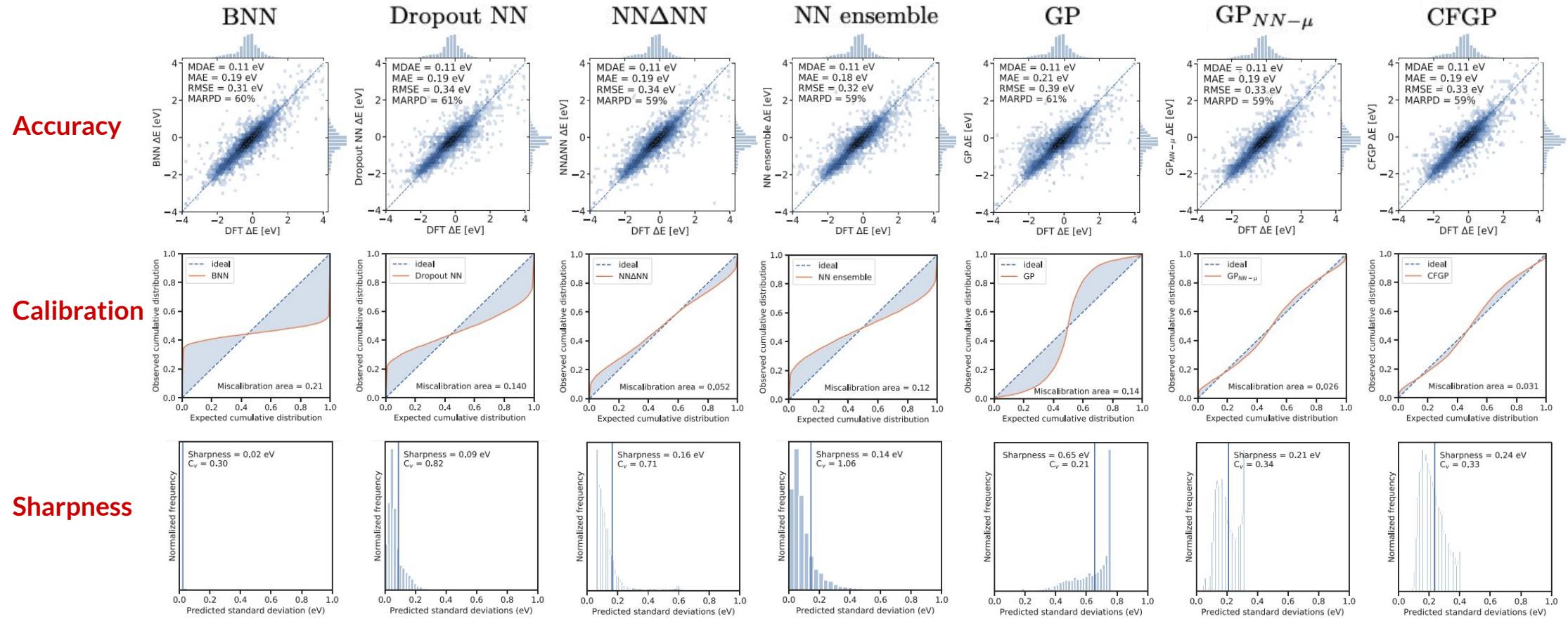
Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", *Tran, *Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

Deep Uncertainty Models – Example: Computational Catalyst Design



"Methods for comparing uncertainty quantifications for material property predictions", *Tran, *Neiswanger, et al., MLST. 2020

"Computational catalyst discovery: Active classification through myopic multiscale sampling", *Tran, *Neiswanger, et al., J. Chem. Phys. 2021

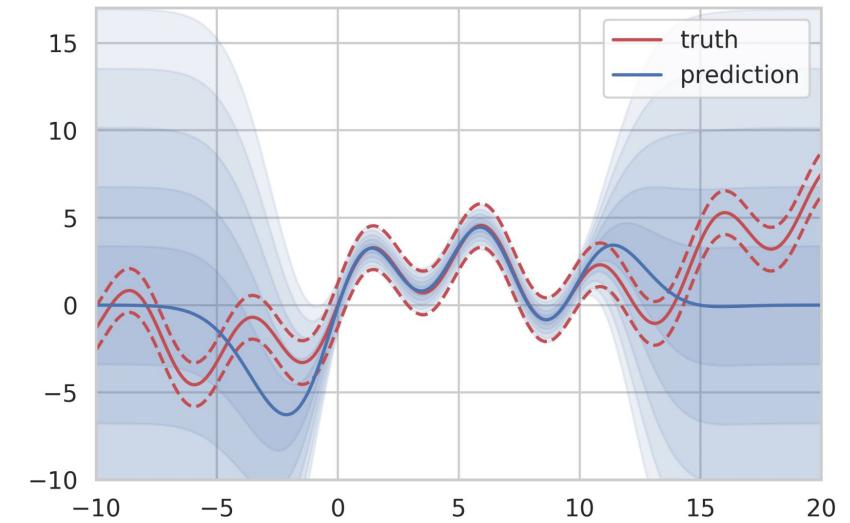
Deep Uncertainty Models – Decision Making Under Uncertainty

Many applications of these models in
decision making under uncertainty.
(Especially in the sciences and engineering).

Deep Uncertainty Models – Decision Making Under Uncertainty

Many applications of these models in decision making under uncertainty.
(Especially in the sciences and engineering).

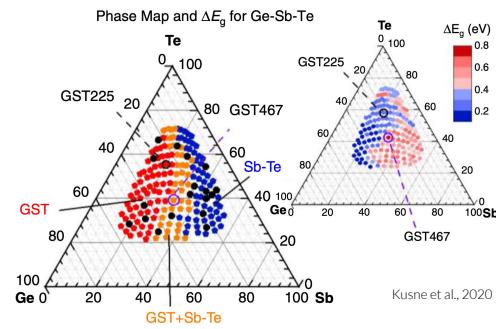
E.g., using probabilistic models for optimization, experimental design, and active learning.



Deep Uncertainty Models – Decision Making Under Uncertainty

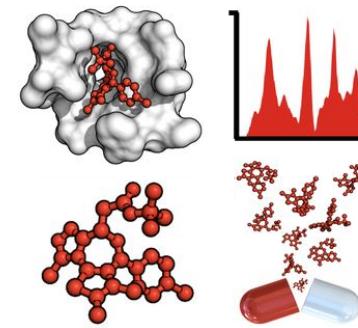
SCIENCE

Materials design



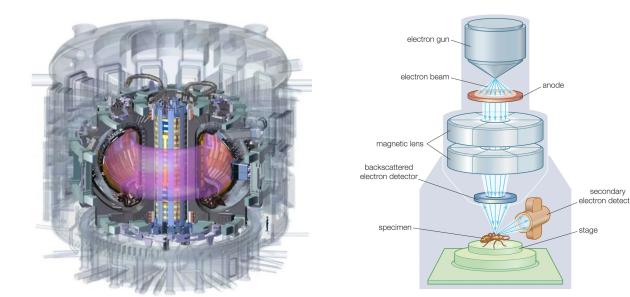
Materials design over large spaces of compounds.

Drug discovery



Target identification, lead optimization.

Large Science Machines

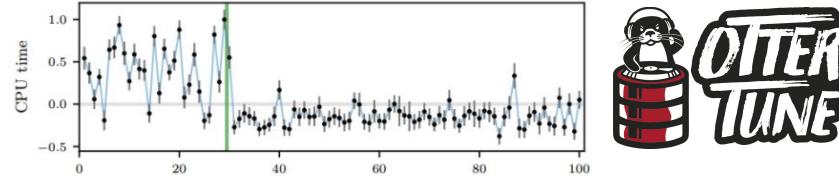


Optimize settings & controls for science use-cases.

Deep Uncertainty Models – Decision Making Under Uncertainty

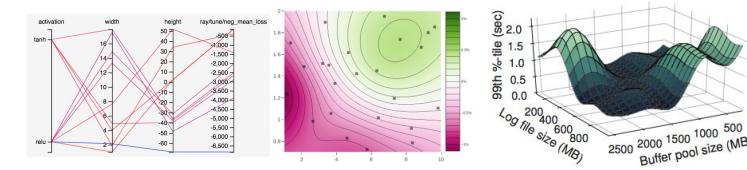
ENGINEERING

Computer Systems



Monitor and tune configurations for performance (latency, throughput).

Machine Learning



Find and assess new methods, models, architectures, hyperparameters.

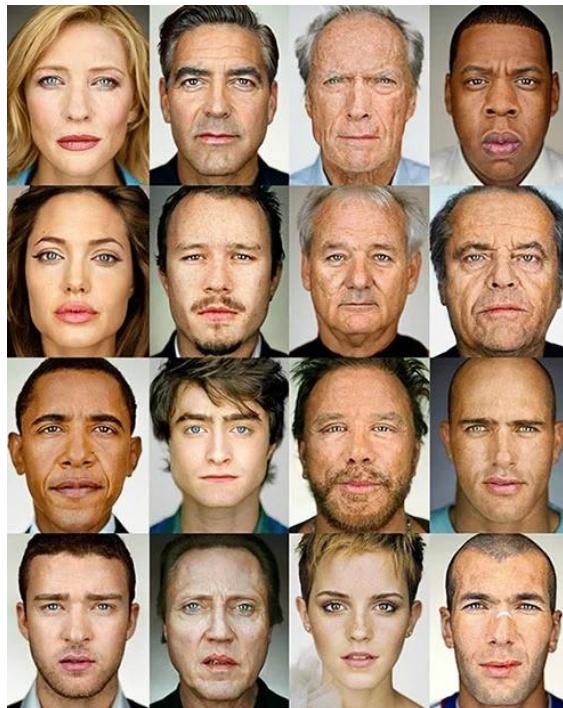
Probabilistic Modeling – High Dimensions

Often we want to model distributions of **high-dimensional data**.

Probabilistic Modeling – High Dimensions

Often we want to model distributions of **high-dimensional data**.

E.g., joint distributions over **image** or **video** data.



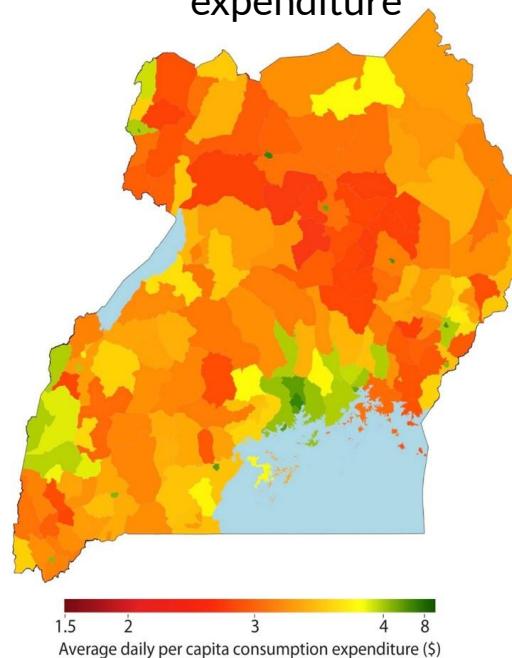
Kevin Beaty, "We made gif portraits of everyone we met on Colfax for 5 hours"

Probabilistic Modeling – High Dimensions

Often we want to model distributions of **high-dimensional data**.

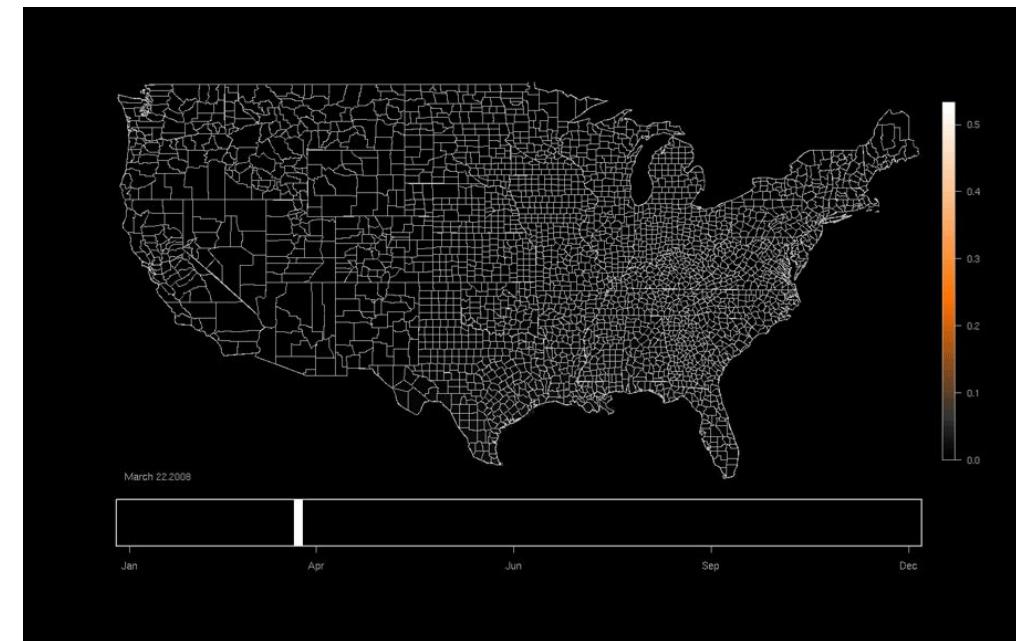
E.g., joint distributions over **spatiotemporal data**.

Uganda, est. daily per-capita expenditure



Source: Stefano Ermon, Stanford University

Seasonal Bird Migrations Trends (from eBird)

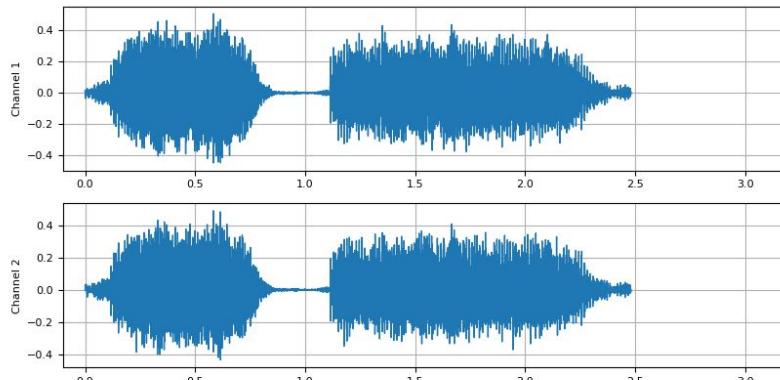


Source: Cornell Lab of Ornithology

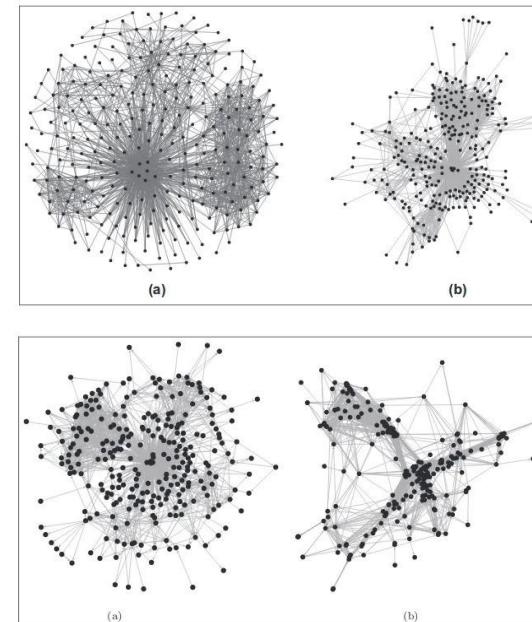
Probabilistic Modeling – High Dimensions

Often we want to model distributions of **high-dimensional data**.

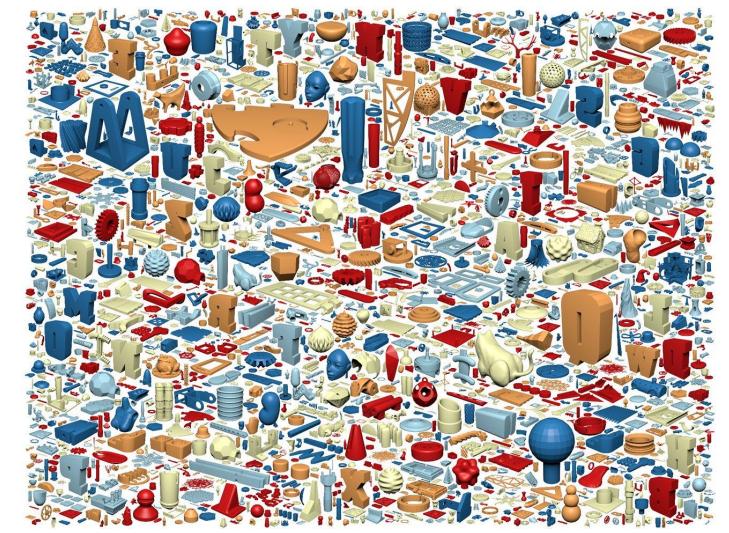
E.g., joint distributions over **audio**, **graph**, **3d structure** data, and more.



Source: Torchaudio package



Source: Sviatoslav Kovalev, "Large Graph Visualization Tools and Approaches"



Source: Thingi10K dataset.

Probabilistic Modeling – High Dimensions

A few tasks for high-dimensional probabilistic modeling include:

Probabilistic Modeling – High Dimensions

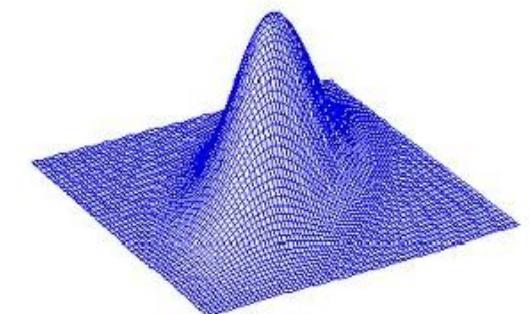
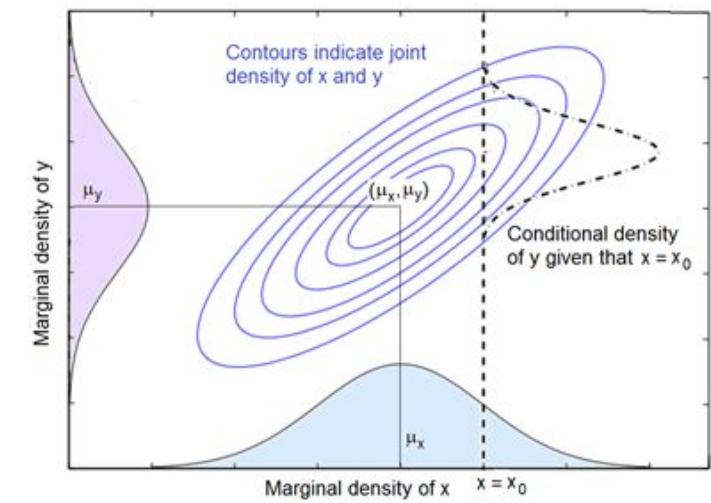
A few tasks for high-dimensional probabilistic modeling include:

1. Representation: What is the best way to represent or approximate a distribution? Exact specification? Simplified approximation? Samples from distribution?

Probabilistic Modeling – High Dimensions

A few tasks for high-dimensional probabilistic modeling include:

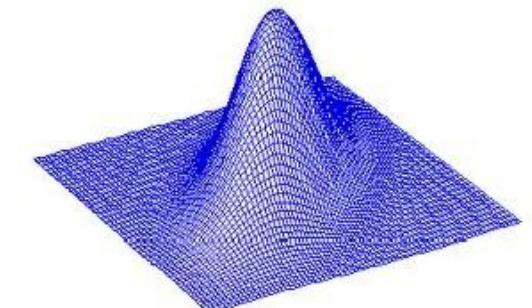
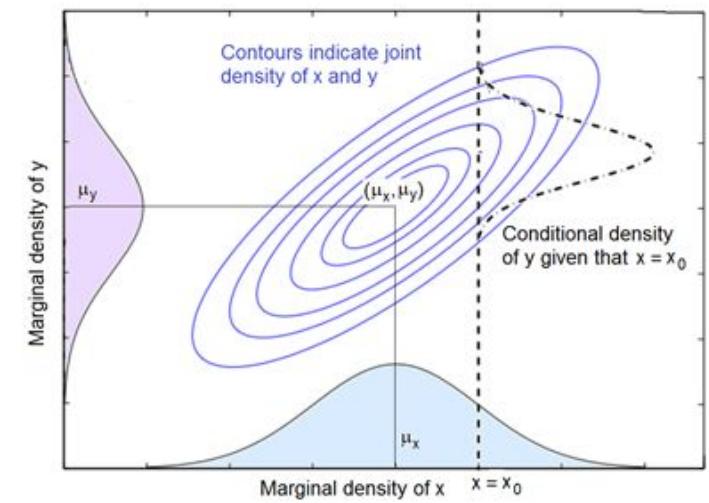
1. Representation: What is the best way to represent or approximate a distribution? Exact specification? Simplified approximation? Samples from distribution?
2. Inference: How do we infer marginal or conditional distributions? Or statistics of these distributions?



Probabilistic Modeling – High Dimensions

A few tasks for high-dimensional probabilistic modeling include:

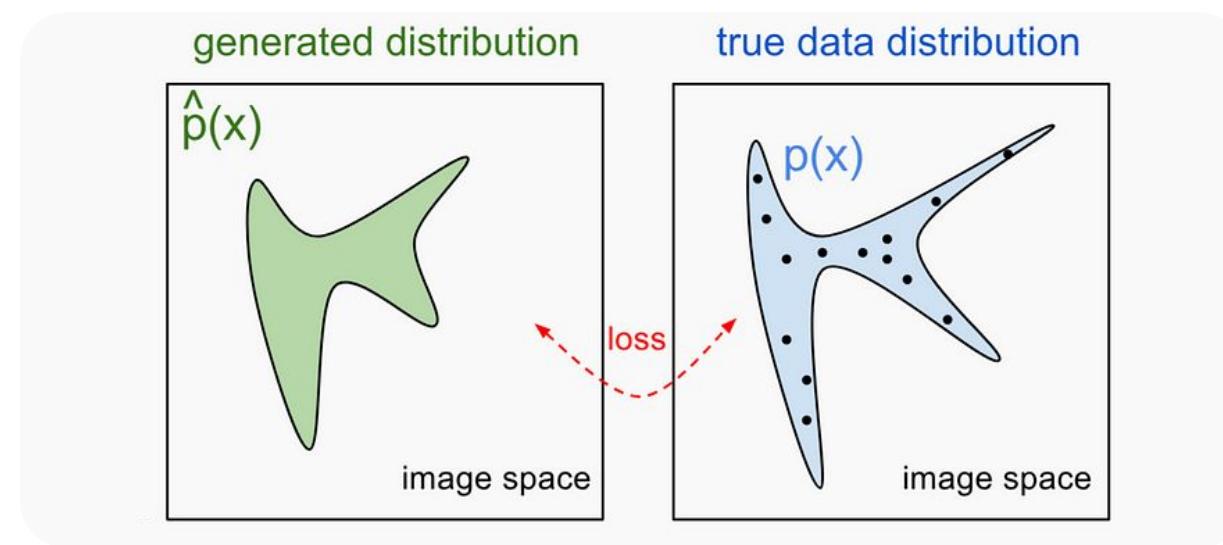
1. Representation: What is the best way to represent or approximate a distribution? Exact specification? Simplified approximation? Samples from distribution?
2. Inference: How do we infer marginal or conditional distributions? Or statistics of these distributions?
3. Learning: How can we fit a model to a high-dimensional dataset?



Probabilistic Models → Generative Models

4. Generation: After learning a model (distribution) from data, how can generate new samples from this distribution.

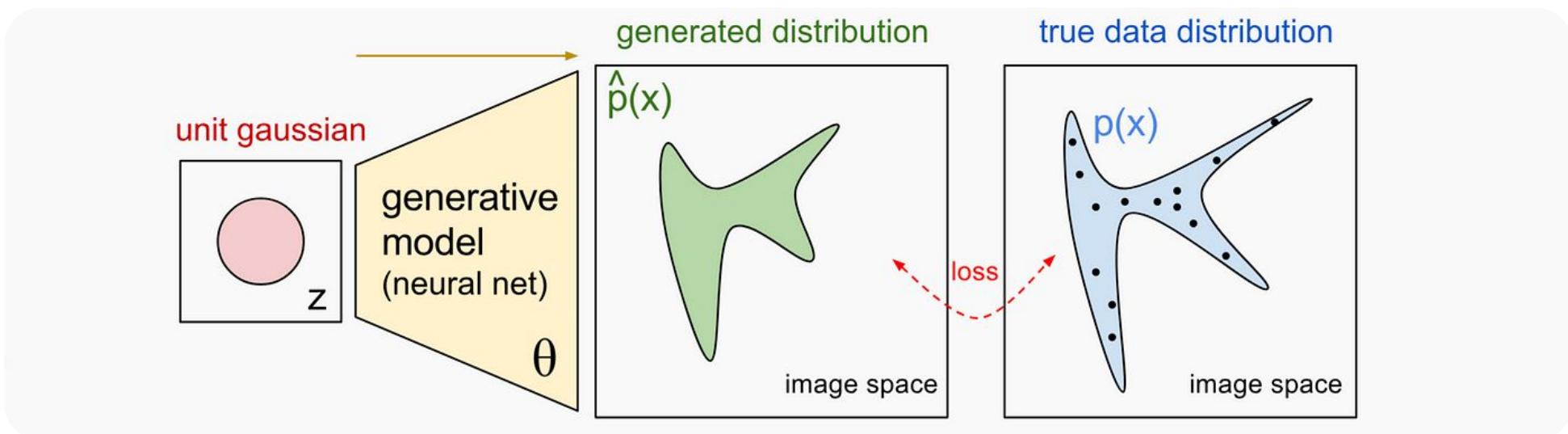
This leads us to the task of generative modeling.



Probabilistic Models → Generative Models

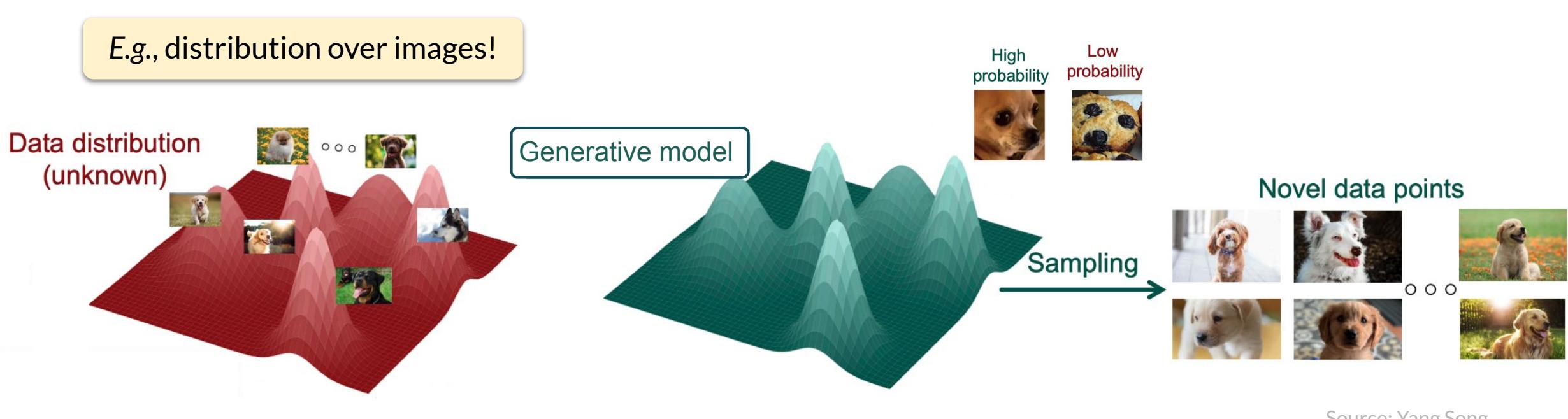
4. Generation: After learning a model (distribution) from data, how can generate new samples from this distribution.

This leads us to the task of ~~generative modeling~~ deep generative modeling.



Probabilistic Models → Generative Models

4. Generation: After learning a model (distribution) from data, how can generate new samples from this distribution.



Probabilistic Models → Generative Models

4. Generation: After learning a model (distribution) from data, how can generate new samples from this distribution.

Sometimes we aim to model/sample from:

- A marginal distribution, $p(x)$ – e.g., a distribution over images.
- A condition distribution, $p(x|y)$ – e.g., a distribution given a label y .

All probabilistic models!

Generative Models – Example Types

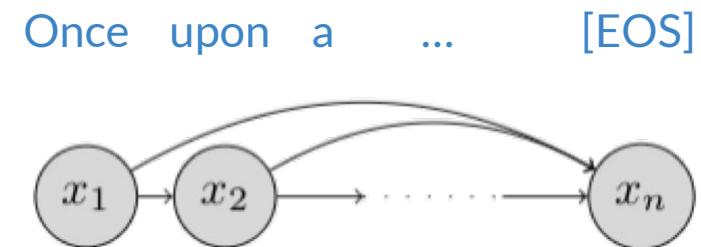
Generative Models Include:

Generative Models – Example Types

Generative Models Include:

$$\begin{aligned} p(x_1) \\ p(x_2 | x_1) \\ p(x_3 | x_2, x_1) \\ p(x_4 | x_3, x_2, x_1) \\ \dots \\ p(x_n | x_{n-1}, \dots, x_2, x_1) \end{aligned}$$

Autoregressive Models
(e.g., LLMs)



Generative Models – Example Types

Generative Models Include:

$$\begin{aligned}\mathcal{N}(x_1, x_2, x_3, \dots, x_n \mid \mathbf{0}, \mathbf{I}) \\ p_1(x_1, x_2, x_3, \dots, x_n) \\ p_2(x_1, x_2, x_3, \dots, x_n) \\ \dots \\ p_n(x_1, x_2, x_3, \dots, x_n) \\ p_{\text{data}}(x_1, x_2, x_3, \dots, x_n)\end{aligned}$$

**Diffusion and
Score-based Models**

Flow-based Models
(e.g. NF, CNF, FM)



Generative Modeling – Intuition for the Difficulty

It's interesting to consider the difficulty of the problem at hand.

Generative Modeling – Intuition for the Difficulty

It's interesting to consider the difficulty of the problem at hand.

Typical phone photo: 700 x 1400 pixels, each pixel in [0, 255], and (R, G, B) channels.

Generative Modeling – Intuition for the Difficulty

It's interesting to consider the difficulty of the problem at hand.

Typical phone photo: 700 x 1400 pixels, each pixel in [0, 255], and (R, G, B) channels.

⇒ Number of possible images: $256^{(700 \times 1400 \times 3)} \approx 10^{(7,000,000)}$ **WOW**

Generative Modeling – Intuition for the Difficulty

It's interesting to consider the difficulty of the problem at hand.

Typical phone photo: 700 x 1400 pixels, each pixel in [0, 255], and (R, G, B) channels.

⇒ Number of possible images: $256^{(700 \times 1400 \times 3)} \approx 10^{(7,000,000)}$ **WOW**

Number of atoms in universe is just 10^{80} :-).

Generative Modeling – Intuition for the Difficulty

It's interesting to consider the difficulty of the problem at hand.

Typical phone photo: 700 x 1400 pixels, each pixel in [0, 255], and (R, G, B) channels.

⇒ Number of possible images: $256^{(700 \times 1400 \times 3)} \approx 10^{(7,000,000)}$ **WOW**

Number of atoms in universe is just 10^{80} :-).

More importantly, number of images in ImageNet is 10^7 , and LAION is 10^9 .

Generative Modeling – Intuition for the Difficulty

It's interesting to consider the difficulty of the problem at hand.

Typical phone photo: 700 x 1400 pixels, each pixel in [0, 255], and (R, G, B) channels.

⇒ Number of possible images: $256^{(700 \times 1400 \times 3)} \approx 10^{(7,000,000)}$ **WOW**

Number of atoms in universe is just 10^{80} :-).

More importantly, number of images in ImageNet is 10^7 , and LAION is 10^9 .

⇒ Learning a generative model is a **highly underdetermined problem!**

Generative Models – Many Applications

And yet we have methods that work!

Generative Models – Many Applications

Applications include:

Generative Models – Many Applications

Applications include: **Image Generation**

Generative Models – Many Applications

Applications include: **Image Generation**



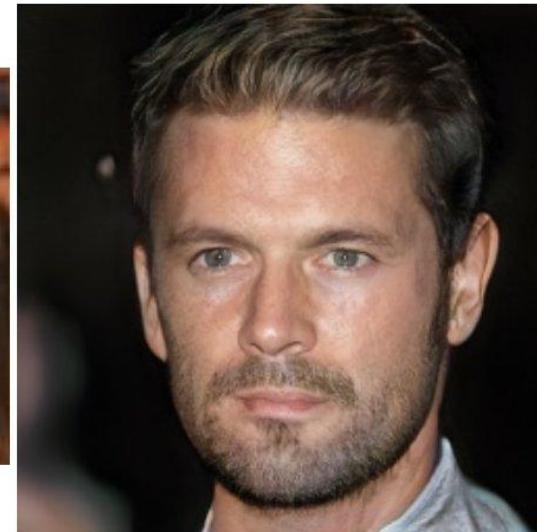
2014



2015



2016



2017



2018

$p(\text{images})$

Source: Ian Goodfellow, 2019, Twitter

https://x.com/goodfellow_ian/status/1084973596236144640

Generative Models – Many Applications

Applications include: **Image Generation** ... and in 2024:



Sources: Flux (Black Forest Labs) and Midjourney

Generative Models – Many Applications

Applications include: **Image Generation** ... and in 2024:



Sources: Flux (Black Forest Labs) and Midjourney

Generative Models – Many Applications

Applications include: **Video Generation** (more recently)

Generative Models – Many Applications

Applications include: **Video Generation** (more recently)



Sources: Sora (OpenAI) and Pika Labs.

$p(\text{next image} \mid \text{previous images})$

Generative Models – Many Applications

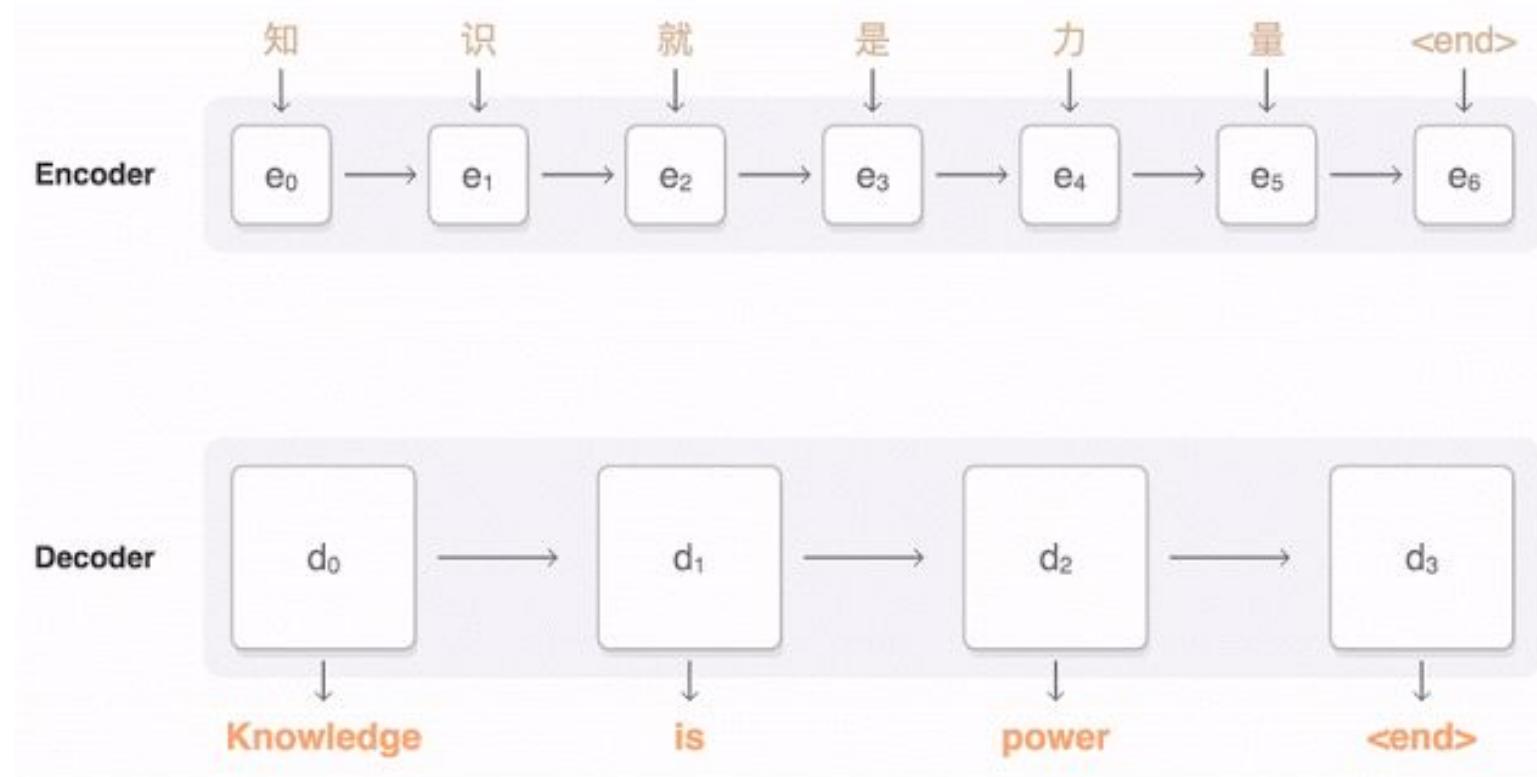
Applications include: **Text Generation**

Generative Models – Many Applications

Applications include: **Text Generation**

$p(\text{next word} \mid \text{previous words})$

2017: Machine Translation with Seq2seq via RNN and attention mechanism

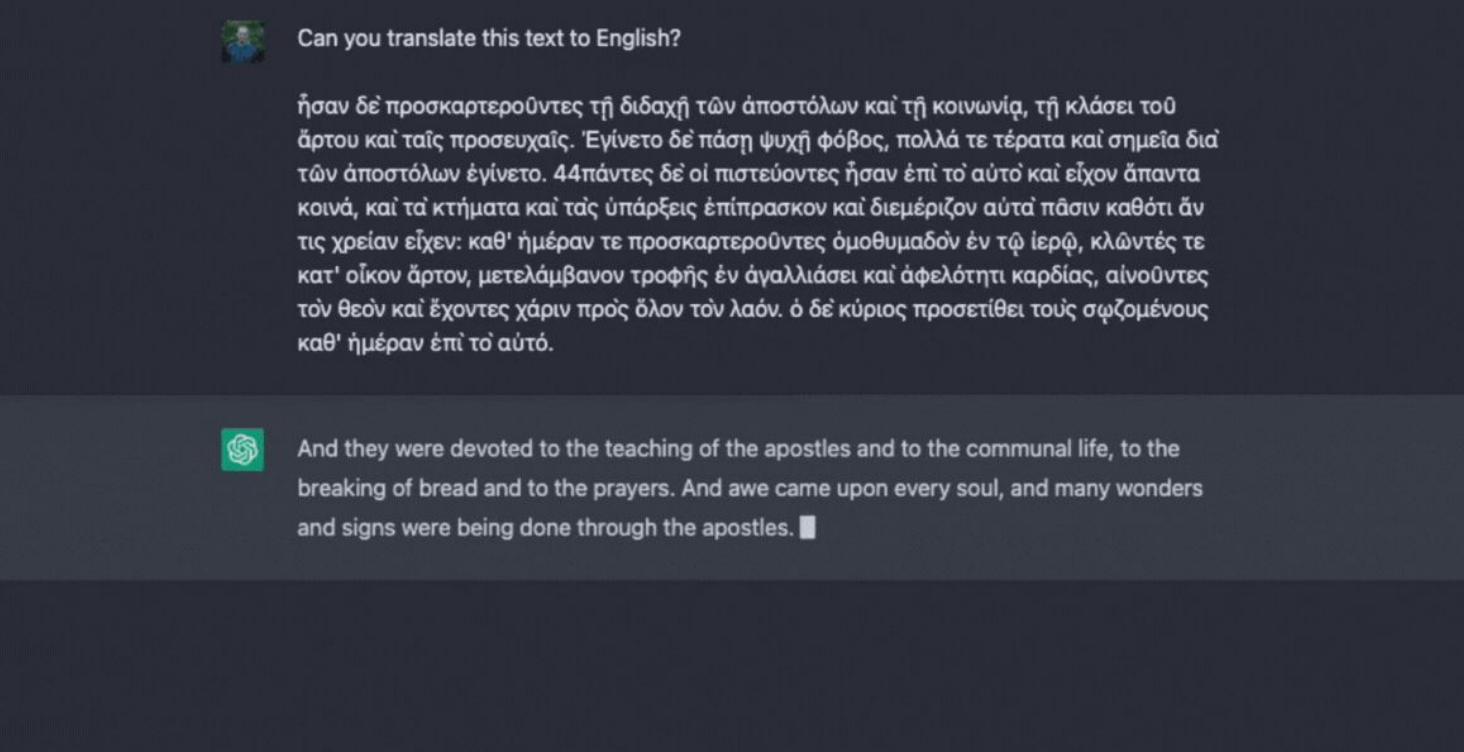


Generative Models – Many Applications

Applications include: **Text Generation**

$p(\text{next word} \mid \text{previous words})$

2024: ChatGPT from OpenAI



The screenshot shows a conversation between a user and ChatGPT. The user asks, "Can you translate this text to English?" Below is the Greek text followed by the English translation generated by ChatGPT.

ἡσαν δὲ προσκαρτεροῦντες τῇ διδαχῇ τῶν ἀποστόλων καὶ τῇ κοινωνίᾳ, τῇ κλάσει τοῦ ἄρτου καὶ ταῖς προσευχαῖς. Ἐγίνετο δὲ πάσῃ ψυχῇ φόβος, πολλά τε τέρατα καὶ σημεῖα δια' τῶν ἀποστόλων ἐγίνετο. 44πάντες δὲ οἱ πιστεύοντες ἡσαν ἐπὶ τὸ αὐτὸν καὶ εἶχον ἄπαντα κοινά, καὶ τὰ κτήματα καὶ τὰς ὑπάρξεις ἐπίπρασκον καὶ διεμέριζον αὐτὰ πᾶσιν καθότι ἂν τις χρείαν εἶχεν: καθ' ἡμέραν τε προσκαρτεροῦντες ὁμοθυμαδὸν ἐν τῷ ιερῷ, κλώντες τε κατ' οἶκον ἄρτον, μετελάμβανον τροφῆς ἐν ἀγαλλιάσει καὶ ἀφελότητι καρδίας, αἰνοῦντες τὸν θεόν καὶ ἔχοντες χάριν πρὸς ὅλον τὸν λαόν. ὁ δὲ κύριος προσετίθει τοὺς σωζομένους καθ' ἡμέραν ἐπὶ τὸ αὐτό.

And they were devoted to the teaching of the apostles and to the communal life, to the breaking of bread and to the prayers. And awe came upon every soul, and many wonders and signs were being done through the apostles. ■

Source: Bureauworks

Generative Models – Many Applications

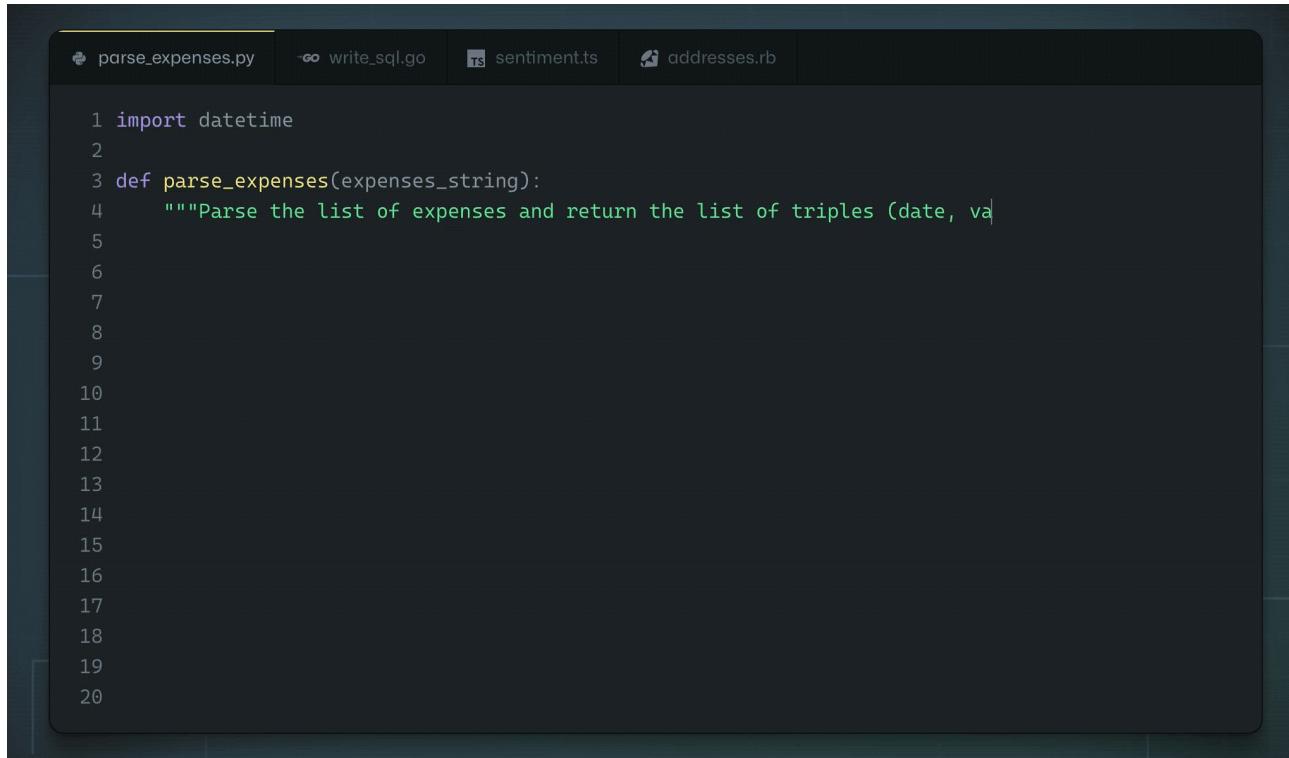
Applications include: **Code Generation**

Generative Models – Many Applications

Applications include: **Code Generation**

$p(\text{next word} \mid \text{previous words})$

2021: OpenAI Codex → GitHub Copilot



A screenshot of a dark-themed code editor showing a Python file named `parse_expenses.py`. The code defines a function `parse_expenses` that takes a string of expenses and returns a list of triples. The code is partially visible, ending with a placeholder at the end of the string parameter.

```
1 import datetime
2
3 def parse_expenses(expenses_string):
4     """Parse the list of expenses and return the list of triples (date, va
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
```

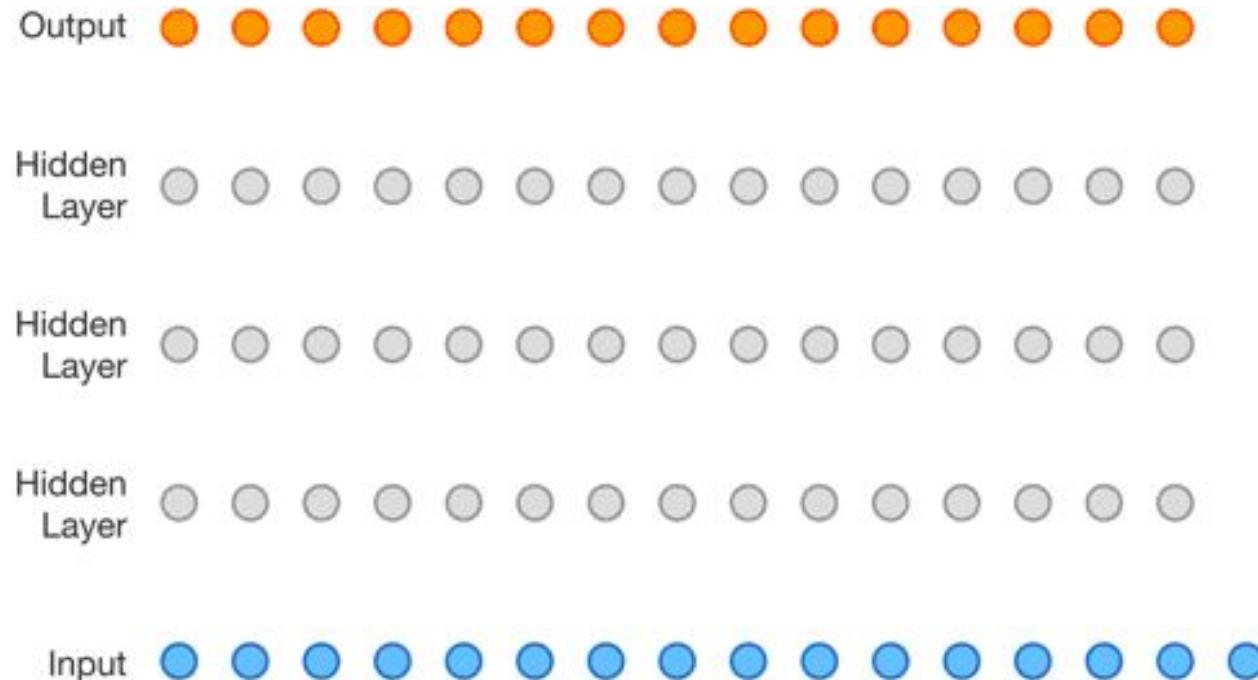
Source: OpenAI/Github Codex

Generative Models – Many Applications

Applications include: **Speech Generation**

Generative Models – Many Applications

Applications include: **Speech Generation** (Text to Speech) $p(\text{speech} \mid \text{text})$



WaveNet, 2016 (Google DeepMind)

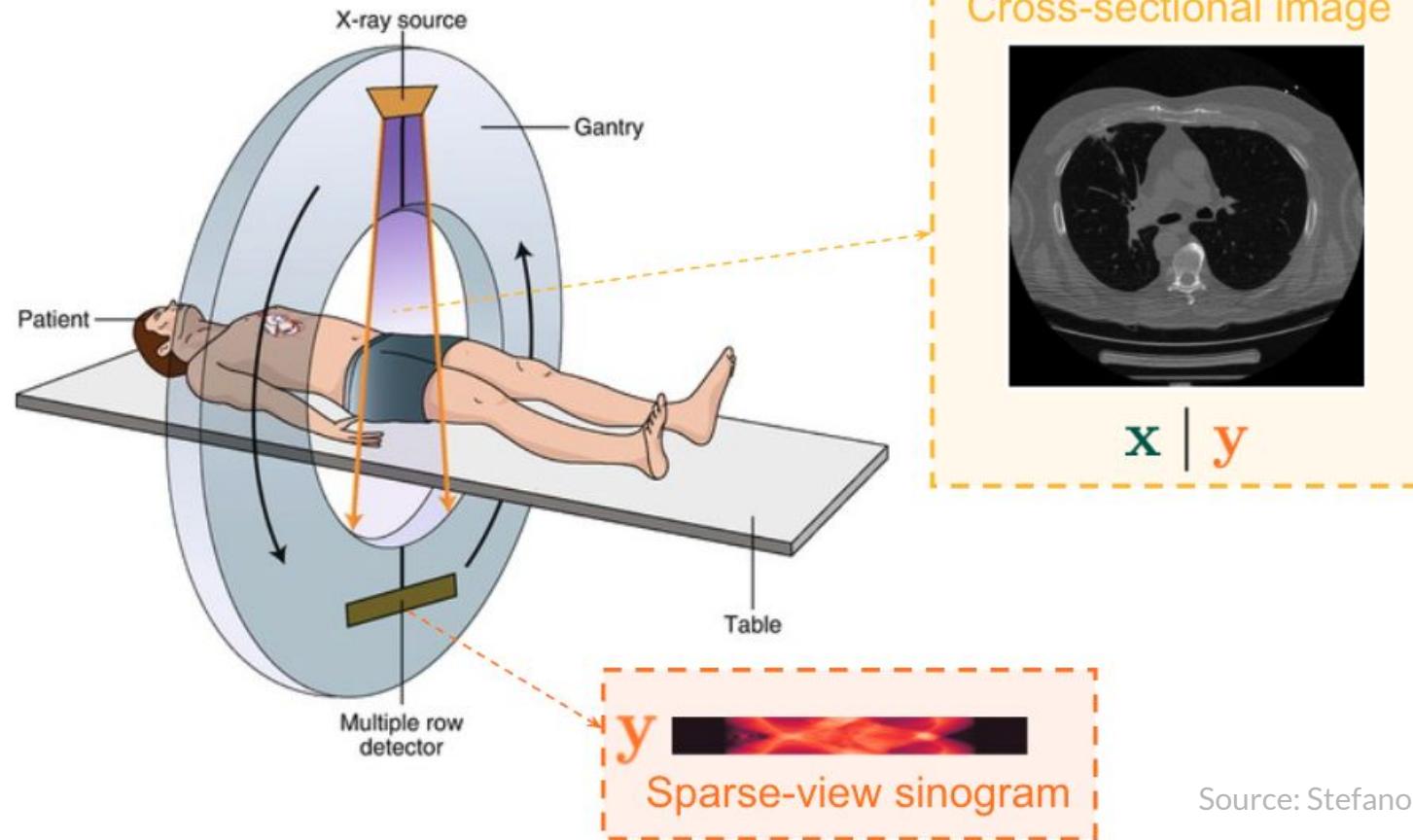
Generative Models – Many Applications

Applications include: **Medical Image Reconstruction**

Generative Models – Many Applications

Applications include: **Medical Image Reconstruction**

$p(\text{image} \mid \text{X-ray signal})$



Source: Stefano Ermon, Deep Generative Models

Generative Models – Many Applications

Applications include: **Science**

Generative Models – Many Applications

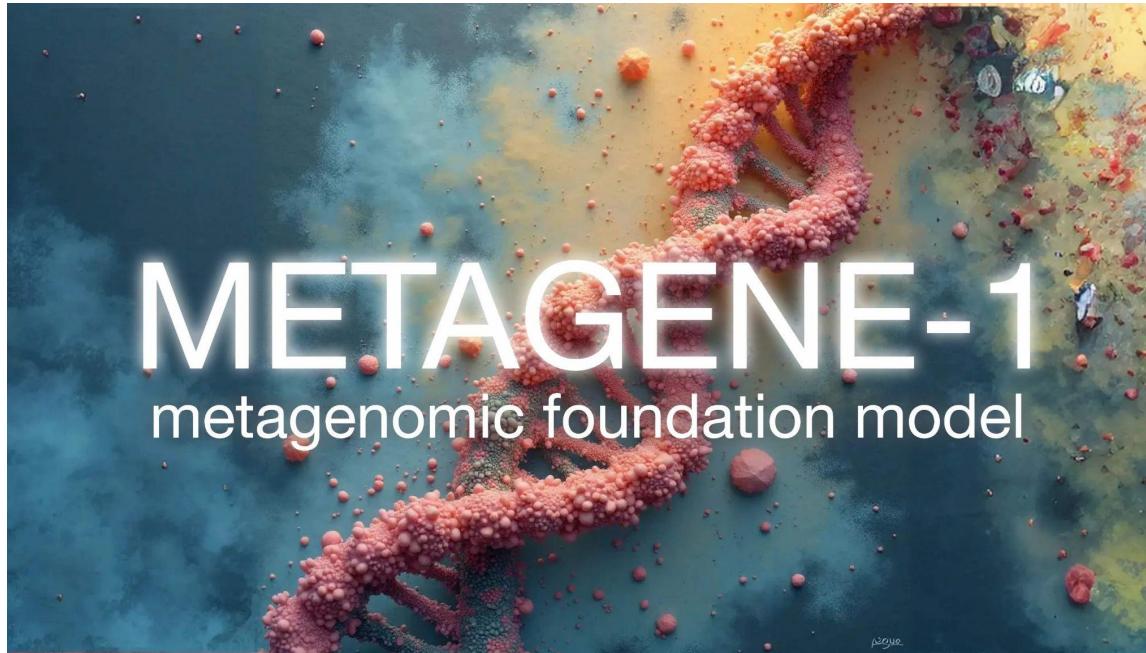
Applications include: **Science** – DNA, RNA, etc.

$p(\text{next DNA} \mid \text{previous DNA})$

Generative Models – Many Applications

Applications include: Science – DNA, RNA, etc.

$p(\text{next DNA} \mid \text{previous DNA})$



<https://metagene.ai>

Trained on DNA/RNA.

Autoregressive model - like LLMs.

7B parameters (~small LLMs).

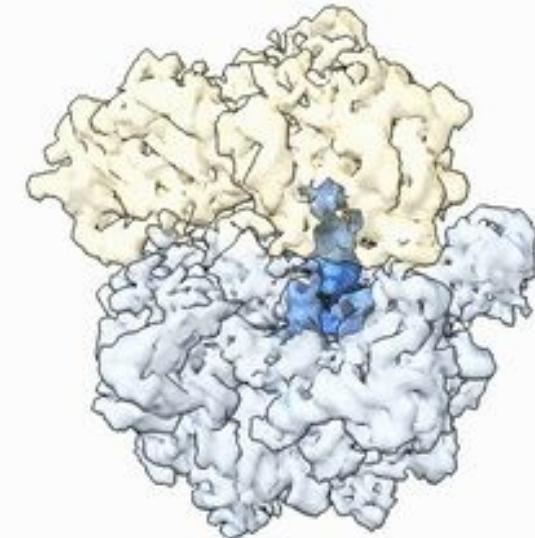
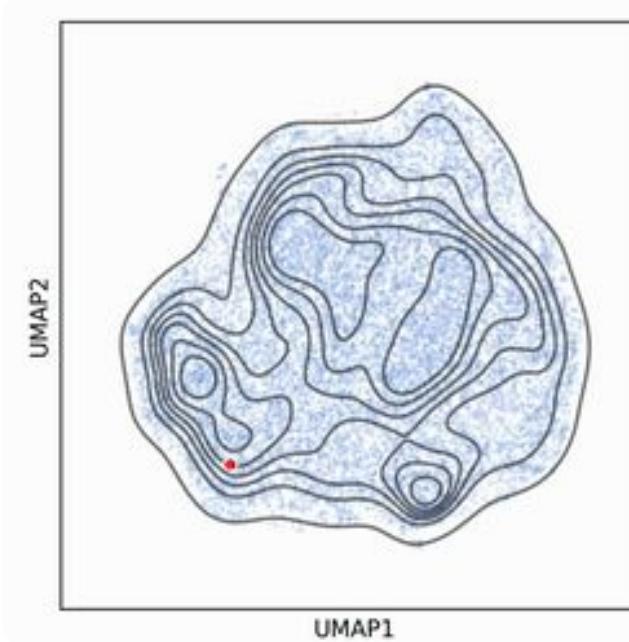
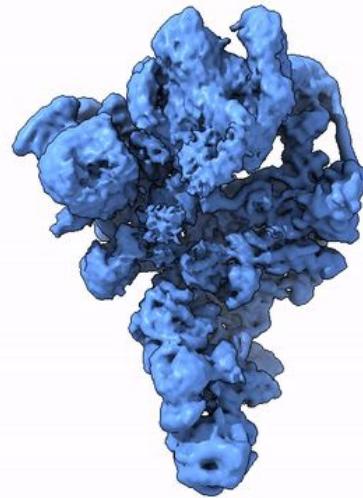
Trained on sequenced wastewater.

Dataset of 1.5 trillion base pairs.

Source: METAGENE-1: Metagenomic Foundation Model for Pandemic Monitoring

Generative Models – Many Applications

Applications include: **Science** – Cryo-EM for 3D Protein Structure Modeling



$p(3D \text{ structure} \mid \text{cryo-EM image})$

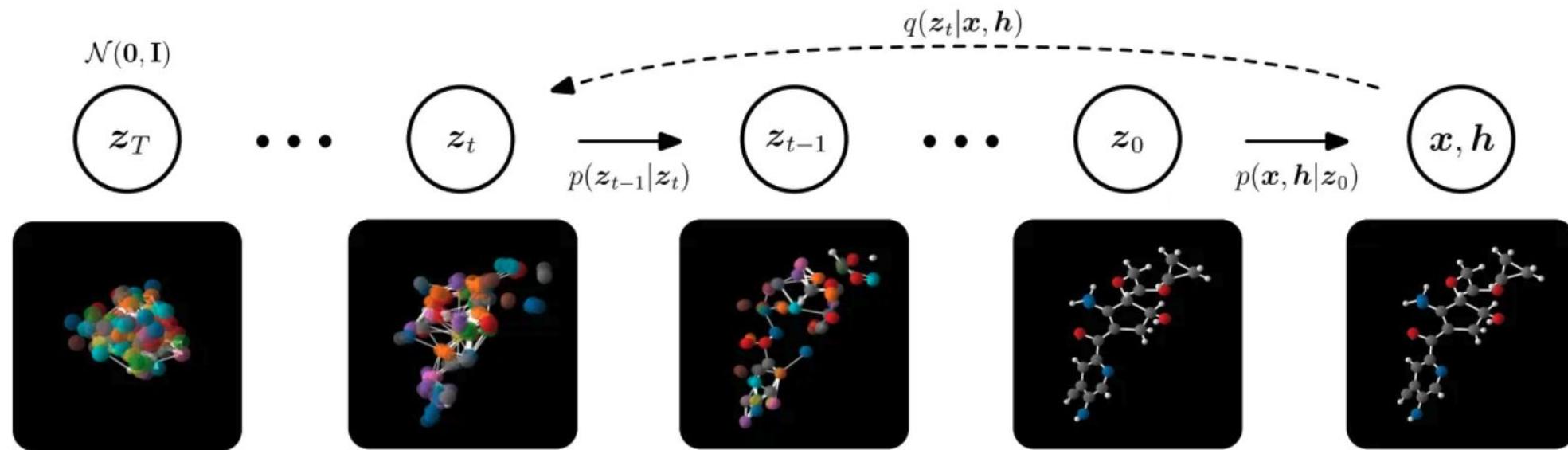
Sources:

CryoDRGN: reconstruction of heterogeneous cryo-EM structures using neural networks.

CryoDRGN-ET: deep reconstructing generative networks for visualizing dynamic biomolecules inside cells.

Generative Models – Many Applications

Applications include: Science – Molecular generation in chemistry



$p(\text{molecules})$

Sources: Equivariant Diffusion for Molecule Generation in 3D

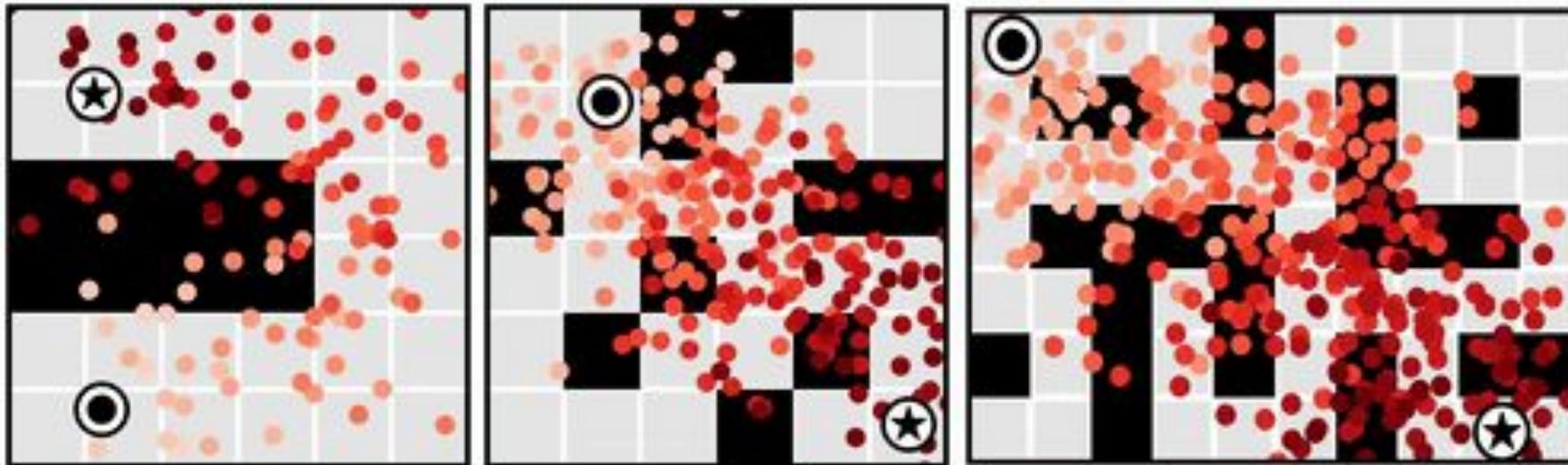
Generative Models – Many Applications

Applications include: **Actions & Planning**

Generative Models – Many Applications

Applications include: **Actions & Planning**

Diffusion for planning $p(\text{action sequence} \mid \text{environment map})$



Source: Planning with Diffusion for
Flexible Behavior Synthesis

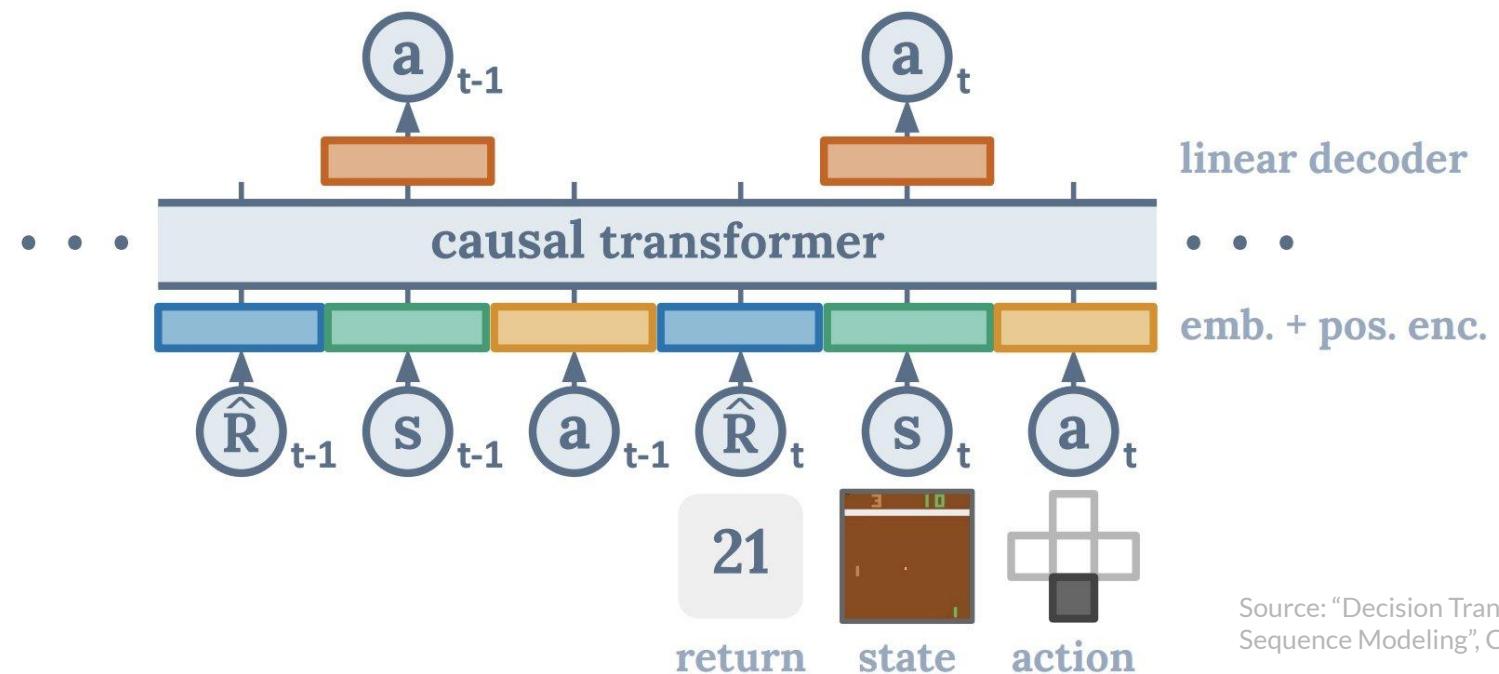
Generative Models – Many Applications

Applications include: **Decision Making**

Generative Models – Many Applications

Applications include: **Decision Making**

Decision Transformer – autoregressive model: $p(\text{states}, \text{actions}, \text{returns})$



Source: "Decision Transformer: Reinforcement Learning via Sequence Modeling", Chen et al., 2021

Generative Models – Many Applications

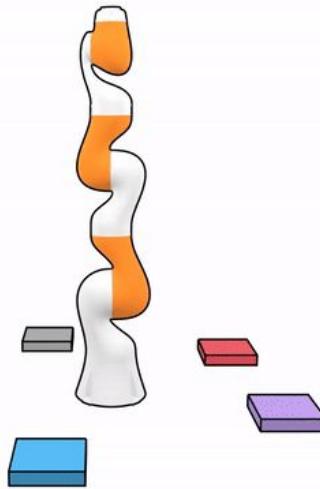
Applications include: **Robotics**

Generative Models – Many Applications

Applications include: **Robotics**

$p(\text{actions} \mid \text{images, text})$

Diffusion planning for robotics



Source: Planning with Diffusion for Flexible Behavior Synthesis.

Vision-language-action model



Source: Physical Intelligence π0.

Generative Models – Many Applications

Applications include: **Conditional Generative Models** – for other tasks...

Generative Models – Many Applications

Applications include: **Conditional Generative Models** – for other tasks...

$$p(\text{zebra images} \mid \text{horse images})$$



Source: Andrej Risteski

Generative Models – Many Applications

Applications include: **Style Transfer** – e.g., on images

Generative Models – Many Applications

Applications include: **Style Transfer** – e.g., on images

$p(\text{Monet painting} \mid \text{input image})$ $p(\text{Van Gogh painting} \mid \text{input image})$



Input Image



Monet



Van Gogh

Source: Andrej Risteski

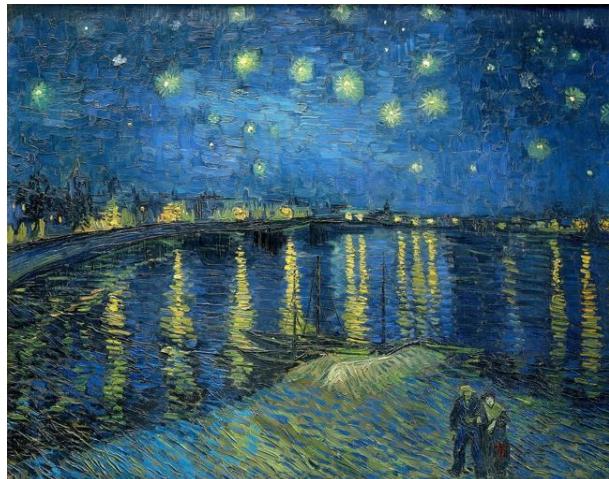
Generative Models – Many Applications

Applications include: **Actions for Art Generation**

Generative Models – Many Applications

Applications include: **Actions for Art Generation**

Inverse painting: $p(\text{paint strokes} \mid \text{painting})$



Source: Inverse painting project

SketchODE: $p(\text{pen strokes} \mid \text{text})$



Source: SketchODE project.

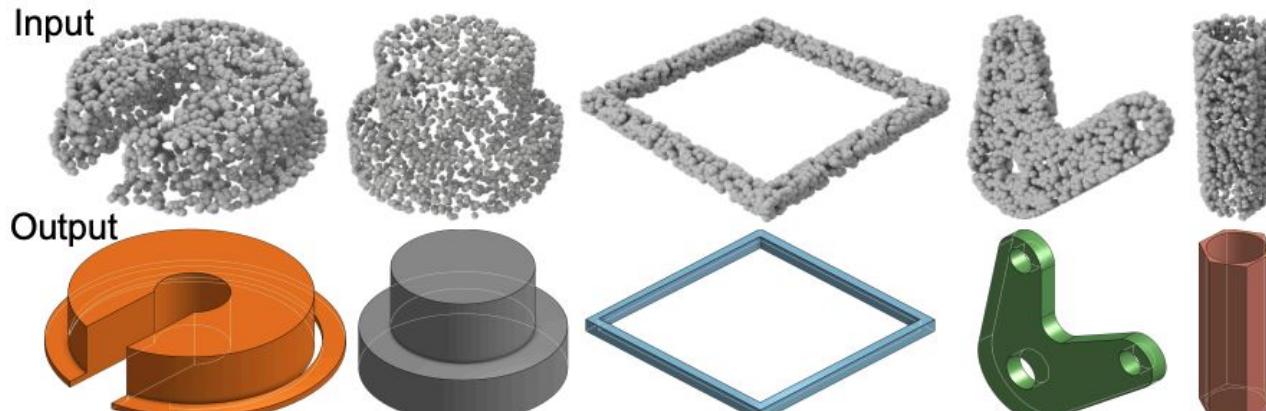
Generative Models – Many Applications

Applications include: **Generation for CAD (computer-aided design)**

Generative Models – Many Applications

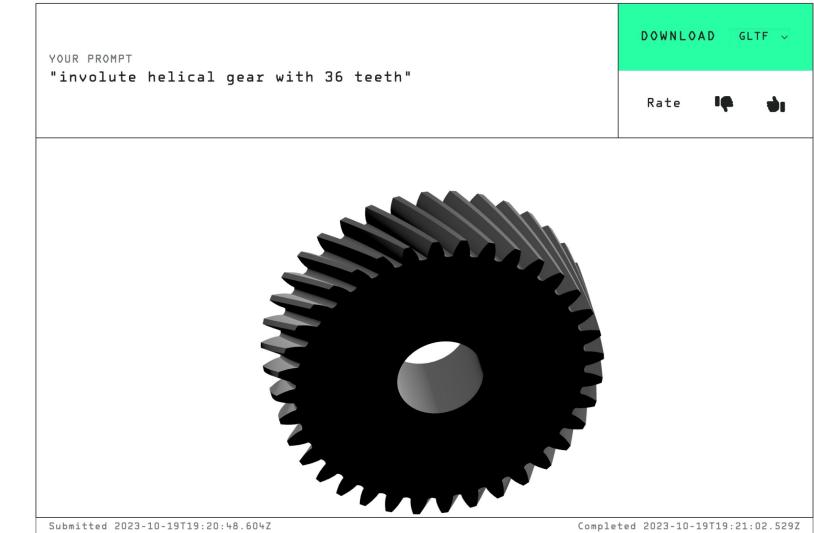
Applications include: **Generation for CAD (computer-aided design)**

Point-cloud to CAD



Source: DeepCAD: A Deep Generative Network for Computer-Aided Design Models

Text to CAD



Source: zoo.dev

Generative Models – Downstream Tasks

After learning a generative model, often want to carry out tasks such as:

Generative Models – Downstream Tasks

After learning a generative model, often want to carry out tasks such as:

1. **Sampling:** How can we generate novel data from the model distribution, i.e., $x_{\text{new}} \sim p_{\theta}(x)$?
2. **Density estimation:** Given a datapoint x , what is the probability density assigned by the model, $p_{\theta}(x)$.
3. **Unsupervised representation learning:** How can we learn meaningful feature representations for a datapoint x ? (And then use them!)

We will cover these topics in this class.

Instructor and Teaching Assistant



Willie Neiswanger
Instructor



Oliver Liu
Teaching Assistant

Instructor and Teaching Assistant



Oliver Liu
Teaching Assistant



Oliver Liu

(My PhD student).

Works on:

- Large multimodal models,
- scientific discovery,
- reasoning and decision-making.

About me – some background context



Willie Neiswanger
Instructor

About me – some background context

In PhD and postdoc: *Probabilistic modeling with application to science and engineering.*

About me – some background context

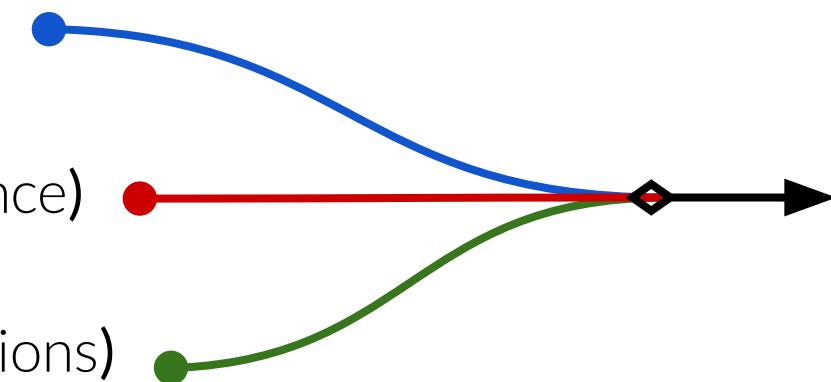
In PhD and postdoc: *Probabilistic modeling with application to science and engineering.*

Use techniques from:

Experimental Design (statistics)

Active Learning (computer science)

Bayesian Opt & Bandits (operations)



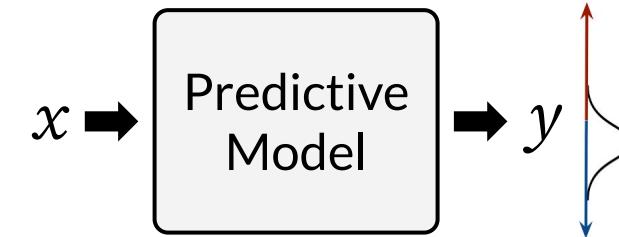
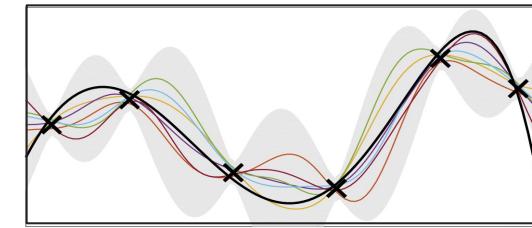
AI-Driven
Optimization and
Experimental Design

Focus on combining probabilistic machine learning with decision making.

About me – some background context

Probabilistic machine learning

- Classic probabilistic models
- Generative models



About me – some background context

Probabilistic machine learning

- Classic probabilistic models
- Generative models

Applications to science and engineering:

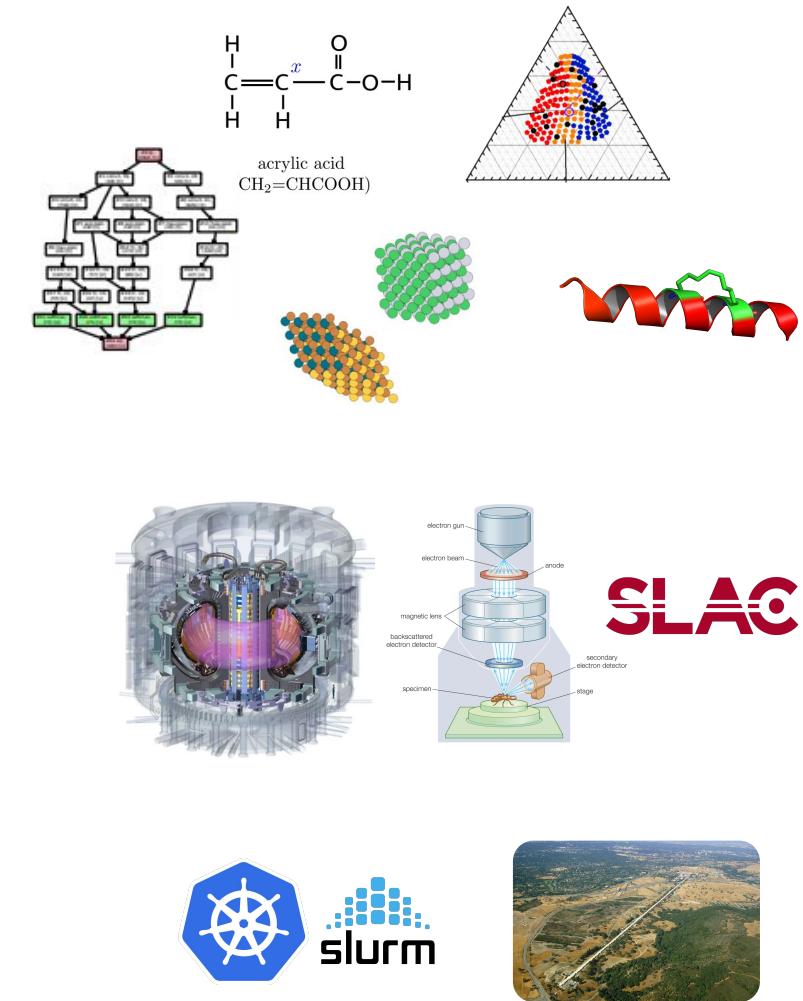
About me – some background context

Probabilistic machine learning

- Classic probabilistic models
- Generative models

Applications to science and engineering:

- Black-box optimization & experimental design:
 - Materials science (high-throughput screening).
 - Scientific machines (particle accelerators, tokamaks).
 - Computer/ML systems (config. tuning, hyperparameter opt.).



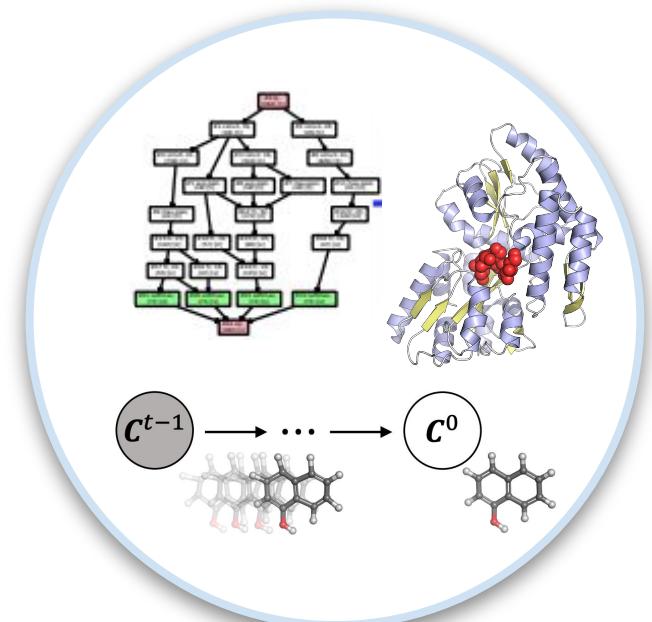
About me – some background context

Probabilistic machine learning

- Classic probabilistic models
- Generative models

Applications to science and engineering:

- Black-box optimization & experimental design:
 - Materials science (high-throughput screening).
 - Scientific machines (particle accelerators, tokamaks).
 - Computer/ML systems (config. tuning, hyperparameter opt.).
- Generative + probabilistic modeling:
 - Chemical / molecular design
 - Neural architecture design



At the end of my postdoc (~mid 2023)

Some collaborators began a project on open source reproductions of LLMs ...



At the end of my postdoc (~mid 2023)

Some collaborators began a project on open source reproductions of LLMs ...



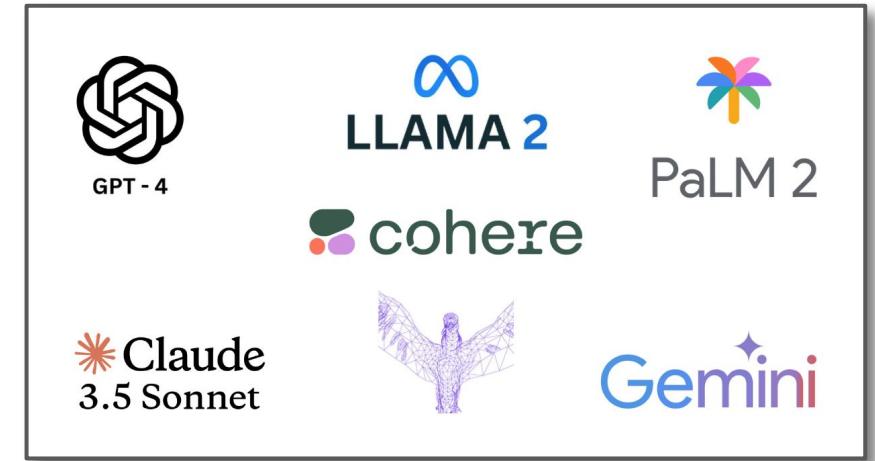
Through support from some universities/companies they got resources to do LLM pretraining.

- ⇒ I thought it would be a good opportunity to learn
- ⇒ Started working on large generative models.



At the end of my postdoc (~mid 2023)

Some collaborators began a project on open source reproductions of LLMs ...



Through support from some universities/companies they got resources to do LLM pretraining.

- ⇒ I thought it would be a good opportunity to learn
- ⇒ Started working on large generative models.



Increasingly working on LLM & large generative modeling projects for past ~2 years.

- Including LLM/foundation model pretraining, more recently diffusion/flow models.

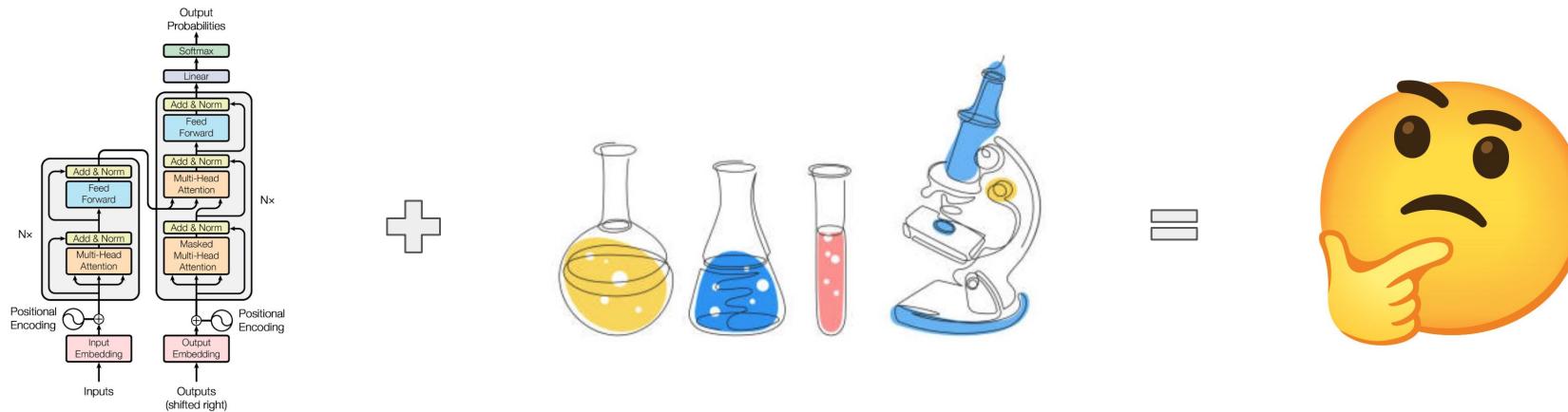
My current research

⇒ It influenced my research!

I still do probabilistic modeling and decision making for science, but...

I also combine this with large generative models, e.g.

- Generative/foundation models for scientific data (e.g., in biology, materials).
- LLMs for (sequential) decision making, and its variants.



Syllabus: Goals, Assignments, Grading, and More

Course Description & Learning Objectives

This course focuses on probabilistic models and their central role within modern machine learning and generative modeling.

We aim to cover topics including:

Course Description & Learning Objectives

This course focuses on probabilistic models and their central role within modern machine learning and generative modeling.

We aim to cover topics including:

- **Probabilistic graphical models** (e.g., Bayesian networks) and **approximate inference** algorithms (e.g., MCMC and variational inference).

Course Description & Learning Objectives

This course focuses on probabilistic models and their central role within modern machine learning and generative modeling.

We aim to cover topics including:

- **Probabilistic graphical models** (e.g., Bayesian networks) and **approximate inference** algorithms (e.g., MCMC and variational inference).
- **Deep generative models** (e.g., autoregressive, score-matching, diffusion, and flow-based approaches).

Course Description & Learning Objectives

This course focuses on probabilistic models and their central role within modern machine learning and generative modeling.

We aim to cover topics including:

- **Probabilistic graphical models** (e.g., Bayesian networks) and **approximate inference** algorithms (e.g., MCMC and variational inference).
- **Deep generative models** (e.g., autoregressive, score-matching, diffusion, and flow-based approaches).
- **Uncertainty quantification** in supervised/deep learning (Gaussian processes, Bayesian neural networks, ensembles, Monte Carlo dropout, UQ in LLMs)

Course Description & Learning Objectives

This course focuses on probabilistic models and their central role within modern machine learning and generative modeling.

We aim to cover topics including:

- **Probabilistic graphical models** (e.g., Bayesian networks) and **approximate inference** algorithms (e.g., MCMC and variational inference).
- **Deep generative models** (e.g., autoregressive, score-matching, diffusion, and flow-based approaches).
- **Uncertainty quantification** in supervised/deep learning (Gaussian processes, Bayesian neural networks, ensembles, Monte Carlo dropout, UQ in LLMs)
- Methods for probabilistic **model-based sequential decision-making** (e.g., Bayesian optimization, and information-based experimental design)

Preparation and Prerequisites

This course is designed for students currently pursuing research, or who wish to pursue research, in probabilistic machine learning or deep generative models!

- A major component will be carrying out a project relevant to your research interests.
- Also: reading and presenting research papers on probabilistic/generative models.

Preparation and Prerequisites

This course is designed for students currently pursuing research, or who wish to pursue research, in probabilistic machine learning or deep generative models!

- A major component will be carrying out a project relevant to your research interests.
- Also: reading and presenting research papers on probabilistic/generative models.

Students are expected to be comfortable with reading and presenting machine learning conference papers (e.g., NeurIPS, ICLR, ICML).

Beneficial to have familiarity with machine learning, algorithms, and probability.
However, there are no formal (required) class prerequisites.

Typical Class Structure

First half (1.5 - 2 hours) – Lecture from me on course materials.

- Core lesson on the topic of the day.
- [Periodically] “Lab session”: implementation or programming lesson related to topic.
- Also aim to review some classic papers on the topic (+ give some history of the field).

Typical Class Structure

Second half – Activities related to paper presentations, course project, and more.

- Student paper presentations on probabilistic/generative ML papers.
- Small-group meetings for feedback on projects.
- Project pitches, and project final presentations.

Laptops in Class – Recommended

You are welcome to use a laptop in this class (though not required), e.g., for

- Taking notes during lectures.
- (If interested) following along with implementation/lab session.
- Sometimes will be useful for roles during paper presentations (e.g., scribe).
- Sometimes will be useful for other tasks during course projects (e.g., sign-up for presentation schedule).

Class Website

Course website URL:

willieneis.github.io/probgen-spring2025

Please check it out for:

- Class news and updates.
- Lecture schedule & topics.
- Reminders of assignments, grading policies.

CSCI 699: Probabilistic and Generative Models, Spring 2025

Spring 2025, Fridays at 2:00-5:20pm in KAP 144

Instructor: Willie Neiswanger

This course focuses on probabilistic models and their central role within modern machine learning and generative modeling. With probabilistic methods increasingly driving advancements in AI, this course will explore its applications across a range of topics, including approximate inference algorithms (MCMC, variational inference), deep generative models (autoregressive, score-matching, diffusion, and flow-based models), and model-based sequential decision making.

Course Staff



Willie Neiswanger

Instructor

Office hours: Friday 5-6:30pm

Location: TBD



Oliver Liu

Teaching Assistant

Office hours: Day-Time and Day-Time

Location: TBD

Logistics

- **Assignments:** Submit all written assignments, including all project-related write-ups, on Brightspace. Grades and feedback will also be provided on Brightspace.
- **General discussion:** Please use the official course Slack channel for general questions.
- **Other discussion:** Email Willie and Oliver (neiswang@ and zliu2898@) or come to office hours to discuss individual matters, such as project ideas or grading.

Prerequisites

This course is designed for students currently pursuing research, or who wish to pursue research, in probabilistic machine learning or deep generative models. It will be beneficial to have familiarity with machine learning (at the level of CSCI 567), algorithms (at the level of CSCI 570), and probability (at the level of MATH 505a). Students are expected to be comfortable with reading and presenting modern machine learning conference papers.

Schedule

Course Submission Site

We'll be using Brightspace - learning management system.

URL: <https://brightspace.usc.edu/>

Mainly used for

- Submitting assignments (reports for class project, scribe notes, etc),
- Internal course content, documents, etc.
- In general: anything internal to the class that we do not intend to (immediately) make public.

The screenshot shows the USC Viterbi School of Engineering Brightspace course page for CSCI 699: Probabilistic and Generative Models, Spring 2025. The top navigation bar includes links for Home, Announcements, Content, Activities, My Grades, Help, Course Tools, and Library Resources. A banner image of water droplets on leaves is visible. Below the banner, there's a 'Slim Announcements Widget' showing a single announcement about the course's private nature. To the right, there's a 'Calendar' section showing an upcoming event for Thursday, January 16, 2025, and an 'Activity Feed' section with a placeholder for creating a post.

Communications with Instructor/TA

For individual matters, send us an **email**:
neiswang@usc.edu, zliu2898@usc.edu

- (You are also welcome to direct message us on **Slack** via your USC Slack account)

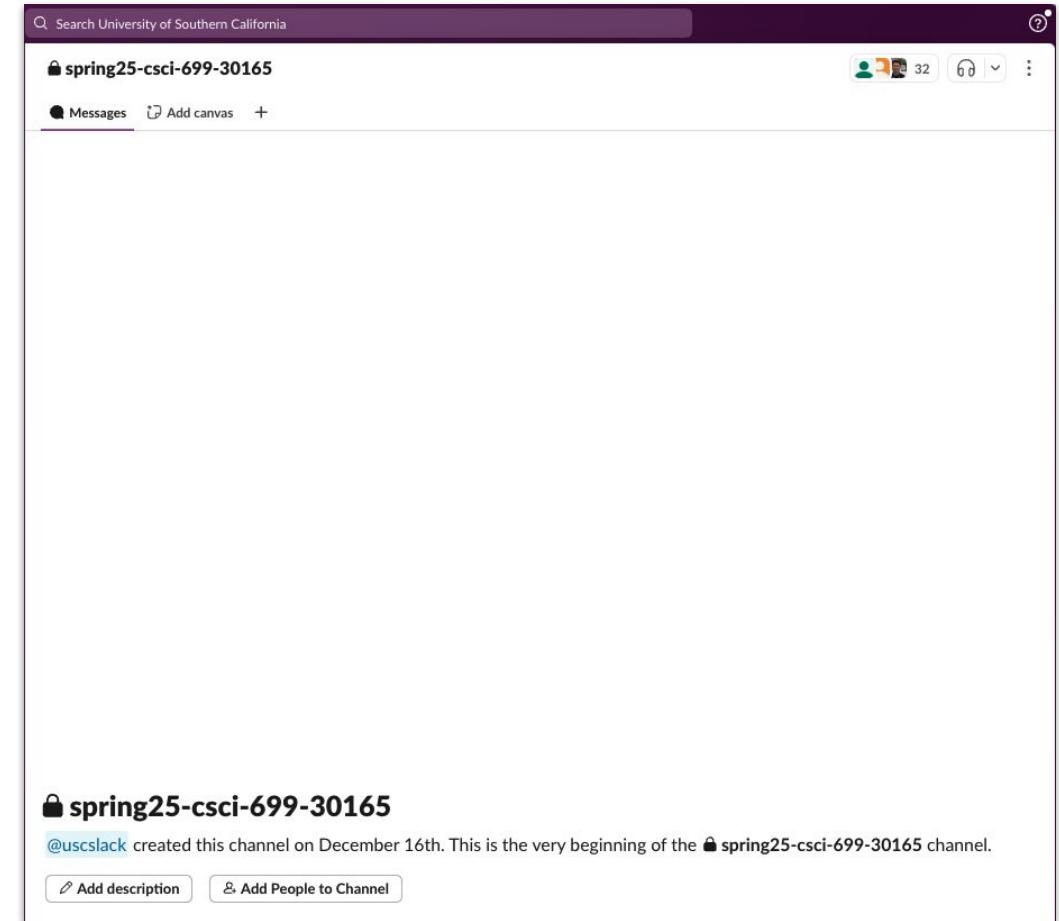
Communications with Instructor/TA

For individual matters, send us an **email**:
neiswang@usc.edu, zliu2898@usc.edu

- (You are also welcome to direct message us on **Slack** via your USC Slack account)

For general discussion, or broader questions relevant to the class, there is a Slack channel that can be used:
spring25-csci-699-30165

- You should be already added here – please let me know if you aren't!



Assignments and Grading – Grading Breakdown

Assignments and Grading – Grading Breakdown

<u>Assignment</u>	<u>% of Grade</u>
1. Paper Presentation	20%
2. In-class Participation and Discussion	
2a. Role 1 – Discussion Lead 1	8%
2a. Role 2 – Discussion Lead 2	8%
2a. Role 3 – Scribe	9%
3. Course Project	
3a. Project Pitch	8%
3b. Midway Report	10%
3c. Final Presentation	12%
3d. Final Report	25%

Paper Presentations

Paper Presentations – Goal

- During the semester, each student gives **one 20-minute presentation** on a paper relevant to probabilistic and generative models.
- Ideally from a modern machine learning conference (e.g., NeurIPS, ICML, ICLR, AAAI, etc)
 - Could be a “classic” (older) paper relevant to modern models, as well.
- This will accomplish a few things:
 - Gives the class a chance to see a broad set of interesting probabilistic/generative modeling papers.
 - Gives each student (more) experience with distilling key content from a paper & presenting it.

Paper Presentations – Which papers to present?

There are a few options:

Paper Presentations – Which papers to present?

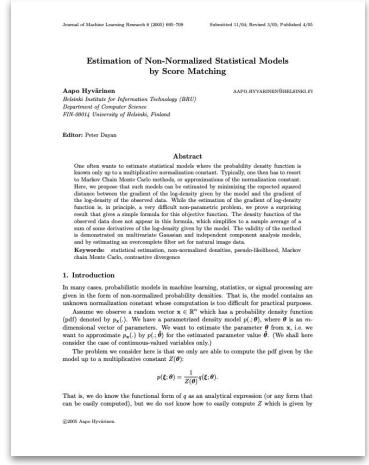
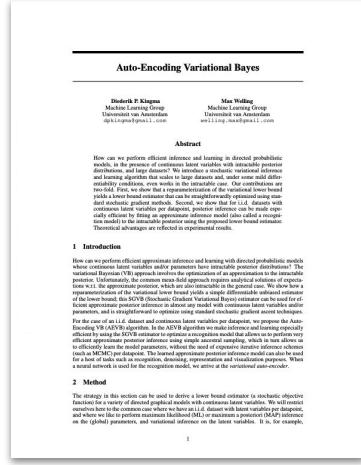
There are a few options:

- You are encouraged to choose a relevant paper to present, ideally which ties into your own research interests.
 - We will have a spreadsheet where you will list the paper you have chosen, so instructors can confirm, and also so that we don't have overlapping papers.

Paper Presentations – Which papers to present?

There are a few options:

- You are encouraged to choose a relevant paper to present, ideally which ties into your own research interests.
 - We will have a spreadsheet where you will list the paper you have chosen, so instructors can confirm, and also so that we don't have overlapping papers.
- I will share a list of interesting probabilistic/generative modeling papers for inspiration— you are welcome to present on one of these.



Paper Presentations – Content of each paper presentation?

Ideally I'd like you to include:

Paper Presentations – Content of each paper presentation?

Ideally I'd like you to include:

1. Motivation (main problem the paper is trying to solve).
2. Brief background on prior work (*i.e.*, related methods, relevant papers).
3. Method and key equations.
 - It'll be important to boil the method down to a concise explanation given ~20 minute presentation.
4. Empirical results – how well does the method work in practice?
5. Your opinion on the paper, its (potential) impact/promise, etc.

Paper Presentations – Schedule

Paper Presentations – Schedule

Starting at the **end of February** we will go through a few (~4) presentations per class (2nd half of class).

Paper Presentations – Schedule

Starting at the **end of February** we will go through a few (~4) presentations per class (2nd half of class).

In the spreadsheet, which I will share, students will sign up for a time slot during the semester.

To make it fair, students who sign up for presentations on the first two dates will get slightly more-lenient grading (a point of extra credit on this assignment).

A	B	C	D	E	F	G	H
Date	Presentation ID	Presenter Name (sign up!)	Paper Title	Link/url to paper	Discussion Lead 1 (sign up!)	Discussion Lead 2 (sign up!)	Scribe (sign up!)
February 28							
	1						
	2						
	3						
	4						
March 7							
	5						

In-Class Participation and Discussion

There is a grade for in-class participation/discussion (during student presentations).
Students will sign up for one of three roles (again in the spreadsheet).
Each student has to do each role once during semester.

In-Class Participation and Discussion

There is a grade for in-class participation/discussion (during student presentations). Students will sign up for one of three roles (again in the spreadsheet). Each student has to do each role once during semester.

Roles 1 and 2 – Discussion Leads

- Read the paper before the class; Responsible for formulating a short list (~5 questions) for discussion; bring up a couple of these during class; submit all after class.

Role 3 – Scribe

- Read the paper before the class; take notes on paper during presentation; write up a ~1 page summary of the presentation.

Course Project

Course Project – Overview & Goals

Another important part of this class will be a semester-long course project.

Aims to give hands-on experience in probabilistic & generative models.

I'd like you to choose a topic that is relevant to your personal research interests (or a topic you've wanted the chance to learn more about).

- So that you can use the methods from this class in your research.

Course Project – Group Project

This will be a **group project** – groups of 3-4 students.

- Aiming for ~10 groups total (due to timing constraints)
- We will help facilitate this during class.
- E.g., second half of next class everyone will introduce themselves and describe research interests, which we will write/share, to help in matching.
- Need to aim to form teams and select project idea by roughly end of this month.

Course Project – Guidance & Expectations

What does this project entail?

Course Project – Guidance & Expectations

What does this project entail?

I want people to use a probabilistic or generative model in some way!

- Application of prob/gen models from this class on a novel task or dataset.
- Algorithmic improvements in learning, inference, or evaluation of prob/gen models.
- Theoretical analysis of any aspect of existing prob/gen models.

Course Project – Guidance & Expectations

What does this project entail?

I want people to use a probabilistic or generative model in some way!

- Application of prob/gen models from this class on a novel task or dataset.
- Algorithmic improvements in learning, inference, or evaluation of prob/gen models.
- Theoretical analysis of any aspect of existing prob/gen models.

Goal is to complete a small-scale implementation or pilot study during the class.
I'm more focused on interesting conceptual ideas, rather than on performance/results.

Aim to connect it to the research you are focusing on outside of this class!

Course Project – Guidance & Expectations

The project **does not** need to be at the level of submission to a major machine learning conference.

(But it's great if works out to be the initial seed of a paper that gets submitted – you can let me know if you are interested in help on this).

Course Project – Guidance & Expectations

The project **does not** need to be at the level of submission to a major machine learning conference.

(But it's great if works out to be the initial seed of a paper that gets submitted – you can let me know if you are interested in help on this).

Possible topics:

PGMs, MCMC/VI, deep generative models (VAE, diffusion, flow matching), autoregressive models/LLMs, deep uncertainty quantification, probabilistic predictive model, decision making (BayesOpt).

Possible application areas:

NLP/language, computer vision/image/video, tabular data, discrete data, scientific data (biology, chemistry, physics), connect to other areas of CS? Other fields? Art?

Course Project – Guidance & Expectations

This is a chance to be creative with projects!

Try something creative, perhaps try to model a new data type.

Take a risk, see if you can make something work.

If you put in effort/give an honest shot, we will not grade harshly if the project doesn't work out.

Goal is to just to complete a small-scale implementation or pilot student during the class.



Source: [sketchode project](#)

Course Project – Assignments

The project will be worth a substantial portion of the grade, and consist of four main assignments:

Course Project – Assignments

The project will be worth a substantial portion of the grade, and consist of four main assignments:

- **Project pitch:** Each group will come up with a project idea, make a few slides, and share their idea with the class for feedback (**10 minutes long, on Feb 7 & Feb 14**)

Course Project – Assignments

The project will be worth a substantial portion of the grade, and consist of four main assignments:

- **Project pitch:** Each group will come up with a project idea, make a few slides, and share their idea with the class for feedback (**10 minutes long, on Feb 7 & Feb 14**)
- **Midway report:** Each group will write a short (4 page) midway report for their project, focusing on a literature review, implementation plan, and any initial experiments.
 - Latex template will be provided.

Course Project – Assignments

The project will be worth a substantial portion of the grade, and consist of four main assignments:

- **Project pitch:** Each group will come up with a project idea, make a few slides, and share their idea with the class for feedback (**10 minutes long, on Feb 7 & Feb 14**)
- **Midway report:** Each group will write a short (4 page) midway report for their project, focusing on a literature review, implementation plan, and any initial experiments.
 - Latex template will be provided.
- **Final presentation:** Each group will give a presentation to the class on their final project (**30 minutes long, on Apr 25 & May 2**)

Course Project – Assignments

The project will be worth a substantial portion of the grade, and consist of four main assignments:

- **Project pitch:** Each group will come up with a project idea, make a few slides, and share their idea with the class for feedback (**10 minutes long, on Feb 7 & Feb 14**)
- **Midway report:** Each group will write a short (4 page) midway report for their project, focusing on a literature review, implementation plan, and any initial experiments.
 - Latex template will be provided.
- **Final presentation:** Each group will give a presentation to the class on their final project (**30 minutes long, on Apr 25 & May 2**)
- **Final report:** Each group will submit a final report for their project, describing all details, background, prior work, and results (**8-10 pages long, due May 9**).
 - Latex template will be provided.

Course Project – Good Course Project Examples

From some previous PGM and DGM courses:

Stanford CS236 Deep Generative Models, course project topics and titles:

- [link, more recent](#)

CMU 10-708 Probabilistic Graphical Models, course projects and reports:

- [link](#)

Course Project – Compute Resources

Depending on the type of model you choose to focus on, external compute resources (in particular, GPU resources), may be useful

Keeping in mind we are aiming for a small-scale/pilot implementation for these projects!

Course Project – Compute Resources

Depending on the type of model you choose to focus on, external compute resources (in particular, GPU resources), may be useful

Keeping in mind we are aiming for a small-scale/pilot implementation for these projects!

CARC Allocation:

- I've requested a CARC allocation for this project, will give access to some GPU resources for training models.

Google Colab

- Free resources for students to run machine learning code.
- Implement training in notebook. Examples (e.g., for diffusion models):
 - [link \(HF diffusers\)](#), [link \(diffusion from scratch\)](#), [link \(simple diffusion\)](#)

LightningAI Studio Student Account

- Some GPU resources for students, [link](#)

Assignments and Grading – To Summarize...

<u>Assignment</u>	<u>% of Grade</u>
1. Paper Presentation (<i>Individual</i>)	20%
2. In-class Participation and Discussion (<i>Individual</i>)	
2a. Role 1 – Discussion Lead 1	8%
2a. Role 2 – Discussion Lead 2	8%
2a. Role 3 – Scribe	9%
3. Course Project (<i>Group</i>)	
3a. Project Pitch	8%
3b. Midway Report	10%
3c. Final Presentation	12%
3d. Final Report	25%

Academic Integrity

This course will follow the expectations for academic integrity as stated in the [USC Student Handbook](#).

This is primarily a paper presentation & (group) project class, so less risk of any academic integrity violations.

Please **ask the instructor/TA** if you are unsure about what constitutes unauthorized assistance, or what information requires citation and/or attribution!

Use of Generative AI

... e.g., in your paper presentations, group project (midway report, final report, etc).

Use of Generative AI

... e.g., in your paper presentations, group project (midway report, final report, etc).

Short answer, **yes** you can use it – this is a generative models class, after all.

And learning to use AI in an effective way is an important emerging skill.

Use of Generative AI

... e.g., in your paper presentations, group project (midway report, final report, etc).

Short answer, **yes** you can use it – this is a generative models class, after all.

And learning to use AI in an effective way is an important emerging skill.

However, keep in mind the following caveats:

- AI tools are permitted to help you brainstorm topics or revise work you have already written.
- If you provide minimum-effort prompts, you will get low-quality results. You will need to refine your prompts to get good outcomes. This will take work.
- Proceed with caution when using AI, don't assume the information provided is accurate.
- You need to **acknowledge that you used AI**. Please include sentences at the end of any assignment that used AI, explaining how and why you used AI.

Optional Readings – and Supplementary Materials

The following resources are useful for many of the topics covered in this class:

1. Kevin Murphy, “Machine Learning: A Probabilistic Perspective”, 2012 ([link](#)).
2. Kevin Murphy, “Probabilistic Machine Learning: Advanced Topics”, 2023 ([link](#)).
3. Chris Bishop, “Pattern Recognition and Machine Learning”, 2006 ([link](#)).
4. Chris Bishop, “Deep Learning - Foundations and Concepts”, 2024 ([link](#)).
5. Stefano Ermon, Deep Generative Models, Course Notes ([link](#)).

A Tour of the Class

A Tour of the Class

Want to roughly go through the topics we'll learn in this class.

And the schedule* over the coming semester.

*subject to change or adjustment if needed 😎.

Lecture 1 (Today, January 17) – Class Overview and Probability Review

Lecture 1 (Today, January 17) – Class Overview and Probability Review

We are here  currently going through an overview of the class.

In a few minutes we'll do a probability review of terminology/notation useful for various models & methods we'll learn in this class.

Probability Review – Probability Space

Definition. Probability Space (Ω, \mathcal{F}, P)

Event Space $\mathcal{F} \subseteq 2^\Omega$

- A set whose elements A ("events") are subsets of Ω .
- Satisfies three properties:
 - contains empty set
 - closed under complements
 - closed under countable unions

An example event:
The coin lands on heads and the die lands on an odd number.

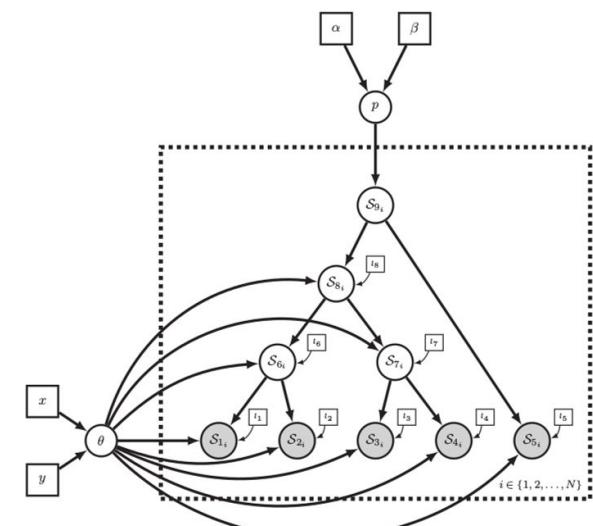
$$\begin{aligned}\emptyset &\in \mathcal{F} \\ A \in \mathcal{F} &\implies \Omega \setminus A \in \mathcal{F} \\ A_1, A_2, \dots \in \mathcal{F} &\implies \cup_i A_i \in \mathcal{F}\end{aligned}$$

Lecture 2 (January 24) – Probabilistic Graphical Models (PGMs)

Lecture 2 (January 24) – Probabilistic Graphical Models (PGMs)

Lecture: PGMs – using ideas from both probability and graph theory to derive efficient machine learning algorithms

- Directed vs undirected graphical models. Bayes networks vs Markov random fields.
- Plate notation, generative process notation.
- Inference vs learning in graphical models.
- Famous/classic PGMs (HMMs, GMMs, LDA, VAE, etc.)
- Some probabilistic programming.



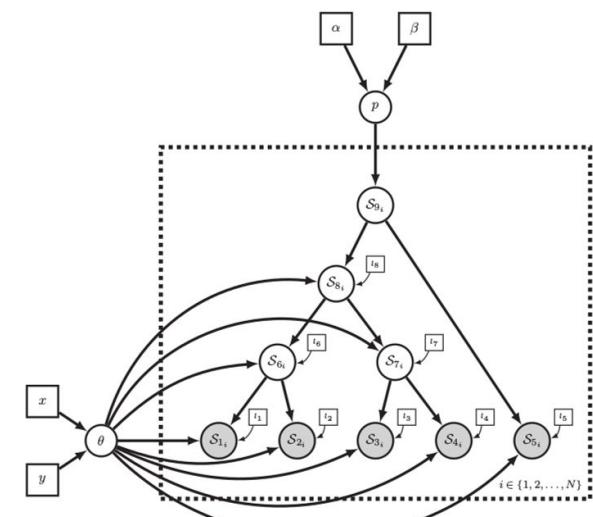
Lecture 2 (January 24) – Probabilistic Graphical Models (PGMs)

Lecture: PGMs – using ideas from both probability and graph theory to derive efficient machine learning algorithms

- Directed vs undirected graphical models. Bayes networks vs Markov random fields.
- Plate notation, generative process notation.
- Inference vs learning in graphical models.
- Famous/classic PGMs (HMMs, GMMs, LDA, VAE, etc.)
- Some probabilistic programming.

After:

- Class introductions and research interests, set up for group project.
- Compute resources.



Source: "Probabilistic Graphical Model Representation in Phylogenetics", Höhna 2014

Lecture 3 (January 31) – PGMs: Classic Algorithms

Lecture 3 (January 31) – PGMs: Classic Algorithms

Lecture: Classic algorithms in probabilistic graphical models (PGMs) for exact and approximate inference & learning.

- Forward-backward algorithm in HMMs.
- Viterbi algorithm in HMMs.
- Exact inference and variable elimination algorithm.
- Belief Propagation (message passing) algorithm.
- Expectation-maximization (EM) algorithm.



Source: USC Viterbi Magazine, "The Viterbi Algorithm at 50"

Lecture 3 (January 31) – PGMs: Classic Algorithms

Lecture: Classic algorithms in probabilistic graphical models (PGMs) for exact and approximate inference & learning.

- Forward-backward algorithm in HMMs.
- Viterbi algorithm in HMMs.
- Exact inference and variable elimination algorithm.
- Belief Propagation (message passing) algorithm.
- Expectation-maximization (EM) algorithm.

After:

- Lab session on probabilistic programming (Pyro, Stan, NumPyro, GPyTorch, etc).
- Check-in on group formation and project ideation.



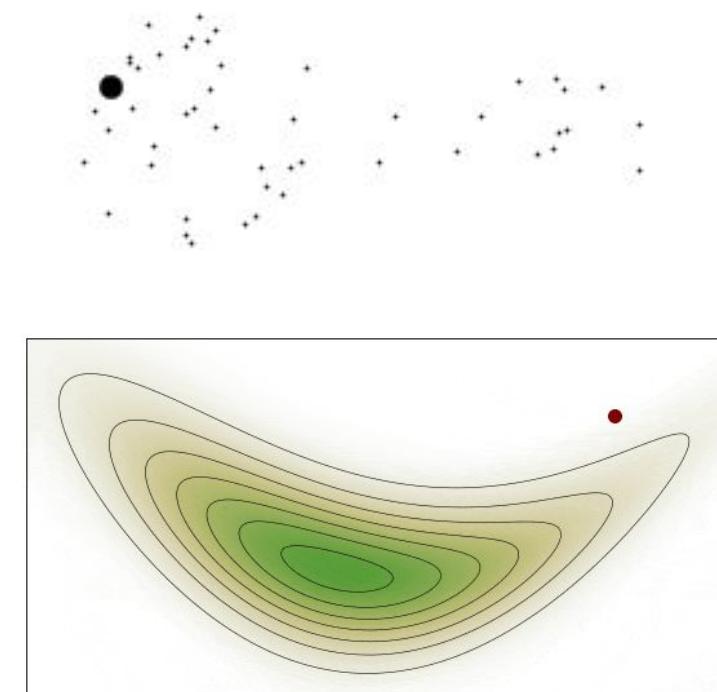
Source: USC Viterbi Magazine, "The Viterbi Algorithm at 50"

Lecture 4 (February 7) – Approximate Bayesian Inference

Lecture 4 (February 7) – Approximate Bayesian Inference

Lecture: Approximate inference algorithms in probabilistic models, including the classics – Markov chain Monte Carlo (MCMC), and variational inference (VI).

- Markov chain Monte Carlo (MCMC) methods, including:
 - Metropolis–Hastings algorithm, Gibbs sampling, slice sampling.
 - Gradient-based methods: Langevin Monte Carlo, Hamiltonian Monte Carlo.
- Variational inference (VI) methods, including
 - Evidence lower bound (ELBO), Stochastic VI, black-box VI.



Source: "Creating animations with MCMC", Krepl 2018,
Wikimedia commons,

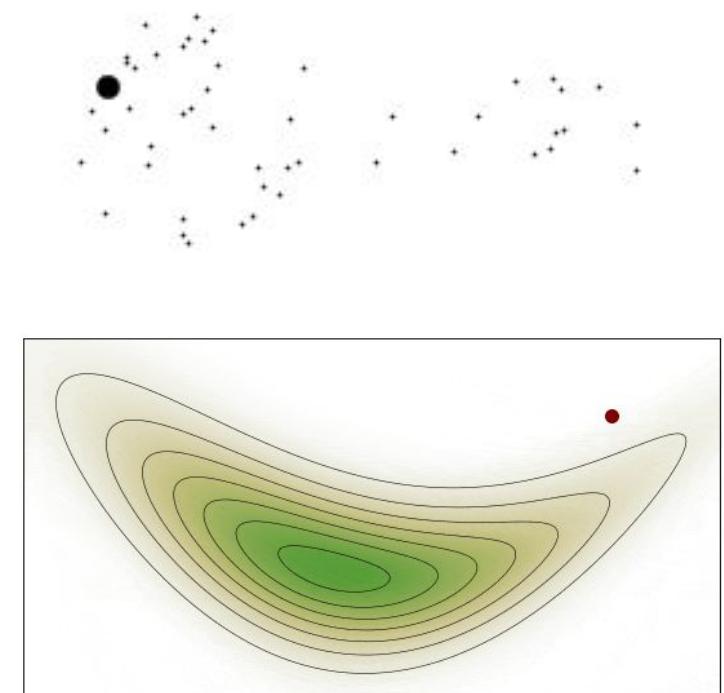
Lecture 4 (February 7) – Approximate Bayesian Inference

Lecture: Approximate inference algorithms in probabilistic models, including the classics – Markov chain Monte Carlo (MCMC), and variational inference (VI).

- Markov chain Monte Carlo (MCMC) methods, including:
 - Metropolis–Hastings algorithm, Gibbs sampling, slice sampling.
 - Gradient-based methods: Langevin Monte Carlo, Hamiltonian Monte Carlo.
- Variational inference (VI) methods, including
 - Evidence lower bound (ELBO), Stochastic VI, black-box VI.

After:

- Project Pitches #1: Groups 1-5



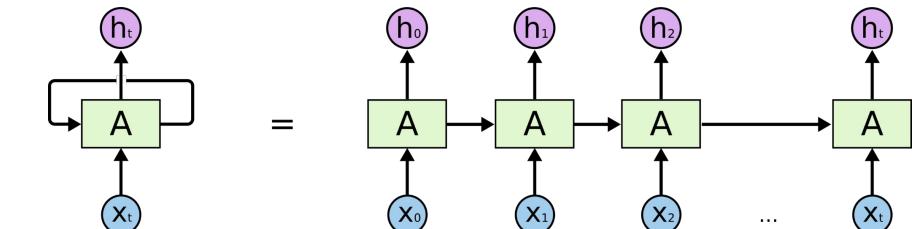
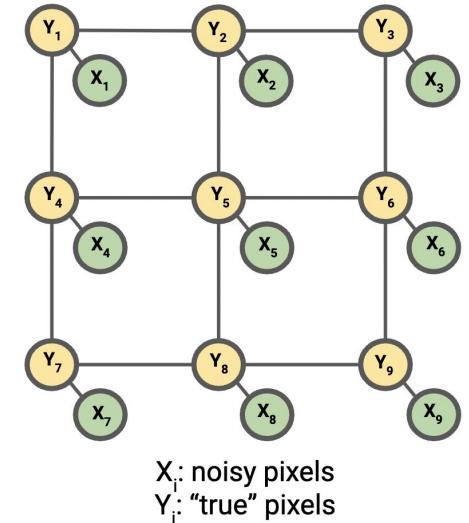
Source: "Creating animations with MCMC", Krepl 2018,
Wikimedia commons,

Lecture 5 (February 14) – Intro to Deep Generative Models

Lecture 5 (February 14) – Intro to Deep Generative Models

Lecture: Introduction to deep (probabilistic) generative models, model paradigms, learning in these models, optimal representations and objective functions.

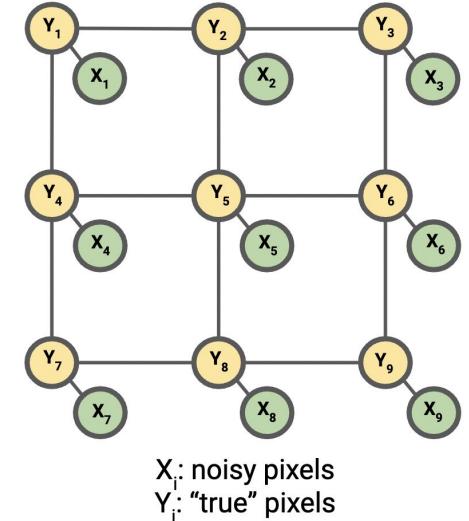
- Modeling paradigms: autoregressive vs diffusion vs flow models.
- Maximum-likelihood / Kullback-Leibler divergence-based Learning.
- Energy-based models
- Early autoregressive models.



Lecture 5 (February 14) – Intro to Deep Generative Models

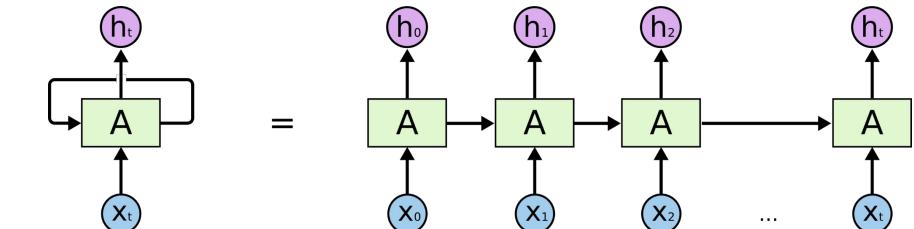
Lecture: Introduction to deep (probabilistic) generative models, model paradigms, learning in these models, optimal representations and objective functions.

- Modeling paradigms: autoregressive vs diffusion vs flow models.
- Maximum-likelihood / Kullback-Leibler divergence-based Learning.
- Energy-based models
- Early autoregressive models.



After:

- Project Pitches #2: Groups 6-10

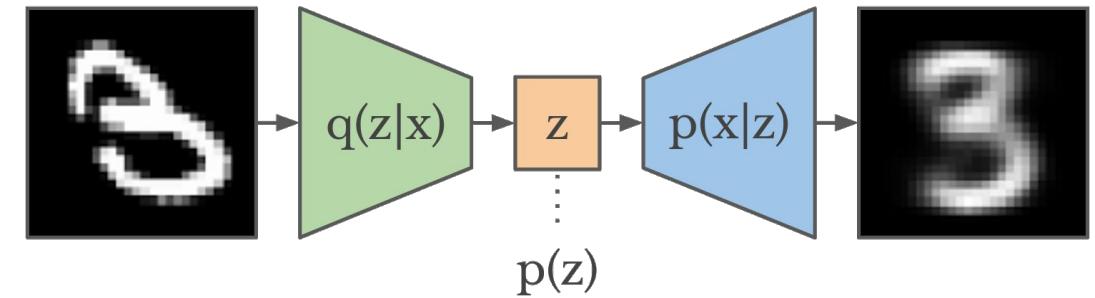


Lecture 6 (February 21) – VAEs and GANs

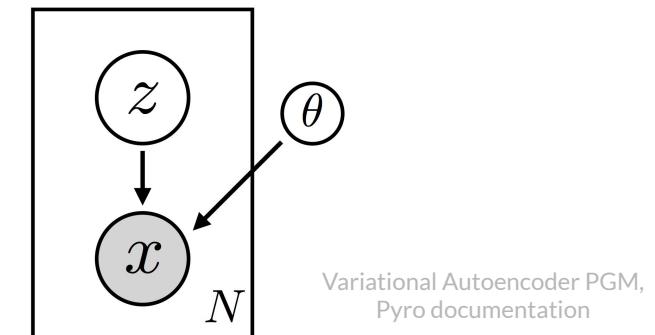
Lecture 6 (February 21) – VAEs and GANs

Lecture: Two historically iconic deep generative models – variational autoencoders (VAEs) and generative adversarial networks (GANs).

- Variational autoencoders
 - Relation to other graphical models.
 - Relation to (classic) autoencoders.
- Generative adversarial networks
 - Extensions in recent times.



"Building Variational Auto-Encoders in TensorFlow", Danijar Hafner

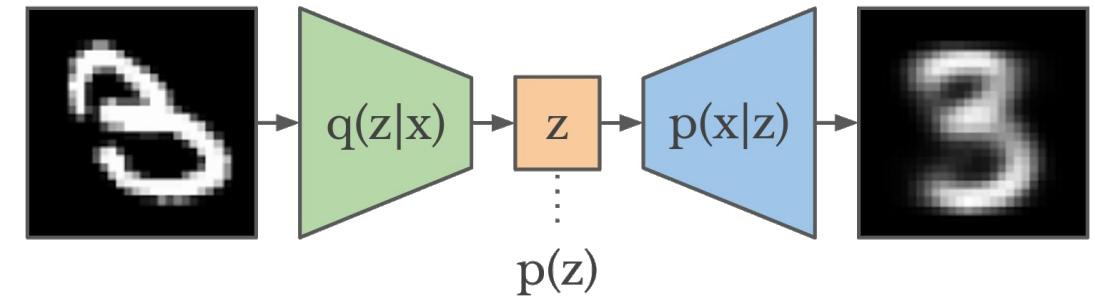


Variational Autoencoder PGM,
Pyro documentation

Lecture 6 (February 21) – VAEs and GANs

Lecture: Two historically iconic deep generative models – variational autoencoders (VAEs) and generative adversarial networks (GANs).

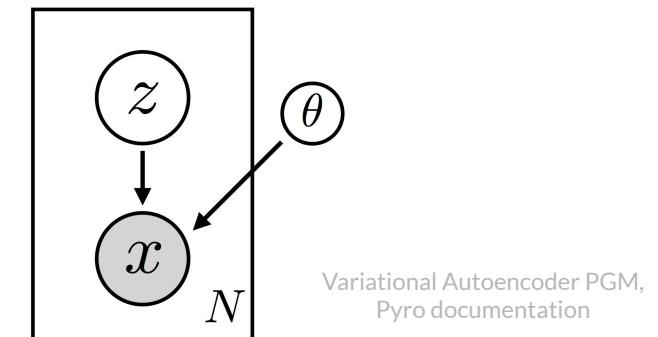
- Variational autoencoders
 - Relation to other graphical models.
 - Relation to (classic) autoencoders.
- Generative adversarial networks
 - Extensions in recent times.



"Building Variational Auto-Encoders in TensorFlow", Danijar Hafner

After:

- Project feedback/discussion via short in-class meetings with each group.



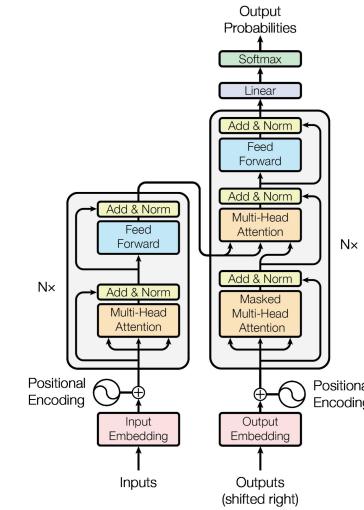
Variational Autoencoder PGM,
Pyro documentation

Lecture 7 (February 28) – Autoregressive Models

Lecture 7 (February 28) – Autoregressive Models

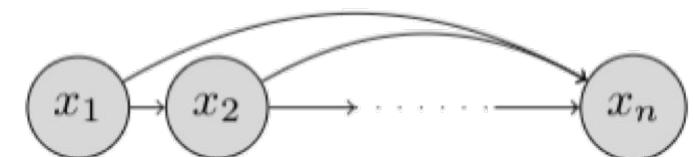
Lecture: Autoregressive generative models, the transformer neural architecture, and its use in large language modeling.

- Neural autoregressive density estimation
- Recurrent neural networks
- Attention-based models
- Transformer models
- Modern autoregressive transformer LLMs



"Attention Is All You Need", Vaswani et al., 2017

Once upon a ... [EOS]



Source: Stefano Ermon, Deep Graphical Models

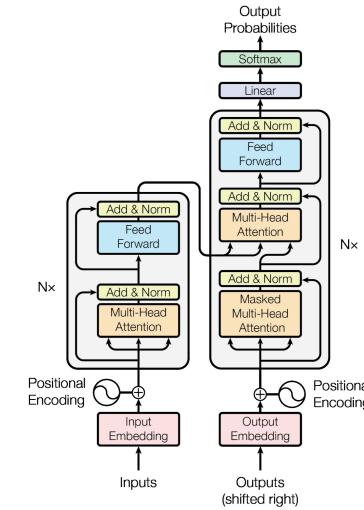
Lecture 7 (February 28) – Autoregressive Models

Lecture: Autoregressive generative models, the transformer neural architecture, and its use in large language modeling.

- Neural autoregressive density estimation
- Recurrent neural networks
- Attention-based models
- Transformer models
- Modern autoregressive transformer LLMs

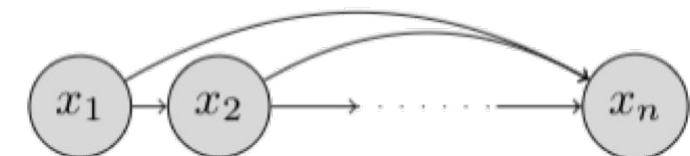
After:

- Paper presentations from students.



"Attention Is All You Need", Vaswani et al., 2017

Once upon a ... [EOS]



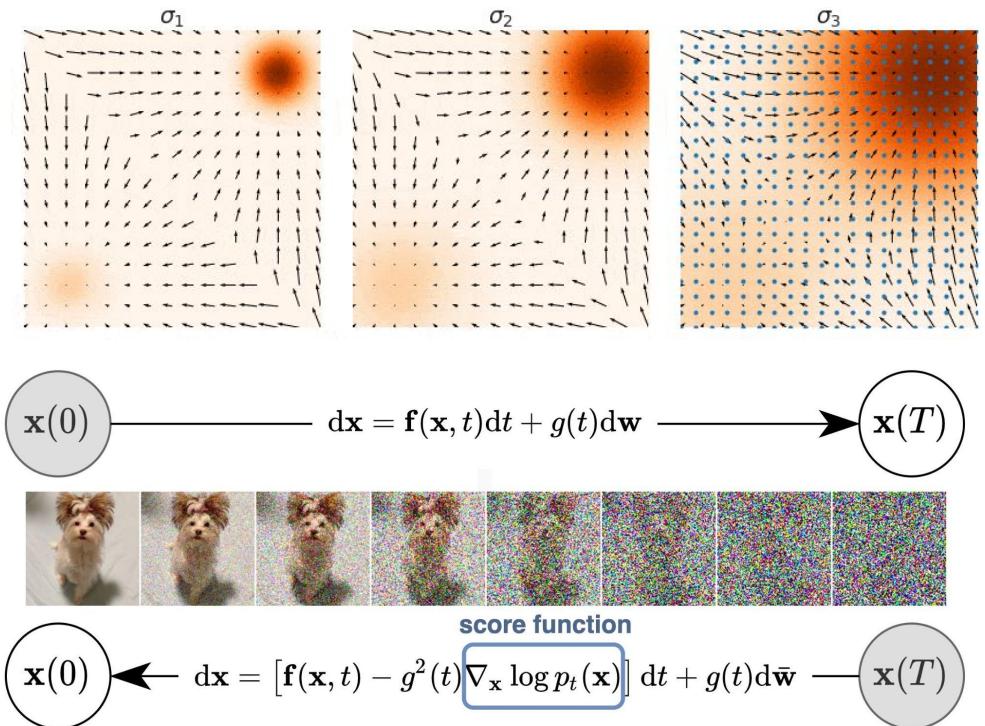
Source: Stefano Ermon, Deep Graphical Models

Lecture 8 (March 7) – Diffusion Models

Lecture 8 (March 7) – Diffusion Models

Lecture: Score-based generative models and diffusion models.

- The (Stein) score function, score matching.
- Score-based generative models.
- Sampling with Langevin Monte Carlo.
- Denoising diffusion probabilistic models.
- Learning with variational inference.



Source: Yang Song, "Generative Modeling by Estimating Gradients of the Data Distribution"

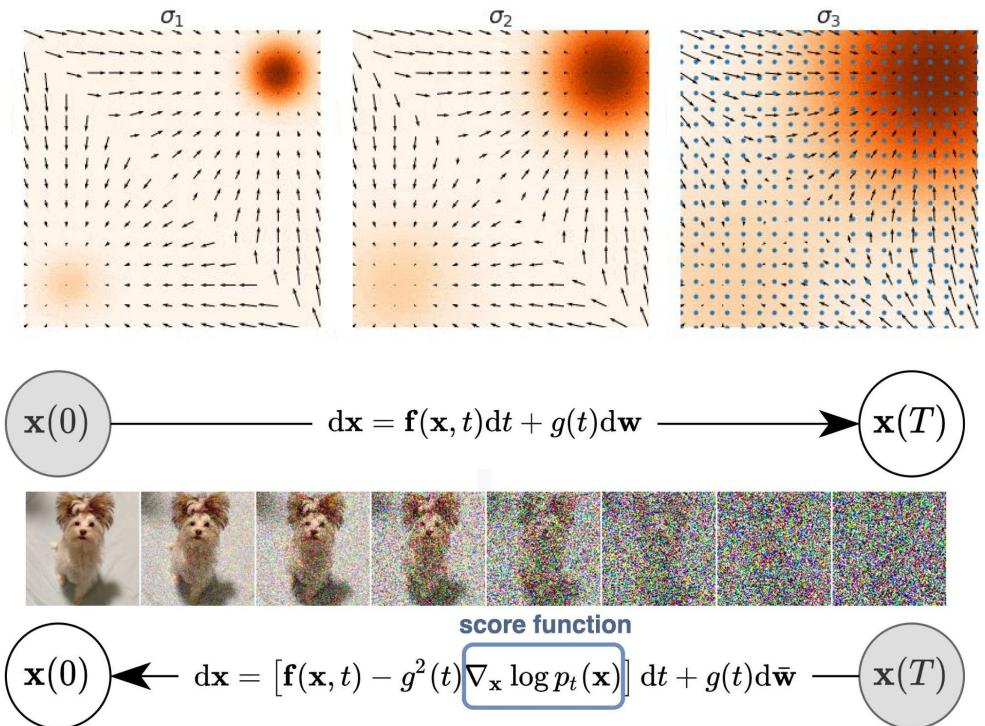
Lecture 8 (March 7) – Diffusion Models

Lecture: Score-based generative models and diffusion models.

- The (Stein) score function, score matching.
- Score-based generative models.
- Sampling with Langevin Monte Carlo.
- Denoising diffusion probabilistic models.
- Learning with variational inference.

After:

- Paper presentations from students.



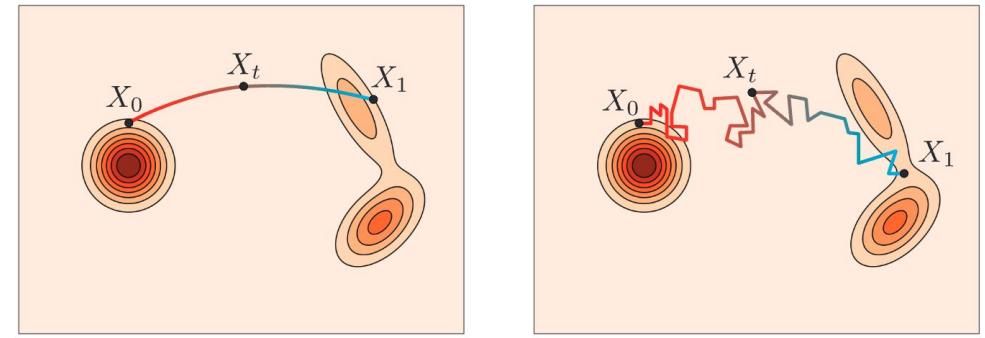
Source: Yang Song, "Generative Modeling by Estimating Gradients of the Data Distribution"

Lecture 9 (March 14) – Flow-based Models

Lecture 9 (March 14) – Flow-based Models

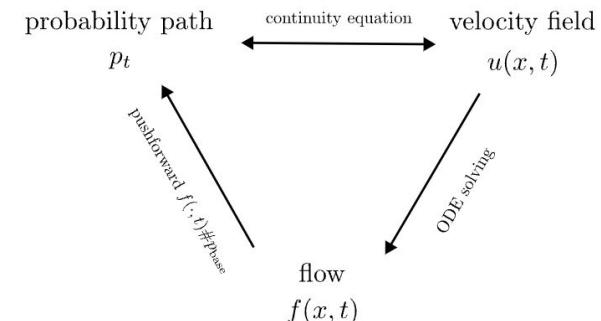
Lecture: Normalizing flows, continuous normalizing flows, neural ODEs, flow matching, and more.

- Normalizing flows
- Continuous normalizing flows
- Neural ordinary differential equations
- Conditional flow matching, and flow matching
- Extensions: discrete flow matching, generator matching.



(a) Flow

(b) Diffusion



"Flow Matching Guide and Code",
Lipman et al., 2024
"A Visual Dive into Conditional Flow Matching", 2024

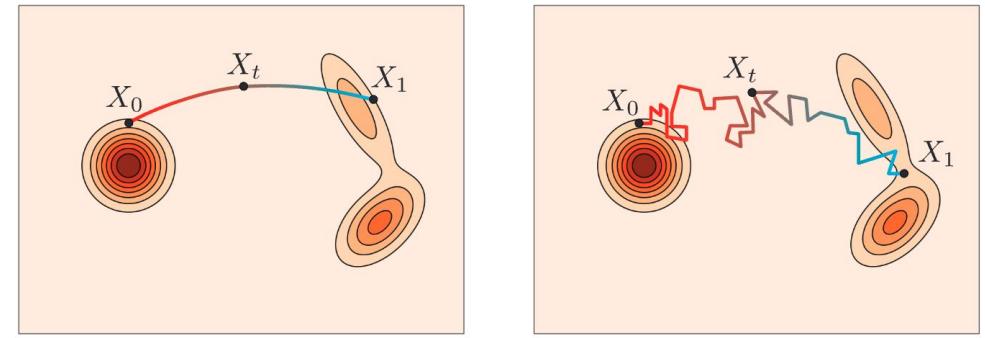
Lecture 9 (March 14) – Flow-based Models

Lecture: Normalizing flows, continuous normalizing flows, neural ODEs, flow matching, and more.

- Normalizing flows
- Continuous normalizing flows
- Neural ordinary differential equations
- Conditional flow matching, and flow matching
- Extensions: discrete flow matching, generator matching.

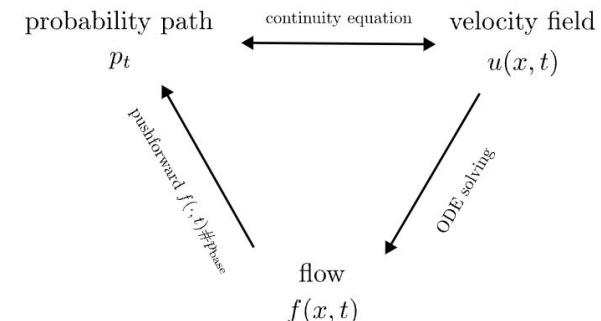
After:

- Paper presentations from students.



(a) Flow

(b) Diffusion



"Flow Matching Guide and Code",
Lipman et al., 2024
"A Visual Dive into Conditional Flow Matching", 2024

Note: March 21st is Spring Break

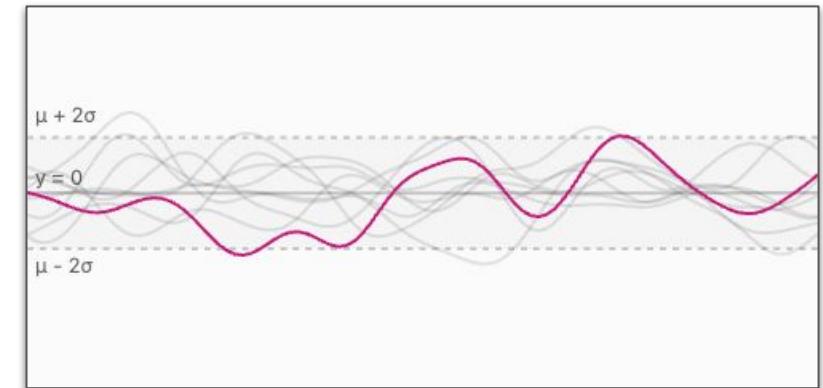


Lecture 10 (March 28) – Predictive Uncertainty Quantification

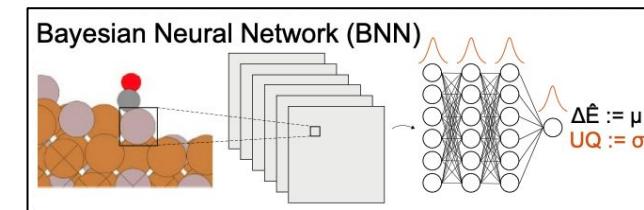
Lecture 10 (March 28) – Predictive Uncertainty Quantification

Lecture: Predictive UQ models, including Gaussian processes (GPs), Bayesian neural networks (BNNs), ensembles/dropout, UQ in LLMs, and relevant metrics.

- Classic predictive UQ: Bayesian parametric models, Gaussian processes (GPs).
- Neural/deep predictive UQ: Bayesian neural networks (BNN), ensembles, MC-dropout.
- Uncertainty quantification in LLMs.
- Relevant metrics: calibration, proper scoring rules.



A Visual Exploration of Gaussian Processes, distill.pub

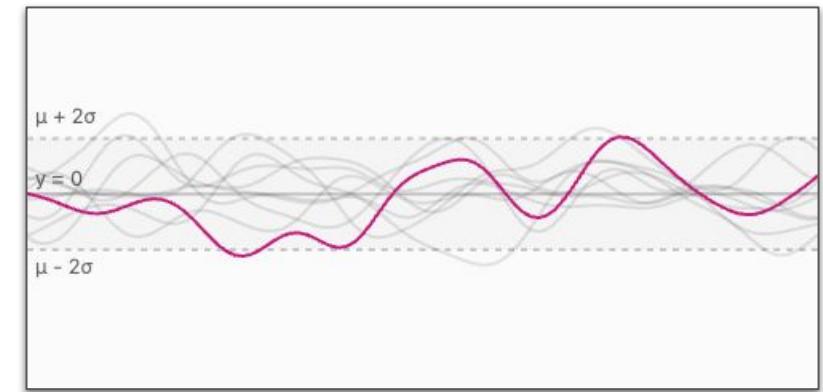


Tran, Neiswanger,
et al., MLST. 2020

Lecture 10 (March 28) – Predictive Uncertainty Quantification

Lecture: Predictive UQ models, including Gaussian processes (GPs), Bayesian neural networks (BNNs), ensembles/dropout, UQ in LLMs, and relevant metrics.

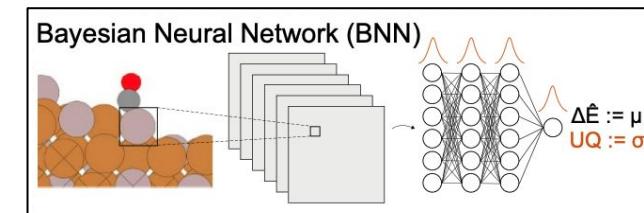
- Classic predictive UQ: Bayesian parametric models, Gaussian processes (GPs).
- Neural/deep predictive UQ: Bayesian neural networks (BNN), ensembles, MC-dropout.
- Uncertainty quantification in LLMs.
- Relevant metrics: calibration, proper scoring rules.



A Visual Exploration of Gaussian Processes, distill.pub

After:

- Paper presentations from students.



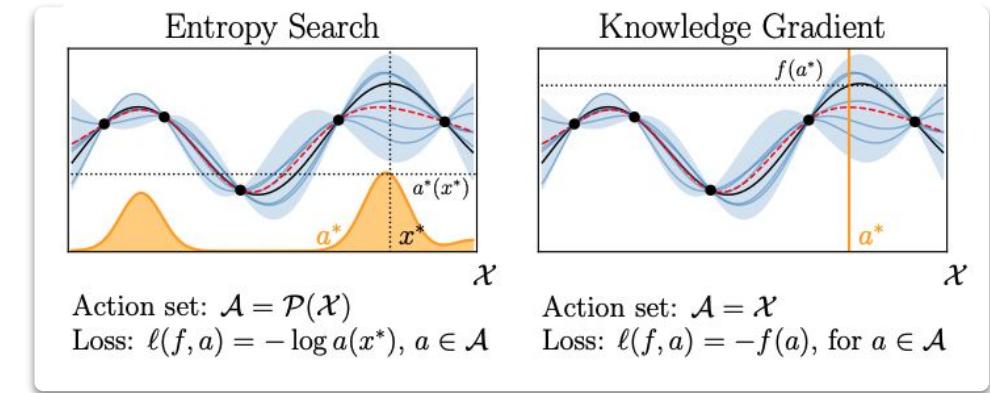
Tran, Neiswanger,
et al., MLST. 2020

Lecture 11 (April 4) – Active Learning and Optimization

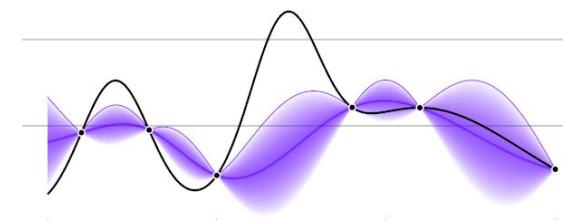
Lecture 11 (April 4) – Active Learning and Optimization

Lecture: Using predictive UQ models for active learning and sequential decision making, including Bayesian optimization and optimal experimental design.

- Decision making under uncertainty.
- Active learning.
- Bayesian optimization.
- Optimal experimental design.
- Extensions: Bayesian algorithm execution.



"Generalizing Bayesian Optimization with Decision-theoretic Entropies", 2022

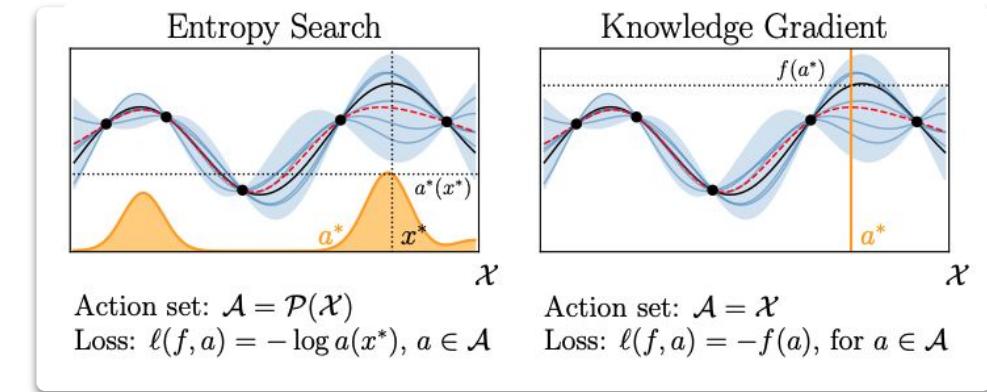


Bayesian
Optimization,
Wikipedia

Lecture 11 (April 4) – Active Learning and Optimization

Lecture: Using predictive UQ models for active learning and sequential decision making, including Bayesian optimization and optimal experimental design.

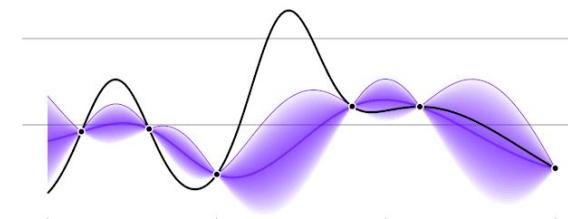
- Decision making under uncertainty.
- Active learning.
- Bayesian optimization.
- Optimal experimental design.
- Extensions: Bayesian algorithm execution.



"Generalizing Bayesian Optimization with Decision-theoretic Entropies", 2022

After:

- Paper presentations from students.



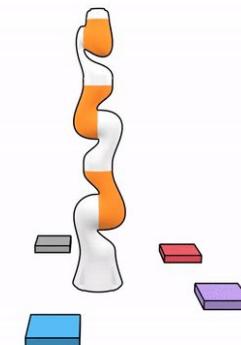
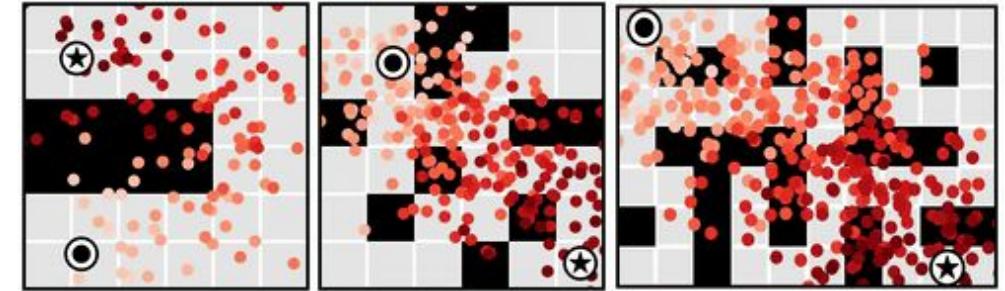
Bayesian
Optimization,
Wikipedia

Lecture 12 (April 11) – Generative Models in Decision Making

Lecture 12 (April 11) – Generative Models in Decision Making

Lecture: Using generative models (diffusion, VAE, LLMs, etc), in decision making and optimization procedures.

- High-dimensional BayesOpt with VAE.
- Diffusion planning.
- Vision-language-action models.
- Decision transformer.
- Decision making with LLMs.

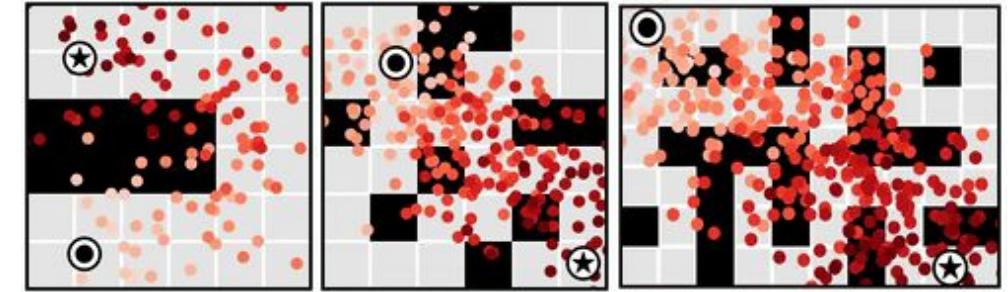


Source: Planning with
Diffusion for Flexible
Behavior Synthesis.

Lecture 12 (April 11) – Generative Models in Decision Making

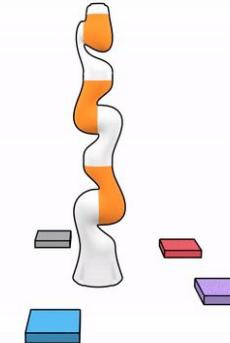
Lecture: Using generative models (diffusion, VAE, LLMs, etc), in decision making and optimization procedures.

- High-dimensional BayesOpt with VAE.
- Diffusion planning.
- Vision-language-action models.
- Decision transformer.
- Decision making with LLMs.



After:

- Paper presentations from students.



Source: Planning with
Diffusion for Flexible
Behavior Synthesis.

Lecture 13 (April 18) – Mystery Class

Lecture: Special topic, to be announced :-).



Lecture 13 (April 18) – Mystery Class

Lecture: Special topic, to be announced :-).

After:

- Finish all student paper presentations.



Lecture 14 (April 25) – Final Presentations #1

Full day of final presentations of class projects! (Teams 1-5)

- Each team will give a presentation.
- ~30 minutes long.
- Describing their full project and results.

Lecture 14 (May 2) – Final Presentations #2

Full day of final presentations of class projects! (Teams 6-10)

- Each team will give a presentation.
- ~30 minutes long.
- Describing their full project and results.

Probability Review

Probability Review

Goal: give a brief background/review of probability that will be useful in this class.

Probability Review

Goal: give a brief background/review of probability that will be useful in this class.

Terminology & notation:

- From elementary probability theory (with a few mentions of measure-theoretic probability theory).
- Which are used throughout modern machine learning work.

Probability Review

Goal: give a brief background/review of probability that will be useful in this class.

Terminology & notation:

- From elementary probability theory (with a few mentions of measure-theoretic probability theory).
- Which are used throughout modern machine learning work.

Will review:

- Probability space, conditional probability, chain rule, independence (for events).
- Random variables, random vectors.
- Cumulative distribution function (CDF), probability density/mass functions (PDF, PMF)
- Joint probability distributions, expectations, chain rule, etc.

Probability Review – Probability Space

Definition. Probability Space (Ω, \mathcal{F}, P)

Probability Review – Probability Space

Definition. Probability Space (Ω, \mathcal{F}, P)

Sample Space Ω

Event Space \mathcal{F}

Probability Function (Measure) P

Probability Review – Probability Space

Definition. Probability Space (Ω, \mathcal{F}, P)



Experiment: “roll a die and flip a coin”

Probability Review – Probability Space

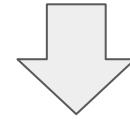
Definition. Probability Space (Ω, \mathcal{F}, P)

Sample Space Ω

- The set of all outcomes of a random experiment



Experiment: “roll a die and flip a coin”



Ω H-1, H-2, H-3,
H-4, H-5, H-6,
T-1, T-2, T-3,
T-4, T-5, T-6,

Probability Review – Probability Space

Definition. Probability Space (Ω, \mathcal{F}, P)

Event Space $\mathcal{F} \subseteq 2^\Omega$

- A set whose elements A (“events”) are subsets of Ω .

An example event:

The coin lands on heads and the die lands on an odd number.

Probability Review – Probability Space

Definition. Probability Space (Ω, \mathcal{F}, P)

Event Space $\mathcal{F} \subseteq 2^\Omega$

- A set whose elements A (“events”) are subsets of Ω .
- Satisfies three properties:
 - contains empty set
 - closed under complements
 - closed under countable unions

An example event:

The coin lands on heads and the die lands on an odd number.

$$\emptyset \in \mathcal{F}$$

$$A \in \mathcal{F} \implies \Omega \setminus A \in \mathcal{F}$$

$$A_1, A_2, \dots \in \mathcal{F} \implies \cup_i A_i \in \mathcal{F}$$

Probability Review – Probability Space

Definition. Probability Space (Ω, \mathcal{F}, P)

Probability Function (Measure) P

- A function $P : \mathcal{F} \rightarrow \mathbb{R}$.

Probability Review – Probability Space

Definition. Probability Space (Ω, \mathcal{F}, P)

Probability Function (Measure) P

- A function $P : \mathcal{F} \rightarrow \mathbb{R}$.
- Satisfies three properties:
 - Non-negativity
 - Normalization
 - Finite-additivity

“Axioms of Probability”

$$P(A) \geq 0, \text{ for all } A \in \mathcal{F}$$

$$P(\Omega) = 1$$

$$\text{If } A_1, A_2, \dots \text{ are disjoint} \implies P(\cup_i A_i) = \sum_i P(A_i)$$

Probability Review – Probability Space

Definition. Probability Space (Ω, \mathcal{F}, P)

Measure Theoretic Brief

\mathcal{F} is a σ -algebra.

(Ω, \mathcal{F}) is a measurable space.

P is a probability measure.

Probability Function (Measure) P

- A function $P : \mathcal{F} \rightarrow \mathbb{R}$.
- Satisfies three properties:
 - Non-negativity
 - Normalization
 - Finite-additivity

“Axioms of Probability”

$$P(A) \geq 0, \text{ for all } A \in \mathcal{F}$$

$$P(\Omega) = 1$$

$$\text{If } A_1, A_2, \dots \text{ are disjoint} \implies P(\cup_i A_i) = \sum_i P(A_i)$$

Probability Review – Probability Space

Some properties of a probability space (Ω, \mathcal{F}, P) :

Probability Review – Probability Space

Some properties of a probability space (Ω, \mathcal{F}, P) :

- $A \subseteq B \implies P(A) \leq P(B)$
- $P(A \cap B) \leq \min(P(A), P(B))$
- Union Bound: $P(A \cup B) \leq P(A) + P(B)$
- $P(\Omega - A) = 1 - P(A)$
- Law of Total Probability: If A_1, \dots, A_k are disjoint, and $\bigcup_{i=1}^k A_i = \Omega$, then $\sum_{i=1}^k P(A_i) = 1$

Probability Review – Conditional Probability (Events)

Definition. Conditional probability of any event A given event B:
(where B has non-zero probability)

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

“Probability measure of the event A, assuming the occurrence of event B”

Probability Review – Conditional Probability (Events)

Definition. Conditional probability of any event A given event B:
(where B has non-zero probability)

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

Example:

Probability that the die lands on an odd number, given that the die lands on >2.

“Probability measure of the event A, assuming the occurrence of event B”

Probability Review – Chain Rule (*Events*)

Definition. Let A_1, \dots, A_k be events, where $P(A_i) > 0$. Then the chain rule states that:

$$\begin{aligned} & P(A_1 \cap A_2 \cap \dots \cap A_k) \\ &= P(A_1)P(A_2|A_1)P(A_3|A_2 \cap A_1) \cdots P(A_k|A_1 \cap A_2 \cap \dots \cap A_{k-1}) \end{aligned}$$

⇒ derived by applying conditional probability multiple times!

Probability Review – Independence (*Events*)

Definition. Two events A and B are called independent if

$$P(A \cap B) = P(A)P(B)$$

or

$$P(A | B) = P(A)$$

“Observing B does not have any effect on the probability of A”

Probability Review – Random Variables

Probability Review – Random Variables

Suppose we just care about a (real-valued) function of outcomes $\omega \in \Omega$.

E.g. if Ω is outcomes of 5 coin flips, we might care about “number of heads that occur”, or “length of longest run of tails” \Rightarrow both are real-valued functions.



And we care about the induced distribution over a numeric/real space, given by this function.

Probability Review – Random Variables

Definition. A Random Variable is a (measurable) function $X : \Omega \rightarrow \mathbb{R}$. We can define the probability of X in the following way.

Note: “measurable function” ensures that this is an event in \mathcal{F} with a well-defined probability

Probability Review – Random Variables

Definition. A Random Variable is a (measurable) function $X : \Omega \rightarrow \mathbb{R}$. We can define the probability of X in the following way.

For discrete random variables:

Note: “measurable function” ensures that this is an event in \mathcal{F} with a well-defined probability

$$P(X = k) = P(\{\omega : X(\omega) = k\})$$

And for continuous random variables:

$$P(a \leq X \leq b) = P(\{\omega : a \leq X(\omega) \leq b\})$$

Probability Review – Random Variables

Definition. A Random Variable is a (measurable) function $X : \Omega \rightarrow \mathbb{R}$. We can define the probability of X in the following way.

Doob, J.

Quoted in Statistical Science

“While writing my book [Stochastic Processes] I had an argument with Feller. He asserted that everyone said ‘random variable’ and I asserted that everyone said ‘chance variable.’

We obviously had to use the same name in our books, so we decided the issue by a stochastic procedure. That is, we tossed for it and he won.”

Note: “measurable function” ensures that this is an event in \mathcal{F} with a well-defined probability

$$\{\omega : X(\omega) = k\})$$

$$\{\omega : a \leq X(\omega) \leq b\})$$

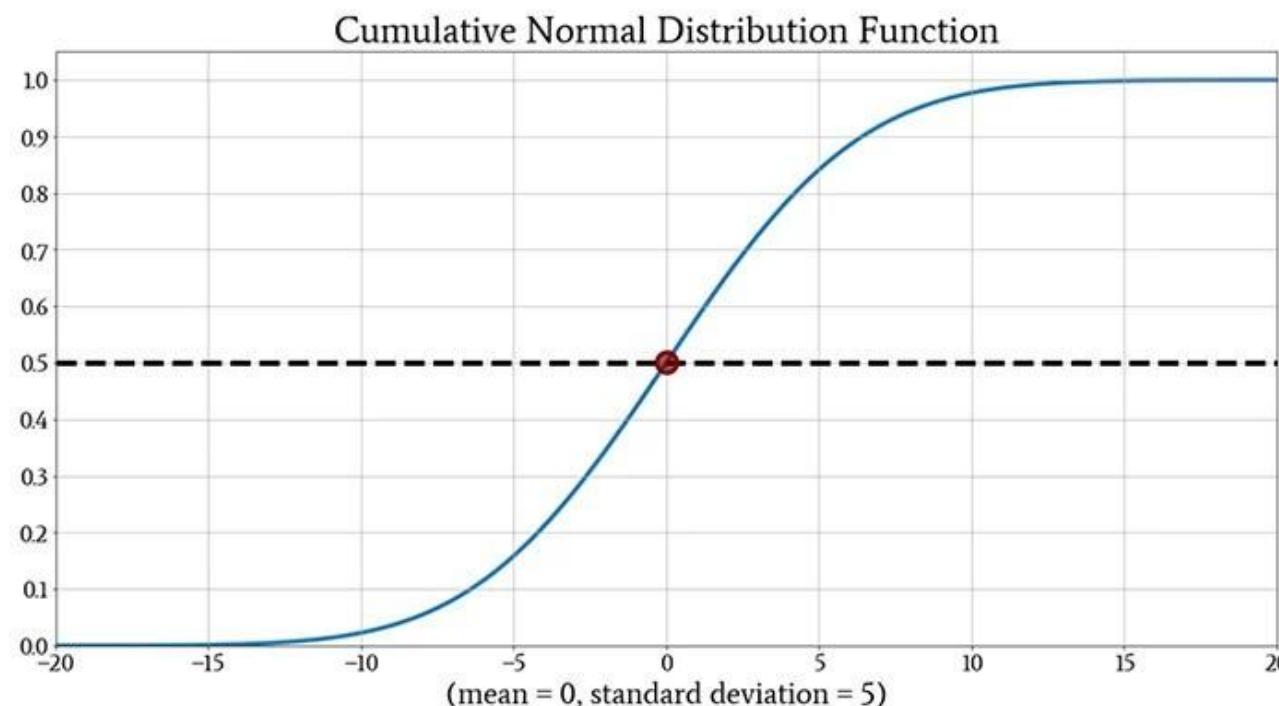
Probability Review – CDF, PMF, PDF

It can be useful to define alternative functions to specify probability measures of random variables – e.g. CDF, PMF, and PDF.

Probability Review – Cumulative Distribution Function (CDF)

Probability Review – Cumulative Distribution Function (CDF)

Definition. Cumulative Distribution Function (CDF) is a function $F_X : \mathbb{R} \rightarrow [0, 1]$ which specifies a probability measure as $F_X(x) = P(X \leq x)$.



Probability Review – Cumulative Distribution Function (CDF)

Definition. Cumulative Distribution Function (CDF) is a function $F_X : \mathbb{R} \rightarrow [0, 1]$ which specifies a probability measure as $F_X(x) = P(X \leq x)$.

- $0 \leq F_X(x) \leq 1$
- $\lim_{x \rightarrow -\infty} F_X(x) = 0$
- $\lim_{x \rightarrow \infty} F_X(x) = 1$
- $x \leq y \implies F_X(x) \leq F_X(y)$

Probability Review – Probability Mass Function (PMF)

Probability Review – Probability Mass Function (PMF)

Definition. When X is a discrete random variable, a Probability Mass Function (PMF) is a function $p_X : \Omega \rightarrow \mathbb{R}$ such that $p_X(x) = P(X = x)$.

Note: $\text{Val}(X)$ is the set of possible sample-space values that X can take on.

Properties:

- $0 \leq p_X(x) \leq 1$
- $\sum_{x \in \text{Val}(X)} p_X(x) = 1$
- $\sum_{x \in A} p_X(x) = P(X \in A)$

Probability Review – Probability Density Function (PDF)

Definition. When X is a continuous random variable, and the CDF is differentiable everywhere, a Probability Density Function (PDF) is $f_X : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$f_X(x) = \frac{dF_X(x)}{dx}$$

For very small Δx , we have $P(x \leq X \leq x + \Delta x) \approx f_X(x)\Delta x$

Probability Review – Probability Density Function (PDF)

Definition. When X is a continuous random variable, and the CDF is differentiable everywhere, a Probability Density Function (PDF) is $f_X : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$f_X(x) = \frac{dF_X(x)}{dx}$$

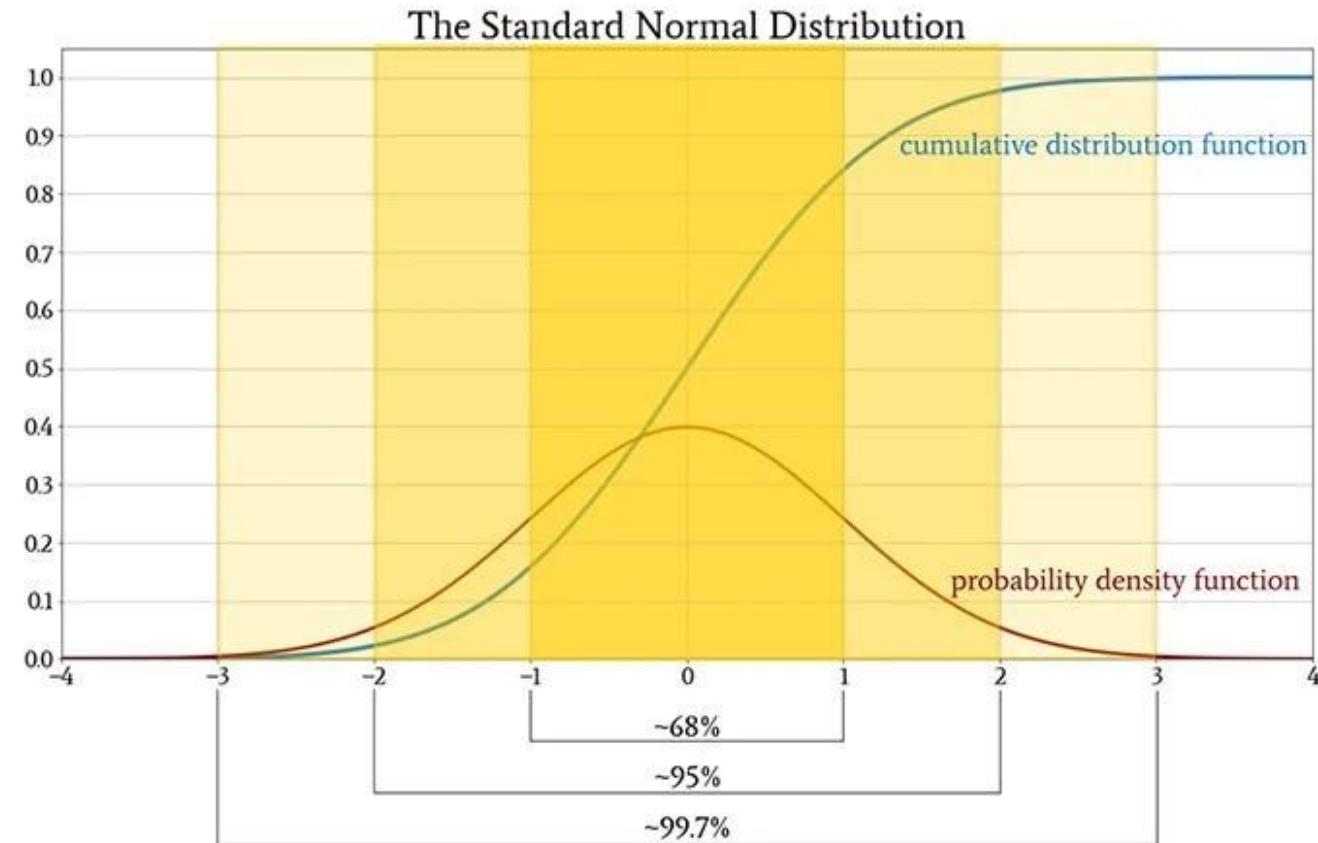
For very small Δx , we have $P(x \leq X \leq x + \Delta x) \approx f_X(x)\Delta x$

Properties:

- $f_X(x) \geq 0$
- $\int_{-\infty}^{\infty} f_X(x)dx = 1$
- $\int_{x \in A} f_X(x)dx = P(X \in A)$

Probability Review – Probability Density Function (PDF)

Definition. When X is a continuous random variable, and the CDF is differentiable everywhere, a Probability Density Function (PDF) is $f_X : \mathbb{R} \rightarrow \mathbb{R}$ such that



Probability Review – Example Distributions

We will encounter many common distributions in this class, including:

Continuous

- Normal (Gaussian) distribution
- Uniform distribution (continuous)
- Gamma distribution
- Beta distribution
- Exponential distribution
- Chi-squared distribution
- Fisher-Snedecor distribution

Discrete

- Poisson distribution
- Binomial distribution
- Geometric distribution
- Hypergeometric distribution
- Discrete uniform distribution

Probability Review – Expectation

Probability Review – Expectation

Definition. Let $g(X)$ be an arbitrary test function $g : \mathbb{R} \rightarrow \mathbb{R}$.

Then for a discrete random variable

$$\mathbb{E}[g(X)] = \sum_{x \in \text{Val}(X)} g(x)p_X(x)$$

And for a continuous random variable

$$\mathbb{E}[g(X)] = \int_{-\infty}^{\infty} g(x)f_X(x)dx$$

Probability Review – Expectation

Definition. Let $g(X)$ be an arbitrary test function $g : \mathbb{R} \rightarrow \mathbb{R}$.

Then for a discrete random variable

$$\mathbb{E}[g(X)] = \sum_{x \in \text{Val}(X)} g(x)p_X(x)$$

And for a continuous random variable **(alternative notation)**

$$\mathbb{E}_{x \sim f_X}[g(x)] = \int_{-\infty}^{\infty} g(x)f_X(x)dx$$

Probability Review – Variance

Probability Review – Variance

Definition. The variance is a measure of how concentrated a random variable is around its mean, and is defined as

$$\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2]$$

We can also write the variance as

$$\begin{aligned}\mathbb{E}[(X - \mathbb{E}[X])^2] &= \mathbb{E}[X^2 - 2\mathbb{E}[X]X + \mathbb{E}[X]^2] \\ &= \mathbb{E}[X^2] - 2\mathbb{E}[X]\mathbb{E}[X] + \mathbb{E}[X]^2 \\ &= \mathbb{E}[X^2] - \mathbb{E}[X]^2\end{aligned}$$

Probability Review – Properties of Expectation and Variance

A few properties of note for the expectation and variance:

Probability Review – Properties of Expectation and Variance

A few properties of note for the expectation and variance:

- $\mathbb{E}[a] = a$ for any constant $a \in \mathbb{R}$
- $\mathbb{E}[af(X)] = a\mathbb{E}[f(X)]$ for any constant $a \in \mathbb{R}$
- Linearity of Expectation: $\mathbb{E}[f(X) + g(X)] = \mathbb{E}[f(X)] + \mathbb{E}[g(X)]$
- $\text{Var}[a] = 0$ for any constant $a \in \mathbb{R}$
- $\text{Var}[af(X)] = a^2\text{Var}[f(X)]$ for any constant $a \in \mathbb{R}$
- If $X \perp Y$, then $\text{Var}[f(X) + g(Y)] = \text{Var}[f(X)] + \text{Var}[g(Y)]$

Probability Review – Joint Distributions (Two or More Random Variables)

Suppose we have two random variables, X and Y , and want to model their outcomes *simultaneously* during a random experiment.

Probability Review – Joint CDF

Probability Review – Joint CDF

Definition. The Joint Cumulative Distribution Function (joint CDF) of random variables X and Y is defined to be

$$F_{XY}(x, y) = P(X \leq x, Y \leq y)$$

By knowing the joint CDF, any event involving X and Y can be calculated!

Probability Review – Joint CDF

Definition. The Joint Cumulative Distribution Function (joint CDF) of random variables X and Y is defined to be

$$F_{XY}(x, y) = P(X \leq x, Y \leq y)$$

The joint CDF is related to the **marginal CDF** for each random variable by:

$$F_X(x) = \lim_{y \rightarrow \infty} F_{XY}(x, y)dy$$

$$F_Y(y) = \lim_{x \rightarrow \infty} F_{XY}(x, y)dx$$

Probability Review – Joint PDF

Probability Review – Joint PDF

Definition. Let X and Y be random variables with joint CDF $F_{XY}(x, y)$. Then the joint probability density function (joint PDF) is defined to be

$$f_{XY}(x, y) = \frac{\partial^2 F_{XY}(x, y)}{\partial x \partial y}$$

And, similarly to the one-dimensional case, we have that

$$\int_{x,y \in A} f_{XY}(x, y) dx dy = P((X, Y) \in A)$$

Probability Review – Marginal PDFs

Probability Review – Marginal PDFs

Definition. Let X and Y be random variables with joint PDF $f_{XY}(x, y)$. Then the marginal probability density functions (marginal PDFs) are defined to be

$$f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x, y) dy$$

$$f_Y(y) = \int_{-\infty}^{\infty} f_{XY}(x, y) dx$$

Often say “marginalize out y” (top), or “marginalize out x” (bottom).

Probability Review – Conditional Distributions

Probability Review – Conditional Distributions

Seeks to answer: “What is the probability distribution over X when we know that Y has a certain value?”

Definition. The conditional distribution of X given that $Y = y$ is defined to be

$$f_{X|Y}(x | y) = \frac{f_{XY}(x, y)}{f_Y(y)}$$

assuming that $f_Y(y) \neq 0$.

Probability Review – Conditional Distributions

Seeks to answer: “What is the probability distribution over X when we know that Y has a certain value?”

Definition. The conditional distribution of X given that $Y = y$ is defined to be

$$f_{X_1|X_2,\dots,X_n}(x_1 \mid x_2, \dots, x_n) = \frac{f_{X_1,X_2,\dots,X_n}(x_1, x_2, \dots, x_n)}{f_{X_2,\dots,X_n}(x_2, \dots, x_n)}$$

assuming that the denominator is not equal to 0.

For higher dimensions!

Probability Review – Conditional Distributions

Seeks to answer: “What is the probability distribution over X when we know that Y has a certain value?”

Definition. The conditional distribution of X given that $Y = y$ is defined to be

$$p(x_1 \mid x_2, \dots, x_n) = \frac{p(x_1, x_2, \dots, x_n)}{p(x_2, \dots, x_n)}$$

assuming that the denominator is not equal to 0.

Common notation for PDFs!

Probability Review – Conditional Distributions

Seeks to answer: “What is the probability distribution over X when we know that Y has a certain value?”

Note that the conditional PDF implies the following factorization:

$$f_{X|Y}(x | y) = \frac{f_{XY}(x, y)}{f_Y(y)} \iff$$

$$f_{XY}(x, y) = f_{X|Y}(x | y)f_Y(y) = f_{Y|X}(y | x)f_X(x)$$

This will show up a lot!

Probability Review – Bayes Rule



An important expression of conditional distribution of random variables!

Wikipedia

Probability Review – Bayes Rule



An important expression of conditional distribution of random variables!

Wikipedia

Definition. Bayes Rule is defined to be

$$f_{Y|X}(y | x) = \frac{f_{XY}(x, y)}{f_X(x)} = \frac{f_{X|Y}(x | y)f_Y(y)}{\int_{-\infty}^{\infty} f_{X|Y}(x | y')f_Y(y')dy'}$$

Probability Review – Bayes Rule



An important expression of conditional distribution of random variables!

Wikipedia

Definition. Bayes Rule is defined to be

$$f_{Y|X}(y | x) = \frac{f_{XY}(x, y)}{f_X(x)} = \frac{f_{X|Y}(x | y)f_Y(y)}{\int_{-\infty}^{\infty} f_{X|Y}(x | y')f_Y(y')dy'}$$

Likelihood

Prior

Marginal / Evidence

Probability Review – Expectation (for Joint PDF)

Probability Review – Expectation (for Joint PDF)

Suppose we have two random variables and a test function $g : \mathbb{R}^2 \rightarrow \mathbb{R}$. Then the expected value can be written

$$\mathbb{E}_{x,y \sim f_{XY}} [g(x, y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f_{XY}(x, y) dx dy$$

Notation often shortened to:

$$\mathbb{E}_{x,y} [g(x, y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f_{XY}(x, y) dx dy$$

Probability Review – Covariance

Probability Review – Covariance

Definition. The covariance of two random variables X and Y is defined to be

$$\text{Cov}[X, Y] = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

Describes whether two variables tend to move in the same direction (positive covariance) or opposite directions (negative covariance).

Probability Review – Covariance

Definition. The covariance of two random variables X and Y is defined to be

$$\text{Cov}[X, Y] = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

Which alternatively can be written

$$\begin{aligned}\text{Cov}[X, Y] &= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY - X\mathbb{E}[Y] - Y\mathbb{E}[X] + \mathbb{E}[X]\mathbb{E}[Y]] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] - \mathbb{E}[Y]\mathbb{E}[X] + \mathbb{E}[X]\mathbb{E}[Y] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]\end{aligned}$$

Probability Review – Chain Rule (for Joint PDFs)

Probability Review – Chain Rule (for Joint PDFs)

Similar to the expression above, for a joint pdf we can write the chain rule as:

Using more-common notation for PDFs

Joint PDF factorization

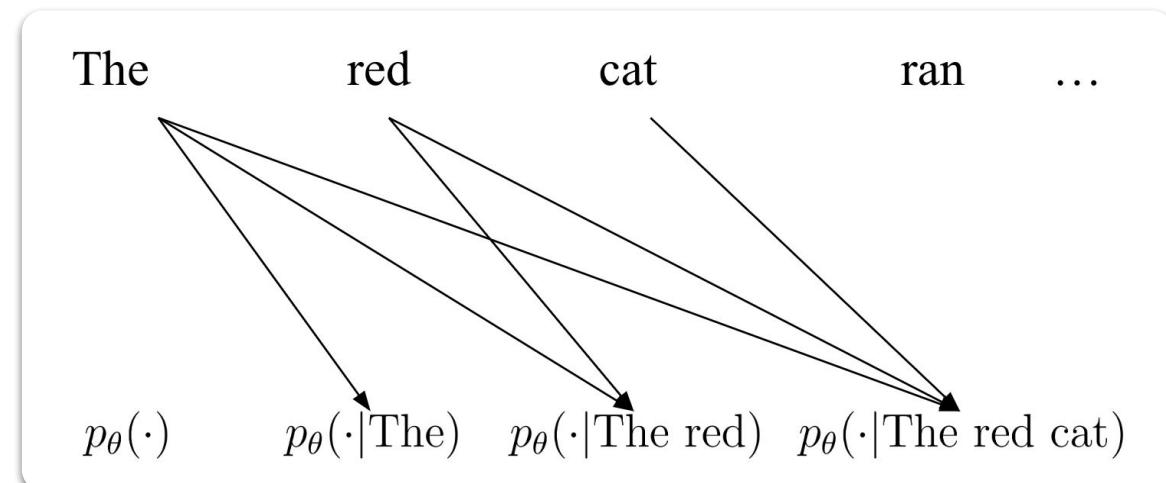
$$\begin{aligned} p(x_1, x_2, \dots, x_n) &= p(x_n \mid x_1, x_2, \dots, x_{n-1})p(x_1, x_2, \dots, x_{n-1}) && \text{(repeatedly!)} \\ &= p(x_n \mid x_1, x_2, \dots, x_{n-1})p(x_{n-1} \mid x_1, x_2, \dots, x_{n-2})p(x_1, x_2, \dots, x_{n-2}) \\ &= \dots \\ &= p(x_1) \prod_{i=2}^n p(x_i \mid x_1, \dots, x_{i-1}) \end{aligned}$$

Full joint PDF factorization

Probability Review – Chain Rule (for Joint PDFs)

$$\begin{aligned} p(x_1, x_2, \dots, x_n) &= p(x_n \mid x_1, x_2, \dots, x_{n-1})p(x_1, x_2, \dots, x_{n-1}) \\ &= p(x_n \mid x_1, x_2, \dots, x_{n-1})p(x_{n-1} \mid x_1, x_2, \dots, x_{n-2})p(x_1, x_2, \dots, x_{n-2}) \\ &= \dots \\ &= p(x_1) \prod_{i=2}^n p(x_i \mid x_1, \dots, x_{i-1}) \end{aligned}$$

Autoregressive Models (LLMs)!



Probability Review – Independence (for Joint PDFs)

Probability Review – Independence (for Joint PDFs)

Using more-common notation for PDFs

We say that random variables X_1, X_2, \dots, X_n are independent if

$$p(x_1, x_2, \dots, x_n) = p(x_1)p(x_2) \cdots p(x_n)$$

Probability Review – Independence (for Joint PDFs)

Using more-common notation for PDFs

We say that random variables X_1, X_2, \dots, X_n are independent if

$$p(x_1, x_2, \dots, x_n) = p(x_1)p(x_2) \cdots p(x_n)$$

Contrast this with the previous factorization (no independence assumption):

$$p(x_1, x_2, \dots, x_n) = p(x_1) \prod_{i=2}^n p(x_i \mid x_1, \dots, x_{i-1})$$

⇒ intuitively, independence means that knowing the value of one variable does not have any effect on the conditional probability distribution of the other variable.

Probability Review – Random Vectors

When working with n random variables, often convenient to describe them as a **random vector**

$$X : \Omega \rightarrow \mathbb{R}^n \quad X = [X_1, X_2, \dots, X_n]^\top$$

Probability Review – Random Vectors

When working with n random variables, often convenient to describe them as a **random vector**

$$X : \Omega \rightarrow \mathbb{R}^n \quad X = [X_1, X_2, \dots, X_n]^\top$$

Where we can write the expected value, given a test function $g : \mathbb{R}^n \rightarrow \mathbb{R}$, as

$$\begin{aligned} E[g(X)] &= E_x[g(x)] = \int_{\mathbb{R}^n} g(x_1, x_2, \dots, x_n) p(x_1, x_2, \dots, x_n) dx_1 dx_2 \cdots dx_n \\ &= \int_{\mathbb{R}^n} g(x) p(x) dx, \end{aligned} \quad x = (x_1, \dots, x_n) \in \mathbb{R}^n$$

Probability Review – Covariance Matrix

Probability Review – Covariance Matrix

A covariance matrix is the $n \times n$ square matrix with entries $\Sigma_{ij} = \text{Cov}[X_i, X_j]$ written as

$$\Sigma = \begin{bmatrix} \text{Cov}[X_1, X_1] & \cdots & \text{Cov}[X_1, X_n] \\ \vdots & \ddots & \vdots \\ \text{Cov}[X_n, X_1] & \cdots & \text{Cov}[X_n, X_n] \end{bmatrix}$$

$$= \begin{bmatrix} E[X_1^2] - E[X_1]E[X_1] & \cdots & E[X_1X_n] - E[X_1]E[X_n] \\ \vdots & \ddots & \vdots \\ E[X_nX_1] - E[X_n]E[X_1] & \cdots & E[X_n^2] - E[X_n]E[X_n] \end{bmatrix}$$

$$= E[XX^\top] - E[X]E[X]^\top = \cdots = E((X - E[X])(X - E[X])^\top).$$

Next Class

Next Class

Lecture: PGMs – using ideas from both probability and graph theory to derive efficient machine learning algorithms

After:

- Class introductions and research interests.
- Everyone should share their name, and a couple of research interests. We will record this in a doc and share it with the class.
 - The goal is to help facilitate matching for group projects!