

Understanding Diffusion Models: A Unified Perspective

Calvin Luo

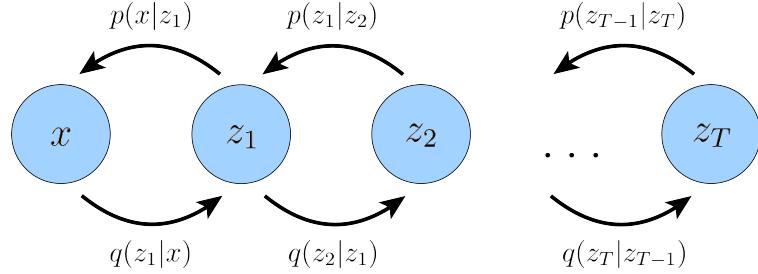
Google Research, Brain Team

calvinluo@google.com

August 26, 2022

Contents

Introduction: Generative Models	1
Background: ELBO, VAE, and Hierarchical VAE	2
Evidence Lower Bound	2
Variational Autoencoders	4
Hierarchical Variational Autoencoders	5
Variational Diffusion Models	6
Learning Diffusion Noise Parameters	14
Three Equivalent Interpretations	15
Score-based Generative Models	17
Guidance	20
Classifier Guidance	21
Classifier-Free Guidance	21
Closing	22



Markovi anHVAE:

$$p(\mathbf{x}, \mathbf{z}_{1:T}) = p(\mathbf{z}_T)p_{\theta}(\mathbf{x}|\mathbf{z}_1) \prod_{t=2}^T p_{\theta}(\mathbf{z}_{t-1}|\mathbf{z}_t) \quad (23)$$

$$q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}) = q_{\phi}(\mathbf{z}_1|\mathbf{x}) \prod_{t=2}^T q_{\phi}(\mathbf{z}_t|\mathbf{z}_{t-1}) \quad (24)$$

==> ELBO :

$$\begin{aligned} \log p(\mathbf{x}) &= \log \int p(\mathbf{x}, \mathbf{z}_{1:T}) d\mathbf{z}_{1:T} \\ &\geq \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})} \right] \end{aligned} \quad (25)$$

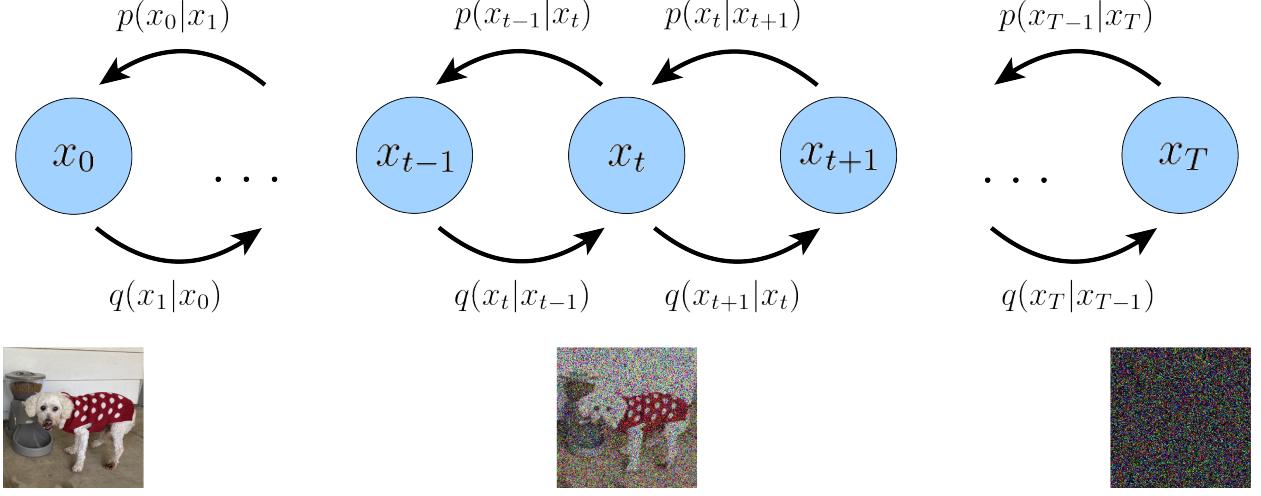
==>

$$\mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})} \right] = \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})} \left[\log \frac{p(\mathbf{z}_T)p_{\theta}(\mathbf{x}|\mathbf{z}_1) \prod_{t=2}^T p_{\theta}(\mathbf{z}_{t-1}|\mathbf{z}_t)}{q_{\phi}(\mathbf{z}_1|\mathbf{x}) \prod_{t=2}^T q_{\phi}(\mathbf{z}_t|\mathbf{z}_{t-1})} \right] \quad (29)$$

Variational Diffusion Models

a VDM [4, 5, 6] i sa MHVAE wi th

- Th el latent di mensi on sexa ctly equal to th edata di mensi on
- Th estructure of th el latent encoder a teachti mestepi not learned; i ti spre-defined as a li nea Gaussian model. In oth erwords, i ti a Gaussian and stri b uti centered around the output of the previous mestep
- Th eGaussi a para meters of th el latent encoders vary over time in such a way that the stri b uti ofn th el latent at final ti mestep T i sa standa rdGaussian an



Furthermore, we explore the hierarchical Markov property between hierarchical transitions from a standard HVAE.

Let us expand on the implications of these assumptions. From the first restriction, we can now represent both true data samples and latent variables \mathbf{x}_t , where $t = 0$ represents true data samples and $t \in [1, T]$ represents corresponding latent variables indexed by t . The VIM posterior is the same as the eHVAE posterior (Eq24):

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}) \quad (30)$$

From the second assumption, we know that the distribution of each latent variable in the encoder is a Gaussian centered around its previous hierarchical latent. Unlike a Markovian HVAE, the structure of the encoder at each time step is not learned; it is fixed as a linear Gaussian model, where the mean and standard deviation can be set beforehand as hyperparameters learned as parameters [6]. We parameterize the Gaussian encoder with $\mu(\mathbf{x}_t) = \sqrt{\alpha_t} \mathbf{x}_{t-1}$, and variance $\Sigma_t(\mathbf{x}_t) = (1 - \alpha_t) \mathbf{I}$, where the form of the coefficients are chosen such that the latent variables stay at a similar scale; in other words, the encoding process is *variance-preserving*. Note that alternative Gaussian parameterizations are allowed, and lead to similar derivations. The main takeaway is that α_t is a (potentially learnable) coefficient that can vary with the hierarchy. Mathematically, encoder transitions are denoted as:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t} \mathbf{x}_{t-1}, (1 - \alpha_t) \mathbf{I}) \quad (31)$$

From the third assumption, we know that evolves over time according to a fixed or learnable schedule structured such that the distribution of the final \mathbf{x}_T is a standard Gaussian. We can then update the joint distribution of a Markovian HVAE (Eq23) to write the joint distribution for a VIM as:

$$p(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \quad (32)$$

where,

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I}) \quad (33)$$

Collectively, what this set of assumptions describes is a steady noise injection of an image input over time; we progressively corrupt an image by adding Gaussian noise until it becomes completely identical to pure Gaussian noise. Visually, this process is depicted in Figure

Note that our encoder distribution $q(\mathbf{z}_t | \mathbf{x}_{t-1})$ are no longer parameterized by ϕ , as they are completely modeled as Gaussian and with defined mean and variance parameters at each timestep. Therefore, in a VIM we are only interested in learning conditional probability $p_{\theta}(\mathbf{z}_{t-1} | \mathbf{x}_t)$, so that we can simulate new data. After optimizing the VIM the sampling procedure is as simple as sampling Gaussian noise \mathbf{z}_{t-1} and iteratively running the denoising transition $t_{\theta}(\mathbf{z}_{t-1} | \mathbf{x}_t)$ for T steps to generate a novel \mathbf{x}_0 .

Like any HVAE, the eVIM can be optimized by maximizing the eELBO:

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right] \quad (37)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right] \quad (38)$$

$$= \mathbb{E}_{q(\mathbf{x}_1 | \mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0 | \mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{T-1}, \mathbf{x}_T | \mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T | \mathbf{x}_{T-1})} \right] + \sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_t, \mathbf{x}_{t+1} | \mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})}{q(\mathbf{x}_t | \mathbf{x}_{t-1})} \right] \quad (44)$$

$$\begin{aligned} &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1 | \mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0 | \mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{\mathbb{E}_{q(\mathbf{x}_{T-1} | \mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_{T-1}) \| p(\mathbf{x}_T))]}_{\text{prior matching term}} \\ &\quad - \sum_{t=1}^{T-1} \underbrace{\mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1} | \mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_t | \mathbf{x}_{t-1}) \| p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1}))]}_{\text{consistency term}} \end{aligned} \quad (45)$$

The detailed view of the eELBO can be interpreted

1. $\mathbb{E}_{q(\mathbf{x}_1 | \mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0 | \mathbf{x}_1)]$ can be interpreted as *reconstruction term*, predicting the log probability of the original data sample given the first-step latent. This term also appears in a vanilla VAE, and can be trained similarly.
2. $\mathbb{E}_{q(\mathbf{x}_{T-1} | \mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_{T-1}) \| p(\mathbf{x}_T))]$ is *a prior matching term*; it is minimized when the final latent distribution matches the Gaussian prior. This term requires no optimization, as it has no trainable parameters; furthermore, as we have assumed a large enough T such that the final distribution is Gaussian, this term effectively becomes zero.
3. $\mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1} | \mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_t | \mathbf{x}_{t-1}) \| p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1}))]$ is a *consistency term*; it endeavors to make the distribution at \mathbf{x}_t consistent, from both forward and backward processes. That is, a denoising step from a noisy image should match the corresponding denoising step from a cleaner image, for every intermediate timestep; this is reflected mathematically by the KLD divergence. This term is minimized when we train $p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})$ to match the Gaussian distribution $q(\mathbf{x}_t | \mathbf{x}_{t-1})$, which is defined in Equation 31.

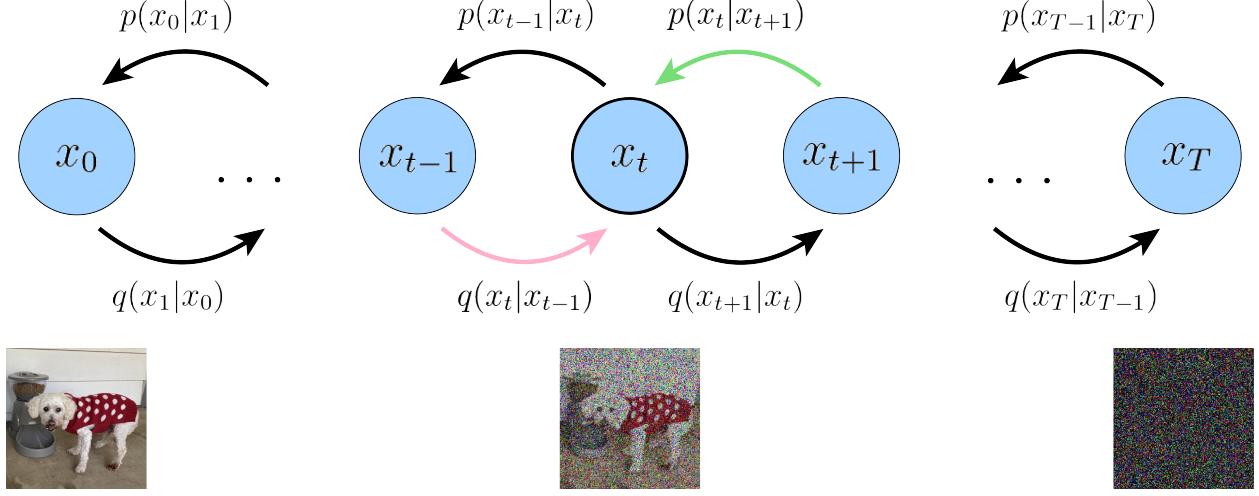


Figure 4: Under our first derivation, a VIM can be optimized by ensuring that for every intermediate posterior or from the latent above $p(\mathbf{x}_t|\mathbf{x}_{t+1})$ matches the Gaussian corruption of the latent before it $q(\mathbf{x}_t|\mathbf{x}_{t-1})$. In this figure, for each intermediate state we minimize the difference between the distributions represented by the pink and green arrows.

Usually, this interpretation of the ELBO is depicted in **4.1. FTlge** of optimizing a VIM is primarily dominated by the third term, since we must optimize over all time steps:

Under this derivation, all terms of the ELBO are computed as expectations, and can therefore be approximated using Monte Carlo estimates. However, actually optimizing the ELBO using these terms we just derived might be suboptimal; because the consistency term is computed as an expectation over two random variables $\{\mathbf{x}_{t-1}, \mathbf{x}_{t+1}\}$ for every time step, the variance of its Monte Carlo estimate could potentially be higher than a term that is estimated using only one random variable per time step. As it is computed by summing up consistency terms, the final estimated value of the ELBO may have higher variance. **Toradure**

Let us instead derive a form for our ELBO where each term is computed as an expectation over only one variable. This is given by the encoder transition $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)$, where the extra condition is superfluous due to the Markov property. Then according to Bayes rule:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) = \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \quad (46)$$

Armed with (46), we can rewrite the derivation using from the ELBO in Eq 37:

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (47)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} \right] \quad (53)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_T|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (54)$$

$$= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t, \mathbf{x}_{t-1}|\mathbf{x}_0)} \left[\log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (57)$$

$$= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}} \quad (58)$$

We have successfully derived a interpretation of the eELBO that can be estimated with lower variance than the standard ELBO. The term can be approximated and optimized using a Monte Carlo estimate.

- $\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]$ can be interpreted as a reconstruction term; it is analogous in the ELBO of a variational VAE, this term can be approximated and optimized using a Monte Carlo estimate.
- $D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))$ represents how close the distribution of the final noisy input is to the standard Gaussian prior or the latent parameters and = 0 under our assumptions.
- $\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]$ is a denoising matching term. We learn desired denoising transitions step $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ as an approximation to tractable, ground-truth denoising transition step $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$. The $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ transition step can act as a ground-truth signal, since it defines how to denoise noisy images with access to what the final, completely denoised image should be. This term is therefore minimized when the two denoising steps match as closely as possible, as measured by the KL divergence.

As we denote, one of the steps in the process of both ELBO derivation (Equations 15 and 58), only the Markov assumption is used; as a result this formula is only true for a very simple Markovian HVAE.

When $T=1$, both of the ELBO interpretations for a VDM exactly recreate the ELBO equation from a variational VAE, as in Eq.(19).

In the derivation of the eELBO, the bulk of the optimization cost once again lies in the summand terms, which dominate the reconstruction term. Whereas each $D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))$ is difficult to minimize for arbitrary posteriors in arbitrarily complex MHVAEs due to the added complexity of simultaneously learning the encoder, in a VDM we can leverage the Gaussian transition assumption to make optimization tractable

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

As we already know that $q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) = q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1-\alpha_t)\mathbf{I})$ from our assumption regarding encoder transitions (31), what remains is deriving the forms of $q(\mathbf{x}_t|\mathbf{x}_0)$ and $q(\mathbf{x}_{t-1}|\mathbf{x}_0)$.

Fortunately, there are also tractable ways utilizing the fact that the encoder transitions of a VDM are linear Gaussian models. Recall that under the reparameterization trick, samples $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_{t-1})$ can be:

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1-\alpha_t}\boldsymbol{\epsilon} \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I}) \quad (59)$$

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1-\alpha_{t-1}}\boldsymbol{\epsilon} \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I}) \quad (60)$$

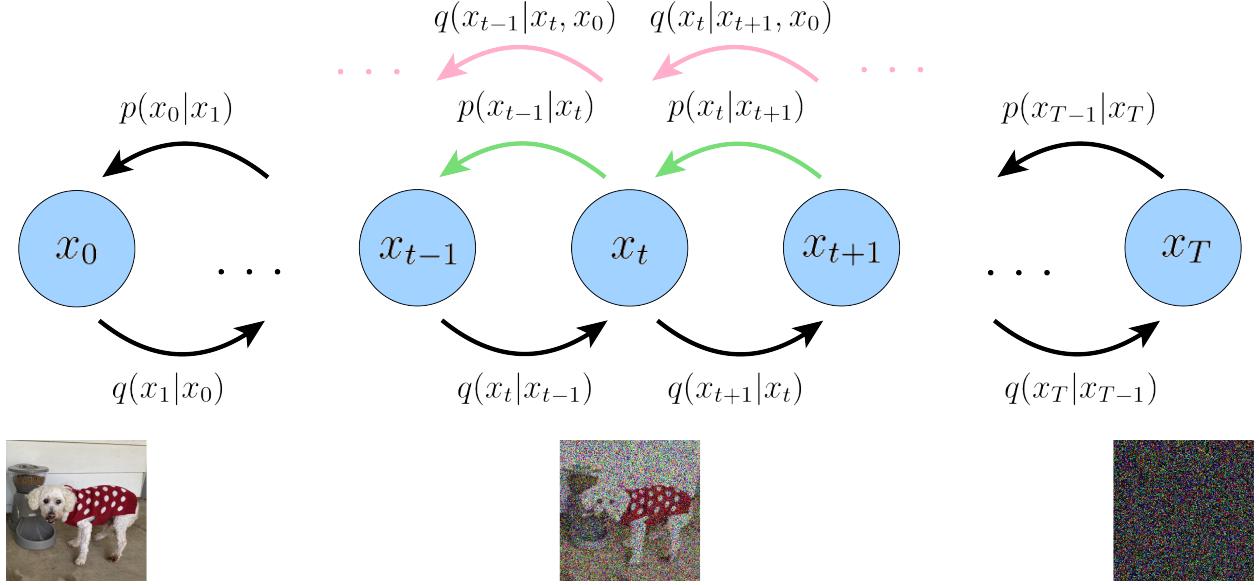


Figure 5: Detailed is an alternate, lower-variance method to optimize a VIM we compute the form of ground-truth denoising step $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ using Bayes rule, and minimize it's KL Divergence with our approximate denoising step $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$. This is since again denoted visually by matching the distributions represented by the green arrows with those of the pink arrows. Artificial liberty is at play here; in the full picture, each pink arrow must also stem from, as it is also a conditioning term.

Then, the form of $q(\mathbf{x}_t|\mathbf{x}_0)$ can be recursively derived through repeated applications of the reparameterization trick. Suppose that we have access to random noise variables $\{\mathbf{e}_t^*, \mathbf{e}_t\}_{t=0}^T \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}; \mathbf{I})$. Then, for an arbitrary sample $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0)$, we can rewrite it as:

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \mathbf{e}_t^* \quad (61)$$

$$= \sqrt{\alpha_t} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \mathbf{e}_{t-2}^* \right) + \sqrt{1 - \alpha_t} \mathbf{e}_t^* \quad (62)$$

$$(63)$$

$$= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \mathbf{e}_{t-2} \quad (64)$$

$$= \dots \quad (65)$$

$$= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \mathbf{e}_0 \quad (66)$$

$$\sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \quad (67)$$

$$(68)$$

$$(69)$$

$$(70)$$

where in Equation 64 we have utilized the fact that ~~the sum of two independent Gaussian random variables remains a Gaussian with mean being the sum of the two means, and variance being the sum of the two variances.~~ Interpreting $\sqrt{1 - \alpha_t} \mathbf{e}_t^*$ as a sample from $\mathcal{N}(\mathbf{0}, (1 - \alpha_t) \mathbf{I})$, and $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \mathbf{e}_{t-2}^*$ as a sample from $\mathcal{N}(\mathbf{0}, (\alpha_t - \alpha_t \alpha_{t-1}) \mathbf{I})$, we can then treat their sum as a random variable sampled from $\mathcal{N}(\mathbf{0}, (1 - \alpha_t + \alpha_t - \alpha_t \alpha_{t-1}) \mathbf{I}) = \mathcal{N}(\mathbf{0}, (1 - \alpha_t \alpha_{t-1}) \mathbf{I})$. A sample from this distribution can then be represented using the reparameterization trick $\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \mathbf{e}_0$, as in Equation 66.

We have therefore derived the Gaussian form $q(\mathbf{x}_t | \mathbf{x}_0)$. This derivation can be modified to also yield the Gaussian parameterization described $q(\mathbf{x}_t | \mathbf{x}_0)$. Now, knowing the forms of both $q(\mathbf{x}_t | \mathbf{x}_0)$ and $q(\mathbf{x}_{t-1} | \mathbf{x}_0)$, we can proceed to calculate the form of $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ by substituting it into the Bayes rule expansion:

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1} | \mathbf{x}_0)}{q(\mathbf{x}_t | \mathbf{x}_0)} \quad (71)$$

$$\propto \exp \left\{ - \left[\frac{(\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{x}_{t-1})^2}{2(1 - \alpha_t)} + \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0)^2}{2(1 - \bar{\alpha}_{t-1})} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0)^2}{2(1 - \bar{\alpha}_t)} \right] \right\} \quad (73)$$

$$\propto \exp \left\{ - \frac{1}{2} \left[- \frac{2\sqrt{\alpha_t} \mathbf{x}_t \mathbf{x}_{t-1}}{1 - \alpha_t} + \frac{\alpha_t \mathbf{x}_{t-1}^2}{1 - \alpha_t} + \frac{\mathbf{x}_{t-1}^2}{1 - \bar{\alpha}_{t-1}} - \frac{2\sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_{t-1} \mathbf{x}_0}{1 - \bar{\alpha}_{t-1}} \right] \right\} \quad (76)$$

$$= \exp \left\{ - \frac{1}{2} \left[\frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \mathbf{x}_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t} \mathbf{x}_t}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0}{1 - \bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \quad (80)$$

$$= \exp \left\{ - \frac{1}{2} \left(\frac{1}{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}} \right) \left[\mathbf{x}_{t-1}^2 - 2 \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1}) \mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t) \mathbf{x}_0}{1 - \bar{\alpha}_t} \mathbf{x}_{t-1} \right] \right\} \quad (83)$$

$$\propto \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1}) \mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t) \mathbf{x}_0}{1 - \bar{\alpha}_t}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I}}_{\Sigma_q(t)}) \quad (84)$$

where in Equation 55, $C(\mathbf{x}_t, \mathbf{x}_0)$ is a constant term with respect to \mathbf{x}_{t-1} computed as a combination of only \mathbf{x}_t , \mathbf{x}_0 , and α values; this term is implicitly returned in Eq 84 to complete the square.

We have therefore shown that at each step $\mathbf{x}_t \sim q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ is normally distributed, with mean $\mu_q(\mathbf{x}_t, \mathbf{x}_0)$ that is a function of \mathbf{x}_t and \mathbf{x}_0 , and variance $\Sigma_q(t)$ as a function of coefficients. These α coefficients are known and fixed at each timestep; they are either set permanently when modeled as hyperparameters, or treated as the current inference output of a network that seeks to model them. Following Equation 84, we can rewrite our variance equation $\Sigma_q(t) = \sigma_q^2(t) \mathbf{I}$, where:

$$\sigma_q^2(t) = \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \quad (85)$$

In order to match approximate denoising transition operator $q(\mathbf{x}_t | \mathbf{x}_0)$ to ground-truth denoising transition step $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ as closely as possible, we can also model it as a Gaussian. Furthermore, α terms are known to be frozen at each timestep, we can immediately construct the variance of the approximate denoising transition step to also $\Sigma_q(t) = \sigma_q^2(t) \mathbf{I}$. We must parameterize its mean $\mu_q(\mathbf{x}_t, t)$ as a function of \mathbf{x}_t , however, since $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$ does not condition on \mathbf{x}_0 .

Recall KL D vergence b etween two Gaussi anli stri b uti ons

$$D_{\text{KL}}(\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x) \| \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)) = \frac{1}{2} \left[\log \frac{|\boldsymbol{\Sigma}_y|}{|\boldsymbol{\Sigma}_x|} - d + \text{tr}(\boldsymbol{\Sigma}_y^{-1} \boldsymbol{\Sigma}_x) + (\boldsymbol{\mu}_y - \boldsymbol{\mu}_x)^T \boldsymbol{\Sigma}_y^{-1} (\boldsymbol{\mu}_y - \boldsymbol{\mu}_x) \right] \quad (86)$$

==>

$$\begin{aligned} & \arg \min_{\boldsymbol{\theta}} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} | \mathbf{x}_t)) \\ &= \arg \min_{\boldsymbol{\theta}} D_{\text{KL}}(\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \| \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_q(t))) \end{aligned} \quad (87)$$

$$= \arg \min_{\boldsymbol{\theta}} \frac{1}{2} [(\boldsymbol{\mu}_{\boldsymbol{\theta}} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1} (\boldsymbol{\mu}_{\boldsymbol{\theta}} - \boldsymbol{\mu}_q)] \quad (90)$$

$$= \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} [\|\boldsymbol{\mu}_{\boldsymbol{\theta}} - \boldsymbol{\mu}_q\|_2^2] \quad (92)$$

where we have written $\boldsymbol{\mu}_q$ as shorthand and $\boldsymbol{\mu}_{\boldsymbol{\theta}}$ as shorthand for brevity. In other words, we want to optimise $\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)$ that matches $\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0)$, which from our derived Equation 84 on takes the form:

$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t} \quad (93)$$

As $\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)$ also conditions on \mathbf{x}_t , we can match $\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0)$ closely by setting it to the following form:

$$\boxed{\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{\mathbf{x}}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t}} \quad (94)$$

where $\hat{\mathbf{x}}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)$ is parameterised by a neural network that seeks to predict from noisy image and time index. Then, the optimisation problem simplifies to:

$$\begin{aligned} & \arg \min_{\boldsymbol{\theta}} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} | \mathbf{x}_t)) \\ &= \arg \min_{\boldsymbol{\theta}} D_{\text{KL}}(\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \| \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_q(t))) \end{aligned} \quad (95)$$

$$= \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{\mathbf{x}}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t} - \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t} \right\|_2^2 \right] \quad (96)$$

$$= \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)^2} [\|\hat{\mathbf{x}}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \mathbf{x}_0\|_2^2] \quad (99)$$

Therefore, optimising a VIMbois down to learning a neural network to predict the original ground truth image from an arbitrary noisy version. Furthermore, minimising the summation term of our derived ELBO objective (Equation 58) across all noise levels can be approximated by minimising the expectation over all time steps:

$$\boxed{\arg \min_{\boldsymbol{\theta}} \mathbb{E}_{t \sim U\{2, T\}} [\mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} | \mathbf{x}_t))]]} \quad (100)$$

which can then be optimised using stochastic samples over time steps.

Learning Diffusion Noise Parameters

Let us investigate how the noise parameters of a VIM can be jointly learned. One potential approach is to model α_t using a neural network $\hat{x}_\theta(t)$ with parameters θ . However, this is inefficient as inference must be performed multiple times at each time step to compute $\bar{\alpha}_t$. Whereas each computation incurs computational cost, we can also derive an alternative way to learn the diffusion noise parameters. By substituting our variance equation from Equation 85 into our derived per-step objective in Equation 99, we can reduce:

$$\frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}(1-\alpha_t)^2}{(1-\bar{\alpha}_t)^2} [\|\hat{x}_\theta(\mathbf{x}_t, t) - \mathbf{x}_0\|_2^2] = \frac{1}{2} \frac{\bar{\alpha}_{t-1}(1-\alpha_t)^2}{(1-\bar{\alpha}_t)^2} [\|\hat{x}_\theta(\mathbf{x}_t, t) - \mathbf{x}_0\|_2^2] \quad (101)$$

$$= \frac{1}{2} \left(\frac{\bar{\alpha}_{t-1}}{1-\bar{\alpha}_{t-1}} - \frac{\bar{\alpha}_t}{1-\bar{\alpha}_t} \right) [\|\hat{x}_\theta(\mathbf{x}_t, t) - \mathbf{x}_0\|_2^2] \quad (108)$$

$\text{SNR} = \frac{\mu^2}{\sigma^2}$, we can write the eSNR at each time step t :

$$\text{SNR}(t) = \frac{\bar{\alpha}_t}{1-\bar{\alpha}_t} \quad (109)$$

Then Eq 108 (and 99) ==>

Precise(t-1) - Precise(t)

$$\frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}(1-\alpha_t)^2}{(1-\bar{\alpha}_t)^2} [\|\hat{x}_\theta(\mathbf{x}_t, t) - \mathbf{x}_0\|_2^2] = \frac{1}{2} (\text{Precise}(t-1) - \text{Precise}(t)) [\|\hat{x}_\theta(\mathbf{x}_t, t) - \mathbf{x}_0\|_2^2] \quad (110)$$

As the name implies, the SNR represents the ratio between the original signal and the amount of noise present; a higher SNR represents more signal and a lower SNR represents more noise. In a diffusion model, we require the SNR to monotonically decrease as time increases; this formalizes the notion that perturbed input \mathbf{x}_t becomes increasingly noisy over time, until it becomes identical to a standard Gaussian at

Following the simplification of the objective in Equation 99 directly parameterize the SNR at each time step using a neural network, and learn it jointly along with the diffusion model. As the SNR must monotonically decrease over time, we can represent it as:

$$\text{SNR}(t) = \exp(-\omega_\eta(t)) \quad (111)$$

where $\omega_\eta(t)$ is modeled as a monotonically increasing neural network with parameters η ; results in a monotonically decreasing function, whereas the exponential forces the resulting term to be positive. Note that the objective in Equation 108 is now optimized over η as well. By combining our parameterization of SNR in Equation 111 with our definition of SNR in Equation 109, we can also explicitly derive elegant forms for the value of $\bar{\alpha}_t$ as well as for the value of $-\bar{\alpha}_t$:

$$\frac{\bar{\alpha}_t}{1-\bar{\alpha}_t} = \exp(-\omega_\eta(t)) \quad (112)$$

$$\therefore \bar{\alpha}_t = \text{sigmoid}(-\omega_\eta(t)) \quad (113)$$

$$\therefore 1-\bar{\alpha}_t = \text{sigmoid}(\omega_\eta(t)) \quad (114)$$

These terms are necessary for a variety of computations; for example, during optimization, they are used to create arbitrary noise from input \mathbf{x}_0 using the reparameterization trick, as derived in Equation

Three Equivalent Interpretations

As we previously proved, a VDM can be trained using a neural network to predict the conditional probability $p_{\theta}(\mathbf{x}_t | \mathbf{x}_0)$ given a natural language \mathbf{x}_0 from an arbitrary training set consisting of \mathbf{x}_t and its next. However, \mathbf{x}_0 has two other equivalent parameterizations, which lead to two further interpretations of a VDM.

1, we can notice that the parameterization in equation (69) implies

$$\mathbf{x}_0 = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_0}{\sqrt{\bar{\alpha}_t}} \quad (115)$$

Plugging this into our previous derivation yields the true denoising transition $\mu_q(\mathbf{x}_t | \mathbf{x}_0)$, which can be derived as:

$$\mu_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t} \quad (116)$$

$$= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t) \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_0}{\sqrt{\bar{\alpha}_t}}}{1 - \bar{\alpha}_t} \quad (117)$$

$$= \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \boldsymbol{\epsilon}_0 \quad (119)$$

2, we can set our approximated denoising transition mean $\mu_{\theta}(\mathbf{x}_t, t)$ as:

$$\mu_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t) \quad (125)$$

and the corresponding optimization problem becomes:

$$\begin{aligned} & \arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) \\ &= \arg \min_{\theta} D_{\text{KL}}(\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \| \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}, \boldsymbol{\Sigma}_q(t))) \end{aligned} \quad (126)$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)\alpha_t} [\|\boldsymbol{\epsilon}_0 - \hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t)\|_2^2] \quad (130)$$

Here, $\hat{\epsilon}_{\theta}(\mathbf{x}_t, t)$ is a neural network that learns to predict the source $\text{noiseN}(\epsilon; \mathbf{0}, \mathbf{I})$ that determines \mathbf{x}_t from \mathbf{x}_0 . We have already seen that a VIM by predicting the original image \mathbf{x}_0 is equivalent to learning to predict the noise ϵ itself. However, some works have found that predicting the noise has resulted in better performance [5, 7].

3 common interpretation of VIMs, we appeal to Tweedie's Formula [8]. In English, Tweedie's Formula states that the true mean of an EFD, given samples drawn from it can be estimated by the MLE of the samples (assuming the samples are drawn from a normal distribution) plus some correction terms involving the score of the estimate in the case of justified sampling. The empirical mean is just the sample itself, commonly used to estimate the sample bias, which is often served as a simple one-end-of-the-underlying-distribution shift that reflects the score of the sample and corrects the main MLE of the samples towards the true mean.

for $\mathbf{z} \sim \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$, Tweedie's Formula:

$$\mathbb{E}[\boldsymbol{\mu}_z | \mathbf{z}] = \mathbf{z} + \boldsymbol{\Sigma}_z \nabla_{\mathbf{z}} \log p(\mathbf{z})$$

apply it to predict the true posterior mean of \mathbf{x}_t given test samples. From Eq70, we know

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

Then Tweedie's Formula \Rightarrow

$$\mathbb{E}[\boldsymbol{\mu}_{x_t} | \mathbf{x}_t] = \mathbf{x}_t + (1 - \bar{\alpha}_t) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \quad (131)$$

According to Tweedie's Formula, the

best estimate for the true mean that \mathbf{x}_t is generated from, $\boldsymbol{\mu}_{x_t} = \sqrt{\bar{\alpha}_t} \mathbf{x}_0$, is defined as:

$$\sqrt{\bar{\alpha}_t} \mathbf{x}_0 = \mathbf{x}_t + (1 - \bar{\alpha}_t) \nabla \log p(\mathbf{x}_t) \quad (132)$$

$$\therefore \mathbf{x}_0 = \frac{\mathbf{x}_t + (1 - \bar{\alpha}_t) \nabla \log p(\mathbf{x}_t)}{\sqrt{\bar{\alpha}_t}} \quad (133)$$

Then we can plug Eq133 into our ground-truth denoising transformation mean $\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0)$ once again

$$\Rightarrow \boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t} \quad (134)$$

$$= \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t) \frac{\mathbf{x}_t + (1 - \bar{\alpha}_t) \nabla \log p(\mathbf{x}_t)}{\sqrt{\bar{\alpha}_t}}}{1 - \bar{\alpha}_t} \quad (135)$$

$$= \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \nabla \log p(\mathbf{x}_t) \quad (137)$$

\Rightarrow approximating denoising mean

$$\mu_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} s_{\theta}(\mathbf{x}_t, t) \quad (143)$$

\Rightarrow

$$\begin{aligned} & \arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) \\ & \quad (144) \end{aligned}$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} s_{\theta}(\mathbf{x}_t, t) - \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \nabla \log p(\mathbf{x}_t) \right\|_2^2 \right] \quad (145)$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{\alpha_t} \left[\|s_{\theta}(\mathbf{x}_t, t) - \nabla \log p(\mathbf{x}_t)\|_2^2 \right] \quad (148)$$

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

The astute reader will notice that the score function $\nabla \log p(\mathbf{x}_t)$ looks remarkably similar in form to the source noise ϵ_0 . This is an interesting coincidence combining Tweedie's Formula (133) with the reparameterization trick (115):

$$\mathbf{x}_0 = \frac{\mathbf{x}_t + (1 - \bar{\alpha}_t) \nabla \log p(\mathbf{x}_t)}{\sqrt{\bar{\alpha}_t}} = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_0}{\sqrt{\bar{\alpha}_t}} \quad (149)$$

$$\therefore (1 - \bar{\alpha}_t) \nabla \log p(\mathbf{x}_t) = -\sqrt{1 - \bar{\alpha}_t} \epsilon_0 \quad (150)$$

$$\nabla \log p(\mathbf{x}_t) = -\frac{1}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_0 \quad (151)$$

As it turns out, the two terms are off by a constant factor that scales with time! The score function measures how to move in data space to maximize the log probability; intuitively, since the source noise is added to a natural image to corrupt it, moving in its opposite direction "denoises" the image and would be the best update to increase the subsequent log probability. Our mathematical proof justifies this intuition; we have explicitly shown that learning to model the score function is equivalent to modeling the negative of the source noise (up to a scaling factor).

We have therefore derived three equivalent objectives to optimize a VIM learning a neural network to predict the original image \mathbf{x}_0 , the source noise ϵ_0 , or the score of the image at an arbitrary level $\nabla \log p(\mathbf{x}_t)$. The VIM can be easily trained by stochastic sampling time steps t and minimizing the norm of the prediction against the ground truth target.

Score-based Generative Models

We have shown that VIM can be learned simply by optimizing a neural network $s_{\theta}(\mathbf{x}_t, t)$ to predict the score function $\nabla \log p(\mathbf{x}_t)$. However, in our derivation of the score term arrived from an application of Tweedie's Formula; this doesn't necessarily hold for the general case. Fortunately, we can look to another class of generative models, SBGMs [9, 10, 11], for exactly this reason. As it turns out, we can show that the VIM formula is equivalent to the SBGM formula on average, provided that the two interpretations are consistent with each other.

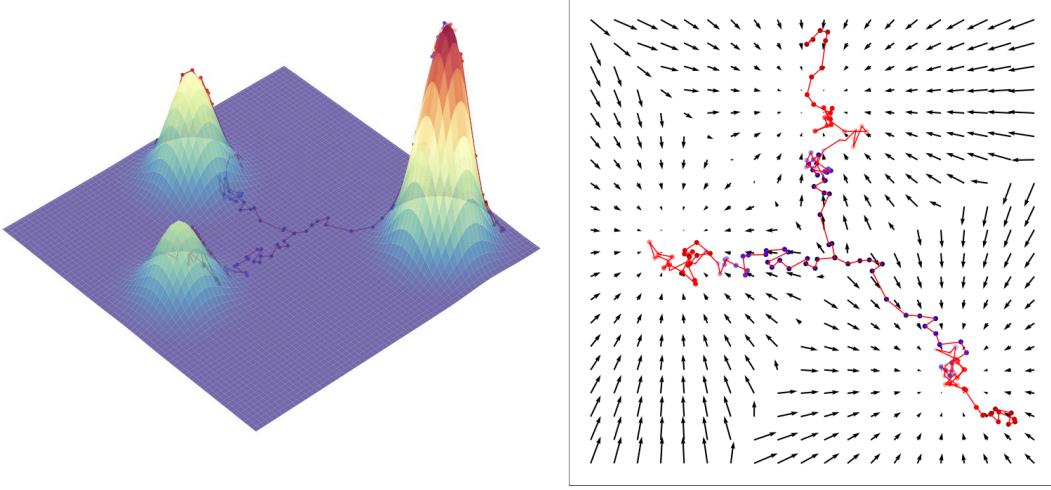


Figure 6: Visualizations of generated samples from MoGs. The left figure shows a 3D contour plot of a bimodal distribution with a blue dashed line indicating a trajectory. The right figure shows a 2D vector field plot with a grid of arrows representing gradients, and a red dashed line indicating a trajectory starting from a point and moving towards one of the modes.

To better understand why optimizability makes sense, we take a detour and revisit energy-based models [12, 13]:

$$p_{\theta}(\mathbf{x}) = \frac{1}{Z_{\theta}} e^{-f_{\theta}(\mathbf{x})} \quad (152)$$

use a neural network $s_{\theta}(\mathbf{x})$ to learn the score function $\nabla \log p(\mathbf{x})$:

$$\begin{aligned} \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) &= -\nabla_{\mathbf{x}} f_{\theta}(\mathbf{x}) \\ &\approx s_{\theta}(\mathbf{x}) \end{aligned} \quad (156)$$

which can be freely represented as a neural network without involving any normalization constants. The score model can be optimized directly using the Divergence with the ground truth score function

$$\mathbb{E}_{p(\mathbf{x})} \left[\|s_{\theta}(\mathbf{x}) - \nabla \log p(\mathbf{x})\|_2^2 \right] \quad (157)$$

What does the score function represent? For every taking the gradient of its logarithm with respect to \mathbf{x} essentially describes what direction in data space to move in order to further increase its likelihood.

Intuitively, then, the score function defines a vector field over the entire space that points towards the modes. Usually, this is depicted in the right plot of Figure 6. Then, by learning the score function of the true data distribution, we can generate samples by starting at any arbitrary point in the same space and iteratively following the score until a mode is reached. This sampling procedure is known as Langevin dynamics, and is mathematically described as:

$$\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + c \nabla \log p(\mathbf{x}_i) + \sqrt{2c\epsilon}, \quad i = 0, 1, \dots, K \quad (158)$$

where \mathbf{x}_0 is randomly sampled from a prior distribution (such as uniform) $\sim \text{Unif}(\mathbf{0}, \mathbf{I})$. It is an extra noise term to ensure that the generated samples do not always collapse onto a mode, but hover around it for diversity. Furthermore, because the learned score function is deterministic, sampling with the noise term involved adds stochasticity to the generation process, allowing us to avoid deterministic trajectories. This is particularly useful when sampling is interested from a position that lies between multiple modes. A visualization of Langevin dynamics sampling and the benefits of the noise term is shown in Figure 6.

Note that the objective in Eq 157 relies on having access to the ground truth score function, which is unavailable for complex distributions such as scene modeling natural images. Fortunately, alternative techniques known as SM [14–17] have been developed to minimize Fish divergence with respect to the ground truth score, and can be optimized with SGB.

Collectively, learning to represent a distribution as a score function and using it to generate samples through MC techniques, such as Langevin dynamics, is known as SbGM [9, 10, 11].

Theorem 2 reexamines problems with vanilla score matching, as detailed by Song and Ermon [9]. Firstly, the score function is ill-defined when it lies on a low-dimensional manifold in a high-dimensional space. This can be seen mathematically; all points not on the low-dimensional manifold would have probability zero, the elong of which is undefined. This is particularly inconvenient when trying to learn a generative model over natural images, which is known to lie on a low-dimensional manifold of the entire ambient space.

Secondly, the estimated score function trained via vanilla score matching will not be accurate in low density regions. This is evident from the objective we minimize in Equation 157. Because it is an expectation over $p(\mathbf{x})$, and explicitly trained on samples from it, the model will not receive an accurate learning signal for rarely seen or unseen examples. This is problematic, since our sampling strategy involves starting from a random location in the high-dimensional space, which is most likely random noise, and moving according to the learned score function. Since we are following a noisy or inaccurate score estimate, the final generated samples may be suboptimal as well, or require many more iterations to converge on an accurate output.

Lastly, LDs sampling may not mix, even if it is performed using the ground truth scores. Suppose that the true data distribution is a mixture:

$$p(\mathbf{x}) = c_1 p_1(\mathbf{x}) + c_2 p_2(\mathbf{x}) \quad (159)$$

The entropy score is computed, the semantically meaningful coefficients are lost, since the operation splits the coefficients from the distribution and the gradient operation zeros it out. To visualize this, note that the ground truth score function shown in the figure 6 is agnostic of the different weights between the two distributions; LD sampling from the depicted initialization point has a rough likelihood arriving at each mode, despite the bottleneck mode having a high weight than the actual MoGs.

It turns out that the retraining procedure successfully addresses multiple levels of Gaussian noise to the data. 1, as the support of a Gaussian noise distribution is the entire space, a perturbation sample will no longer be confined to a low-dimensional manifold. 2, adding large Gaussian noise will increase the mode covers in the distribution, adding more training signal in low density regions. 3, adding multiple levels of Gaussian noise will result in intermediate distributions that respect the ground truth mixing coefficients.

Formally, we can choose a positive sequence of noise levels $\{\sigma_t\}_{t=1}^T$ and define a sequence of progressively perturbed data distributions:

$$p_{\sigma_t}(\mathbf{x}_t) = \int p(\mathbf{x}) \mathcal{N}(\mathbf{x}_t; \mathbf{x}, \sigma_t^2 \mathbf{I}) d\mathbf{x} \quad (160)$$

Then a neural network $s_\theta(\mathbf{x}, t)$ is learned using SM to learn the score function for all noise levels simultaneously:

$$\arg \min_{\theta} \sum_{t=1}^T \lambda(t) \mathbb{E}_{p_{\sigma_t}(\mathbf{x}_t)} \left[\|s_\theta(\mathbf{x}, t) - \nabla \log p_{\sigma_t}(\mathbf{x}_t)\|_2^2 \right] \quad (161)$$

where $\lambda(t)$ is a positive weight function that conditions on the noise level t . Note that this objective almost exactly matches the objective derived in Equation 18 to train a Variational Diffusion Model. Furthermore, the authors propose an annealed Langevin dynamics sampling as a generative procedure, in which samples are produced by running Langevin dynamics for $t = T, T-1, \dots, 2, 1$ in sequence. This initialization is chosen from some fixed prior (such as uniform), and each subsequent sampling step starts from the final samples of the previous simulation. Because the noise levels steadily decrease over time steps, we reduce the step size over time, the samples eventually converge into a true mode. This is directly analogous to the sampling procedure performed in the MHVAE interpretation of a VDM, where a randomly initialized data vector is iteratively refined over decreasing noise levels.

Therefore we have established an explicit connection between VDMs and SbGMs, both in their training objectives and sampling procedures.

One question is how to naturally generalize diffusion models to an infinite number of time steps. Under the MHVAE view, this can be interpreted as extending the number of hierarchies to infinity $T \rightarrow \infty$. It is clearer to represent this from the equivalent SbGM perspective; under an infinite number of noise scales, the perturbation of an image over continuous time can be represented as a stochastic process, and therefore described by a SDE. Sampling is then performed by reversing the SDE, which naturally requires estimating the score function at each continuous-valued noise level [10]. Different parameterizations of the SDE essentially describe different perturbations over time, enabling flexible modeling of the denoising procedure [6].

Guidance

So far, we have focused on modeling just the data distribution $p(\mathbf{x})$. However, we are often also interested in learning conditional distributions $p(\mathbf{x}|y)$, which would enable us to explicitly control the data we generate through conditioning information y . This forms the backbone of image super-resolution models such as Cascaded Diffusion Models [18], as well as SOTA image-text models such as DALL-E 2 [19] andImagen [7].

A natural way to add conditioning information is simply alongside the timestep information, at each iteration. Recall our joint distribution from Eq 32:

$$p(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$$

Then to turn this into a conditional DM, we can simply add arbitrary conditioning information y at each transition step as

$$p(\mathbf{x}_{0:T}|y) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t, y) \quad (162)$$

For example, y could be a text encoding in image-text generation, or a low-resolution image to perform super-resolution. We are thus able to learn the core neural networks of a VDM as before, by predicting $\hat{\mathbf{x}}_\theta(\mathbf{x}_t, t, y) \approx \mathbf{x}_0$, $\hat{\mathbf{e}}_\theta(\mathbf{x}_t, t, y) \approx \mathbf{e}_0$, or $s_\theta(\mathbf{x}_t, t, y) \approx \nabla \log p(\mathbf{x}_t|y)$ for each desired interpretation and implementation.

A caveat of this vanilla formulation is that a conditional diffusion model trained in this way may potentially learn to ignore or downplay any given conditional information. Guidance is therefore proposed as a way to more explicitly control the amount of weight the model gives to the conditional information, at the cost of sample diversity. The two most popular forms of guidance are known as Classifier-Free Guidance [0, 20] and Classifier-Free Guidance [21].

Classifier Guidance

Let us begin with the score-based formulation of a diffusion model, where our goal is $\nabla \log p(\mathbf{x}_t|y)$, the score of the conditional model, at arbitrary noise level. Recall that ∇ is shorthand for $\nabla_{\mathbf{x}_t}$ in the interest of brevity. By Bayes rule, we can derive the following equivalent form:

$$\nabla \log p(\mathbf{x}_t|y) = \nabla \log \left(\frac{p(\mathbf{x}_t)p(y|\mathbf{x}_t)}{p(y)} \right) \quad (163)$$

$$= \nabla \log p(\mathbf{x}_t) + \nabla \log p(y|\mathbf{x}_t) - \nabla \log p(y) \quad (164)$$

$$= \underbrace{\nabla \log p(\mathbf{x}_t)}_{\text{unconditional score}} + \underbrace{\nabla \log p(y|\mathbf{x}_t)}_{\text{adversarial gradient}} \quad (165)$$

where $\text{grad } \log p(y|\mathbf{x}_t) = 0$.

Our final derived result can be interpreted as learning an unconditional score function combined with the adversarial gradient of a classifier $p(y|\mathbf{x}_t)$. Therefore in Classifier Guidance [10, 20], the score of an unconditional diffusion model is learned as previously derived, alongside a classifier that takes in arbitrary noisy \mathbf{x}_t and attempts to predict conditional information y . Then during the sampling procedure, the overall conditional score function used for annealed Langevin dynamics is computed as the sum of the unconditional score function and the adversarial gradient of the noisy classifier.

In order to introduce fine-grained control to either encourage or discourage the model to consider the conditioning information, Classifier Guidance scales the adversarial gradient of the noisy classifier by a hyperparameter γ . The score function learned under Classifier Guidance is:

$$\boxed{\nabla \log p(\mathbf{x}_t|y) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(y|\mathbf{x}_t)} \quad (166)$$

Intuitively, when $\gamma = 0$ the conditional DM learns to ignore the conditional information entirely, and when γ is large the conditional DM learns to produce samples that heavily adhere to the conditional information. This would come at the cost of sample diversity, as it would only produce data that would be easy to regenerate the provided conditional information from, even at noisy levels.

One noted drawback of Classifier Guidance is its reliance on a separately learned classifier. Because the classifier must handle arbitrary noisy inputs, while most existing pretrained classification models are not optimized to do, it must be learned alongside the diffusion model.

Classifier-Free Guidance

In CFG [21], the authors distinguish between the separate classifier model in favor of an unconditional DM and a conditional DM. To derive the score function under CFG, we can first rearrange Eq 165 to show that:

$$\nabla \log p(y|\mathbf{x}_t) = \nabla \log p(\mathbf{x}_t|y) - \nabla \log p(\mathbf{x}_t) \quad (167)$$

Then Eq 166, ==>

$$\nabla \log p(\mathbf{x}_t|y) = \nabla \log p(\mathbf{x}_t) + \gamma (\nabla \log p(\mathbf{x}_t|y) - \nabla \log p(\mathbf{x}_t)) \quad (168)$$

$$= \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{x}_t|y) - \gamma \nabla \log p(\mathbf{x}_t) \quad (169)$$

$$= \underbrace{\gamma \nabla \log p(\mathbf{x}_t|y)}_{\text{conditional score}} + \underbrace{(1 - \gamma) \nabla \log p(\mathbf{x}_t)}_{\text{unconditional score}} \quad (170)$$

Once again, γ is a term that controls how much our learned conditional model cares about the conditional information. When $\gamma = 0$, the learned conditional model completely ignores the conditional information and learns an unconditional diffusion model. When $\gamma > 1$, the model explicitly learns the vanilla conditional distribution without guidance. When $\gamma < 1$, the diffusion model not only prioritizes the conditional score function, but also moves in the direction away from the unconditional score function. In other words, it reduces the probability of generating samples that do not use conditional information, in favor of the samples that explicitly do. This also has the effect of decreasing sample diversity at the cost of generating samples that accurately match the conditional information.

Because learning two separate diffusion models is expensive, we can learn both the conditional and unconditional diffusion models together as a singular conditional model; the unconditional diffusion model can be queried by replacing the conditional information with fixed constant values, such as zeros. This is especially useful for performing random dropout on the conditional information. This fier-Free Guidance is elegant because it enables us greater control over our conditional generation procedure while requiring nothing beyond the training of a singular conditional diffusion model.

Closing

Allow us to recapitulate our findings over the course of our explorations. First, we derive Variational Diffusion Models as a specific case of a Markovian Hierarchical Variational Autoencoder, where three key assumptions enable tractable computation and scalable optimization of the ELBO. We then prove that optimizing a VIM boils down to learning a neural network to predict one of three potential objectives: the original source image from any arbitrary noise distribution or the original source noise from any arbitrary noise level. Then, we dive deeper into what it means to learn the score function, and connect it explicitly with the perspective of Score-based Generation Modeling. Lastly, we cover how to learn a conditional distribution using diffusion models.

In summary, diffusion models have their own unique capabilities as generative models; indeed, they power the current state-of-the-art models on text-conditional image generation such as Imagen and DALL-E 2. Furthermore, the mathematical tools that enable these models are exceedingly elegant. However, there still remain a few drawbacks to consider:

- It is unlikely that this is somehow, as humans, naturally model and generate data; we do not generate samples as random noise that we interpretively denoise.
- The VIM does not produce interpretable latents. Whereas a VAE would hopefully learn a structured latent space through the optimization of its encoder, in a VIM the encoder at each timestep is already given as a learned Gaussian model and cannot be optimized flexibly. Therefore, the intermediate latents are restricted as just noisy versions of the original input.
- The latents are restricted to the same dimensionality as the original input, further frustrating efforts to learn meaningful, compressed latent structure.
- Sampling is an expensive procedure, as multiple denoising steps must be run under both formulations. Recall that one of the restrictions is that a large enough number of timesteps is needed to ensure the final latent is completely Gaussian noise; during sampling we must iterate over all these timesteps to generate a sample.

As a final note, the success of diffusion models highlights the power of Hierarchical VAEs as a generative model. We have shown that when we generalize *infinite* latent hierarchical es, even if the encoder is trivial and the latent dimension is fixed and Markovian transitions are assumed, we are still able to learn powerful models of data. This suggests that further performance gains can be achieved in the case of general, deep HVAEs, where complex encoders and semantically meaningful latent spaces can be potentially learned.