
- Understanding Diffusion Models: A Unified Perspective

2023.05.03

Jong Jun Won

Preprint

Understanding Diffusion Models: A Unified Perspective

Calvin Luo

Google Research, Brain Team

calvinluo@google.com

August 26, 2022

C Luo 저술 · 2022 · 24회 인용

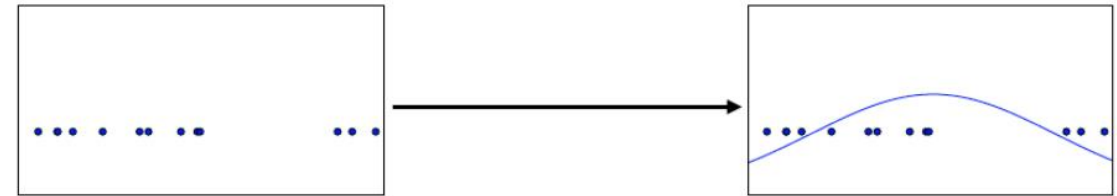
1. Introduction : Generative Models
2. Background – ELBO, VAE, Hierarchical VAE
3. Variational Diffusion Models
4. Three equivalent Interpretations

1. Introduction

Generative Models

- Given observed samples \mathbf{x} , the goal is to learn to model its **true data distribution** $P(\mathbf{x})$
- Data distribution : **the statistical characteristics and patterns** present in the real data
It specifies the probabilities of all events – Introduction to probability
- Once learned, we can generate new samples from our approximate model.
- GAN
- Likelihood-based : **Variational Autoencoders**

- Density estimation



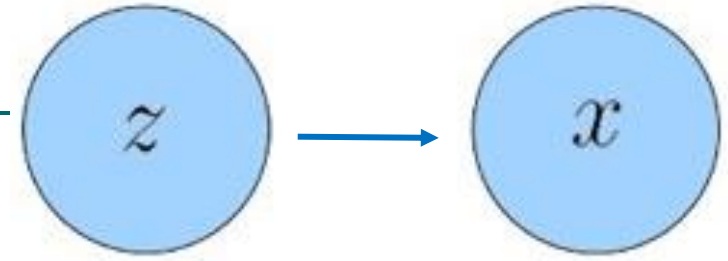
- Sample generation



Training examples

Model samples

2. Background



Latent variable

- We can think of The data(X) as **generated** by an associated unseen latent variable(Z).
- We generally seek to learn **latent representations (low-dimension)** rather than higher-dimensional ones.
- Learning data distribution in lower-dimension is more easier than high-dimension (sparsity, complexity)
- Learning lower-dimensional latents can also be seen as a form of compression
- The best intuition for expressing this idea is through Plato's Allegory of the Cave.

Latent space

🌐 2 languages

[Article](#) [Talk](#)

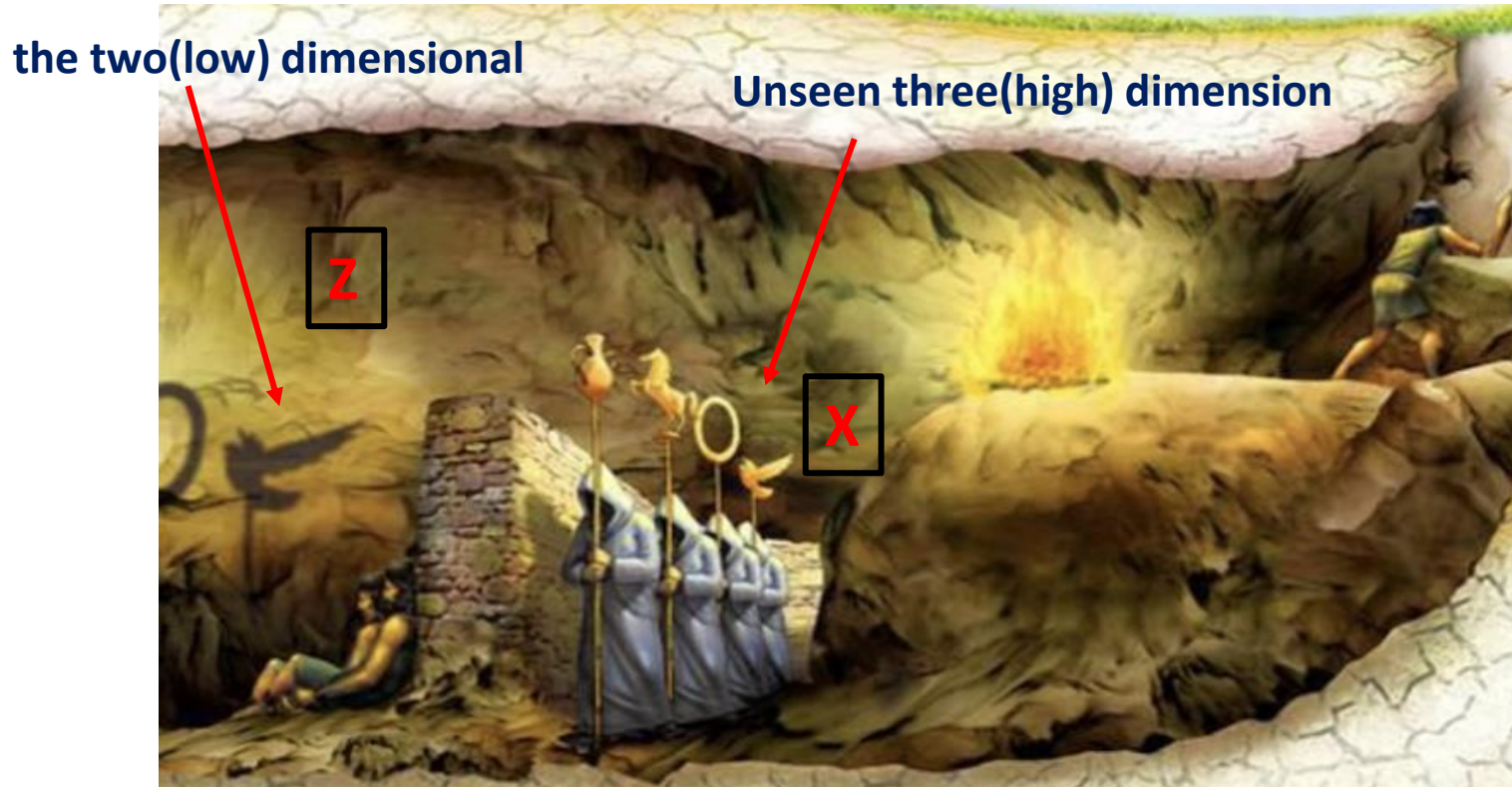
[Read](#) [Edit](#) [View history](#) [Tools](#)

From Wikipedia, the free encyclopedia

A **latent space**, also known as a **latent feature space** or **embedding space**, is an **embedding** of a set of items within a **manifold** in which items resembling each other are positioned closer to one another in the latent space. Position within the latent space can be viewed as being defined by a set of **latent variables** that emerge from the resemblances from the objects.

2. Background

Latent variable



1. Feature Extraction, Compression from dataset ($X \rightarrow Z$)

2. Generation from latent variable ($Z \rightarrow X$)

(Generated) Object : encapsulate abstract properties (size, shape, and more)

The cave people can never see the hidden objects

They can still reason and draw(**generate**) inferences about them.

In a similar way, we can approximate latent representations that describe the data we observe.

2. Background

Evidence Lower Bound

- Likelihood-based : to learn a model to maximize the likelihood $P(x)$
- The latent variables and the data as modeled by a joint distribution $P(x, z)$

$$p(x) = \int p(x, z) dz$$

$$p(x) = \frac{p(x, z)}{p(z|x)}$$

2. Background

Evidence Lower Bound

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

- maximizing the likelihood $p(\mathbf{x})$ is difficult because it either involves integrating out **all latent variables \mathbf{z}** , which is **intractable** for complex models,

in practice *any* solution takes too many resources to be useful,

$$p(\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})}$$

- it involves having access to a ground truth latent encoder $p(\mathbf{z}|\mathbf{x})$

2. Background

Evidence Lower Bound

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

$q_{\phi}(z|x)$: a parameterizable ϕ model that is learned to estimate the true distribution over latent variables for given x

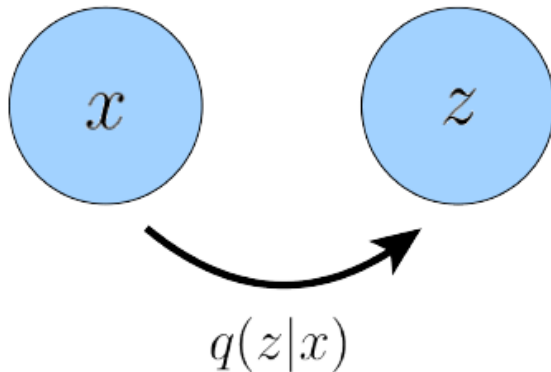
2. Background

Evidence Lower Bound

- the **Evidence Lower Bound** (ELBO) is a lower bound of the log likelihood of $P(x)$ (the evidence).
- Then, **maximizing the ELBO** becomes a proxy objective with which to **optimize a latent variable model**
- $q_{\phi}(z|x)$: parameters ϕ that we seek to optimize true latent variables for given observations x

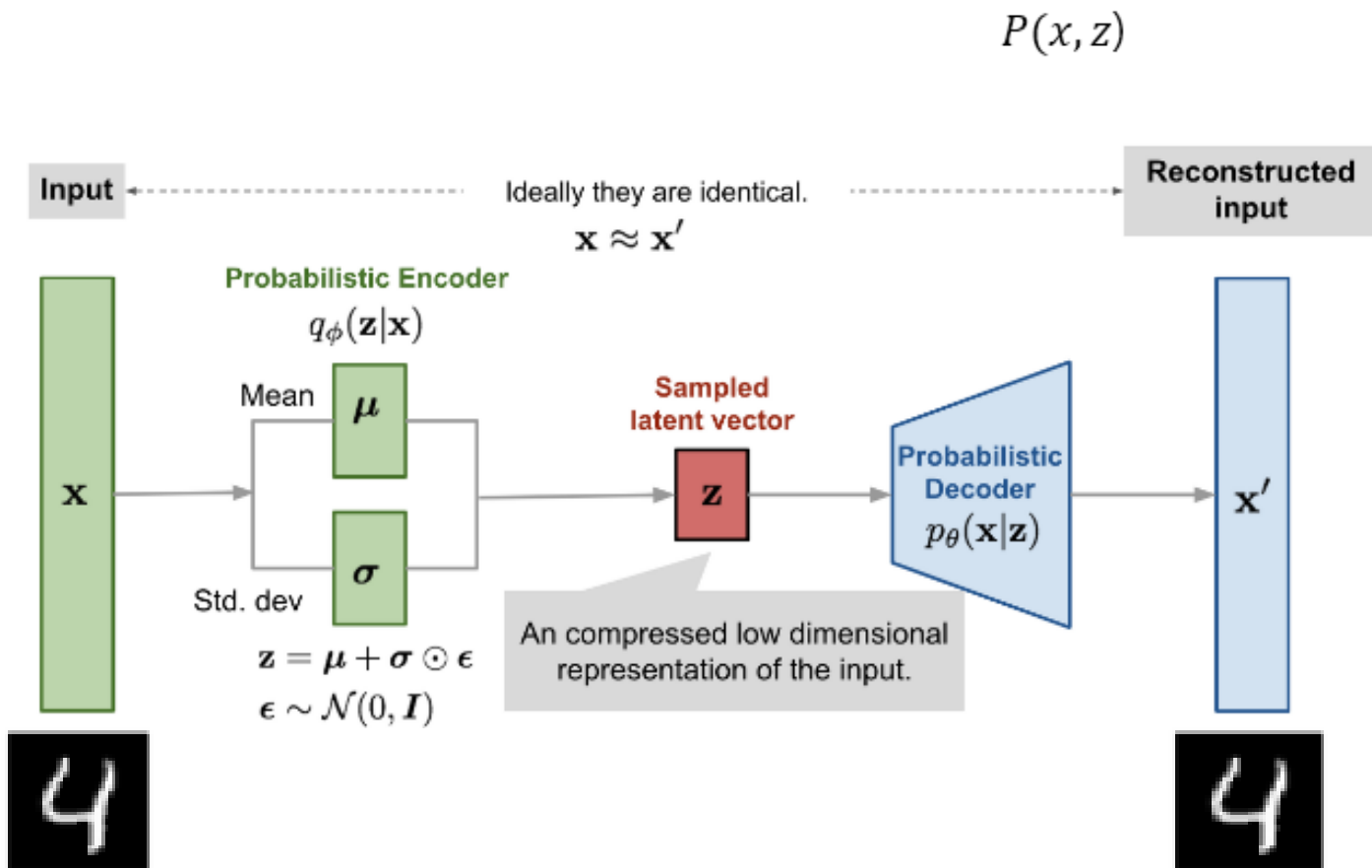
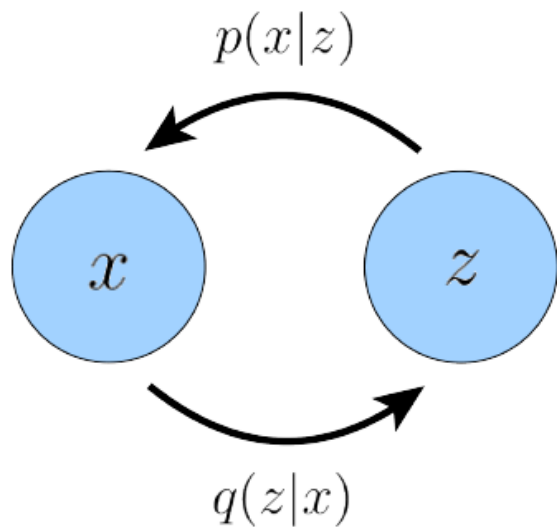
$$\log p(x) \geq \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x, z)}{q_{\phi}(z|x)} \right] \xrightarrow[\text{approximation}]{p(x) = \frac{p(x, z)}{p(z|x)}}$$

- Our goal is to learn this underlying latent structure that describes our observed data



2. Background

Variational Autoencoder

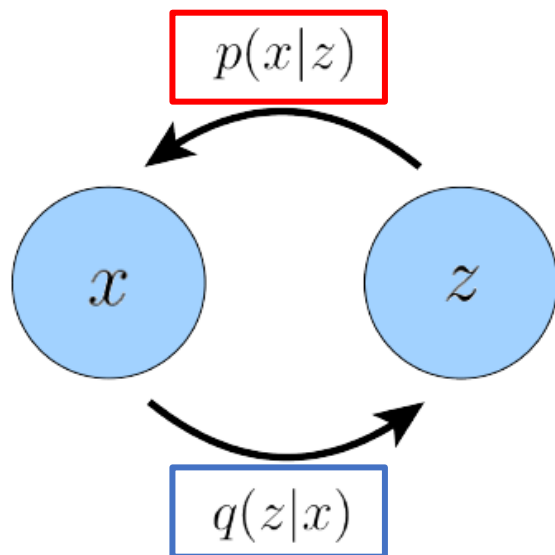


2. Background

Variational Autoencoder

- The ELBO is optimized jointly over parameters ϕ and θ

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

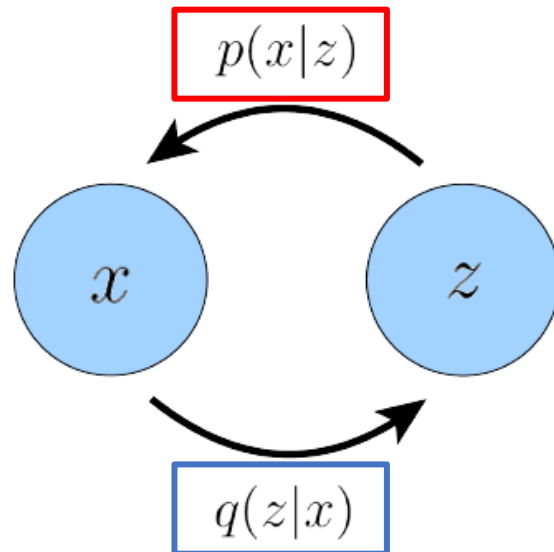


2. Background

Variational Autoencoder

- The encoder : a multivariate Gaussian with diagonal covariance
- The latent : a standard multivariate Gaussian

$$\underbrace{\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_{\phi}(z|x) \parallel p(z))}_{\text{prior matching term}}$$



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

$$z = \mu + \sigma \epsilon \quad \text{with } \epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I})$$

$$q_{\phi}(z|x) = \mathcal{N}(z; \mu_{\phi}(x), \sigma_{\phi}^2(x) \mathbf{I})$$

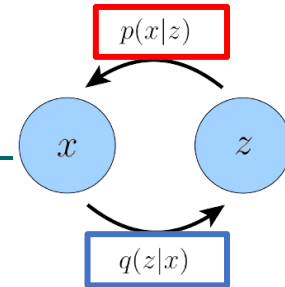
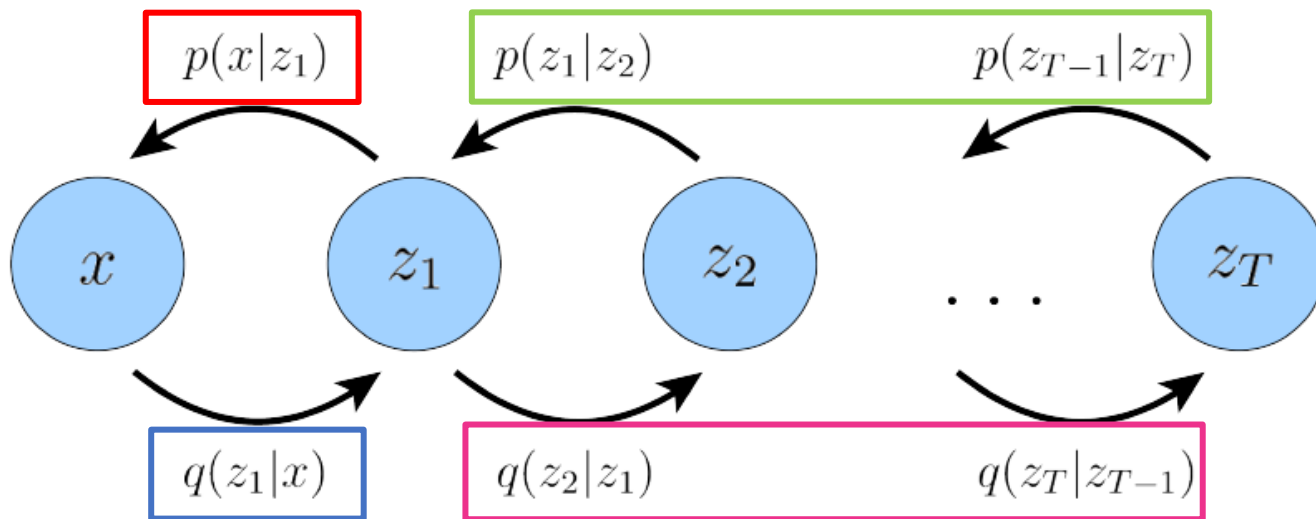
$$z = \mu_{\phi}(x) + \sigma_{\phi}(x) \odot \epsilon \quad \text{with } \epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I})$$

Dot product

2. Background

Hierarchical Variational Autoencoder

- A Hierarchical Variational Autoencoder (HAVE) is a generalization of a VAE that extends to multiple hierarchies over latent variables.
- Markovian HVAE
- In a MHVAE, the generative process is a **Markov chain**,
- **Each latent z_t is generated only from the previous latent.**



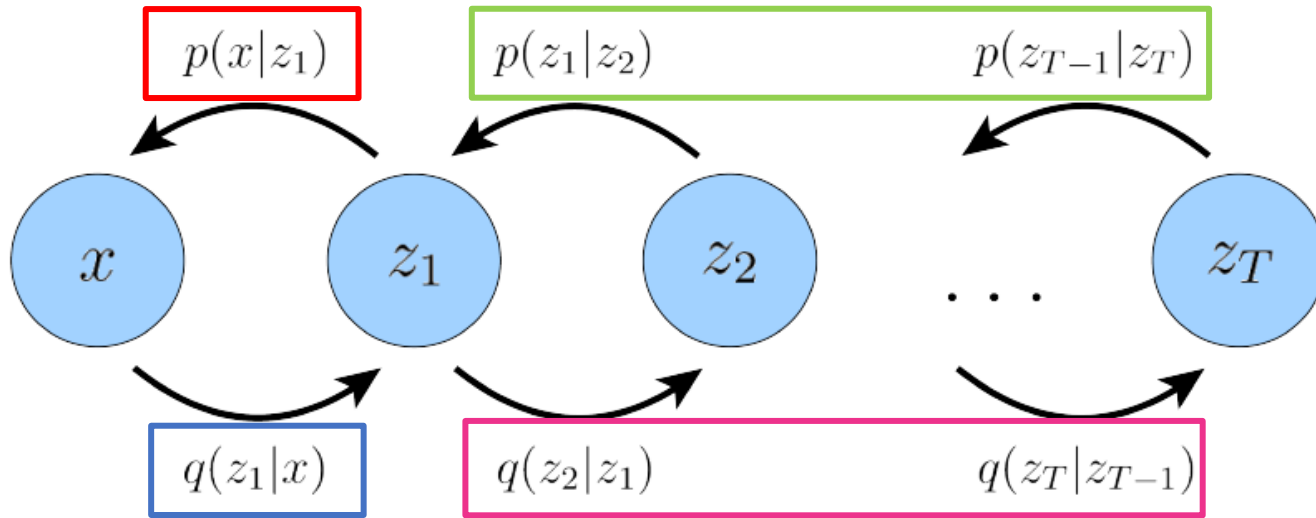
$$\mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x, z)}{q_{\phi}(z|x)} \right]$$

$$p(x, z_{1:T}) = p(z_T) p_{\theta}(x|z_1) \prod_{t=2}^T p_{\theta}(z_{t-1}|z_t)$$

$$q_{\phi}(z_{1:T}|x) = q_{\phi}(z_1|x) \prod_{t=2}^T q_{\phi}(z_t|z_{t-1})$$

2. Background

Hierarchical Variational Autoencoder



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

$$p(x, z_{1:T}) = p(z_T) p_\theta(x|z_1) \prod_{t=2}^T p_\theta(z_{t-1}|z_t)$$

$$q_\phi(z_{1:T}|x) = q_\phi(z_1|x) \prod_{t=2}^T q_\phi(z_t|z_{t-1})$$

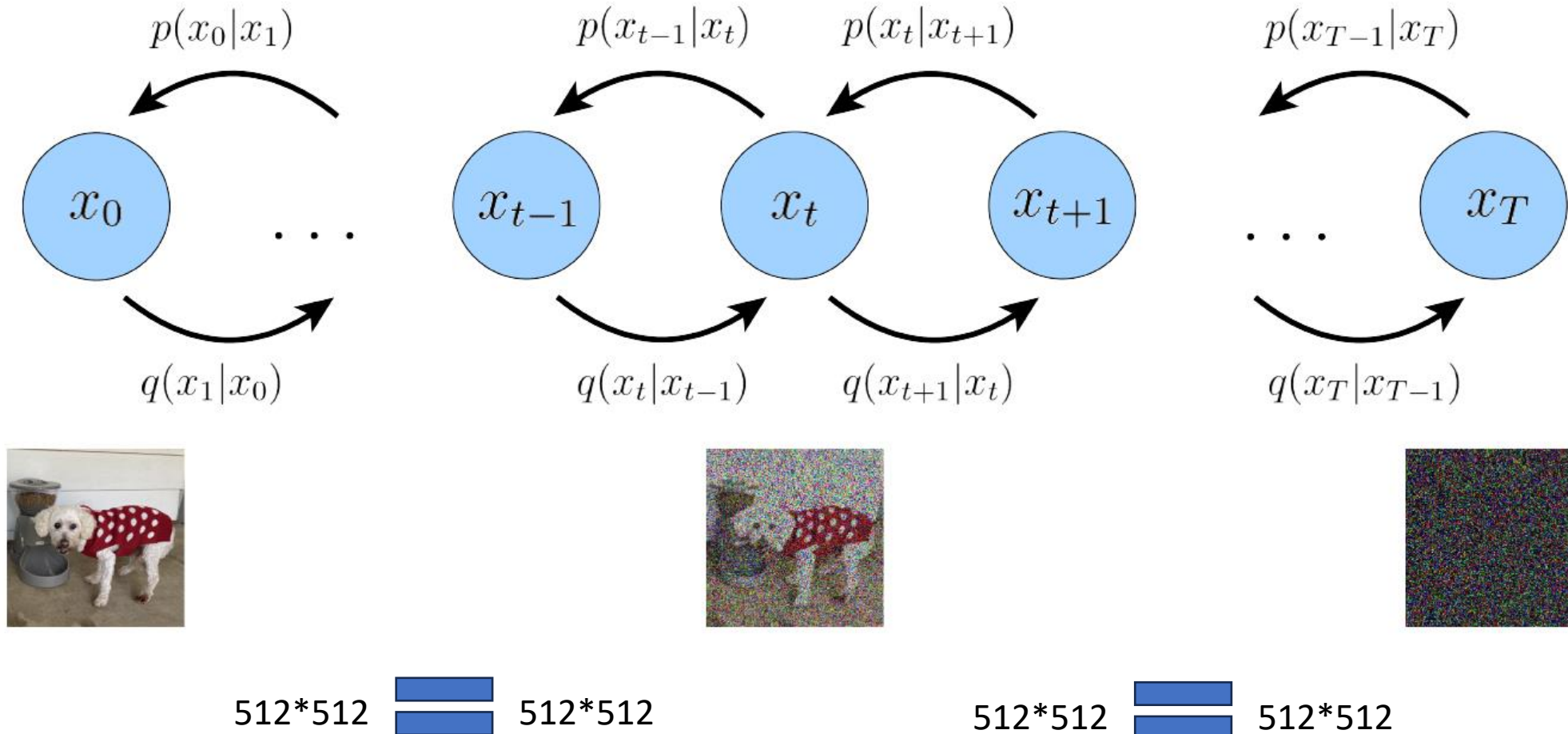
$$\mathbb{E}_{q_\phi(z_{1:T}|x)} \left[\log \frac{p(x, z_{1:T})}{q_\phi(z_{1:T}|x)} \right] = \mathbb{E}_{q_\phi(z_{1:T}|x)} \left[\log \frac{p(z_T) p_\theta(x|z_1) \prod_{t=2}^T p_\theta(z_{t-1}|z_t)}{q_\phi(z_1|x) \prod_{t=2}^T q_\phi(z_t|z_{t-1})} \right]$$

3. Variational Diffusion Models

VDM is as a Markovian HVAE with three key restrictions

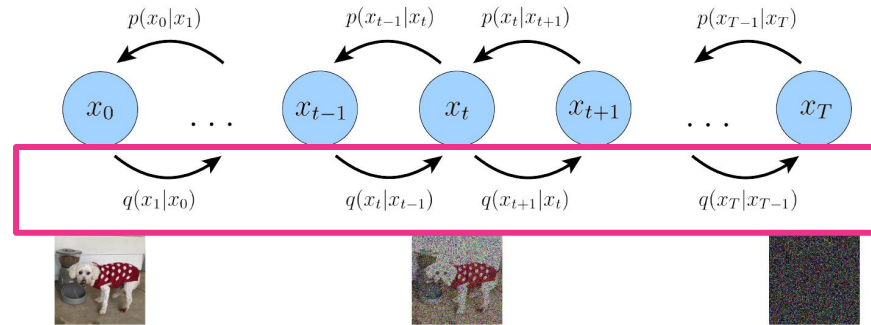
3. Variational Diffusion Models

First, The latent dimension is exactly equal to the data dimension



3. Variational Diffusion Models

Second, The structure of **the latent encoder** at each timestep is **not learned**



$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$$

- it is pre-defined as a linear Gaussian model.



- It is a Gaussian distribution centered around the output of the previous timestep.

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I})$$

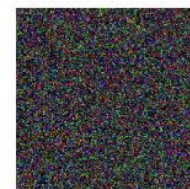
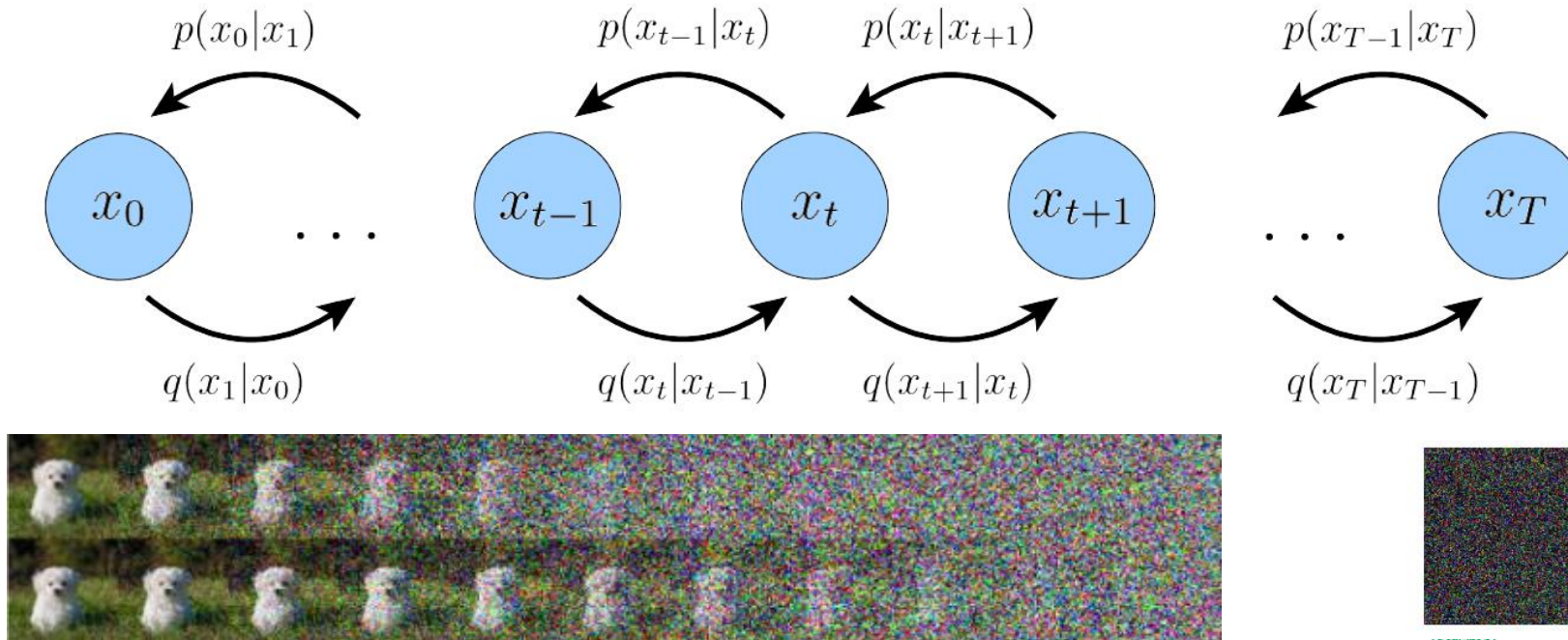
- α_t is coefficient that can vary with the hierarchical depth t
- Signal to noise ratio(SNR) must monotonically decrease over time

3. Variational Diffusion Models

Third, The Gaussian parameters of the latent encoders vary over time

- In such a way that the distribution of the latent at final timestep T is a standard Gaussian

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I})$$



#. HOW Diffusion, WHY Gaussian?

Referred From

Deep Unsupervised Learning using Nonequilibrium Thermodynamics

J Sohl-Dickstein 저술 · 2015 ·

- inspired by **non-equilibrium statistical physics**, is to systematically and slowly destroy structure in a data distribution through an *iterative forward diffusion process*.
- We then *learn a reverse diffusion process* that restores structure in data, yielding a highly flexible and tractable generative model of the data

		<i>Gaussian</i>	<i>Binomial</i>
Well behaved (analytically tractable) distribution	$\pi(\mathbf{x}^{(T)}) =$	$\mathcal{N}(\mathbf{x}^{(T)}; \mathbf{0}, \mathbf{I})$	$\mathcal{B}(\mathbf{x}^{(T)}; 0.5)$
Forward diffusion kernel	$q(\mathbf{x}^{(t)} \mathbf{x}^{(t-1)}) =$	$\mathcal{N}(\mathbf{x}^{(t)}; \mathbf{x}^{(t-1)}\sqrt{1-\beta_t}, \mathbf{I}\beta_t)$	$\mathcal{B}(\mathbf{x}^{(t)}; \mathbf{x}^{(t-1)}(1-\beta_t) + 0.5\beta_t)$
Reverse diffusion kernel	$p(\mathbf{x}^{(t-1)} \mathbf{x}^{(t)}) =$	$\mathcal{N}(\mathbf{x}^{(t-1)}; \mathbf{f}_\mu(\mathbf{x}^{(t)}, t), \mathbf{f}_\Sigma(\mathbf{x}^{(t)}, t))$	$\mathcal{B}(\mathbf{x}^{(t-1)}; \mathbf{f}_b(\mathbf{x}^{(t)}, t))$

#. HOW Diffusion, WHY Gaussian?

Referred From

ON THE THEORY OF STOCHASTIC PROCESSES, WITH PARTICULAR REFERENCE TO APPLICATIONS

W. FELLER

CORNELL UNIVERSITY

W Feller 저술 - 1949

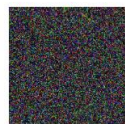
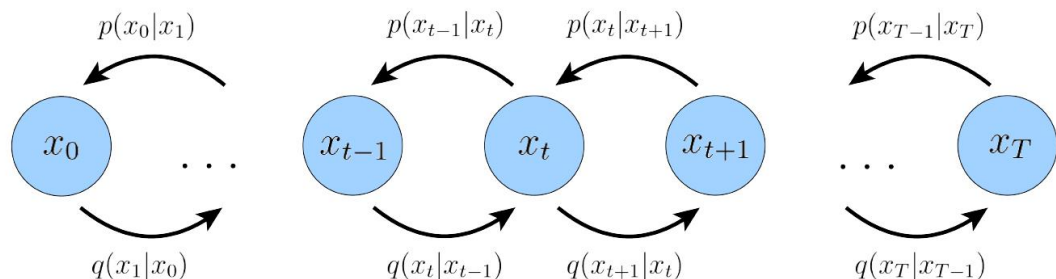
(ONR. Project for Research in Probability)

HOW?

Markov process : differential

Forward, Backward의 확률 분포가 같다.

-> 원본 분포의 추정이 가능하다.



Markov processes leading to ordinary differential equations

Under these conditions²¹ $u(\tau, \xi; t, x)$ satisfies the “forward equation”

$$u_t(\tau, \xi; t, x) = \frac{1}{2} [a(t, x)u(\tau, \xi; t, x)]_{xx} + [b(t, x)u(\tau, \xi; t, x)]_x,$$

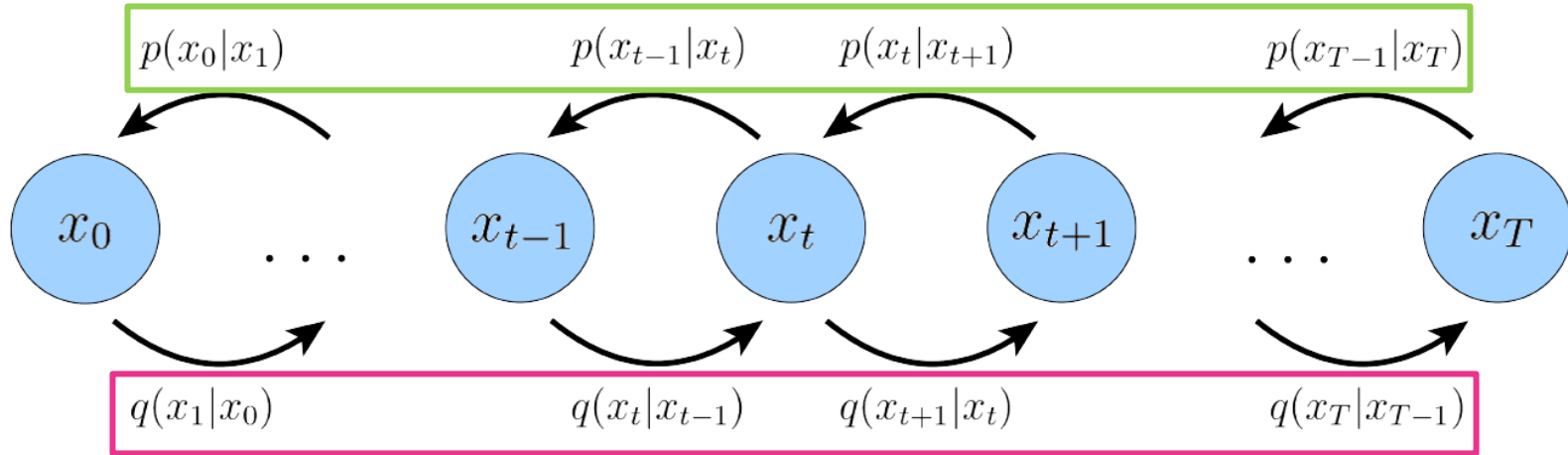
and the “backward equation”

$$u_\tau(\tau, \xi; t, x) = \frac{1}{2} a(\tau, \xi)u_{\xi\xi}(\tau, \xi; t, x) + b(\tau, \xi)u_\xi(\tau, \xi; t, x);$$

3. Variational Diffusion Models

Basic Formula

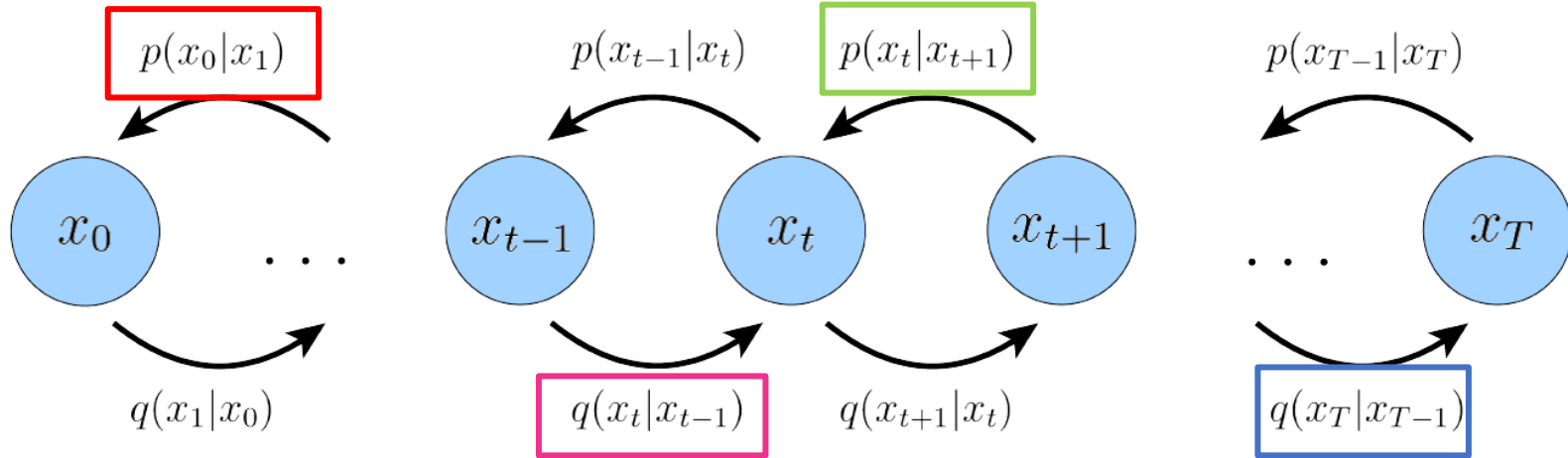
$$\begin{aligned} & \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \end{aligned}$$



3. Variational Diffusion Models

Basic Formula

$$\begin{aligned} & \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \end{aligned}$$



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

3. Variational Diffusion Models

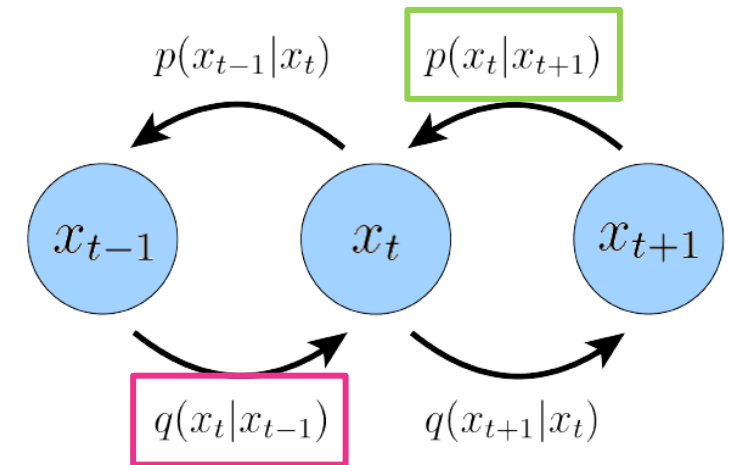
Basic Formula

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

However, actually optimizing the ELBO using the terms we just derived might be suboptimal; because the consistency term is computed as an expectation over two random variables.

the variance of its Monte Carlo estimate could potentially be higher than a term that is estimated using only one random variable per timestep.

Let us try to derive a form for our ELBO where each term is computed as an expectation over only one random variable at a time



3. Variational Diffusion Models

Basic Formula

The key insight is that we can rewrite encoder transitions

$$q(x_t|x_{t-1}) \Leftarrow$$

$$= q(x_t|x_{t-1}, x_{t-2}) \Leftarrow$$

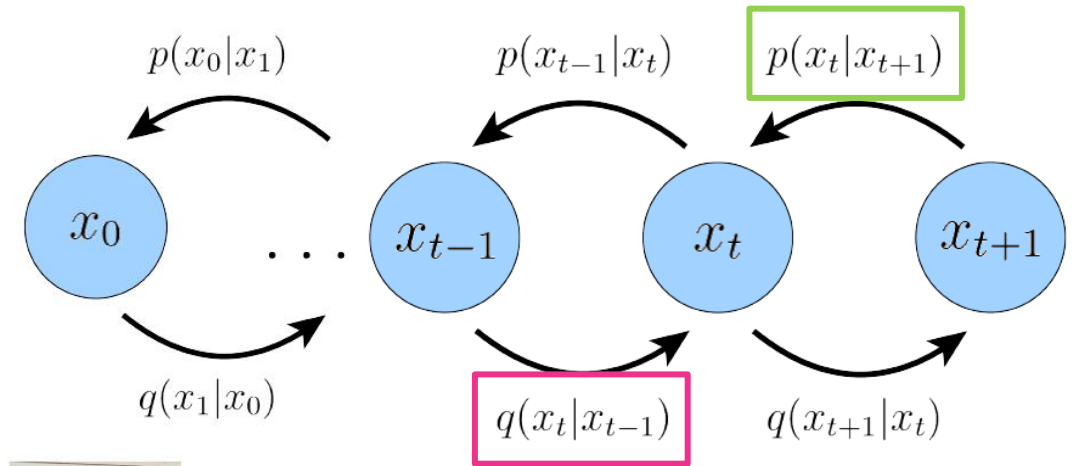
$$= q(x_t|x_{t-1}, x_{t-2}, x_1) \Leftarrow$$

$$= q(x_t|x_{t-1}, x_1, x_0) \Leftarrow$$

$$= q(x_t|x_{t-1}, x_0) \Leftarrow$$

$$q(x_t|x_{t-1}, x_0) = \frac{q(x_{t-1}|x_t, x_0)q(x_t|x_0)}{q(x_{t-1}|x_0)}$$

$$\mathbb{E}_{q(x_{1:T}|x_0)} \left[\log \frac{p(x_T)p_\theta(x_0|x_1) \prod_{t=2}^T p_\theta(x_{t-1}|x_t)}{q(x_1|x_0) \prod_{t=2}^T q(x_t|x_{t-1}, x_0)} \right]$$



3. Variational Diffusion Models

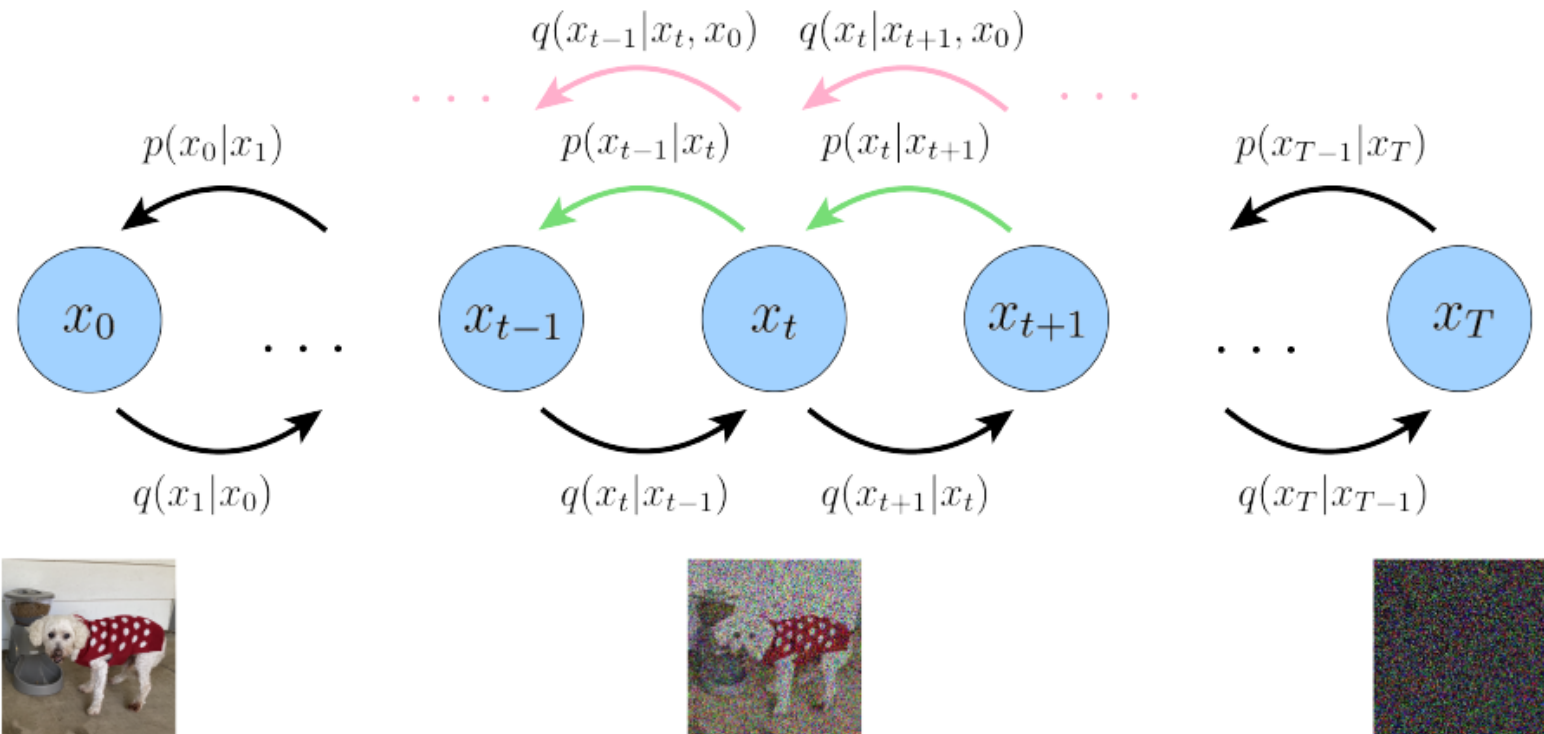
Basic Formula

$$q(x_t|x_{t-1}, x_0) = \frac{q(x_{t-1}|x_t, x_0)q(x_t|x_0)}{q(x_{t-1}|x_0)}$$

$$\mathbb{E}_{q(x_{1:T}|x_0)} \left[\log \frac{p(x_T)p_\theta(x_0|x_1) \prod_{t=2}^T p_\theta(x_{t-1}|x_t)}{q(x_1|x_0) \prod_{t=2}^T q(x_t|x_{t-1}, x_0)} \right]$$

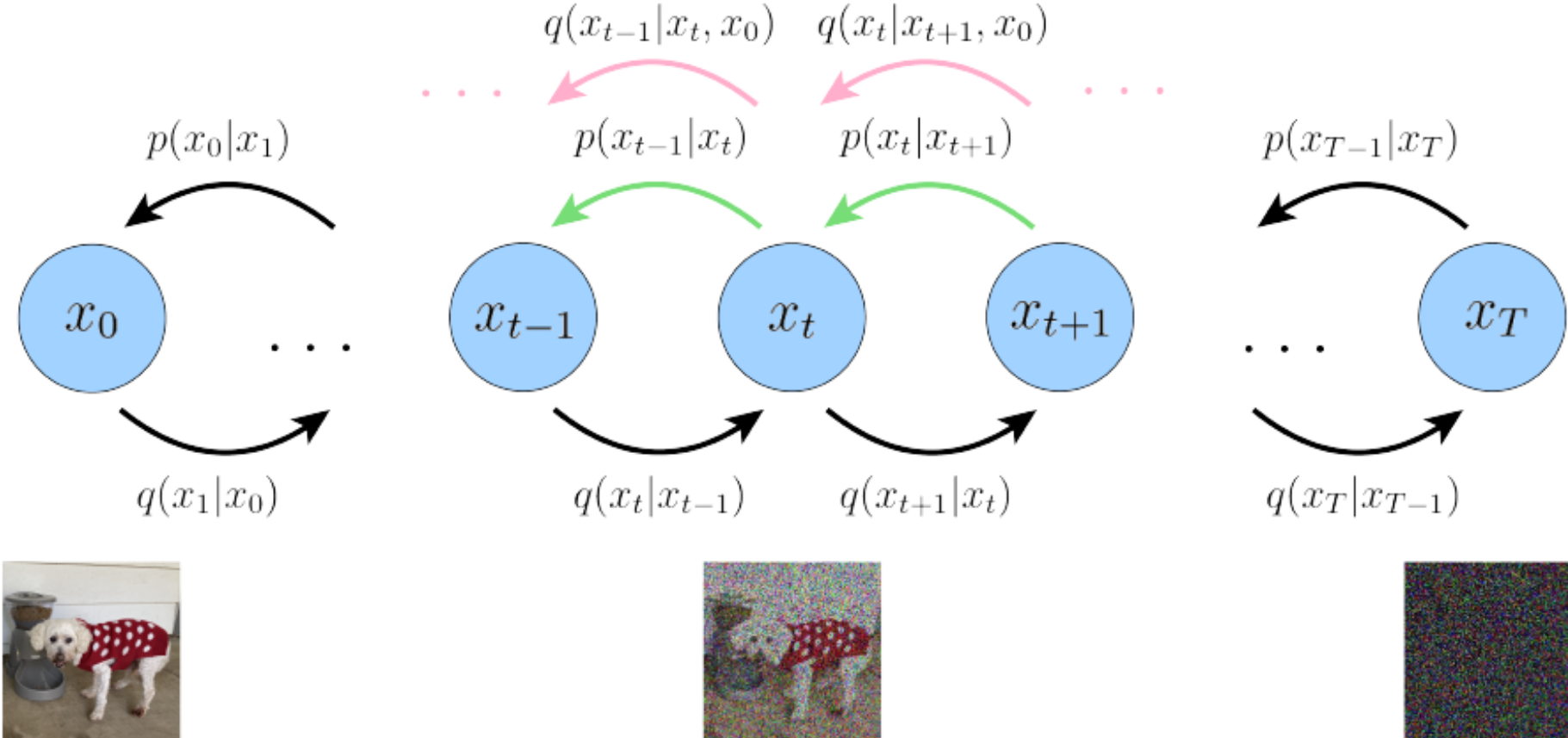
$$= \underbrace{\mathbb{E}_{q(x_1|x_0)} [\log p_\theta(x_0|x_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(x_T|x_0) \parallel p(x_T))}_{\text{prior matching term}}$$

$$- \sum_{t=2}^T \underbrace{\mathbb{E}_{q(x_t|x_0)} [D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_\theta(x_{t-1}|x_t))]}_{\text{denoising matching term}}$$



3. Variational

Basic Formula



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

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

3. Variational Diffusion Models

Formula

$$[D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t))]$$

$$q(x_t|x_{t-1}, x_0) = \frac{q(x_{t-1}|x_t, x_0)q(x_t|x_0)}{q(x_{t-1}|x_0)}$$

$$\propto \mathcal{N}(x_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}}_{\mu_q(x_t, x_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I}}_{\Sigma_q(t)})$$

$$p(x^{(t-1)}|x^{(t)}) = \mathcal{N}(x^{(t-1)}; f_{\mu}(x^{(t)}, t), f_{\Sigma}(x^{(t)}, t))$$

$$u(\tau, \xi; t, x) = \frac{1}{2\{\pi(t - \tau)\}^{1/2}} \exp\left\{-\frac{(x - \xi)^2}{4(t - \tau)}\right\}$$

Gaussian distribution.

3. Variational Diffusion Models

Formula

$$u(\tau, \xi; t, x) = \frac{1}{2\{\pi(t - \tau)\}^{1/2}} \exp \left\{ -\frac{(x - \xi)^2}{4(t - \tau)} \right\}$$

Gaussian distribution.

$$[D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t))]$$

$$p(x^{(t-1)}|x^{(t)}) = \mathcal{N}(x^{(t-1)}; f_{\mu}(x^{(t)}, t), f_{\Sigma}(x^{(t)}, t))$$

$$\propto \mathcal{N}(x_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}}_{\mu_q(x_t, x_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I}}_{\Sigma_q(t)})$$

$$D_{\text{KL}}(\mathcal{N}(x; \mu_x, \Sigma_x) \parallel \mathcal{N}(y; \mu_y, \Sigma_y)) = \frac{1}{2} \left[\log \frac{|\Sigma_y|}{|\Sigma_x|} - d + \text{tr}(\Sigma_y^{-1} \Sigma_x) + (\mu_y - \mu_x)^T \Sigma_y^{-1} (\mu_y - \mu_x) \right]$$

$$\arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\mu_{\theta} - \mu_q\|_2^2 \right]$$

μ_{θ} as shorthand for $\mu_{\theta}(x_t, t)$

μ_q as shorthand for $\mu_q(x_t, x_0)$

3. Variational Diffusion Models

Formula

$$[D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t))]$$

$$\arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\mu_{\theta} - \mu_q\|_2^2 \right] \quad \begin{array}{l} \mu_{\theta} \text{ as shorthand for } \mu_{\theta}(x_t, t) \\ \mu_q \text{ as shorthand for } \mu_q(x_t, x_0) \end{array}$$

$$\mu_q(x_t, x_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}$$

$$\mu_{\theta}(x_t, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{x}_{\theta}(x_t, t)}{1 - \bar{\alpha}_t}$$

$\hat{x}_{\theta}(x_t, t)$ is parameterized by a neural network that seeks to predict original image from noisy image and time index.

3. Variational Diffusion Models

Formula

$$\mu_q(x_t, x_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t} \quad \mu_q \text{ as shorthand for } \mu_q(x_t, x_0)$$

$$\mu_\theta(x_t, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{x}_\theta(x_t, t)}{1 - \bar{\alpha}_t} \quad \mu_\theta \text{ as shorthand for } \mu_\theta(x_t, t)$$

$\hat{x}_\theta(x_t, t)$ is parameterized by a neural network that seeks to predict original image from noisy image and time index.

$$[D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_\theta(x_{t-1}|x_t))]$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\|\mu_\theta - \mu_q\|_2^2 \right]$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)^2} \left[\|\hat{x}_\theta(x_t, t) - x_0\|_2^2 \right] \quad \longrightarrow$$

Optimizing VDM boils down to **learning a neural network to predict the original ground truth image** from an arbitrarily noisified version of it

4. Three equivalent Interpretations

Formula

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}(1-\alpha_t)^2}{(1-\bar{\alpha}_t)^2} \left[\|\hat{x}_{\theta}(x_t, t) - x_0\|_2^2 \right]$$

$$\mu_q(x_t, x_0) = \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)x_0}{1-\bar{\alpha}_t}$$

$$x_0 = \frac{x_t - \sqrt{1-\bar{\alpha}_t}\epsilon_0}{\sqrt{\bar{\alpha}_t}}$$

$$= \frac{1}{\sqrt{\alpha_t}}x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}}\epsilon_0$$

$$\mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha_t}}x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}}\hat{\epsilon}_{\theta}(x_t, t)$$

$\hat{\epsilon}_{\theta}(x_t, t)$ is parameterized by a neural network that learns to predict the source noise

$$[D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t))]$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1-\alpha_t)^2}{(1-\bar{\alpha}_t)\alpha_t} \left[\|\epsilon_0 - \hat{\epsilon}_{\theta}(x_t, t)\|_2^2 \right] \longrightarrow$$

learning a VDM by predicting the original image is equivalent to learning to predict the noise

4. Three equivalent Interpretations

Formula

$$\mu_q(x_t, x_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t} = \frac{1}{\sqrt{\alpha_t}}x_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}}\nabla \log p(x_t)$$

$$\mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha_t}}x_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}}s_{\theta}(x_t, t)$$

$$[D_{\text{KL}}(q(x_{t-1}|x_t, x_0) || p_{\theta}(x_{t-1}|x_t))]$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{\alpha_t} \left[\|s_{\theta}(x_t, t) - \nabla \log p(x_t)\|_2^2 \right]$$

Here, $s_{\theta}(x_t, t)$ is a neural network that learns to predict the score function $\nabla_{x_t} \log p(x_t)$, which is the gradient of x_t in data space, for any arbitrary noise level t .

4. Three equivalent Interpretations

Formula

Mathematically, for a Gaussian variable $z \sim \mathcal{N}(z; \mu_z, \Sigma_z)$.

Tweedie's Formula states that : $\mathbb{E} [\mu_z | z] = z + \Sigma_z \nabla_z \log p(z)$

$$q(x_t | x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

$$\mathbb{E} [\mu_{x_t} | x_t] = x_t + (1 - \bar{\alpha}_t) \nabla_{x_t} \log p(x_t)$$

$$\sqrt{\bar{\alpha}_t} x_0 = x_t + (1 - \bar{\alpha}_t) \nabla \log p(x_t)$$

$$\therefore x_0 = \frac{x_t + (1 - \bar{\alpha}_t) \nabla \log p(x_t)}{\sqrt{\bar{\alpha}_t}}$$

$$\mu_q(x_t, x_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t} = \frac{1}{\sqrt{\alpha_t}}x_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \nabla \log p(x_t)$$

5. Conclusion

Formula

Three equivalent objectives to optimize a VDM:

1. To predict the original image
2. To predict the source noise
3. To predict $\nabla \log p(x_t)$ at an arbitrary noise level