

# Lecture 01 - Deep Generative Modeling

**Baojian Zhou**

**Topics on Diffusion Models**

**School of Data Science, Fudan University**

# The book

<https://arxiv.org/abs/2510.21890>

## The Principles of Diffusion Models

## From Origins to Advances



Chieh-Hsin Lai  
Sony AI

Yang Song  
OpenAI

Dongjun Kim  
Stanford University

**Yuki Mitsufuji**  
Sony Corporation, Sony AI

**Stefano Ermon**  
Stanford University

## ① Chieh-Hsin (Jesse) Lai

- Research Scientist at Sony AI's Music Foundation Model Team
- Ph.D. in Mathematics from the University of Minnesota, Twin Cities
- Latest Publications

① Chieh-Hsin (Jesse) Lai

- Research Scientist at Sony AI's Music Foundation Model Team
  - Ph.D. in Mathematics from the University of Minnesota, Twin Cities
  - [Latest Publications](#)

② Song Yang

- Research Principal at Meta Superintelligence Labs
  - Ph.D. in CS from Stanford University, advised by Stefano Ermon
  - [Latest Publications](#)

① Chieh-Hsin (Jesse) Lai

- Research Scientist at Sony AI's Music Foundation Model Team
  - Ph.D. in Mathematics from the University of Minnesota, Twin Cities
  - [Latest Publications](#)

② Song Yang

- Research Principal at Meta Superintelligence Labs
  - Ph.D. in CS from Stanford University, advised by Stefano Ermon
  - [Latest Publications](#)

③ Dongjun Kim

- Postdoctoral scholar in Prof. Stefano Ermon's group at Stanford
  - Ph.D at KAIST
  - [Latest Publications](#)

① Chieh-Hsin (Jesse) Lai

- Research Scientist at Sony AI's Music Foundation Model Team
  - Ph.D. in Mathematics from the University of Minnesota, Twin Cities
  - [Latest Publications](#)

② Song Yang

- Research Principal at Meta Superintelligence Labs
  - Ph.D. in CS from Stanford University, advised by Stefano Ermon
  - [Latest Publications](#)

③ Dongjun Kim

- Postdoctoral scholar in Prof. Stefano Ermon's group at Stanford
  - Ph.D at KAIST
  - [Latest Publications](#)

## ④ Yuki Mitsufuji

- Distinguished Engineer, Sony Group Corporation
  - PhD in Information Science & Technology, the University of Tokyo
  - [Latest Publications](#)

## ① Chieh-Hsin (Jesse) Lai

- Research Scientist at Sony AI's Music Foundation Model Team
- Ph.D. in Mathematics from the University of Minnesota, Twin Cities
- [Latest Publications](#)

## ② Song Yang

- Research Principal at Meta Superintelligence Labs
- Ph.D. in CS from Stanford University, advised by Stefano Ermon
- [Latest Publications](#)

## ③ Dongjun Kim

- Postdoctoral scholar in Prof. Stefano Ermon's group at Stanford
- Ph.D at KAIST
- [Latest Publications](#)

## ④ Yuki Mitsufuji

- Distinguished Engineer, Sony Group Corporation
- PhD in Information Science & Technology, the University of Tokyo
- [Latest Publications](#)

## ⑤ Stefano Ermon

- Associate Professor, Department of CS, Stanford University
- Computer Science, Cornell University
- [Latest Publications](#)

## Main content

- The principles that have shaped the development of diffusion models
  - Diffusion modeling starts by defining a forward process that gradually corrupts data into noise, linking the data distribution to a simple prior through a continuum of intermediate distributions.
  - The main goal is to learn a reverse process that transforms noise back into data while recovering the same intermediates.

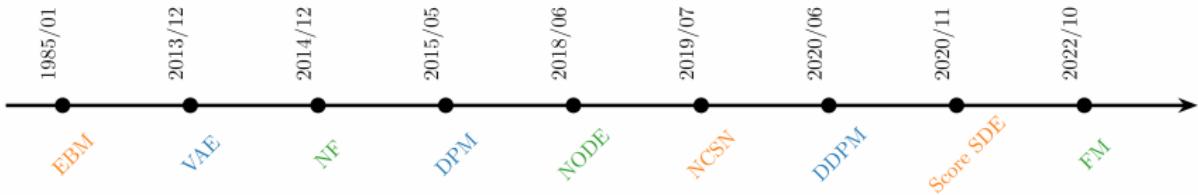
# Main content

- The principles that have shaped the development of diffusion models
- Diffusion modeling starts by defining a forward process that gradually corrupts data into noise, linking the data distribution to a simple prior through a continuum of intermediate distributions.
- The main goal is to learn a reverse process that transforms noise back into data while recovering the same intermediates.
- Three complementary ways to formalize the above two processes:
  - Variational view (Variational autoencoders)
  - Score-based view (EBMs)
  - Flow-based view (Normalizing flows)

They share a common backbone: a learned time-dependent velocity field whose flow transports a simple prior to the data.

- Sampling amounts to solving a differential equation that evolves noise into data along a continuous generative trajectory.

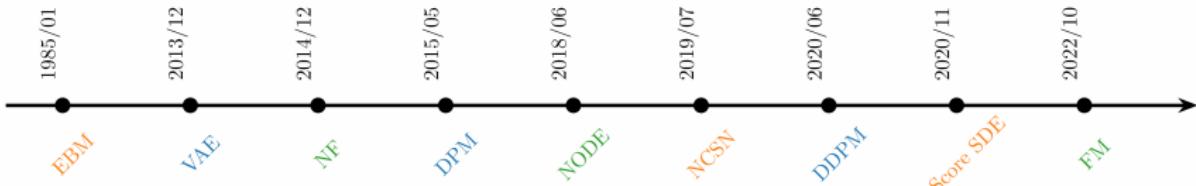
# Timeline of diffusion model perspectives



## ① Score-based view

- Energy-Based Model (EBM) [Ackley et al., 1985], Noise Conditional Score Network (NCSN) [Song and Ermon, 2019], Score SDE [Song et al., 2020]

# Timeline of diffusion model perspectives



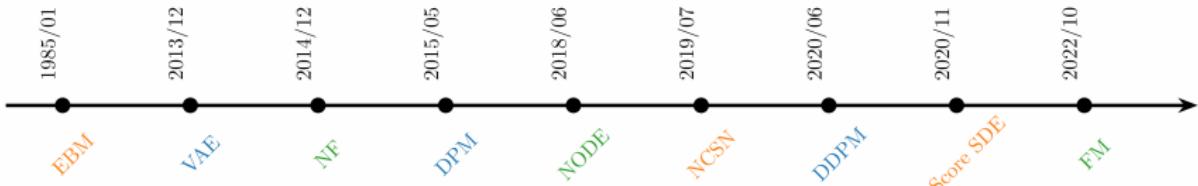
## ① Score-based view

- Energy-Based Model (EBM) [Ackley et al., 1985], Noise Conditional Score Network (NCSN) [Song and Ermon, 2019], Score SDE [Song et al., 2020]

## ② Variational view

- Variational Autoencoder (VAE) [Kingma and Welling, 2013], Diffusion Probabilistic Models (DPM) [Sohl-Dickstein et al., 2015], DDPM [Ho et al., 2020]

# Timeline of diffusion model perspectives



## ① Score-based view

- Energy-Based Model (EBM) [Ackley et al., 1985], Noise Conditional Score Network (NCSN) [Song and Ermon, 2019], Score SDE [Song et al., 2020]

## ② Variational view

- Variational Autoencoder (VAE) [Kingma and Welling, 2013], Diffusion Probabilistic Models (DPM) [Sohl-Dickstein et al., 2015], DDPM [Ho et al., 2020]

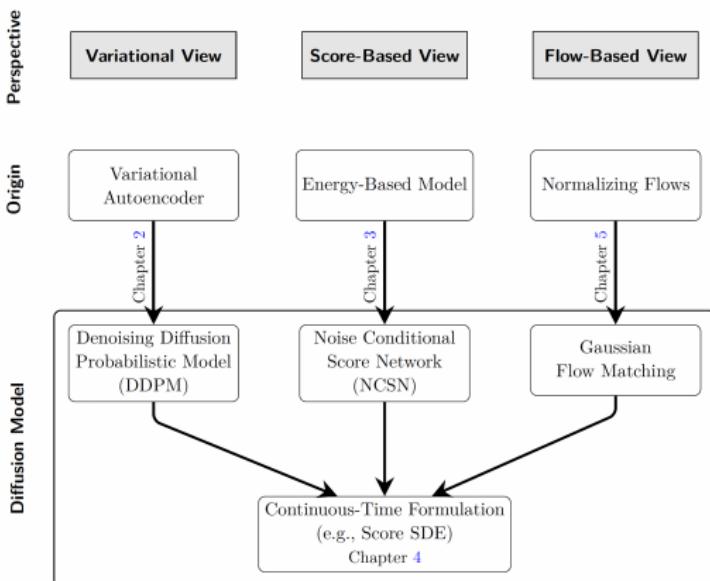
## ③ Flow-based view

- Normalizing Flow (NF) [Rezende and Mohamed, 2015], Neural ODE (NODE) [Chen et al., 2018], Flow-matching (FM) [Lipman et al., 2022]

# First 6 Chapters

Chapter 1

## Overview of Deep Generative Modeling



## **Unifying Principles**

Chapter 6

- Conditional Strategy
  - Fokker-Planck Equation

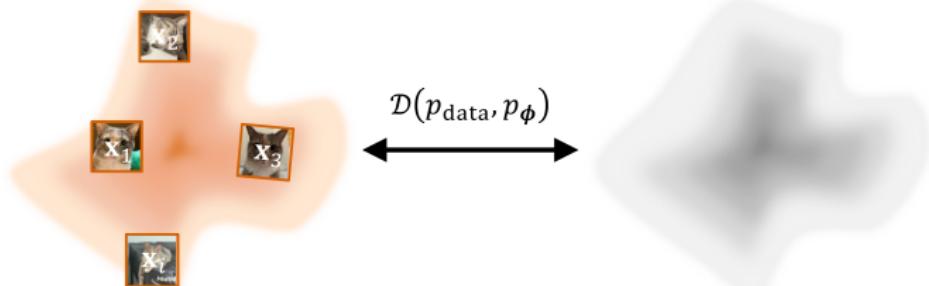
# Roadmap of Our Seminar

- ① Introduction to DGMs ([Lecture 01](#), this lecture)
- ② Core Perspectives on Diffusion Models
  - ① Variational View ([Lecture 02](#))
  - ② Score-Based View ([Lecture 03](#), Yu Xiang)
  - ③ Diffusion Models Today ([Lecture 04](#), Zhang Hanxing)
  - ④ Flow-Based View ([Lecture 05](#), Yuxiang Wang)
- ③ Unified and Systematic Lens on Diffusion Models ([Lecture 06](#), Tang Jiyang)
- ④ Diffusion Models and Optimal Transport ([Lecture 07](#), Jieran Lu)
- ⑤ Guidance and Controllable Generation ([Lecture 08](#), Li Mingze, Kailin Han)
- ⑥ Sophisticated Solvers for Fast Sampling ([Lecture 09](#), Huang binbin, Chenghao Yang)
- ⑦ Distillation-Based Methods for Fast Sampling ([Lecture 10](#), Zhou Xinyu)
- ⑧ Learning Fast Generators from Scratch ([Lecture 11](#))

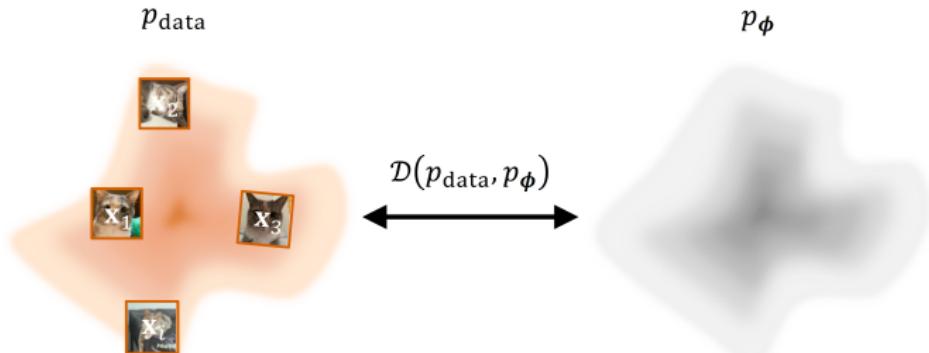
## Deep Generative Modeling

$p_{\text{data}}$

$$p_\phi$$



# Deep Generative Modeling



- Deep generative models (DGMs) are neural networks (NNs) that learn a probability distribution over high-dimensional data (e.g., images, text, audio) so they can generate new examples that resemble the dataset.
- Training a DGM is essentially minimizing the discrepancy between the model distribution  $p_{\phi}$  and the unknown data distribution  $p_{\text{data}}$ .

# Deep Generative Modeling

The goal of learning DGMs are twofold:

- **Realistic generation:** to generate novel, realistic samples indistinguishable from real data.
- **Controllable generation:** to enable fine-grained and interpretable control over the generative process.

# Deep Generative Modeling

The goal of learning DGMs are twofold:

- **Realistic generation:** to generate novel, realistic samples indistinguishable from real data.
- **Controllable generation:** to enable fine-grained and interpretable control over the generative process.

Assumptions of  $p_\phi$  and  $p_{\text{data}}$ :

- Take a collection of examples (e.g., images, text, audio) drawn from  $p_{\text{data}}$  as input: A finite set of samples  $D_{\text{tr}} := \{\mathbf{x}^{(i)}\}_{i=1}^N \stackrel{\text{i.i.d.}}{\sim} p_{\text{data}}$
- The true distribution  $p_{\text{data}}$  is an unknown and complex
- The learned model  $p_\phi$  is commonly called **generative model**

## Deep Generative Modeling

- $p_\phi$ : the probability distribution where  $\phi$  are learnable NN parameters
  - $p_{\text{data}}$ : the true data distribution we want to approximate
  - $\mathbf{x} \in \mathbb{R}^d$ : a data point (an image, a sequence of text, an audio)
  - $\mathcal{D}(p, q)$ : measure the distance between two distributions



# Deep Generative Modeling

- $p_\phi$ : the probability distribution where  $\phi$  are learnable NN parameters
  - $p_{\text{data}}$ : the true data distribution we want to approximate
  - $\mathbf{x} \in \mathbb{R}^d$ : a data point (an image, a sequence of text, an audio)
  - $\mathcal{D}(p, q)$ : measure the distance between two distributions

Given training samples  $D_{\text{tr}}^T$ , we fit  $\phi$  by minimizing  $\mathcal{D}(p_{\text{data}}, p_\phi)$

$$\phi^* \in \arg \min_{\phi} \mathcal{D}(p_{\text{data}}(\mathbf{x}), p_{\phi}(\mathbf{x})). \quad (1.1)$$

Hopefully, we have

$$p_{\phi^*}(\mathbf{x}) \approx p_{\text{data}}(\mathbf{x}). \quad (1.2)$$

Once  $p_{\phi^*}$  is available, we can

- ① Generate  $\mathbf{x}$ : draw  $\mathbf{x} \sim p_\phi$  via the model's sampling methods
  - ② Evaluate  $p_\phi(\mathbf{x}')$ : compute the probability or likelihood of any sample  $\mathbf{x}'$

## What is $\mathcal{D}$ ?

## Choices of $\mathcal{D}$ - KL divergence

A standard choice is the (forward) Kullback-Leibler (KL) divergence

$$\begin{aligned} \mathcal{D}_{\text{KL}}(p_{\text{data}} \| p_{\phi}) &:= \int p_{\text{data}}(\mathbf{x}) \cdot \log \frac{p_{\text{data}}(\mathbf{x})}{p_{\phi}(\mathbf{x})} d\mathbf{x} \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\text{data}}(\mathbf{x}) - \log p_{\phi}(\mathbf{x})]. \end{aligned} \quad (1.3)$$

# Choices of $\mathcal{D}$ - KL divergence

A standard choice is the (forward) Kullback-Leibler (KL) divergence

$$\begin{aligned}\mathcal{D}_{\text{KL}}(p_{\text{data}} \| p_{\phi}) &:= \int p_{\text{data}}(\mathbf{x}) \cdot \log \frac{p_{\text{data}}(\mathbf{x})}{p_{\phi}(\mathbf{x})} d\mathbf{x} \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\text{data}}(\mathbf{x}) - \log p_{\phi}(\mathbf{x})].\end{aligned}\quad (1.3)$$

- ➊  $\mathcal{D}_{\text{KL}}(p \| q) \geq 0$  and  $\mathcal{D}_{\text{KL}}(p \| q) = 0$  iff  $p = q$
- ➋  $\mathcal{D}_{\text{KL}}(p, q)$  is asymmetric:  $\mathcal{D}_{\text{KL}}(p_{\text{data}} \| p_{\phi}) \neq \mathcal{D}_{\text{KL}}(p_{\phi} \| p_{\text{data}})$
- ➌  $\mathcal{D}_{\text{KL}}$  is called **relative entropy** in information theory, where it measures how much information is lost when  $p_{\phi}$  is used to approximate  $p_{\text{data}}$
- ➍ Minimizing  $\mathcal{D}_{\text{KL}}(p_{\text{data}} \| p_{\phi})$  encourages **mode covering**:  
If there exists a set of positive measure  $A$  with  $p_{\text{data}}(A) > 0$  but  $p_{\phi}(\mathbf{x}) = 0$  for  $\mathbf{x} \in A$ , then the integrand contains  $\log(p_{\text{data}}(\mathbf{x})/0) = +\infty$  on  $A$ , so  $\mathcal{D}_{\text{KL}} = +\infty$ . Thus, minimizing KL forces the model to assign probability wherever the data has support.

# Choices of $\mathcal{D}$ - KL divergence

- ⑤ Minimizing KL = Minimizing cross-entropy

$$\begin{aligned}\mathcal{D}_{\text{KL}}(p_{\text{data}} \| p_{\phi}) &= \underbrace{-\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\phi}(\mathbf{x})]}_{\mathcal{H}(p_{\text{data}}, p_{\phi})} - \underbrace{(-\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\text{data}}(\mathbf{x})])}_{\mathcal{H}(p_{\text{data}})} \\ &= \mathcal{H}(p_{\text{data}}, p_{\phi}) - \mathcal{H}(p_{\text{data}})\end{aligned}$$

- $\mathcal{H}(p_{\text{data}})$  is the **entropy** of  $p_{\text{data}}$
- $\mathcal{H}(p_{\text{data}}, p_{\phi})$  is the **cross-entropy** of  $p_{\text{data}}$  relative to  $p_{\phi}$

# Choices of $\mathcal{D}$ - KL divergence

- ⑤ Minimizing KL = Minimizing cross-entropy

$$\begin{aligned}\mathcal{D}_{\text{KL}}(p_{\text{data}} \| p_{\phi}) &= \underbrace{-\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\phi}(\mathbf{x})]}_{\mathcal{H}(p_{\text{data}}, p_{\phi})} - \underbrace{(-\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\text{data}}(\mathbf{x})])}_{\mathcal{H}(p_{\text{data}})} \\ &= \mathcal{H}(p_{\text{data}}, p_{\phi}) - \mathcal{H}(p_{\text{data}})\end{aligned}$$

- $\mathcal{H}(p_{\text{data}})$  is the **entropy** of  $p_{\text{data}}$
- $\mathcal{H}(p_{\text{data}}, p_{\phi})$  is the **cross-entropy** of  $p_{\text{data}}$  relative to  $p_{\phi}$

- ⑥ Minimizing KL  $\Leftrightarrow$  Maximum likelihood estimation (MLE)

$$\min_{\phi} \mathcal{D}_{\text{KL}}(p_{\text{data}} \| p_{\phi}) \Leftrightarrow \max_{\phi} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\phi}(\mathbf{x})]$$

Given  $D_{tr}^N := \{\mathbf{x}^{(i)}\}_{i=1}^N \stackrel{\text{i.i.d.}}{\sim} p_{\text{data}}$ , yielding the empirical MLE

$$\phi^* \in \arg \min_{\phi} \left\{ \hat{\mathcal{L}}_{\text{MLE}}(\phi) := -\frac{1}{N} \sum_{i=1}^N \log p_{\phi}(\mathbf{x}^{(i)}) \right\}.$$

# Choices of $\mathcal{D}$ - Fisher Divergence

**Fisher Divergence:** another important concept for (score-based) diffusion modeling. For two distributions  $p_{\text{data}}$  and  $p_\phi$ , it is defined as

$$\mathcal{D}_F(p_{\text{data}} \| p_\phi) := \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[ \|\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_\phi(\mathbf{x})\|_2^2 \right] \quad (1.4)$$

---

<sup>1</sup>A **vector field** is an assignment of a vector to each point in  $\mathbb{R}^d$ .

# Choices of $\mathcal{D}$ - Fisher Divergence

**Fisher Divergence:** another important concept for (score-based) diffusion modeling. For two distributions  $p_{\text{data}}$  and  $p_\phi$ , it is defined as

$$\mathcal{D}_F(p_{\text{data}} \| p_\phi) := \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[ \|\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_\phi(\mathbf{x})\|_2^2 \right] \quad (1.4)$$

- ①  $\mathcal{D}_F(p_{\text{data}} \| p_\phi) \geq 0$ ,  $\mathcal{D}_F(p_{\text{data}} \| p_\phi) = 0$  if and only if  $p_{\text{data}} = p_\phi$ .
- ② The Fisher divergence  $\mathcal{D}_F$  measures the discrepancy between the **score functions**  $\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$  and  $\nabla_{\mathbf{x}} \log q_\phi(\mathbf{x})$ , which are vector fields pointing toward regions of higher probability.<sup>1</sup> (**score matching**)
- ③ It is invariant to **normalization constants**, since scores depend only on gradients of log-densities, and it forms the basis of score matching.
- ④ The model  $p_\phi$  is trained to align its score field with that of the data.

---

<sup>1</sup>A **vector field** is an assignment of a vector to each point in  $\mathbb{R}^d$ .

# Choices of $\mathcal{D}$ - $f$ divergence

Different divergences capture different geometric or statistical notions of discrepancy, which in turn affect the optimization dynamics of learning algorithms. A broad family is the  $f$ -divergences:

$$\mathcal{D}_f(p\|q) = \int q(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) d\mathbf{x}, \quad f(1) = 0, \quad (1.5)$$

where  $f : \mathbb{R}_+ \rightarrow \mathbb{R}$  is a convex function. Examples of  $f$ -divergence:

- ① When  $f(u) = u \log u$ , then  $\mathcal{D}_f = \mathcal{D}_{\text{KL}}$

$$f(u) = u \log u \quad \Rightarrow \quad \mathcal{D}_f = \mathcal{D}_{\text{KL}}(p\|q)$$

- ② Jensen-Shannon (JS) divergence (extended version see [Nielsen \[2025\]](#))

$$f(u) = \frac{1}{2} \left[ u \log u - (u + 1) \log \frac{1+u}{2} \right] \quad \Rightarrow \quad \mathcal{D}_f = \mathcal{D}_{\text{JS}}(p\|q)$$

# Practical choices of $\mathcal{D}$ - $f$ divergence

- ③ Total Variation distance (see [Tao et al., 2024])

$$f(u) = \frac{1}{2}|u - 1| \quad \Rightarrow \quad \mathcal{D}_f = \mathcal{D}_{\text{TV}}(p, q) = \frac{1}{2} \int |p(\mathbf{x}) - q(\mathbf{x})| d\mathbf{x}$$

Properties of these  $f$ -divergences

- $\mathcal{D}_{\text{JS}}$  is symmetric and its explicit form is

$$\mathcal{D}_{\text{JS}}(p\|q) = \frac{1}{2}\mathcal{D}_{\text{KL}}\left(p\|\frac{1}{2}(p+q)\right) + \frac{1}{2}\mathcal{D}_{\text{KL}}\left(q\|\frac{1}{2}(p+q)\right).$$

- $\mathcal{D}_{\text{JS}}$  is theoretically related to GAN (see Goodfellow et al. [2014])
- The explicit form of  $\mathcal{D}_{\text{TV}}$  is

$$\mathcal{D}_{\text{TV}}(p, q) = \frac{1}{2} \int_{\mathbb{R}^D} |p(\mathbf{x}) - q(\mathbf{x})| d\mathbf{x} = \sup_{A \subset \mathbb{R}^D} |p(A) - q(A)|.$$

# Challenges in modeling distributions

To model  $p_{\text{data}}$ , we can parameterize  $p_{\text{data}}$  using NNs with parameters  $\phi$ , creating a model we denote as  $p_\phi$ . For  $p_\phi$  to be a valid probability density function, it must satisfy two fundamental properties:

- ① **Non-negativity:**  $p_\phi(\mathbf{x}) \geq 0$  for all  $\mathbf{x}$  in the domain.
- ② **Normalization:**  $p_\phi$  is normalized in domain, i.e.,  $\int p_\phi(\mathbf{x}) d\mathbf{x} = 1$ .

# Challenges in modeling distributions

To model  $p_{\text{data}}$ , we can parameterize  $p_{\text{data}}$  using NNs with parameters  $\phi$ , creating a model we denote as  $p_\phi$ . For  $p_\phi$  to be a valid probability density function, it must satisfy two fundamental properties:

- ① **Non-negativity:**  $p_\phi(\mathbf{x}) \geq 0$  for all  $\mathbf{x}$  in the domain.
- ② **Normalization:**  $p_\phi$  is normalized in domain, i.e.,  $\int p_\phi(\mathbf{x}) d\mathbf{x} = 1$ .

A network can naturally produce a real scalar  $E_\phi(\mathbf{x}) \in \mathbb{R}$  for input  $\mathbf{x}$ . To interpret this output as a valid density, it must be transformed to satisfy the above two conditions. A practical alternative is to view  $E_\phi : \mathbb{R}^D \rightarrow \mathbb{R}$  as defining an unnormalized density and then enforce these properties explicitly.

# Challenges in modeling distributions

**Step 1: Ensuring Non-Negativity.** By applying a positive function to  $E_\phi(\mathbf{x})$ , such as  $|E_\phi(\mathbf{x})|$ ,  $E_\phi^2(\mathbf{x})$ , or  $\exp(E_\phi(\mathbf{x}))$ . A standard and convenient choice is the exponential function. It gives us  $\tilde{p}_\phi(\mathbf{x})$ :

$$\tilde{p}_\phi(\mathbf{x}) = \exp(E_\phi(\mathbf{x}))$$

# Challenges in modeling distributions

**Step 1: Ensuring Non-Negativity.** By applying a positive function to  $E_\phi(\mathbf{x})$ , such as  $|E_\phi(\mathbf{x})|$ ,  $E_\phi^2(\mathbf{x})$ , or  $\exp(E_\phi(\mathbf{x}))$ . A standard and convenient choice is the exponential function. It gives us  $\tilde{p}_\phi(\mathbf{x})$ :

$$\tilde{p}_\phi(\mathbf{x}) = \exp(E_\phi(\mathbf{x}))$$

**Step 2: Enforcing Normalization.** To make  $\tilde{p}_\phi(\mathbf{x})$  a valid probability density, we must divide it by its integral over the entire space as follows:

$$p_\phi(\mathbf{x}) = \frac{\tilde{p}_\phi(\mathbf{x})}{\int \tilde{p}_\phi(\mathbf{x}') d\mathbf{x}'} = \frac{\exp(E_\phi(\mathbf{x}))}{Z(\phi) := \int \exp(E_\phi(\mathbf{x}')) d\mathbf{x}'},$$

where  $Z(\phi)$  is known as the **normalizing constant** or **partition function**

$$Z(\phi) := \int \exp(E_\phi(\mathbf{x}')) d\mathbf{x}'.$$

# Challenges in modeling distributions

**Step 1: Ensuring Non-Negativity.** By applying a positive function to  $E_\phi(\mathbf{x})$ , such as  $|E_\phi(\mathbf{x})|$ ,  $E_\phi^2(\mathbf{x})$ , or  $\exp(E_\phi(\mathbf{x}))$ . A standard and convenient choice is the exponential function. It gives us  $\tilde{p}_\phi(\mathbf{x})$ :

$$\tilde{p}_\phi(\mathbf{x}) = \exp(E_\phi(\mathbf{x}))$$

**Step 2: Enforcing Normalization.** To make  $\tilde{p}_\phi(\mathbf{x})$  a valid probability density, we must divide it by its integral over the entire space as follows:

$$p_\phi(\mathbf{x}) = \frac{\tilde{p}_\phi(\mathbf{x})}{\int \tilde{p}_\phi(\mathbf{x}') d\mathbf{x}'} = \frac{\exp(E_\phi(\mathbf{x}))}{Z(\phi) := \int \exp(E_\phi(\mathbf{x}')) d\mathbf{x}'},$$

where  $Z(\phi)$  is known as the **normalizing constant** or **partition function**

$$Z(\phi) := \int \exp(E_\phi(\mathbf{x}')) d\mathbf{x}'.$$

**The computation of  $Z(\phi)$  is intractable.** This is a central problem that motivates the development of many different families of DGMs.

# Example 1: Energy-Based Models (EBMs)

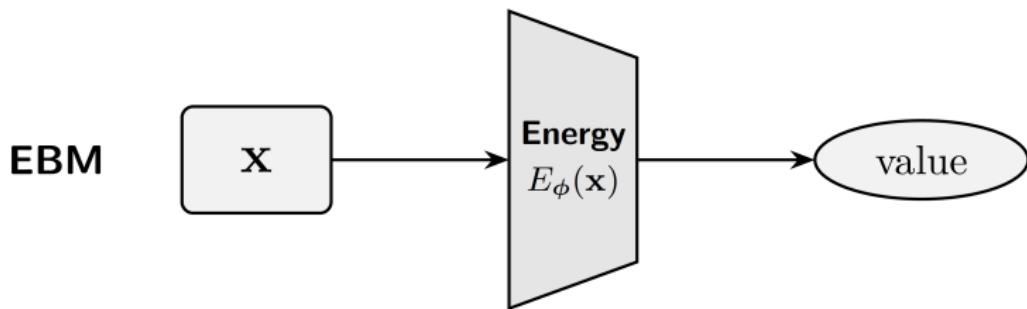


Figure 1.1: EBM maps an input  $x$  to a scalar energy

## Example 1: Energy-Based Models (EBMs)

EBMs define  $p_\phi(\mathbf{x})$  through an energy function  $E_\phi(\mathbf{x})$  that *assigns lower energy to more probable data points*. The probability  $p_\phi(\mathbf{x})$  is defined as:

$$p_\phi(\mathbf{x}) := \frac{\exp(-E_\phi(\mathbf{x}))}{Z(\phi)}, \text{ where } Z(\phi) = \int \exp(-E_\phi(\mathbf{x})) d\mathbf{x}. \quad (1.6)$$

## Example 1: Energy-Based Models (EBMs)

EBMs define  $p_\phi(x)$  through an energy function  $E_\phi(x)$  that *assigns lower energy to more probable data points*. The probability  $p_\phi(x)$  is defined as:

$$p_\phi(x) := \frac{\exp(-E_\phi(x))}{Z(\phi)}, \text{ where } Z(\phi) = \int \exp(-E_\phi(x)) dx. \quad (1.6)$$

- ① Ideas [Ackley et al., 1985, LeCun et al., 2006]
- ②  $p_\phi(x)$  is called Boltzmann distribution
- ③  $Z(\phi)$  is the partition function and computationally intractable
- ④ To avoid the calculation of  $Z(\phi)$ , diffusion models offer an alternative by generating data from the gradient of the log density, which does not depend on  $Z(\phi)$ . (More details in [Lecture 03](#))
- ⑤ Training EBMs usually involves maximizing the MLE of the data

# Example 1: Energy-Based Models (EBMs)

How to train your energy-based models [Song and Kingma, 2021]?

- The standard for learning  $p_\phi$  from i.i.d. data is MLE. We fit  $p_\phi(\mathbf{x})$  to  $p_{\text{data}}(\mathbf{x})$  by maximizing the expected log-likelihood function over  $p_{\text{data}}$ , defined by

$$\arg \max_{\phi} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log p_\phi(\mathbf{x})] \Leftrightarrow$$

$$\arg \min_{\phi} -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log p_\phi(\mathbf{x})] \Leftrightarrow$$

$$\arg \min_{\phi} \left\{ \mathcal{D}_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \| p_\phi(\mathbf{x})) - \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log p_{\text{data}}(\mathbf{x})] \right\} \Leftrightarrow$$

$$\arg \min_{\phi} \mathcal{D}_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \| p_\phi(\mathbf{x}))$$

- We cannot directly compute the likelihood of an EBM, but we can still estimate the gradient of the log-likelihood with MCMC approaches.

## Example 1: Energy-Based Models (EBMs)

**Why gradient information?** If two continuously differentiable real-valued functions  $f(\mathbf{x})$  and  $g(\mathbf{x})$  have equal first derivatives everywhere, then  $f(\mathbf{x}) \equiv g(\mathbf{x}) + C$ . When  $f(\mathbf{x})$  and  $g(\mathbf{x})$  are log probability density functions (PDFs) with equal first derivatives, and therefore  $f(\mathbf{x}) \equiv g(\mathbf{x})$  (see Claim 2).

## Example 1: Energy-Based Models (EBMs)

**Why gradient information?** If two continuously differentiable real-valued functions  $f(\mathbf{x})$  and  $g(\mathbf{x})$  have equal first derivatives everywhere, then  $f(\mathbf{x}) \equiv g(\mathbf{x}) + C$ . When  $f(\mathbf{x})$  and  $g(\mathbf{x})$  are log probability density functions (PDFs) with equal first derivatives, and therefore  $f(\mathbf{x}) \equiv g(\mathbf{x})$  (see Claim 2).

The gradient of the log-probability of an EBM in Equ. (1.6) decomposes as a sum of two terms:

$$\nabla_{\phi} \log p_{\phi}(\mathbf{x}) = -\nabla_{\phi} E_{\phi}(\mathbf{x}) - \nabla_{\phi} \log Z(\phi) \quad (1.7)$$

- The first term,  $-\nabla_{\phi} E_{\phi}(\mathbf{x})$ , is straightforward to evaluate with NNs.
- The second term,  $\nabla_{\phi} \log Z(\phi)$  is **intractable to compute exactly**.
- Score matching (Lecture 03):

$$\mathcal{L}_{\text{SM}}(\phi) = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \|\nabla_{\mathbf{x}} \log p_{\phi}(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})\|_2^2. \quad (1.8)$$

## Example 2: Autoregressive Models

AR

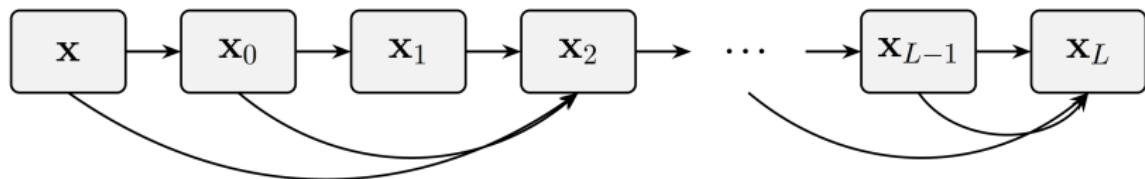


Figure 1.2: AR generates a sequence  $\{x_\ell\}$  left to right with causal dependencies

## Example 2: Autoregressive Models

Deep autoregressive (AR) models factorize the joint data distribution  $p_\phi$  into a product of conditional probabilities using the chain rule of probability:

$$p_\phi(\mathbf{x}) = \prod_{i=1}^D p_\phi(x_i | \mathbf{x}_{<i}), \quad (1.9)$$

where  $\mathbf{x} = (x_1, \dots, x_D)$  and  $\mathbf{x}_{<i} = (x_1, \dots, x_{i-1})$ .

- Each conditional  $p_\phi(x_i | \mathbf{x}_{<i})$  is parameterized by a neural network
- Each  $p_\phi(x_i | \mathbf{x}_{<i})$  is normalized by design (e.g., via softmax for discrete or parameterized Gaussian for continuous variables)
- Training proceeds by maximizing the exact likelihood, or equivalently, minimizing the negative log-likelihood
- A foundational class of likelihood-based generative models and key approaches in modern research

# Example 3: Variational Autoencoders (VAEs)

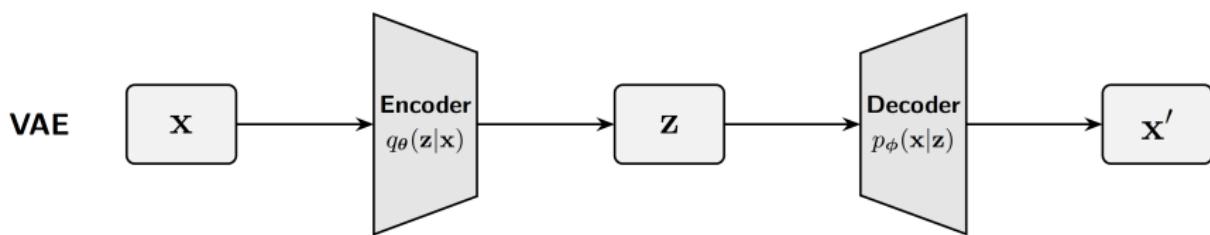


Figure 1.3: VAE encodes  $x$  to a latent  $z$  and decodes to a reconstruction  $x'$

## Example 3: Variational Autoencoders (VAEs)

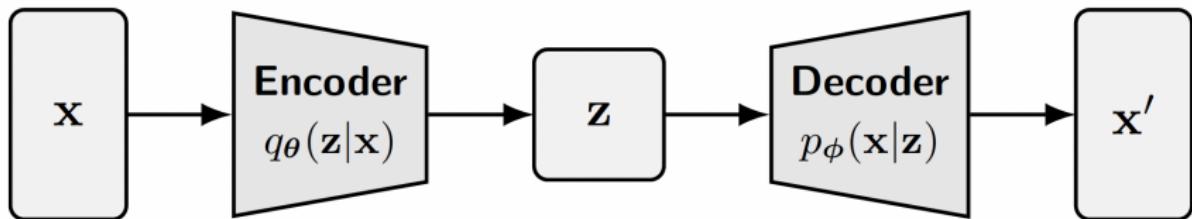


Figure 1.4: VAE consists of a stochastic encoder  $q_\theta(z | x)$  that maps data  $x$  to a latent variable  $z$ , and a decoder  $p_\phi(x | z)$  that reconstructs data from the latent.

## Example 3: Variational Autoencoders (VAEs)

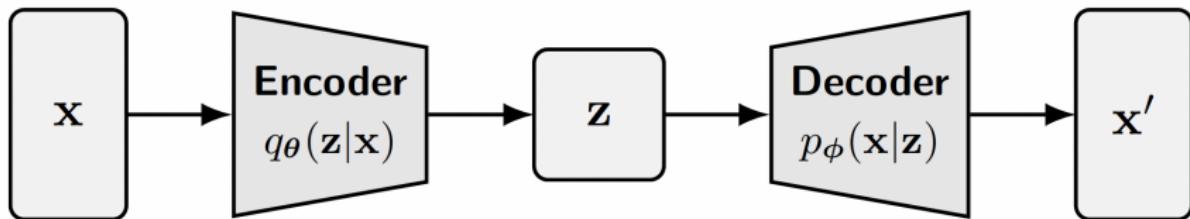


Figure 1.4: VAE consists of a stochastic encoder  $q_{\theta}(z | x)$  that maps data  $x$  to a latent variable  $z$ , and a decoder  $p_{\phi}(x | z)$  that reconstructs data from the latent.

- Ideas [Hinton and Salakhutdinov, 2006, Kingma and Welling, 2013]
- From a probabilistic view, VAEs learn both an **encoder** (inference network),  $q_{\theta}(z | x)$ , which approximates the unknown distribution of latent variables given the data, and a **decoder** (the generator),  $p_{\phi}(x | z)$ , which reconstructs data from these latent variables.

## Example 3: Variational Autoencoders (VAEs)

To define a computable training objective, since we cannot directly optimize  $\log p_\phi(x)$ , we can maximize a lower bound on it, that is, the Evidence Lower Bound (ELBO):

## Example 3: Variational Autoencoders (VAEs)

To define a computable training objective, since we cannot directly optimize  $\log p_\phi(\mathbf{x})$ , we can maximize a lower bound on it, that is, the Evidence Lower Bound (ELBO):

- **Step 1:** The likelihood  $p_\phi(\mathbf{x})$  can be treated as a latent-variable generative model through the marginal likelihood:

$$p_\phi(\mathbf{x}) = \int p_\phi(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z} \quad (1.10)$$

## Example 3: Variational Autoencoders (VAEs)

To define a computable training objective, since we cannot directly optimize  $\log p_\phi(\mathbf{x})$ , we can maximize a lower bound on it, that is, the Evidence Lower Bound (ELBO):

- **Step 1:** The likelihood  $p_\phi(\mathbf{x})$  can be treated as a latent-variable generative model through the marginal likelihood:

$$p_\phi(\mathbf{x}) = \int p_\phi(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z} \quad (1.10)$$

- **Step 2:** To connect our intractable generator to real data, consider the reverse question: given an observation  $\mathbf{x}$ , what latent codes  $\mathbf{z}$  could have produced it? By Bayes' rule, the posterior distribution is

$$p_\phi(\mathbf{z} | \mathbf{x}) = \frac{p_\phi(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{p_\phi(\mathbf{x})} = \frac{p_\phi(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{\int_{\mathbf{z}'} p_\phi(\mathbf{x}, \mathbf{z}') d\mathbf{z}'}. \quad (1.11)$$

The denominator is intractable.

## Example 3: Variational Autoencoders (VAEs)

- Step 3: The *variational* step in VAEs addresses Step 3 by replacing  $p_\phi(z | x)$  with a tractable approximation  $q_\theta(z | x)$ , which is an encoder (or inference network) parameterized by  $\theta$ :

$$q_\theta(z | x) \approx p_\phi(z | x). \quad (1.12)$$

In practice, the encoder maps each observed data point  $x$  to a distribution over latent codes.

## Example 3: Variational Autoencoders (VAEs)

- **Step 3:** The *variational* step in VAEs addresses Step 3 by replacing  $p_\phi(z | x)$  with a tractable approximation  $q_\theta(z | x)$ , which is an encoder (or inference network) parameterized by  $\theta$ :

$$q_\theta(z | x) \approx p_\phi(z | x). \quad (1.12)$$

In practice, the encoder maps each observed data point  $x$  to a distribution over latent codes.

- **Step 4:** The log likelihood can then be treated as marginal of the density function  $p_\phi(x, z)$  over  $z$

$$\begin{aligned} \log p_\phi(x) &= \log \int p_\phi(x, z) dz = \log \int q_\theta(z | x) \frac{p_\phi(x, z)}{q_\theta(z | x)} dz \\ &= \log \mathbb{E}_{z \sim q_\theta(z | x)} \left[ \frac{p_\phi(x, z)}{q_\theta(z | x)} \right] \geq \underbrace{\mathbb{E}_{z \sim q_\theta(z | x)} \left[ \log \frac{p_\phi(x, z)}{q_\theta(z | x)} \right]}_{\text{Evidence Lower Bound (ELBO)}}, \end{aligned}$$

where the inequality is due to Jensen's inequality in Claim 3.

# Example 3: Variational Autoencoders (VAEs)

## Theorem 1 (Evidence Lower Bound (ELBO))

For any data point  $x$ , the log-likelihood satisfies:

$$\log p_\phi(x) \geq \mathcal{L}_{\text{ELBO}}(x; \theta, \phi), \quad (1.13)$$

where the ELBO is given by

$$\mathcal{L}_{\text{ELBO}}(x; \theta, \phi) = \underbrace{\mathbb{E}_{z \sim q_\theta(z|x)} [\log p_\phi(x | z)]}_{\text{Reconstruction Term}} - \underbrace{\mathcal{D}_{\text{KL}}(q_\theta(z | x) \| p(z))}_{\text{Latent Regularization}}$$

## Example 3: Variational Autoencoders (VAEs)

### Theorem 1 (Evidence Lower Bound (ELBO))

For any data point  $x$ , the log-likelihood satisfies:

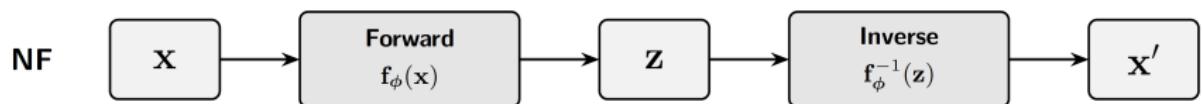
$$\log p_\phi(x) \geq \mathcal{L}_{\text{ELBO}}(x; \theta, \phi), \quad (1.13)$$

where the ELBO is given by

$$\mathcal{L}_{\text{ELBO}}(x; \theta, \phi) = \underbrace{\mathbb{E}_{z \sim q_\theta(z|x)} [\log p_\phi(x | z)]}_{\text{Reconstruction Term}} - \underbrace{\mathcal{D}_{\text{KL}}(q_\theta(z | x) \| p(z))}_{\text{Latent Regularization}}$$

- The first term encourages accurate reconstruction of the data
- The second regularizes the latent variables by keeping them close to a simple prior distribution  $p_{\text{prior}}(z) := p(z)$  (often Gaussian)
- VAEs provide a principled way to combine neural networks with latent variable models and remain one of the most widely used approaches

# Example 4: Normalizing Flows



## Example 4: Normalizing Flows



- NF first applies an invertible map  $f_\phi$  between  $x$  and  $z$
- It then uses  $f_\phi^{-1}$  to produce  $x'$

## Example 4: Normalizing Flows

- Classic flow-based models, such as Normalizing Flows (NFs) [Rezende and Mohamed, 2015] and Neural Ordinary Differential Equations (NODEs) [Chen et al., 2018] aim to learn a bijective mapping  $f_\phi$  between a simple latent distribution  $z$  and a complex data distribution  $x$  via an invertible operator.
- This is achieved either through a sequence of bijective transformations (in NFs) or by modeling the transformation as an Ordinary Differential Equation (in NODEs).
- These models leverage the “change-of-variable formula for densities”, enabling MLE training:

$$\log p_\phi(x) = \log p(z) + \log \left| \det \frac{\partial f_\phi^{-1}(x)}{\partial z} \right| \quad (1.14)$$

# Example 5: Generative Adversarial Networks

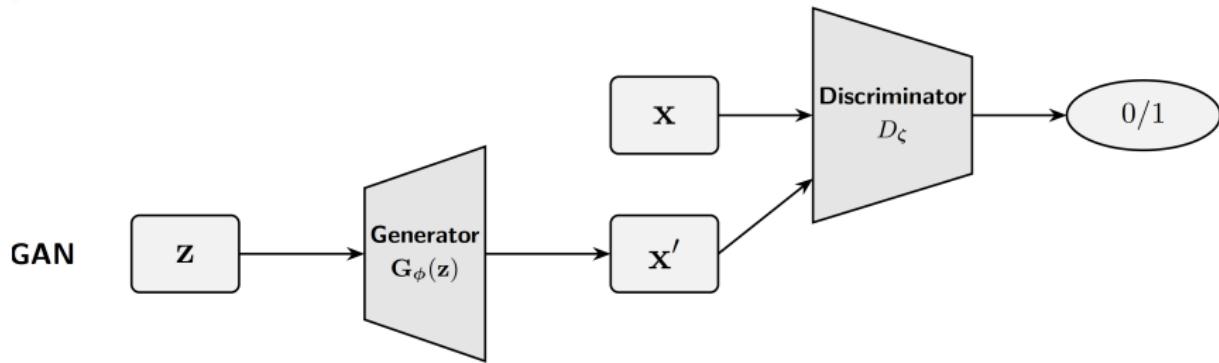


Figure 1.5: GAN transforms noise  $z$  to a sample  $x'$  that is judged against real  $x$  by a discriminator  $D_\zeta$

## Example 5: Generative Adversarial Networks

- Generative Adversarial Networks (GANs) consist of two NNs, a generator  $G_\phi$  and a discriminator  $D_\zeta$ , that compete against each other.
- The generator  $G_\phi$  aims to create realistic samples  $G_\phi(\mathbf{z})$  from random noise  $\mathbf{z} \sim p_{\text{prior}}$ , while the discriminator attempts to distinguish between real samples  $\mathbf{x}$  and generated samples  $G_\phi(\mathbf{z})$ .
- The objective function for GANs can be formulated as:

$$\min_{G_\phi} \max_{D_\zeta} \underbrace{\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D_\zeta(\mathbf{x})]}_{\text{real}} + \underbrace{\mathbb{E}_{\mathbf{z} \sim p_{\text{prior}}(\mathbf{z})} [\log (1 - D_\zeta(G_\phi(\mathbf{z})))]}_{\text{fake}}.$$

- GANs do not define an explicit density function and therefore bypass likelihood estimation entirely. They focus on generating samples that closely mimic the data distribution.

## Example 5: Generative Adversarial Networks

From a divergence perspective, the discriminator  $D_\zeta$  implicitly measures the discrepancy between the true data distribution  $p_{\text{data}}$  and the generator distribution  $p_{G_\phi}$ , where  $p_{G_\phi}$  denotes the distribution of generated samples  $G_\phi(z)$  obtained from noise  $z \sim p_{\text{prior}}$ . With an optimal discriminator for a fixed generator  $G_\phi$  computed as

$$\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_{G_\phi}(x)}, \quad (1.15)$$

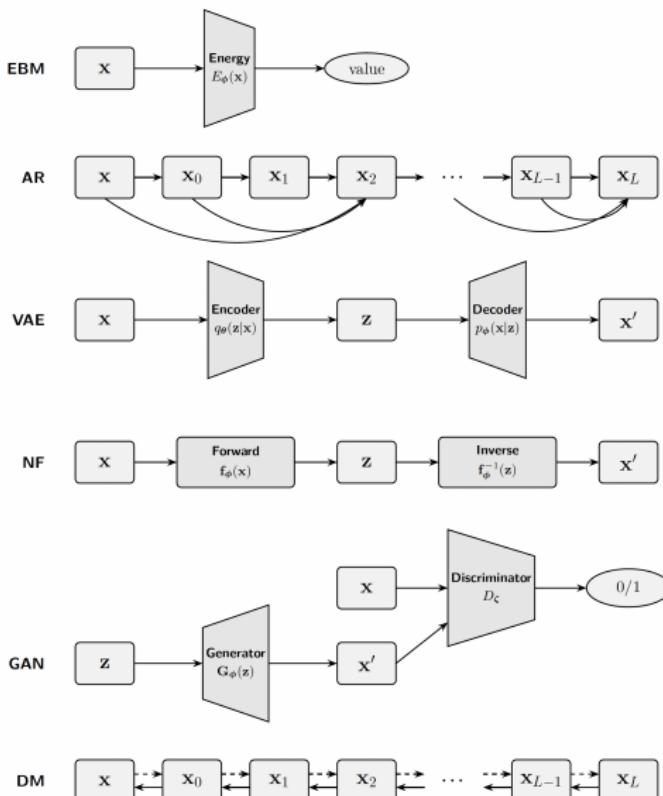
the generator's minimization reduces to

$$\min_{G_\phi} 2\mathcal{D}_{\text{JS}}(p_{\text{data}} \| p_{G_\phi}) - \log 4. \quad (1.16)$$

Here,  $\mathcal{D}_{\text{JS}}$  denotes the Jensen-Shannon divergence, defined as

$$\mathcal{D}_{\text{JS}}(p \| q) := \frac{1}{2}\mathcal{D}_{\text{KL}}\left(p\|\frac{p+q}{2}\right) + \frac{1}{2}\mathcal{D}_{\text{KL}}\left(q\|\frac{p+q}{2}\right). \quad (1.17)$$

# Prominent DGMs



# Quick Summary of DGMs

DGMs span a wide spectrum of modeling strategies. A fundamental distinction lies in how these models parameterize the underlying data distribution, that is, whether they specify  $p_\phi(\mathbf{x})$  explicitly or only implicitly, irrespective of the training objective.

- **Explicit Models:** These models directly parameterize a probability distribution  $p_\phi(\mathbf{x})$  via a tractable or approximately tractable density or mass function. Examples include ARs, NFs, VAEs, and DMs, all of which define  $p_\phi(\mathbf{x})$  either exactly or through a tractable bound.
- **Implicit Models:** These models specify a distribution only through a sampling procedure, typically of the form  $\mathbf{x} = G_\phi(\mathbf{z})$  for some noise variable  $\mathbf{z} \sim p_{\text{prior}}$ . In this case,  $p_\phi(\mathbf{x})$  is not available in closed form and may not be defined at all.

# Closing Remarks

- We begin by defining the primary objective: to learn a tractable model distribution  $p_\phi$  parametrized by  $\phi$  that approximates an unknown, complex data distribution  $p_{\text{data}}$ .
- A central challenge is the computational intractability of the normalizing constant, or partition function  $Z(\phi)$ .
- To circumvent this problem, various families of DGMs have been developed, each employing a distinct strategy. We surveyed several prominent approaches, including EBMs, ARs, VAEs, NFs, and GANs.
- While each of these classical frameworks is significant, three in particular serve as the conceptual origins for the diffusion models that are the focus of this monograph: VAEs, EBMs, and NFs. We will trace the evolution of diffusion models from these three foundational paradigms.

## Claim 1 (Equal derivatives means differences by constant)

If two continuous functions  $f$  and  $g$  have the same derivatives in the same domain, then  $f(\mathbf{x}) = g(\mathbf{x}) + C$  where  $C$  is a constant.

### Proof.

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  and  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  be continuously differentiable and suppose

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \nabla_{\mathbf{x}} g(\mathbf{x}) \quad \text{for all } \mathbf{x}.$$

Define  $h(\mathbf{x}) = f(\mathbf{x}) - g(\mathbf{x})$ . Then

$$\nabla_{\mathbf{x}} h(\mathbf{x}) = \nabla_{\mathbf{x}} f(\mathbf{x}) - \nabla_{\mathbf{x}} g(\mathbf{x}) = 0.$$

That is,  $h(\mathbf{x})$  must be constant. This leads to

$$h(\mathbf{x}) \equiv C \quad \Rightarrow \quad f(\mathbf{x}) = g(\mathbf{x}) + C.$$



## Claim 2 (Equal derivatives mean equal functions)

Suppose  $f$  and  $g$  are log densities:  $p(\mathbf{x}) = e^{f(\mathbf{x})}$ ,  $q(\mathbf{x}) = e^{g(\mathbf{x})}$ , and both are normalized PDFs, i.e.,  $\int e^{f(\mathbf{x})} d\mathbf{x} = 1$ ,  $\int e^{g(\mathbf{x})} d\mathbf{x} = 1$ . If the gradients of  $f$  and  $g$  are equal everywhere, then  $f \equiv g$ .

### Proof.

If the gradients of  $f$  and  $g$  are equal everywhere, by Claim 1,  $f(\mathbf{x}) = g(\mathbf{x}) + C$ . This means  $e^{f(\mathbf{x})} = e^{g(\mathbf{x})+C} = e^C e^{g(\mathbf{x})}$ . Integrate both sides:

$$1 = \int e^{f(\mathbf{x})} d\mathbf{x} = \int e^C e^{g(\mathbf{x})} d\mathbf{x} = e^C \int e^{g(\mathbf{x})} d\mathbf{x} = e^C \cdot 1 = e^C.$$

The only choice is  $C = 0$ . Hence  $f(\mathbf{x}) = g(\mathbf{x})$  for all  $\mathbf{x}$ , and therefore  $p(\mathbf{x}) = q(\mathbf{x})$ . ■

# Jensen's Inequality

Jensen's inequality generalizes the statement that the secant line of a convex function lies above the graph of the function: the secant line consists of weighted means of the convex function (for  $t \in [0, 1]$ ).

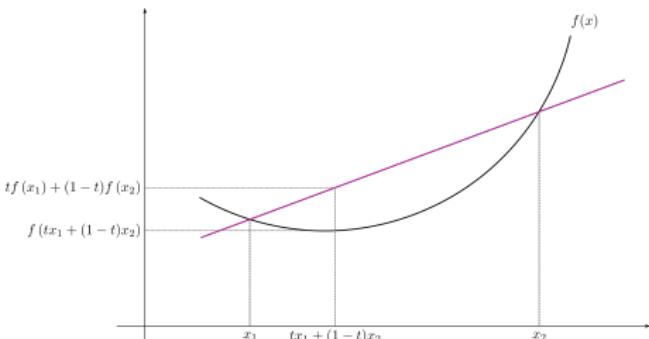


Figure 1.6: A secant line of a convex function lies above its graph.

- Weighted means of the convex function:  $tf(x_1) + (1 - t)f(x_2)$
- The convex function of the weighted means:  $f(tx_1 + (1 - t)x_2)$

In this case, Jensen's inequality provides

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2).$$

# Jensen's Inequality

- **Finite form:** For a real **convex** function  $\varphi$ , numbers  $x_1, x_2, \dots, x_n$  in its domain, and positive weights  $a_i$ , Jensen's inequality can be stated as:

$$\varphi\left(\frac{\sum a_i x_i}{\sum a_i}\right) \leq \frac{\sum a_i \varphi(x_i)}{\sum a_i}$$

and the inequality is reversed if  $\varphi$  is **concave**, which is

$$\varphi\left(\frac{\sum a_i x_i}{\sum a_i}\right) \geq \frac{\sum a_i \varphi(x_i)}{\sum a_i}.$$

Equality holds if and only if  $x_1 = x_2 = \dots = x_n$  or  $\varphi$  is linear on a domain containing  $x_1, x_2, \dots, x_n$ .

- **Example:** The function  $\log(x)$  is concave, we have the following inequality

$$\log\left(\frac{\sum_{i=1}^n x_i}{n}\right) \geq \frac{\sum_{i=1}^n \log(x_i)}{n}$$

# Jensen's Inequality

## Claim 3 (Probabilistic form of Jensen's Inequality)

*In the probability theory setting, let  $(\Omega, \mathcal{F}, P)$  be a probability space,  $X$  an integrable real-valued random variable and  $\varphi$  a convex function. Then*

$$\varphi(\mathbb{E}[X]) \leq \mathbb{E}[\varphi(X)]. \quad (1.18)$$

- Note that the equality holds if and only if  $\varphi$  is a linear function on some convex set  $A$  such that  $P(X \in A) = 1$ .

David H Ackley, Geoffrey E Hinton, and Terrence J Sejnowski. A learning algorithm for boltzmann machines. *Cognitive science*, 9(1):147–169, 1985.

Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018.

Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.

Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *science*, 313(5786):504–507, 2006.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, Fujie Huang, et al. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006.

Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.

Frank Nielsen. Two tales for a geometric jensen–shannon divergence. *arXiv preprint arXiv:2508.05066*, 2025.

Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR, 2015.

Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. pmlr, 2015.

Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.

Yang Song and Diederik P Kingma. How to train your energy-based models. *arXiv preprint arXiv:2101.03288*, 2021.

Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.

Lan Tao, Shirong Xu, Chi-Hua Wang, Namjoon Suh, and Guang Cheng. Discriminative estimation of total variation distance: A fidelity auditor for generative data. *arXiv preprint arXiv:2405.15337*, 2024.