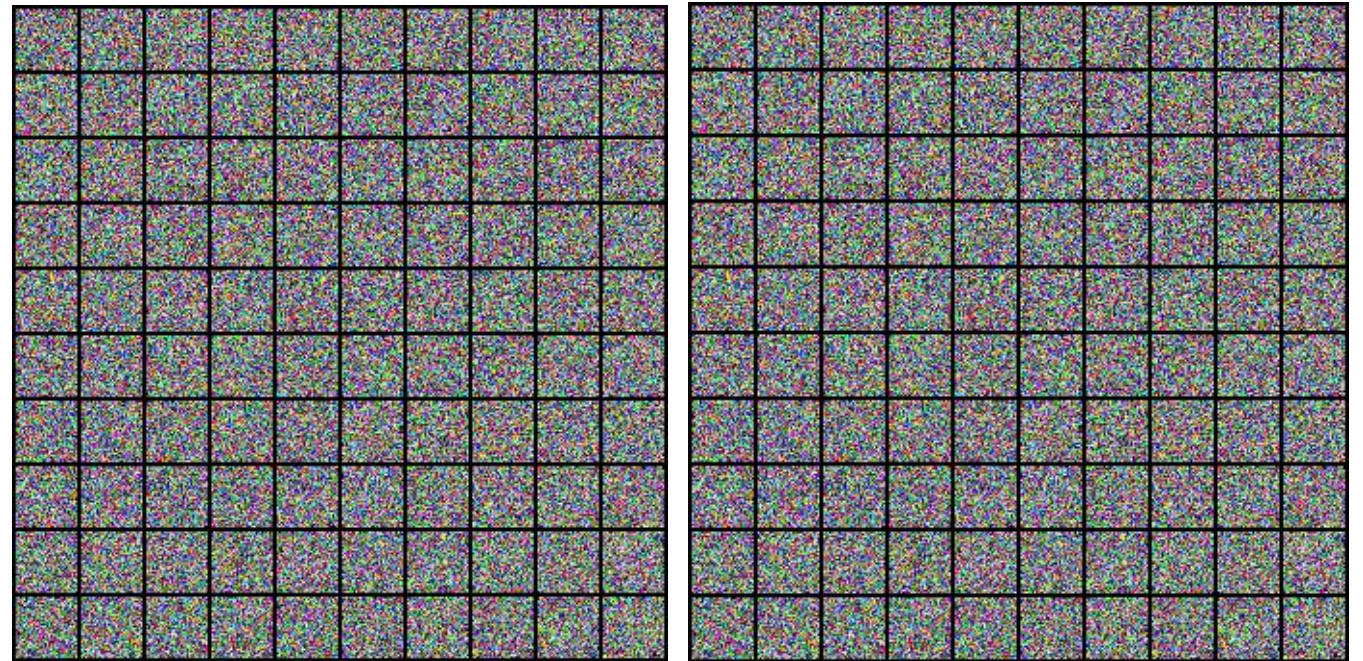




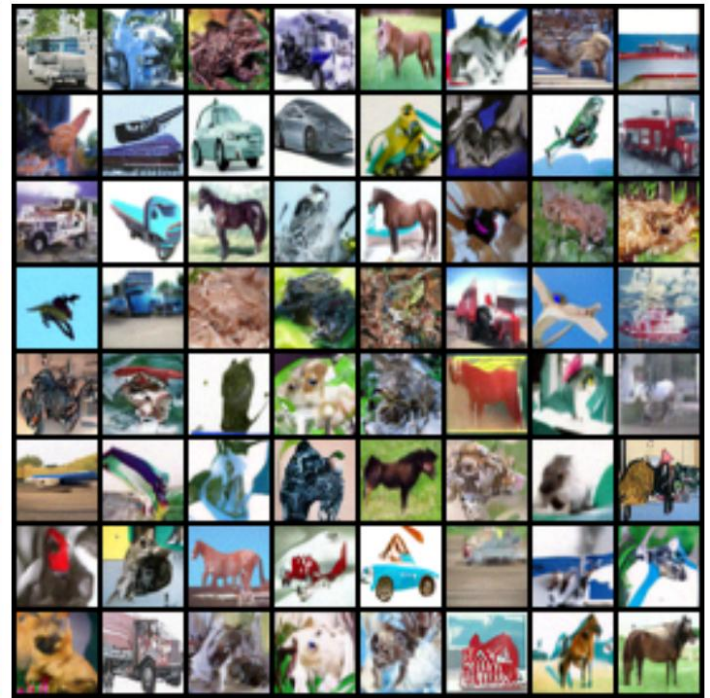
The Art of Diffusion Model
- Discrete Perspective

Bangyan Liao

liaobangyan@westlake.edu.cn



3 0 6 8 0 4 9 4
2 2 6 3 5 8 9 2
0 0 3 3 6 3 4 2
0 6 0 4 7 3 4 4
5 5 2 4 0 2 7 8
5 2 2 8 9 3 6 1
2 1 6 7 7 2 6 0
6 8 6 4 6 5 8 6

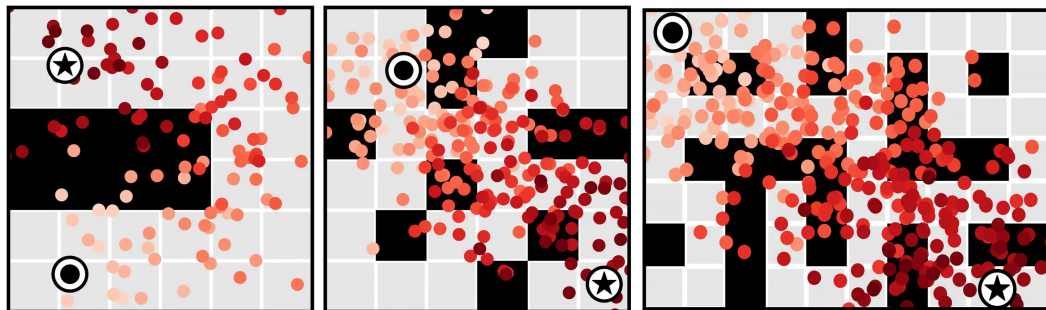


Motivation and Outline

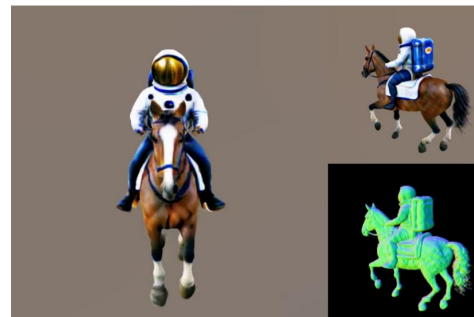
Various Applications with Diffusion Model

Motivation

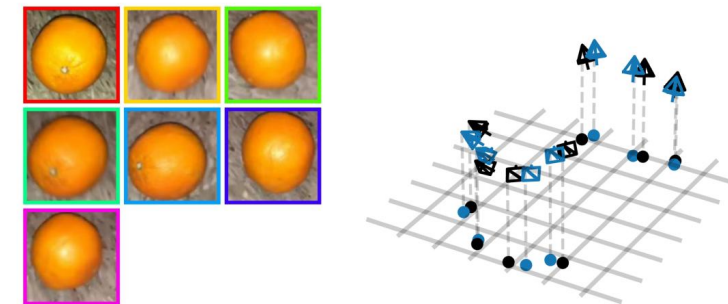
Motion Planning^[1]



3D Generation^[3]



Camera Pose Estimation^[5]



Reinforcement Learning^[2]

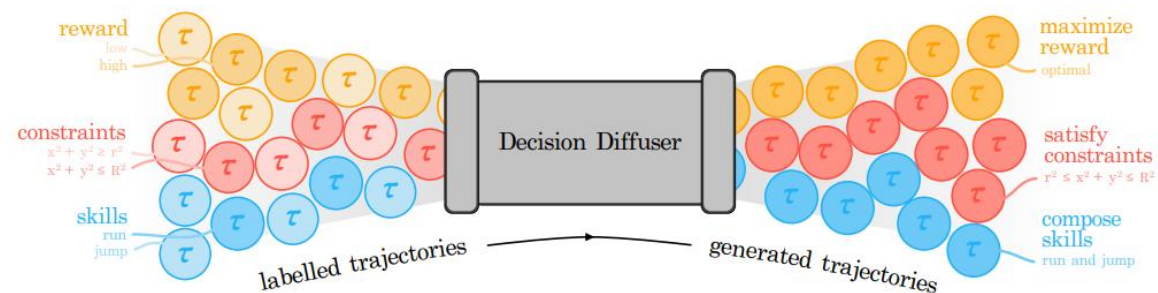
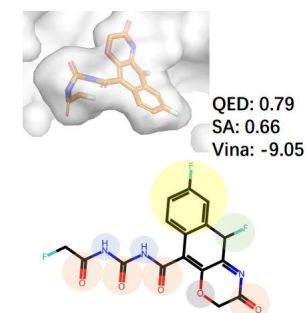


Image inpainting^[4]



Protein Structure^[6]



- [1] Janner M, Du Y, Tenenbaum J B, et al. Planning with diffusion for flexible behavior synthesis[J]. arXiv preprint arXiv:2205.09991, 2022.
- [2] Ajay A, Du Y, Gupta A, et al. Is conditional generative modeling all you need for decision-making?[J]. arXiv preprint arXiv:2211.15657, 2022.
- [3] Shi Y, Wang P, Ye J, et al. Mvdream: Multi-view diffusion for 3d generation[J]. arXiv preprint arXiv:2308.16512, 2023.
- [4] Lugmayr A, Danelljan M, Romero A, et al. Repaint: Inpainting using denoising diffusion probabilistic models[C]//Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2022: 11461-11471.
- [5] Zhang J Y, Lin A, Kumar M, et al. Cameras as rays: Pose estimation via ray diffusion[J]. arXiv preprint arXiv:2402.14817, 2024.
- [6] Huang Z, Yang L, Zhou X, et al. Interaction-based Retrieval-augmented Diffusion Models for Protein-specific 3D Molecule Generation[C]//Forty-first International Conference on Machine Learning.

Three Perspectives on Discrete Version Diffusion Model

Motivation

Soce Matching

Noise Conditioned Soce Matching Network (NCSN)^[1]

Hierarchical VAE

Denoising Diffusion Probabilistic Model (DDPM)^[2]

General Signal-Noise-Ratio

Variational Diffusion Model (VDM)^[3]

Motivation

Noise Conditioned Score Matching Network (NCSN)^[1]

- Revisiting Naive Score Matching
- Score Matching perspective of Diffusion Model
- Objective and Noise Schedule Improvement

Denoising Diffusion Probabilistic Model (DDPM)^[2]

Variational Diffusion Model (VDM)^[3]

- General Signal-Noise-Ratio perspective of Diffusion Model
- Neural Network Parameterized Signal-Noise-Ratio

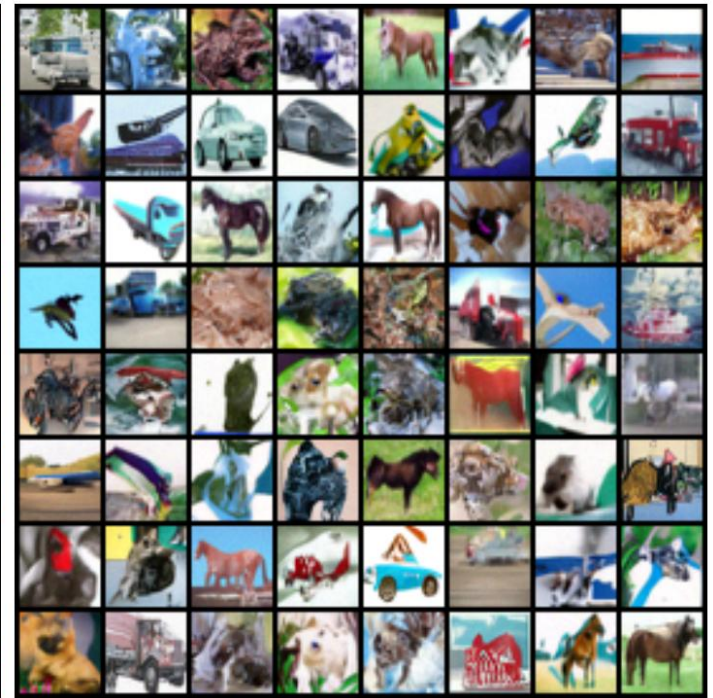
Hierarchical VAE perspective of DDPM^[2]

- Revisiting VAE and Evidence Lower Bound
- Hierarchical VAE perspective of Diffusion Model
- Reparameterization Triangle

Further Readings

Reference

Outline



Noise Conditioned Score Matching Network (NCSN)^[1]

- Revisiting Naive Score Matching
- Score Matching perspective of Diffusion Model
- Objective and Noise Schedule Improvement

Why score matching?

NCSN^[1] - Revisiting Naive Score Matching

Score matching can be used to approximate the **intractable distribution q** with **tractable distribution p** .

$$\min_p F(p, q) = \frac{1}{2} \mathbb{E}_q [\|\nabla_x \log p(x) - \nabla_x \log q(x)\|_2^2] \quad (\text{score matching})$$

Example of intractable distribution q

Energy-based Model	$q(x) = \frac{1}{Z} \exp[-E(x)]$	$\nabla_x \log q(x) = -\nabla_x E(x)$
Unnormalized Model	$q(x) = \frac{q^*(x)}{Z}$	$\nabla_x \log q(x) = \nabla_x \log q^*(x)$
Bayesian Posterior	$q(x D) = \frac{q(x,D)}{\int q(x,D)}$	$\nabla_x \log q(x) = \nabla_x \log q(x, D)$

Example of tractable distribution p

Gaussian Distribution	$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$	$\nabla_x \log p(x) = -\frac{(x-\mu)}{\sigma^2}$
Normalizing Flow	$p(x) = MLP(x), x \sim \mathcal{N}(0, 1)$	$\nabla_x \log p(x) = \frac{\nabla_x MLP(x)}{MLP(x)}$

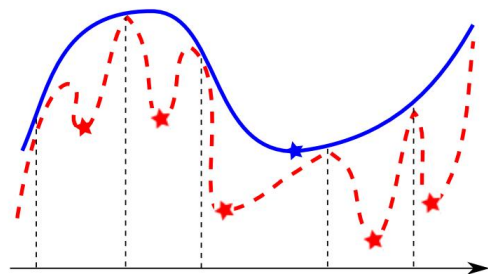
Why score matching?

NCSN^[1] - Revisiting Naive Score Matching

Score matching can be used to approximate the **intractable distribution q** with **tractable distribution p** .

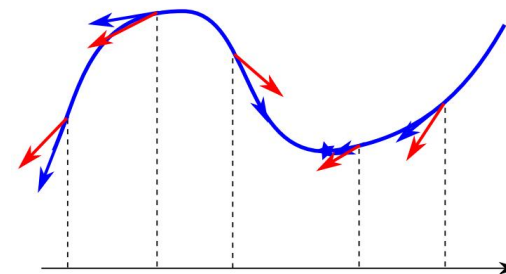
$$\min_p F(p, q) = \frac{1}{2} \mathbb{E}_q [\|\nabla_x \log p(x) - \nabla_x \log q(x)\|_2^2] \quad (\text{score matching})$$

Connection with VI (Variational Inference)



$$\min_p KL [q | p] = \mathbb{E}_q [\log q(x) - \log p(x)]$$

Variational Inference



Score Matching

Compared to VI, score matching leads to more stable score estimation and can get rid of the intractable normalization constant.

Score matching with data is not trivial

NCSN^[1] - Revisiting Naive Score Matching

Score matching can be used to approximate the intractable target data distribution q_{data} with **tractable neural network p_{θ}** .

$$\min_{\theta} F(p_{\theta}, q_{data}) = \frac{1}{2} \mathbb{E}_{q_{data}} [\|\nabla_x \log p_{\theta}(x) - \nabla_x \log q_{data}(x)\|_2^2]$$

the exact score function is unknown

samples are available in this data distribution



Two possible choices for learning the target data score function

- Integration by parts trick
- Denoising Score Matching (DSM)

One choice for sample according to the learned score function

- Langevin Dynamics

Integration by parts trick

NCSN^[1] - Revisiting Naive Score Matching

Integration by parts trick

$$\begin{aligned} & \frac{1}{2} \mathbb{E}_{\mathbf{q}_{data}} [\|\nabla_x \log \mathbf{p}_\theta(x) - \nabla_x \log \mathbf{q}_{data}(x)\|_2^2] \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{q}_{data}} \left[\frac{1}{2} (\nabla_x \log \mathbf{p}_\theta(x))^2 + \text{tr}(\nabla_x^2 \log \mathbf{p}_\theta(x)) \right] \end{aligned}$$

Proof.

$$\begin{aligned} & \min_{\theta} F(\mathbf{q}_{data}, \mathbf{p}_\theta) \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{q}_{data}} [\|\nabla_x \log \mathbf{q}_{data}(x) - \nabla_x \log \mathbf{p}_\theta(x)\|_2^2] \\ &= \frac{1}{2} \int \mathbf{q}_{data}(x) (\nabla_x \log \mathbf{q}_{data}(x) - \nabla_x \log \mathbf{p}_\theta(x))^2 dx \\ &= \frac{1}{2} \int \mathbf{q}_{data}(x) (\nabla_x \log \mathbf{q}_{data}(x))^2 dx + \frac{1}{2} \int \mathbf{q}_{data}(x) (\nabla_x \log \mathbf{p}_\theta(x))^2 dx - \int \mathbf{q}_{data}(x) \nabla_x \log \mathbf{q}_{data}(x) \nabla_x \log \mathbf{p}_\theta(x) dx \\ & \quad \underbrace{\hspace{10em}}_{\text{const}} \quad \underbrace{\mathbf{q}_{data}(x) \nabla_x \log \mathbf{q}_{data}(x) = \nabla_x \mathbf{q}_{data}(x)} \\ &= \frac{1}{2} \int \mathbf{q}_{data}(x) (\nabla_x \log \mathbf{p}_\theta(x))^2 dx - \frac{\mathbf{q}_{data}(x) \nabla_x \log \mathbf{p}_\theta(x) \Big|_{-\infty}^{+\infty}}{\mathbf{q}_{data}(x) \rightarrow 0 \text{ when } |x| \rightarrow \infty} + \int \mathbf{q}_{data}(x) \text{tr}(\nabla_x^2 \log \mathbf{p}_\theta(x)) dx \\ &= \frac{1}{2} \int \mathbf{q}_{data}(x) (\nabla_x \log \mathbf{p}_\theta(x))^2 dx + \int \mathbf{q}_{data}(x) \text{tr}(\nabla_x^2 \log \mathbf{p}_\theta(x)) dx \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{q}_{data}} \left[\frac{1}{2} (\nabla_x \log \mathbf{p}_\theta(x))^2 + \text{tr}(\nabla_x^2 \log \mathbf{p}_\theta(x)) \right] \end{aligned}$$

Denoising Score Matching (DSM)

NCSN^[1] - Revisiting Naive Score Matching

Denoising Score Matching

$$\frac{1}{2} \mathbb{E}_{\mathbf{q}_{data}} [\|\nabla_x \log \mathbf{p}_\theta(x) - \nabla_x \log \mathbf{q}_{data}(x)\|_2^2] \quad (\text{score matching})$$



$$\frac{1}{2} \mathbb{E}_{\mathbf{q}_\sigma(\tilde{x}|x)\mathbf{q}_{data}} [\|\nabla_x \log \mathbf{p}_\theta(\tilde{x}) - \nabla_{\tilde{x}} \log \mathbf{q}_\sigma(\tilde{x}|x)\|_2^2] \quad (\text{denoising score matching})$$

$\mathbf{q}_\sigma(\tilde{x}) = \int \mathbf{q}_\sigma(\tilde{x}|x)\mathbf{q}_{data}(x) dx$ represents the noised data

$\mathbf{q}_\sigma(\tilde{x}) \approx \mathbf{q}_{data}(x)$ when the noise is small enough

Denoising Score Matching for Gaussian Noise

Instead of approximate $\mathbf{q}_{data}(x)$ directly, noised $\mathbf{q}_\sigma(\tilde{x}|x)$ can be specifically defined (Gaussian Noise)

$$\mathbf{q}_\sigma(\tilde{x}|x) = \mathcal{N}(\tilde{x}; x, \sigma^2 \mathbf{I})$$
$$\nabla_{\tilde{x}} \log \mathbf{q}_\sigma(\tilde{x}|x) = \frac{\tilde{x} - x}{\sigma^2}$$

$$\frac{1}{2} \mathbb{E}_{\mathbf{q}_\sigma(\tilde{x}|x)\mathbf{q}_{data}} \left[\left\| \nabla_x \log \mathbf{p}_\theta(\tilde{x}) - \frac{\tilde{x} - x}{\sigma^2} \right\|_2^2 \right] \quad (\text{denoising score matching})$$

Recall the Langevin Dynamics

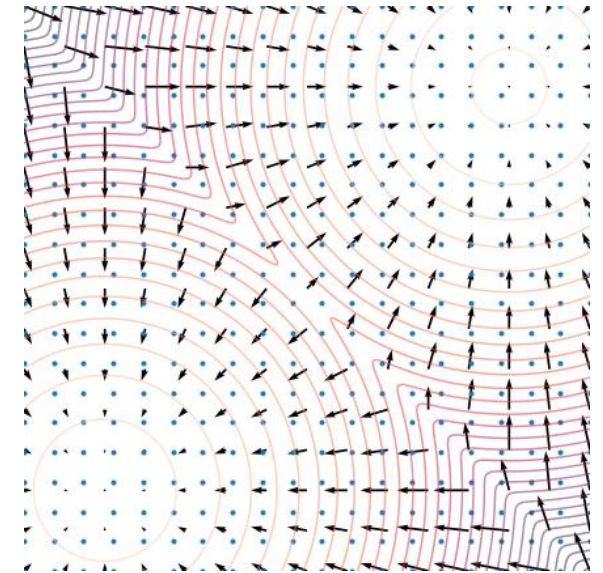
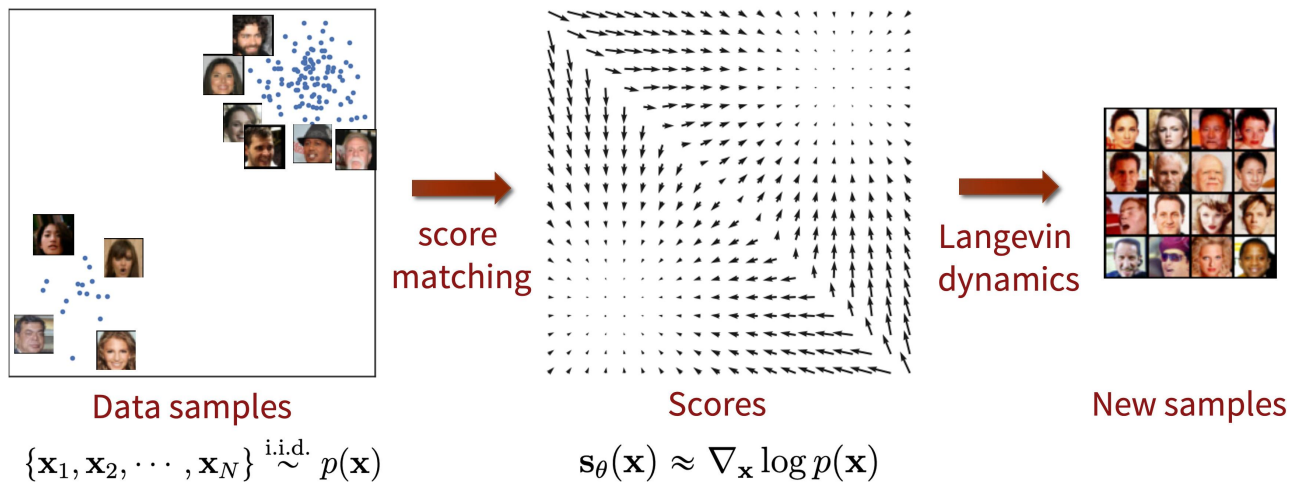
NCSN^[1] - Revisiting Naive Score Matching

Langevin Dynamics

$$x_t = x_{t-1} + \frac{\sigma}{2} \nabla_x \log p_\theta(x_{t-1}) + \sqrt{\sigma} \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

x_0 follows some prior distribution. The distribution of x_T converges to $p_\theta(x)$ when $T \rightarrow \infty$

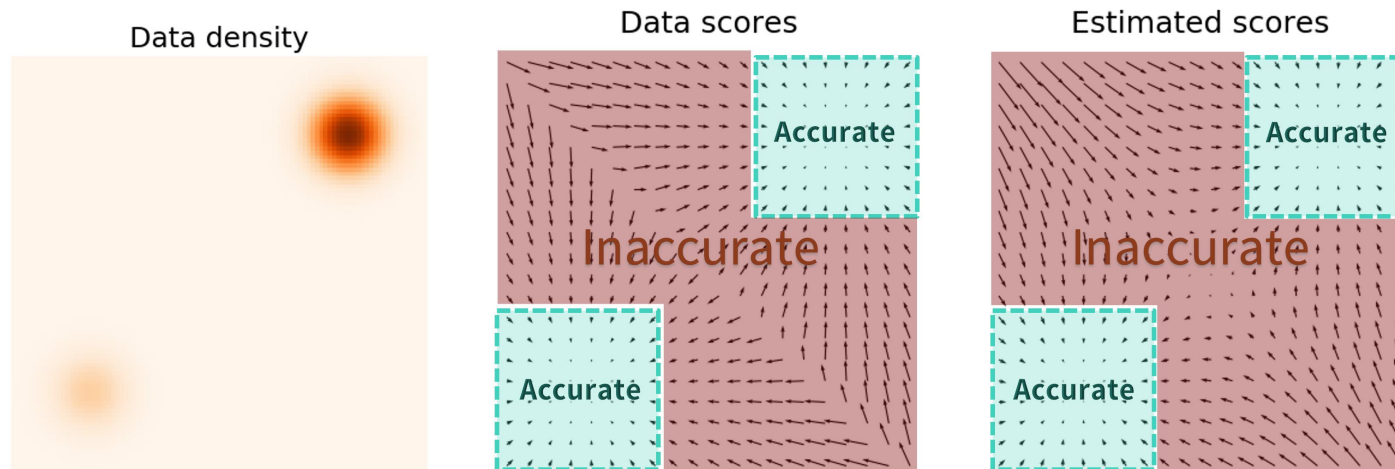
The full pipeline of score matching



Example of Langevin Dynamics

Challenges of score-based generative modeling

NCSN^[1] - Score Matching perspective of Diffusion Model



Example of inaccurate region of learned score

The manifold hypothesis

The manifold hypothesis states that data in the real world tend to concentrate on low dimensional manifolds embedded in a high dimensional space

Low data density regions

The scarcity of data in low density regions can cause difficulties for both score estimation with score matching and MCMC sampling with Langevin dynamics.

Learning score with multi-level noise

NCSN^[1] - Score Matching perspective of Diffusion Model

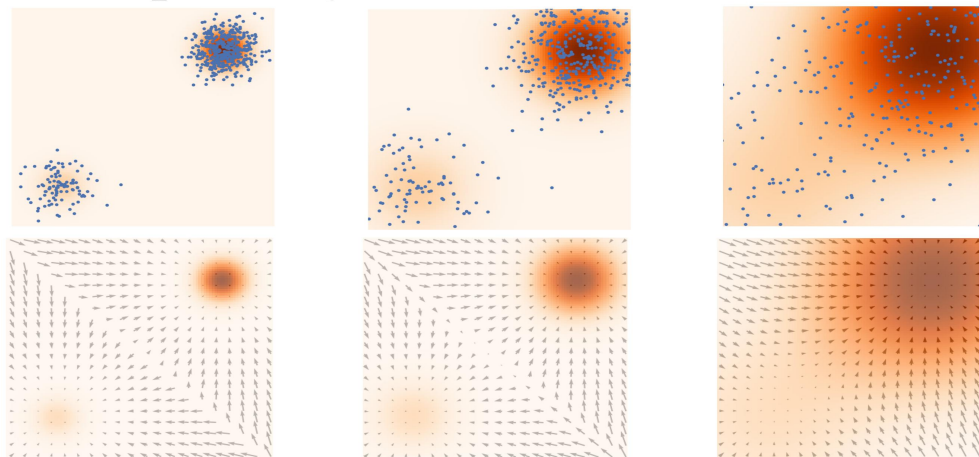
Noise Conditioned Score Matching (NCSN)

$$\mathcal{L}(\theta, \{\sigma_i\}_{i=1}^L) = \sum_{i=1}^L \lambda(\sigma_i) \ell(\theta, \sigma_i)$$

$$\begin{aligned} \ell(\theta, \sigma) &= \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \mathbf{q}_{data}(\mathbf{x})} \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbf{q}_\sigma(\tilde{\mathbf{x}}|\mathbf{x})} [\|\nabla_{\mathbf{x}} \log \mathbf{p}_\theta(\tilde{\mathbf{x}}) - \nabla_{\tilde{\mathbf{x}}} \log \mathbf{q}_\sigma(\tilde{\mathbf{x}}|\mathbf{x})\|_2^2] \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \mathbf{q}_{data}(\mathbf{x})} \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{x}, \sigma^2 \mathbf{I})} \left[\left\| \mathbf{s}_\theta(\tilde{\mathbf{x}}, \sigma) - \frac{\tilde{\mathbf{x}} - \mathbf{x}}{\sigma^2} \right\|_2^2 \right] \end{aligned}$$

Let $\{\sigma_i\}_{i=1}^L$ be a positive geometric sequence that satisfies $\epsilon = \frac{\sigma_1}{\sigma_2} = \dots = \frac{\sigma_{L-1}}{\sigma_L} > 1$

Example of Noise Conditioned Score Matching



$\sigma_3 < \sigma_2 < \sigma_1$

Different noise level of data density and samples.

Different noise level of learned score results.

Annealed Langevin Sampling

NCSN^[1] - Score Matching perspective of Diffusion Model

Annealed Langevin Sampling

Algorithm 1 Annealed Langevin dynamics.

Require: $\{\sigma_i\}_{i=1}^L, \epsilon, T$.

1: Initialize $\tilde{\mathbf{x}}_0$

2: **for** $i \leftarrow 1$ to L **do**

3: $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$ $\triangleright \alpha_i$ is the step size.

4: **for** $t \leftarrow 1$ to T **do**

5: Draw $\mathbf{z}_t \sim \mathcal{N}(0, I)$

6: $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_i}{2} \mathbf{s}_\theta(\tilde{\mathbf{x}}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$

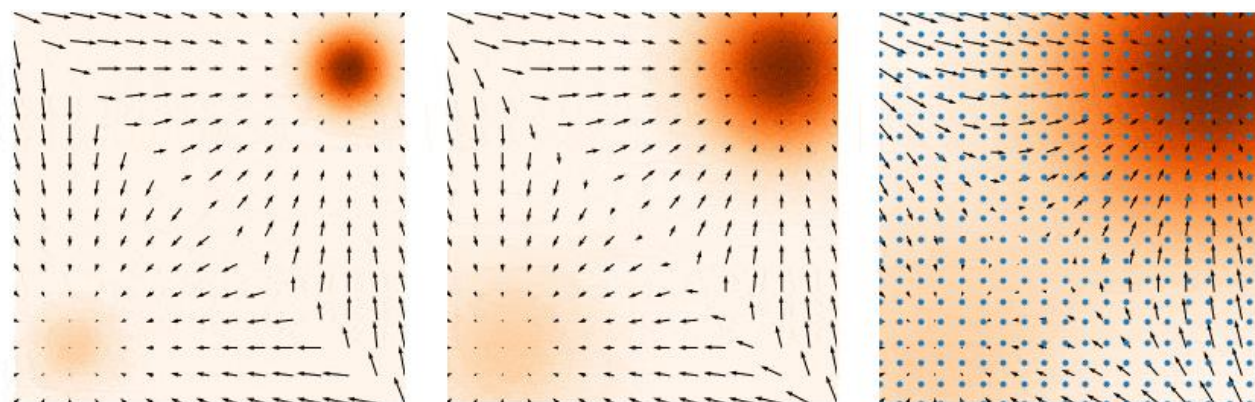
7: **end for**

8: $\tilde{\mathbf{x}}_0 \leftarrow \tilde{\mathbf{x}}_T$

9: **end for**

return $\tilde{\mathbf{x}}_T$

Example of Annealed Langevin Sampling



σ_3

σ_2

σ_1

Annealed Denoised Iteration

Langevin Dynamics Sampling

Objective Improvement

NCSN^[1] - Objective and Noise Schedule Improvement

Score-based Objective

$$\mathcal{L}(\theta, \{\sigma_i\}_{i=1}^L) = \sum_{i=1}^L \lambda(\sigma_i) \ell(\theta, \sigma_i)$$

$$\ell(\theta, \sigma) = \frac{1}{2} \mathbb{E}_{x \sim q_{data}(x)} \mathbb{E}_{\tilde{x} \sim \mathcal{N}(x, \sigma^2 \mathbf{I})} \left[\left\| \mathbf{s}_{\theta}(\tilde{x}, \sigma) - \frac{\tilde{x} - x}{\sigma^2} \right\|_2^2 \right]$$

This score neural network $\mathbf{s}_{\theta}(\tilde{x}, \sigma)$ has unnormalized output score

Residual-based Objective

$$\begin{aligned} \ell(\theta, \sigma) &= \frac{1}{2} \mathbb{E}_{x \sim q_{data}(x)} \mathbb{E}_{\tilde{x} \sim \mathcal{N}(x, \sigma^2 \mathbf{I})} \left[\left\| \mathbf{s}_{\theta}(\tilde{x}, \sigma) - \frac{\tilde{x} - x}{\sigma^2} \right\|_2^2 \right] \\ &= \frac{1}{2} \mathbb{E}_{x \sim q_{data}(x)} \mathbb{E}_{\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})} \left[\frac{1}{\sigma^2} \left\| \boldsymbol{\varepsilon}_{\theta}(\tilde{x}, \sigma) - \boldsymbol{\varepsilon} \right\|_2^2 \right] \end{aligned}$$

Proof.

$$\begin{aligned} q(\tilde{x}|x) &= \mathcal{N}(x, \sigma^2 \mathbf{I}) \\ \tilde{x} &= x + \sigma \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I}) \end{aligned}$$

■

This residual neural network $\boldsymbol{\varepsilon}_{\theta}(\tilde{x}, \sigma)$ has normalized output noise

Noise Schedule Improvement

NCSN^[1] - Objective and Noise Schedule Improvement

Variance-Exploding Noise Schedule

$$\mathcal{L}(\theta, \{\sigma_i\}_{i=1}^L) = \sum_{i=1}^L \lambda(\sigma_i) \ell(\theta, \sigma_i)$$

$$\ell(\theta, \sigma) = \frac{1}{2} \mathbb{E}_{x \sim q_{data}(x)} \mathbb{E}_{\tilde{x} \sim \mathcal{N}(x, \sigma^2 \mathbf{I})} \left[\left\| \mathbf{s}_\theta(\tilde{x}, \sigma) - \frac{\tilde{x} - x}{\sigma^2} \right\|_2^2 \right]$$

This can be viewed as the Variance-Exploding Noise Schedule as

$$\mathbf{x}_t = \mathbf{x}_0 + \epsilon^t \sigma_0 \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$$

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \mathbf{x}_0, \epsilon^t \sigma_0 \mathbf{I})$$

and this variance grows exponentially

Variance-Preserving Noise Schedule

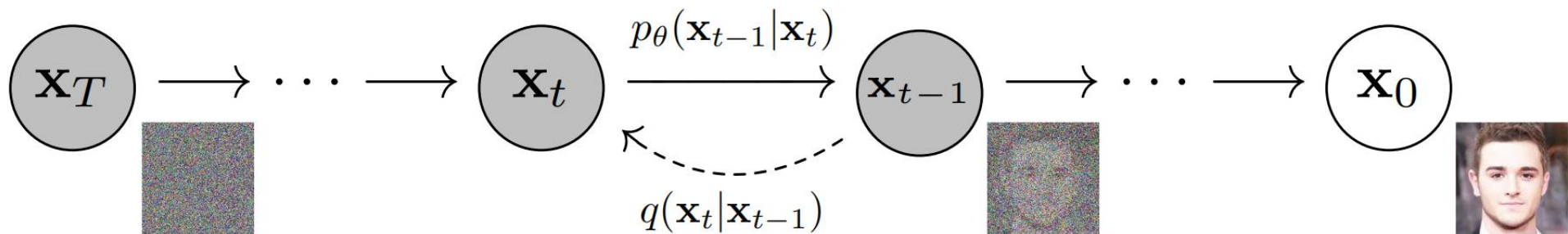
$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$$

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

This variance preserving noise schedule leads more stable learning process

Denosing Diffusion Probabilistic Model (DDPM)^[2]

$$\mathcal{L}(\theta) = \frac{1}{2} \mathbb{E}_{t \sim [0,1]} \mathbb{E}_{x \sim q_{data}(x)} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})} \left[\left\| \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}) - \boldsymbol{\epsilon} \right\|_2^2 \right]$$



Denoising Diffusion Probabilistic Model (DDPM)^[2]

VARIATIONAL BOUND FOR TRAINING

DDPM^[2]

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

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

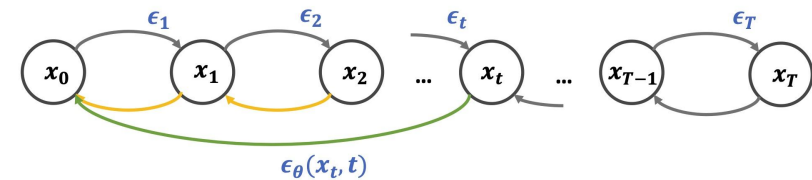
$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}) \quad \alpha_t := 1 - \beta_t \quad \bar{\alpha}_t := \prod_{s=1}^t \alpha_s$$

Efficient training:

$$\mathbb{E}_q \left[\underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{L_T} + \sum_{t>1} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}} - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) \right]$$

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

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



VARIATIONAL BOUND FOR TRAINING

DDPM^[2]

Gaussian stochastic process:

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

$$= \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[\frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon} \right) - \boldsymbol{\mu}_\theta(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}), t) \right\|^2 \right]$$

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \tilde{\boldsymbol{\mu}}_t \left(\mathbf{x}_t, \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t)) \right) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right)$$

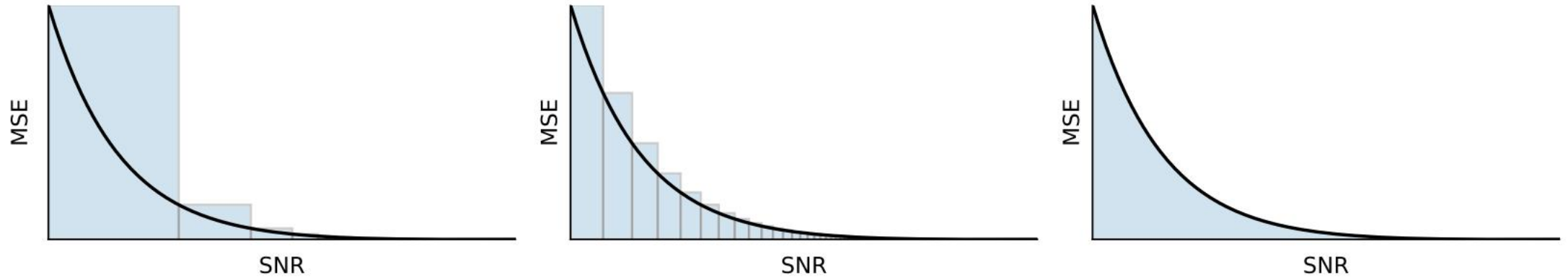
$$\mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[\frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1 - \bar{\alpha}_t)} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2 \right]$$

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
 $\nabla_{\boldsymbol{\theta}} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$
 - 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
-



Variational Diffusion Model (VDM)^[3]

- General Signal-Noise-Ratio perspective of Diffusion Model
- Neural Network Parameterized Signal-Noise-Ratio

Variance-Exploding Noise Schedule (NCSN^[1])

$$\mathbf{x}_t = \mathbf{x}_0 + \epsilon^t \sigma_0 \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I})$$

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \mathbf{x}_0, \epsilon^t \sigma_0 \mathbf{I})$$

Variance-Preserving Noise Schedule (DDPM^[2])

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

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

Variance-Nueral Noise Schedule (VDM^[3])

$$\mathbf{x}_t = NN(t) \mathbf{x}_0 + NN'(t) \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I})$$

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; NN(t) \mathbf{x}_0, [NN'(t)]^2 \mathbf{I})$$

This will be introduced in this section.

Neural Network Parameterized Signal-Noise-Ratio

VDM^[3] - Neural Network Parameterized Signal-Noise-Ratio

Variance-Exploding Noise Schedule (NCSN^[1])

$$\mathbf{x}_t = \mathbf{x}_0 + \epsilon^t \sigma_0 \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I})$$
$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \mathbf{x}_0, \epsilon^t \sigma_0 \mathbf{I})$$

Variance-Preserving Noise Schedule (DDPM^[2])

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I})$$
$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

Signal-Noise-Ratio must be non-negative and monotonically decreasing.

$$\text{SNR}(t) = \frac{\alpha_t^2}{\sigma_t^2} = \exp(-\gamma_\eta(t))$$

$$\sigma_t^2 = \text{sigmoid}(\gamma_\eta(t))$$

$$\alpha_t^2 = \text{sigmoid}(-\gamma_\eta(t))$$

Neural Network Parameterized Noise Schedule (VDM^[3])

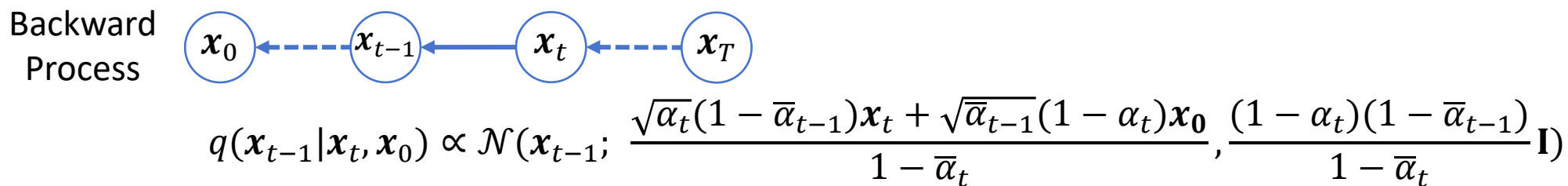
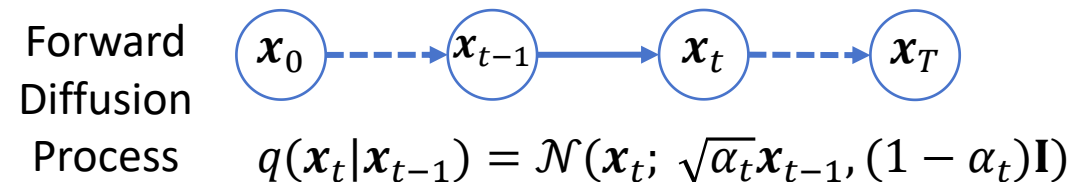
$$\mathbf{x}_t = \sqrt{\text{sigmoid}(-\gamma_\eta(t))} \mathbf{x}_0 + \sqrt{\text{sigmoid}(\gamma_\eta(t))} \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I})$$

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\text{sigmoid}(-\gamma_\eta(t))} \mathbf{x}_0, \text{sigmoid}(\gamma_\eta(t)) \mathbf{I})$$

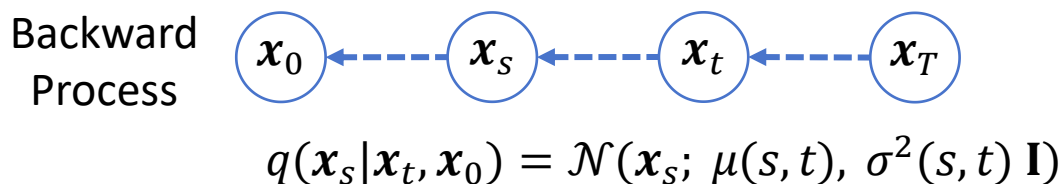
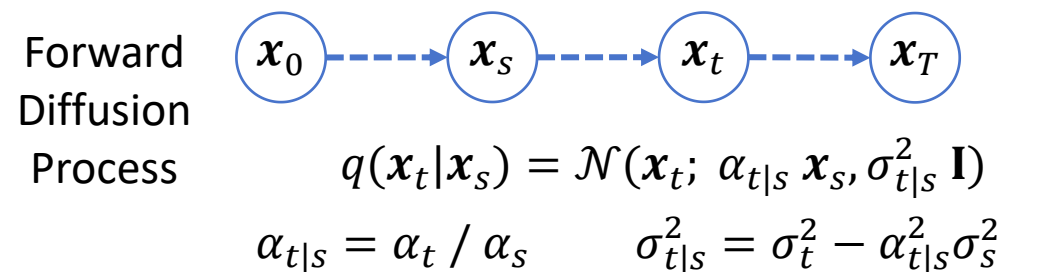
General Signal-Noise-Ratio

VDM^[3] - General Signal-Noise-Ratio perspective of Diffusion Model

Markov Process in DDPM^[2]



General Markov Process



ground truth	neural approximation
$\mu(s, t) = \frac{\alpha_{t s} \sigma_s^2}{\sigma_t^2} \mathbf{x}_t + \frac{\alpha_s \sigma_{t s}^2}{\sigma_t^2} \mathbf{x}_0$	$\mu_\theta(s, t) = \frac{\alpha_{t s} \sigma_s^2}{\sigma_t^2} \mathbf{x}_t + \frac{\alpha_s \sigma_{t s}^2}{\sigma_t^2} \mathbf{x}_\theta(\mathbf{x}_t, t)$
$\sigma^2(s, t) = \frac{\sigma_{t s}^2 \sigma_s^2}{\sigma_t^2}$	$\sigma^2(s, t) = \frac{\sigma_{t s}^2 \sigma_s^2}{\sigma_t^2}$

Objective function in VDM

VDM^[3] - General Signal-Noise-Ratio perspective of Diffusion Model

consistency/diffusion term

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

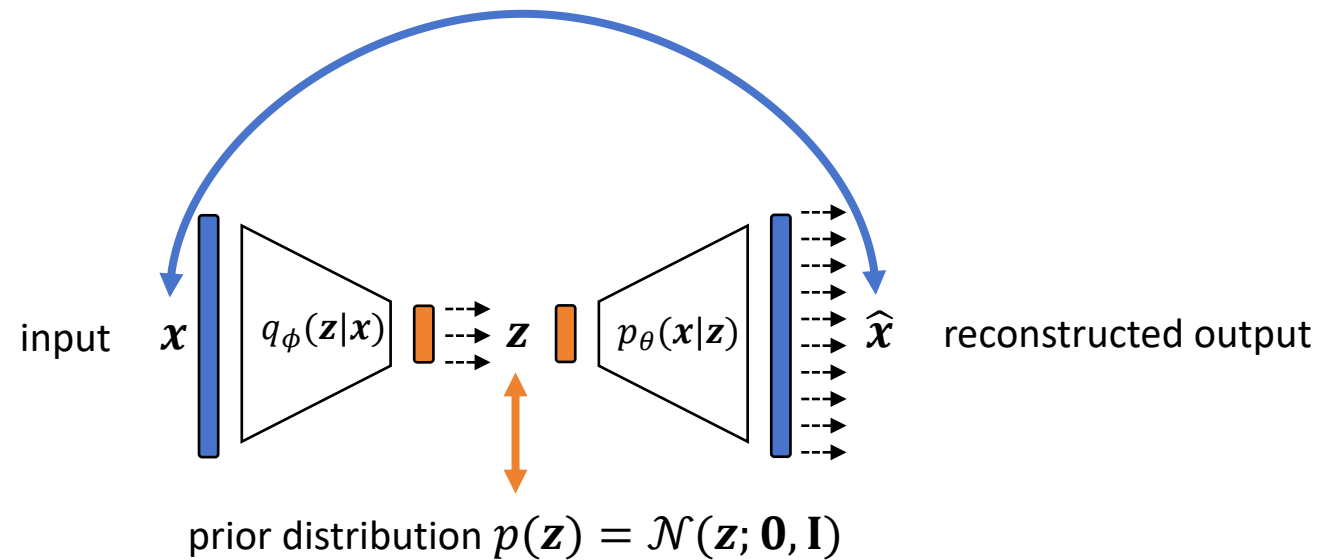
$$\begin{aligned} & D_{KL}(q(\mathbf{x}_s|\mathbf{x}_t, \mathbf{x}_0) | p_{\theta}(\mathbf{x}_s|\mathbf{x}_t)) \\ &= \frac{1}{2 \sigma^2(s, t)} \|\mu(s, t) - \mu_{\theta}(s, t)\|_2^2 \\ &= \frac{1}{2} \left(\frac{\alpha_s^2}{\sigma_s^2} - \frac{\alpha_t^2}{\sigma_t^2} \right) \|\mathbf{x}_0 - \mathbf{x}_{\theta}\|_2^2 \\ &= \frac{1}{2} (\text{SNR}(s) - \text{SNR}(t)) \|\mathbf{x}_0 - \mathbf{x}_{\theta}(\mathbf{x}_t, t)\|_2^2 \end{aligned}$$

$$\frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}), t \sim \mathcal{U}(0, T)} [(\text{SNR}(t-1) - \text{SNR}(t)) \|\mathbf{x}_0 - \mathbf{x}_{\theta}(\alpha_t \mathbf{x}_0 + \sigma_t \epsilon, t)\|_2^2]$$

reparameterization trick

$$\mathbf{x}_{\theta}(\mathbf{x}_t, t) = (\mathbf{x}_t - \sigma_t \epsilon_{\theta}(\mathbf{x}_t, t)) / \alpha_t$$

$$\frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}), t \sim \mathcal{U}(0, T)} [(\text{SNR}(t-1) - \text{SNR}(t)) (\sigma_t / \alpha_t) \|\epsilon - \epsilon_{\theta}(\alpha_t \mathbf{x}_0 + \sigma_t \epsilon, t)\|_2^2]$$



Hierarchical VAE perspective of DDPM^[2]

- Revisiting VAE and Evidence Lower BOund
- Hierarchical VAE perspective of Diffusion Model
- Reparameterization Triangle

Evidence Lower BOund of VI

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

Proof.

$$\begin{aligned} \log p(\mathbf{x}; \theta) &= \int q_\phi(\mathbf{z}) \log p(\mathbf{x}; \theta) dz \\ &= \int q_\phi(\mathbf{z}) (\log p(\mathbf{x}, \mathbf{z}; \theta) - \log p(\mathbf{z}|\mathbf{x}; \theta)) dz \\ &= \int q_\phi(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \theta)}{q_\phi(\mathbf{z})} dz - \int q_\phi(\mathbf{z}) \log \frac{p(\mathbf{z}|\mathbf{x}; \theta)}{q_\phi(\mathbf{z})} dz \\ &= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}; \theta)}{q_\phi(\mathbf{z})} \right] + D_{KL}(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x}; \theta)) \\ &\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}; \theta)}{q_\phi(\mathbf{z})} \right] \quad \blacksquare \end{aligned}$$

Evidence Lower BOund of VAE

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

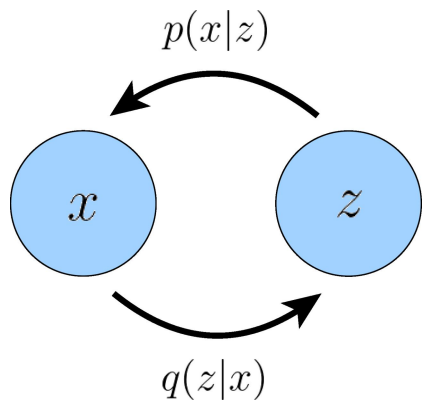
Proof.

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

Variational Auto Encoder

DDPM^[2] - Revisiting VAE and Evidence Lower Bound

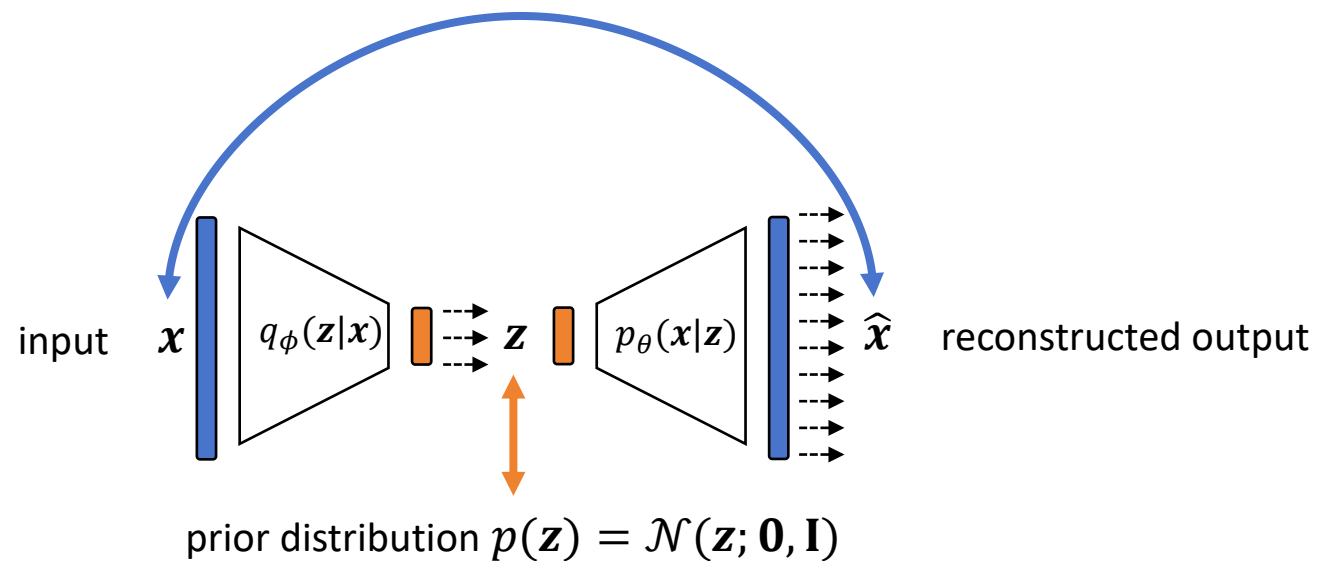
Graph Structure



ELBO Objective of VAE

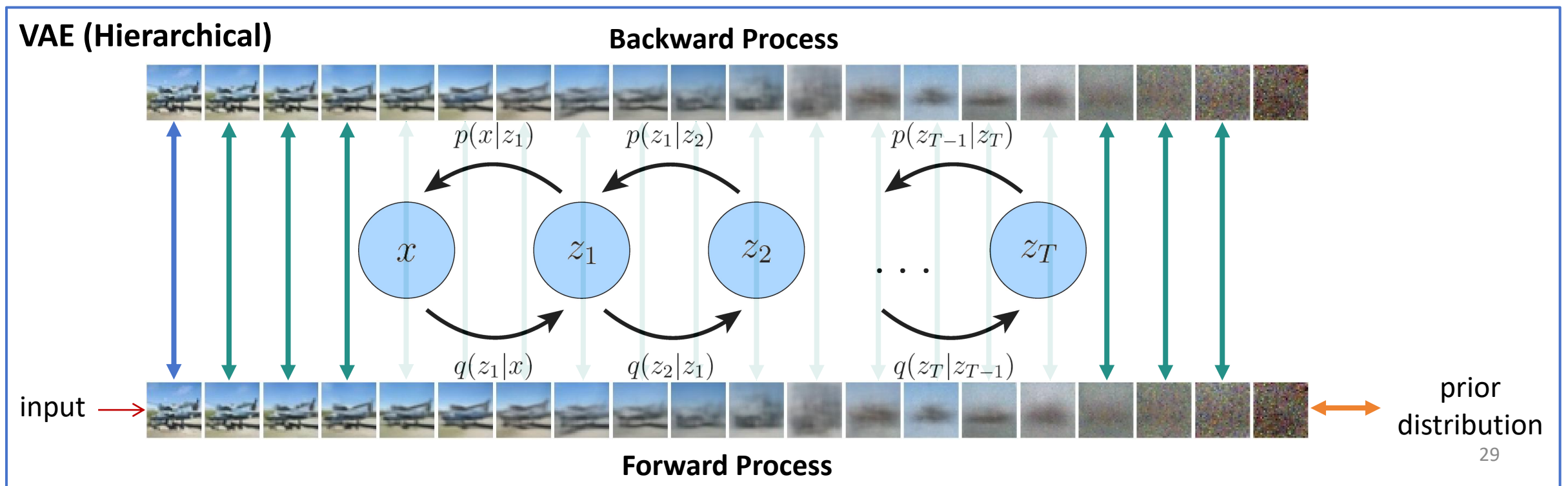
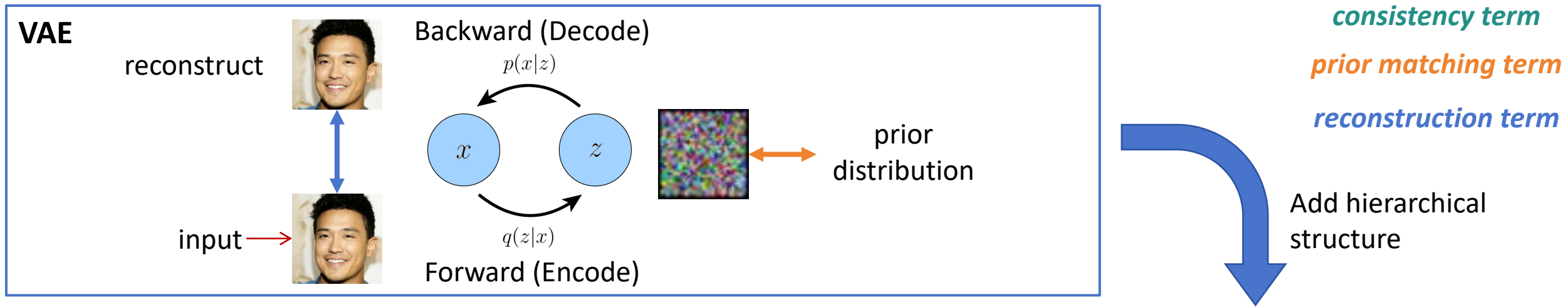
$$\mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \frac{p(x, z; \theta)}{q_\phi(z|x)} \right] = \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x) \parallel p(z))$$

reconstruction term *prior matching term*



From VAE to Hierarchical VAE

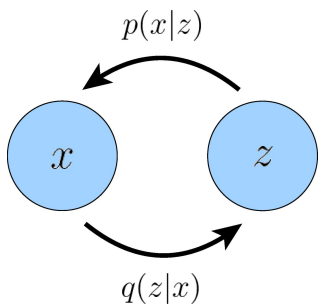
DDPM^[2] - Hierarchical VAE perspective of Diffusion Model



From VAE to Hierarchical VAE

DDPM^[2] - Hierarchical VAE perspective of Diffusion Model

Graph Structure



ELBO Objective of VAE

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

reconstruction term

prior matching term

ELBO Objective of Hierarchical VAE

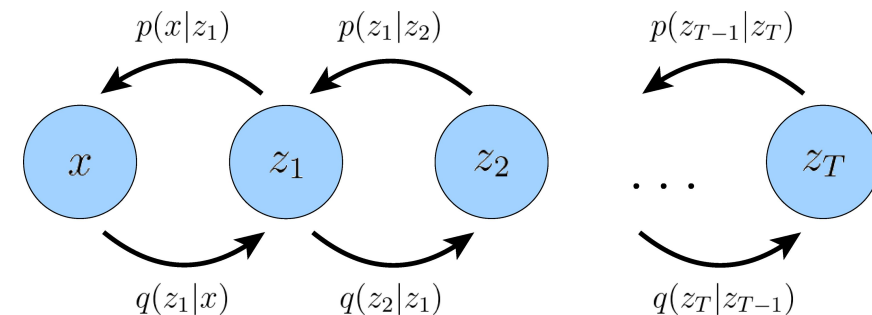
$$\begin{aligned} & \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}_{1:T}; \theta)}{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})} \left[\log p_{\theta}(\mathbf{x}|\mathbf{z}_1) + \log \frac{p(\mathbf{z}_T)}{q_{\phi}(\mathbf{z}_T|\mathbf{z}_{T-1})} + \log \frac{\prod_{t=1}^{T-1} p_{\theta}(\mathbf{z}_t|\mathbf{z}_{t+1})}{\prod_{t=1}^{T-1} q_{\phi}(\mathbf{z}_t|\mathbf{z}_{t-1})} \right] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}_1|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z}_1)] - \mathbb{E}_{q_{\phi}(\mathbf{z}_{T-1}|\mathbf{x})} [D_{KL}(q_{\phi}(\mathbf{z}_T|\mathbf{z}_{T-1}) \parallel p(\mathbf{z}_T))] - \sum_{t=1}^{T-1} \mathbb{E}_{q_{\phi}(\mathbf{z}_{t-1}, \mathbf{z}_{t+1}|\mathbf{x})} [D_{KL}(q_{\phi}(\mathbf{z}_t|\mathbf{z}_{t-1}) \parallel p_{\theta}(\mathbf{z}_t|\mathbf{z}_{t+1}))] \end{aligned}$$

reconstruction term

prior matching term

consistency/diffusion term

Graph Structure



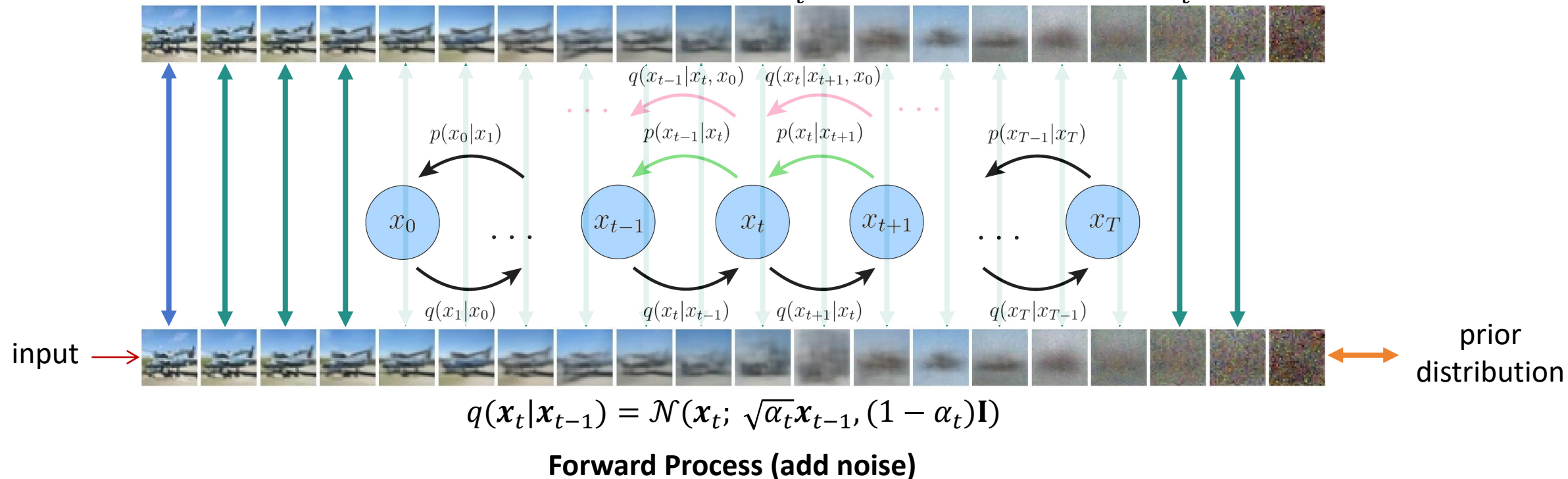
From Hierarchical VAE to Diffusion Model

DDPM^[2] - Hierarchical VAE perspective of Diffusion Model

Backward Process (Denoising)

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \propto \mathcal{N}(\mathbf{x}_{t-1}; \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}, \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I})$$

consistency term
prior matching term
reconstruction term



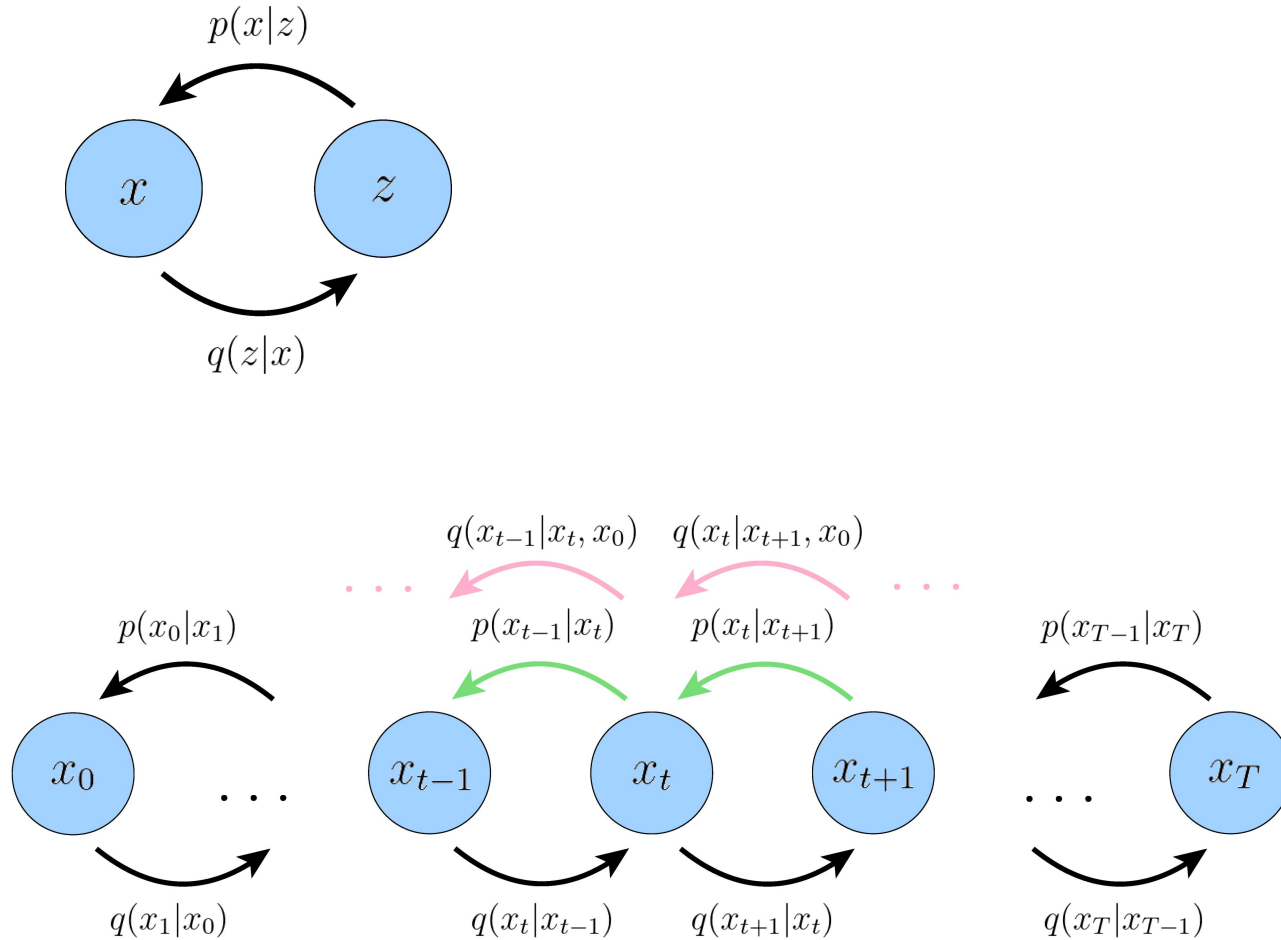
Benefits.

The forward process is “deterministic” without neural network. Instead of minimizing the $D_{KL}(q_\phi(\mathbf{z}_t|\mathbf{z}_{t-1}) | p_\theta(\mathbf{z}_t|\mathbf{z}_{t+1}))$, diffusion model minimizes $D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) | p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))$ which have more supervised signal (more stable to train).

Let us derive the objective function

DDPM^[2] - Hierarchical VAE perspective of Diffusion Model

Graph Structure



ELBO Objective of VAE

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



$z \rightarrow x_{1:T}$

1. Replace single variable distribution with the hierarchical joint distribution

$$\mathbb{E}_{q(x_{1:T}|x_0)} \left[\log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right]$$

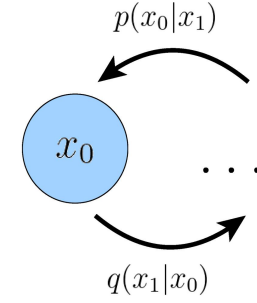
2. Decompose the joint distribution

Let us derive the objective function

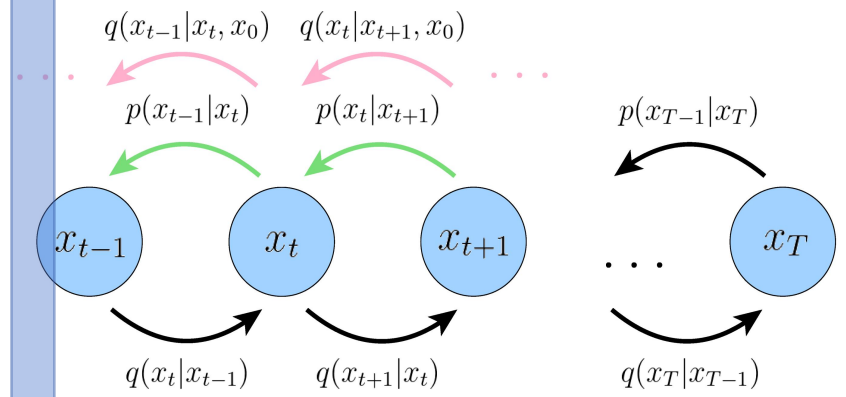
DDPM^[2] - Hierarchical VAE perspective of Diffusion Model

Joint Distribution Decomposite

$$\begin{aligned}
 & \mathbb{E}_{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] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} \frac{\prod_{t=2}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} \frac{\prod_{t=2}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=2}^T \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{q(\mathbf{x}_T|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{\prod_{t=2}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=2}^T q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \left[\log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] - D_{KL}(q(\mathbf{x}_T|\mathbf{x}_0) | p(\mathbf{x}_T)) - \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) | p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]
 \end{aligned}$$



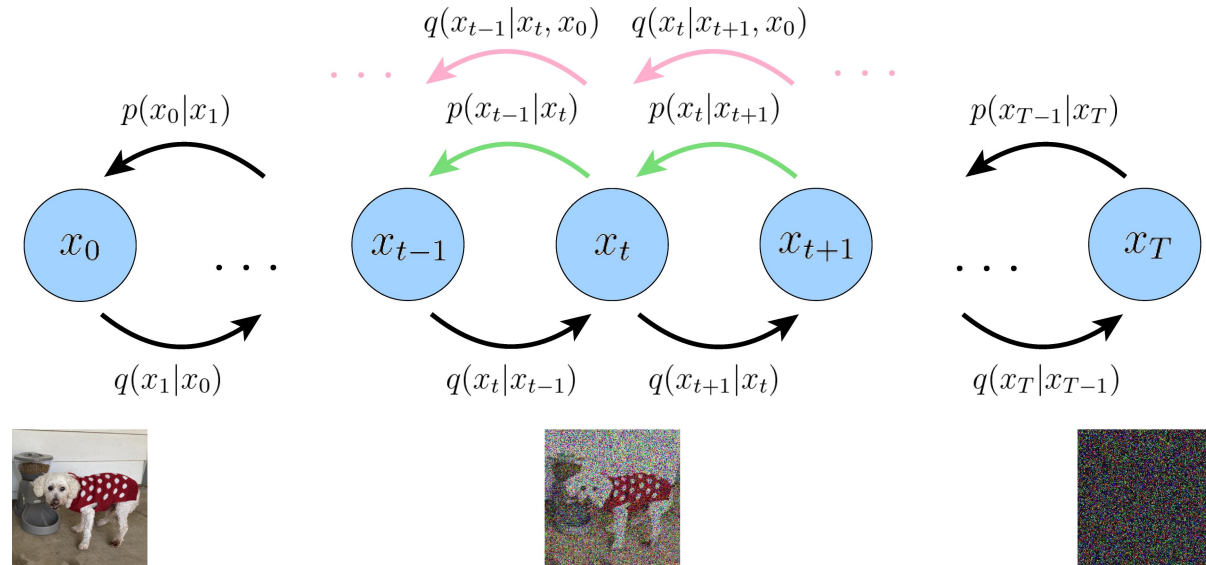
Graph Structure



Let us derive the objective function

DDPM^[2] - Hierarchical VAE perspective of Diffusion Model

Graph Structure



$$\mathbb{E}_{q(x_1|x_0)}[\log p_\theta(x_0|x_1)] - D_{KL}(q(x_T|x_0) | p(x_T))$$

Can be ignored

reconstruction term *prior matching term*

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

consistency/diffusion term

3. Combine with the deterministic noise schedule.

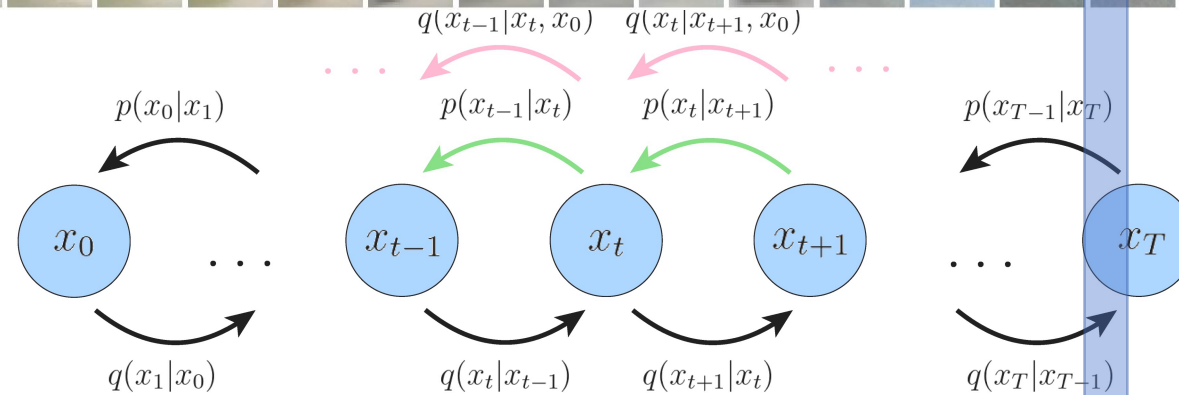
Let us derive the objective function

DDPM^[2] - Hierarchical VAE perspective of Diffusion Model

Deterministic Noise Schedule

Backward Process (Denoising)

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



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

Forward Process (add noise)

Let us derive the objective function

DDPM^[2] - Hierarchical VAE perspective of Diffusion Model

ground-truth

$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}$$

neural approximation

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

consistency/diffusion term

$$\begin{aligned} & \arg \min_{\theta} \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))] \\ &= \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} \left[\frac{1}{2\sigma_q^2(t)} \left[\|\boldsymbol{\mu}_q - \boldsymbol{\mu}_\theta\|_2^2 \right] \right] \\ &= \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} \left[\frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)^2} \left[\|\mathbf{x}_0 - \hat{\mathbf{x}}_\theta(\mathbf{x}_t, t)\|_2^2 \right] \right] \end{aligned}$$

4. Rewrite the objective with ε - prediction.

reparameterization trick

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_0 \\ q(\mathbf{x}_t | \mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \\ \mathbf{x}_0 &= \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_0}{\sqrt{\bar{\alpha}_t}} \end{aligned}$$

ground-truth

$$\begin{aligned} \boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t} \\ &= \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \epsilon_0 \end{aligned}$$

neural approximation

$$\begin{aligned} \boldsymbol{\mu}_\theta(\mathbf{x}_t, t) &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{\mathbf{x}}_\theta(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t} \\ &= \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\epsilon}_\theta(\mathbf{x}_t, t) \end{aligned}$$

Let us derive the objective function

DDPM^[2] - Hierarchical VAE perspective of Diffusion Model

ground-truth

$$\begin{aligned} \boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t} \\ &= \frac{1}{\sqrt{\alpha_t}}\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}}\boldsymbol{\epsilon}_0 \end{aligned}$$

neural approximation

$$\begin{aligned} \boldsymbol{\mu}_\theta(\mathbf{x}_t, t) &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{\mathbf{x}}_\theta(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t} \\ &= \frac{1}{\sqrt{\alpha_t}}\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}}\hat{\boldsymbol{\epsilon}}_\theta(\mathbf{x}_t, t) \end{aligned}$$

See proof in [4]

consistency/diffusion term

$$\begin{aligned} &\arg \min_{\theta} \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) | p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))] \\ &= \arg \min_{\theta} \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\frac{1}{2\sigma_q^2(t)} \left\| \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}} (\boldsymbol{\epsilon}_0 - \hat{\boldsymbol{\epsilon}}_\theta(\mathbf{x}_t, t)) \right\|_2^2 \right] \\ &= \arg \min_{\theta} \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)\alpha_t} \|\boldsymbol{\epsilon}_0 - \hat{\boldsymbol{\epsilon}}_\theta(\mathbf{x}_t, t)\|_2^2 \right] \end{aligned}$$

Can be ignored

Denoising Diffusion Probabilistic Model (DDPM)^[2]

$$\mathcal{L}(\theta) = \frac{1}{2} \mathbb{E}_{t \sim [0,1]} \mathbb{E}_{\mathbf{x} \sim q_{data}(\mathbf{x})} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})} \left[\left\| \boldsymbol{\epsilon}_\theta(\sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}) - \boldsymbol{\epsilon} \right\|_2^2 \right]$$

Training and sampling in DDPM

DDPM^[2] - Hierarchical VAE perspective of Diffusion Model

Algorithm Definition

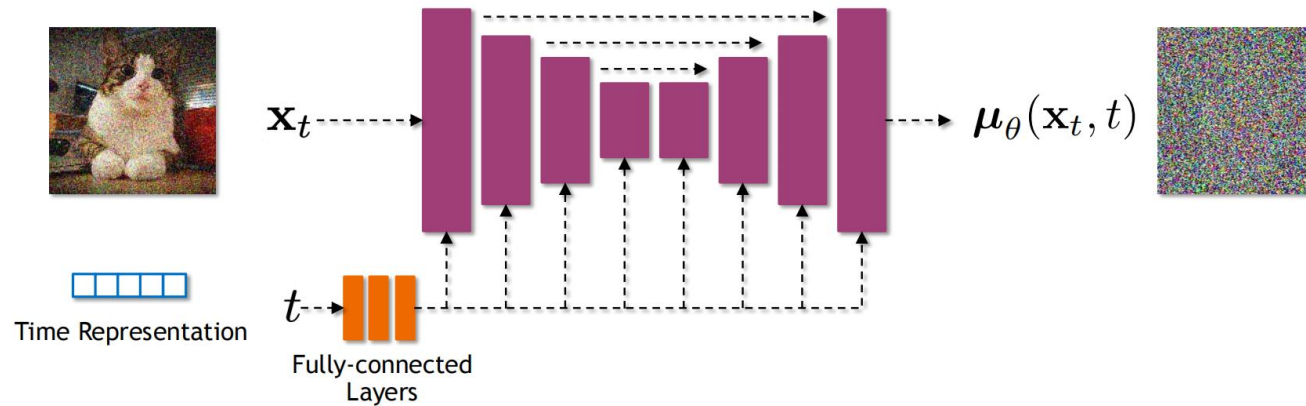
Algorithm 1 Training

- 1: **repeat**
- 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on $\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon; t)\|^2$
- 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: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
- 5: **end for**
- 6: **return** \mathbf{x}_0

Diffusion models often use U-Net architectures with ResNet blocks, skip connections, and self-attention layers.



Time representation: sinusoidal positional embeddings or random Fourier features.

Time features are fed to residual blocks using either simple spatial addition or adaptive group normalization layers.

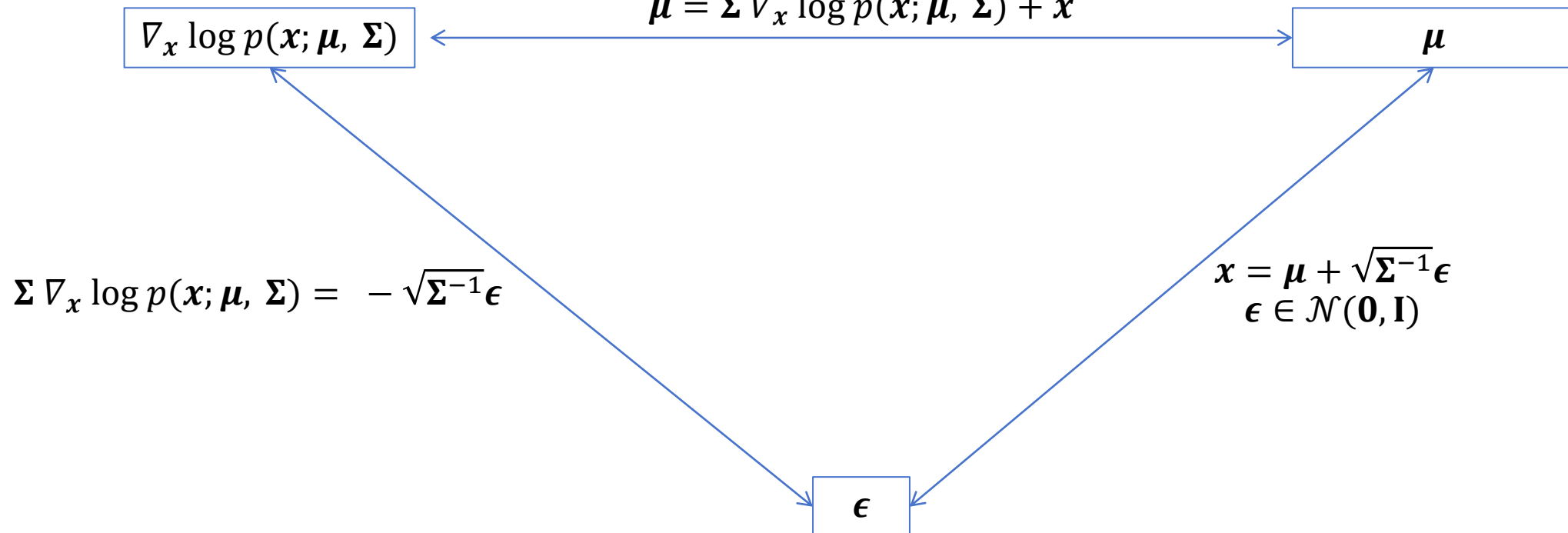
Nueral Network Structure

Reparameterization Triangle

DDPM^[2] - Reparameterization Triangle

$$p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^N \det \boldsymbol{\Sigma}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \nabla_{\mathbf{x}} \left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) = -\boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})$$
$$\boldsymbol{\mu} = \boldsymbol{\Sigma} \nabla_{\mathbf{x}} \log p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \mathbf{x}$$



Generative Modeling by Estimating Gradients of the Data Distribution

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

Variational Diffusion Models 2.0

Understanding Diffusion Model Objectives as the ELBO with Simple Data Augmentation

Durk Kingma Ruiqi Gao



<https://slideslive.com/39006359/>

Understanding Diffusion Models: A Unified Perspective

Calvin Luo
Google Research, Brain Team
calvinluo@google.com

August 26, 2022

<https://arxiv.org/abs/2208.11970>

- [1] Song Y, Ermon S. Generative modeling by estimating gradients of the data distribution[J]. Advances in neural information processing systems, 2019, 32.
- [2] Ho J, Jain A, Abbeel P. Denoising diffusion probabilistic models[J]. Advances in neural information processing systems, 2020, 33: 6840-6851.
- [3] Kingma D, Salimans T, Poole B, et al. Variational diffusion models[J]. Advances in neural information processing systems, 2021, 34: 21696-21707.
- [4] Luo C. Understanding diffusion models: A unified perspective[J]. arXiv preprint arXiv:2208.11970, 2022.