

Diffusion Model Journal Club

Kingma & Gao: Understanding Diffusion Objectives as the ELBO with Simple Data Augmentation

Johannes Zenn

Tübingen AI Center, University of Tübingen, IMPRS-IS

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



Tübingen AI Center

imprs-is

Creating Data From Noise

Last Week: Modeling the Reverse Stochastic Differential Equation



- forward SDE (from data to noise)

$$dz = f(z, t)dt + g(t)d\mathbf{w} \quad (1)$$

- the forward SDE can be reversed (Anderson, 1982) resulting in the reverse SDE

$$dz = [f(z, t) - g(t)^2 \nabla_x \log q(x, t)]dt + g(t)d\mathbf{w} \quad (2)$$

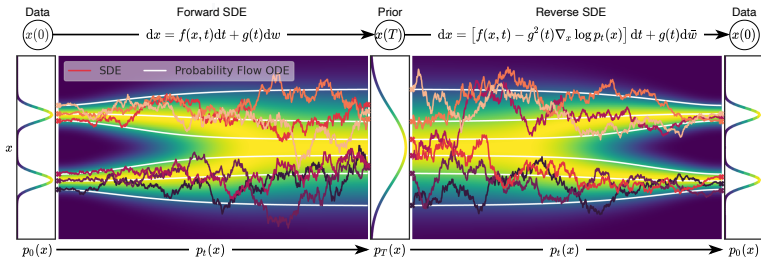


Figure taken from Song et al. (2020)

Two Stochastic Differential Equations

Last Week: Variance-Preserving and Variance-Exploding Forward SDEs



Variance-Preserving (VP) SDE

$$dz = \underbrace{-\frac{1}{2} \left(\frac{d}{dt} \log(1 + e^{-\lambda_t}) \right)}_{\text{drift}} z dt + \underbrace{\sqrt{\frac{d}{dt} \log(1 + e^{-\lambda_t})}}_{\text{diffusion}} dw \quad (3)$$

Variance-Exploding (VE) SDE

$$dz = 0 dt + \sqrt{\frac{d}{dt} \log(1 + e^{-\lambda_t})} dw \quad (4)$$

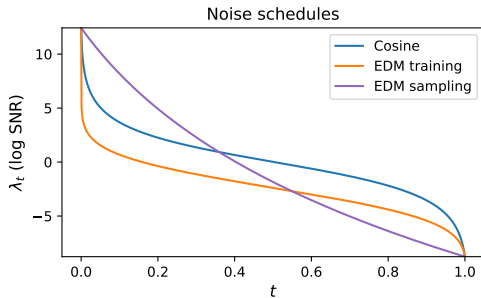


Figure taken from Kingma & Gao (2024)

Denoising Score Matching

Following Hyvärinen & Dayan (2005) and Vincent (2011)



- ▶ for generative modeling, we need to approximate the score (of the data) $\nabla_{\mathbf{x}} \log q(\mathbf{x}, t)$ in the reverse SDE

$$d\mathbf{z} = [\mathbf{f}(\mathbf{z}, t) - g(t)^2 \nabla_{\mathbf{x}} \log q(\mathbf{x}, t)] dt + g(t) d\mathbf{w} \quad (5)$$

- ▶ following TODO: Hyvärinen 2005, Vincent 2011 we want to minimize the following objective

$$\mathbb{E}_{q(\mathbf{x})} [\|\nabla_{\mathbf{x}} \log q(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x})\|_2^2] \quad (6)$$

- ▶ we have

$$\begin{aligned} & \mathbb{E}_{q(\mathbf{x})} [\|\nabla_{\mathbf{x}} \log q(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x})\|_2^2] \\ &= \mathbb{E}_{q(\mathbf{x})} \left[\left\| \frac{\partial \log q(\mathbf{x})}{\partial \mathbf{x}} \right\|^2 \right] - \mathbb{E}_{q(\mathbf{x})} \left[\left\langle \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}), \frac{\partial \log q(\mathbf{x})}{\partial \mathbf{x}} \right\rangle \right] \\ & \quad + \mathbb{E}_{q(\mathbf{x})} [\|\nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x})\|^2] \end{aligned} \quad (7)$$

Denoising Score Matching

Following Hyvärinen & Dayan (2005) and Vincent (2011)



- and, after some manipulations,

$$\mathbb{E}_{q(\mathbf{x})} \left[\left\langle \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}), \frac{\partial \log q(\mathbf{x})}{\partial \mathbf{x}} \right\rangle \right] = \mathbb{E}_{q(\mathbf{x}, \tilde{\mathbf{x}})} \left[\left\langle \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}), \frac{\partial \log q(\mathbf{x} | \tilde{\mathbf{x}})}{\partial \mathbf{x}} \right\rangle \right] \quad (8)$$

where $q(\mathbf{x}, \tilde{\mathbf{x}}) = q(\mathbf{x} | \tilde{\mathbf{x}})q_0(\tilde{\mathbf{x}})$ and $q(\mathbf{x} | \tilde{\mathbf{x}})$ is a noising distribution

- therefore,

$$\begin{aligned} \mathbb{E}_{q(\mathbf{x})} [\|\nabla_{\mathbf{x}} \log q(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x})\|_2^2] \\ = \mathbb{E}_{q(\mathbf{x}, \tilde{\mathbf{x}})} [\|\nabla_{\mathbf{x}} \log q(\mathbf{x} | \tilde{\mathbf{x}}) - \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x})\|_2^2] \end{aligned} \quad (9)$$

- we parameterize the score of the model by

$$\begin{aligned} \mathbb{E}_{q(\mathbf{x}, \tilde{\mathbf{x}})} [\|\nabla_{\mathbf{x}} \log q(\mathbf{x} | \tilde{\mathbf{x}}) - \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x})\|_2^2] \\ = \mathbb{E}_{q(\mathbf{x}, \tilde{\mathbf{x}})} [\|\nabla_{\mathbf{x}} \log q(\mathbf{x} | \tilde{\mathbf{x}}) - s_{\theta}(\mathbf{x})\|_2^2] \end{aligned} \quad (10)$$

Denoising Score Matching

Following Hyvärinen & Dayan (2005) and Vincent (2011)



- ▶ now, introduce multiple noise levels λ

$$\mathbf{z}_\lambda = \alpha_\lambda \mathbf{x} + \sigma_\lambda \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- ▶ we recover the denoising score matching (DSM) objective

$$\mathcal{L}_{\text{DSM}}(\mathbf{x}) = \mathbb{E}_{t, \boldsymbol{\epsilon}} \left[\left\| s_\theta(\mathbf{z}_\lambda; \lambda) - \frac{\mathbf{z}_\lambda - \alpha_\lambda \mathbf{x}}{\sigma_\lambda^2} \right\|_2^2 \right] \quad (11)$$

- ▶ where we used the gradient of $\log q(\mathbf{z}_\lambda | \mathbf{x})$

$$\nabla_{\mathbf{z}_\lambda} \log q(\mathbf{z}_\lambda | \mathbf{x}) = -\frac{\mathbf{z}_\lambda - \alpha_\lambda \mathbf{x}}{\sigma_\lambda^2} \quad (12)$$

Parameterization of the Score Network

Last Week: Following Karras et al. (2022)



- now, we can approximate the score (of the data) $\nabla_{\mathbf{x}} \log q(\mathbf{x}, t)$ in the reverse SDE

$$d\mathbf{z} = [\mathbf{f}(\mathbf{z}, t) - g(t)^2 \nabla_{\mathbf{x}} \log q(\mathbf{x}, t)] dt + g(t) d\mathbf{w} \quad (13)$$

by $\nabla_{\mathbf{x}} \log q(\mathbf{x}, t) \approx \mathbf{s}_{\theta}(\mathbf{z}; \lambda)$

- the score network can be parameterized in various ways, based on the relationships between \mathbf{z}_{λ} , \mathbf{x} , and ϵ

$$\mathbf{z}_{\lambda} = \alpha_{\lambda} \mathbf{x} + \sigma_{\lambda} \epsilon \quad (14)$$

$$\mathbf{x} = \alpha_{\lambda}^{-1} (\mathbf{z}_{\lambda} - \sigma_{\lambda} \epsilon) \quad (15)$$

$$\epsilon = \sigma_{\lambda}^{-1} (\mathbf{z}_{\lambda} - \alpha_{\lambda} \mathbf{x}) \quad (16)$$

$$\mathbf{s}_{\theta}(\mathbf{z}; \lambda) = -\nabla_{\mathbf{z}} E_{\theta}(\mathbf{z}, \lambda) \quad (17)$$

$$\mathbf{s}_{\theta}(\mathbf{z}; \lambda) = -\hat{\epsilon}_{\theta}(\mathbf{z}; \lambda) / \sigma_{\lambda} \quad (18)$$

$$\mathbf{s}_{\theta}(\mathbf{z}; \lambda) = -(\mathbf{z} - \alpha_{\lambda} \hat{\mathbf{x}}_{\theta}(\mathbf{z}; \lambda)) / \sigma_{\lambda}^2 \quad (19)$$

F -Prediction (Karras et al., 2022)

Last Week: Variance Exploding (VE) SDE Special Case



- F -prediction model: In Karras et al., the F -prediction model is used to parameterize the score network under the variance-exploding (VE) SDE
- F -prediction Formula for VE SDE

$$\mathbf{x} = \frac{\tilde{\sigma}_{\text{data}}^2}{e^{-\lambda} + \tilde{\sigma}_{\text{data}}^2} \mathbf{z}_{\lambda} + \frac{e^{-\lambda/2} \tilde{\sigma}_{\text{data}}}{\sqrt{e^{-\lambda} + \tilde{\sigma}_{\text{data}}^2}} \mathbf{F} \quad (20)$$

$$\mathbf{F} = \frac{\sqrt{e^{-\lambda} + \tilde{\sigma}_{\text{data}}^2}}{e^{-\lambda/2} \tilde{\sigma}_{\text{data}}} \mathbf{x} - \frac{\tilde{\sigma}_{\text{data}} \alpha_{\lambda}}{e^{-\lambda/2} \sqrt{e^{-\lambda} + \tilde{\sigma}_{\text{data}}^2}} \mathbf{z}_{\lambda} \quad (21)$$

where $\tilde{\sigma}_{\text{data}} = 0.5$.

Variational Diffusion Models

Starting with Kingma et al. (2021)



Let's get probabilistic.

Evidence Lower Bound

Slowly approaching the paper for today



- ▶ the log evidence (marginal likelihood) of the data \mathbf{x} under the model is

$$\log p(\mathbf{x}) = \log \int_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \quad (22)$$

where $p(\mathbf{x}, \mathbf{z})$ is the joint distribution of the data \mathbf{x} and latent variables \mathbf{z} .

- ▶ we can derive a lower bound (Jensen's inequality) on the log-evidence, introducing the variational distribution $q(\mathbf{z}|\mathbf{x})$

$$\log p(\mathbf{x}) = \log \int_{\mathbf{z}} q(\mathbf{z}|\mathbf{x}) \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})} d\mathbf{z} \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})} \right] \quad (23)$$

This gives us the evidence lower bound (ELBO) (Blei et al., 2017)

- ▶ we rewrite the ELBO

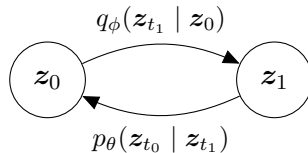
$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] - D_{KL}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \quad (24)$$

From Variational Autoencoders to Diffusion Models

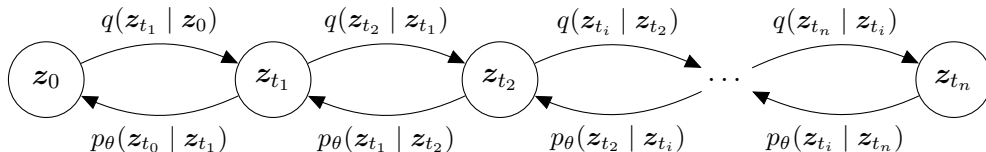
Originally by Kingma & Welling (2014) and Rezende et al. (2014)



Variational Autoencoder



Diffusion Model





- ▶ We aim to learn a generative model $p_{\theta}(\mathbf{x})$ that approximates $q(\mathbf{x})$, where \mathbf{x} is drawn from a dataset. (think: \mathbf{x} image or latent representation)
- ▶ The model incorporates a sequence of latent variables \mathbf{z}_t for $t \in [0, 1]$, where $\mathbf{z}_{0,\dots,1} := \mathbf{z}_0, \dots, \mathbf{z}_1$.
- ▶ forward process: defines a conditional distribution $q(\mathbf{z}_{0,\dots,1}|\mathbf{x})$
- ▶ generative model: defines a joint distribution $p(\mathbf{z}_{0,\dots,1})$



- ▶ forward diffusion process generates a sequence of increasingly noisy versions of the data \mathbf{x} , latent variables at time $t \in [0, 1]$ are denoted \mathbf{z}_t

$$q(\mathbf{z}_t|\mathbf{x}) = \mathcal{N}(\alpha_t \mathbf{x}, \sigma_t^2 \mathbf{I}) \quad \Rightarrow \quad \mathbf{z}_t = \alpha_t \mathbf{x} + \sigma_t \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (25)$$

where α_t and σ_t^2 are strictly positive scalar functions of time.

- ▶ Signal-to-Noise Ratio (SNR)

$$\text{SNR}(t) = \frac{\alpha_t^2}{\sigma_t^2} \quad (26)$$

- ▶ variance-preserving (VP): $\alpha_t = \sqrt{1 - \sigma_t^2}$, variance-exploding (VE): $\alpha_t^2 = 1$
- ▶ noise schedule σ_t^2 is learned via a monotonic neural network $\gamma(t)$

$$\sigma_t^2 = \text{sigmoid}(\gamma_{\boldsymbol{\eta}}(t)), \quad \alpha_t^2 = \text{sigmoid}(-\gamma_{\boldsymbol{\eta}}(t)), \quad \text{SNR}(t) = \exp(-\gamma_{\boldsymbol{\eta}}(t)) \quad (27)$$

- ▶ reverse diffusion process is the (hierarchical) generative model

$$p(\mathbf{x}) = \int_{\mathbf{z}} p(\mathbf{z}_1) p(\mathbf{x}|\mathbf{z}_0) \prod_{i=1}^T p(\mathbf{z}_{s(i)}|\mathbf{z}_{t(i)}), \quad (28)$$

where $s(i) = (i - 1)/T$, and $t(i) = i/T$.

- ▶ for $t = T$, with VP diffusion and sufficiently small SNR(1),

$$p(\mathbf{z}_1) = \mathcal{N}(\mathbf{z}_1; 0, \mathbf{I}), \quad (29)$$

- ▶ generative model approximates the (unknown) conditional distribution $q(\mathbf{x}|\mathbf{z}_0)$

$$p(\mathbf{x}|\mathbf{z}_0) = \prod_i p(x_i|z_{0,i}), \quad (30)$$

where $p(x_i|z_{0,i}) \propto q(z_{0,i}|x_i)$, remember: $q(\mathbf{z}_t|\mathbf{x}) = \mathcal{N}(\alpha_t \mathbf{x}, \sigma_t^2 \mathbf{I})$



- conditional model distributions

$$p(\mathbf{z}_s | \mathbf{z}_t) = q(\mathbf{z}_s | \mathbf{z}_t, \mathbf{x} = \hat{\mathbf{x}}_\theta(\mathbf{z}_t; t)), \quad (31)$$

where $\hat{\mathbf{x}}_\theta(\mathbf{z}_t; t)$ is the noise-prediction parameterization

- we see

$$q(\mathbf{z}_t | \mathbf{z}_s) = \mathcal{N}(\alpha_{t|s} \mathbf{z}_s, \sigma_{t|s}^2 \mathbf{I}), \quad \alpha_{t|s} = \frac{\alpha_t}{\alpha_s}, \quad \sigma_{t|s}^2 = \sigma_t^2 - \alpha_{t|s}^2 \sigma_s^2 \quad (32)$$

- posterior $q(\mathbf{z}_s | \mathbf{z}_t, \mathbf{x}) \propto q(\mathbf{z}_s | \mathbf{x}) q(\mathbf{z}_t | \mathbf{z}_s)$, we get

$$q(\mathbf{z}_s | \mathbf{z}_t, \mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_Q(\mathbf{z}_t, \mathbf{x}; s, t), \sigma_Q^2(s, t) \mathbf{I}), \quad (33)$$

where

$$\sigma_Q^2(s, t) = \frac{\sigma_s^2 \sigma_{t|s}^2}{\sigma_t^2}, \quad \boldsymbol{\mu}_Q(\mathbf{z}_t, \mathbf{x}; s, t) = \frac{\alpha_{t|s} \sigma_s^2}{\sigma_t^2} \mathbf{z}_t + \frac{\alpha_s \sigma_{t|s}^2}{\sigma_t^2} \mathbf{x}. \quad (34)$$

- for the generative model we replace \mathbf{x} by the predicted $\hat{\mathbf{x}}_\theta(\mathbf{z}_t; t)$

$$p(\mathbf{z}_s | \mathbf{z}_t) = \mathcal{N}(\boldsymbol{\mu}_\theta(\mathbf{z}_t; s, t), \sigma_Q^2(s, t) \mathbf{I}), \quad (35)$$

- ▶ remember: ELBO

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] - \text{KL}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \quad (36)$$

- ▶ minimize the negative ELBO

$$\begin{aligned} -\log p(\mathbf{x}) &\leq -\text{VLB}(\mathbf{x}) \\ &= \underbrace{D_{KL}(q(\mathbf{z}_1|\mathbf{x})||p(\mathbf{z}_1))}_{\text{Prior loss}} + \underbrace{\mathbb{E}_{q(\mathbf{z}_0|\mathbf{x})} [-\log p(\mathbf{x}|\mathbf{z}_0)]}_{\text{Reconstruction loss}} + \underbrace{\mathcal{L}_T}_{\text{Diffusion loss}} \end{aligned} \quad (37)$$

- ▶ for finite T , the diffusion loss can be expressed as

$$\mathcal{L}_T = \sum_{i=1}^T \mathbb{E}_{q(\mathbf{z}_{t(i)}|\mathbf{x})} D_{KL}[q(\mathbf{z}_{s(i)}|\mathbf{z}_{t(i)}, \mathbf{x})||p_{\theta}(\mathbf{z}_{s(i)}|\mathbf{z}_{t(i)})] \quad (38)$$

- we can write

$$\mathcal{L}_T = \sum_{i=1}^T \mathbb{E}_{q(\mathbf{z}_{t(i)}|\mathbf{x})} D_{KL}[q(\mathbf{z}_{s(i)}|\mathbf{z}_{t(i)}, \mathbf{x}) || p_{\theta}(\mathbf{z}_{s(i)}|\mathbf{z}_{t(i)})] \quad (39)$$

$$= \frac{T}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}), i \sim U\{1, T\}} [(\text{SNR}(s) - \text{SNR}(t)) \|\mathbf{x} - \hat{\mathbf{x}}_{\theta}(\mathbf{z}_t; t)\|_2^2] \quad (40)$$

- here, \mathcal{L}_T is an unbiased Monte Carlo estimator, where
 - $U\{1, T\}$ is the uniform distribution over $\{1, \dots, T\}$,
 - $s = (i-1)/T, t = i/T$,
 - $\mathbf{z}_t = \alpha_t \mathbf{x} + \sigma_t \epsilon$, where $\epsilon \sim \mathcal{N}(0, \mathbf{I})$.

Variational Lower Bound

A Closer Look on the Connection to the SNR



► **KL Divergence between Gaussians:**

$$D_{KL}(q(z_s|z_t, \mathbf{x})||p_\theta(z_s|z_t)) = \frac{1}{2\sigma_Q^2(s, t)} \|\boldsymbol{\mu}_Q - \boldsymbol{\mu}_\theta\|_2^2 \quad (41)$$

► **Mean Difference $\|\boldsymbol{\mu}_Q - \boldsymbol{\mu}_\theta\|_2^2$:**

$$\boldsymbol{\mu}_Q - \boldsymbol{\mu}_\theta = \frac{\alpha_s^2}{\sigma_s^2} (\mathbf{x} - \hat{\mathbf{x}}_\theta(z_t; t)) \quad (42)$$

► **Simplification of KL Divergence:**

$$D_{KL}(q(z_s|z_t, \mathbf{x})||p_\theta(z_s|z_t)) = \frac{1}{2} \frac{\alpha_s^2}{\sigma_s^2} \left(1 - \frac{\alpha_t^2}{\sigma_t^2}\right) \|\mathbf{x} - \hat{\mathbf{x}}_\theta(z_t; t)\|_2^2 \quad (43)$$

► **Final Expression with SNR:**

$$D_{KL}(q(z_s|z_t, \mathbf{x})||p_\theta(z_s|z_t)) = \frac{1}{2} (\text{SNR}(s) - \text{SNR}(t)) \|\mathbf{x} - \hat{\mathbf{x}}_\theta(z_t; t)\|_2^2 \quad (44)$$

- ▶ in the limit as $T \rightarrow \infty$, the diffusion loss \mathcal{L}_T can be written as a function of $\tau = 1/T$

$$\mathcal{L}_T = \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}), i \sim U\{1, T\}} \left[\frac{\text{SNR}(t - \tau) - \text{SNR}(t)}{\tau} \|\mathbf{x} - \hat{\mathbf{x}}_{\theta}(\mathbf{z}_t; t)\|_2^2 \right] \quad (45)$$

- ▶ as $\tau \rightarrow 0$ and $T \rightarrow \infty$, we get

$$\mathcal{L}_{\infty} = -\frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}), t \sim \mathcal{U}[0, 1]} [\text{SNR}'(t) \|\mathbf{x} - \hat{\mathbf{x}}_{\theta}(\mathbf{z}_t; t)\|_2^2] \quad (46)$$

- ▶ The continuous-time formulation of the diffusion loss is then:

$$\mathcal{L}_{\infty} = -\frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \int_0^1 \text{SNR}'(t) \|\mathbf{x} - \hat{\mathbf{x}}_{\theta}(\mathbf{z}_t; t)\|_2^2 dt \quad (47)$$



- ▶ we know: the signal-to-noise ratio (SNR) function is invertible: $t = \text{SNR}^{-1}(v)$.
- ▶ therefore, by a change of variables (t to v) we can express the diffusion loss as

$$\mathcal{L}_\infty = \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \int_{\text{SNR}_{\min}}^{\text{SNR}_{\max}} \|\mathbf{x} - \tilde{\mathbf{x}}_\theta(\mathbf{z}_v, v)\|_2^2 dv \quad (48)$$

- ▶ shape of SNR function between $t = 0$ and $t = 1$ does not affect the diffusion loss; only the values at the endpoints (SNR_{\min} and SNR_{\max}) matter
- ▶ two different diffusion processes with the same SNR_{\min} and SNR_{\max} will define the same distribution $p(\mathbf{x})$, up to a rescaling of the latents
- ▶ therefore, variance-preserving and variance-exploding diffusion models are equivalent in continuous time



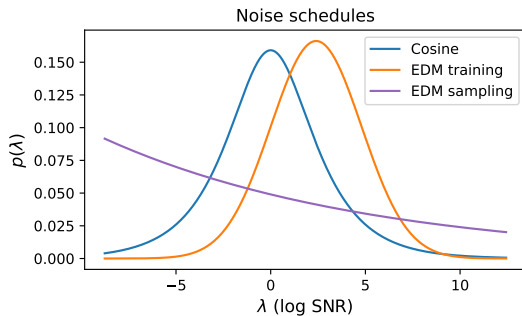
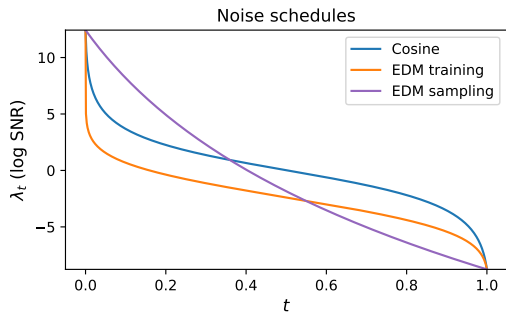
- ▶ The continuous-time diffusion loss can be generalized to a weighted loss (no equivalence)

$$\mathcal{L}_{\infty}(\mathbf{x}, w) = \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \int_{\text{SNR}_{\min}}^{\text{SNR}_{\max}} w(v) \|\mathbf{x} - \tilde{\mathbf{x}}_{\theta}(\mathbf{z}_v, v)\|_2^2 dv \quad (49)$$

- ▶ $w(v)$ is a weighting function that can emphasize different noise levels
- ▶ different objectives in diffusion models are special cases of this weighted diffusion loss



- Noise schedule is a strictly monotonically decreasing function $\lambda = f_\lambda(t)$ with endpoints $\lambda_{\max} := f_\lambda(0)$ and $\lambda_{\min} := f_\lambda(1)$.



Noise Schedules

Deriving the PDF from the Noise Schedule



- During training, time t is sampled uniformly: $t \sim \mathcal{U}(0, 1)$. The noise level λ is then computed via $\lambda = f_\lambda(t)$.
- This results in a distribution over noise levels, with the probability density function (PDF) given by:

$$p(\lambda_t) = -\frac{d}{d\lambda} f_\lambda^{-1}(\lambda_t) = -\frac{dt}{d\lambda} = -\frac{1}{f'_\lambda(t)}$$

- The PDF $p(\lambda_t)$ represents the likelihood of encountering a particular noise level λ during training.

| Noise schedule name | $\lambda = f_\lambda(t) = \dots$ | $t = f_\lambda^{-1}(\lambda) = \dots$ | $p(\lambda) = -\frac{d}{d\lambda} f_\lambda^{-1}(\lambda)$ |
|---------------------------|--|---|---|
| Cosine | $-2 \log(\tan(\pi t/2))$ | $(2/\pi) \arctan(e^{-\lambda/2})$ | $\text{sech}(\lambda/2)/(2\pi)$ |
| Shifted Cosine | $-2 \log(\tan(\pi t/2)) + 2s$ | $(2/\pi) \arctan(e^{-\lambda/2-s})$ | $\text{sech}(\lambda/2 - s)/(2\pi)$ |
| EDM (Training) | $-F_{\mathcal{N}}^{-1}(t; 2.4, 2.4^2)$ | $F_{\mathcal{N}}(-\lambda; 2.4, 2.4^2)$ | $\mathcal{N}(\lambda; 2.4, 2.4^2)$ |
| EDM (Sampling) | $-2\rho \log(\sigma_{\max}^{1/\rho} + (1-t)(\sigma_{\min}^{1/\rho} - \sigma_{\max}^{1/\rho}))$ | $1 - \frac{e^{-\lambda/(2\rho)} - \sigma_{\max}^{1/\rho}}{\sigma_{\min}^{1/\rho} - \sigma_{\max}^{1/\rho}}$ | $\frac{e^{-\lambda/(2\rho)}}{2\rho(\sigma_{\max}^{1/\rho} - \sigma_{\min}^{1/\rho})}$ |
| Flow Matching (OT) | $2 \log((1-t)/t)$ | $1/(1 + e^{\lambda/2})$ | $\text{sech}^2(\lambda/4)/8$ |

- ▶ we define the log signal-to-noise ratio (SNR) as

$$\lambda = \log(\alpha_\lambda^2 / \sigma_\lambda^2) \quad (50)$$

- ▶ therefore, we can write \mathcal{L}_∞ as

$$\mathcal{L}_\infty = -\frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}), t \sim \mathcal{U}[0, 1]} [\text{SNR}'(t) \|\mathbf{x} - \hat{\mathbf{x}}_\theta(\mathbf{z}_t; t)\|_2^2] \quad (51)$$

$$= -\frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}), t \sim \mathcal{U}[0, 1]} \left[\frac{d\lambda}{dt} \|\mathbf{x} - \hat{\mathbf{x}}_\theta(\mathbf{z}_t; t)\|_2^2 \right] \quad (52)$$

- ▶ we can similarly use the noise prediction parameterization

$$-\text{ELBO}(\mathbf{x}) = \frac{1}{2} \mathbb{E}_{t \sim \mathcal{U}(0, 1), \epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[-\frac{d\lambda}{dt} \cdot \|\hat{\epsilon}_\theta(\mathbf{z}_t; \lambda_t) - \epsilon\|_2^2 \right] + c \quad (53)$$

Negative ELBO and Further Connections

Reiterating Various Parameterizations



- noise prediction parameterization

$$-\text{ELBO}(\mathbf{x}) = \frac{1}{2} \mathbb{E}_{t \sim \mathcal{U}(0,1), \epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[-\frac{d\lambda}{dt} \cdot \|\hat{\epsilon}_{\theta}(\mathbf{z}_t; \lambda_t) - \epsilon\|_2^2 \right] + c \quad (54)$$

- using

$$\mathbf{s}_{\theta}(\mathbf{z}_t; \lambda_t) = -\frac{\hat{\epsilon}_{\theta}(\mathbf{z}_t; \lambda_t)}{\sigma_{\lambda}}$$

and $\tilde{w}(t) = \sigma_t^2$ we can connect

$$\mathcal{L}_{\text{DSM}}(\mathbf{x}) = \mathbb{E}_{t \sim \mathcal{U}(0,1), \epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[\tilde{w}(t) \cdot \|\mathbf{s}_{\theta}(\mathbf{z}_t, \lambda_t) - \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{x})\|_2^2 \right] \quad (55)$$

to

$$\mathcal{L}_{\epsilon}(\mathbf{x}) = \frac{1}{2} \mathbb{E}_{t \sim \mathcal{U}(0,1), \epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[\|\hat{\epsilon}_{\theta}(\mathbf{z}_t; \lambda_t) - \epsilon\|_2^2 \right] \quad (56)$$



- ▶ general form of the weighted loss is

$$\mathcal{L}_w(\mathbf{x}) = \frac{1}{2} \mathbb{E}_{t \sim \mathcal{U}(0,1), \epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[w(\lambda_t) \cdot \left(-\frac{d\lambda}{dt} \right) \cdot \|\hat{\epsilon}_\theta(\mathbf{z}_t; \lambda_t) - \epsilon\|_2^2 \right] \quad (57)$$

- ▶ weighted loss can be rewritten as an integral

$$\mathcal{L}_w(\mathbf{x}) = \frac{1}{2} \int_{\lambda_{\min}}^{\lambda_{\max}} w(\lambda) \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} [\|\hat{\epsilon}_\theta(\mathbf{z}_\lambda; \lambda) - \epsilon\|_2^2] d\lambda \quad (58)$$

- ▶ loss does not depend on the specific noise schedule λ_t except for the endpoints λ_{\min} and λ_{\max} , noise schedule (only!) affects the variance of the Monte Carlo estimator



- ▶ the noise schedule $p(\lambda)$ acts as an *importance sampling distribution*

$$\mathcal{L}_w(\mathbf{x}) = \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}), \lambda \sim p(\lambda)} \left[\frac{w(\lambda_t)}{p(\lambda)} \|\hat{\epsilon}_\theta(\mathbf{z}_t; \lambda_t) - \epsilon\|_2^2 \right] \quad (59)$$

Theorem

If the weighting function $w(\lambda_t)$ is monotonic, then the weighted diffusion objective is equivalent to the ELBO with data augmentation (additive noise).

Result: Any objective with (implied) monotonic weighting, can be understood as equivalent to the ELBO with simple data augmentation (additive noise)

- ▶ adaptive noise schedule: by lowering the variance of the loss estimator, this often significantly speeds up optimization

Proof of Theorem 1 (Kingma & Gao, 2024)

KL Divergence and Time Derivative



- **KL Divergence:** define $\mathcal{L}(t; \mathbf{x}) = D_{KL}(q(\mathbf{z}_{t,\dots,1}|\mathbf{x})||p(\mathbf{z}_{t,\dots,1}))$ for the KL divergence between $q(\mathbf{z}_{t,\dots,1}|\mathbf{x})$ and $p(\mathbf{z}_{t,\dots,1})$ for timesteps t to 1

$$\mathcal{L}(t; \mathbf{x}) := D_{KL}q(\mathbf{z}_{t,\dots,1}|\mathbf{x})||p(\mathbf{z}_{t,\dots,1})) \quad (60)$$

- **Time Derivative:** show that the time derivative of $\mathcal{L}(t; \mathbf{x})$ is

$$\frac{d}{dt}\mathcal{L}(t; \mathbf{x}) = \frac{1}{2} \frac{d\lambda}{dt} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} [\|\epsilon - \hat{\epsilon}_{\theta}(\mathbf{z}_t; \lambda_t)\|_2^2] \quad (61)$$

- **Rewriting the Weighted Loss:** therefore, the weighted loss is rewritten as

$$\mathcal{L}_w(\mathbf{x}) = - \int_0^1 \frac{d}{dt} \mathcal{L}(t; \mathbf{x}) w(\lambda_t) dt \quad (62)$$

Proof of Theorem 1 (Kingma & Gao, 2024)

Integration by Parts and Monotonic Weighting



- weighted loss

$$\mathcal{L}_w(\mathbf{x}) = - \int_0^1 \frac{d}{dt} \mathcal{L}(t; \mathbf{x}) w(\lambda_t) dt \quad (63)$$

- **Integration by Parts:** integration by parts yields

$$\mathcal{L}_w(\mathbf{x}) = \int_0^1 \frac{d}{dt} w(\lambda_t) \mathcal{L}(t; \mathbf{x}) dt + w(\lambda_{\max}) \mathcal{L}(0; \mathbf{x}) + \text{constant} \quad (64)$$

- **Monotonic Weighting:** assume $w(\lambda_t)$ is monotonically increasing and normalized

$$\mathcal{L}_w(\mathbf{x}) = \mathbb{E}_{p_w(t)} [\mathcal{L}(t; \mathbf{x})] + \text{constant} \quad (65)$$

- **Probability Distribution:** $p_w(t) = \frac{d}{dt} w(\lambda_t)$ is a probability distribution over $t \in [0, 1]$, meaning that $\mathcal{L}_w(\mathbf{x})$ becomes an expected KL divergence.



- we note that

$$\mathcal{L}(t; \mathbf{x}) = D_{KL}(q(\mathbf{z}_{t,\dots,1}|\mathbf{x})||p(\mathbf{z}_{t,\dots,1})) = -\mathbb{E}_{q(\mathbf{z}_t|\mathbf{x})}[\text{ELBO}_t(\mathbf{z}_t)] - \mathcal{H}(q(\mathbf{z}_t|\mathbf{x})) \quad (66)$$

$$\text{ELBO}_t(\mathbf{z}_t) := \mathbb{E}_{q(\tilde{\mathbf{z}}_t|\mathbf{z}_t)}[\log p(\mathbf{z}_t, \tilde{\mathbf{z}}_t) - \log q(\tilde{\mathbf{z}}_t|\mathbf{z}_t)] \leq \log p(\mathbf{z}_t), \tilde{\mathbf{z}}_t := \mathbf{z}_{t+dt,\dots,1} \quad (67)$$

- therefore,

$$\mathcal{L}(t; \mathbf{x}) = D_{KL}(q(\mathbf{z}_{t,\dots,1}|\mathbf{x})||p(\mathbf{z}_{t,\dots,1})) \geq D_{KL}(q(\mathbf{z}_t|\mathbf{x})||p(\mathbf{z}_t)) \quad (68)$$

- $\mathcal{L}(t; \mathbf{x})$ is the expected negative ELBO of noise-perturbed data \mathbf{z}_t

$$\mathcal{L}(t; \mathbf{x}) = -\mathbb{E}_{q(\mathbf{z}_t|\mathbf{x})}[\text{ELBO}_t(\mathbf{z}_t)] + c. \geq -\mathbb{E}_{q(\mathbf{z}_t|\mathbf{x})}[\log p(\mathbf{z}_t)] + c. \quad (69)$$

$$\mathcal{L}_w(\mathbf{x}) = \mathbb{E}_{p_w(t)} [\mathcal{L}(t; \mathbf{x})] + c.$$

$$= - \underbrace{\mathbb{E}_{p_w(t), q(\mathbf{z}_t|\mathbf{x})} [\text{ELBO}_t(\mathbf{z}_t)]}_{\text{ELBO of noise-perturbed data}} + c. \geq - \underbrace{\mathbb{E}_{p_w(t), q(\mathbf{z}_t|\mathbf{x})} [\log p(\mathbf{z}_t)]}_{\text{ll. of noise-perturbed data}} + c. \quad (70)$$

Results

ImageNet 64×64



| Model parameterization | Training noise schedule | Weighting function | Monotonic? | DDPM sampler | | EDM sampler | |
|---------------------------|-------------------------|-------------------------------------|------------|--------------|-----------------------------------|-------------|-----------------------------------|
| | | | | FID ↓ | IS ↑ | FID ↓ | IS ↑ |
| ϵ -prediction | Cosine | $\text{sech}(\lambda/2)$ (Baseline) | | 1.85 | 54.1 ± 0.79 | 1.55 | 59.2 ± 0.78 |
| " | Cosine | $\text{sigmoid}(-\lambda + 1)$ | ✓ | 1.75 | 55.3 ± 1.23 | | |
| " | Cosine | $\text{sigmoid}(-\lambda + 2)$ | ✓ | 1.68 | 56.8 ± 0.85 | 1.46 | 60.4 ± 0.86 |
| " | Cosine | $\text{sigmoid}(-\lambda + 3)$ | ✓ | 1.73 | 56.1 ± 1.36 | | |
| " | Cosine | $\text{sigmoid}(-\lambda + 4)$ | ✓ | 1.80 | 55.1 ± 1.65 | | |
| " | Cosine | $\text{sigmoid}(-\lambda + 5)$ | ✓ | 1.94 | 53.5 ± 1.12 | | |
| " | Adaptive | $\text{sigmoid}(-\lambda + 2)$ | ✓ | 1.70 | 54.8 ± 1.20 | 1.44 | 60.6 ± 1.44 |
| " | Adaptive | EDM-monotonic | ✓ | 1.67 | 56.8 ± 0.90 | 1.44 | 61.1 ± 1.80 |
| EDM (Karras et al., 2022) | EDM (training) | EDM (Baseline) | | | | 1.36 | |
| EDM (our reproduction) | EDM (training) | EDM (Baseline) | | | | 1.45 | 60.7 ± 1.19 |
| " | Adaptive | EDM | | | | 1.43 | 63.2 ± 1.76 |
| " | Adaptive | $\text{sigmoid}(-\lambda + 2)$ | ✓ | | | 1.55 | 63.7 ± 1.14 |
| " | Adaptive | EDM-monotonic | ✓ | | | 1.43 | 63.7 ± 1.48 |
| v -prediction | Adaptive | $\exp(-\lambda/2)$ (Baseline) | ✓ | | | 1.62 | 58.0 ± 1.56 |
| " | Adaptive | $\text{sigmoid}(-\lambda + 2)$ | ✓ | | | 1.51 | 64.4 ± 1.28 |
| " | Adaptive | EDM-monotonic | ✓ | | | 1.45 | 64.6 ± 1.35 |

Results

ImageNet 128×128



| Model parameterization | Training noise schedule | Weighting function | Monotonic? | FID ↓ | | IS ↑ |
|------------------------|-------------------------|---|------------|-------------|-------------|------------------------------------|
| | | | | train | eval | |
| v -prediction | Cosine-shifted | $\exp(-\lambda/2)$ (Baseline) | ✓ | 1.91 | 3.23 | 171.9 ± 2.46 |
| " | Adaptive | $\text{sigmoid}(-\lambda + 2)$ -shifted | ✓ | 1.91 | 3.41 | 183.1 ± 2.20 |
| " | Adaptive | EDM-monotonic-shifted | ✓ | 1.75 | 2.88 | 171.1 ± 2.67 |



| Method | Without guidance | | | With guidance | | |
|---|------------------|-------------|--------------------|----------------|-------------|--------------------|
| | FID ↓ train | eval | IS ↑ | FID ↓ train | eval | IS ↑ |
| 128 × 128 resolution | | | | | | |
| ADM (Dhariwal & Nichol, 2021) | 5.91 | – | – | 2.97 | – | – |
| CDM (Ho et al., 2022) | 3.52 | 3.76 | 128.8 ± 2.5 | – | – | – |
| RIN (Jabri et al., 2023) | 2.75 | – | 144.1 | – | – | – |
| Simple Diffusion (U-Net) (Hooeboom et al., 2023) | 2.26 | 2.88 | 137.3 ± 2.0 | – | – | – |
| Simple Diffusion (U-ViT, L) (Hooeboom et al., 2023) | 1.91 | 3.23 | 171.9 ± 2.5 | 2.05 | 3.57 | 189.9 ± 3.5 |
| VDM++ (Ours) , $w(\lambda) = \text{sigmoid}(-\lambda + 2)$ | 1.91 | 3.41 | 183.1 ± 2.2 | – | – | – |
| VDM++ (Ours) , EDM-monotonic weighting | 1.75 | 2.88 | 171.1 ± 2.7 | 1.78 | 3.16 | 190.5 ± 2.3 |
| 256 × 256 resolution | | | | | | |
| BigGAN-deep (no truncation) (Brock, 2018) | 6.9 | – | 171.4 ± 2.0 | – | – | – |
| MaskGIT (Chang et al., 2022) | 6.18 | – | 182.1 | – | – | – |
| ADM (Dhariwal & Nichol, 2021) | 10.94 | – | – | 3.94 | – | 215.9 |
| CDM (Ho et al., 2022) | 4.88 | 4.63 | 158.7 ± 2.3 | – | – | – |
| RIN (Jabri et al., 2023) | 3.42 | – | 182.0 | – | – | – |
| Simple Diffusion (U-Net) (Hooeboom et al., 2023) | 3.76 | 3.71 | 171.6 ± 3.1 | – | – | – |
| Simple Diffusion (U-ViT, L) (Hooeboom et al., 2023) | 2.77 | 3.75 | 211.8 ± 2.9 | 2.44 | 4.08 | 256.3 ± 5.0 |
| VDM++ (Ours) , EDM-monotonic weighting | 2.40 | 3.36 | 225.3 ± 3.2 | 2.12 | 3.69 | 267.7 ± 4.9 |



| Method | Without guidance | | | With guidance | | |
|--|------------------|-------------|--------------------|----------------|-------------|--------------------|
| | FID ↓ train | eval | IS ↑ | FID ↓ train | eval | IS ↑ |
| <i>Latent diffusion with pretrained VAE:</i> | | | | | | |
| DiT-XL/2 (Peebles & Xie, 2023) | 9.62 | – | 121.5 | 2.27 | – | 278.2 |
| U-ViT (Bao et al., 2023) | – | – | – | 3.40 | – | – |
| Min-SNR- γ (Hang et al., 2023) | – | – | – | 2.06 | – | – |
| MDT (Gao et al., 2023) | 6.23 | – | 143.0 | 1.79 | – | 283.0 |
| 512 × 512 resolution | | | | | | |
| MaskGIT (Chang et al., 2022) | 7.32 | – | 156.0 | – | – | – |
| ADM (Dhariwal & Nichol, 2021) | 23.24 | – | – | 3.85 | – | 221.7 |
| RIN (Jabri et al., 2023) | – | – | – | 3.95 | – | 216.0 |
| Simple Diffusion (U-Net) (Hoogeboom et al., 2023) | 4.30 | 4.28 | 171.0 ± 3.0 | – | – | – |
| Simple Diffusion (U-ViT, L) (Hoogeboom et al., 2023) | 3.54 | 4.53 | 205.3 ± 2.7 | 3.02 | 4.60 | 248.7 ± 3.4 |
| VDM++ (Ours), EDM-monotonic weighting | 2.99 | 4.09 | 232.2 ± 4.2 | 2.65 | 4.43 | 278.1 ± 5.5 |
| <i>Latent diffusion with pretrained VAE:</i> | | | | | | |
| DiT-XL/2 (Peebles & Xie, 2023) | 12.03 | – | 105.3 | 3.04 | – | 240.8 |
| LDM-4 (Rombach et al., 2022) | 10.56 | – | 103.5 ± 1.2 | 3.60 | – | 247.7 ± 5.6 |



Summary of Results:

- ▶ We have demonstrated that the **weighted diffusion loss** generalizes diffusion objectives in the literature.
- ▶ It can be interpreted as a **weighted integral of ELBO objectives**, where each ELBO corresponds to a different noise level.
- ▶ When the weighting function $w(\lambda_t)$ is monotonic, the loss has an interpretation as the ELBO objective with **data augmentation** via noise perturbation.

Implications:

- ▶ The equivalence between **monotonic weighting** and the **ELBO with data augmentation** enables a direct comparison between diffusion models and other **likelihood-based models** such as autoregressive transformers.
- ▶ This opens up avenues for optimizing other model types toward the same objective as monotonically weighted diffusion models.

Discussion