

Embodied Vision

Generative Modeling

Tsung-Wei Ke

Spring 2026



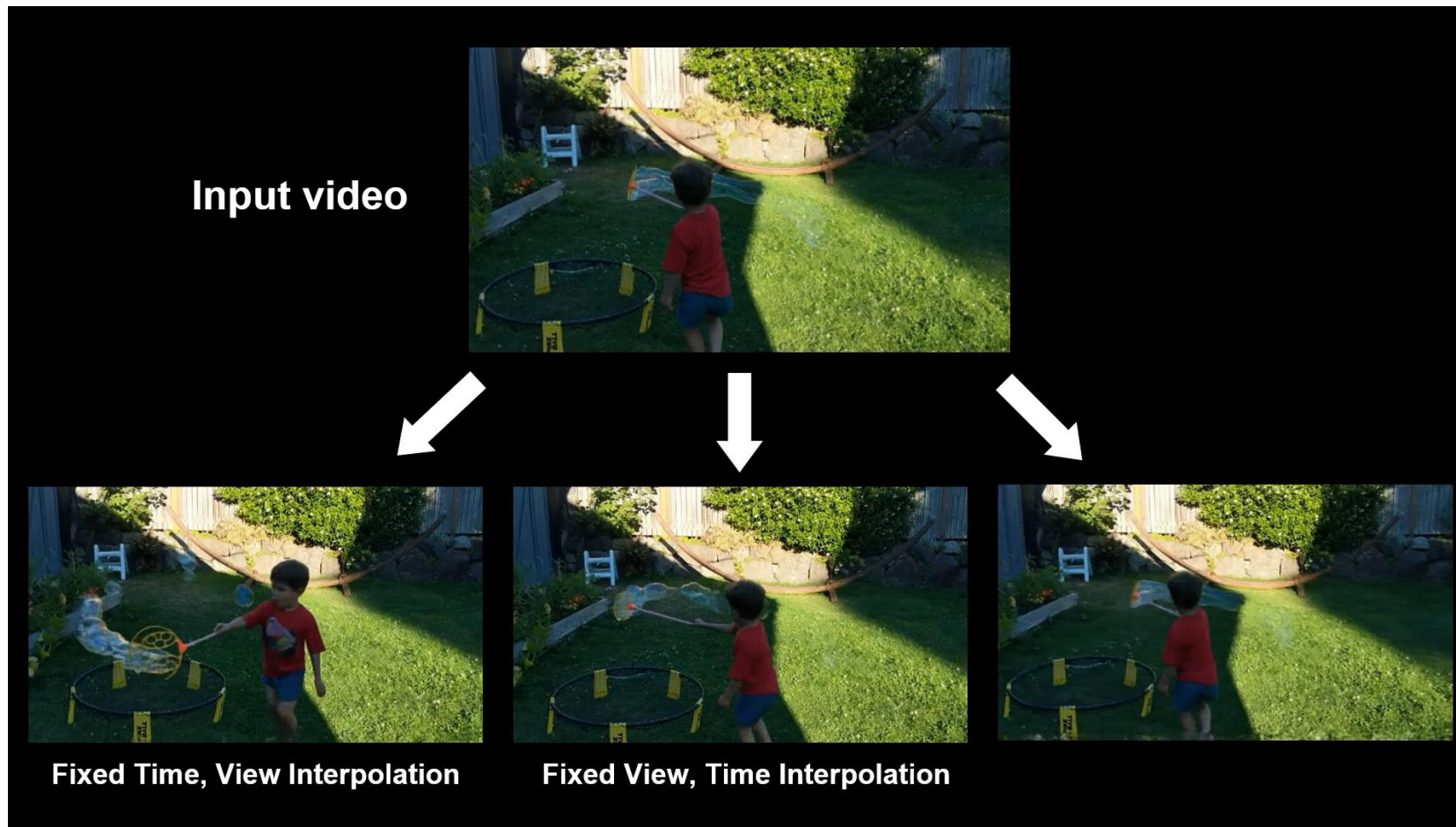
Disclaimer

- This lecture borrows contents heavily from
 - [Deep Generative Models](#) by Kaiming He at MIT

Content

- Generative Modeling
 - Autoregressive Models
 - Variational AutoEncoders
 - Denoising Diffusion Probabilistic Models
 - Flow Matching Models
- 4D Generative Modeling
 - 3D Generation with Multi-view Video Generation
 - 4D Generation with Multi-view Video Generation
 - 3D Geometry Generation

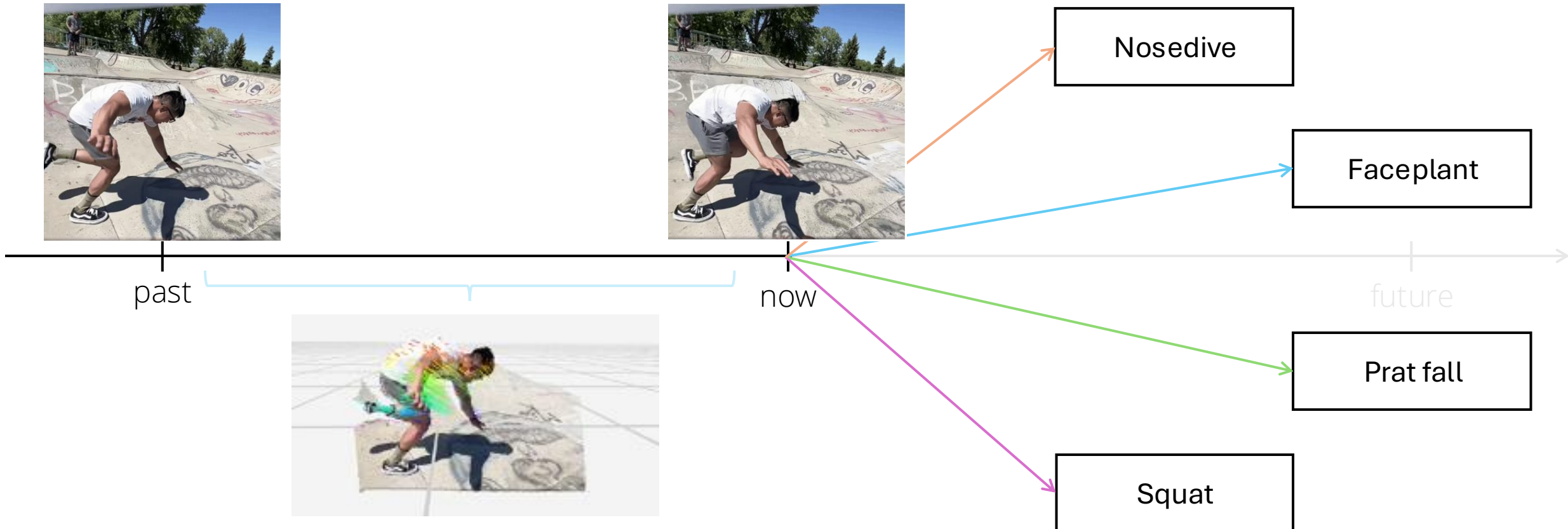
Previously, We Model the 4D World by Observing Monocular Videos



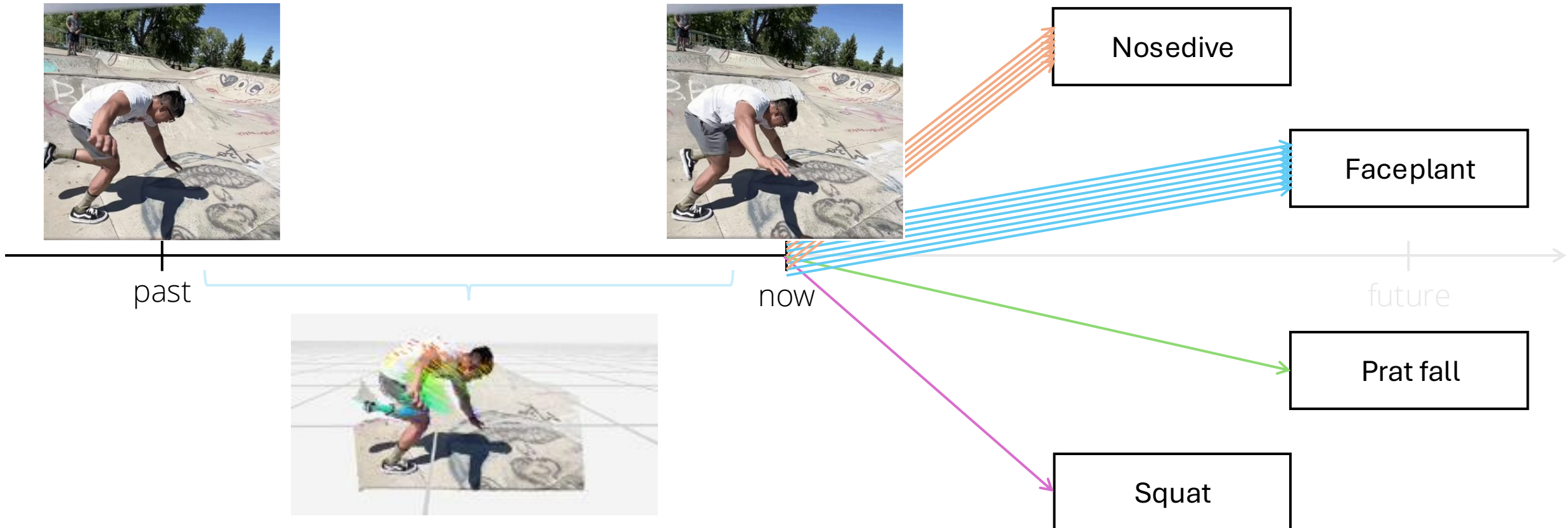
We Inferred the Past Motion of an Observed Object... How About the Future Motion of the Object?



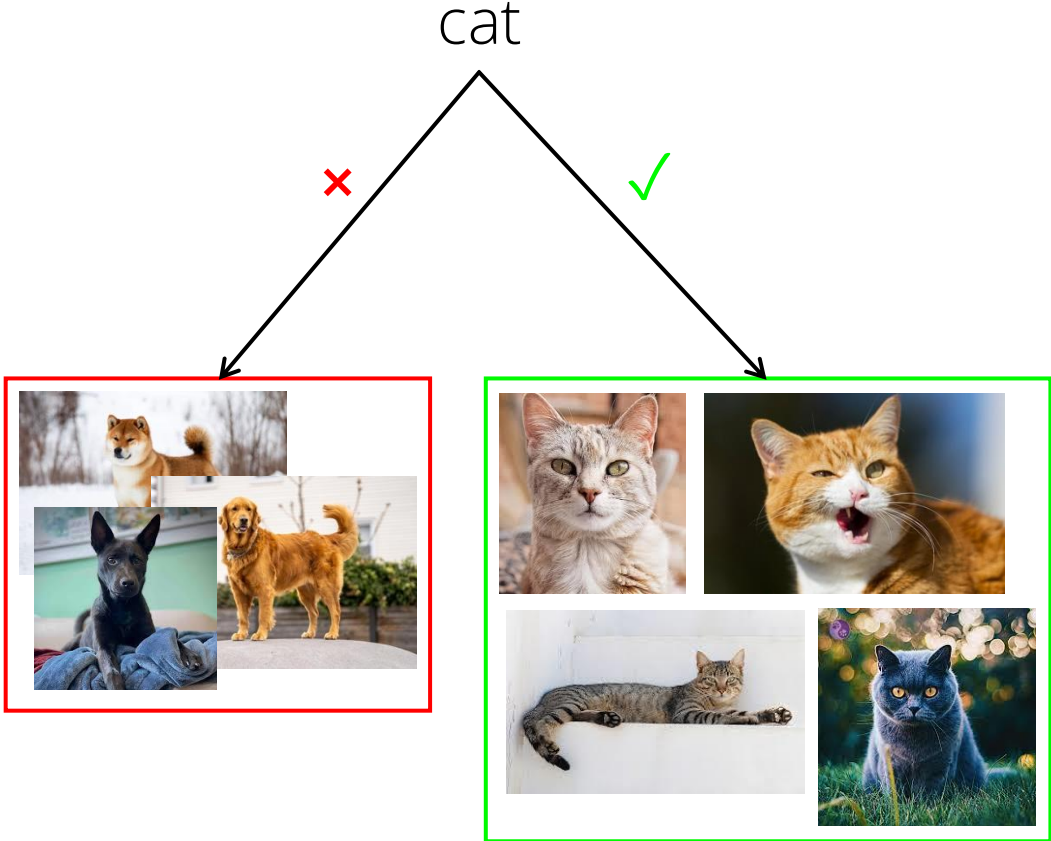
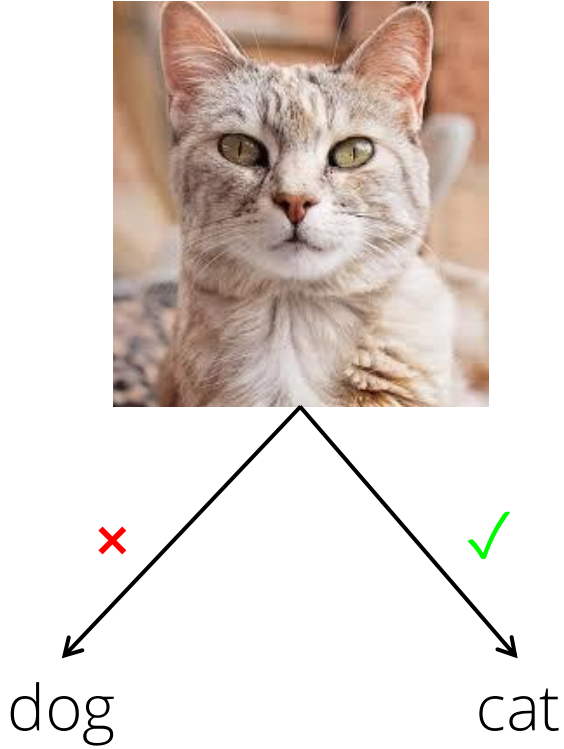
Unfortunately, the Future is Not Unique. We May Have Multiple Plausible Futures



Some Futures are "More Plausible" than Others



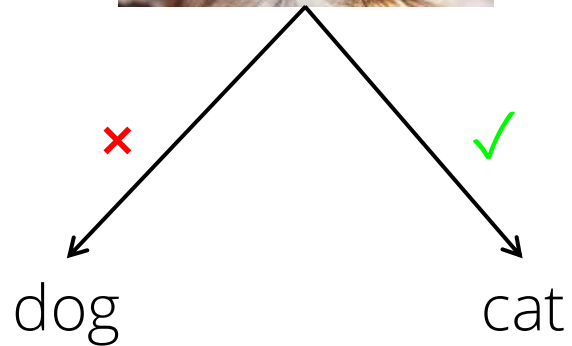
An Analogy of Image Classification and Generation



An Analogy of Image Classification and Generation

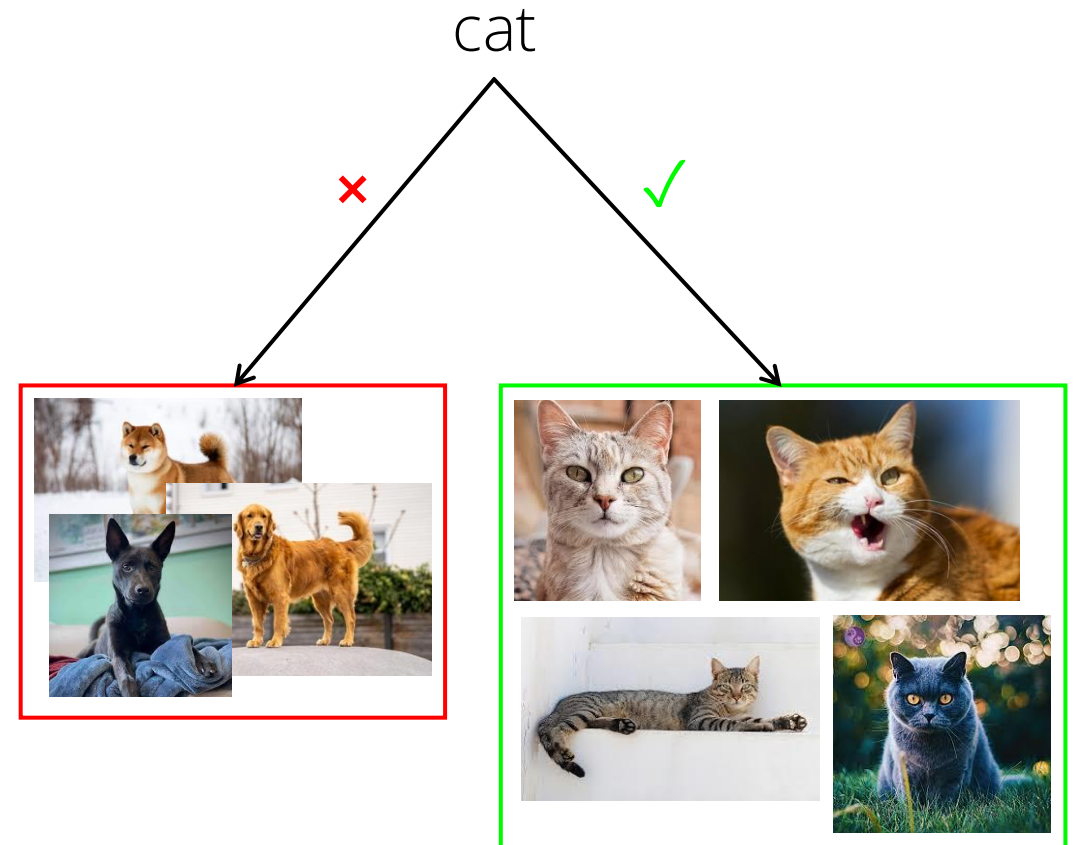
Discriminative

- Sample "x" → Label "y"
- One desired output

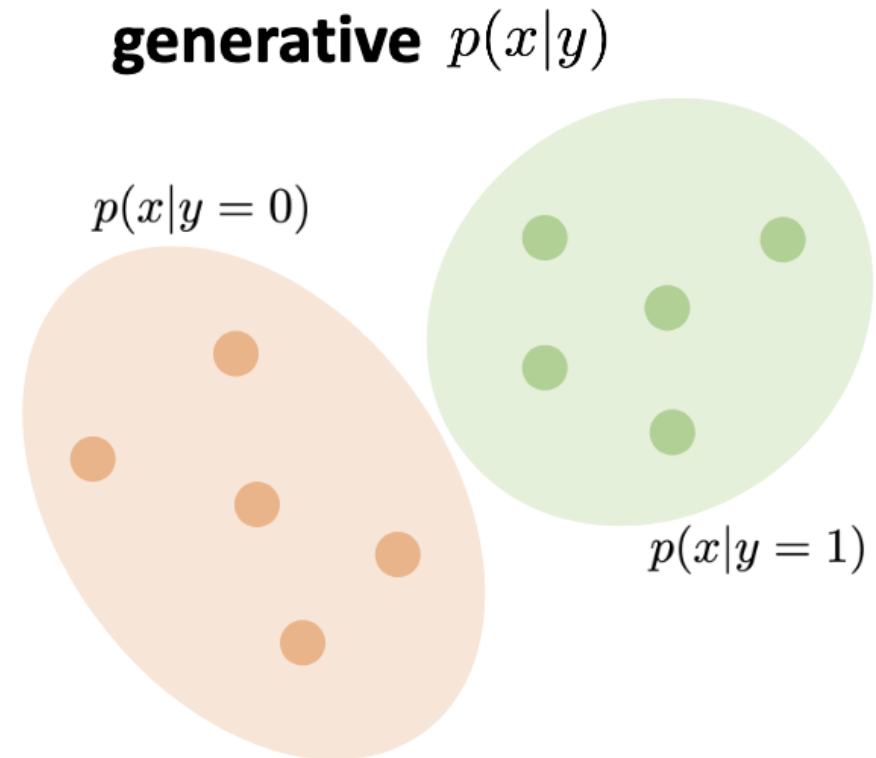
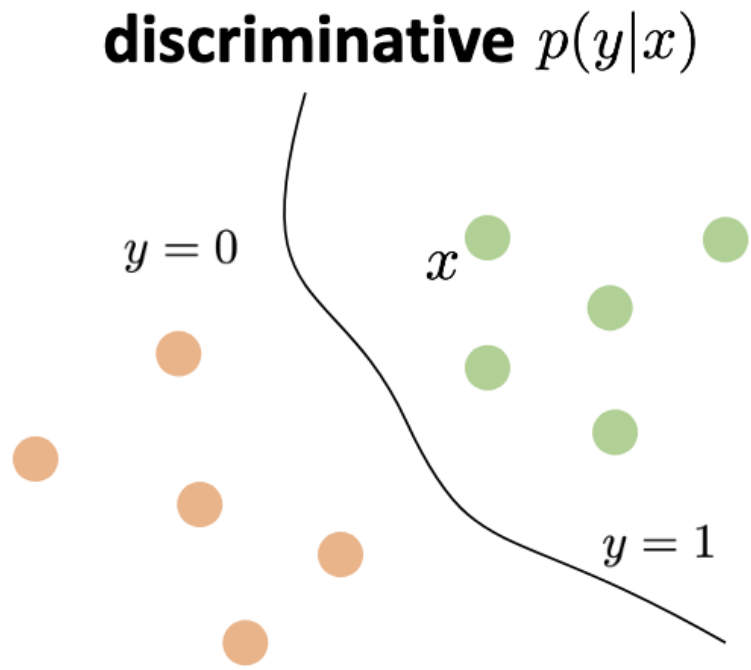


Generative

- Label "y" → Sample "x"
- Multiple possible outputs



Discriminative vs. Generative Models



Discriminative vs. Generative Models

- Discriminative models to be generative: Baye's rule

$$p(y|x) = p(x|y) \frac{p(y)}{p(x)}$$

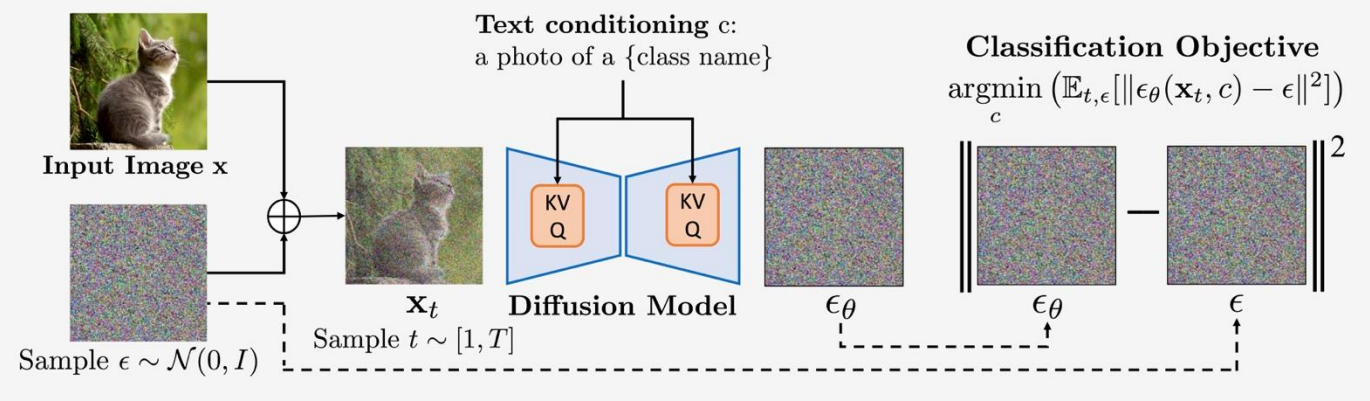
Prior distribution of label "y"
(e.g. uniform distribution)

Prior distribution of sample "x"
(e.g. constant for every "x")

- Generative models to be discriminative:

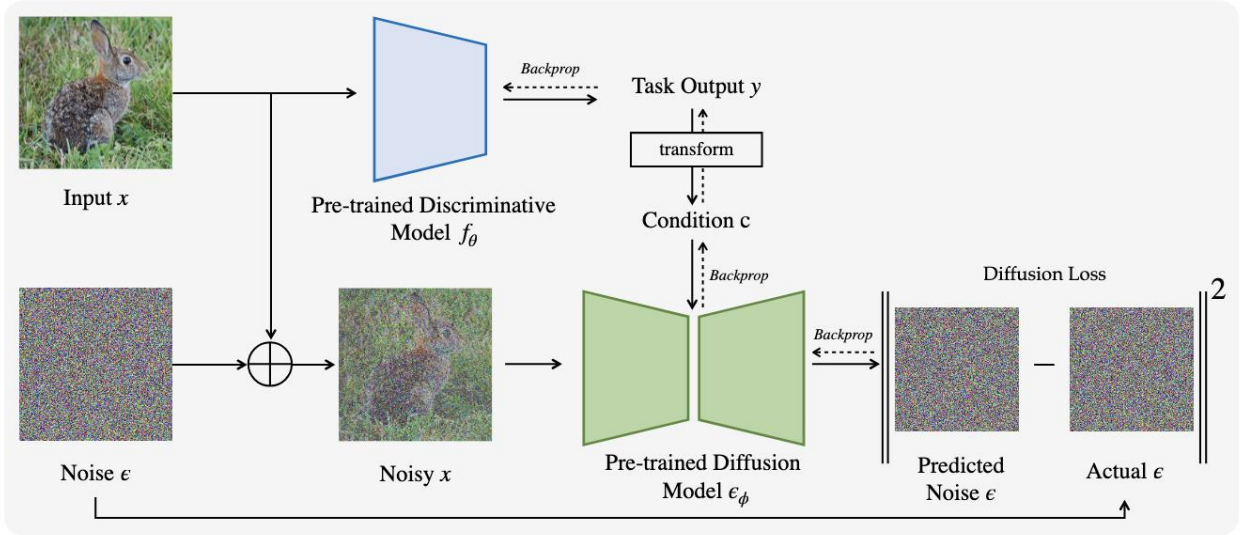
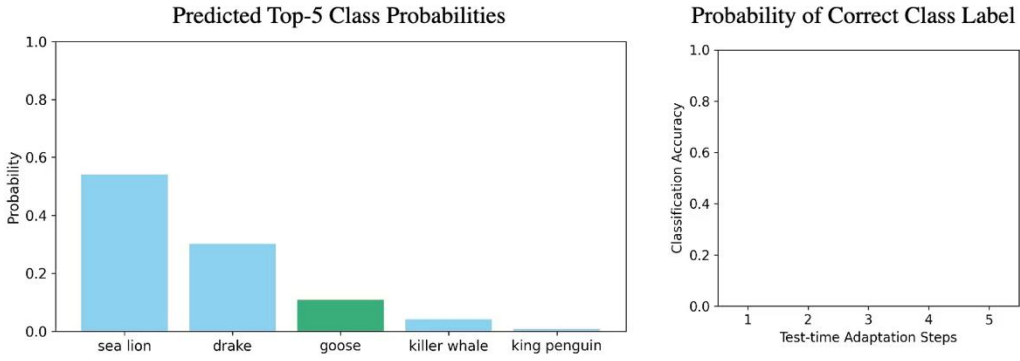
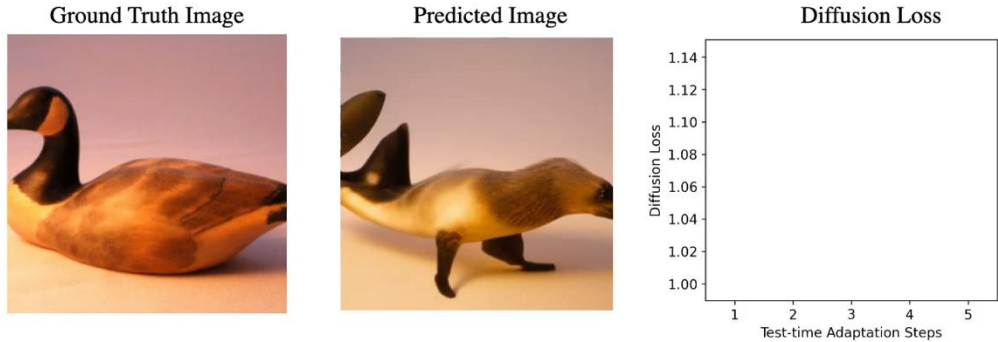
$$p(x|y) = p(y|x) \frac{p(x)}{p(y)}$$

Generative Models are Implicitly Discriminative Models



Your Diffusion Model is Secretly a Zero-Shot Classifier. Li et al.

Improve pre-trained image classifiers by minimizing diffusion loss at test-time



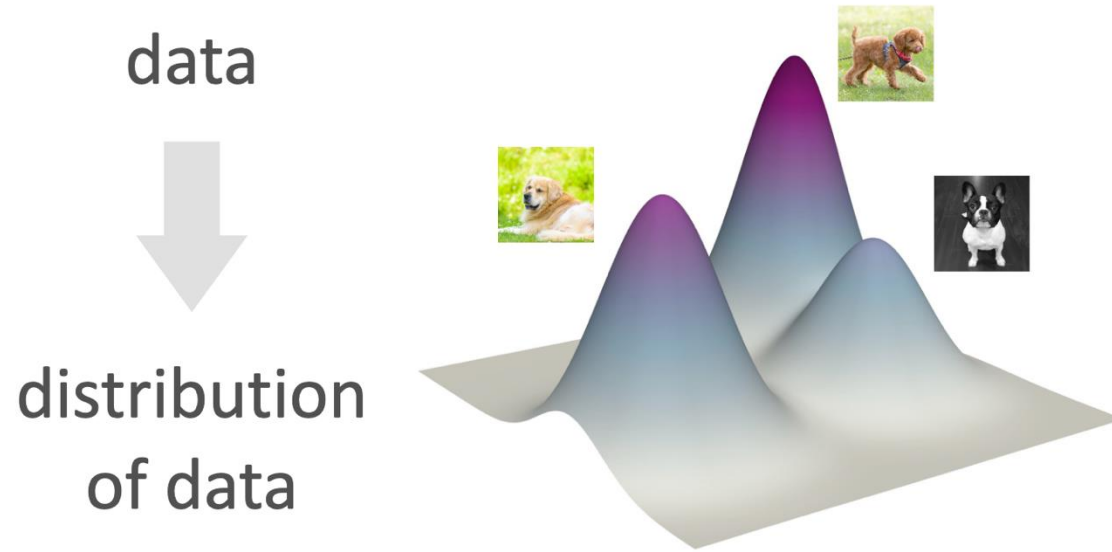
Diffusion-TTA: Test-time Adaptation of Discriminative Models via Generative Feedback. Prabhudesai et al.

Generative models w/ probabilistic modeling

data

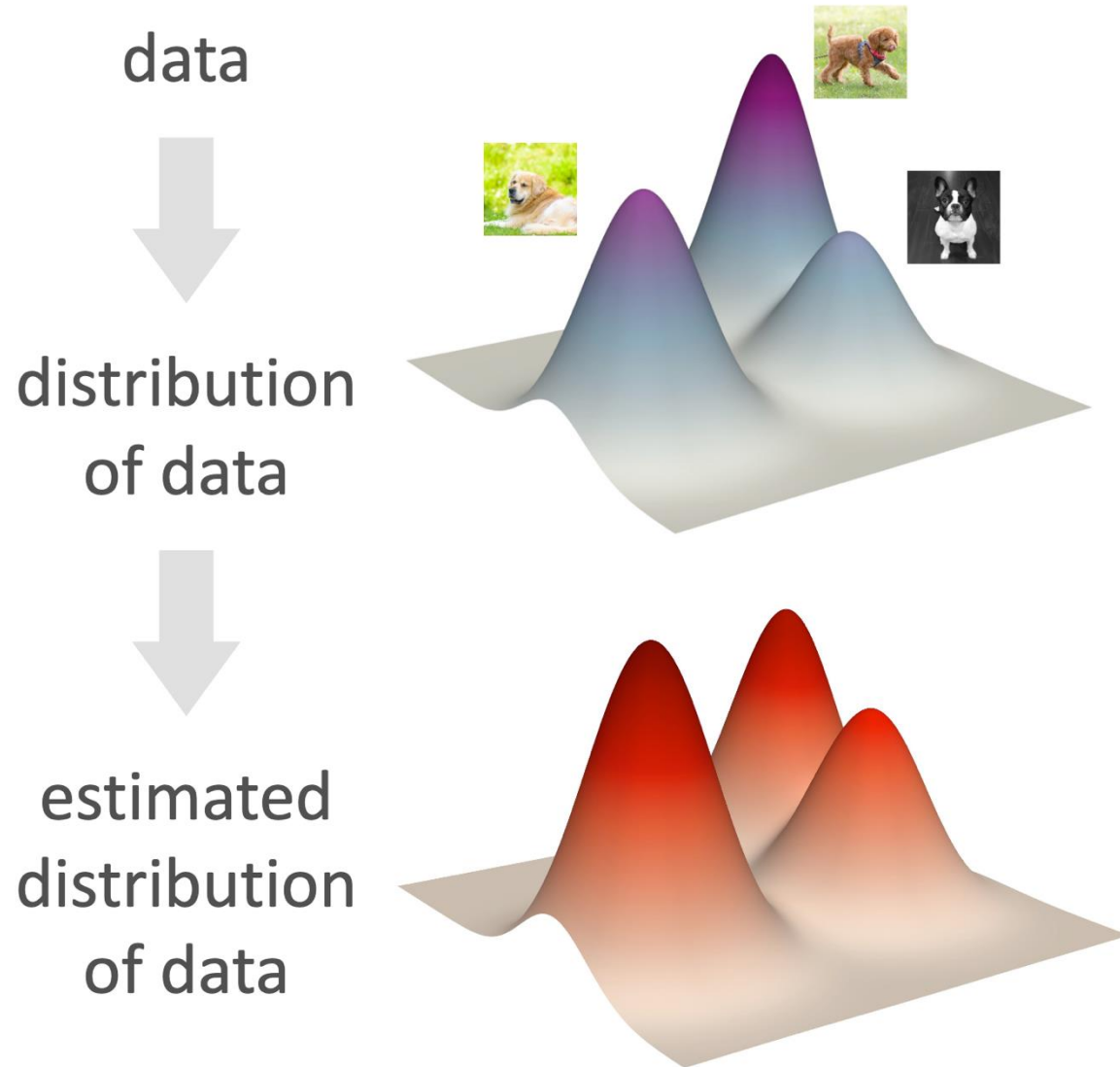


Generative models w/ probabilistic modeling



- Capture the probabilistic distribution of training with the model

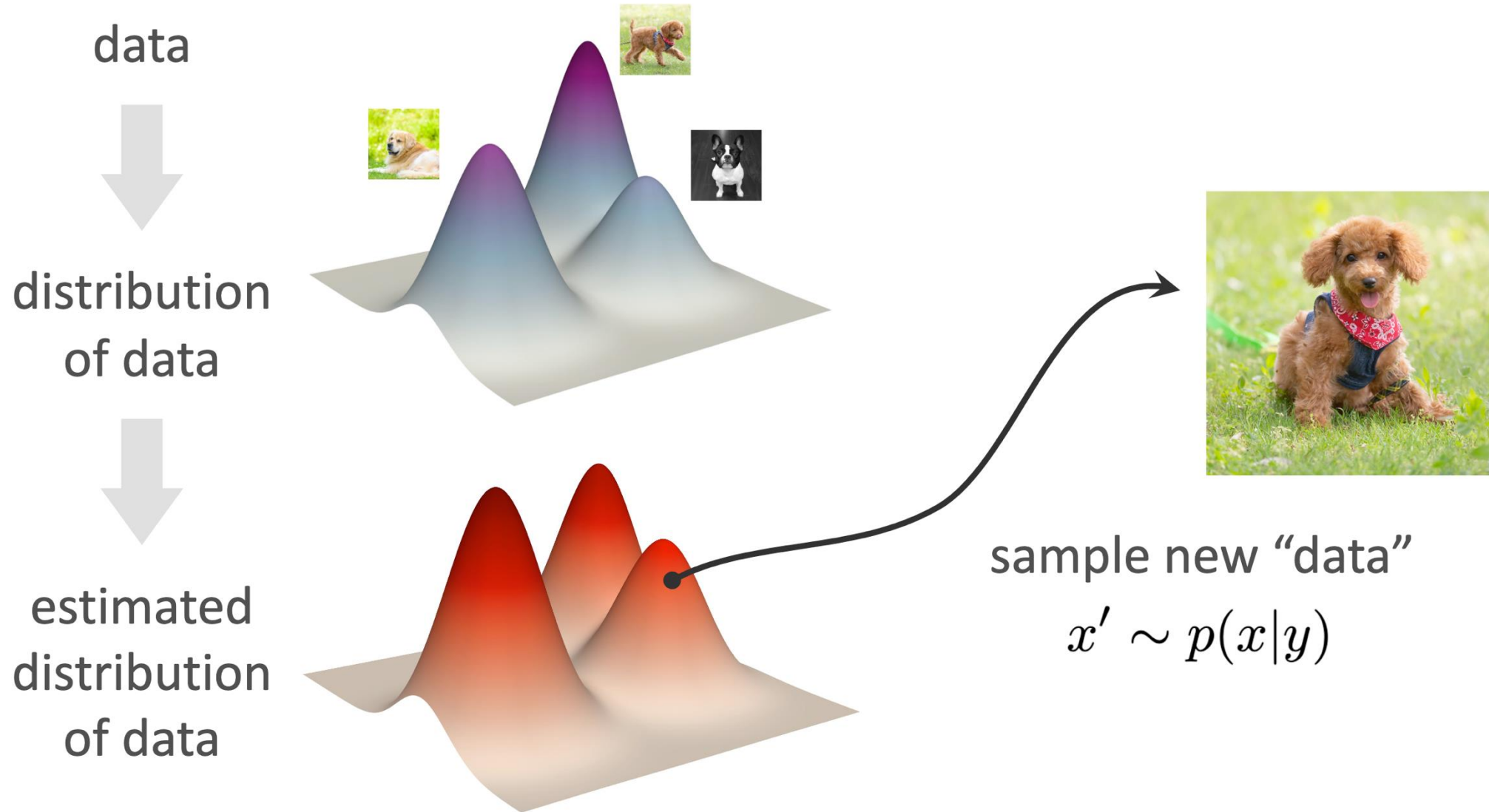
Generative models w/ probabilistic modeling



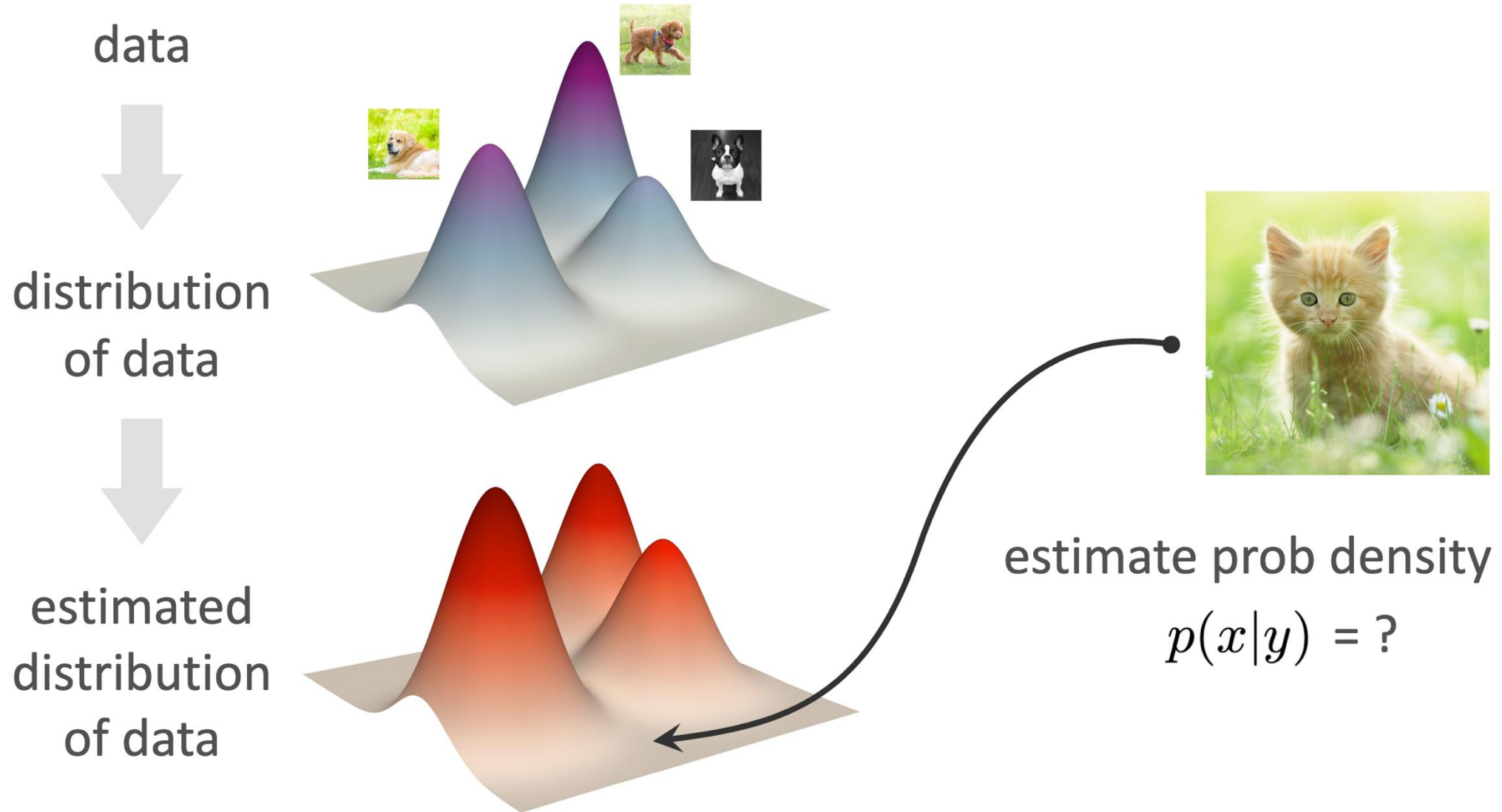
- Optimize a loss function

$$\mathcal{L}(\text{distribution of data}, \text{estimated distribution of data})$$

Generative models w/ probabilistic modeling



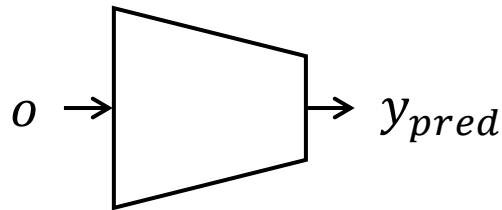
Generative models w/ probabilistic modeling



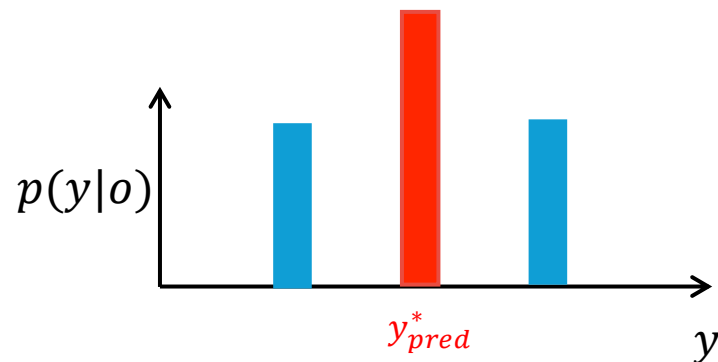
Probabilistic Models Differ From Regression Models

Deterministic

- predicts the same output given the same input

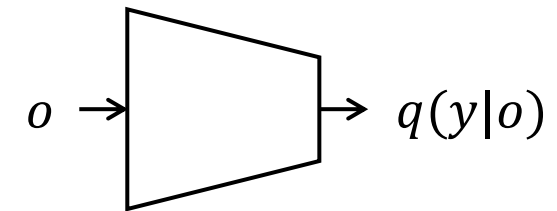


$$y_{pred}^* = \underset{y_{pred}}{\operatorname{argmin}} \frac{1}{N} \sum_{n=1}^N |y_{pred} - y_n|^2$$

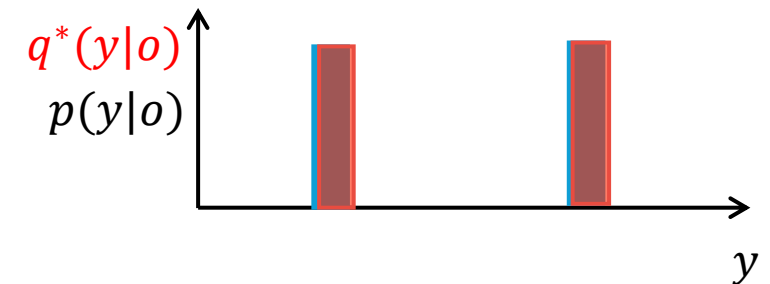


Probabilistic

- predicts the same likelihood of all output given the same input



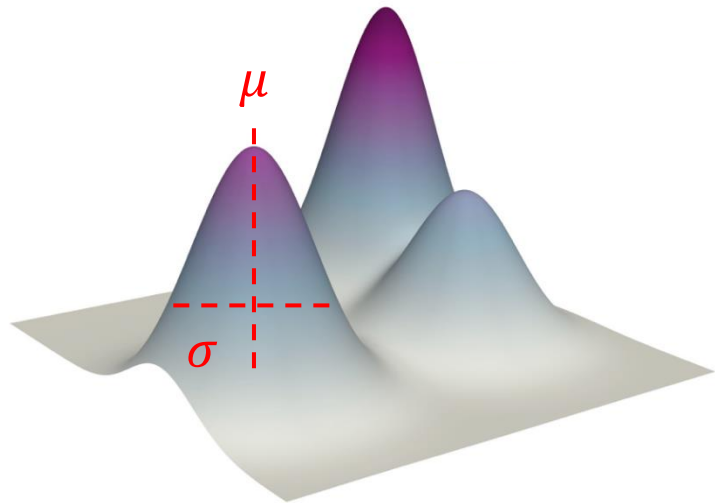
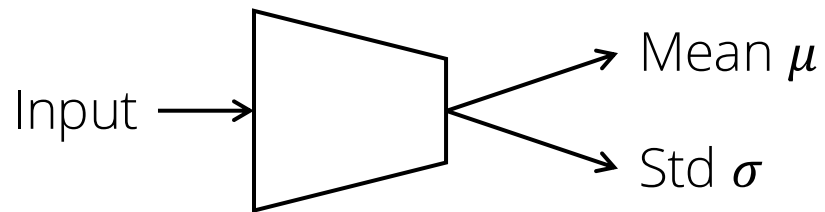
$$q^*(y|o) = \underset{q^*}{\operatorname{argmin}} D_{KL}(q(y|o), p(y|o))$$



How to Model Probabilistic Distribution?

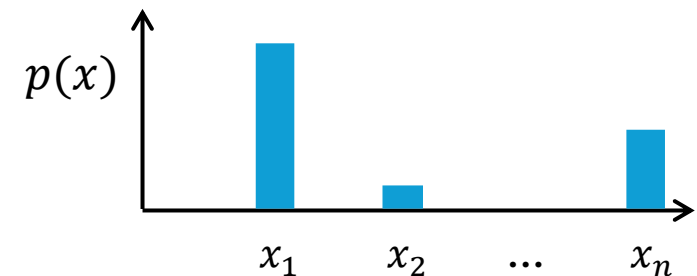
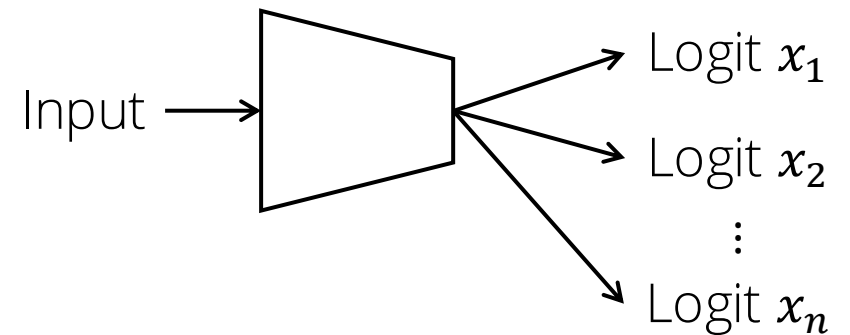
Gaussian

- Predict the mean and standard deviation to model the probability of sample "x"
- VAE, Diffusion models etc



Categorical

- Predict the categorical probability of sample "x"
- Most LLMs etc



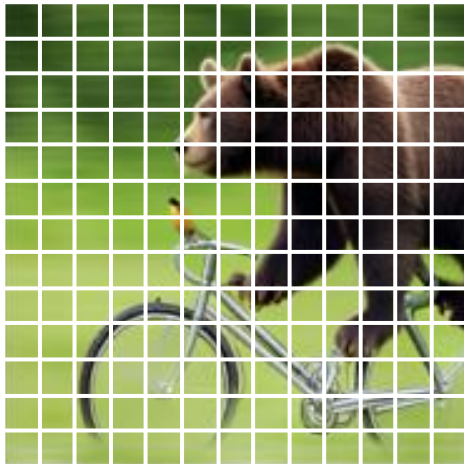
Content

- Generative Modeling
 - Autoregressive Models
 - Variational AutoEncoders
 - Denoising Diffusion Probabilistic Models
 - Flow Matching Models
- 4D Generative Modeling
 - 3D Generation with Multi-view Video Generation
 - 4D Generation with Multi-view Video Generation
 - 3D Geometry Generation

What are Common Generative Models?

Generation Can be Considered as Modeling Joint Distribution

Once upon a time ...



- Text generation:

$$p(\text{word1} = \text{Once}, \text{word2} = \text{upon})$$

- Image generation:

$$p(\text{pixel1} = \text{red}, \text{pixel2} = \text{green})$$

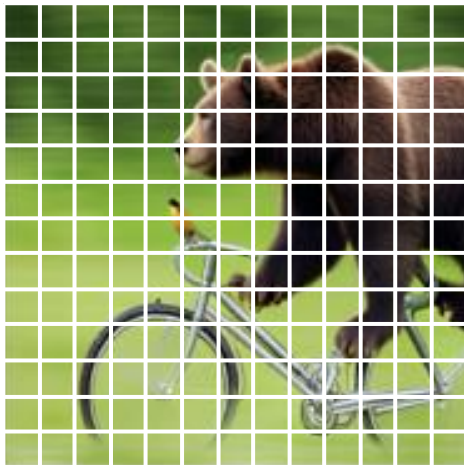
- Joint distribution of data:

$$p(x1 = X1, x2 = X2)$$

➤ How to model this?

Generation Can be Considered as Modeling Joint Distribution

Once upon a time ...



- Text generation:

$$p(\text{word1} = \text{Once}, \text{word2} = \text{upon})$$

- Image generation:

$$p(\text{pixel1} = \text{red}, \text{pixel2} = \text{green})$$

- Joint distribution of data:

$$p(\mathbf{x1} = \mathbf{X1}, \mathbf{x2} = \mathbf{X2})$$

- Modeling as independent distributions

$$p(\mathbf{x1}, \mathbf{x2}) \not\rightarrow p(\mathbf{x1})p(\mathbf{x2})$$

- Modeling as conditional distributions

$$p(\mathbf{x1}, \mathbf{x2}) \rightarrow p(\mathbf{x1})p(\mathbf{x2}|\mathbf{x1})$$

Conditional Distribution Modeling

Chain rule:

Any joint distribution can be written as a product of conditionals

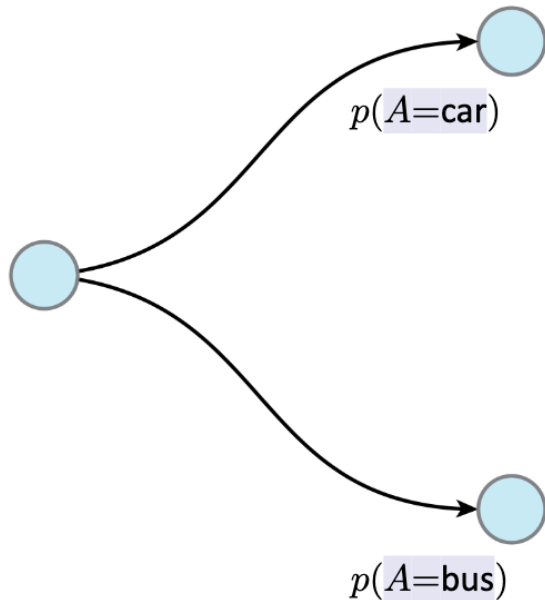
$$p(A, B, C) = p(A)p(B | A)p(C | A, B)$$

Conditional Distribution Modeling

Chain rule:

Any joint distribution can be written as a product of conditionals

$$p(A, B, C) = p(A)p(B | A)p(C | A, B)$$

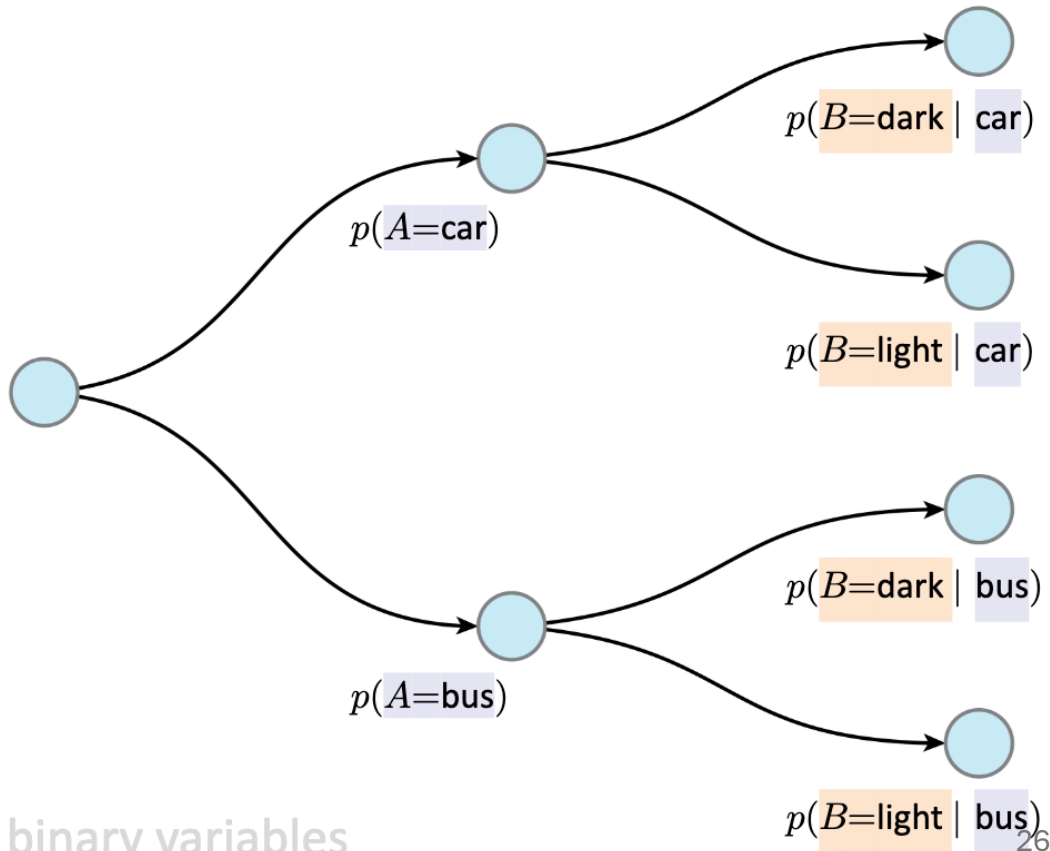


Conditional Distribution Modeling

Chain rule:

Any joint distribution can be written as a product of conditionals

$$p(A, B, C) = p(A)p(B | A)p(C | A, B)$$



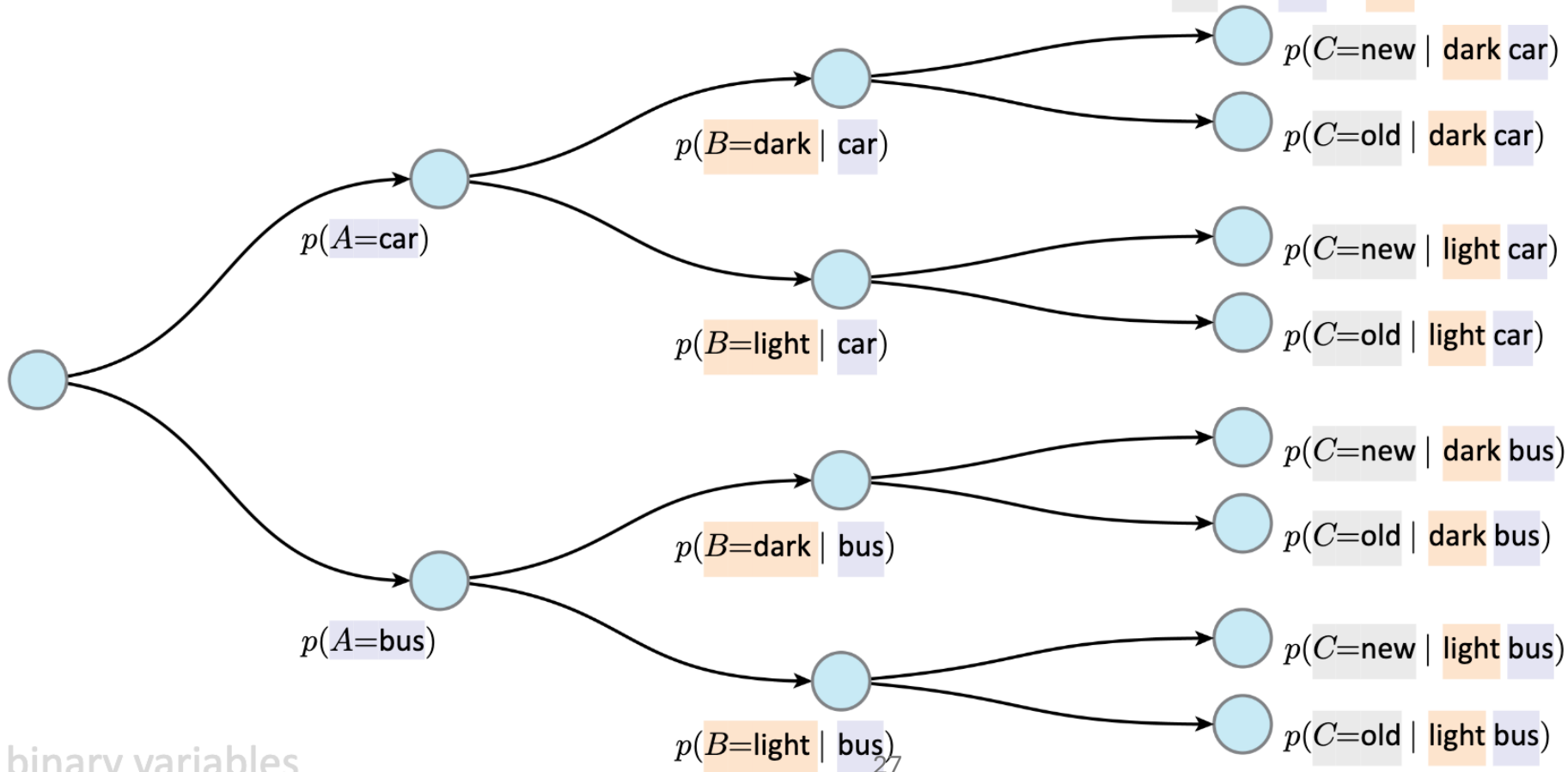
*example: binary variables

Conditional Distribution Modeling

Chain rule:

Any joint distribution can be written as a product of conditionals

$$p(A, B, C) = p(A)p(B | A)p(C | A, B)$$



*example: binary variables

Conditional Distribution Modeling

Chain rule:

Any joint distribution can be written as a product of conditionals

in any order:

$$\begin{aligned} p(A, B, C) &= p(A)p(B | A)p(C | A, B) \\ &= p(A)p(C | A)p(B | A, C) \\ &= p(B)p(A | B)p(C | A, B) \\ &= p(B)p(C | B)p(A | B, C) \\ &= p(C)p(A | C)p(B | A, C) \\ &= p(C)p(B | C)p(A | B, C) \end{aligned}$$

* This is a foundation of Masked Autoregressive (MAR) models:

Li, et al. Autoregressive Image Generation without Vector Quantization, 2024

Conditional Distribution Modeling

Chain rule:

Any joint distribution can be written as a product of conditionals

in any partition:

$$\begin{aligned} p(A, B, C, D) &= p(A, B)p(C, D \mid A, B) \\ &= p(C, D)p(A, B \mid C, D) \\ &= p(A, B, C)p(D \mid A, B, C) \\ &= \dots \end{aligned}$$

* This is a foundation of Masked Autoregressive (MAR) models:

Li, et al. Autoregressive Image Generation without Vector Quantization, 2024

Categorize Generation Models by the Way of Modeling Conditional Distributions

- Condition on data partitions: Autoregressive models

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) = p(\mathbf{x}_1)p(\mathbf{x}_2|\mathbf{x}_1) \cdots p(\mathbf{x}_n|\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n-1})$$

- Condition on latent variables: VAE, Diffusion models

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) = p(\mathbf{z})p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n|\mathbf{z})$$

$$\text{with } p(\mathbf{z}) = p(z_1)p(z_2|z_1) \cdots p(z_n|z_1, z_2, \dots, z_{n-1})$$

Content

- Generative Modeling
 - Autoregressive Models
 - Variational AutoEncoders
 - Denoising Diffusion Probabilistic Models
 - Flow Matching Models
- 4D Generative Modeling
 - 3D Generation with Multi-view Video Generation
 - 4D Generation with Multi-view Video Generation
 - 3D Geometry Generation

AutoRegressive Models

- Condition on data partitions:

$$p(x_1, x_2, \dots, x_n) = p(x_1)p(x_2|x_1) \cdots p(x_n|x_1, x_2, \dots, x_{n-1})$$

AutoRegressive Models

- Condition on data partitions:

$$p(x_1, x_2, \dots, x_n) = p(x_1)p(x_2|x_1) \cdots p(x_n|x_1, x_2, \dots, x_{n-1})$$

$$= p_{\theta}(x_1)p_{\theta}(x_2|x_1) \cdots p_{\theta}(x_n|x_1, x_2, \dots, x_{n-1})$$

An inductive bias that shares models among
different mapping

AutoRegressive Models

- Condition on data partitions:

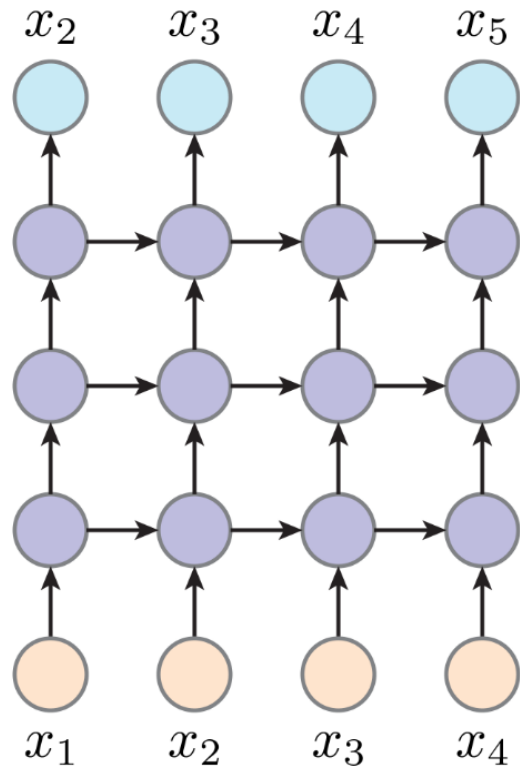
$$p(x_1, x_2, \dots, x_n) = p(x_1)p(x_2|x_1) \cdots p(x_n|x_1, x_2, \dots, x_{n-1})$$

$$= p_{\theta}(x_1)p_{\theta}(x_2|x_1) \cdots p_{\theta}(x_n|x_1, x_2, \dots, x_{n-1})$$

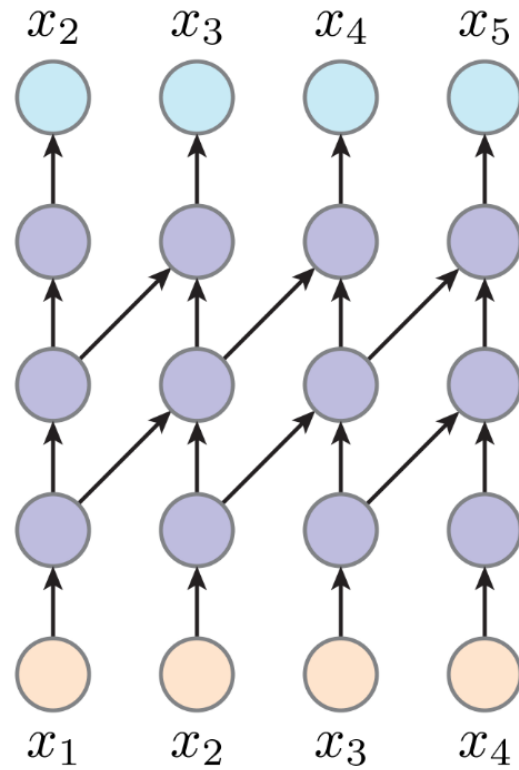
An inductive bias that shares models
among different mappings

What kind of model architectures can
handle these different mappings?

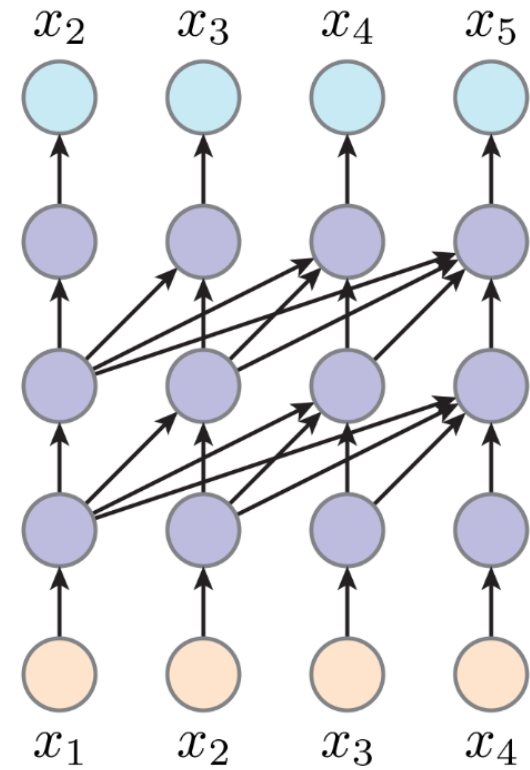
Common Architectures for Autoregression



RNN

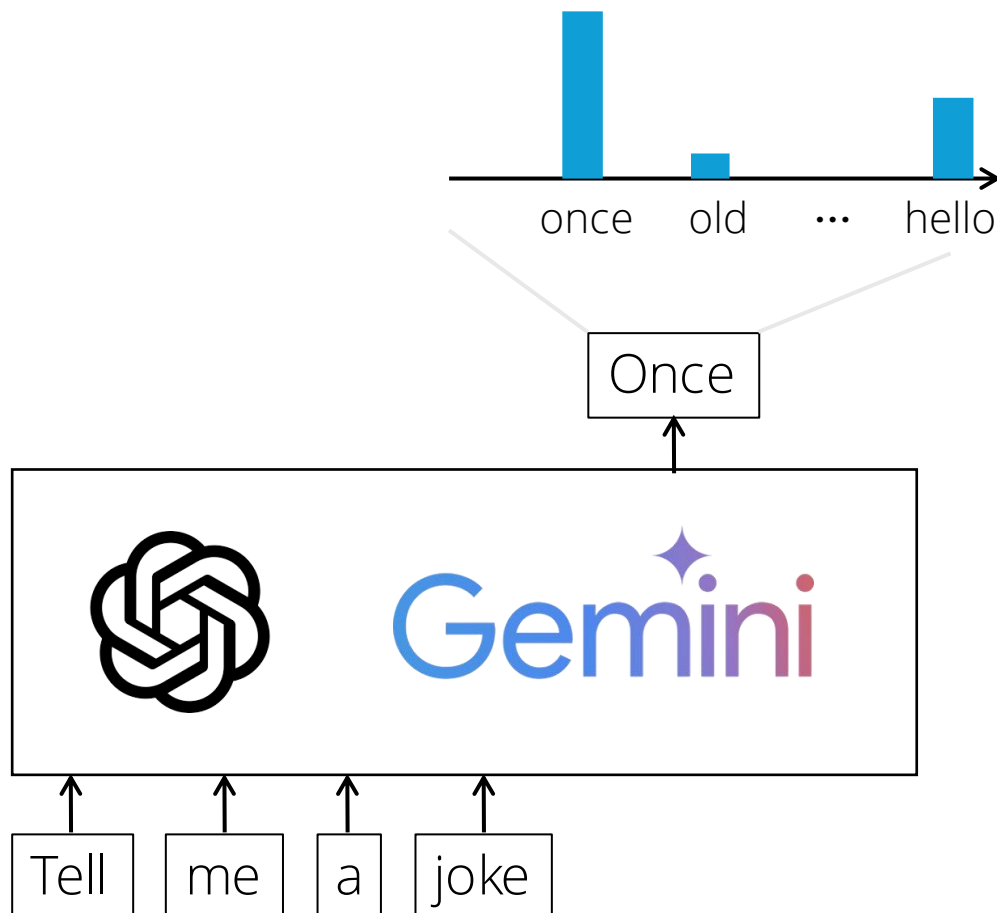


CNN



Attention

Discrete Autoregressive Transformer are Categorical Probabilistic Generative Models



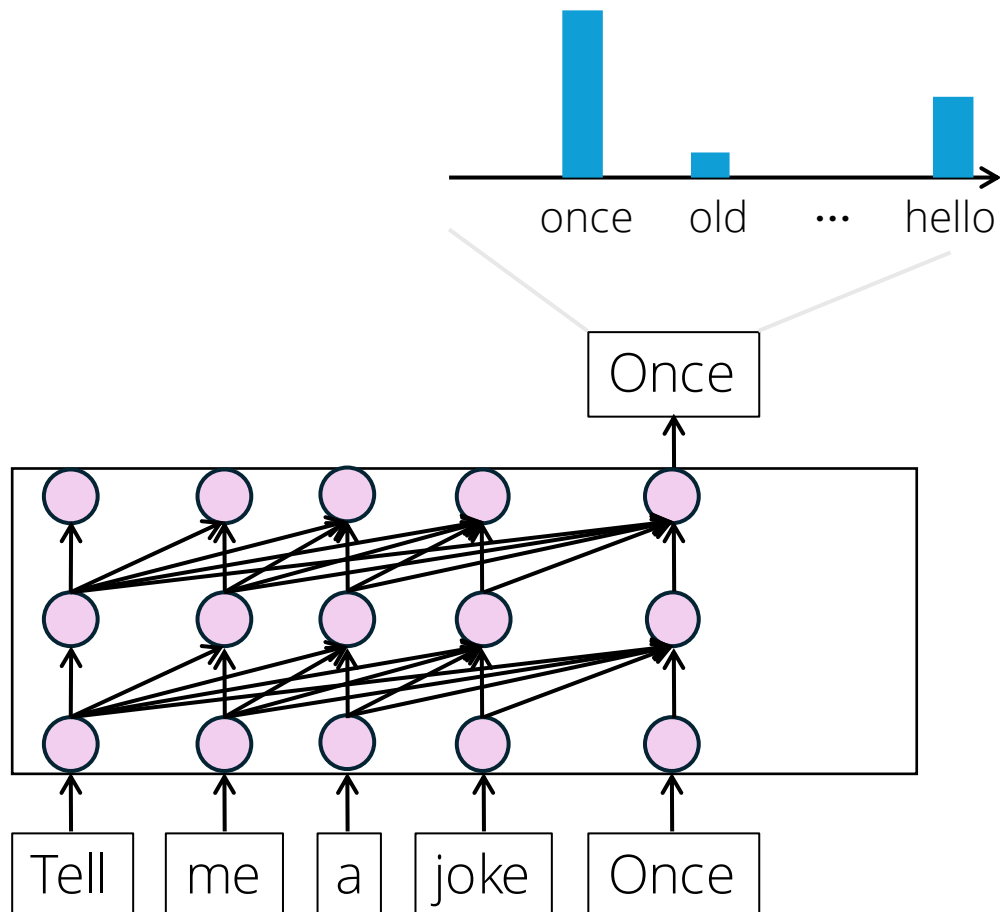
GPT-4.5

I'm going through a tough time after failing a test

Aw, I'm really sorry to hear that. Failing a test can feel pretty tough and discouraging, but remember, it's just one moment—not a reflection of who you are or your capabilities. Want to talk about what happened, or do you just need a distraction? I'm here either way.

<https://openai.com/index/introducing-gpt-4-5/>

Discrete Autoregressive Transformer are Categorical Probabilistic Generative Models



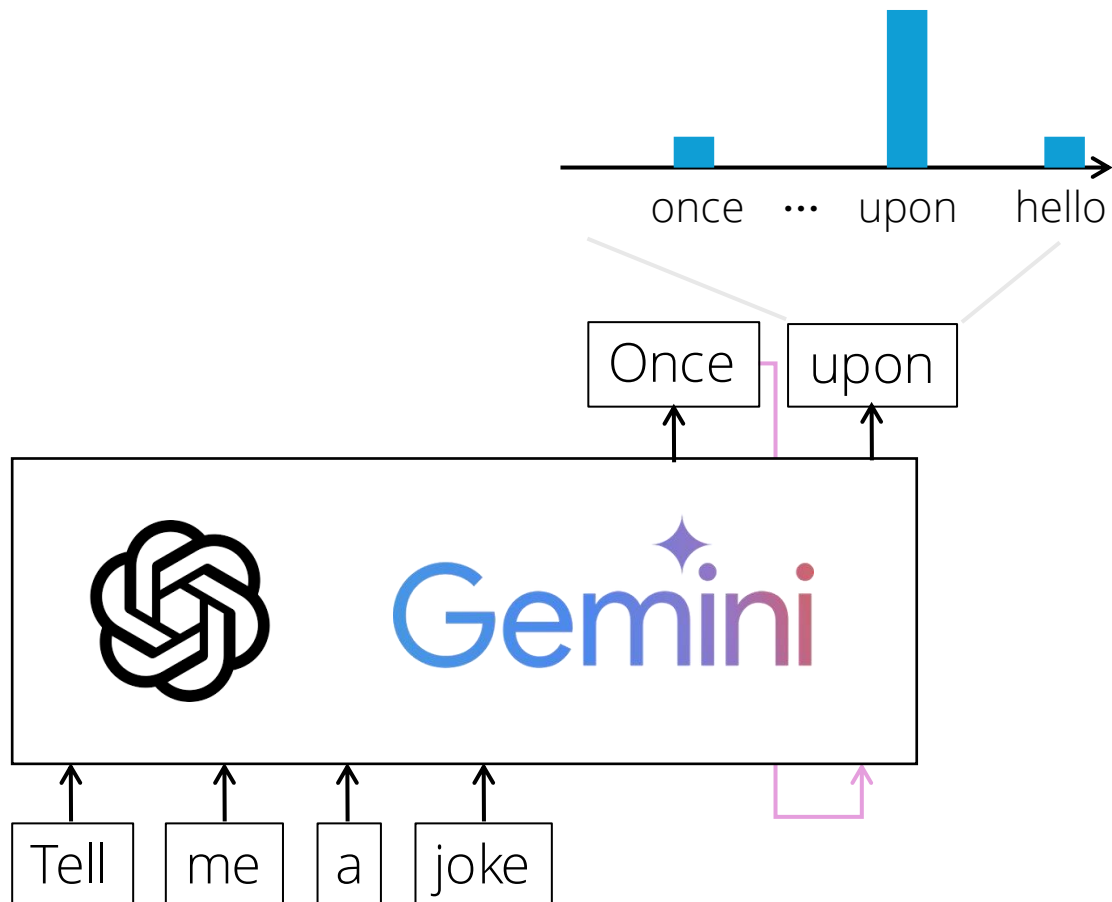
GPT-4.5

I'm going through a tough time after failing a test

Aw, I'm really sorry to hear that. Failing a test can feel pretty tough and discouraging, but remember, it's just one moment—not a reflection of who you are or your capabilities. Want to talk about what happened, or do you just need a distraction? I'm here either way.

<https://openai.com/index/introducing-gpt-4-5/>

Discrete Autoregressive Transformer are Categorical Probabilistic Generative Models



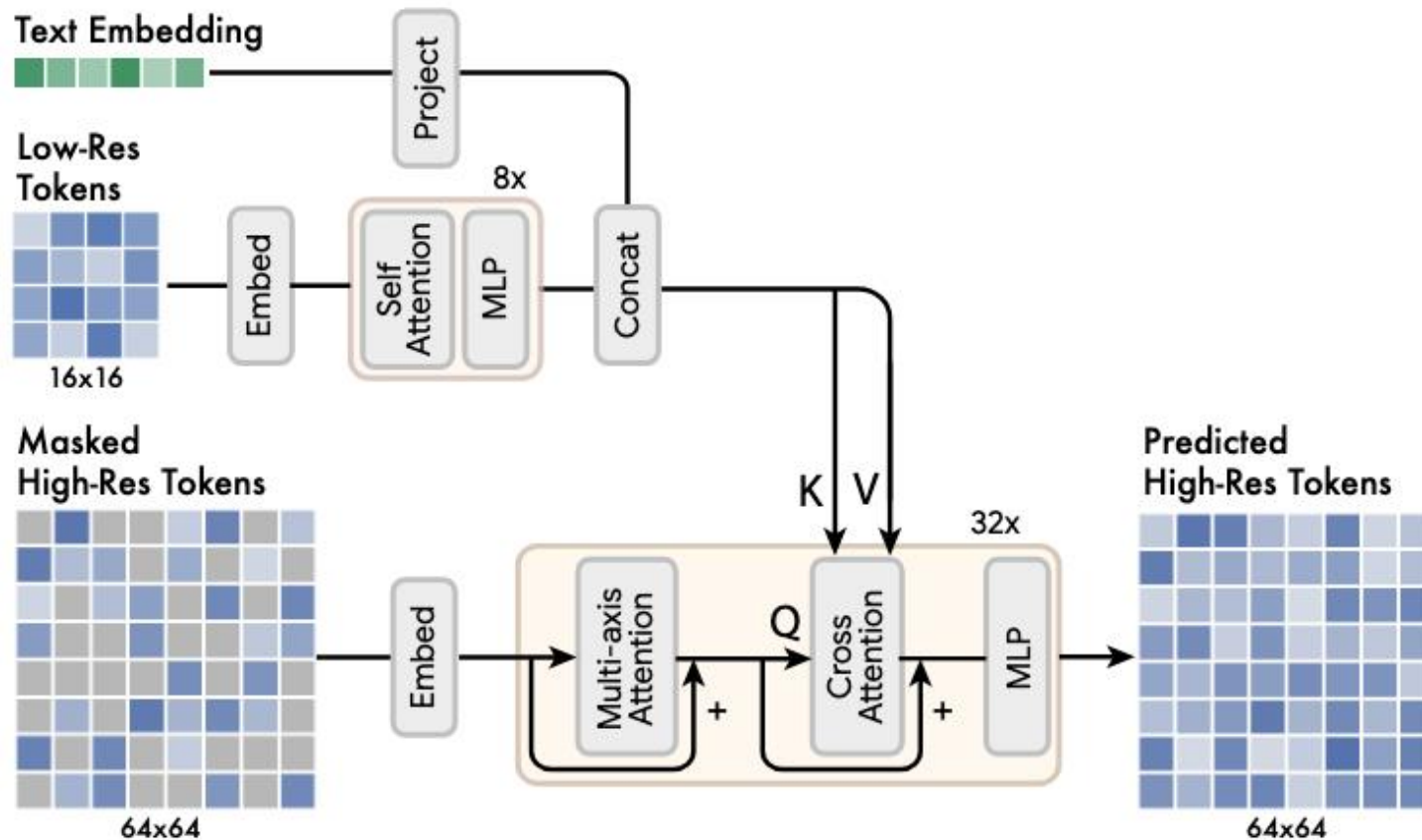
GPT-4.5

I'm going through a tough time after failing a test

Aw, I'm really sorry to hear that. Failing a test can feel pretty tough and discouraging, but remember, it's just one moment—not a reflection of who you are or your capabilities. Want to talk about what happened, or do you just need a distraction? I'm here either way.

<https://openai.com/index/introducing-gpt-4-5/>

Discrete Autoregressive Transformer Can Also Generate Images / Videos



Text

A bear riding a bicycle, with a bird perched on the handlebars.

A high contrast portrait photo of a fluffy hamster wearing an orange beanie and sunglasses holding a sign that says "Let's PAINT!"

LowRes 256x256



HighRes 512x512



Content

- Generative Modeling
 - Autoregressive Models
 - Variational AutoEncoders
 - Denoising Diffusion Probabilistic Models
 - Flow Matching Models
- 4D Generative Modeling
 - 3D Generation with Multi-view Video Generation
 - 4D Generation with Multi-view Video Generation
 - 3D Geometry Generation

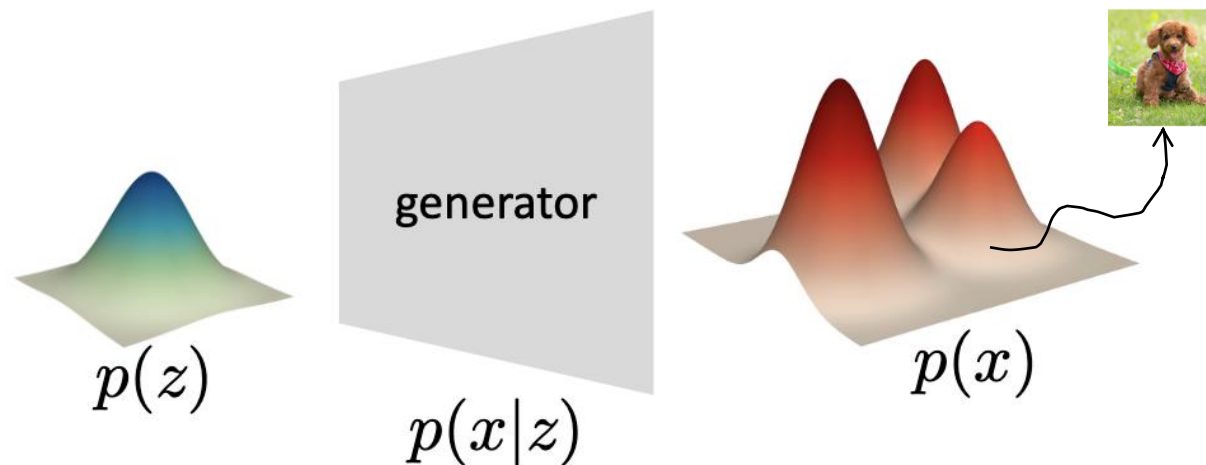
Another Way of Modeling Conditional Distributions

- Condition on latent variables:

$$p(x_1, x_2, \dots, x_n) = p(\mathbf{z})p(x_1, x_2, \dots, x_n|\mathbf{z})$$

You can select a random distribution for variable \mathbf{z} , such as Gaussian.

In fact, Gaussian prior results in minimal variance if the data distribution is unknown [TODO; citation]



In Fact, This Views Generation as Mapping a Prior Distribution $p(\mathbf{z})$ to the Data Distribution $p(\mathbf{x})$



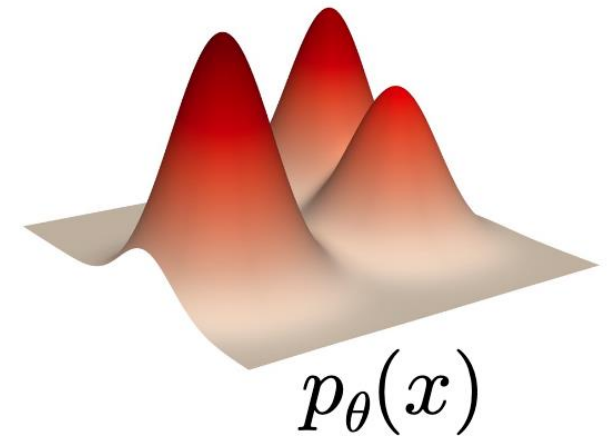
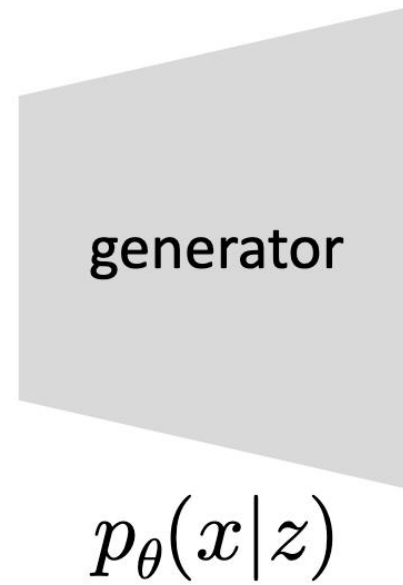
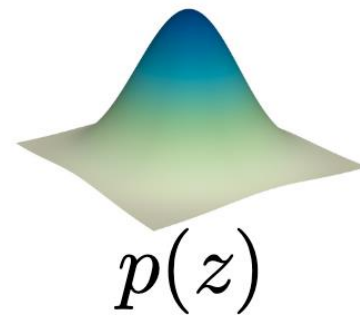
How to model the mapping differentiates generative models, such as VAE, DDPM, Flow Matching Model, and GAN...

VAE as Latent Variable Models

We want to maximize $\mathbb{E}_{x \sim p_{data}} \log p_{\theta}(x)$

with $p_{\theta}(x)$ represented as:

$$p_{\theta}(x) = \int_z p_{\theta}(x|z)p(z)dz$$



VAE as Latent Variable Models

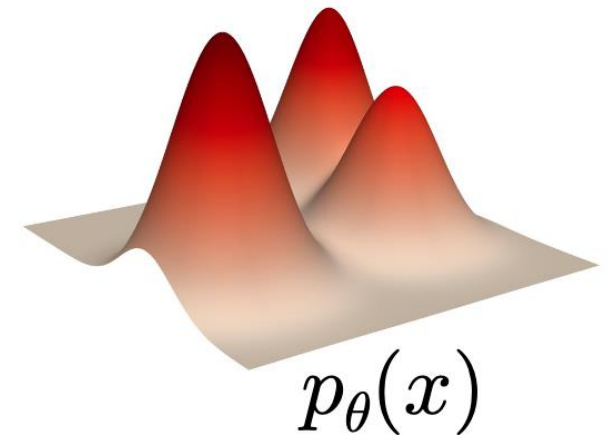
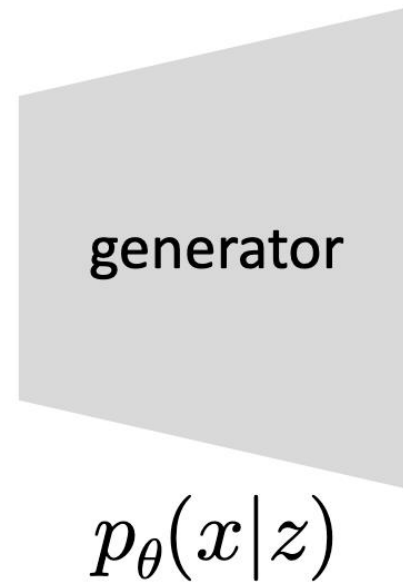
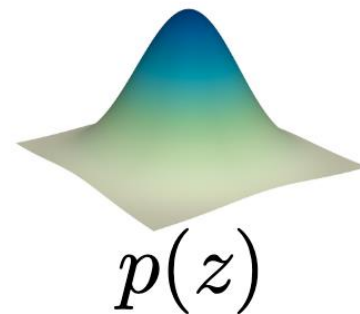
We want to maximize $\mathbb{E}_{x \sim p_{data}} \log p_{\theta}(x)$

with $p_{\theta}(x)$ represented as:

$$p_{\theta}(x) = \int_z p_{\theta}(x|z) p(z) dz$$

Two sets of unknowns:

- We need to optimize for θ
- We can't control **"true"** $p(z)$



Idea: introduce a "controllable" distribution $q(z)$

Latent Variable Models

$$\begin{aligned} & \log p_{\theta}(x) \\ = & \int_z q(z) \log p_{\theta}(x) dz \\ = & \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \right) dz \\ = & \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \frac{q(z)}{q(z)} \right) dz \\ = & \int_z q(z) \left(\log p_{\theta}(x|z) + \log \frac{p_{\theta}(z)}{q(z)} + \log \frac{q(z)}{p_{\theta}(z|x)} \right) dz \\ = & \mathbb{E}_{z \sim q(z)} \left[\log p_{\theta}(x|z) \right] - \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z) \right) + \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z|x) \right) \end{aligned}$$

Rewrite log likelihood by latent z

- for any distribution $q(z)$
- Bayes' rule

Latent Variable Models

$$\begin{aligned} & \log p_{\theta}(x) \\ = & \int_z q(z) \log p_{\theta}(x) dz \\ = & \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \right) dz \\ = & \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \frac{q(z)}{q(z)} \right) dz && \bullet \text{ just algebra} \\ = & \int_z q(z) \left(\log p_{\theta}(x|z) + \log \frac{p_{\theta}(z)}{q(z)} + \log \frac{q(z)}{p_{\theta}(z|x)} \right) dz && \bullet \text{ just algebra} \\ = & \mathbb{E}_{z \sim q(z)} \left[\log p_{\theta}(x|z) \right] - \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z) \right) + \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z|x) \right) \end{aligned}$$

Rewrite log likelihood by latent z

- for any distribution $q(z)$

- Bayes' rule

- just algebra

- just algebra

Latent Variable Models

$$\begin{aligned} & \log p_{\theta}(x) \\ = & \int_z q(z) \log p_{\theta}(x) dz \\ = & \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \right) dz \\ = & \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \frac{q(z)}{q(z)} \right) dz && \bullet \text{ just algebra} \\ = & \int_z q(z) \left(\log p_{\theta}(x|z) + \log \frac{p_{\theta}(z)}{q(z)} + \log \frac{q(z)}{p_{\theta}(z|x)} \right) dz && \bullet \text{ just algebra} \\ = & \mathbb{E}_{z \sim q(z)} \left[\log p_{\theta}(x|z) \right] - \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z) \right) + \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z|x) \right) \end{aligned}$$

Rewrite log likelihood by latent z

- for any distribution $q(z)$

- Bayes' rule

- just algebra

- just algebra

Latent Variable Models

$$\begin{aligned} & \log p_{\theta}(x) \\ = & \int_z q(z) \log p_{\theta}(x) dz \\ = & \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \right) dz \\ = & \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \frac{q(z)}{q(z)} \right) dz \\ = & \int_z q(z) \left(\log p_{\theta}(x|z) + \log \frac{p_{\theta}(z)}{q(z)} + \log \frac{q(z)}{p_{\theta}(z|x)} \right) dz \\ = & \mathbb{E}_{z \sim q(z)} \left[\log p_{\theta}(x|z) \right] - \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z) \right) + \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z|x) \right) \end{aligned}$$

Rewrite log likelihood by latent z

- for any distribution $q(z)$

- Bayes' rule

- just algebra

- just algebra

Latent Variable Models

$$\begin{aligned} & \log p_{\theta}(x) \\ = & \int_z q(z) \log p_{\theta}(x) dz \\ = & \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \right) dz \\ = & \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \frac{q(z)}{q(z)} \right) dz \\ = & \int_z q(z) \left(\log p_{\theta}(x|z) + \log \frac{p_{\theta}(z)}{q(z)} + \log \frac{q(z)}{p_{\theta}(z|x)} \right) dz \\ = & \mathbb{E}_{z \sim q(z)} \left[\log p_{\theta}(x|z) \right] - \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z) \right) + \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z|x) \right) \end{aligned}$$

Rewrite log likelihood by latent z

- for any distribution $q(z)$
- Bayes' rule
- just algebra
- just algebra

Latent Variable Models

$$\begin{aligned} & \log p_{\theta}(x) \\ = & \int_z q(z) \log p_{\theta}(x) dz \\ = & \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \right) dz \\ = & \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \frac{q(z)}{q(z)} \right) dz \\ = & \int_z q(z) \left(\log p_{\theta}(x|z) + \log \frac{p_{\theta}(z)}{q(z)} + \log \frac{q(z)}{p_{\theta}(z|x)} \right) dz \\ = & \mathbb{E}_{z \sim q(z)} \left[\log p_{\theta}(x|z) \right] - \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z) \right) + \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z|x) \right) \end{aligned}$$

Rewrite log likelihood by latent z

- for any distribution $q(z)$
- Bayes' rule
- just algebra
- just algebra

Latent Variable Models

intractable

$$\log p_{\theta}(x)$$

$$= \int_z q(z) \log p_{\theta}(x) dz$$

$$= \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \right) dz$$

$$= \int_z q(z) \log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)} \frac{q(z)}{q(z)} \right) dz$$

$$= \int_z q(z) \left(\log p_{\theta}(x|z) + \log \frac{p_{\theta}(z)}{q(z)} + \log \frac{q(z)}{p_{\theta}(z|x)} \right) dz$$

$$= \mathbb{E}_{z \sim q(z)} \left[\log p_{\theta}(x|z) \right] - \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z) \right) + \mathcal{D}_{\text{KL}} \left(q(z) || p_{\theta}(z|x) \right)$$

tractable

tractable

intractable

Rewrite log likelihood by latent z

- for any distribution $q(z)$
- Bayes' rule

Latent Variable Models

$$\begin{aligned} & \text{intractable} \quad \boxed{\log p_{\theta}(x)} - \boxed{\mathcal{D}_{\text{KL}}(q(z) || p_{\theta}(z|x))} \quad \text{intractable} \\ & = \quad \boxed{\mathbb{E}_{z \sim q(z)} [\log p_{\theta}(x|z)]} - \boxed{\mathcal{D}_{\text{KL}}(q(z) || p_{\theta}(z))} \\ & \qquad \qquad \text{tractable} \qquad \qquad \qquad \text{tractable} \end{aligned}$$

Latent Variable Models

$$\begin{aligned} & \text{intractable} \quad \boxed{\log p_{\theta}(x)} - \boxed{\mathcal{D}_{\text{KL}}(q(z) \parallel p_{\theta}(z|x))} \quad \text{intractable} \\ & = \underbrace{\boxed{\mathbb{E}_{z \sim q(z)} [\log p_{\theta}(x|z)]}}_{\text{tractable}} - \underbrace{\boxed{\mathcal{D}_{\text{KL}}(q(z) \parallel p_{\theta}(z))}}_{\text{tractable}} \end{aligned}$$

- This is called Evidence Lower Bound (ELBO)
- Lower bound of $\log p_{\theta}(x)$
- This equation holds for any distribution $q(z)$

Latent Variable Models

$$\underbrace{\mathbb{E}_{z \sim q(z)} \left[\log p_{\theta}(x|z) \right]}_{\text{tractable}} - \underbrace{\mathcal{D}_{\text{KL}} \left(q(z) \parallel p_{\theta}(z) \right)}_{\text{tractable}}$$

- This is called Evidence Lower Bound (ELBO)
- Lower bound of $\log p_{\theta}(x)$
- This equation holds for any distribution $q(z)$
- Parameterize $q(z)$ by $q_{\phi}(z|x)$

Latent Variable Models

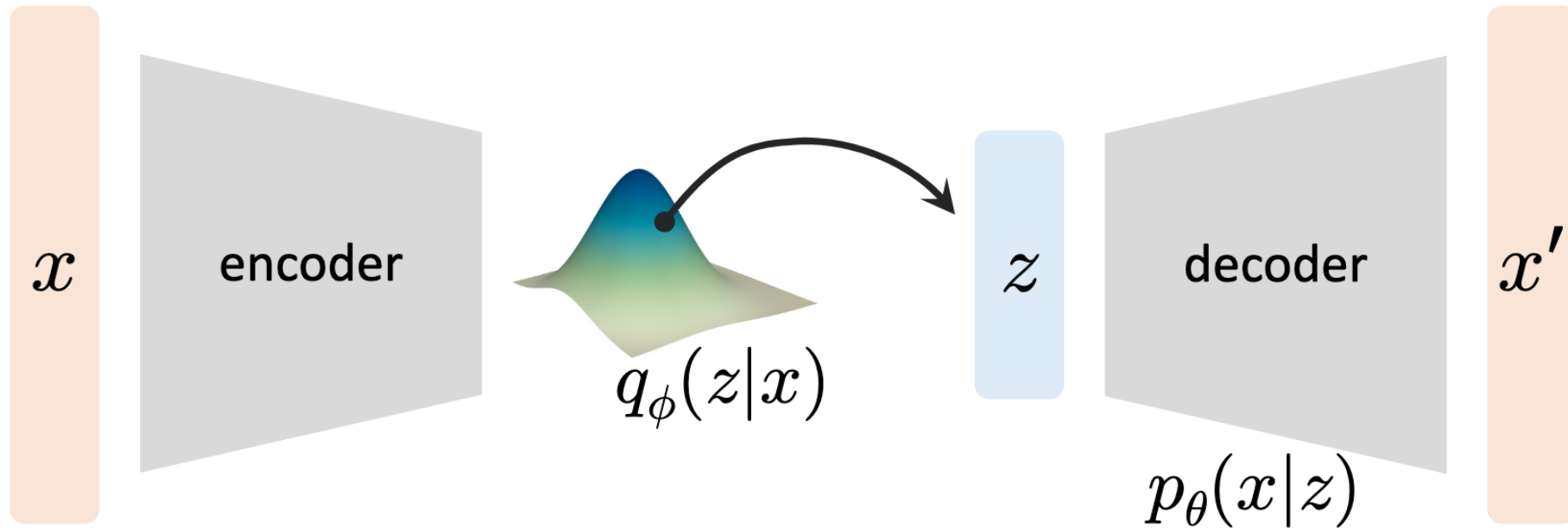
$$\underbrace{\mathbb{E}_{z \sim q_\phi(z)} \left[\log p_\theta(x|z) \right]}_{\text{tractable}} - \underbrace{\mathcal{D}_{\text{KL}} \left(q_\phi(z|x) \parallel p_\theta(z) \right)}_{\text{tractable}}$$

- This is called Evidence Lower Bound (ELBO)
- Lower bound of $\log p_\theta(x)$
- This equation holds for any distribution $q(z)$
- Parameterize $q(z)$ by $q_\phi(z|x)$
- let $p_\theta(z)$ be a simple known prior $p(z)$

Variational Autoencoder

Maximize ELBO \Rightarrow minimize an objective:

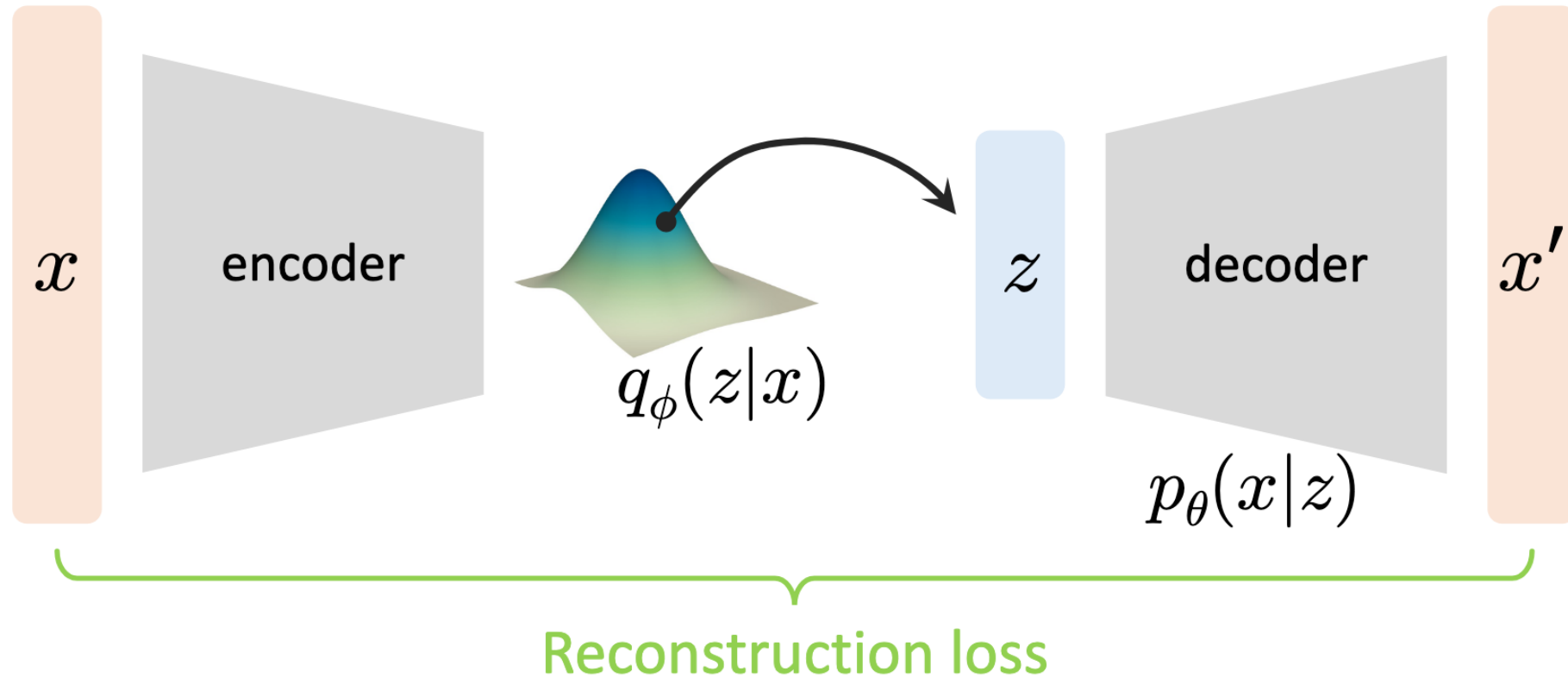
$$\mathcal{L}_{\theta, \phi}(x) = -\mathbb{E}_{z \sim q_{\phi}(z|x)} \left[\log p_{\theta}(x|z) \right] + \mathcal{D}_{\text{KL}} \left(q_{\phi}(z|x) || p(z) \right)$$



Variational Autoencoder

Maximize ELBO \Rightarrow minimize an objective:

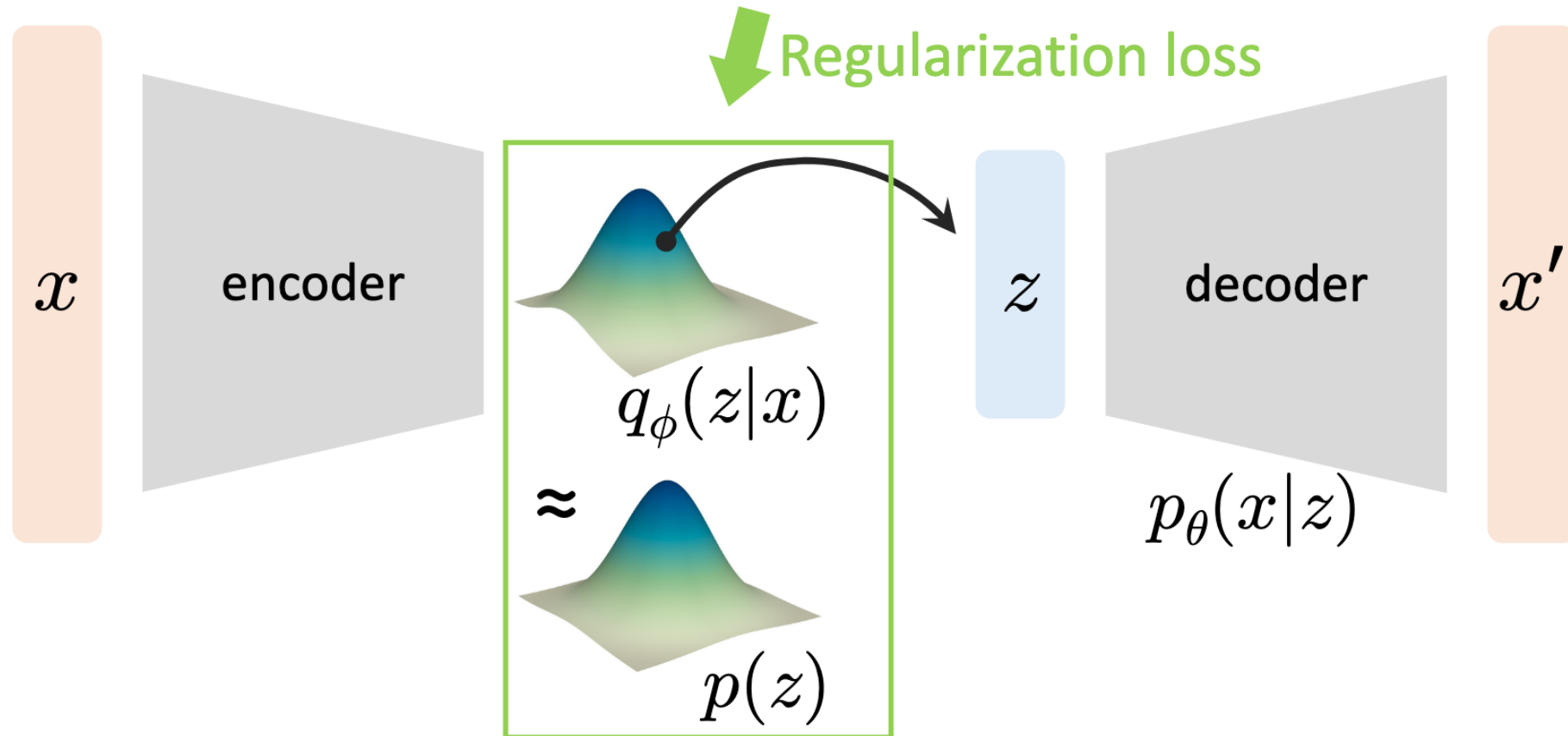
$$\mathcal{L}_{\theta, \phi}(x) = -\mathbb{E}_{z \sim q_{\phi}(z|x)} \left[\log p_{\theta}(x|z) \right] + \mathcal{D}_{\text{KL}} \left(q_{\phi}(z|x) || p(z) \right)$$



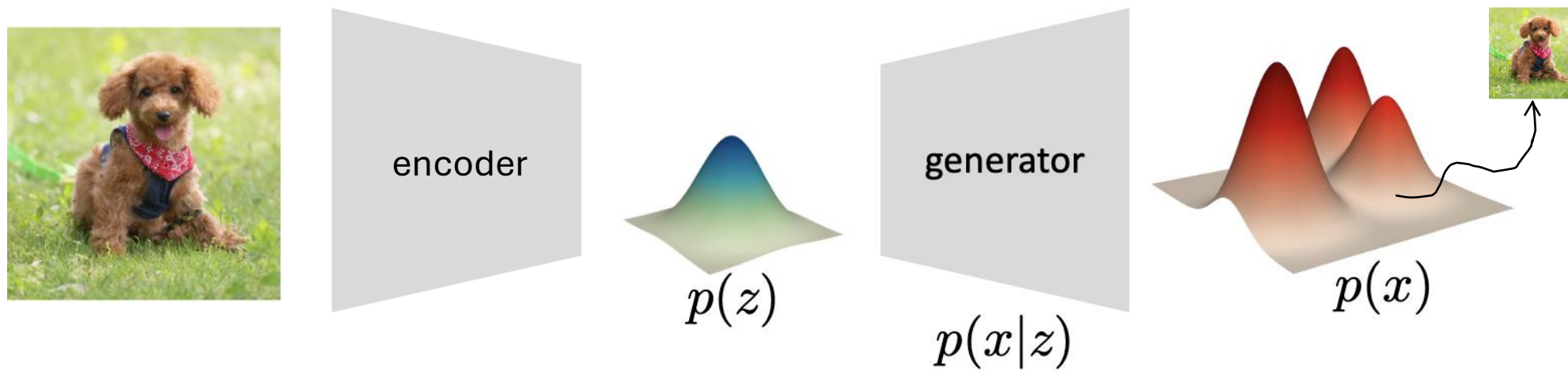
Variational Autoencoder

Maximize ELBO \Rightarrow minimize an objective:

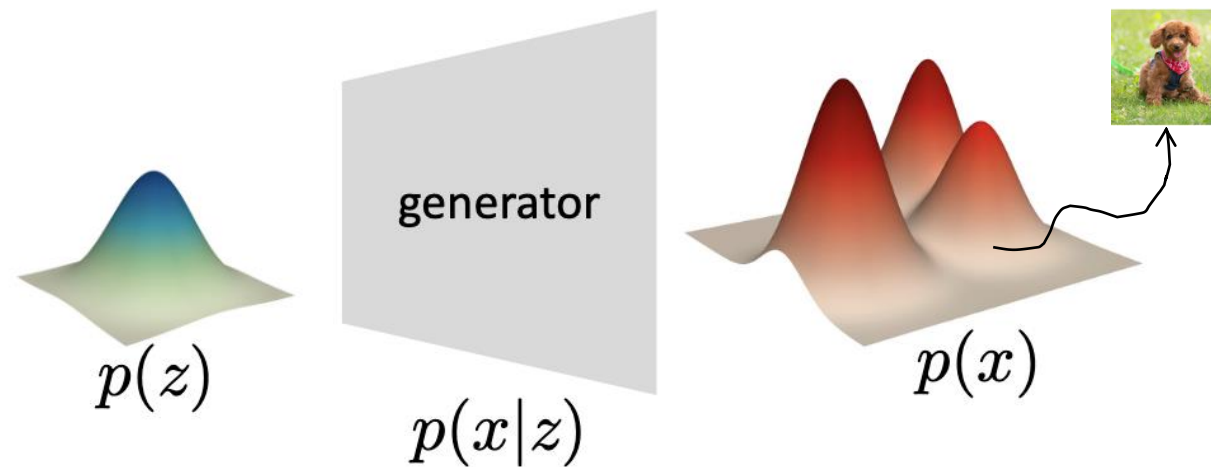
$$\mathcal{L}_{\theta, \phi}(x) = -\mathbb{E}_{z \sim q_{\phi}(z|x)} \left[\log p_{\theta}(x|z) \right] + \mathcal{D}_{\text{KL}} \left(q_{\phi}(z|x) || p(z) \right)$$



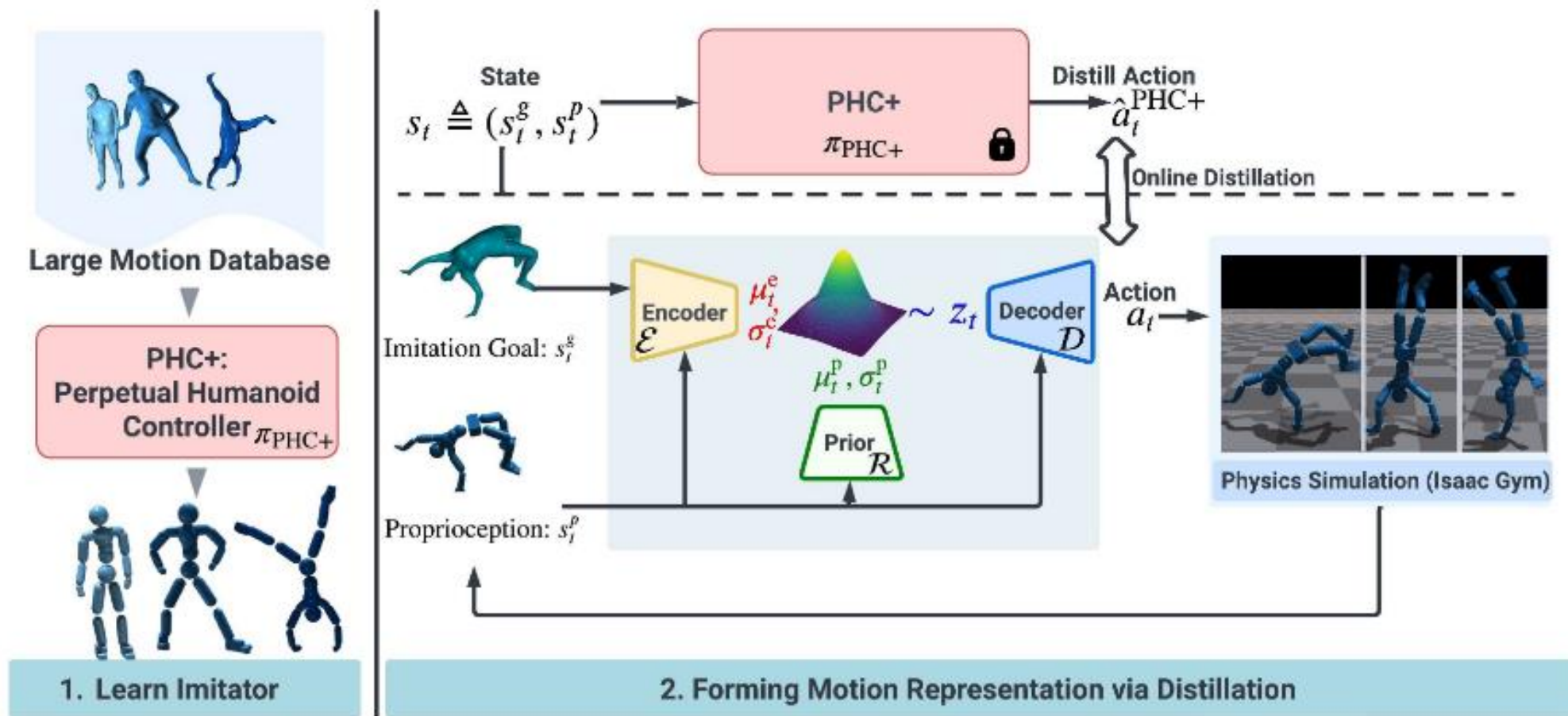
Training:



Testing:



Why Do We Care About VAE? A Common Technique to Compress Data as Priors

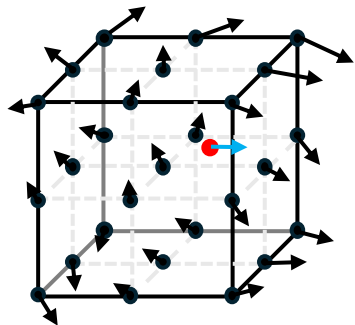


Why Do We Care About VAE? A Common Technique to Compress Data as Priors

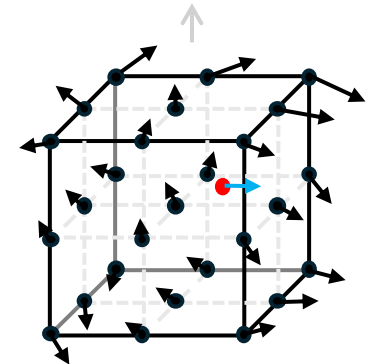
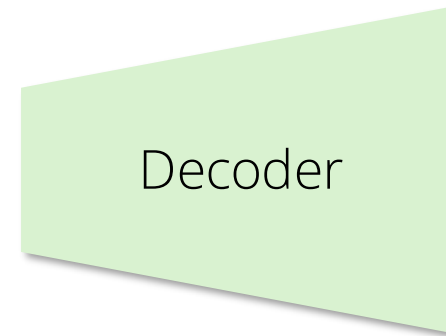
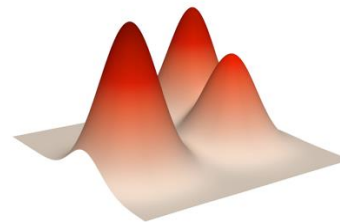
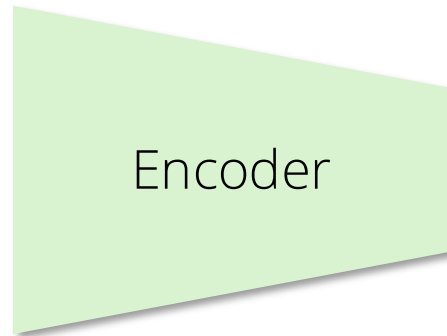
Static Object



Animated Object



Motion Field



Motion Field

Content

- Generative Modeling
 - Autoregressive Models
 - Variational AutoEncoders
 - Denoising Diffusion Probabilistic Models
 - Flow Matching Models
- 4D Generative Modeling
 - 3D Generation with Multi-view Video Generation
 - 4D Generation with Multi-view Video Generation
 - 3D Geometry Generation

Another Way of Modeling Conditional Distributions

- Condition on latent variables:

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) = p(\mathbf{z})p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n | \mathbf{z})$$

$$\text{with } p(\mathbf{z}) = p(\mathbf{z}_1)p(\mathbf{z}_2 | \mathbf{z}_1) \dots p(\mathbf{z}_n | \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_{n-1})$$

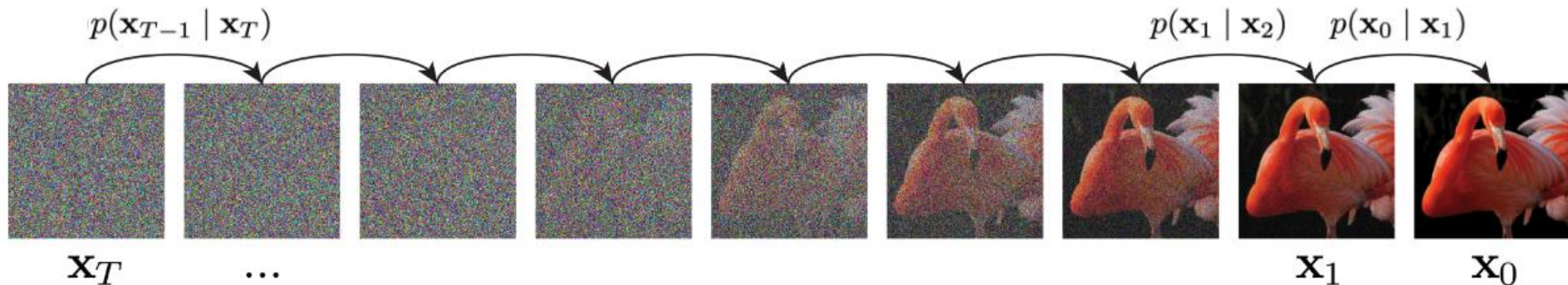
Another Way of Modeling Conditional Distributions

- Condition on latent variables:

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) = p(\mathbf{z})p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n | \mathbf{z})$$

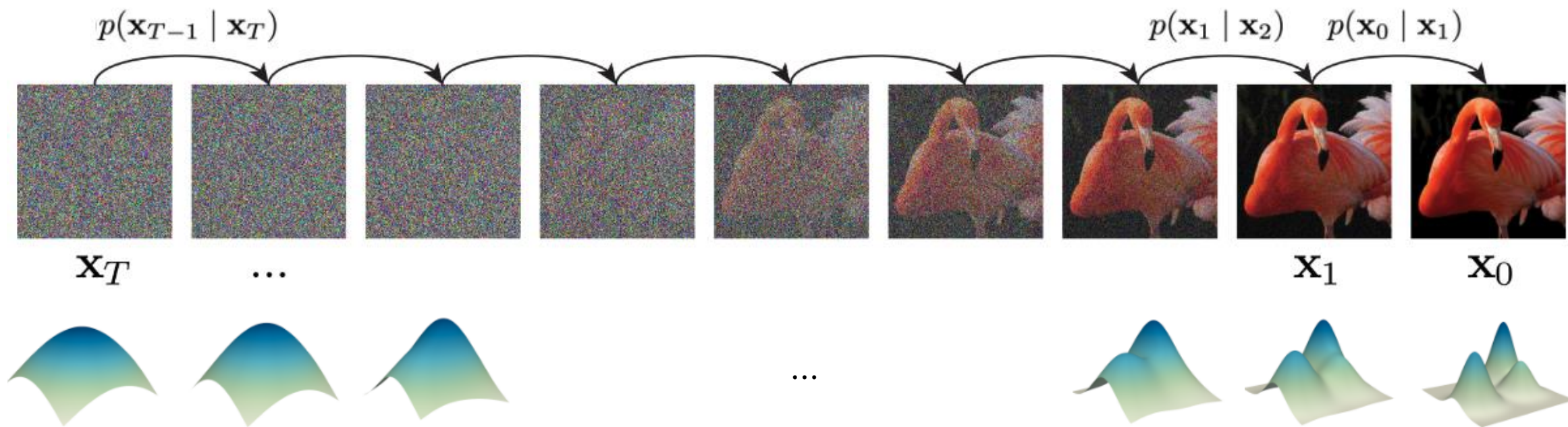
$$\text{with } p(\mathbf{z}) = p(\mathbf{z}_1)p(\mathbf{z}_2 | \mathbf{z}_1) \dots p(\mathbf{z}_n | \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_{n-1})$$

$$p(\mathbf{x}_{0:T}) = p(\mathbf{x}_T)p(\mathbf{x}_{T-1} | \mathbf{x}_T) \dots p(\mathbf{x}_1 | \mathbf{x}_2)p(\mathbf{x}_0 | \mathbf{x}_1)$$



Modeling Transitions of Two Distributions

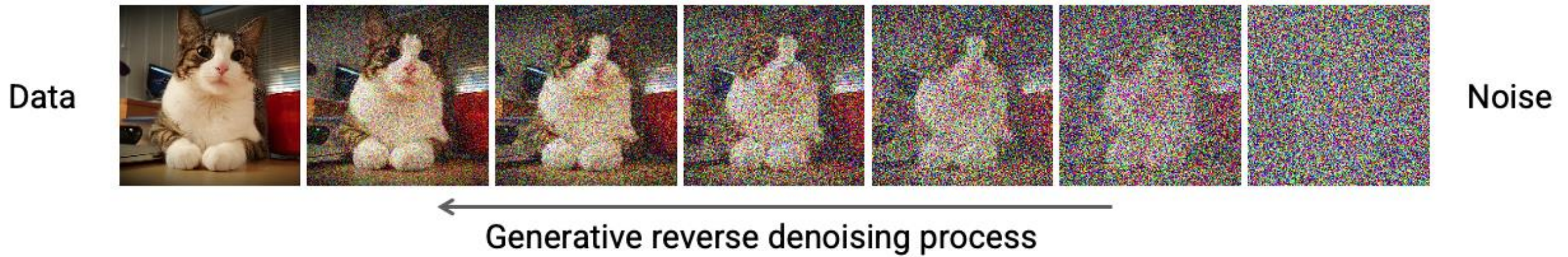
$$p(\mathbf{x}_{0:T}) = p(\mathbf{x}_T)p(\mathbf{x}_{T-1} | \mathbf{x}_T)\dots p(\mathbf{x}_1 | \mathbf{x}_2)p(\mathbf{x}_0 | \mathbf{x}_1)$$



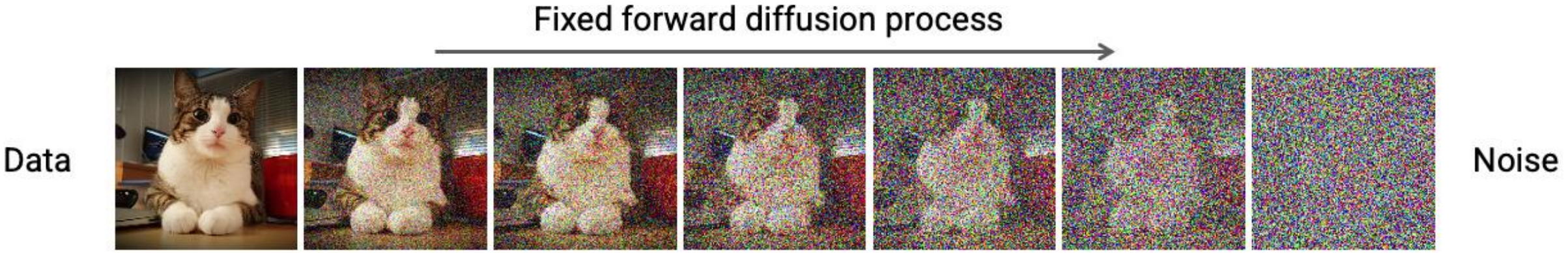
Model How Data is Transported Among the Prior Distribution $p(\mathbf{z})$ and the Target Distribution $p(\mathbf{x})$



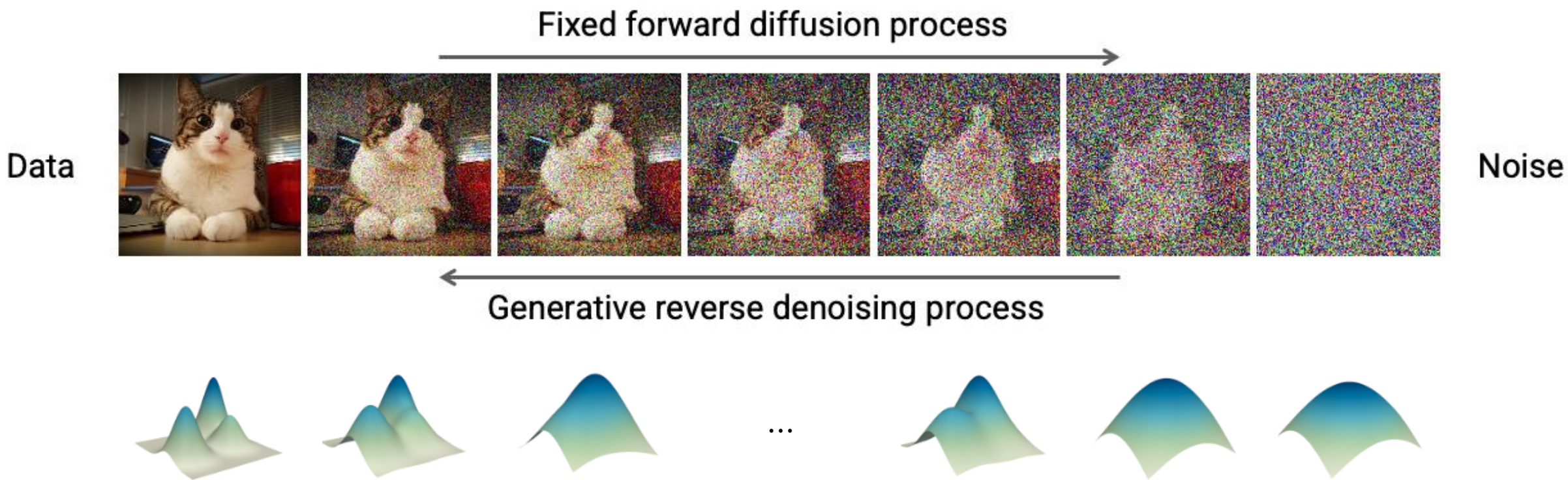
Generation: Map a Prior Distribution $p(\mathbf{z})$ to the Data Distribution $p(\mathbf{x})$



Basic Idea: First, Determine How the Data Distribution $p(\mathbf{x})$ is Mapped to a Prior Distribution $p(\mathbf{z})$

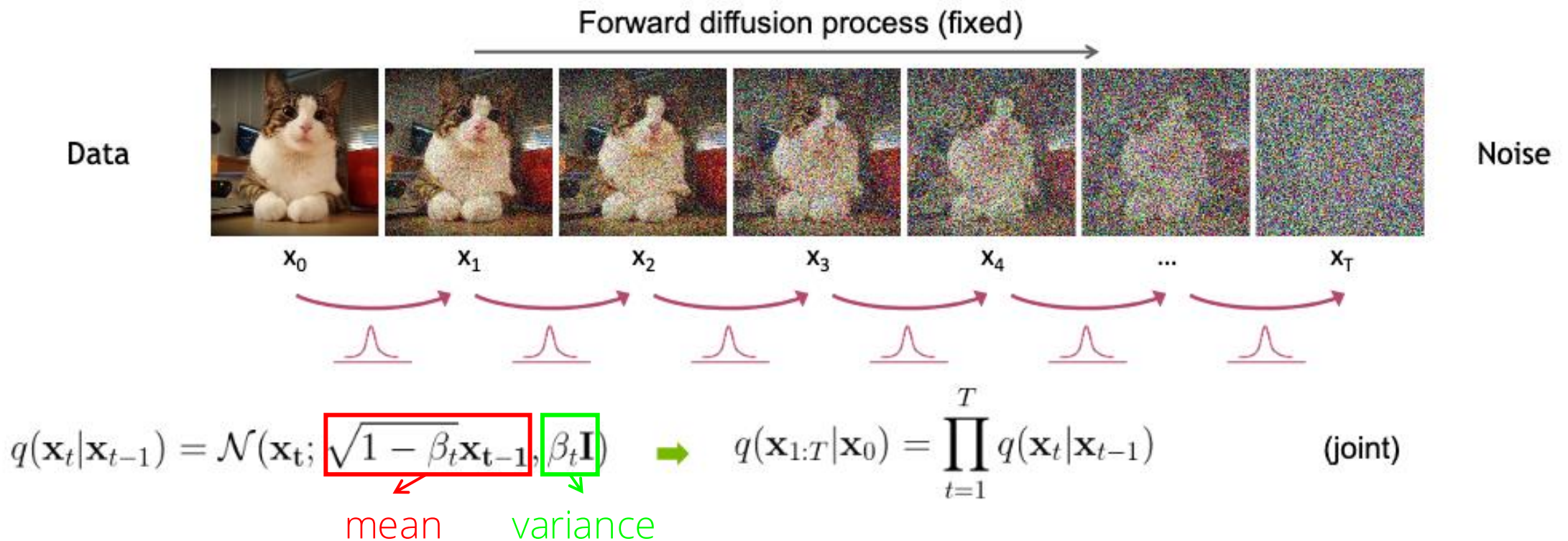


Basic Idea: Then, Learn to Reverse the Mapping



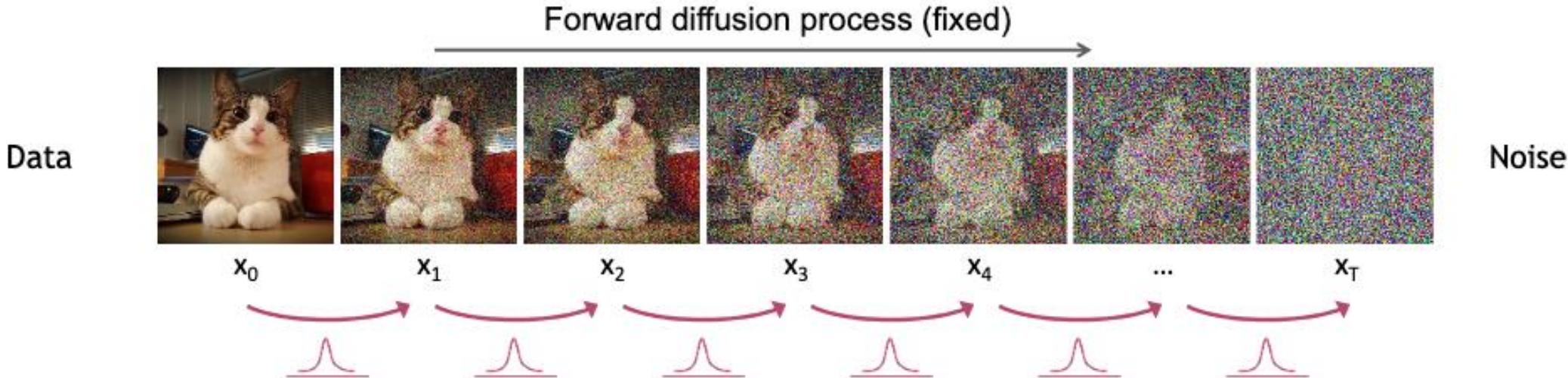
Denoising Diffusion Probabilistic Models Specify the Forward Mapping with Gaussian Diffusion Process

The formal definition of the forward process in T steps:



Denoising Diffusion Probabilistic Models Specify the Forward Mapping with Gaussian Diffusion Process

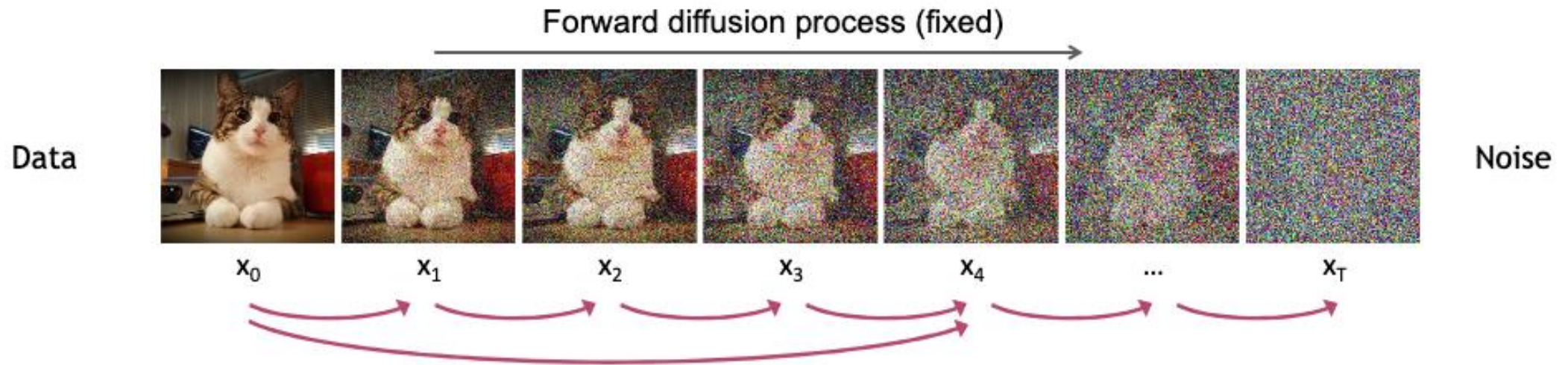
The formal definition of the forward process in T steps:



$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad \rightarrow \quad q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad \text{(joint)}$$

← schedule ← mean ← variance

Denoising Diffusion Probabilistic Models Specify the Forward Mapping with Gaussian Diffusion Process



Define $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$ \longrightarrow $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$ (Diffusion Kernel)

For sampling: $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon$ where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

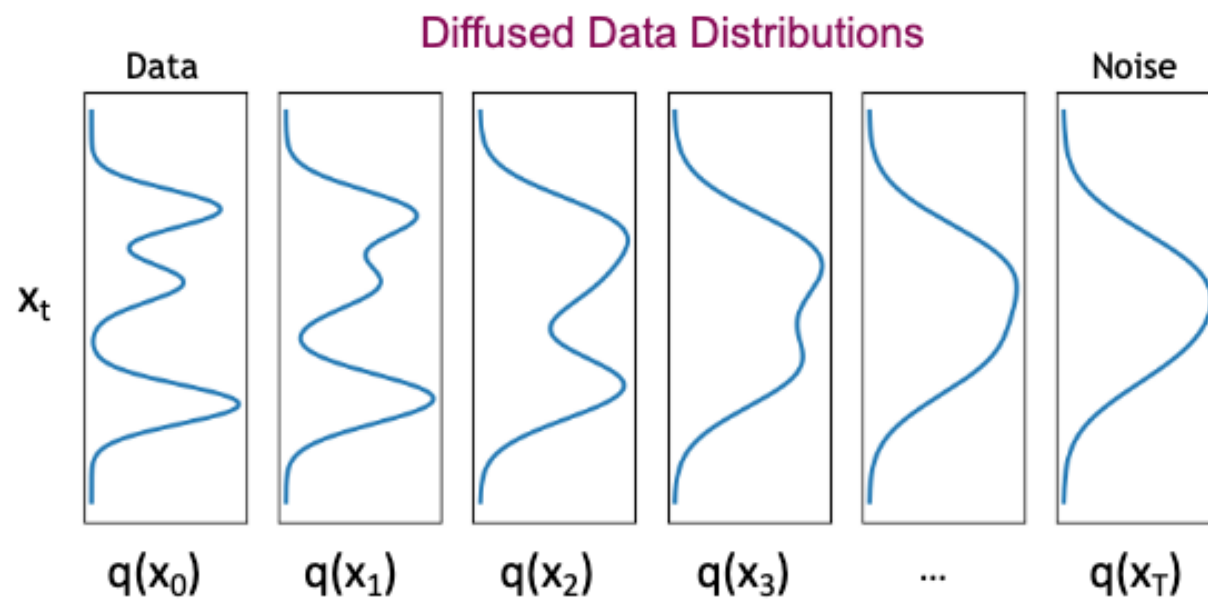
β_t values schedule (i.e., the noise schedule) is designed such that $\bar{\alpha}_T \rightarrow 0$ and $q(\mathbf{x}_T | \mathbf{x}_0) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

What happens to a distribution in the forward diffusion?

So far, we discussed the diffusion kernel $q(\mathbf{x}_t|\mathbf{x}_0)$ but what about $q(\mathbf{x}_t)$?

$$\underbrace{q(\mathbf{x}_t)}_{\text{Diffused data dist.}} = \int \underbrace{q(\mathbf{x}_0, \mathbf{x}_t)}_{\text{Joint dist.}} d\mathbf{x}_0 = \int \underbrace{q(\mathbf{x}_0)}_{\text{Input data dist.}} \underbrace{q(\mathbf{x}_t|\mathbf{x}_0)}_{\text{Diffusion kernel}} d\mathbf{x}_0$$

The diffusion kernel is Gaussian convolution.



We can sample $\mathbf{x}_t \sim q(\mathbf{x}_t)$ by first sampling $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ and then sampling $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0)$ (i.e., ancestral sampling).

Generative Learning by Denoising

Recall, that the diffusion parameters are designed such that $q(\mathbf{x}_T) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

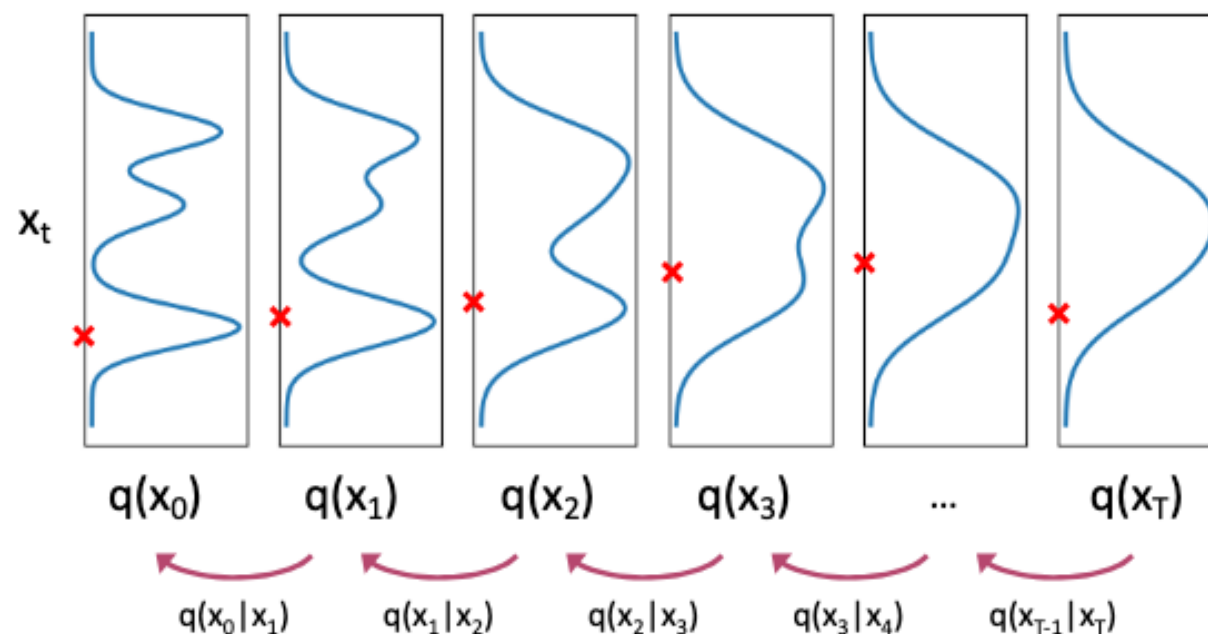
Generation:

Sample $\mathbf{x}_T \sim \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

Iteratively sample $\mathbf{x}_{t-1} \sim \underbrace{q(\mathbf{x}_{t-1}|\mathbf{x}_t)}$

True Denoising Dist.

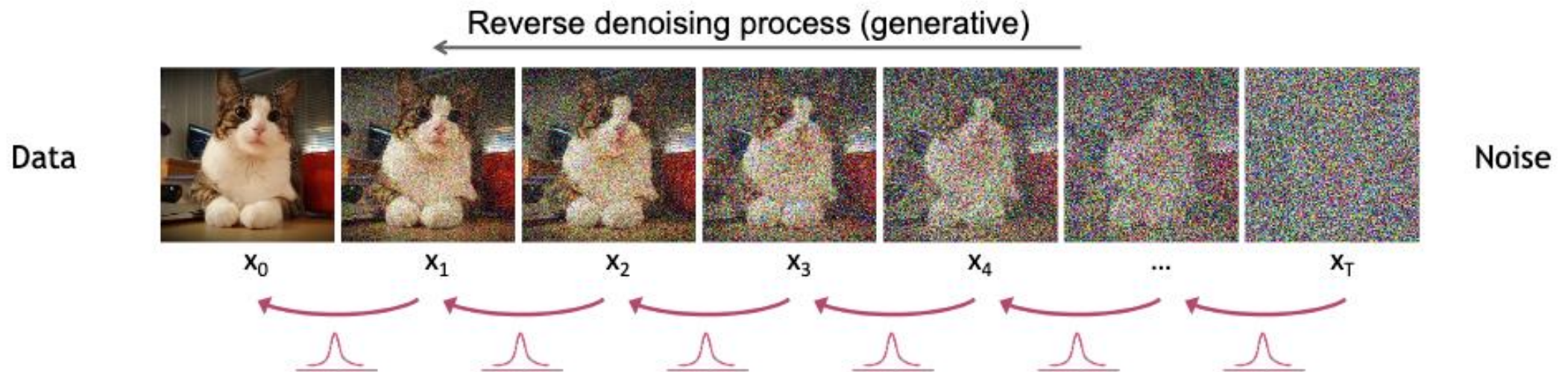
Diffused Data Distributions



Can we approximate $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$? Yes, we can use a **Normal distribution** if β_t is small in each forward diffusion step.

We Directly Have Reverse Mappings of Gaussian Diffusion Process

Formal definition of forward and reverse processes in T steps:

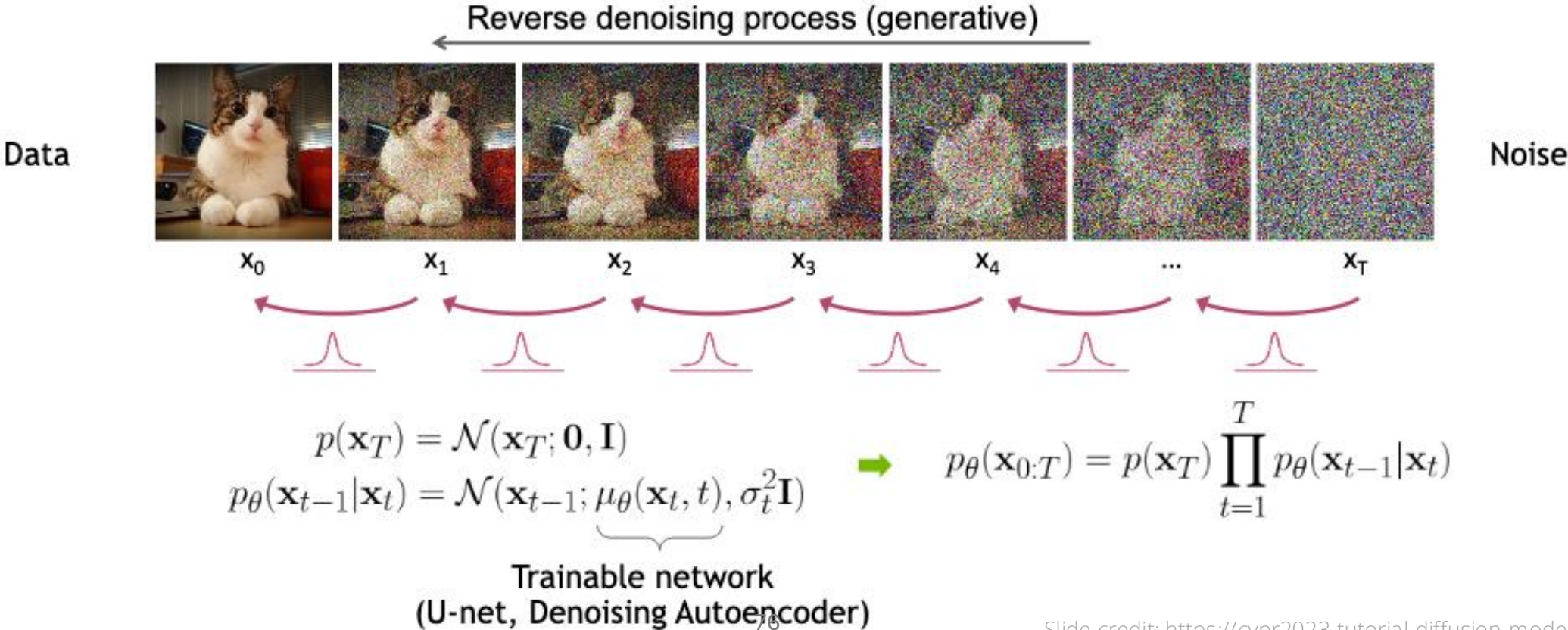


$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\boldsymbol{\beta}}_t \mathbf{I}),$$

$$\text{where } \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) := \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t \quad \text{and} \quad \tilde{\boldsymbol{\beta}}_t := \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

Learn a Model to Reverse Mappings of Gaussian Diffusion Process

Formal definition of forward and reverse processes in T steps:



Learning Denoising Model

Variational upper bound

For training, we can form variational upper bound that is commonly used for training variational autoencoders:

$$\mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] =: L$$

Recall that $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon$. [Ho et al. NeurIPS 2020](#) parameterized the mean of denoising model via:

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{1 - \beta_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right)$$

Using a few simple arithmetic operations, we can write down the variational objective as:

$$L = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), t \sim \mathcal{U}\{1, T\}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\lambda_t \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2 \right]$$

[Ho et al. NeurIPS 2020](#) observe that simply setting λ_t to 1 for all t works best in practice.

Summary

Algorithm 1 Training

- 1: **repeat**
 - 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 4: $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 5: Take gradient descent step on
 $\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$
 - 6: **until** converged
-

Algorithm 2 Sampling

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 2: **for** $t = T, \dots, 1$ **do**
 - 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
 - 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
 - 5: **end for**
 - 6: **return** \mathbf{x}_0
-

Diffusion Models are Awesome

The screenshot displays the AudioBox interface with a timeline from 0:00 to 0:25. The interface includes a top navigation bar with icons for 'Add speech', 'Add sound effect', 'Play', and 'Download'. A purple tip banner contains three instructions: 'Tip 1. Add new audio for your story', 'Tip 2. Click audio clips below to edit, and drag to rearrange', and 'Tip 3. Listen to your audio story and download it'. The main workspace is divided into tracks for 'Jamie', 'Robot Assistant', 'broken engine clanking mak...', and 'A sci fi swoosh space sound...'. Each track contains audio clips with text labels and waveform visualizations. A speaker icon is visible on the right side of the interface.

0:00 0:05 0:10 0:15 0:20 0:25

Jamie

- Oh no, the engine isnt...
- Mm...
- Mm...
- No, I haven't...

Robot Assistant

- Have you tried tighten...
- Have you properly aligned ...

broken engine clanking mak...

- broken engine clanking making loud bangs, followed by a big puff ...

A sci fi swoosh space sound...

- A sci fi swoosh space sound followed by a loud bang and silence

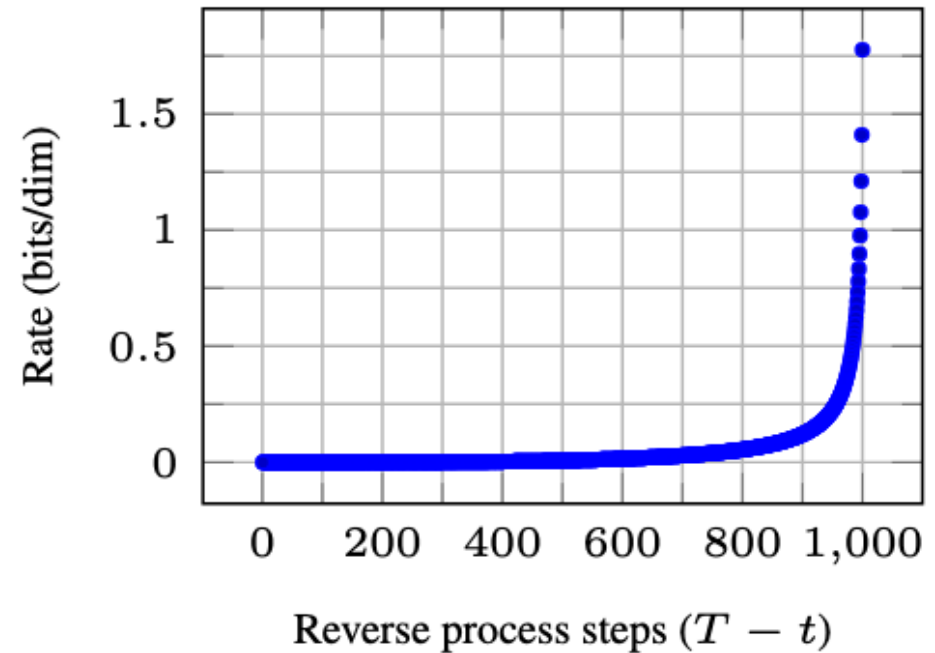
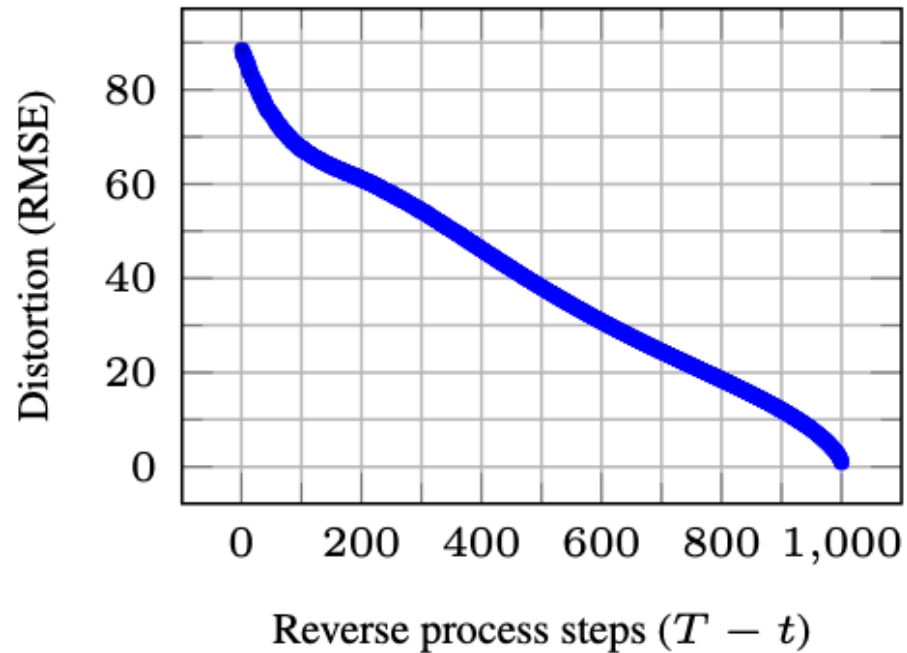
Diffusion Models are Awesome



Diffusion Models are Awesome



Diffusion Models Are Slow...

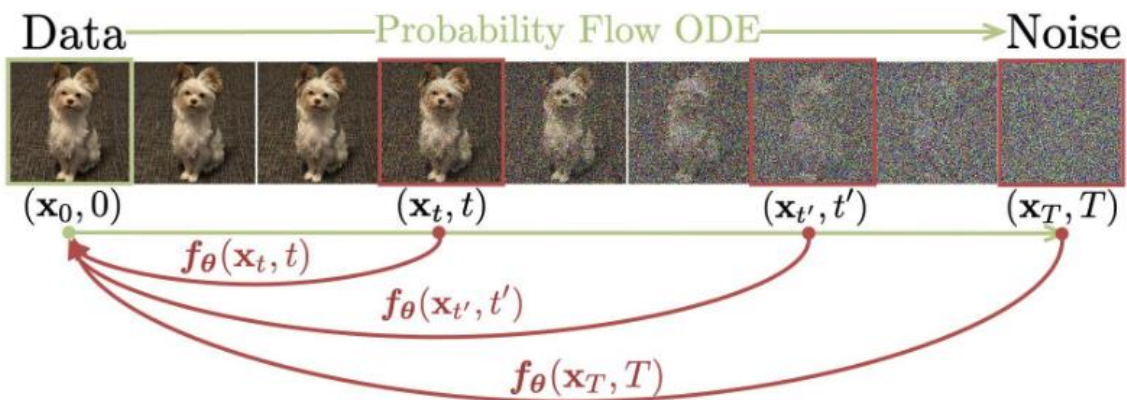


DDPM often requires 1000 samplings steps to generate high-quality data....

Many Techniques Introduced to Accelerate DDPM

Consistency Models

Yang Song¹ Prafulla Dhariwal¹ Mark Chen¹ Ilya Sutskever¹



DENOISING DIFFUSION IMPLICIT MODELS

Jiaming Song, Chenlin Meng & Stefano Ermon

Stanford University

{tsong, chenlin, ermon}@cs.stanford.edu

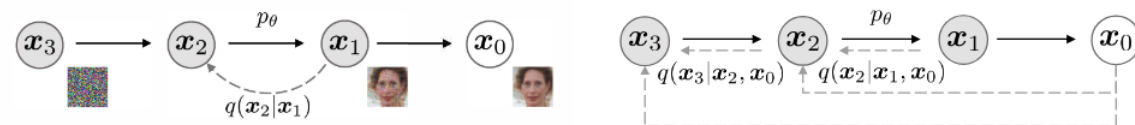


Figure 1: Graphical models for diffusion (left) and non-Markovian (right) inference models.

Content

- Generative Modeling
 - Autoregressive Models
 - Variational AutoEncoders
 - Denoising Diffusion Probabilistic Models
 - Flow Matching Models
- 4D Generative Modeling
 - 3D Generation with Multi-view Video Generation
 - 4D Generation with Multi-view Video Generation
 - 3D Geometry Generation

A Visual Diagram of Distributional Transportation

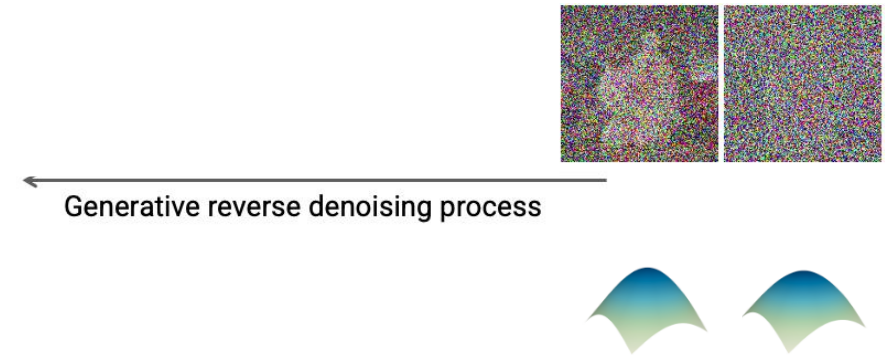
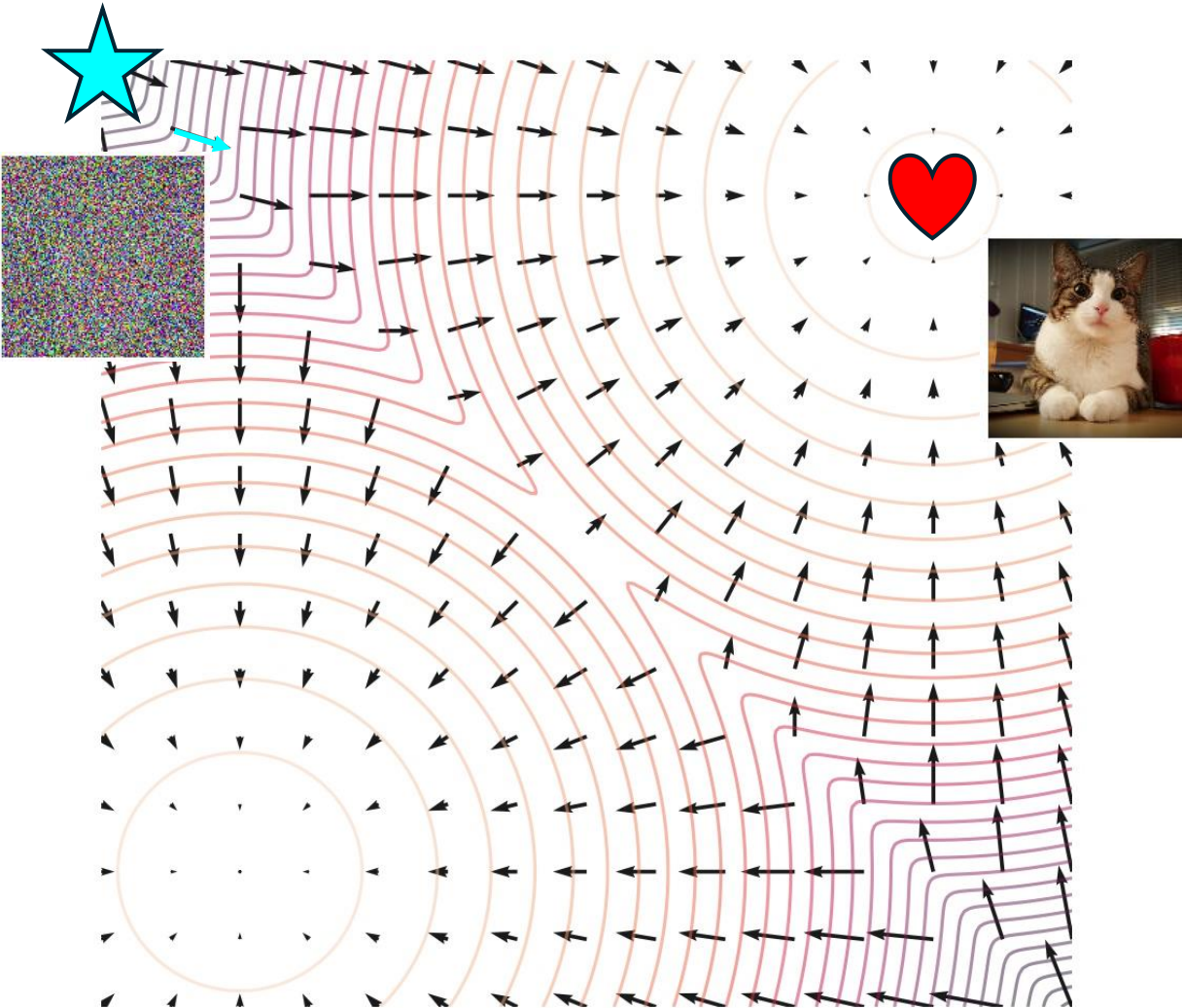


Image credit: <https://yang-song.net/blog/2021/score/>

A Visual Diagram of Distributional Transportation

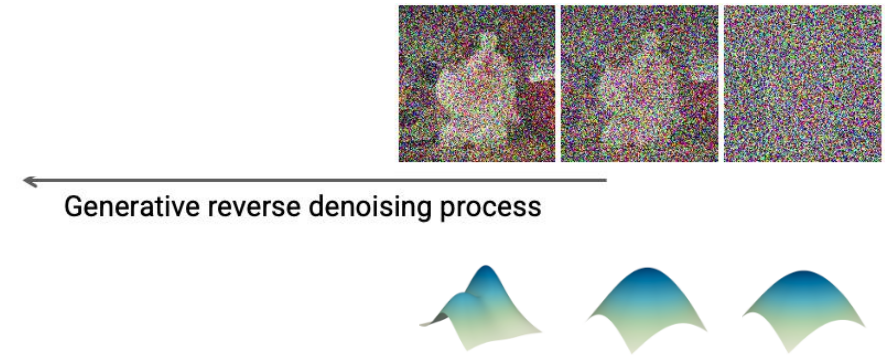
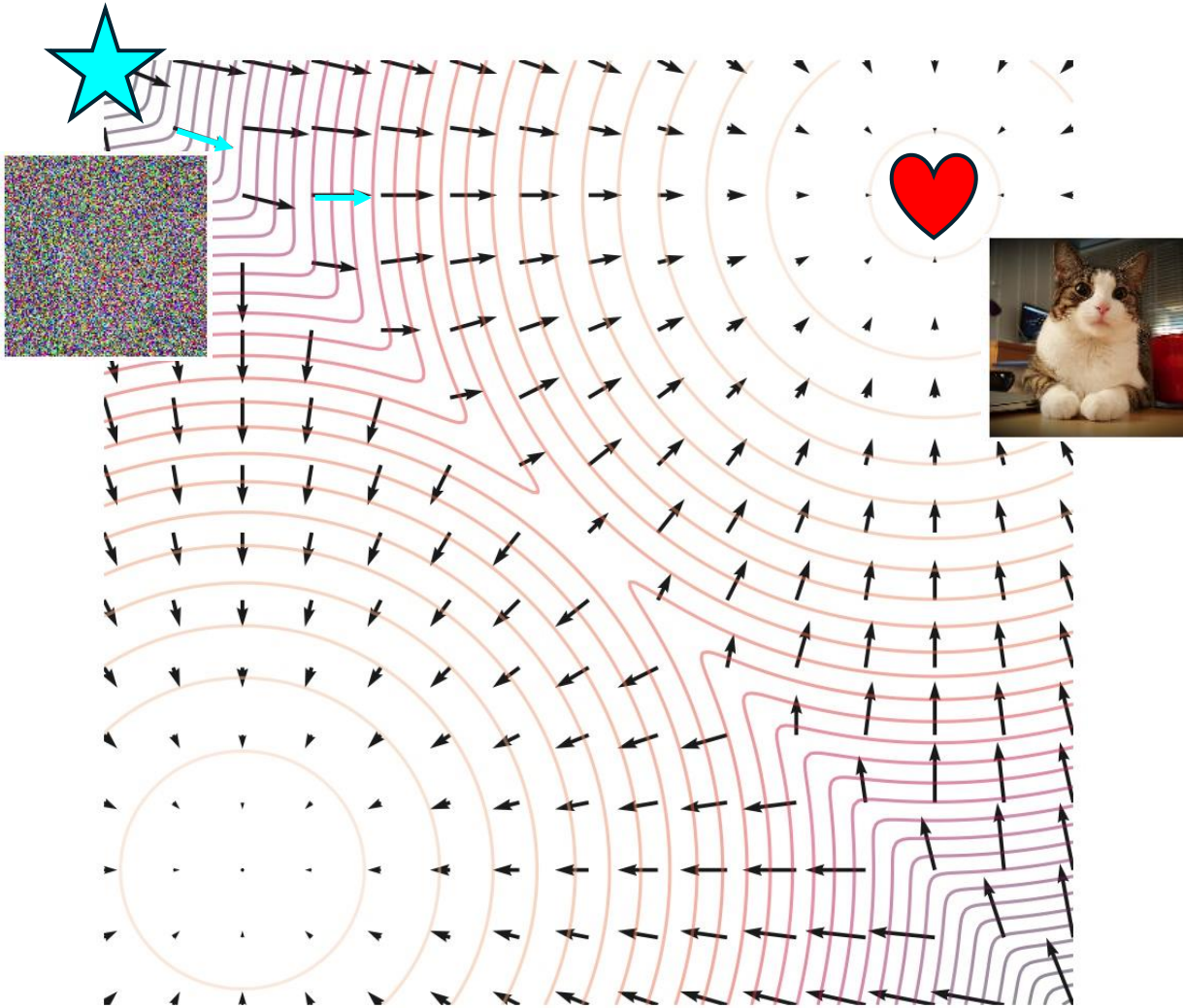


Image credit: <https://yang-song.net/blog/2021/score/>

A Visual Diagram of Distributional Transportation

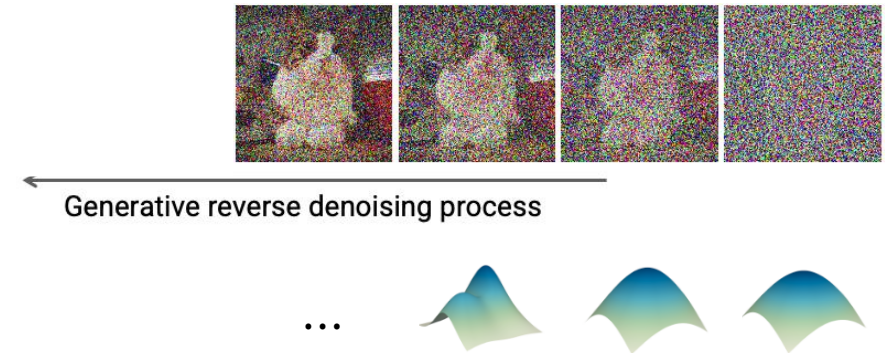
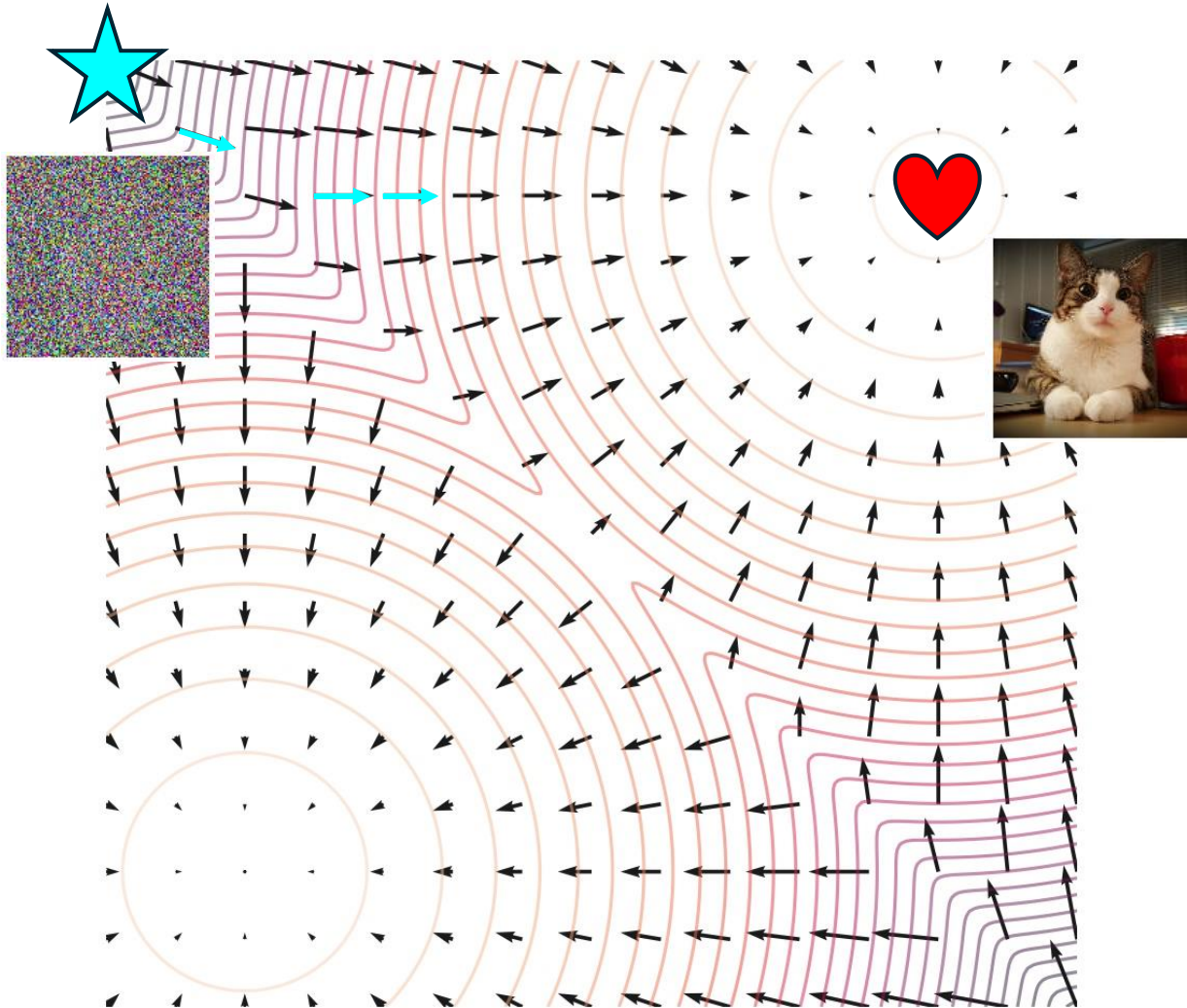


Image credit: <https://yang-song.net/blog/2021/score/>

A Visual Diagram of Distributional Transportation

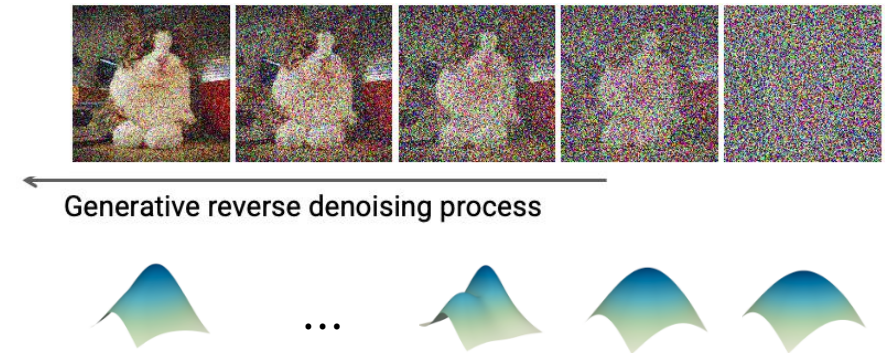
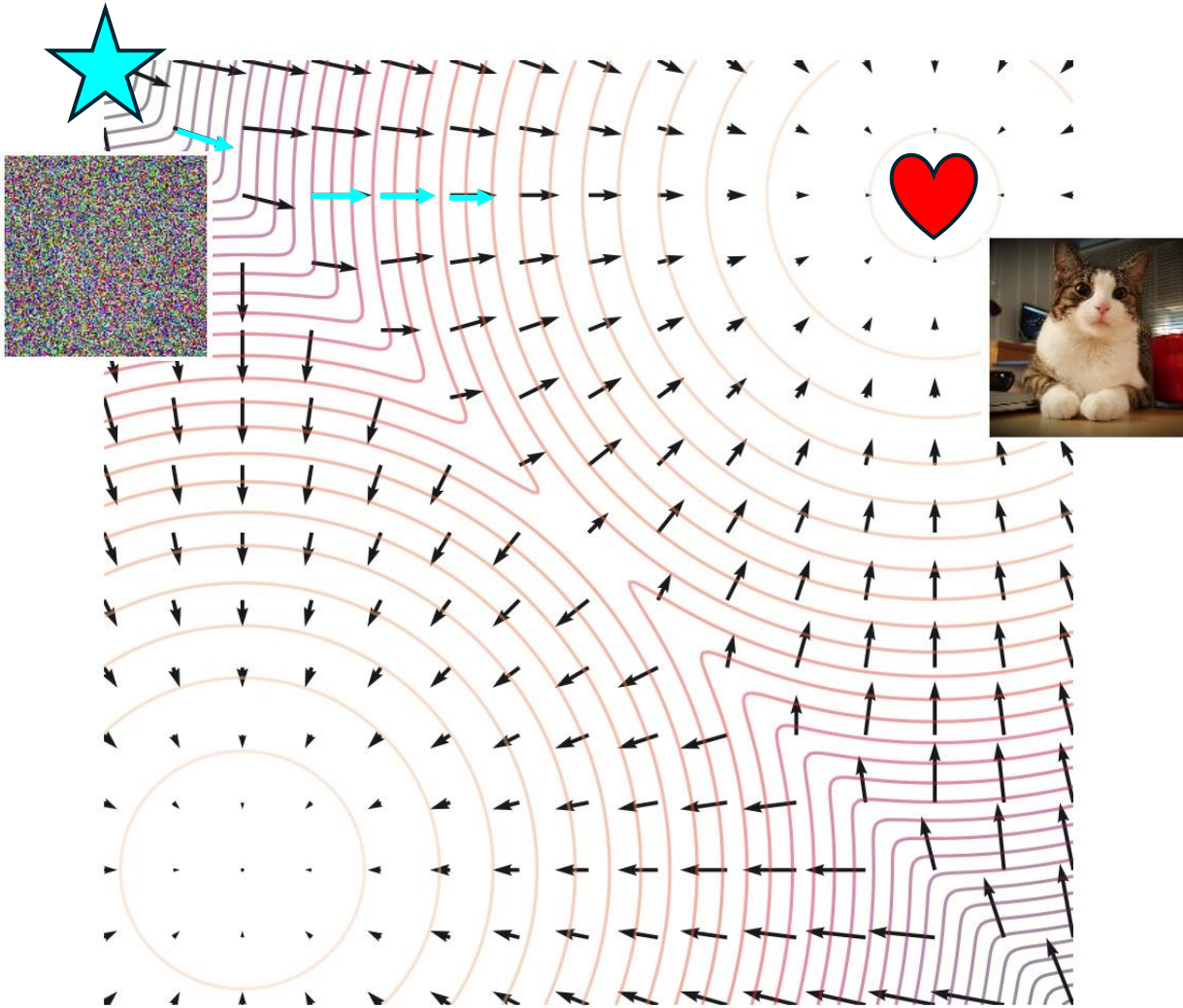
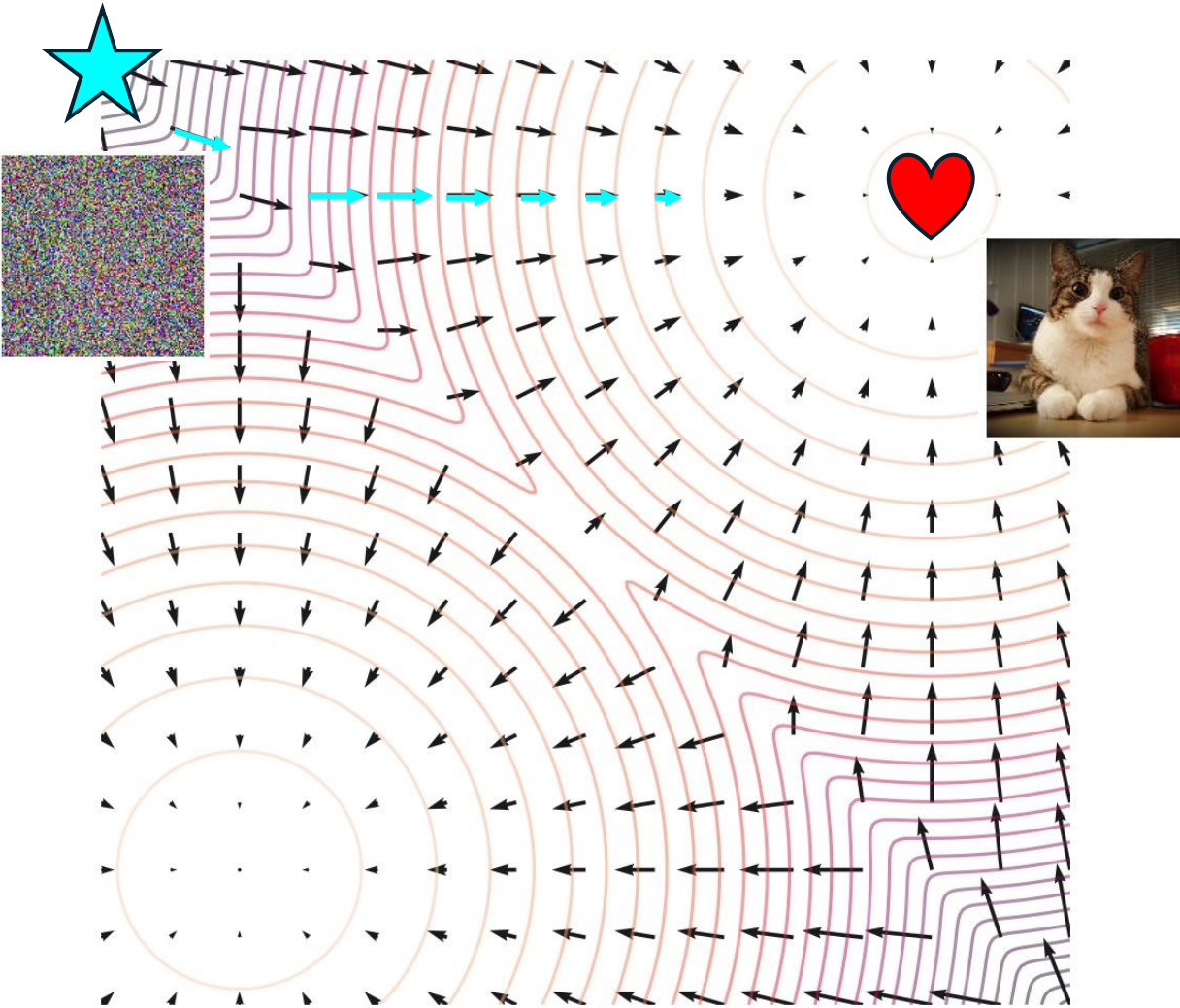


Image credit: <https://yang-song.net/blog/2021/score/>

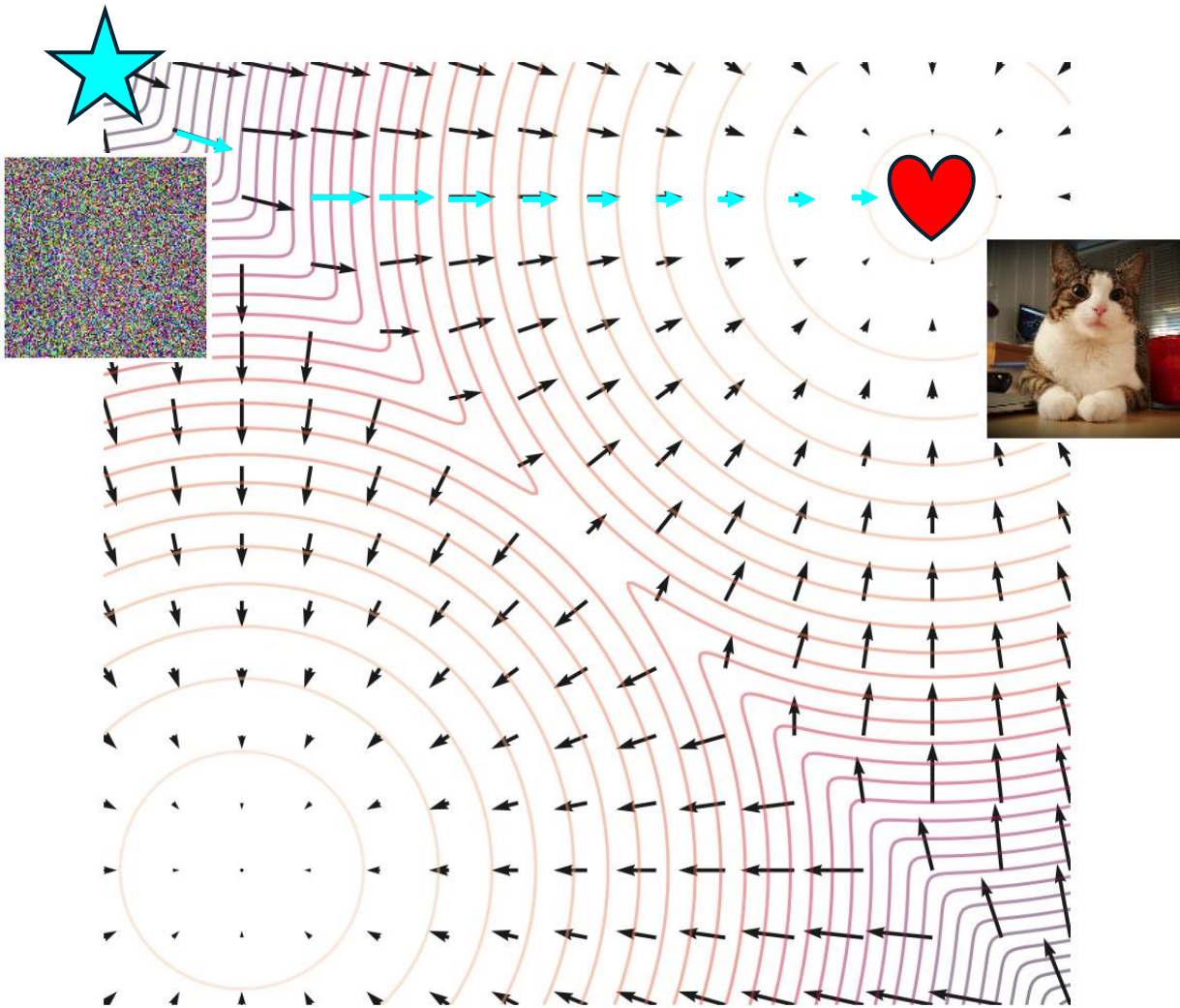
A Visual Diagram of Distributional Transportation



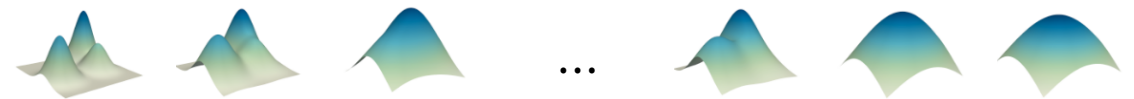
Generative reverse denoising process



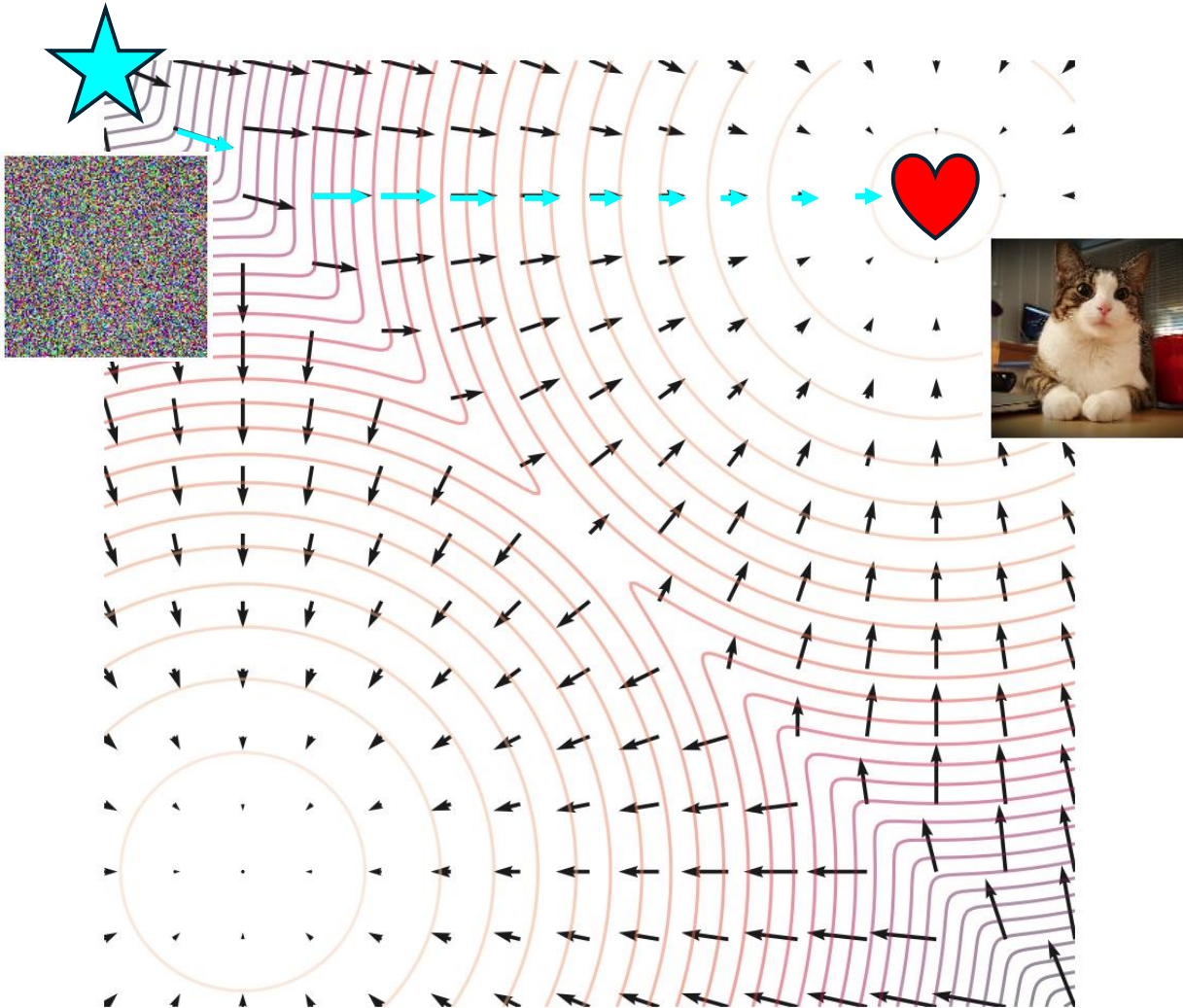
A Visual Diagram of Distributional Transportation



Generative reverse denoising process



A Visual Diagram of Distributional Transportation



← Generative reverse denoising process

- A probabilistic view of generation



$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$$

- A flow view of generation

$$\mathbf{x}_0 = \mathbf{x}_1 + \int_1^0 \mathbf{v}(\mathbf{x}_t, t) dt$$

This is the core idea of flow matching

DDPM vs. Flow Matching

DDPM

- A probabilistic view of generation
- Forward process: Gaussian diffusion process

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

- Backward process: learn a denoising model

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

- Sampling steps are discrete: $\mathbf{t} \in [0, T]$
- Prior distribution is Gaussian

Flow Matching

- A flow view of generation
- Forward process: Any vector field

$$\mathbf{x}_1 = \mathbf{x}_0 - \int_1^0 \mathbf{v}(\mathbf{x}_t, t) dt$$

- Backward process: learn a regression model

$$\mathbf{x}_0 = \mathbf{x}_1 + \int_1^0 \mathbf{u}_\theta(\mathbf{x}_t, t) dt$$

- Sampling steps are continuous: $\mathbf{t} \sim \text{Uniform}(0, 1)$
- Prior distribution is not restricted to Gaussian

DDPM vs. Flow Matching

Flow matching often assign $\mathbf{t} = \mathbf{0}$ to the prior and $\mathbf{t} = \mathbf{1}$ to the data distribution

DDPM

- A probabilistic view of generation
- Forward process: Gaussian diffusion process

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

- Backward process: learn a denoising model

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

- Sampling steps are discrete: $\mathbf{t} \in [0, \mathbf{T}]$
- Prior distribution is Gaussian

Flow Matching

- A flow view of generation
- Forward process: Any vector field

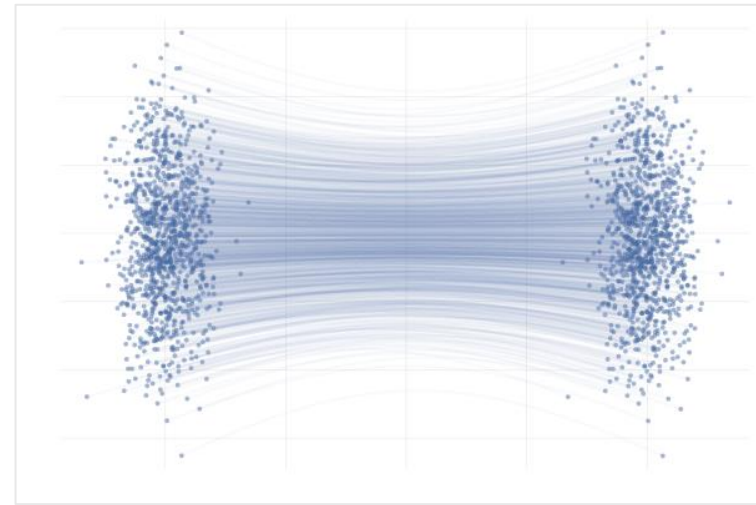
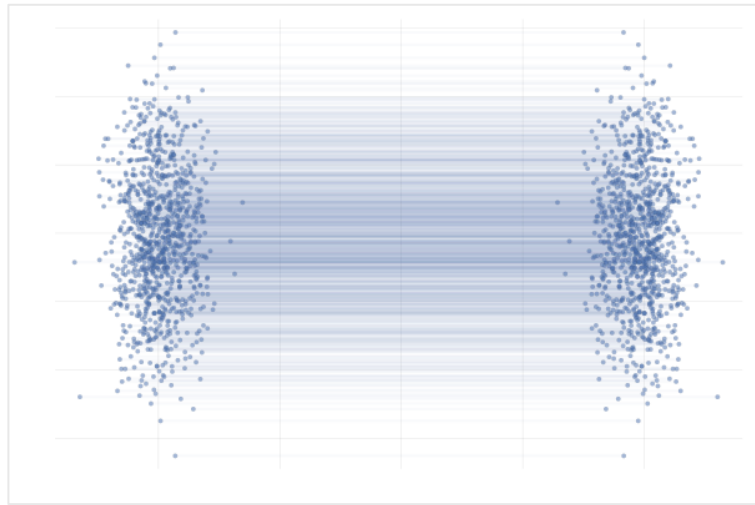
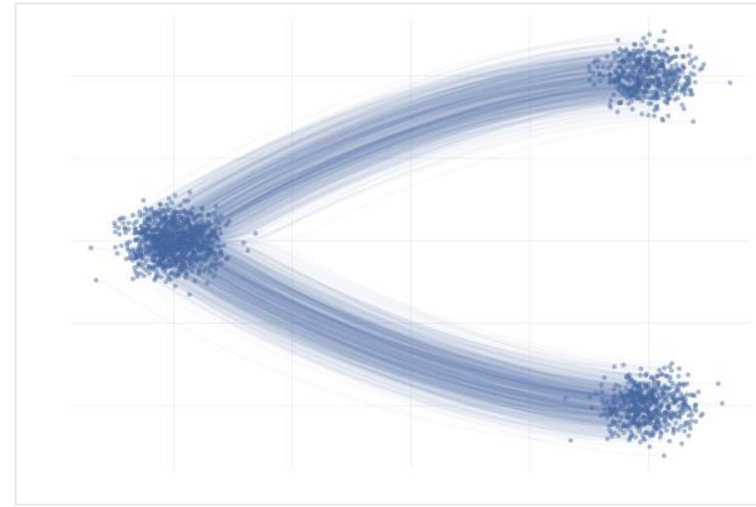
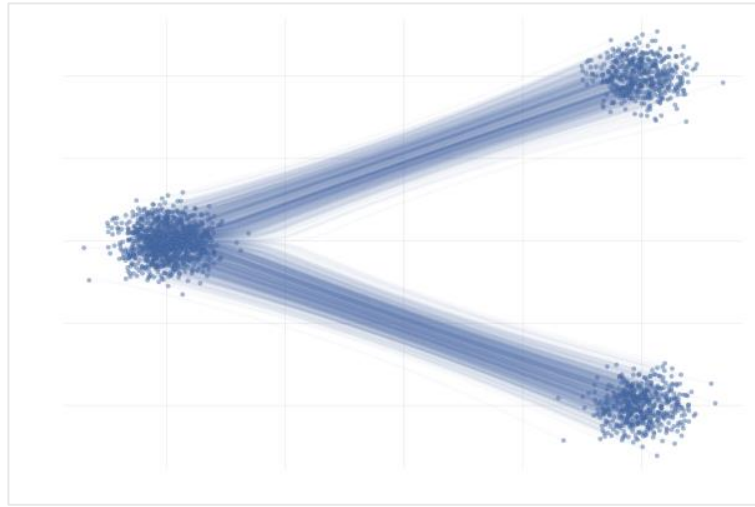
$$\mathbf{x}_0 = \mathbf{x}_1 - \int_0^1 \mathbf{v}(\mathbf{x}_t, t) dt$$

- Backward process: learn a regression model

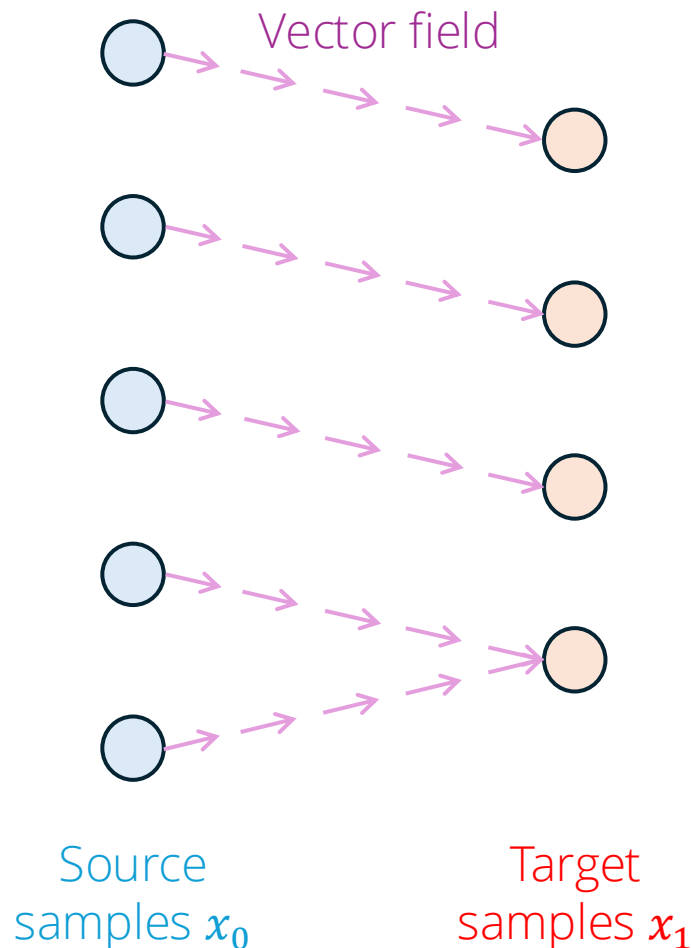
$$\mathbf{x}_1 = \mathbf{x}_0 + \int_0^1 \mathbf{u}_\theta(\mathbf{x}_t, t) dt$$

- Sampling steps are continuous: $\mathbf{t} \sim \text{Uniform}(0, 1)$
- Prior distribution is not restricted to Gaussian

A More General Formulation of Generation: Specify Any Vector Field as the Transportation Map

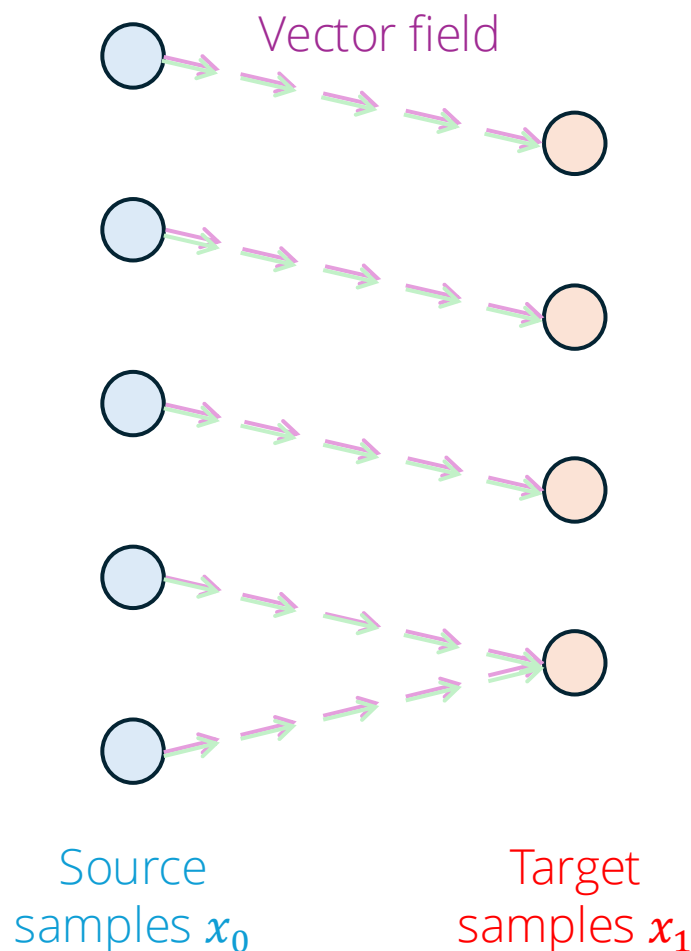


However, Specification of the Vector Field is Intractable...



- A vector field $u_t(x)$ transports source samples to target samples
- The objective of Flow Matching: regresses the vector field with a neural network $v_\theta(t, x)$
$$L_{FM} = \mathbb{E}[\|u_t(x) - v_\theta(t, x)\|^2]$$
- However, $u_t(x)$ is often intractable to obtain
 - The vector field pairs every data with every samples drawn from the prior distribution

Fortunately, the Vector Field can be Composed of a Mixture of Conditional Vector Field



- A vector field $u_t(x)$ transports source samples to target samples
- The objective of Flow Matching: regresses the vector field with a neural network $v_\theta(t, x)$
- Vector field $u_t(x)$ can be composed of a mixture of conditional vector field $u_t(x|x_1)$

$$L_{FM} = \mathbb{E}[\|u_t(x) - v_\theta(t, x)\|^2]$$

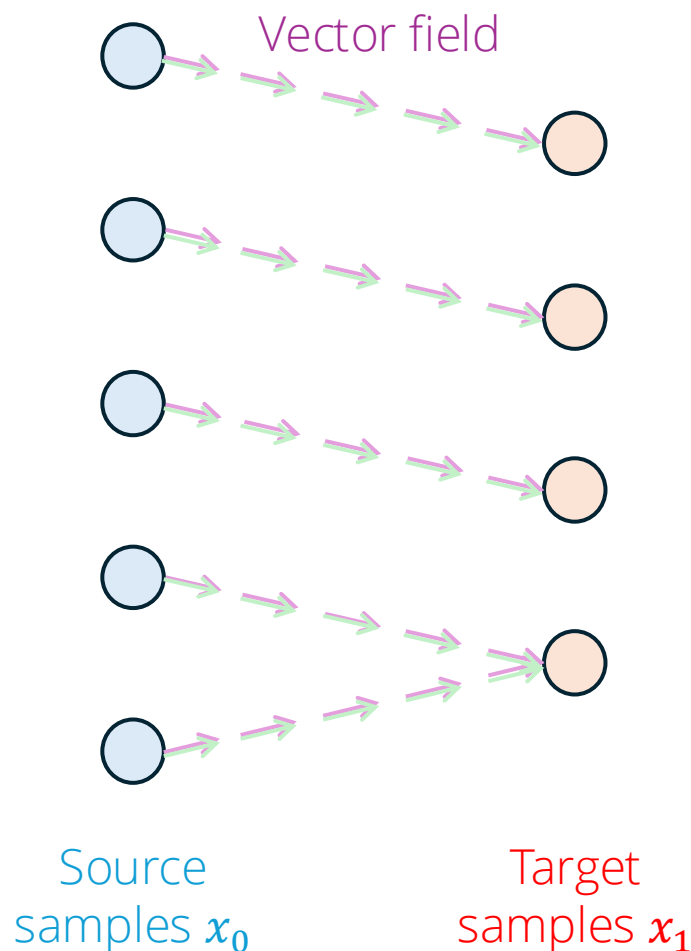
$$u_t(x) = \mathbb{E}_{p(x_1)} \left[\frac{u_t(x|x_1) p_t(x|x_1)}{p_t(x)} \right]$$

$u_t(x|x_1)$: the conditional flow between sample x and data x_1

$p_t(x|x_1)$: the joint distribution of sample x and data x_1

$u_t(x)$: the flow at sample x during time t

Fortunately, the Vector Field can be Composed of a Mixture of Conditional Vector Field



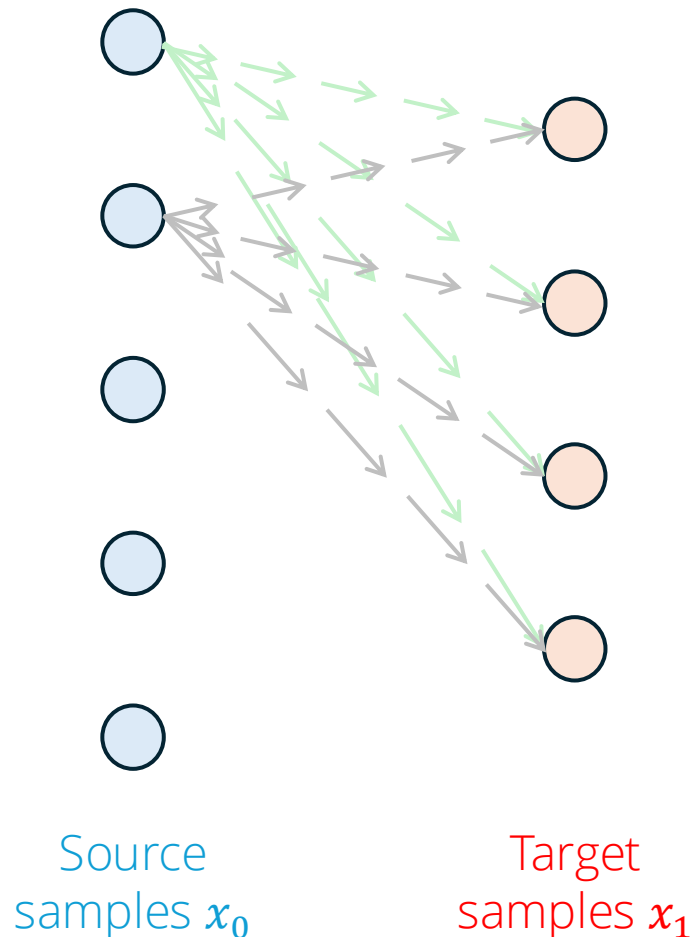
- A vector field $u_t(x)$ transports source samples to target samples
- The objective of Flow Matching: regresses the vector field with a neural network $v_\theta(t, x)$
- Vector field $u_t(x)$ can be composed of a mixture of conditional vector field $u_t(x|x_1)$

$$L_{FM} = \mathbb{E}[\|u_t(x) - v_\theta(t, x)\|^2]$$

$$u_t(x) = \mathbb{E}_{p(x_1)} \left[\frac{u_t(x|x_1) p_t(x|x_1)}{p_t(x)} \right]$$

- We have:
$$\nabla_\theta \mathbb{E}[\|u_t(x) - v_\theta(t, x)\|^2] = \nabla_\theta \mathbb{E}[\|u_t(x|x_1) - v_\theta(t, x)\|^2]$$

A Common Practice to Compose the Vector Field



- In practice, we often
 - Adopt Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{I})$ as the prior distribution
 - Randomly pair source samples and target samples
 - Assign optimal transportation (straight velocity) as the conditional flow for the source-target pair

$$u_t(x) = x_1 - x_0$$
$$x_t = (1 - t)x_0 + tx_1$$

Flow Matching: Learn to Regress the Vector Field

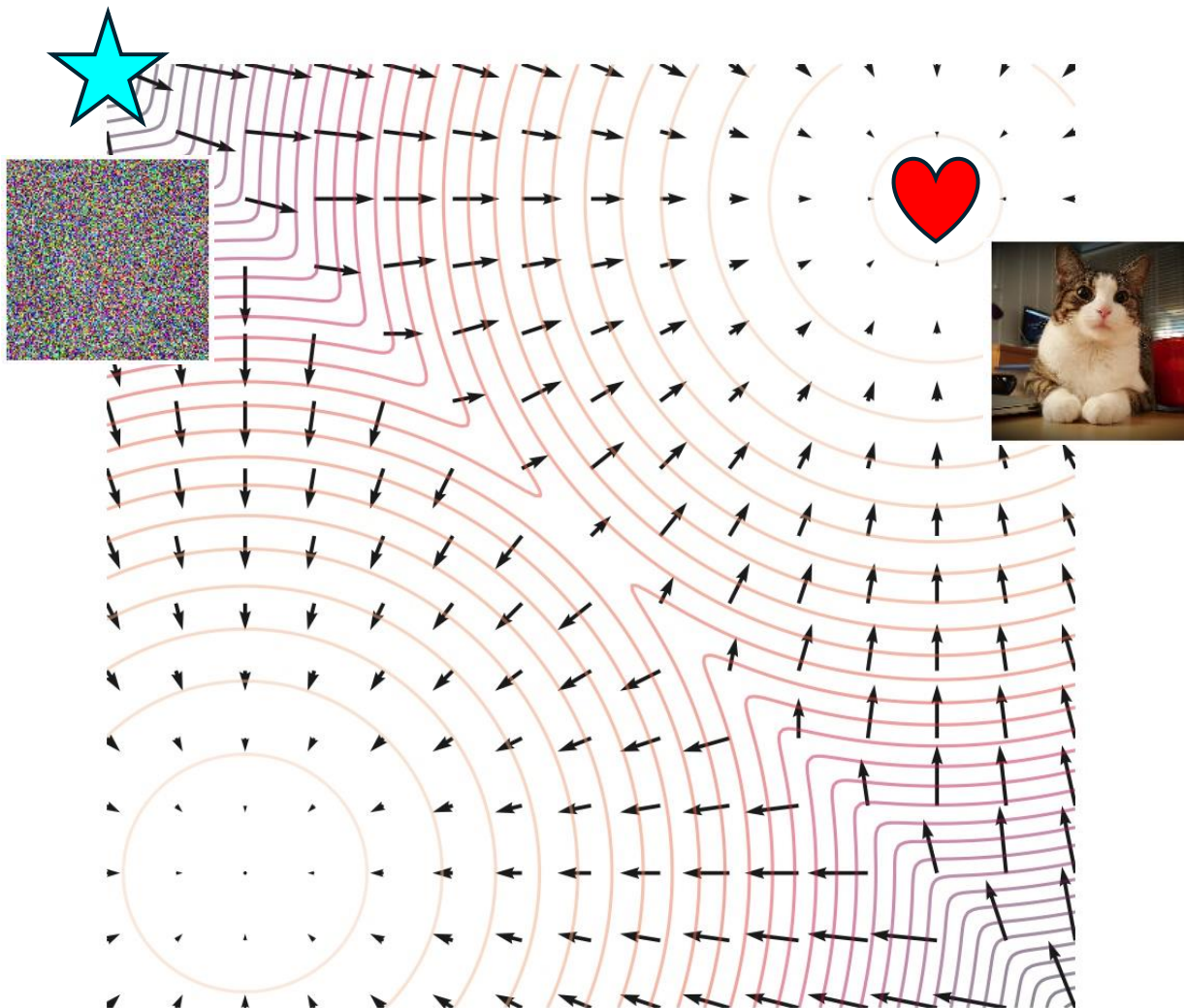
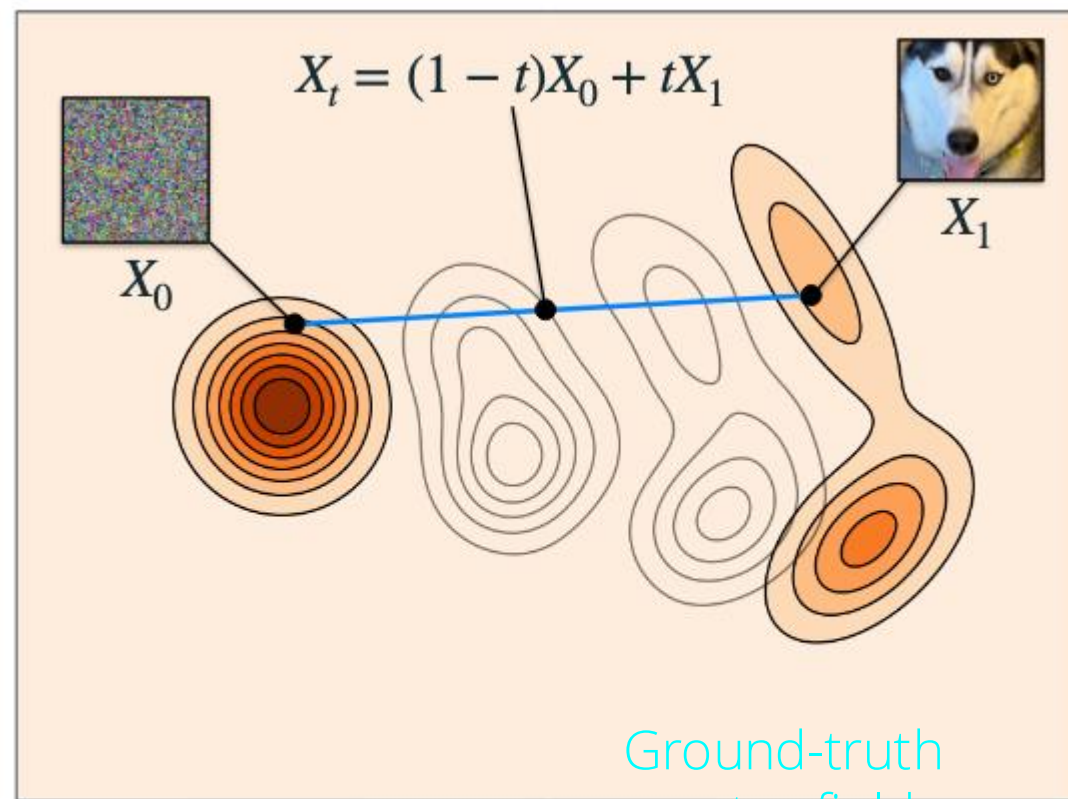


Image credit: <https://yang-song.net/blog/2021/score/>

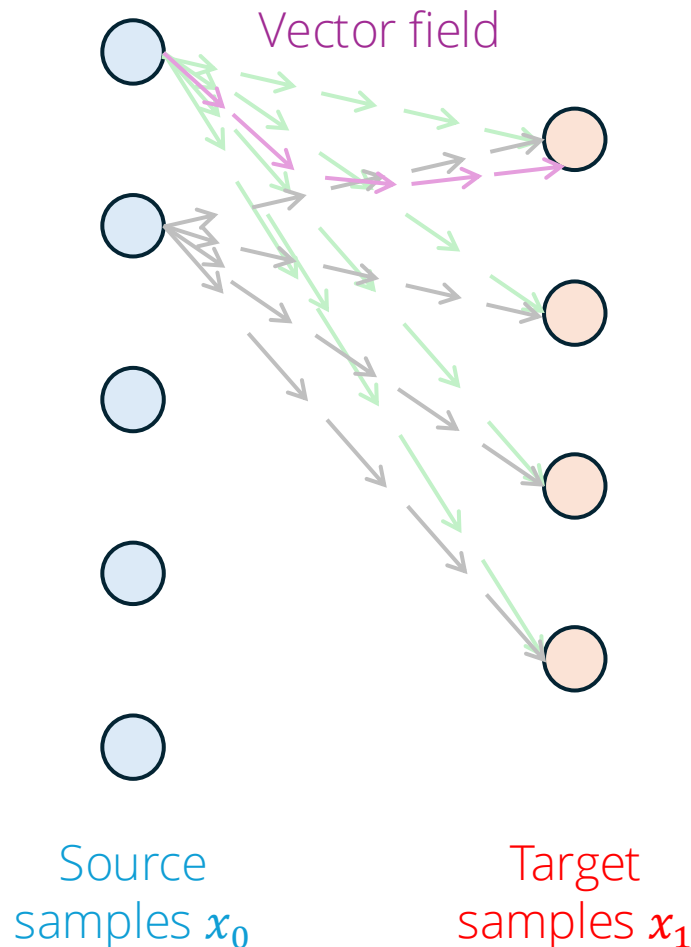


Ground-truth
vector field

$$\mathbb{E}_{t, X_0, X_1} \left\| \boxed{u_t^\theta(X_t)} - \boxed{X_1 - X_0} \right\|^2$$

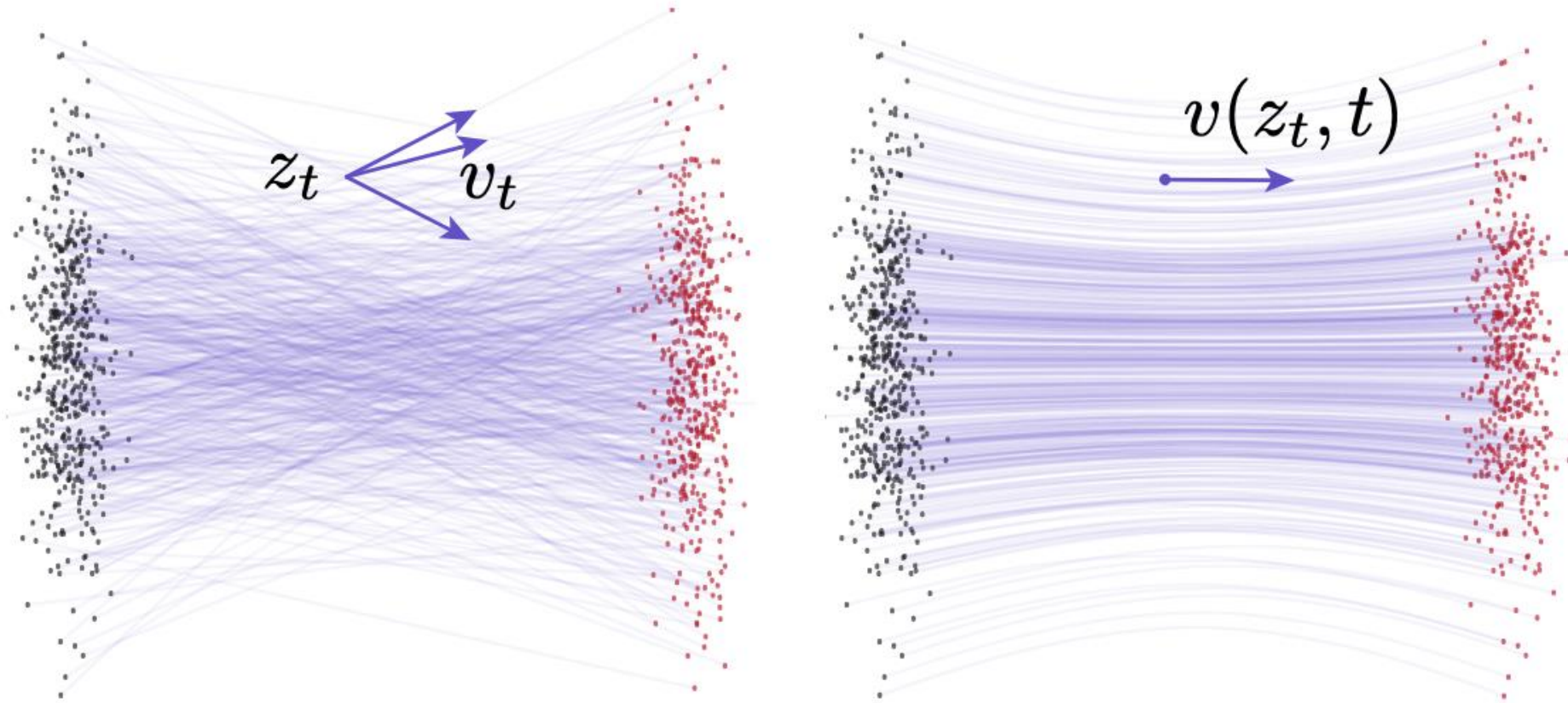
Estimated vector field

A Common Practice to Compose the Vector Field



- In practice, we often
 - Adopt Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{I})$ as the prior distribution
 - Randomly pair source samples and target samples
 - Assign optimal transportation (straight velocity) as the conditional flow for the source-target pair
- Since vector field $u_t(x)$ is composed of a mixture of conditional vector field $u_t(x|x_1)$, which becomes curved under random coupling

The Target Vector Field vs. Learned Vector Field



Flow Matching is a General Formulation that Recovers Diffusion Models if the Vector Field is Specified Correctly

Example I: Diffusion conditional VFs. Diffusion models start with data points and gradually add noise until it approximates pure noise. These can be formulated as stochastic processes, which have strict requirements in order to obtain closed form representation at arbitrary times t , resulting in Gaussian conditional probability paths $p_t(x|x_1)$ with specific choices of mean $\mu_t(x_1)$ and std $\sigma_t(x_1)$ (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2020b). For example, the reversed (noise→data) Variance Exploding (VE) path has the form

$$p_t(x) = \mathcal{N}(x|x_1, \sigma_{1-t}^2 I), \quad (16)$$

where σ_t is an increasing function, $\sigma_0 = 0$, and $\sigma_1 \gg 1$. Next, equation 16 provides the choices of $\mu_t(x_1) = x_1$ and $\sigma_t(x_1) = \sigma_{1-t}$. Plugging these into equation 15 of Theorem 3 we get

$$u_t(x|x_1) = -\frac{\sigma'_{1-t}}{\sigma_{1-t}}(x - x_1). \quad (17)$$

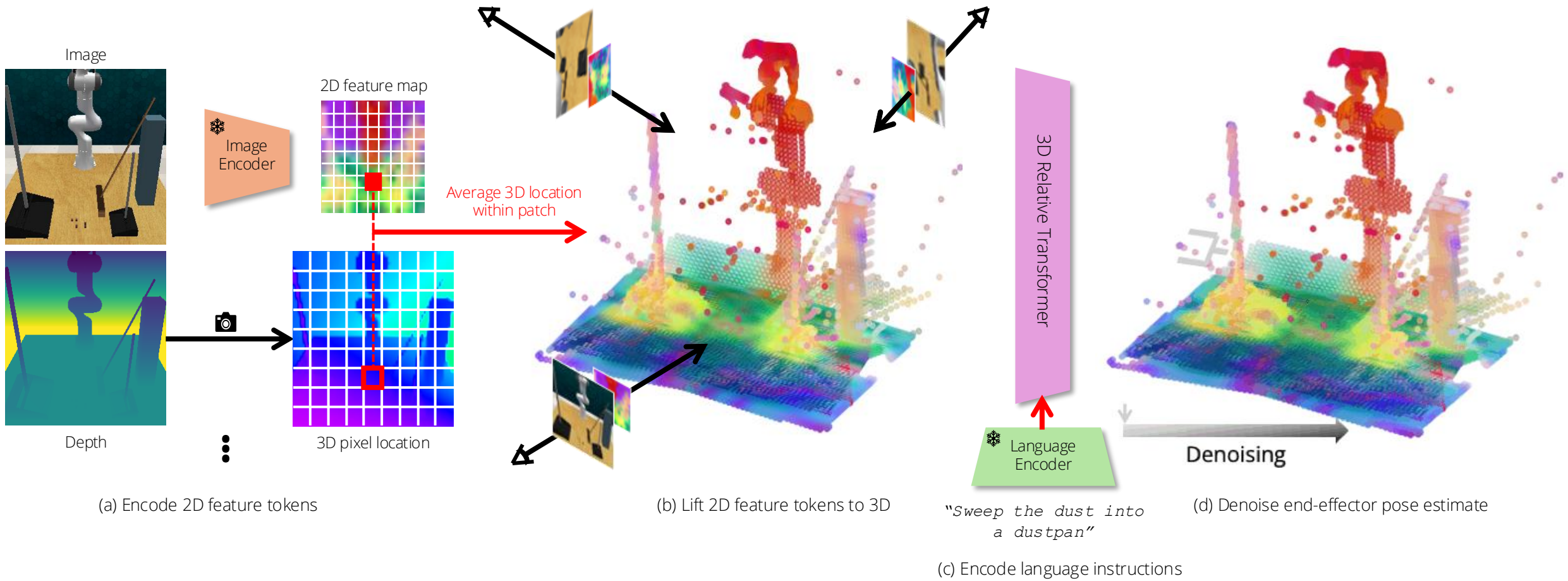
The reversed (noise→data) Variance Preserving (VP) diffusion path has the form

$$p_t(x|x_1) = \mathcal{N}(x | \alpha_{1-t}x_1, (1 - \alpha_{1-t}^2) I), \text{ where } \alpha_t = e^{-\frac{1}{2}T(t)}, T(t) = \int_0^t \beta(s)ds, \quad (18)$$

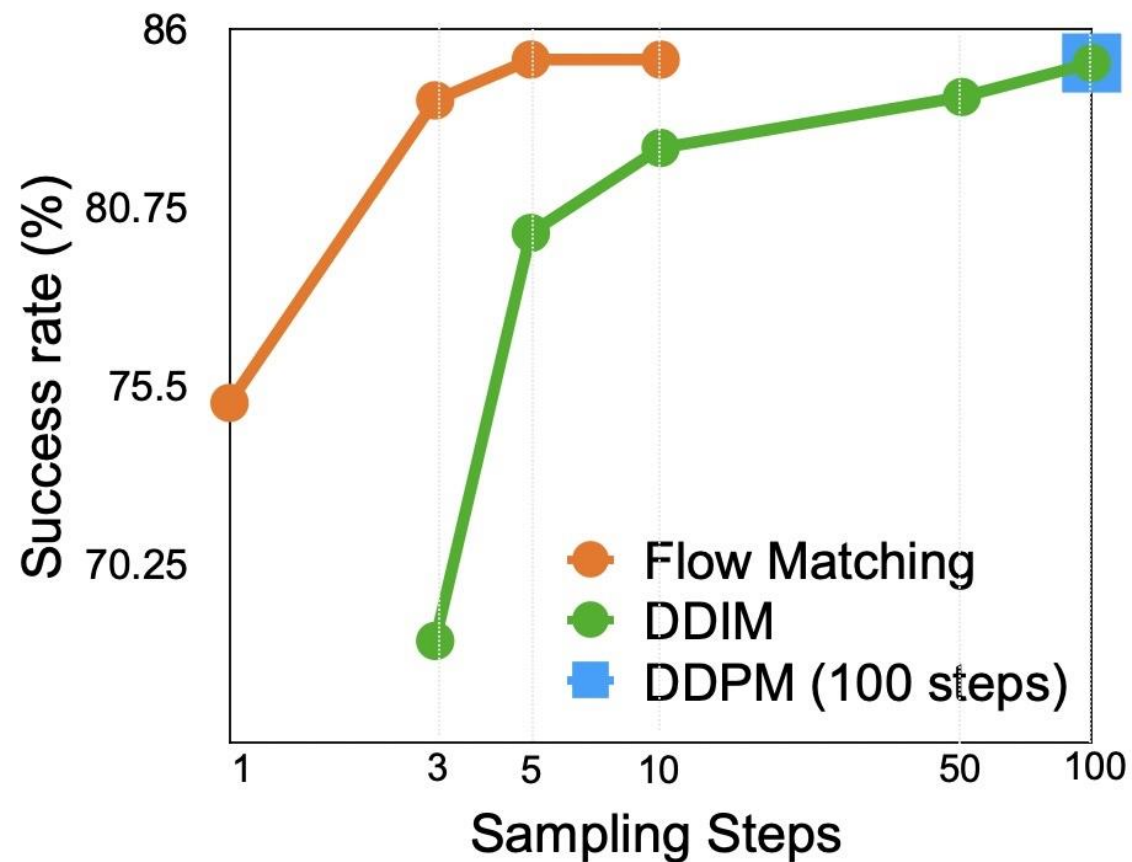
and β is the noise scale function. Equation 18 provides the choices of $\mu_t(x_1) = \alpha_{1-t}x_1$ and $\sigma_t(x_1) = \sqrt{1 - \alpha_{1-t}^2}$. Plugging these into equation 15 of Theorem 3 we get

$$u_t(x|x_1) = \frac{\alpha'_{1-t}}{1 - \alpha_{1-t}^2} (\alpha_{1-t}x - x_1) = -\frac{T'(1-t)}{2} \left[\frac{e^{-T(1-t)}x - e^{-\frac{1}{2}T(1-t)}x_1}{1 - e^{-T(1-t)}} \right]. \quad (19)$$

Why Flow Matching?

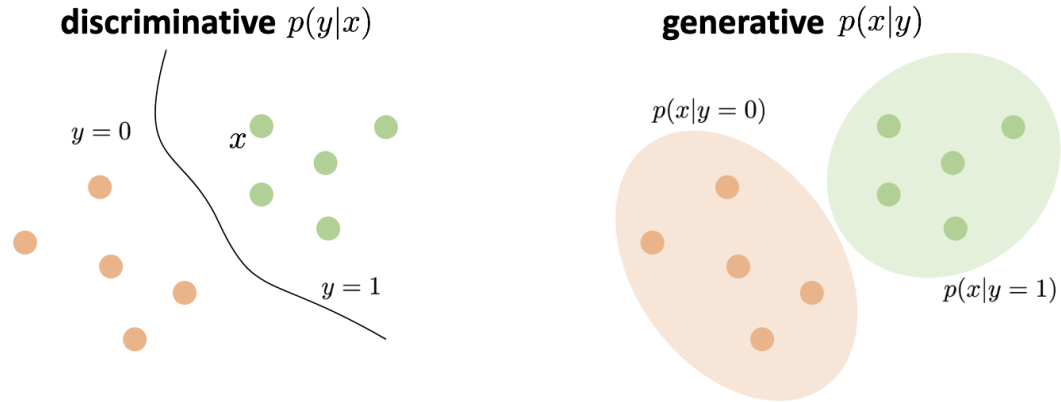


Flow Matching Requires Much Fewer Sampling Steps than DDPM

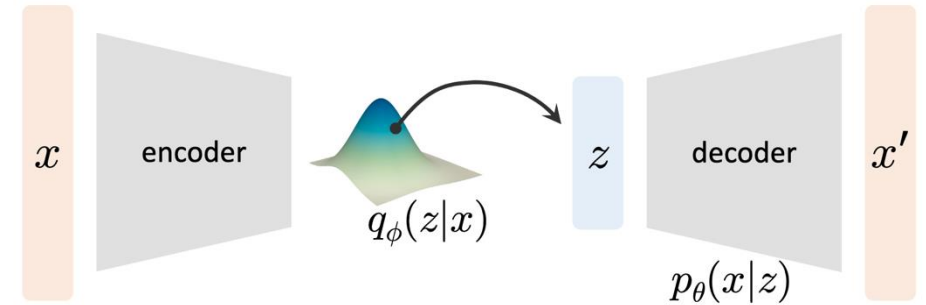


Summary

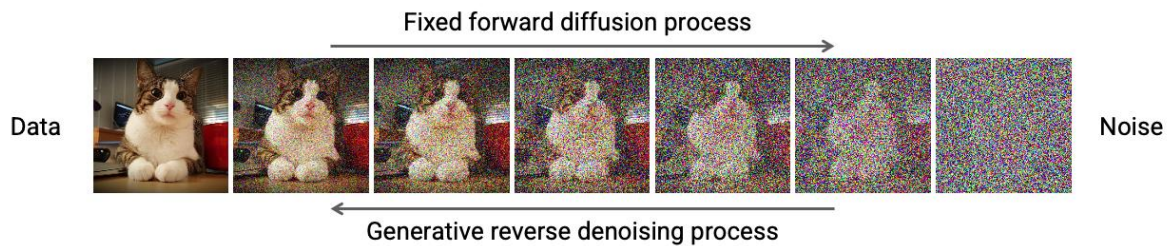
Discriminative Models vs. Generative Models



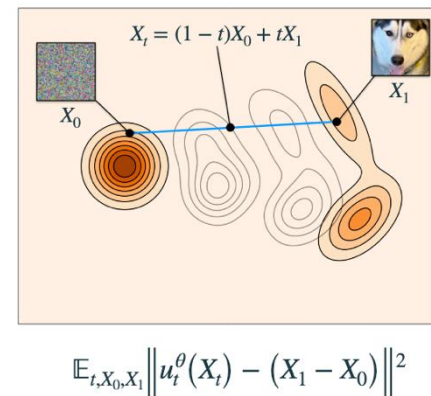
Variational AutoEncoders



Diffusion Models



Flow Matching



Autoregressive Models

