



Cornell Bowers C-IS

College of Computing and Information Science

Deep Learning

Week 7: Diffusion Models

Progress In Generative Modeling

VAEs, 2013



GANs, 2014



PixelCNN, 2016



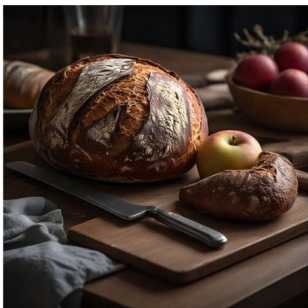
BigGAN, 2019



Imagen, 2022



Text-to-Image Diffusion Models



A bread, an apple, and a knife on a table



a robot cooking dinner in the kitchen



A teddy bear and a stuffed raccoon sitting on a wooden chair side by side



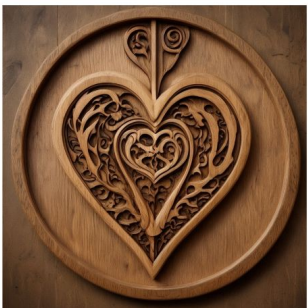
The oil painting shows a cow standing near a tree with red leaves



A traditional tea house in a tranquil garden with blooming cherry blossom trees



a painting of trees near a peaceful lake



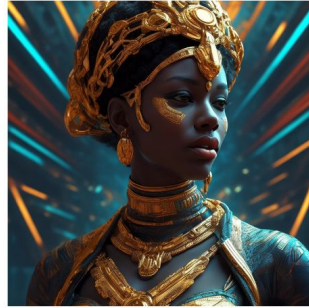
A heart made of wood



an old man with green eyes and a long grey beard



A painting of an adorable rabbit sitting on a colorful splash



an afrofuturist lady wearing gold jewelry



a black basketball shoe with a lightning bolt on it



A cool orange cat wearing sunglasses playing a guitar with a group of dancing bananas

Video Generation



Sora

Video Generation



Sora

Prompt: 3D animation of a small, round, fluffy creature with big, expressive eyes explores a vibrant, enchanted forest.

Video Generation



Prompt: 3D animation of a small, round, fluffy creature with big, expressive eyes explores a vibrant, enchanted forest.

Video Generation



Prompt: 3D animation of a small, round, fluffy creature with big, expressive eyes explores a vibrant, enchanted forest.

Prompt: A cat waking up its sleeping owner demanding breakfast.

Video Generation



Prompt: 3D animation of a small, round, fluffy creature with big, expressive eyes explores a vibrant, enchanted forest.



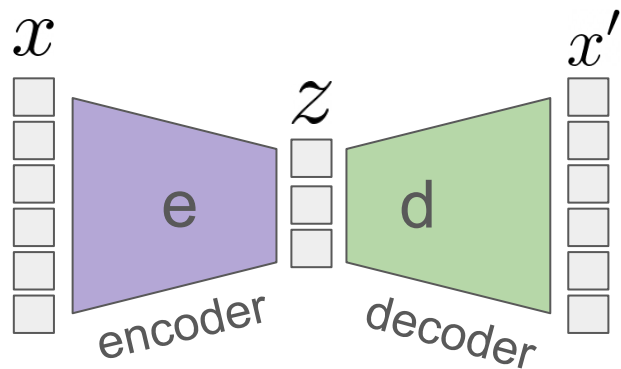
Prompt: A cat waking up its sleeping owner demanding breakfast.

Autoencoders

- Reconstruction loss: mean squared error

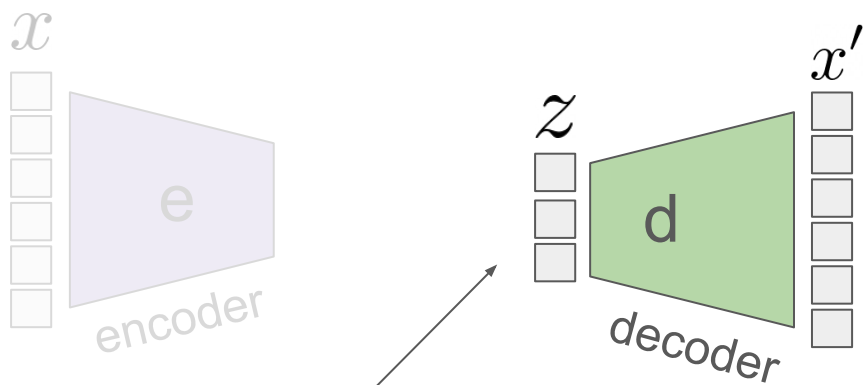
$$\sum_{x \in \mathcal{D}} (x - x')^2$$

where $x' = e(d(x))$



The Result: an Autoencoder.
[Kramer, 1991]

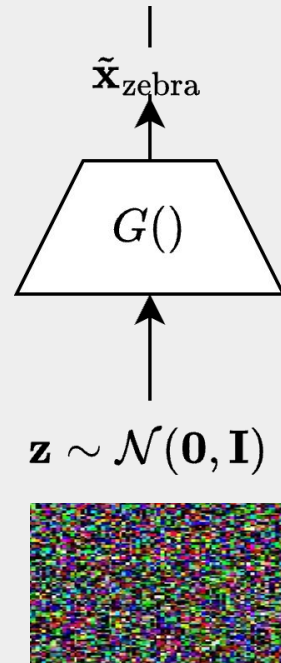
Sampling from an Autoencoder



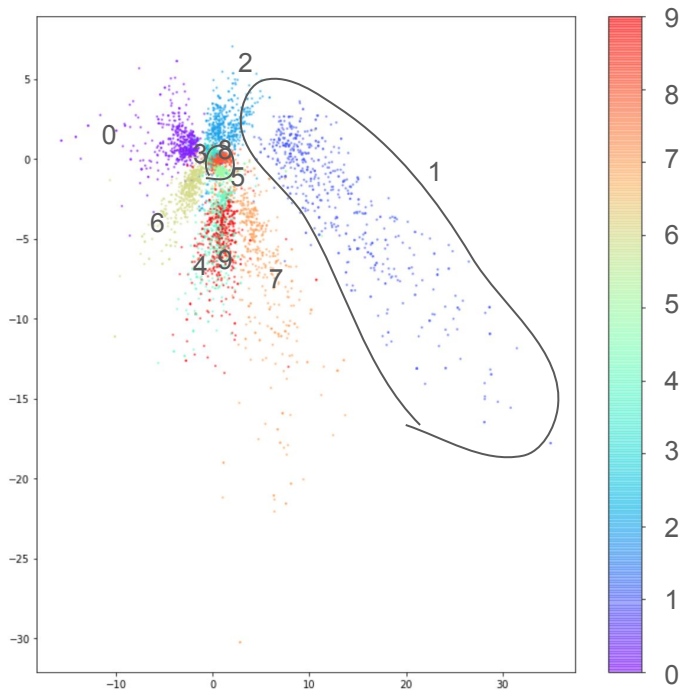
$$z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

feed decoder
(Gaussian) noise?

Recall how we could
sample with GANs...



Autoencoder trained on MNIST: latent space



Not a very nice representation...

- lots of empty space
- no symmetries between digit representations

Question:

What are the implications for sampling?

Figure 3-8. Plot of the latent space, colored by digit



reconstructed sample

$$x' = d(e(x))$$

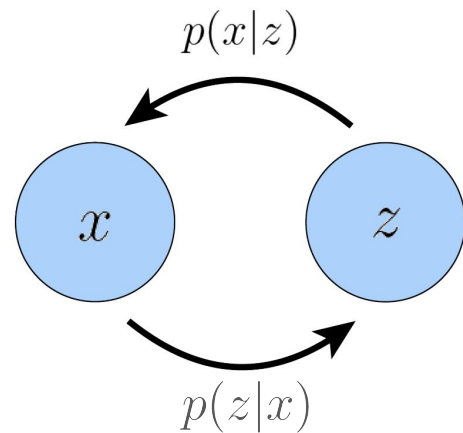


new image?

$$x' = d(\text{noise})$$

Variational Inference

- Have joint model $p(\mathbf{x}, \mathbf{z})$ $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$
- observe x (but not z);
- want to calculate posterior $p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})}$
- which requires $p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}$
 - i.e., the “evidence”.

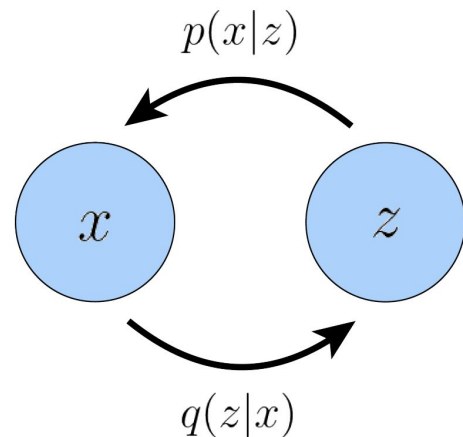


but the integral is often intractable! So, instead ...

Variational Inference

- Introduce a learnable variational approximation of the posterior

$$q_{\phi}(z|\mathbf{x}) \approx p(z|\mathbf{x})$$

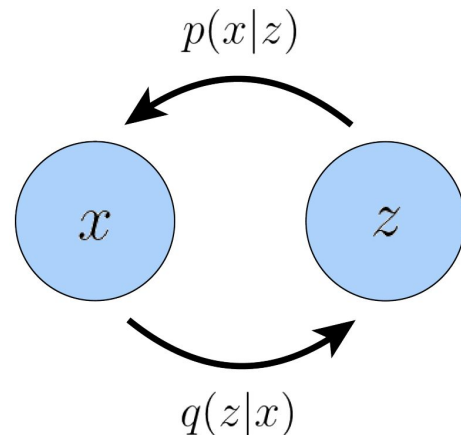


Variational Inference

- Introduce a learnable variational approximation of the posterior

$$q_{\phi}(z|\mathbf{x}) \approx p(z|\mathbf{x})$$

- Bound the likelihood using the variational posterior



$$\log p(\mathbf{x}) = \mathbb{E}_{q_{\phi}(z|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, z)}{q_{\phi}(z|\mathbf{x})} \right] + D_{\text{KL}}(q_{\phi}(z|\mathbf{x}) \parallel p(z|\mathbf{x}))$$

Intractable; Evidence/
Log-likelihood

Tractable; Evidence
Lower Bound (ELBO)

Intractable; Divergence
from true posterior

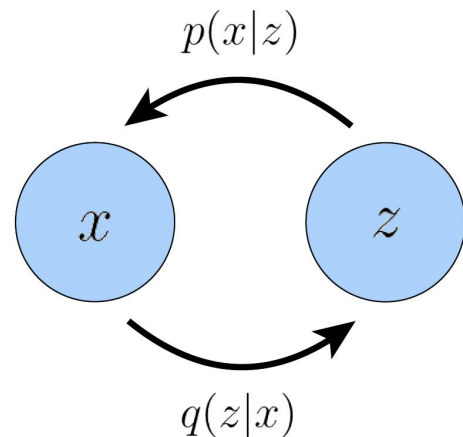
Variational Inference

- Introduce a learnable variational approximation of the posterior

$$q_{\phi}(\mathbf{z}|\mathbf{x}) \approx p(\mathbf{z}|\mathbf{x})$$

- Bound the likelihood with the ELBO

$$\log p(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}))$$

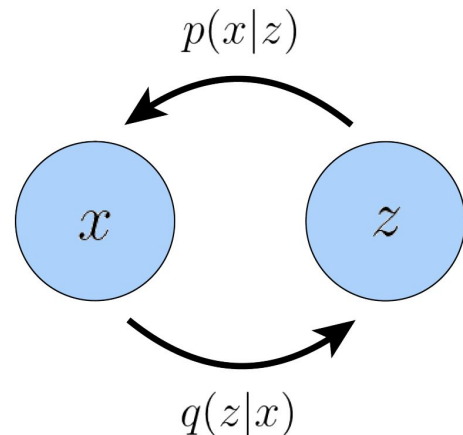


Variational Inference

- Introduce a learnable variational approximation of the posterior

$$q_{\phi}(z|\mathbf{x}) \approx p(z|\mathbf{x})$$

- Bound the likelihood with the ELBO



$$\log p(\mathbf{x}) = \mathbb{E}_{q_{\phi}(z|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, z)}{q_{\phi}(z|\mathbf{x})} \right] + D_{\text{KL}}(q_{\phi}(z|\mathbf{x}) \parallel p(z|\mathbf{x}))$$

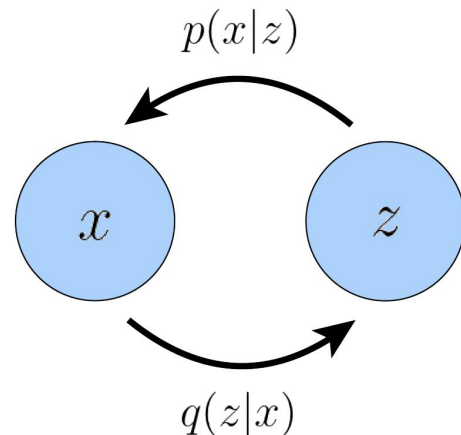
$$\geq \mathbb{E}_{q_{\phi}(z|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, z)}{q_{\phi}(z|\mathbf{x})} \right]$$

Non-negativity
of KL

Tractable; ELBO

The Evidence Lower Bound (ELBO)

- Maximize the ELBO
- Either:
 - Maximizes the likelihood of the observed data
 - Improves the approximation of the unknown posterior



$$\mathbb{E}_{q_{\phi}(z|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, z)}{q_{\phi}(z|\mathbf{x})} \right] = \log p(\mathbf{x}) - D_{\text{KL}}(q_{\phi}(z|\mathbf{x}) \parallel p(z|\mathbf{x}))$$

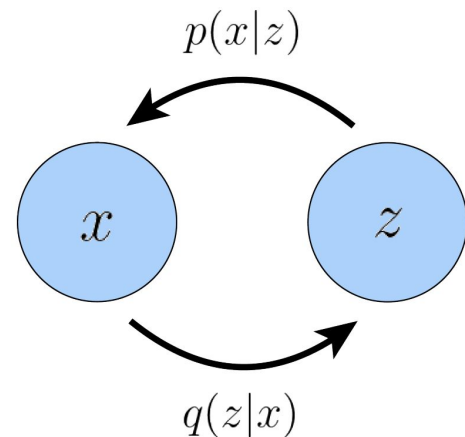
Tractable; ELBO

Intractable;
Evidence

Intractable; Divergence between
approximate and true posterior

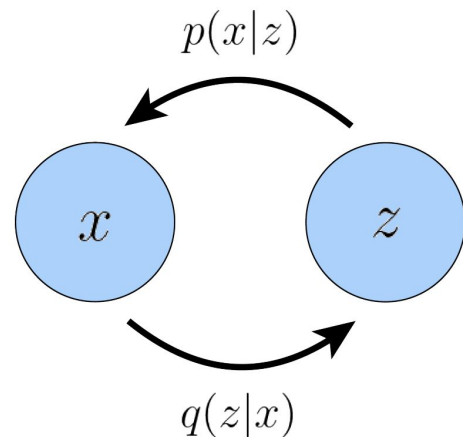
The Evidence Lower Bound (ELBO)

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$$



(Chain Rule of Probability)

The Evidence Lower Bound (ELBO)

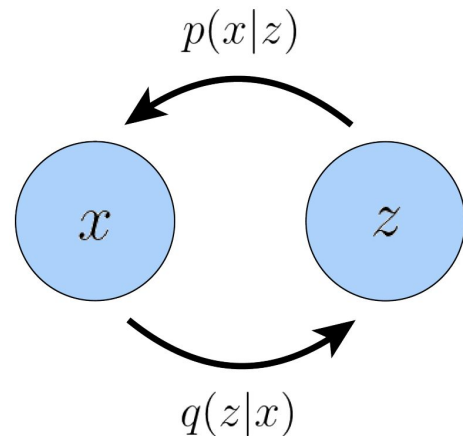


(Chain Rule of Probability)

(Split the Expectation)

$$\begin{aligned}\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]\end{aligned}$$

The Evidence Lower Bound (ELBO)



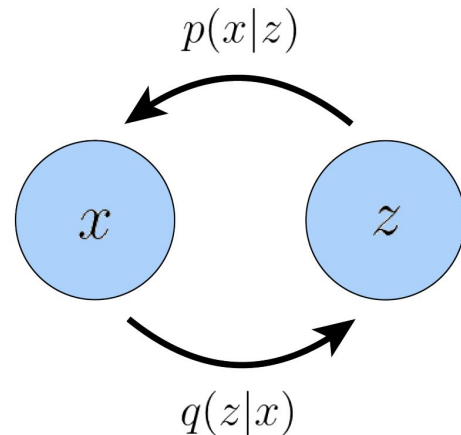
(Chain Rule of Probability)

(Split the Expectation)

(Definition of KL Divergence)

$$\begin{aligned} \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})) \end{aligned}$$

The Evidence Lower Bound (ELBO)



(Chain Rule of Probability)

(Split the Expectation)

(Definition of KL Divergence)

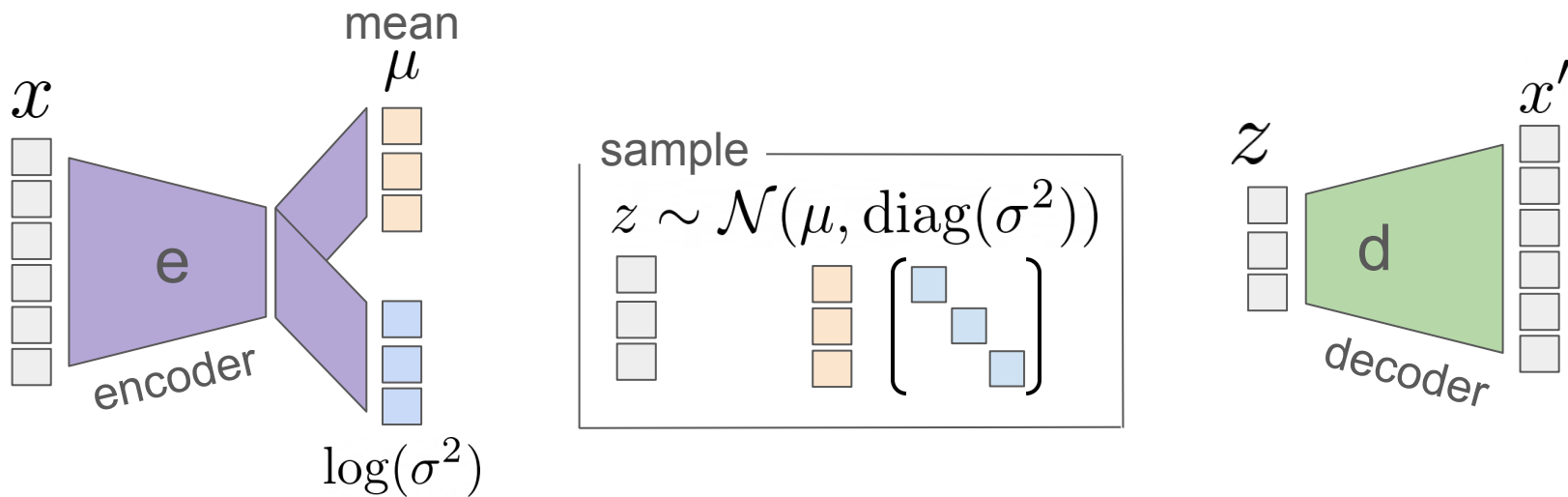
$$p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\begin{aligned} \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \\ &= \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))}_{\text{prior matching term}} \end{aligned}$$

Mean-Squared
Error/Cross-Entropy

Regularization

An Architecture for Gaussians



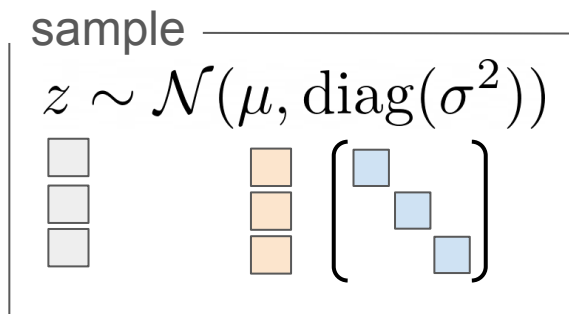
variance

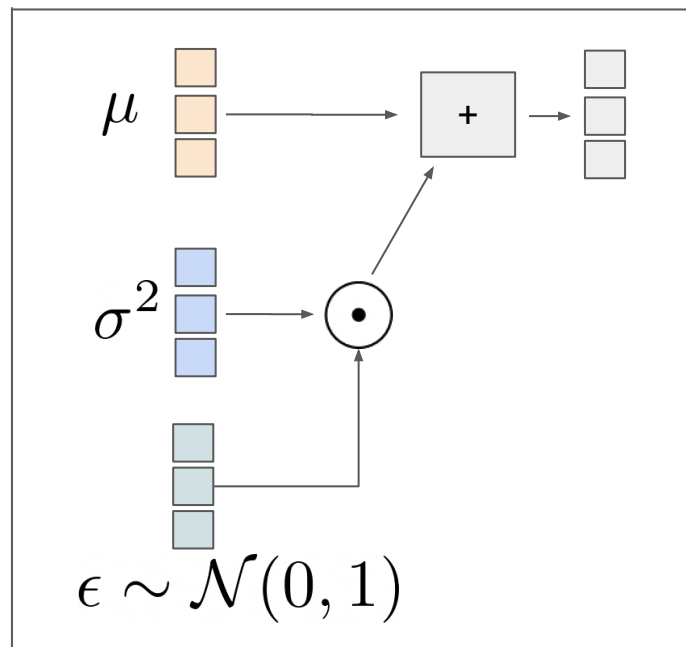
Problem: backpropagation through sampling process?

$$q_{\phi}(z|x)$$

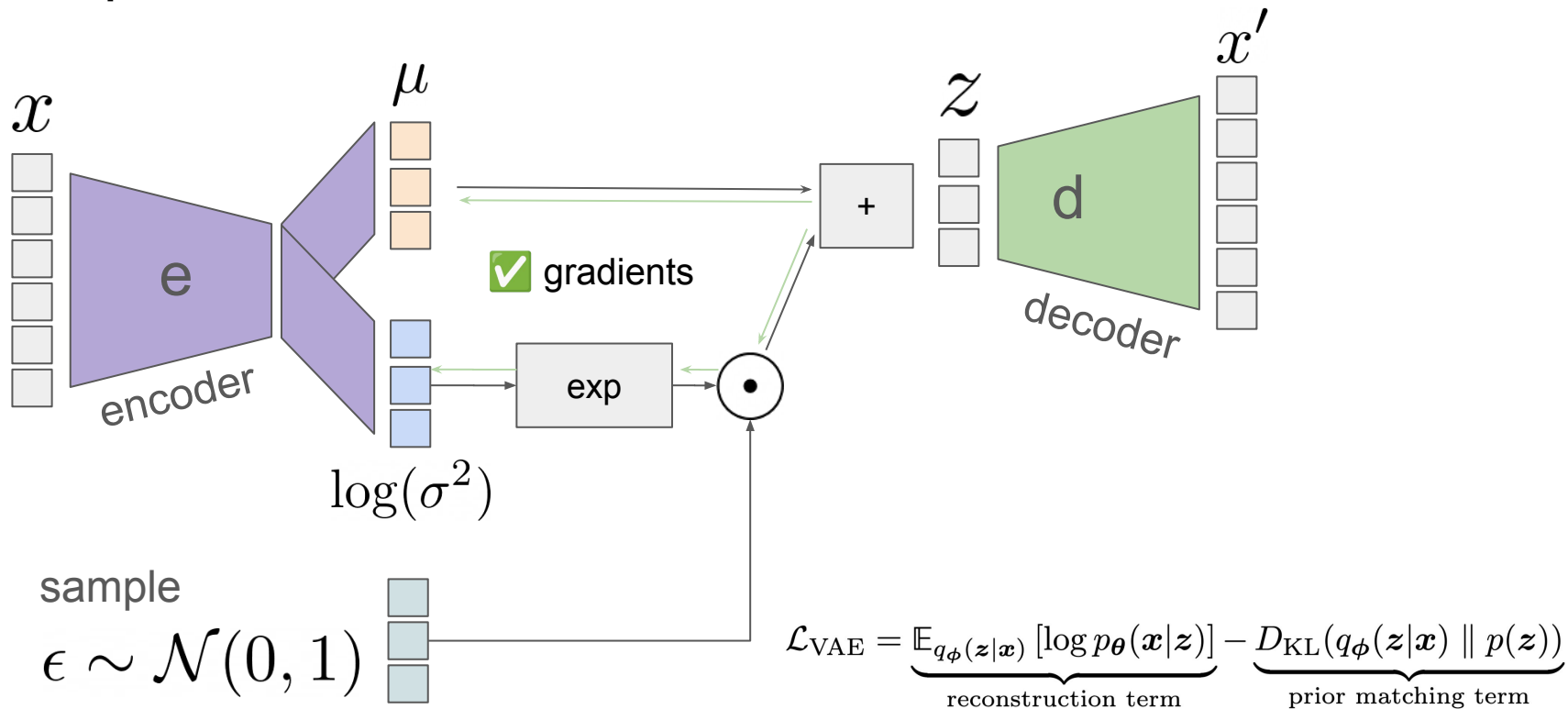
The Reparameterization Trick

$$\mathcal{N}(\mu, \text{diag}(\sigma^2)) = \mu + \sigma^2 \odot \mathcal{N}(0, I)$$

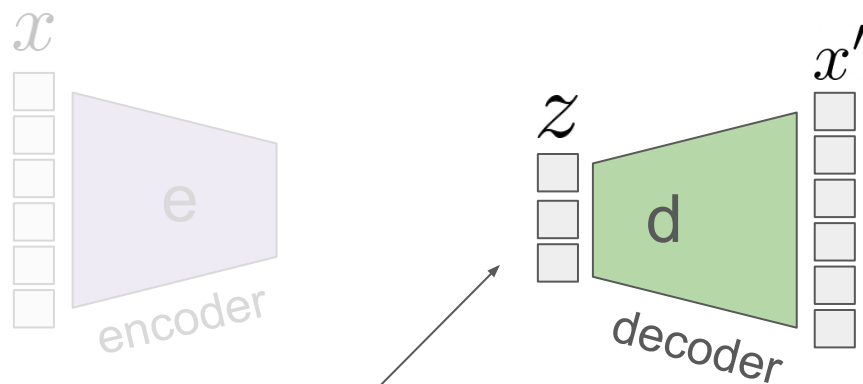


$$=$$


The Reparameterization Trick



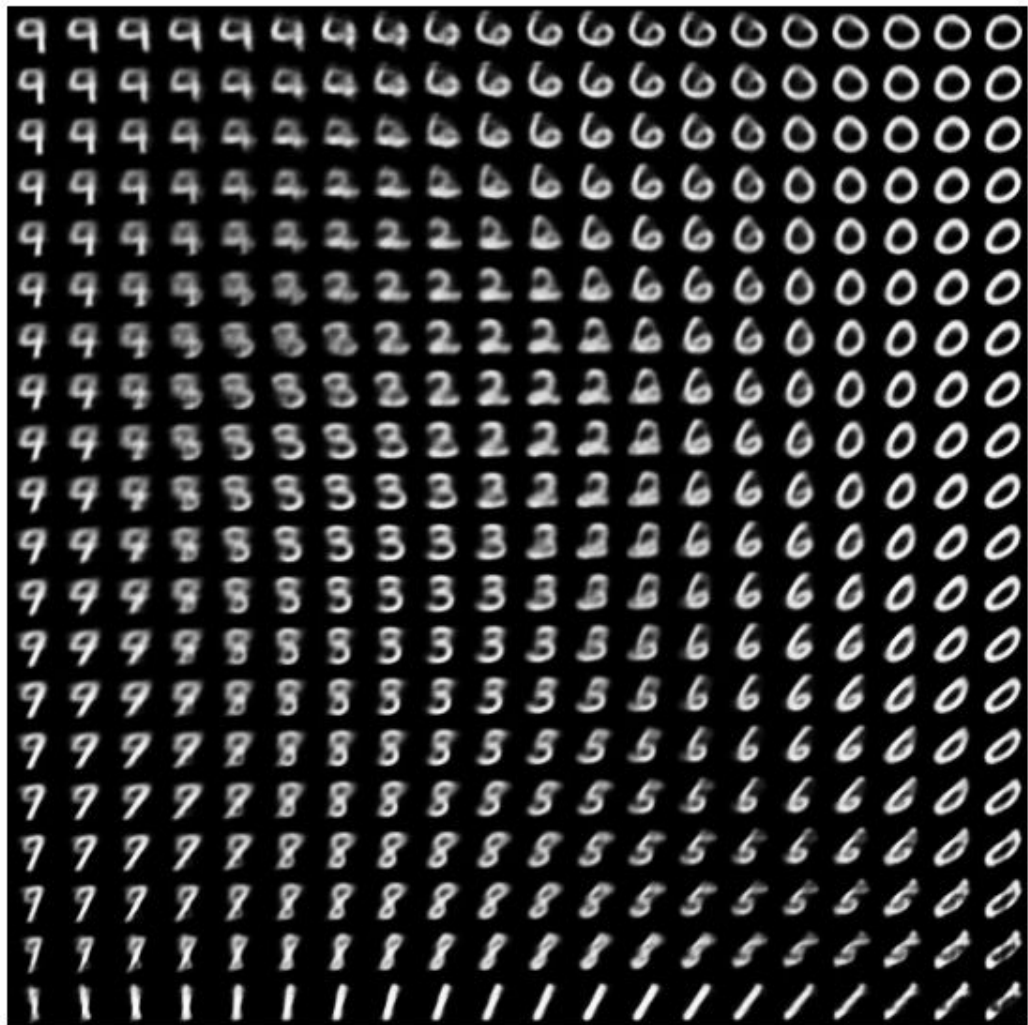
Sampling from a VAE



$$z \sim p(z) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

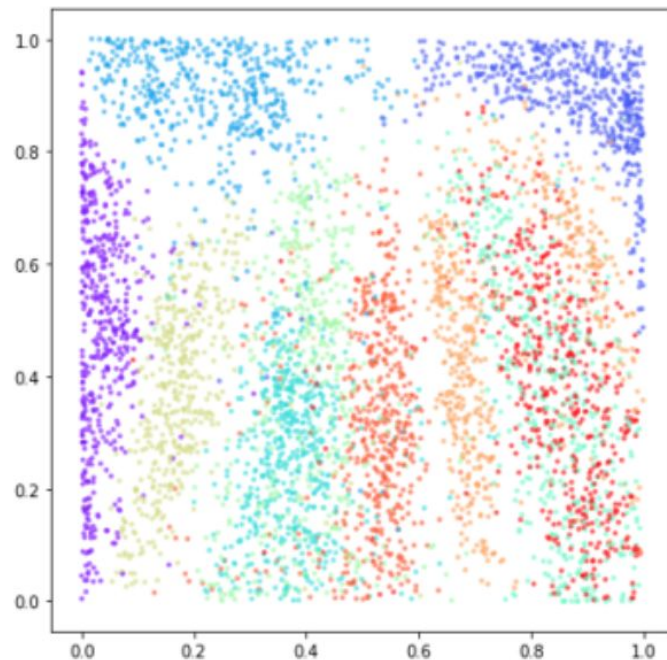
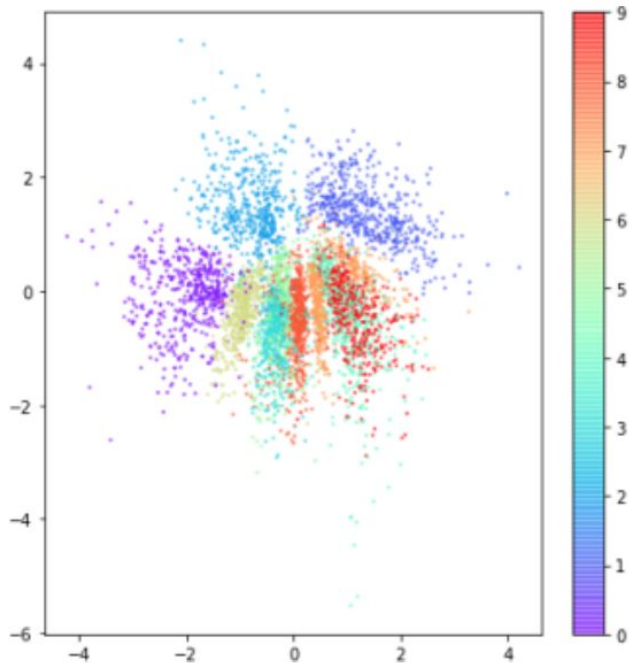
a much nicer space...

can smoothly interpolate digits in
a meaningful, digit-y kind of way



Back to MNIST: Visualizing latent space again

VAE Latent space, note the distribution is centered, and each digit has an equal portion



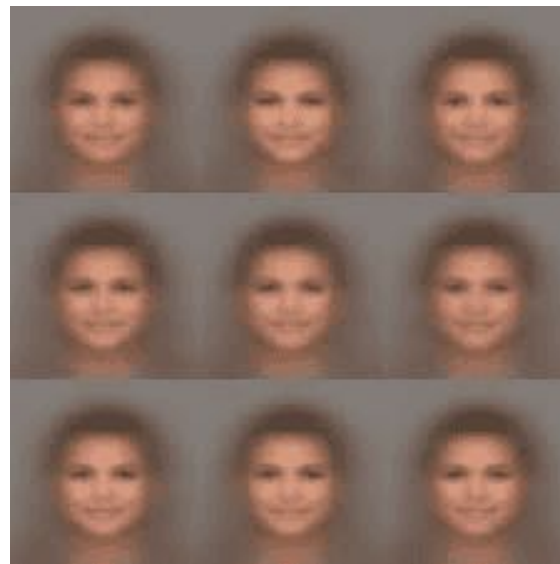
The Biggest Drawback of VAEs

- Out of the box, generated images can be blurry.

Question: Why?



[VAE v. GAN](#)

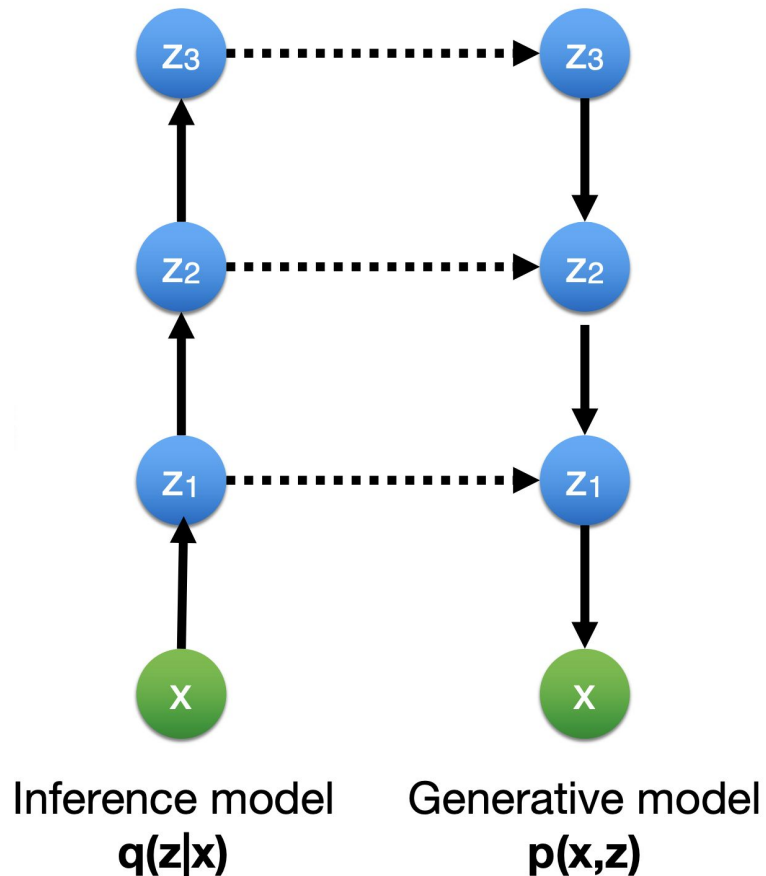


Hierarchical VAEs

- “Flat” VAEs suffer from simple priors
- Define a hierarchical generative process

$$q_{\phi}(\mathbf{z}_{1,2,3}|\mathbf{x}) = q_{\phi}(\mathbf{z}_1|\mathbf{x})q_{\phi}(\mathbf{z}_2|\mathbf{z}_1)q_{\phi}(\mathbf{z}_3|\mathbf{z}_2)$$

$$p_{\theta}(\mathbf{z}_{1,2,3}) = p_{\theta}(\mathbf{z}_3)p_{\theta}(\mathbf{z}_2|\mathbf{z}_3)p_{\theta}(\mathbf{z}_1|\mathbf{z}_2)p_{\theta}(\mathbf{x}|\mathbf{z}_1)$$



Extending the ELBO

- ELBO derivation is unchanged

$$\begin{aligned}\log p(\mathbf{x}) &= \log \int p(\mathbf{x}, \mathbf{z}_{1:T}) d\mathbf{z}_{1:T} \\ &\geq \mathbb{E}_{q_\phi(\mathbf{z}_{1:T}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_\phi(\mathbf{z}_{1:T}|\mathbf{x})} \right]\end{aligned}$$

Extending the ELBO

- Omitting some steps
 - See “Understanding Diffusion Models: A Unified Perspective” by Calvin Luo for a nice walkthrough

$$\begin{aligned}\log p(\mathbf{x}) &= \log \int p(\mathbf{x}, \mathbf{z}_{1:T}) d\mathbf{z}_{1:T} \\ &\geq \mathbb{E}_{q_\phi(\mathbf{z}_{1:T}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_\phi(\mathbf{z}_{1:T}|\mathbf{x})} \right] \\ &= \dots\end{aligned}$$

Extending the ELBO

$$\log p(\mathbf{x}) = \log \int p(\mathbf{x}, \mathbf{z}_{1:T}) d\mathbf{z}_{1:T}$$

$$\geq \mathbb{E}_{q_\phi(\mathbf{z}_{1:T}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_\phi(\mathbf{z}_{1:T}|\mathbf{x})} \right]$$

= ...

$$= \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}_1|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(\mathbf{z}_T|\mathbf{x}) \parallel p(\mathbf{z}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \mathbb{E}_{q_\phi(\mathbf{z}_t|\mathbf{x})} [D_{\text{KL}}(q_\phi(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) \parallel p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t))]$$

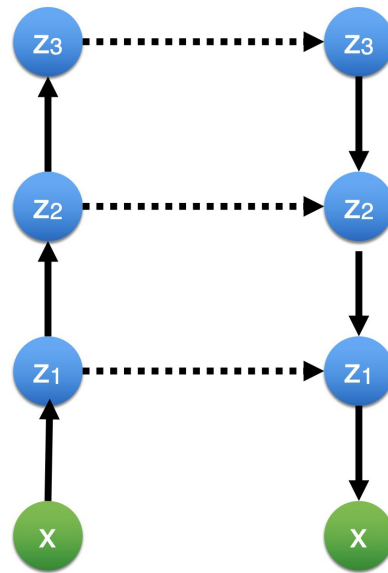
KL-Div between
Gaussians

$$\frac{1}{2} \left\{ \left(\frac{\sigma_0}{\sigma_1} \right)^2 + \frac{(\mu_1 - \mu_0)^2}{\sigma_1^2} - 1 + \ln \frac{\sigma_1^2}{\sigma_0^2} \right\}$$

Consistency Term- Unstable Optimization!

Extending the ELBO

Hard to jointly learn hierarchical encoders and decoders!



Inference model
 $q(\mathbf{z}|\mathbf{x})$

Generative model
 $p(\mathbf{x},\mathbf{z})$

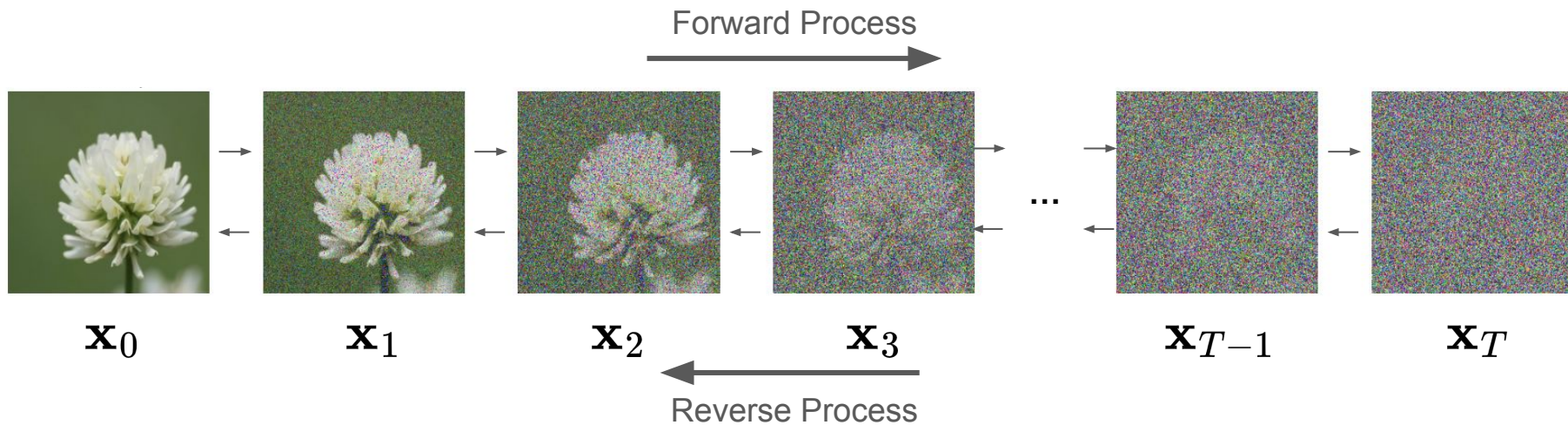
$$= \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}_1|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(\mathbf{z}_T|\mathbf{x}) \parallel p(\mathbf{z}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \mathbb{E}_{q_\phi(\mathbf{z}_t|\mathbf{x})} [D_{\text{KL}}(q_\phi(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) \parallel p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t))]$$

Consistency Term- Unstable Optimization!

Denosing Diffusion Models

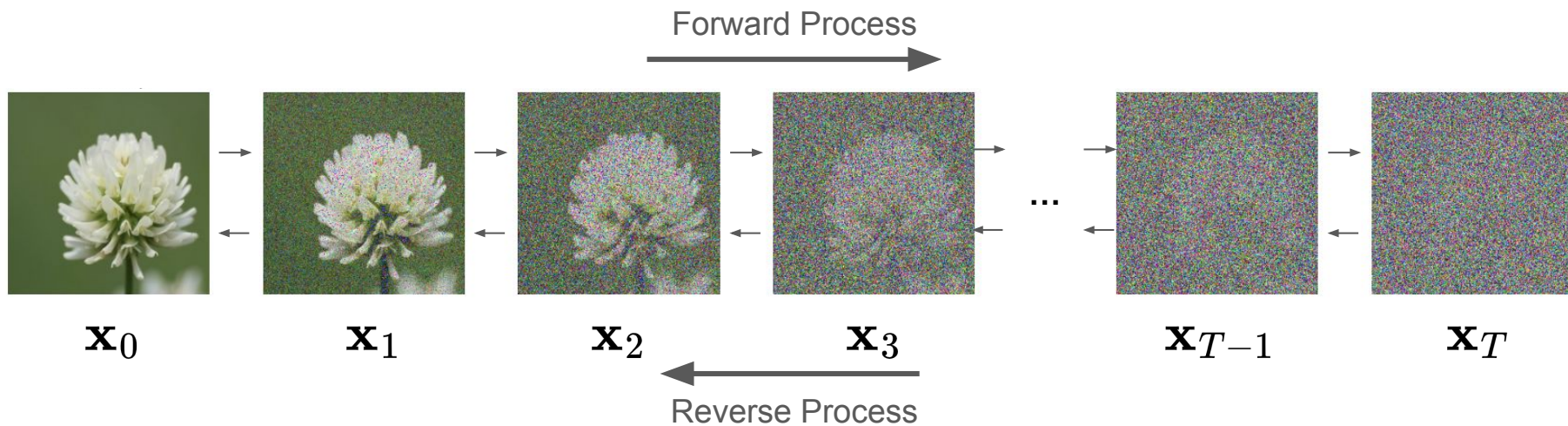
Denosing diffusion models consist of two processes:

- Forward diffusion process that gradually adds noise to input
- Reverse denosing process that learns to generate data by denoising

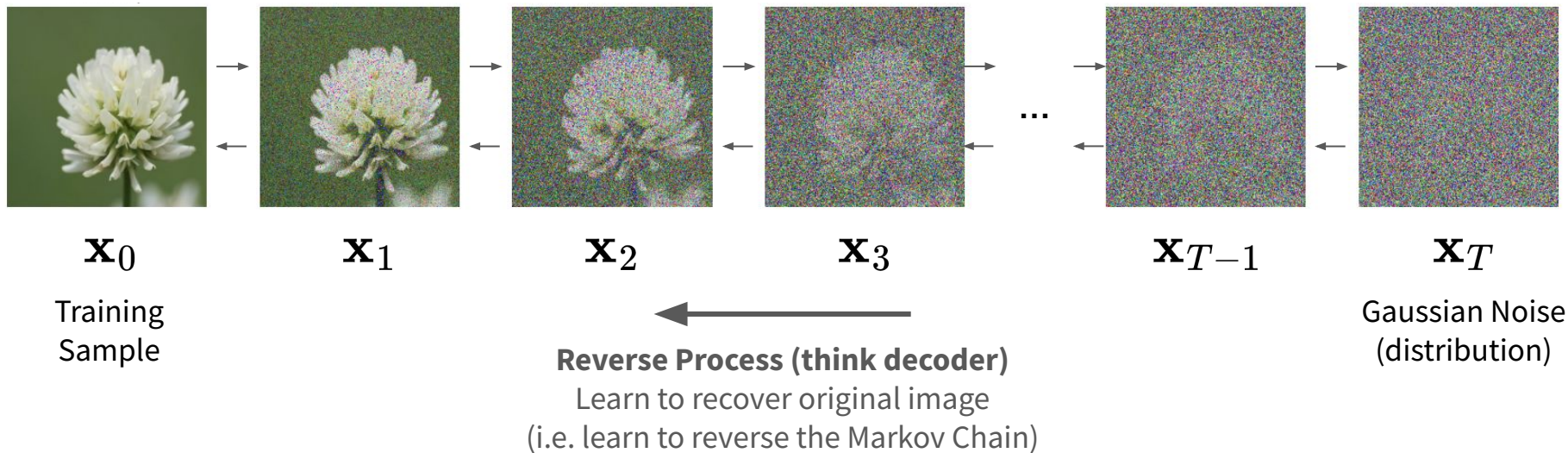


Discuss:

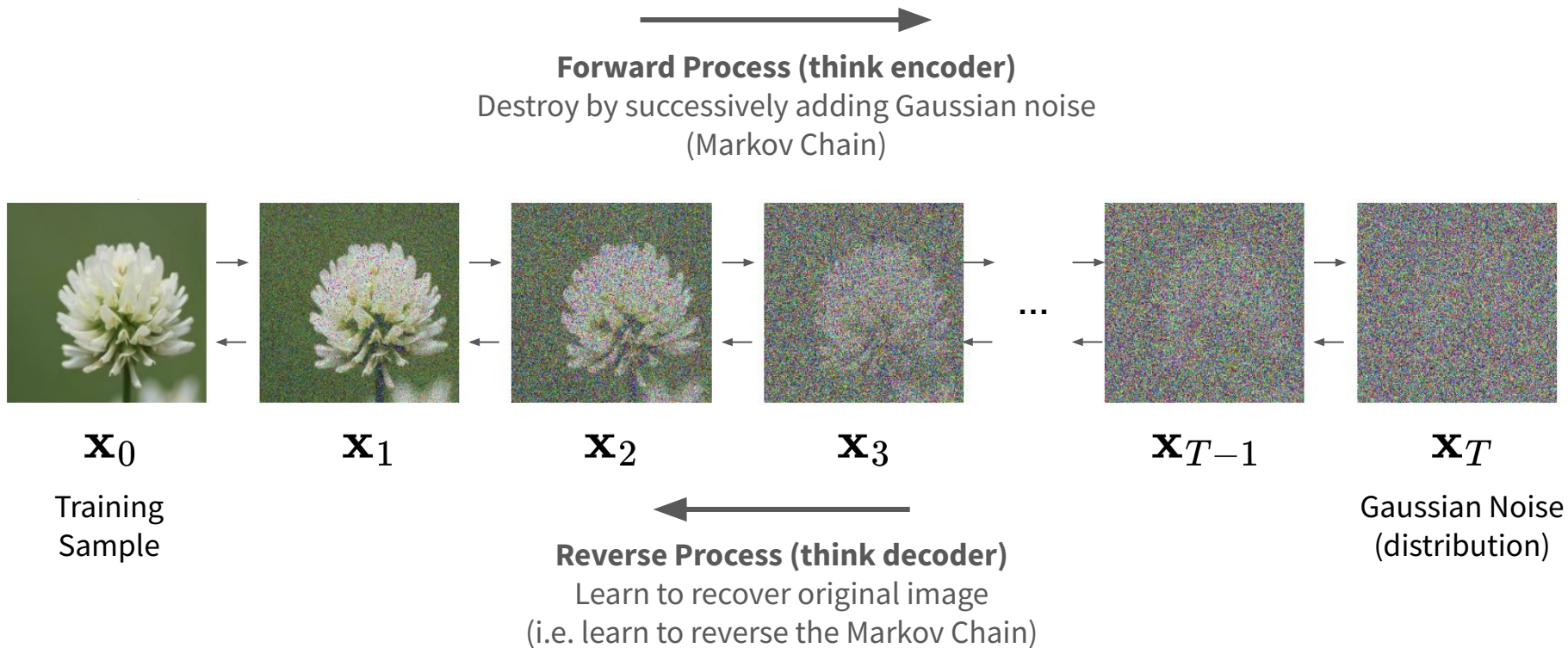
- How to define the forward and reverse directions?



Reverse Process: high level idea

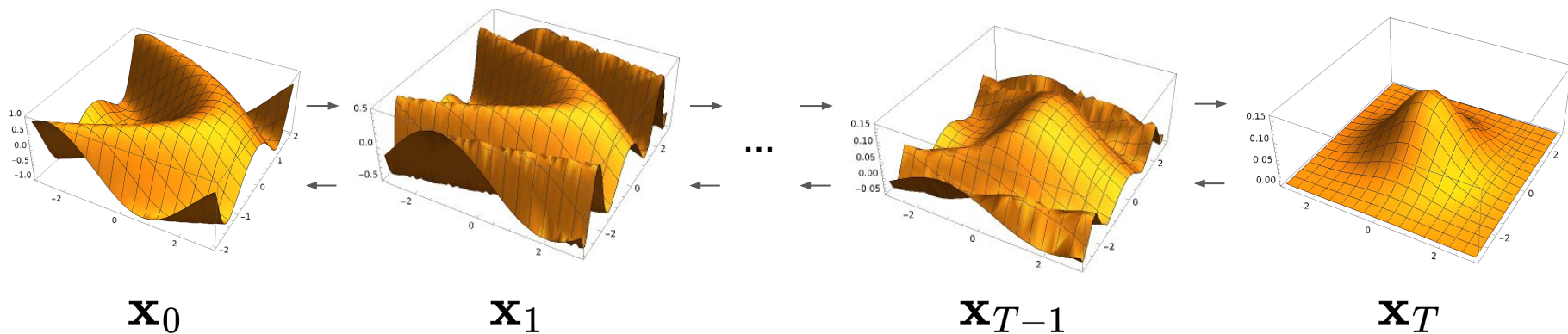


Putting it together



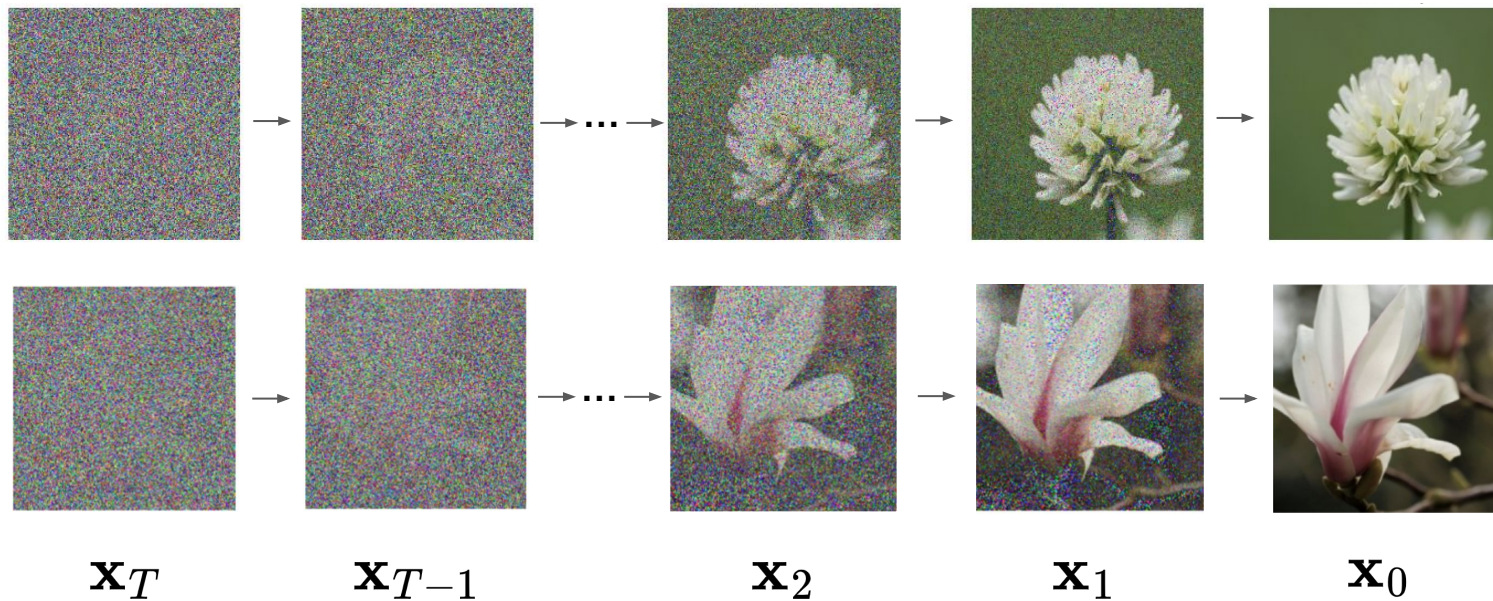
Diffusion Models

We define a mapping to Gaussian noise (forward process)
Want to **learn the reverse mapping to generate data** (reverse process)



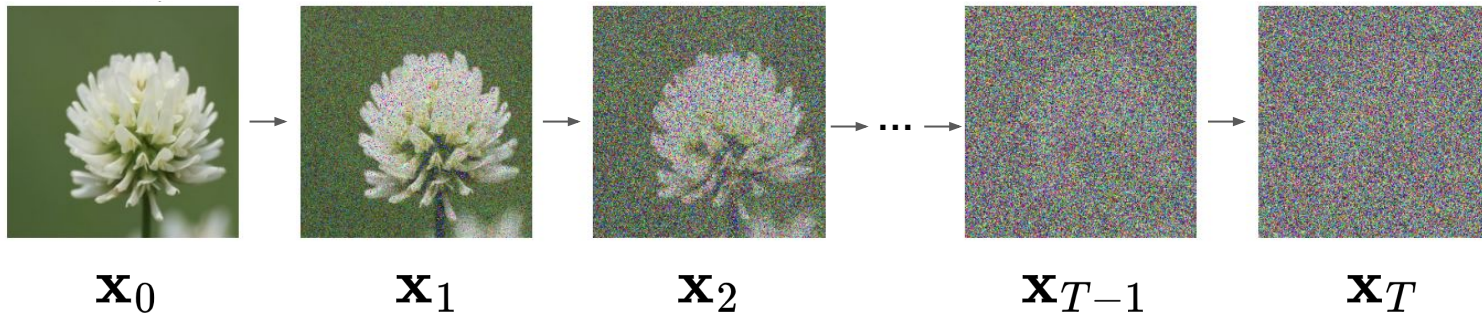
Diffusion Sampling

Different draws of initial noise lead to diverse outputs



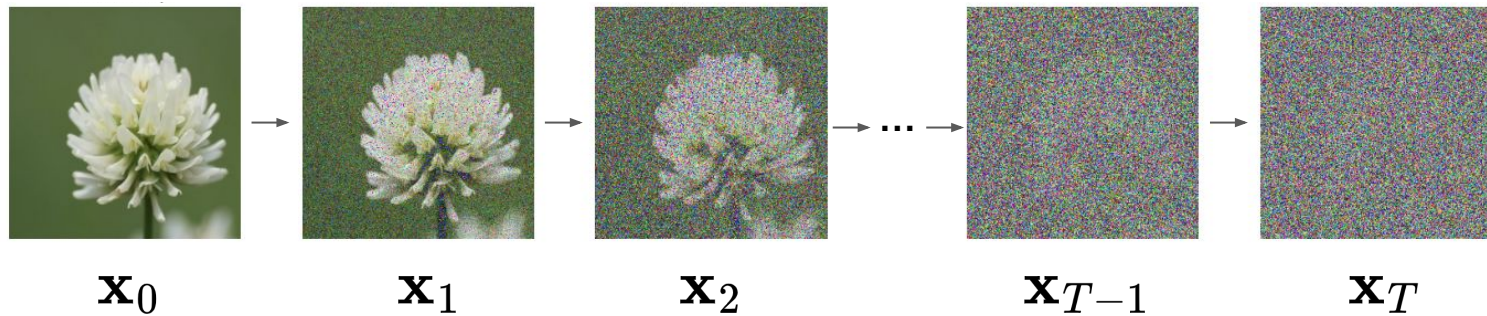
Forward Process Overview

- Destroys original image \mathbf{x}_0 by **successively adding Gaussian noise**
- Desired outcome: At step T , \mathbf{x}_T is a **pure Gaussian noise**
 - i.e. the distribution we map the data manifold to



Details: Forward Process

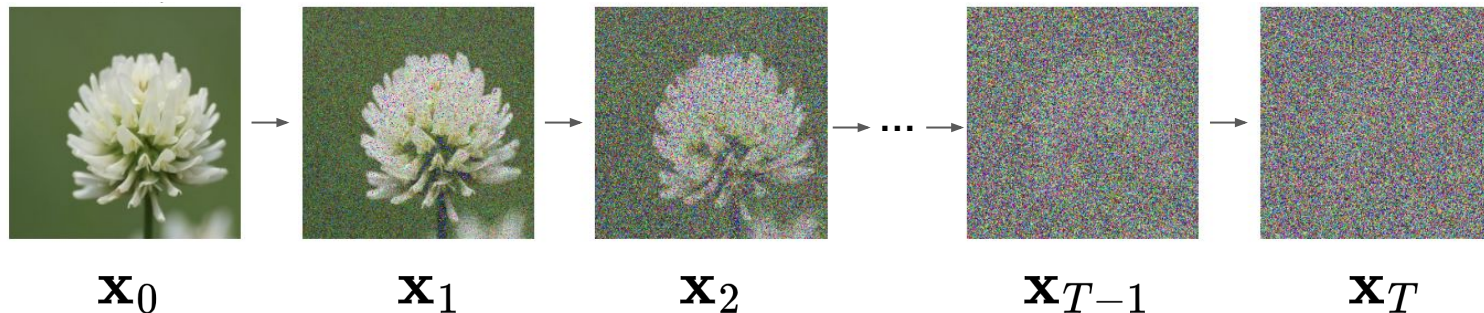
1. \mathbf{x}_0 sampled from some distribution



Details: Forward Process

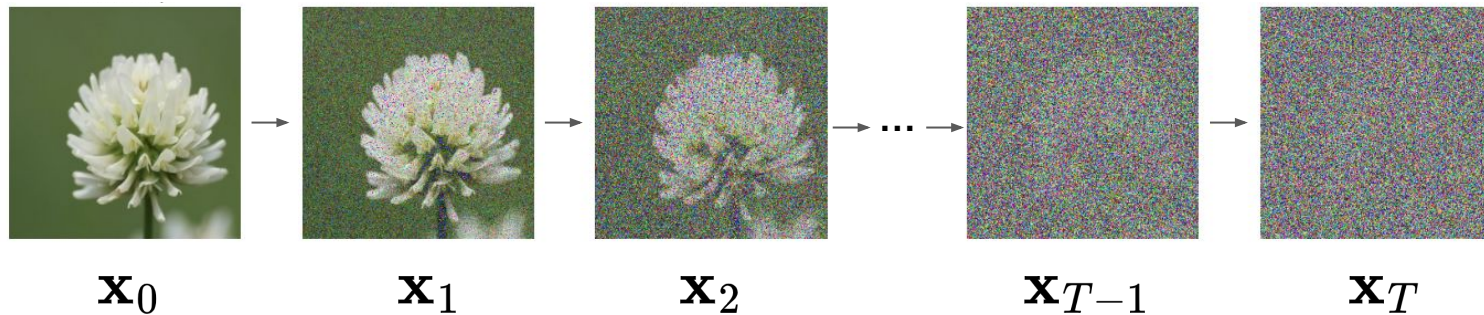
2. \mathbf{x}_t sampled from normal distribution conditioned on \mathbf{x}_{t-1} given by:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad \{\beta_t \in (0, 1)\}_{t=1}^T$$



Details: Forward Process

$\{\beta_t \in (0, 1)\}_{t=1}^T$ is variance schedule (controlling **how** we move toward Gaussian noise)



Details: Forward Process

\mathbf{x}_t sampled from normal distribution conditioned on \mathbf{x}_{t-1} given by:

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

Can we extend this to sampling \mathbf{x}_t in a closed form? We use the re-parametrization trick:

Let $\alpha_t := 1 - \beta_t$, and let $\boldsymbol{\epsilon}_{t-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}$$

Details: Forward Process

Inductively, we can say

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}$$

$$= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \bar{\boldsymbol{\epsilon}}_{t-2}$$

$$= \dots$$

$$= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$$

Merged noise.
epsilon is still $\sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

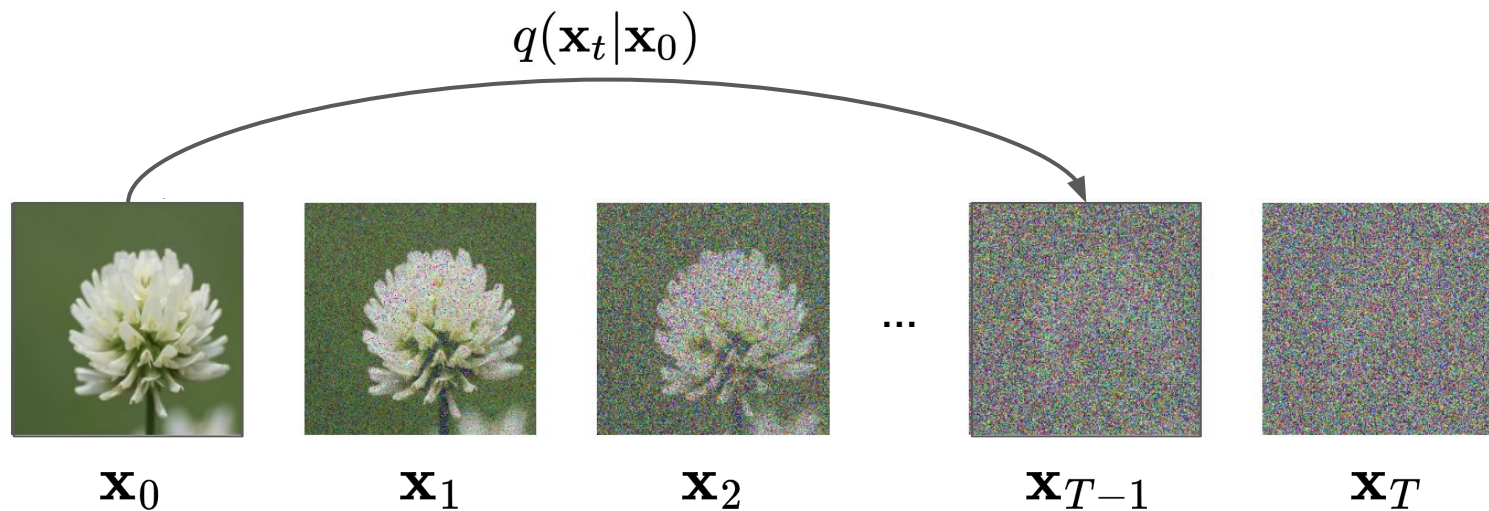
$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

Details: Forward Process

Can sample \mathbf{x}_t in closed-form as $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$

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



Aside: Noise Schedules

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

- Define the noise schedule in terms of $\bar{\alpha}_t \in (0, 1)$
 - Some monotonically decreasing function from 1 to 0
- Cosine Noise schedule:

$$\bar{\alpha}_t = \cos(.5\pi t/T)^2$$

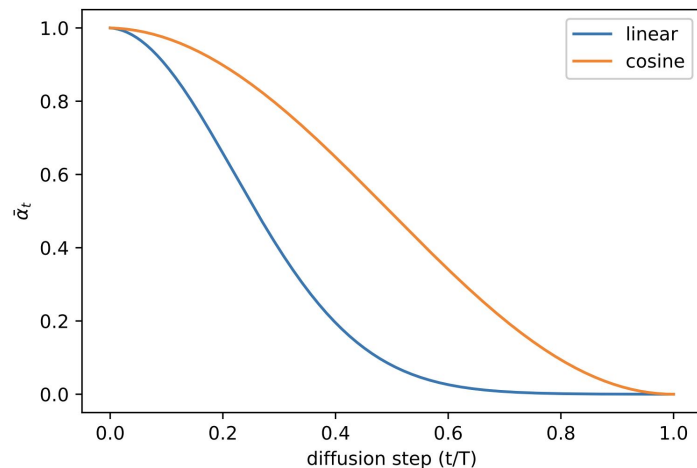
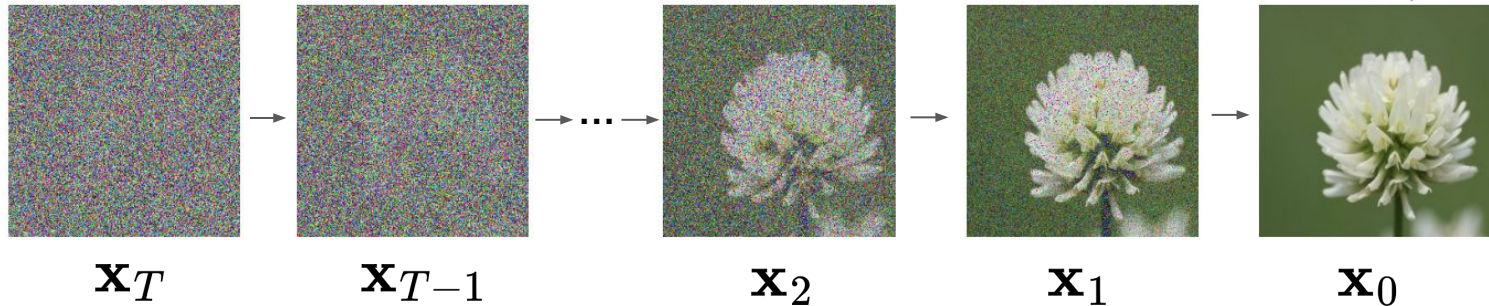


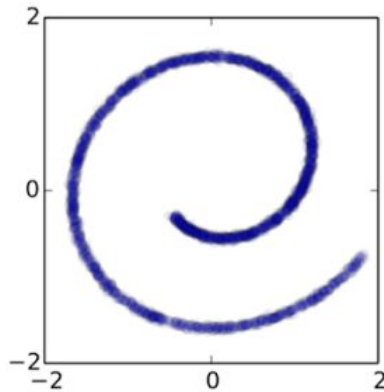
Figure 5. $\bar{\alpha}_t$ throughout diffusion in the linear schedule and our proposed cosine schedule.

Reverse Process Overview

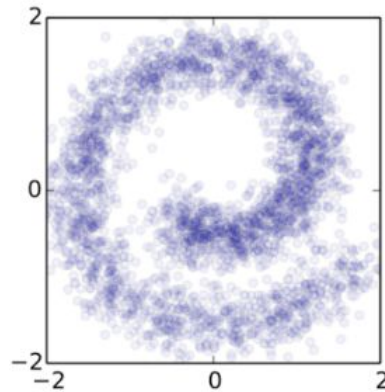
- "Learn to reverse what we just destroyed"
 - Learn time reversal of Markov Chain; we train a model for this
- Desired outcome: some \mathbf{x}_0 close to the original data distribution



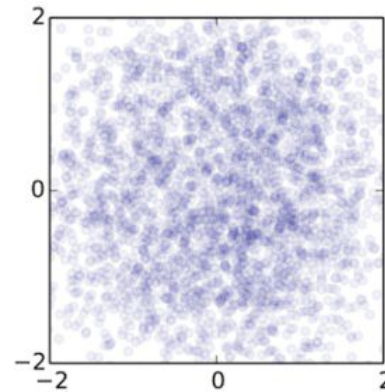
$t = 0$



$t = \frac{T}{2}$



$t = T$

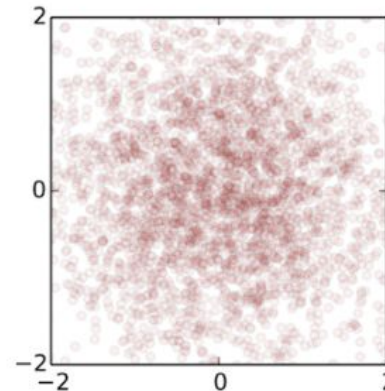
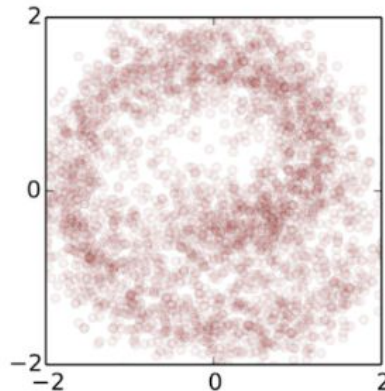
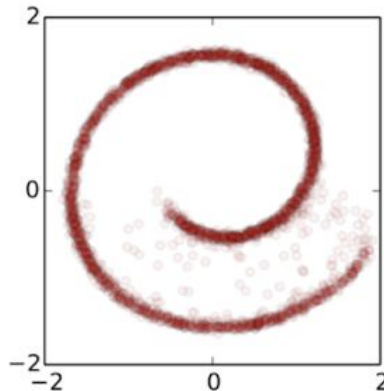


The forward trajectory

$$q(\mathbf{x}_{0:T})$$

The reverse trajectory

$$p_{\theta}(\mathbf{x}_{0:T})$$

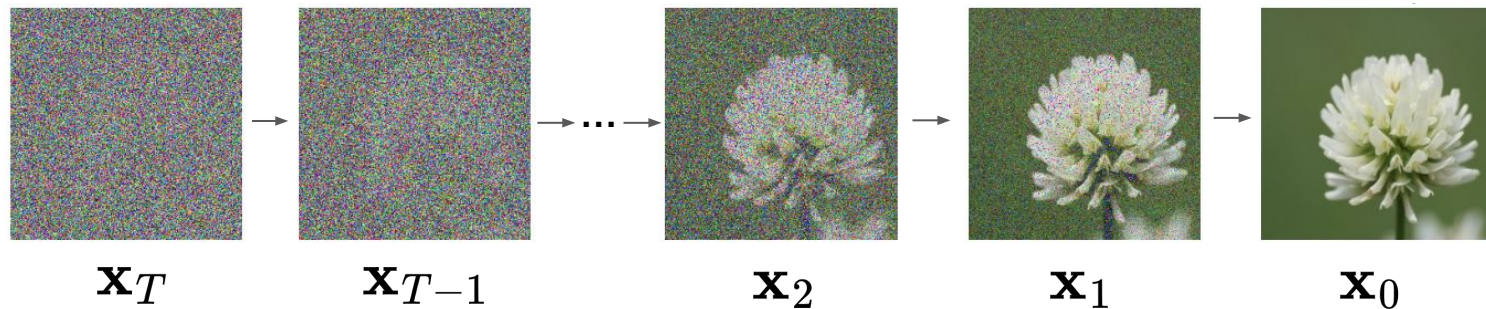


Reverse process:



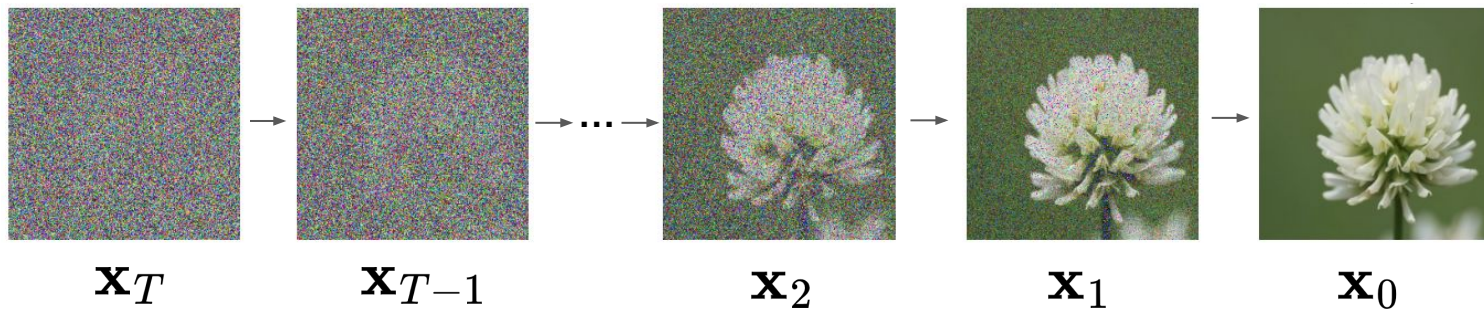
Details: Reverse Process

1. Ideally, sample from reversed conditional distribution $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$



Details: Reverse Process

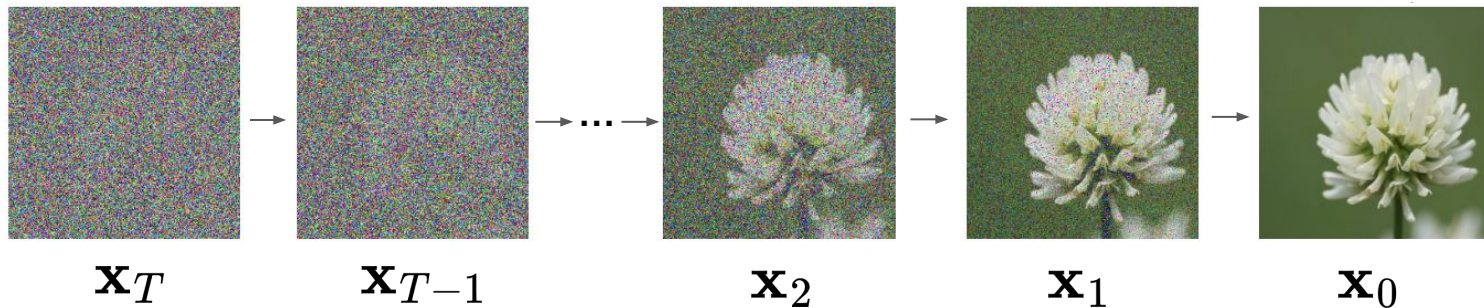
Problem: $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ is **intractable** (can't compute easily)



Details: Reverse Process

Problem: $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ is **intractable** (can't compute easily)

You need to use the entire dataset!

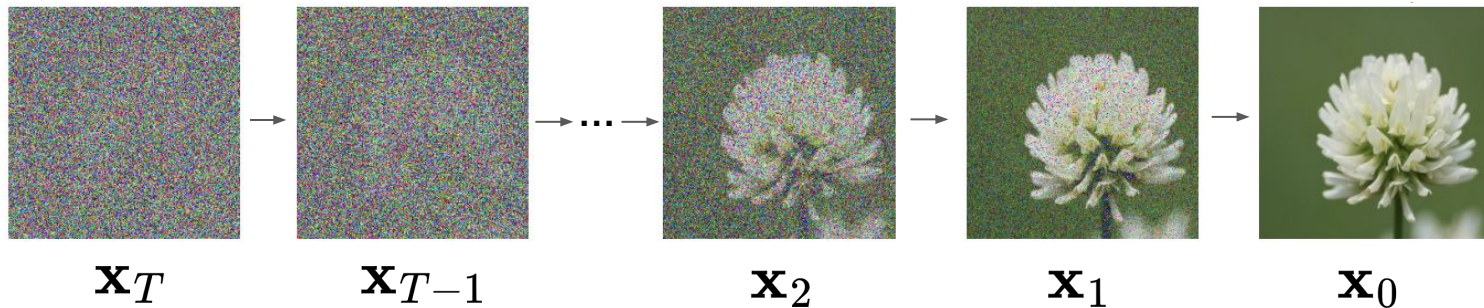


Details: Reverse Process

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

However: $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ is **tractable**

Can reverse the forward process given the original data!



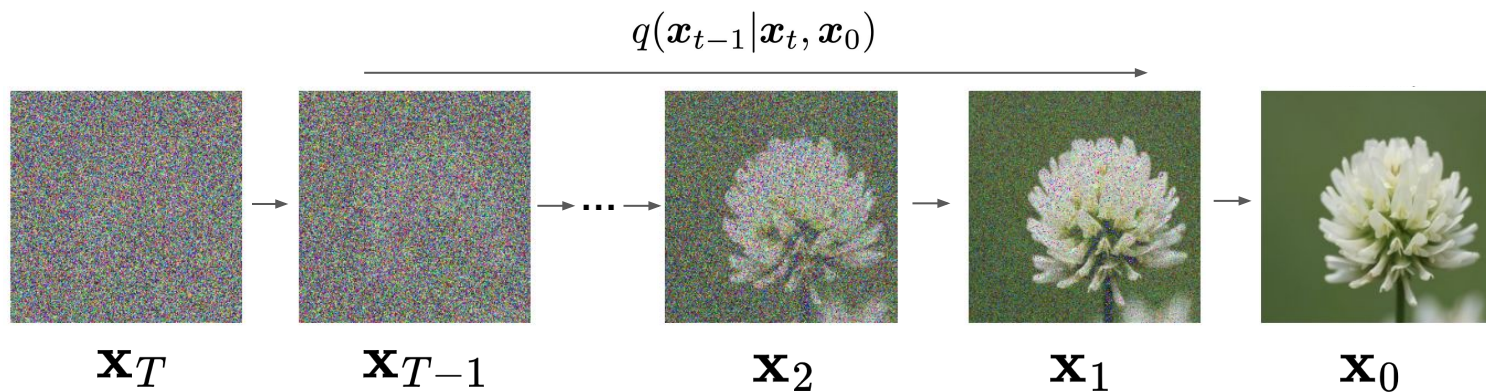
Details: Reverse Process

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

However: $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ is **tractable**

Can reverse the forward process given the original data!

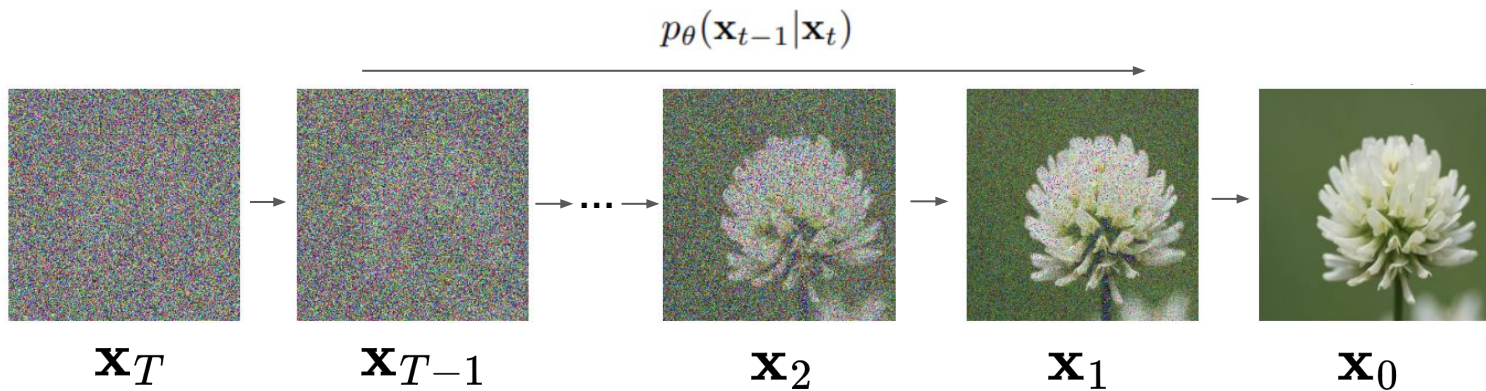
Problem: Don't have any "original data" for inference



Key Idea

We introduce a generative model to approximate the reverse process:

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$$
$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I}) \quad \rightarrow \quad p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$$



Key Idea

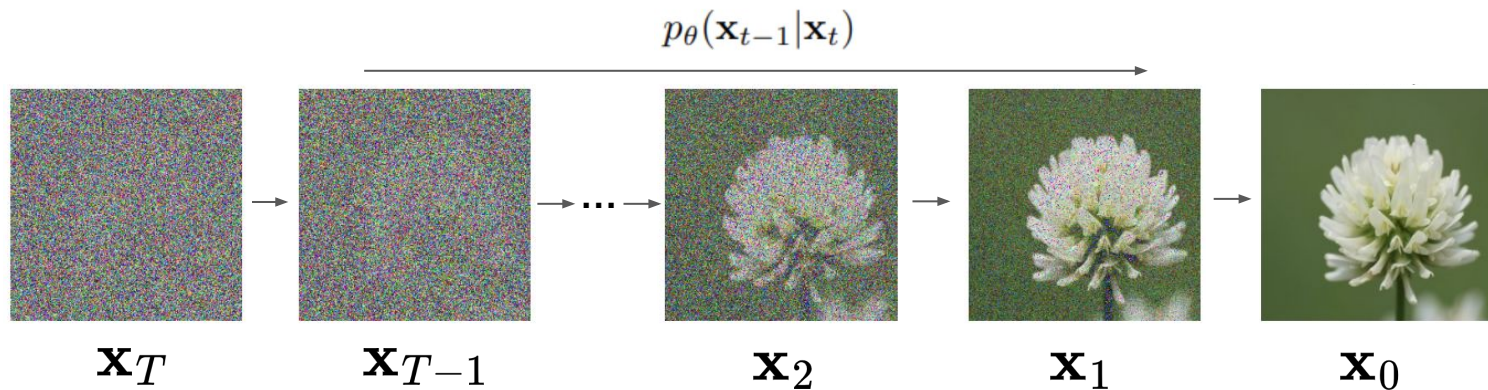
We introduce a generative model to approximate the reverse process:

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I}) \quad \rightarrow \quad p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$$

$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

Learning Objective!

$$\mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t))]$$



Training: Principled Derivation

Find the model that maximizes the likelihood of the training data

i.e. same as VAEs, variational inference; approximate the true posterior

Training Objective

- Bound the likelihood with the ELBO
 - Exactly like hierarchical VAEs

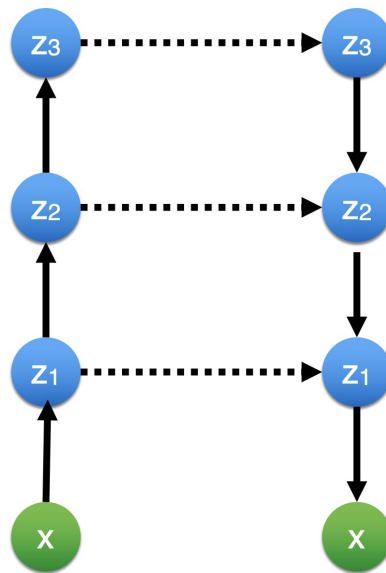
$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]$$

$$= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \underbrace{\sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

Learning Objective!

Discuss:

Differences between diffusion models and hierarchical VAEs?



Inference model
 $q(\mathbf{z}|\mathbf{x})$

Generative model
 $p(\mathbf{x}, \mathbf{z})$

Learning Objective!

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]$$

$$= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} -$$

$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

Training Objective

- Bound the likelihood with the ELBO
 - Exactly like VAEs

$$\begin{aligned}
 \log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\
 &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \underbrace{\sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}
 \end{aligned}$$

Learning Objective!

where $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ is the tractable posterior distribution:

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

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

Parameterizing the Denoising Model

Since both $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ and $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ are Normal distributions, the KL divergence has a simple form:

$$L_{t-1} = D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)) = \mathbb{E}_q \left[\frac{1}{2\sigma_t^2} \|\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) - \mu_\theta(\mathbf{x}_t, t)\|^2 \right] + C$$

Recall that $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon$. [Ho et al. NeurIPS 2020](#) observe that:

$$\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{1 - \beta_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right)$$

They propose to represent the mean of the denoising model using a **noise-prediction network**:

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

With this parameterization

$$L_{t-1} = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\frac{\beta_t^2}{2\sigma_t^2(1 - \beta_t)(1 - \bar{\alpha}_t)} \|\epsilon - \underbrace{\epsilon_\theta(\underbrace{\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}}_{\mathbf{x}_t}, t)}_{\mathbf{x}_t}\|^2 \right] + C$$

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

Training Objective Weighting

ELBO objective leads to a specific regression weight at each time step:

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

Approaches zero!

However, this weight is often very large for small t 's

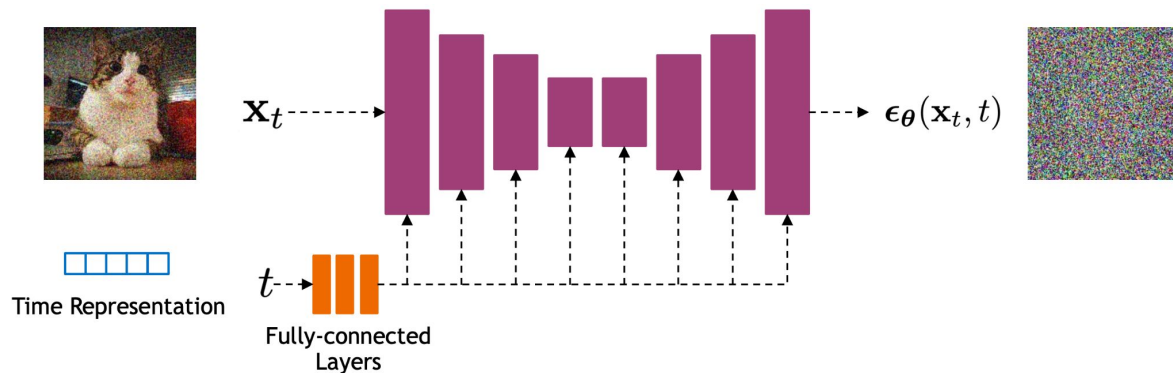
Ho et al., 2020 proposed the following objective to improve perceptual quality:

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

What Network Architecture to Use For ϵ_{θ} ?

People often use U-Nets with residual blocks and self-attention layers at low resolutions

Has same input and output image dimensions



Time representation: sinusoidal positional embeddings

Inject time embedding throughout the network (e.g. additive positional embedding)

Diffusion Results

Outperforms prior generative models when using the **simplified** training objective

ELBO objective performs worse!

Model	IS	FID
Gated PixelCNN [59]	4.60	65.93
Sparse Transformer [7]		
PixelIQN [43]	5.29	49.46
EBM [11]	6.78	38.2
NCSNv2 [56]		31.75
NCSN [55]	8.87 ± 0.12	25.32
SNGAN [39]	8.22 ± 0.05	21.7
SNGAN-DDLS [4]	9.09 ± 0.10	15.42
StyleGAN2 + ADA (v1) [29]	9.74 ± 0.05	3.26
Ours (L , fixed isotropic Σ)	7.67 ± 0.13	13.51
Ours (L_{simple})	9.46 ± 0.11	3.17

ELBO

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

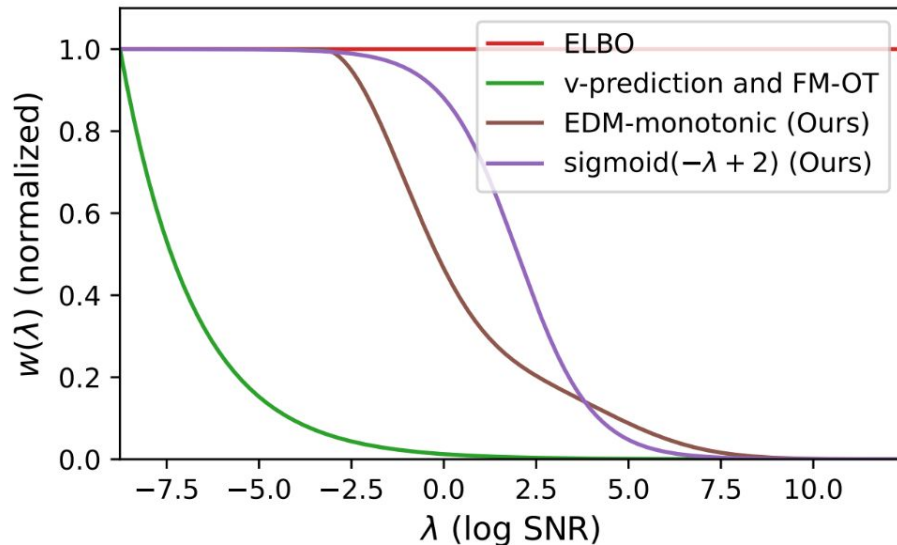
Training Objective Weighting

- ELBO forces the network to model imperceptible details
 - Less modeling capacity dedicated to perceptible details (global image structure, etc.)
- If you care about perceptual quality:
 - Decrease the loss weighting for low noise levels

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

$$\text{SNR}(t) = \bar{\alpha}_t / (1 - \bar{\alpha}_t)$$

$$\log(\text{SNR}(t)) = \log(\bar{\alpha}_t / (1 - \bar{\alpha}_t))$$

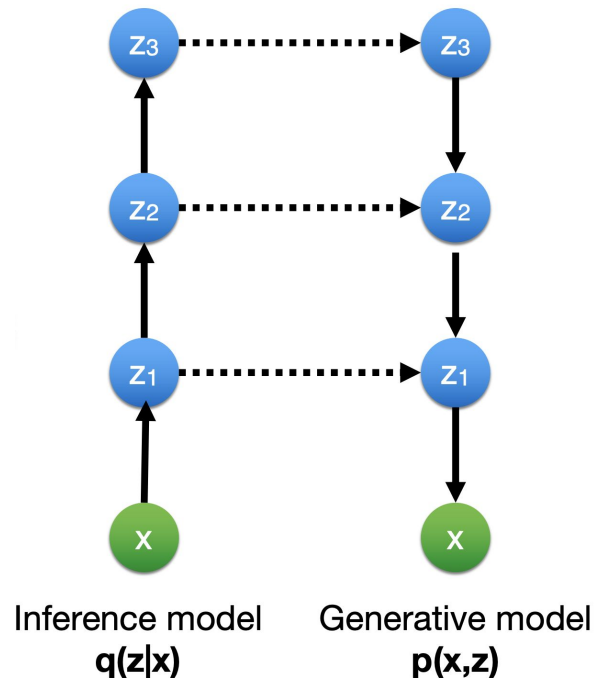


Connection to VAEs

Diffusion models can be considered as a special form of hierarchical VAEs.

However, in diffusion models:

- The inference model is fixed: easier to optimize
- The latent variables have the same dimension as the data.
- The ELBO is decomposed to each time step: fast to train
- Can be made extremely deep (even infinitely deep)
- The model is trained with some reweighting of the ELBO
 - Can trade off likelihood for improved perceptual quality



Alternative Diffusion Parameterization: Data Prediction

Can also view the diffusion network as learning to predict the **original data**

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon \\ \Rightarrow \mathbf{x}_0 &= \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon}{\sqrt{\bar{\alpha}_t}} \\ \Rightarrow \mathbf{x}_\theta(\mathbf{x}_t, t) &= \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}}\end{aligned}$$

Alternative Diffusion Parameterization: Data Prediction

Can also view the diffusion network as learning to predict the **original data**

$$\mathbf{x}_\theta(\mathbf{x}_t, t) = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}}$$

Diffusion training objective: $\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]$

For sampling, want $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$, but don't have access to the original data

Use our estimate of the original data, $\mathbf{x}_\theta(\mathbf{x}_t, t)$, to sample:

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_\theta(\mathbf{x}_t, t)) \approx q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$$

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

Training Algorithm

Repeat until convergence

1. $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ ← Sample original image from image distribution
2. $t \sim U\{1, 2, \dots, T\}$ ← Sample random time step uniformly
3. $\epsilon \sim \mathcal{N}(0, 1)$ ← Sample Gaussian noise
4. Optimizer step on $L(\theta) = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|^2]$
← Model predicts noise applied at time step t and calculate loss

Sampling Algorithm

$\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ← Sample pure Gaussian noise

For $t = T, T - 1 \dots, 1$

$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$ ← Sample Gaussian noise to apply to image

$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ ← Predict noise applied to image and remove that noise

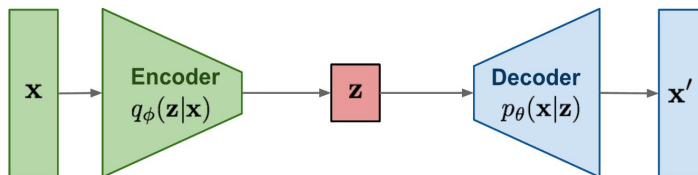
Return \mathbf{x}_0

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = q(\mathbf{x}_{t-1} | \mathbf{x}_t, \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t))$$

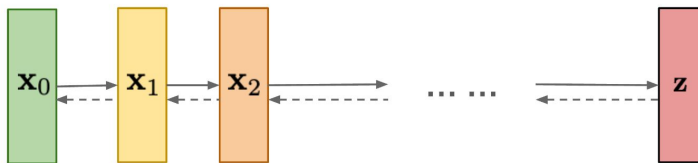


Generative Modeling

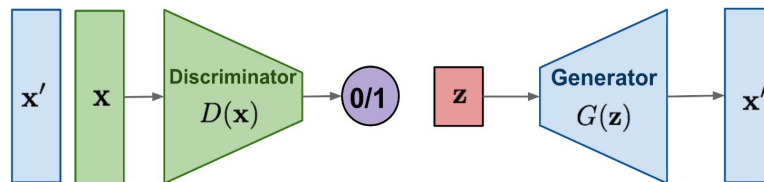
VAE: maximize variational lower bound



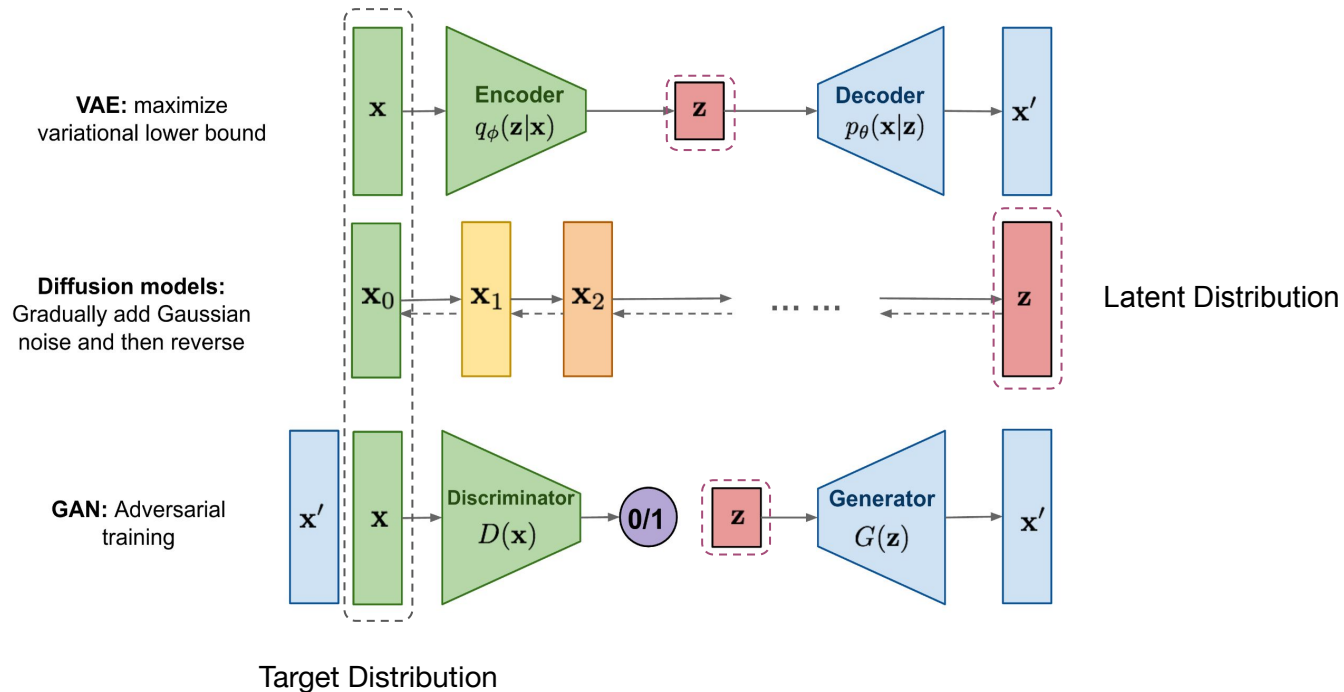
Diffusion models:
Gradually add Gaussian noise and then reverse



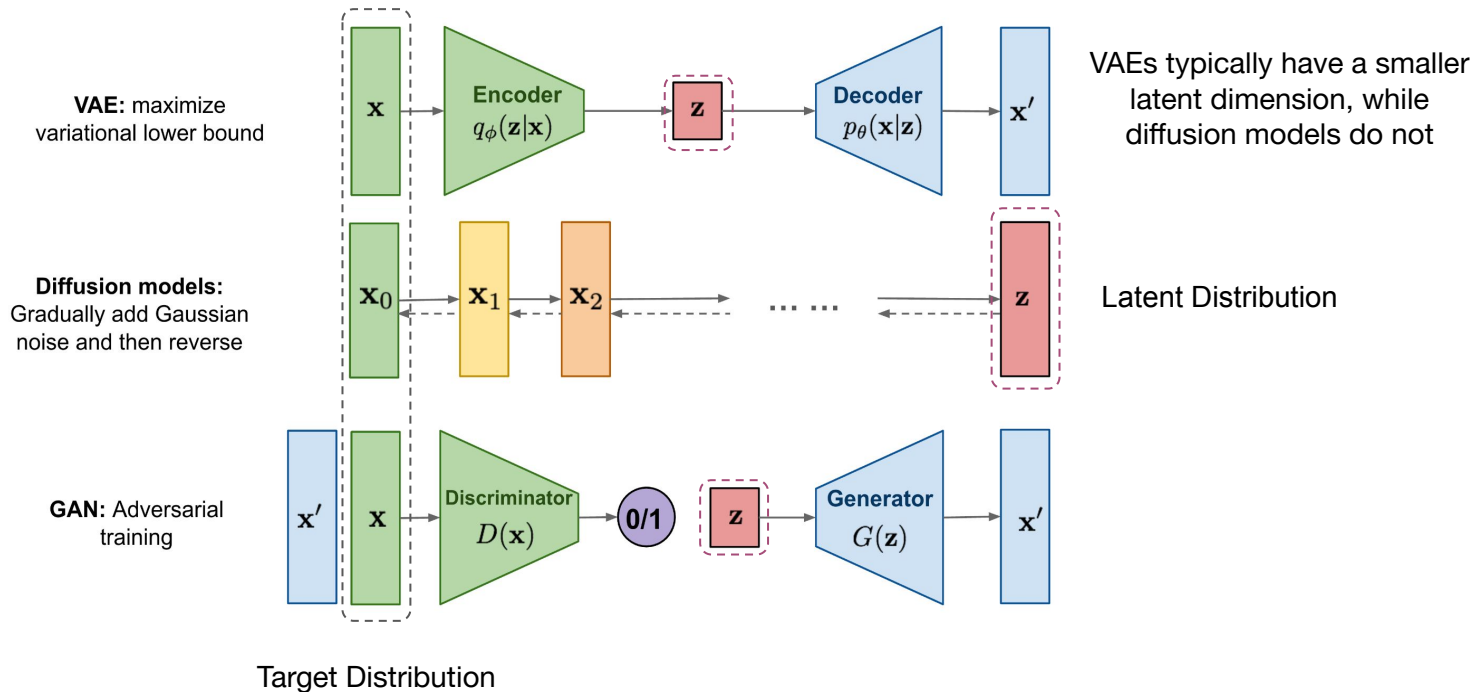
GAN: Adversarial training



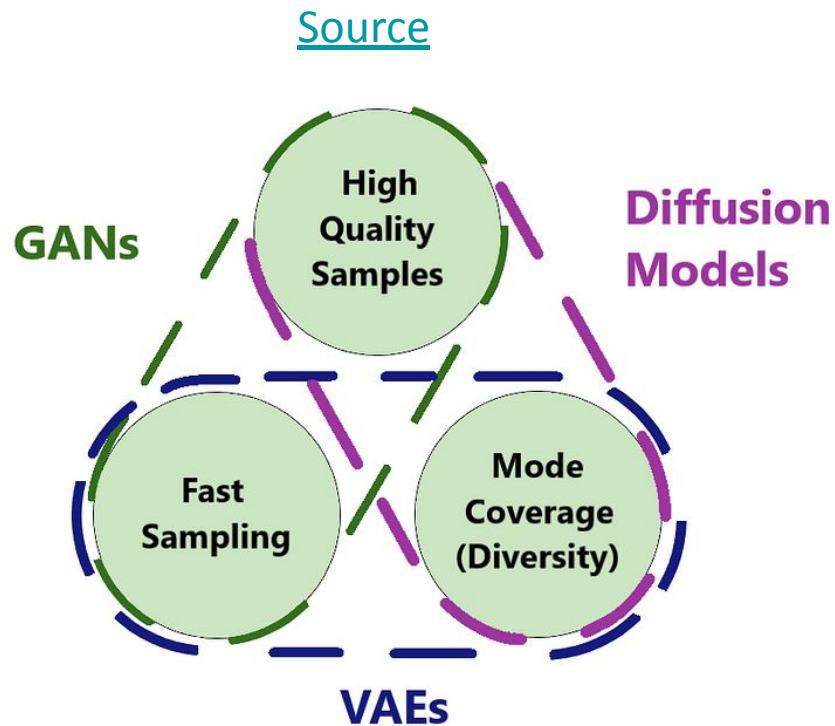
Generative Modeling



Generative Modeling



Diffusion Models vs. VAEs vs. GAN



Cornell Bowers C·IS

Stable Diffusion Demo!

<https://huggingface.co/spaces/stabilityai/stable-diffusion>

Sample input: "messi as a real madrid player"



Recap

- Can bound the likelihood of observed data (i.e. the evidence) with the Evidence Lower Bound (i.e. the ELBO)
- Can learn generative models by maximizing the ELBO
 - VAEs, hierarchical VAEs, Diffusion models
- Diffusion models are a special case of hierarchical VAEs
 - The encoder is fixed to a linear Gaussian model
 - Only learn the decoder
 - Easy to train!
- Learning objective decomposed to each timestep
 - Can be made extremely deep!
 - Can focus on higher noise levels to improve perceptual quality!
- Limitation:
 - Can require many sampling steps for good quality