



Overview Diffusion Model

Suk-Hwan Lee

Artificial Intelligence

Creating the Future

Dong-A University

Division of **C**omputer **E**ngineering &
Artificial **I**ntelligence

References

[Technical Blog]

Lil'Long, "What are Diffusion Models?"

<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/#imagen>

[Technical Blog]

Sergios Karagiannakos, Nikolas Adaloglou, "How diffusion models work: the math from scratch"

<https://theaisummer.com/diffusion-models/>

▼ Table of Contents

- What are Diffusion Models?
 - Forward diffusion process
 - Connection with stochastic gradient Langevin dynamics
 - Reverse diffusion process
 - Parameterization of L_t for Training Loss
 - Simplification
 - Connection with noise-conditioned score networks (NCSN)
 - Parameterization of β_t
 - Parameterization of reverse process variance Σ_θ
- Conditioned Generation
 - Classifier Guided Diffusion
 - Classifier-Free Guidance

- Speed up Diffusion Models
 - Fewer Sampling Steps & Distillation
 - Latent Variable Space
- Scale up Generation Resolution and Quality
- Model Architecture
- Quick Summary
- Citation
- References

[Updated on 2021-09-19: Highly recommend this blog post on [score-based generative modeling](#) by Yang Song (author of several key papers in the references)].

[Updated on 2022-08-27: Added [classifier-free guidance](#), [GLIDE](#), [unCLIP](#) and [Imagen](#).

[Updated on 2022-08-31: Added [latent diffusion model](#).

[Updated on 2024-04-13: Added [progressive distillation](#), [consistency models](#), and the [Model Architecture](#) section.

References

[Technical Blog]

Lil'Long, "What are Diffusion Models?"

<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/#imagen>

- [1] Jascha Sohl-Dickstein et al. "[Deep Unsupervised Learning using Nonequilibrium Thermodynamics](#)." ICML 2015.
- [2] Max Welling & Yee Whye Teh. "**Bayesian learning via stochastic gradient langevin dynamics**." ICML 2011.
- [3] Yang Song & Stefano Ermon. "**Generative modeling by estimating gradients of the data distribution**." NeurIPS 2019.
- [4] Yang Song & Stefano Ermon. "**Improved techniques for training score-based generative models**." NeurIPS 2020.
- [5] Jonathan Ho et al. "**Denoising diffusion probabilistic models**." arxiv Preprint arxiv:2006.11239 (2020). [code]
- [6] Jiaming Song et al. "**Denoising diffusion implicit models**." arxiv Preprint arxiv:2010.02502 (2020). [code]
- [7] Alex Nichol & Prafulla Dhariwal. "**Improved denoising diffusion probabilistic models**" arxiv Preprint arxiv:2102.09672 (2021). [code]

[Technical Blog]

Sergios Karagiannakos, Nikolas Adaloglou, "How diffusion models work: the math from scratch"

<https://theaisummer.com/diffusion-models/>

- [8] Prafulla Dhariwal & Alex Nichol. "**Diffusion Models Beat GANs on Image Synthesis**." arxiv Preprint arxiv:2105.05233 (2021). [code]
- [9] Jonathan Ho & Tim Salimans. "**Classifier-Free Diffusion Guidance**." NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications.
- [10] Yang Song, et al. "**Score-Based Generative Modeling through Stochastic Differential Equations**." ICLR 2021.
- [11] Alex Nichol, Prafulla Dhariwal & Aditya Ramesh, et al. "**GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models**." ICML 2022.
- [12] Jonathan Ho, et al. "**Cascaded diffusion models for high fidelity image generation**." J. Mach. Learn. Res. 23 (2022): 47-1.
- [13] Aditya Ramesh et al. "**Hierarchical Text-Conditional Image Generation with CLIP Latents**." arxiv Preprint arxiv:2204.06125 (2022).

References

[Technical Blog]

Lil'Long, "What are Diffusion Models?"

<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/#imagen>

[Technical Blog]

Sergios Karagiannakos, Nikolas Adaloglou, "How diffusion models work: the math from scratch"

<https://theaisummer.com/diffusion-models/>

[14] Chitwan Saharia & William Chan, et al. "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding." arxiv Preprint arxiv:2205.11487 (2022).

[15] Rombach & Blattmann, et al. "**High-Resolution Image Synthesis with Latent Diffusion Models.**" CVPR 2022.code

[16] Song et al. "**Consistency Models**" arxiv Preprint arxiv:2303.01469 (2023)

[17] Salimans & Ho. "**Progressive Distillation for Fast Sampling of Diffusion Models**" ICLR 2022.

[18] Ronneberger, et al. "U-Net: Convolutional Networks for Biomedical Image Segmentation" MICCAI 2015.

[19] Peebles & Xie. "**Scalable diffusion models with transformers.**" ICCV 2023.

[20] Zhang et al. "**Adding Conditional Control to Text-to-Image Diffusion Models.**" arxiv Preprint arxiv:2302.05543 (2023).

What are Diffusion Models?

- Generative models; GAN, VAE, and Flow-based models. - Great success in generating high-quality samples
- Limitations*
 - GAN** models : Potentially unstable training and less diversity in generation due to their adversarial training nature.
 - VAE** : Relies on a surrogate loss (대리손실함수 - ELBO).
 - Flow models** : Have to use specialized architectures to construct reversible transform
- Diffusion models**
 - Inspired by [non-equilibrium thermodynamics](#).
 - Define a **Markov chain of diffusion steps to slowly add random noise to data** and then **learn to reverse the diffusion process to construct desired data samples from the noise**.
 - Unlike VAE or flow models, diffusion models are **learned with a fixed procedure and the latent variable has high dimensionality** (same as the original data).

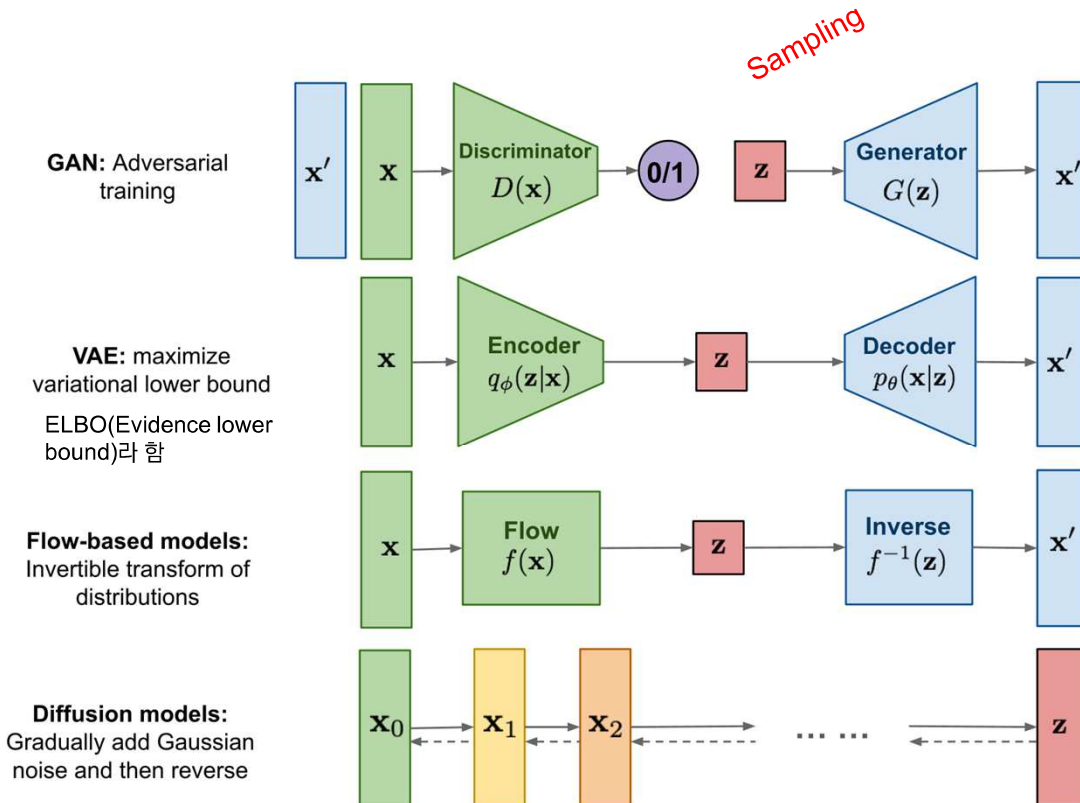


Fig. 1. Overview of different types of generative models.

- Several diffusion-based generative models have been proposed with similar ideas underneath, including [diffusion probabilistic models \(Sohl-Dickstein et al., 2015\)](#), [noise-conditioned score network \(NCSN; Yang & Ermon, 2019\)](#), and [denoising diffusion probabilistic models \(DDPM; Ho et al. 2020\)](#).

What are Diffusion Models?

1. Forward diffusion process

- Given a data point sampled from a real data distribution $\mathbf{x}_0 \sim q(\mathbf{x})$, let us define a *forward diffusion process* in which we add small amount of Gaussian noise to the sample in T steps, producing a sequence of noisy samples $\mathbf{x}_1, \dots, \mathbf{x}_T$. The step sizes are controlled by a variance schedule $\{\beta_t \in (0,1)\}_{t=1}^T$.

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \quad q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

- The data sample \mathbf{x}_0 gradually loses its distinguishable features as the step t becomes larger. Eventually when $T \rightarrow \infty$, \mathbf{x}_T is equivalent to an isotropic Gaussian distribution.

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

- Center of the Noise** $\mu_t = (1 - \beta_t)\mathbf{x}_{t-1}$: Slightly shifting the image by a factor that depends on β_t , with β_t controlling how much we deviate from the original.
- Spread of the Noise** $\Sigma_t = \beta_t\mathbf{I}$ (\mathbf{I} is the identity matrix) : The noise is added independently to each pixel or feature of the image. The term β_t scales this noise.

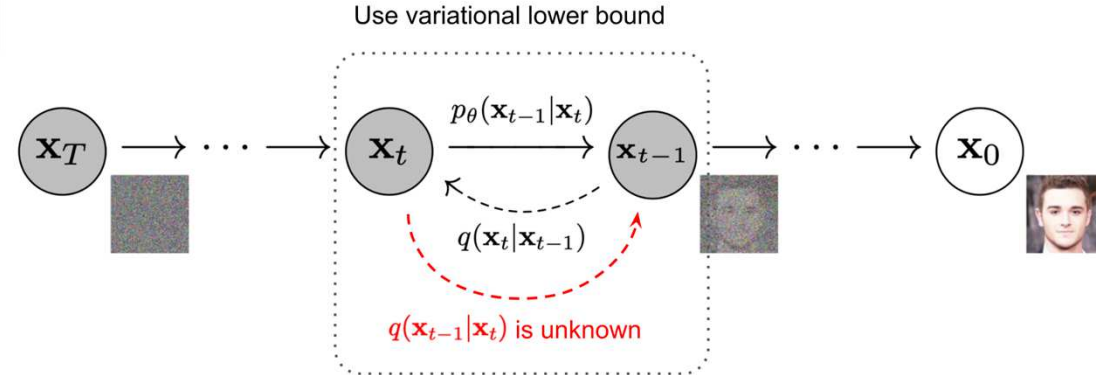


Fig. 2. The Markov chain of forward (reverse) diffusion process of generating a sample by slowly adding (removing) noise. (Image source: Ho et al. 2020 with a few additional annotations)

➤ Iterative Refinement

- At each step t , the image is slightly shifted (according to the mean) and then noise is added around this shifted version (according to the variance). Over multiple steps, this leads to an image that's increasingly noisy. But remember, **it's not random noise; it's structured and based on the original image.**
- The power of the diffusion model lies in **reversing this process. By training on this forward noise process, the model learns to do the reverse: starting from a noisy image, it iteratively refines and 'denoises' it to produce a clear, coherent sample.**

What are Diffusion Models?

1. Forward diffusion process

➤ **Question:** Why do we need the scaling factor $(1 - \beta_t)$ to shift the mean?

- The mean deviation from the image in the diffusion process is essential for achieving the generative capability of the model.

$$\mu_t = (1 - \beta_t)x_{t-1}, \quad \Sigma_t = \beta_t I$$

- The mean, μ_t , in the diffusion process, represents **where we expect our image to be after adding noise at time t .**
- The term $(1 - \beta_t)$ acts as a **tether**, ensuring that our image **doesn't drift too far from its starting point**. However, as β_t increases, our "tether" allows for more slack, meaning the image can drift slightly.
- By controlling this drift using β_t , we ensure that our changes, while making the *vehicle (original)* image different, still retain its essence and are recognized as a *vehicle (original)*.
- Over time, **this controlled deviation** allows the **model** to **explore and generate a diverse set of images**

❖ Introduction to Diffusion Models (Part III. Diffusion Process) : <https://scalexi.medium.com/introduction-to-diffusion-models-part-iii-diffusion-process-cf18bdd36cc4>

➤ Setting the Stage: Multi-dimensional Scenario

- When working in a multi-dimensional space, I denotes the identity matrix, which, in this context, implies each dimension has a consistent standard deviation β_t .
- $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ is a conditional normal distribution described by the mean μ_t and variance Σ_t

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

What are Diffusion Models?

1. Forward diffusion process

➤ Tractability of Progression

- Starting from our initial data \mathbf{x}_0 and progressing to \mathbf{x}_T

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

- The colon symbol (:) means that we are applying the function q across all timesteps from 1 to T .

❖ A challenge arises !!!

- If we wanted to sample \mathbf{x}_t for $t=500$ (before reaching T), we'd have to apply the function 500 times. That's inefficient!

- ❖ Introduction to Diffusion Models (Part III. Diffusion Process) : <https://scalexi.medium.com/introduction-to-diffusion-models-part-iii-diffusion-process-cf18bdd36cc4>

➤ The Reparametrization Trick

- The reparametrization trick simplifies this process.
- Instead of working with β_t , we introduce $\alpha_t = 1 - \beta_t$ and its cumulative product up to time $\bar{\alpha}_t$.
- This trick allows us to express \mathbf{x}_t in terms of the initial point \mathbf{x}_0 and some noise ϵ .

$$x_t \sim q(x_t|x_0) = \mathcal{N}(\bar{\alpha}_t x_0, (1 - \bar{\alpha}_t)I) \quad \Rightarrow \text{Next page}$$

- With β_t pre-defined, we can calculate α_t and $\bar{\alpha}_t$ for all timesteps beforehand. Now, to obtain a sample \mathbf{x}_t at any timestep, just sample the noise ϵ and plug it into the equation above.
- Rather than directly sampling from a distribution, we use another route to achieve our goal for **the efficient computation**.

What are Diffusion Models?

1. Forward diffusion process

A nice property of the above process is that we can sample \mathbf{x}_t at any arbitrary time step t in a closed form using **reparameterization trick**. Let $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$:

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \epsilon_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \bar{\epsilon}_{t-2} \quad \text{where } \bar{\epsilon}_{t-2} \text{ merges two Gaussians (*).} \\ &= \dots \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon\end{aligned}$$

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

(*) Recall that when we merge two Gaussians with different variance, $N(\mathbf{0}, \sigma_1^2 \mathbf{I})$ and $N(\mathbf{0}, \sigma_2^2 \mathbf{I})$, the new distribution is $N(\mathbf{0}, (\sigma_1^2 + \sigma_2^2) \mathbf{I})$. Here the merged standard deviation is $\sqrt{(1 - \alpha_t) + \alpha_t(1 - \alpha_{t-1})} = \sqrt{1 - \alpha_t \alpha_{t-1}}$

Usually, we can afford a larger update step when the sample gets noisier, so $\beta_1 < \beta_2 < \dots < \beta_T$ and therefore $\bar{\alpha}_1 > \dots > \bar{\alpha}_T$

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

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

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} (\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \epsilon_{t-2}) + \sqrt{1 - \alpha_t} \epsilon_{t-1} \\ &= \sqrt{\alpha_t} \sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \underbrace{\sqrt{\alpha_t (1 - \alpha_{t-1})} \epsilon_{t-2}}_{N(0, \alpha_t (1 - \alpha_{t-1}) \mathbf{I})} + \underbrace{\sqrt{1 - \alpha_t} \epsilon_{t-1}}_{N(0, (1 - \alpha_t) \mathbf{I})}\end{aligned}$$

$$N(0, \alpha_t (1 - \alpha_{t-1}) \mathbf{I}) \quad N(0, (1 - \alpha_t) \mathbf{I})$$

Sum of two gaussian distributions with different variance

$$\begin{aligned}N(0, (\alpha_t (1 - \alpha_{t-1}) + (1 - \alpha_t)) \mathbf{I}) \\ = N(0, (1 - \alpha_t \alpha_{t-1}) \mathbf{I})\end{aligned}$$

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

$$\mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0) = N(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \quad \bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

What are Diffusion Models?

1. Forward diffusion process



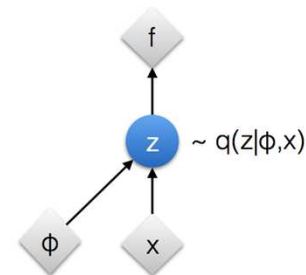
Memo

[Reparameterization trick : VAE](https://lilianweng.github.io/posts/2018-08-12-vae/#reparameterization-trick)

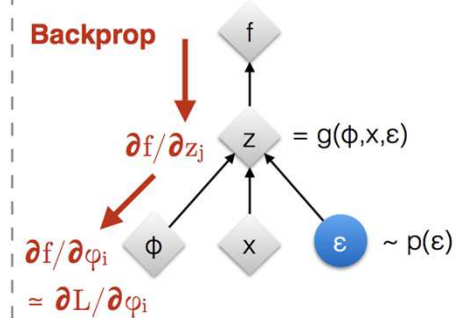
<https://lilianweng.github.io/posts/2018-08-12-vae/#reparameterization-trick>

Reparameterization trick : The expectation term in the loss function invokes generating samples from $\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})$. **Sampling is a stochastic process and therefore we cannot backpropagate the gradient.** To make it trainable, the **reparameterization trick** is introduced: It is often possible to express the random variable \mathbf{z} as a deterministic variable $\mathbf{z} = \mathcal{T}_{\phi}(\mathbf{x}, \epsilon)$, where ϵ is an auxiliary independent random variable, and the transformation function \mathcal{T}_{ϕ} parameterized by ϕ converts ϵ to \mathbf{z} .

Original Form



Reparametrized Form



◇ : Deterministic node

● : Random node

[Kingma, 2013]
[Bengio, 2013]
[Kingma and Welling 2014]
[Rezende et al 2014]

$$\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}^{(i)}, \boldsymbol{\sigma}^{2(i)} \mathbf{I})$$

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}, \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

; Reparameterization trick.

What are Diffusion Models?

Connection with stochastic gradient Langevin dynamics

Langevin dynamics is a concept from physics, developed for statistically modeling molecular systems. Combined with stochastic gradient descent, **stochastic gradient Langevin dynamics** (Welling & Teh 2011) can produce samples from a **probability density** $p(\mathbf{x})$ using only the **gradients** $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ in a Markov chain of updates:

$$\mathbf{x}_t = \mathbf{x}_{t-1} + \frac{\delta}{2} \nabla_{\mathbf{x}} \log p(\mathbf{x}_{t-1}) + \sqrt{\delta} \epsilon_t, \quad \text{where } \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Score : Control Drift

Compared to standard SGD, stochastic gradient Langevin dynamics injects Gaussian noise into the parameter updates to avoid collapses into local minima.

where δ is the step size. When $T \rightarrow \infty$, $\epsilon \rightarrow 0$, \mathbf{x}_T equals to the true probability density $p(\mathbf{x})$

[Stochastic Gradient Langevin Dynamics](https://www.stats.ox.ac.uk/~teh/research/compstats/WelTeh2011a.pdf)

<https://www.stats.ox.ac.uk/~teh/research/compstats/WelTeh2011a.pdf>

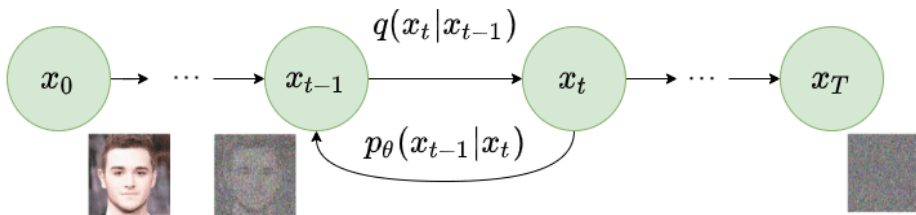
What are Diffusion Models?

2. Reverse diffusion process

➤ Reverse Diffusion: From Noise to Image

You begin with a highly diffused image, represented by \mathbf{x}_T .

The challenge now is to “un-noise” this image, moving sequentially in reverse from \mathbf{x}_T to \mathbf{x}_0 . The process is represented by $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$, indicating the probabilistic rule or mechanism that governs this backward journey.



Reverse Diffusion Process. [Source the AI Summer]

❖ [Introduction to Diffusion Models \(Part III. Diffusion Process\)](#)

- Starting point \mathbf{x}_T : the distorted or noisy image.
- Intermediate states $(\mathbf{x}_t, \mathbf{x}_{t-1}, \dots)$: These are the steps between the completely distorted image and the original, where we systematically remove the noise.
- End point \mathbf{x}_0 : the original, undistorted image.
- $q(\mathbf{x}_t|\mathbf{x}_{t-1})$: This represents the **forward diffusion process**, **guiding how we transition from one state to the next**. Think of it as the rule by which we added noise to our image.
- $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$: This is the heart of **the reverse diffusion**. It's the rule that dictates **how we trace our steps back, effectively removing the noise**.

What are Diffusion Models?

2. Reverse diffusion process

➤ Decoding the Reverse Diffusion Equation

The reverse diffusion process at a specific timestep t

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$$

- $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$: The conditional probability distribution of the image at the previous timestep \mathbf{x}_{t-1} given the current image \mathbf{x}_t . Essentially, it answers the question: “Given the current noisy image \mathbf{x}_t , what is the likelihood of the previous image being \mathbf{x}_{t-1} ?”
- $\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t)$: The mean of the Gaussian distribution, parameterized by θ . It's a function of the current image \mathbf{x}_t and the timestep t . In simpler terms, it provides the expected or “average” image at the previous timestep, given our current image. It's the central point from where the reverse diffusion starts.
- $\boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t)$: Represents the Gaussian distribution's variance (or spread), also parameterized by θ . This variance indicates how “spread out” or varied the possible images at \mathbf{x}_{t-1} could be, given \mathbf{x}_t . A larger variance means that there is a wider range of possible images for \mathbf{x}_{t-1} , while a smaller variance indicates that the possible images are closely packed around the mean.

❖ [Introduction to Diffusion Models \(Part III. Diffusion Process\)](#)

❖ The mean gives the best guess of that previous frame, while the variance tells how confident or uncertain this prediction is.

What are Diffusion Models?

2. Reverse diffusion process

If we can reverse the above process and sample from $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$, we will be able to recreate the true sample from a Gaussian noise input, $\mathbf{x}_T \sim N(\mathbf{0}, I)$.

Note that if β_t is small enough, $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ will also be Gaussian. Unfortunately, we cannot easily estimate $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ because it needs to use the entire dataset and therefore we need to learn a model p_θ to approximate these conditional probabilities in order to run the **reverse diffusion process**.

Reverse diffusion process

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

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

Forward diffusion process

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

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

$$q(\mathbf{x}_t|\mathbf{x}_0) = N(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)I) \quad \bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

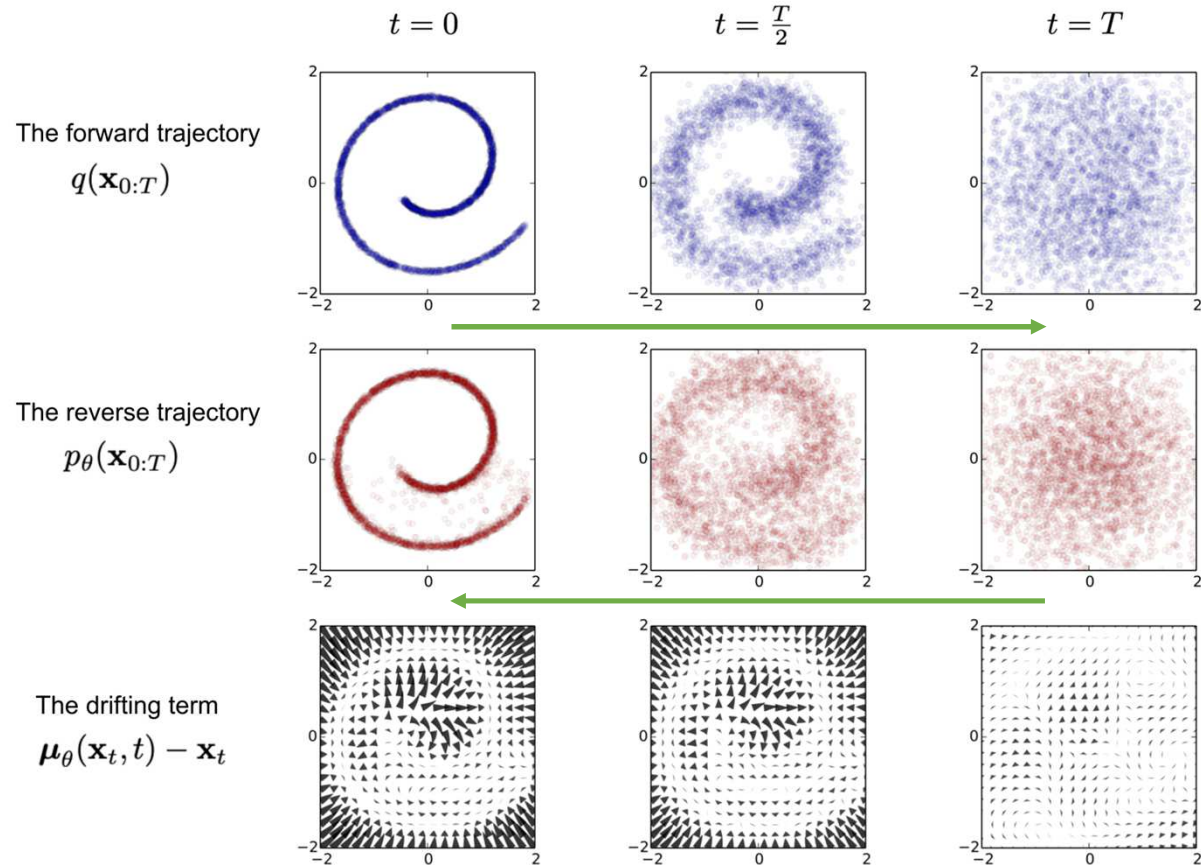


Fig. 3. An example of training a diffusion model for modeling a 2D swiss roll data. (Image source: [Sohl-Dickstein et al., 2015](#))

What are Diffusion Models?

2. Reverse diffusion process

It is noteworthy that the **reverse conditional probability is tractable** when conditioned on \mathbf{x}_0 :

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

Using **Bayes' rule**, we have:

$$\begin{aligned} q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) &= q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \\ &\propto \exp \left(-\frac{1}{2} \left(\frac{(\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{x}_{t-1})^2}{\beta_t} + \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0)^2}{1 - \bar{\alpha}_{t-1}} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0)^2}{1 - \bar{\alpha}_t} \right) \right) \\ &= \exp \left(-\frac{1}{2} \left(\frac{\mathbf{x}_t^2 - 2\sqrt{\alpha_t} \mathbf{x}_t \mathbf{x}_{t-1} + \alpha_t \mathbf{x}_{t-1}^2}{\beta_t} + \frac{\mathbf{x}_{t-1}^2 - 2\sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0 \mathbf{x}_{t-1} + \bar{\alpha}_{t-1} \mathbf{x}_0^2}{1 - \bar{\alpha}_{t-1}} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0)^2}{1 - \bar{\alpha}_t} \right) \right) \\ &= \exp \left(-\frac{1}{2} \left(\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1}^2 - \left(\frac{2\sqrt{\alpha_t}}{\beta_t} \mathbf{x}_t + \frac{2\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} \mathbf{x}_0 \right) \mathbf{x}_{t-1} + C(\mathbf{x}_t, \mathbf{x}_0) \right) \right) \end{aligned}$$

where $C(\mathbf{x}_t, \mathbf{x}_0)$ is some function not involving \mathbf{x}_{t-1} and details are omitted.

📖 Conditional distributions of the multivariate normal distribution :
<https://statproofbook.github.io/P/mvn-cond.html>

📖 Multivariate normal distribution :
https://en.wikipedia.org/wiki/Multivariate_normal_distribution

📖 Bayes Theorem of 3 Events

$$P(A|B \cap C) = \frac{P(B|A \cap C) P(A|C)}{P(B|C)}$$

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$$

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

$$\alpha_t = 1 - \beta_t$$

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

What are Diffusion Models?

2. Reverse diffusion process

Following the standard Gaussian density function, the mean and variance can be parameterized as follows (recall that $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$):

$$\begin{aligned}\tilde{\beta}_t &= 1 / \left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) = 1 / \left(\frac{\alpha_t - \bar{\alpha}_t + \beta_t}{\beta_t(1 - \bar{\alpha}_{t-1})} \right) = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \cdot \beta_t \\ \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) &= \left(\frac{\sqrt{\alpha_t}}{\beta_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} \mathbf{x}_0 \right) / \left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) \\ &= \left(\frac{\sqrt{\alpha_t}}{\beta_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} \mathbf{x}_0 \right) \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \cdot \beta_t \\ &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0\end{aligned}$$

Thanks to the nice property, we can represent $\mathbf{x}_0 = \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_t)$ and plug it into the above equation and obtain:

$$\begin{aligned}\tilde{\mu}_t &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_t) \\ &= \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_t \right)\end{aligned}$$

$$\begin{aligned}q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I}) \\ &\propto \exp \left(-\frac{1}{2} \left(\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1}^2 - \left(\frac{2\sqrt{\alpha_t}}{\beta_t} \mathbf{x}_t + \frac{2\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} \mathbf{x}_0 \right) \mathbf{x}_{t-1} + C(\mathbf{x}_t, \mathbf{x}_0) \right) \right) \\ &\propto \exp \left(-\frac{1}{2\tilde{\beta}_t} (\mathbf{x}_{t-1} - \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0))^2 \right)\end{aligned}$$

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \epsilon_{t-1} && \text{; where } \epsilon_{t-1}, \epsilon_{t-2}, \dots \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \bar{\epsilon}_{t-2} && \text{; where } \bar{\epsilon}_{t-2} \text{ merges two Gaussians (*)}. \\ &= \dots \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon\end{aligned}$$

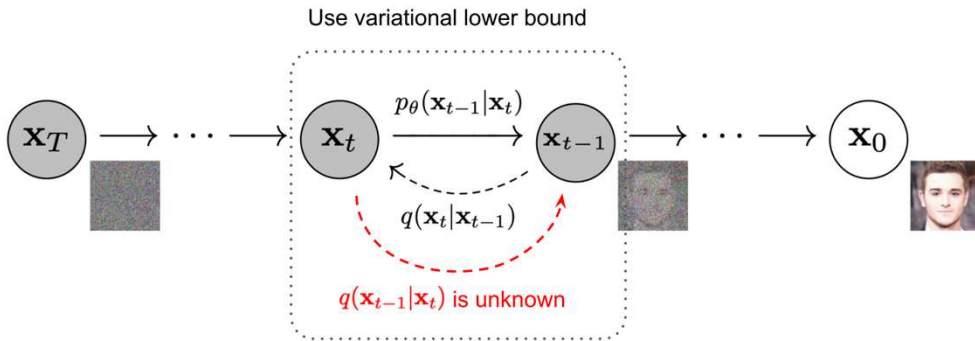
What are Diffusion Models?

2. Reverse diffusion process

As demonstrated in Fig. 2., such a setup is very similar to [VAE](#) and thus we can use the variational lower bound to optimize the negative log-likelihood.

$$\begin{aligned}
 -\log p_\theta(\mathbf{x}_0) &\leq -\log p_\theta(\mathbf{x}_0) + D_{\text{KL}}(q(\mathbf{x}_{1:T}|\mathbf{x}_0) \| p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)) \\
 &= -\log p_\theta(\mathbf{x}_0) + \mathbb{E}_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})/p_\theta(\mathbf{x}_0)} \right] \\
 &= -\log p_\theta(\mathbf{x}_0) + \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} + \log p_\theta(\mathbf{x}_0) \right] \\
 &= \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right]
 \end{aligned}$$

$$\text{Let } L_{\text{VLB}} = \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \geq -\mathbb{E}_{q(\mathbf{x}_0)} \log p_\theta(\mathbf{x}_0)$$



It is also straightforward to get the same result using Jensen's inequality. Say we want to minimize the cross entropy as the learning objective,

$$\begin{aligned}
 L_{\text{CE}} &= -\mathbb{E}_{q(\mathbf{x}_0)} \log p_\theta(\mathbf{x}_0) \\
 &= -\mathbb{E}_{q(\mathbf{x}_0)} \log \left(\int p_\theta(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \right) \\
 &= -\mathbb{E}_{q(\mathbf{x}_0)} \log \left(\int q(\mathbf{x}_{1:T}|\mathbf{x}_0) \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \right) \\
 &= -\mathbb{E}_{q(\mathbf{x}_0)} \log \left(\mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right) \\
 &\leq -\mathbb{E}_{q(\mathbf{x}_{0:T})} \log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \\
 &= \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] = L_{\text{VLB}}
 \end{aligned}$$

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}) \quad p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

What are Diffusion Models?

2. Reverse diffusion process

To convert each term in the equation to be *analytically computable*, the objective can be further rewritten to be a combination of several KL-divergence and entropy terms

See the detailed step-by-step process in [Appendix B in Sohl-Dickstein et al., 2015](#):

Deep Unsupervised Learning using Nonequilibrium Thermodynamics

$$\begin{aligned}
 L_{\text{VLB}} &= \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \\
 &= \mathbb{E}_q \left[\log \frac{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=1}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \quad \text{By Bayes' rule} \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \left(\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} \cdot \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \right) + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} + \log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1)} \right] \\
 &= \mathbb{E}_q \left[\log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{p_\theta(\mathbf{x}_T)} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)} - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) \right] \\
 &= \mathbb{E}_q \left[\underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p_\theta(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \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}} - \underbrace{\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)}_{L_0} \right]
 \end{aligned}$$

What are Diffusion Models?

2. Reverse diffusion process

Let's label each component in the variational lower bound loss separately:

$$L_{\text{VLB}} = L_T + L_{T-1} + \dots + L_0$$

where $L_T = D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p_\theta(\mathbf{x}_T))$

✓ $L_t = D_{\text{KL}}(q(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}))$ for $1 \leq t \leq T-1$

$$L_0 = -\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)$$

- Every KL term in L_{VLB} (except for L_0) compares two Gaussian distributions and therefore they can be computed in closed form.

- L_T is constant and can be ignored during training because q has no learnable parameters and \mathbf{x}_T is a Gaussian noise.
- [Ho et al. 2020](#) (Denoising Diffusion Probabilistic Models) models L_0 using a separate discrete decoder derived from $N(\mathbf{x}_0; \boldsymbol{\mu}_\theta(\mathbf{x}_1, \mathbf{1}), \boldsymbol{\Sigma}_\theta(\mathbf{x}_1, \mathbf{1}))$

Closed form : [Multivariate normal distributions](#)

https://en.wikipedia.org/wiki/Kullback%E2%80%93Leibler_divergence#Multivariate_normal_distributions

What are Diffusion Models?

3. Parameterization of L_t for Training Loss

Recall that we need to learn a neural network to approximate the conditioned probability distributions in the reverse diffusion process, $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = N(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$.

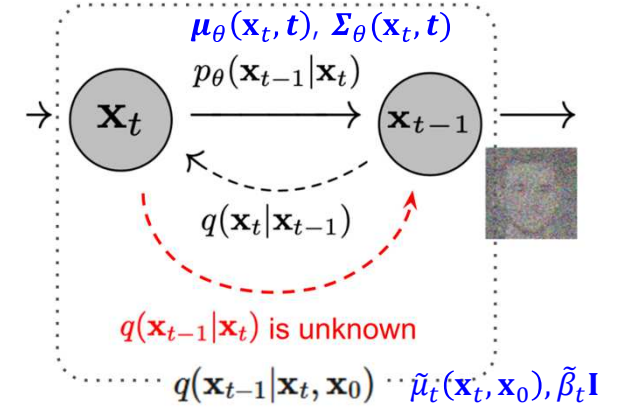
We would like to train $\boldsymbol{\mu}_\theta$ to predict $\tilde{\boldsymbol{\mu}}_t = \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\boldsymbol{\epsilon}_t)$. Because \mathbf{x}_t is available as input at training time, we can reparameterize the Gaussian noise term instead to make it predict $\boldsymbol{\epsilon}_t$ from the input \mathbf{x}_t at time step t :

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

$$\text{Thus } \mathbf{x}_{t-1} = \mathcal{N}(\mathbf{x}_{t-1}; \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

The loss term L_t is parameterized to minimize the difference from $\tilde{\boldsymbol{\mu}}_t$:

$$\begin{aligned} L_t &= \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[\frac{1}{2\|\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)\|_2^2} \|\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) - \boldsymbol{\mu}_\theta(\mathbf{x}_t, t)\|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[\frac{1}{2\|\boldsymbol{\Sigma}_\theta\|_2^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\boldsymbol{\epsilon}_t \right) - \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) \right\|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[\frac{(1-\alpha_t)^2}{2\alpha_t(1-\bar{\alpha}_t)\|\boldsymbol{\Sigma}_\theta\|_2^2} \|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[\frac{(1-\alpha_t)^2}{2\alpha_t(1-\bar{\alpha}_t)\|\boldsymbol{\Sigma}_\theta\|_2^2} \|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\boldsymbol{\epsilon}_t, t)\|^2 \right] \end{aligned}$$



$$\begin{aligned} L_t &= D_{\text{KL}}(q(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})) \\ p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) &= N(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)) \\ q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\boldsymbol{\beta}}_t \mathbf{I}) \end{aligned}$$

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

What are Diffusion Models?

Simplification

Empirically, [Ho et al. \(2020\)](#) (Denoising Diffusion Probabilistic Models) found that training the diffusion model works better with a simplified objective that ignores the weighting term:

$$\begin{aligned} L_t^{\text{simple}} &= \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} \left[\|\epsilon_t - \epsilon_\theta(\mathbf{x}_t, t)\|^2 \right] \\ &= \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} \left[\|\epsilon_t - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_t, t)\|^2 \right] \end{aligned}$$

The final simple objective is:

$$L_{\text{simple}} = L_t^{\text{simple}} + C$$

where C is a constant not depending on θ .

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
        $\nabla_\theta \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2$ 
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:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

Fig. 4. The training and sampling algorithms in DDPM (Image source: Ho et al. 2020)

What are Diffusion Models?

Connection with noise-conditioned score networks (NCSN)

- [Song & Ermon \(2019\)](#) proposed a **score-based generative modeling** method where **samples** are produced via **Langevin dynamics** using **gradients of the data distribution estimated with score matching**.
- The score of each sample \mathbf{x} 's density probability is defined as its gradient $\nabla_{\mathbf{x}} \log q(\mathbf{x})$.
- A score network $\mathbf{s}_{\theta}: \mathbb{R}^D \rightarrow \mathbb{R}^D$ is trained to estimate it, $\mathbf{s}_{\theta}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log q(\mathbf{x})$.
- To make it scalable with high-dimensional data in the deep learning setting, they proposed to use either **denoising score matching** ([Vincent, 2011](#)) or **sliced score matching** (use random projections; [Song et al., 2019](#)).
Denoising score matching adds a pre-specified small noise to the data $q(\tilde{\mathbf{x}}|\mathbf{x})$ and estimates $q(\tilde{\mathbf{x}})$ with score matching.
- Recall that Langevin dynamics can sample data points from a probability density distribution using only the **score** $\nabla_{\mathbf{x}} \log q(\mathbf{x})$ in an iterative process..

☞ By notation in Song : Score function – Score Network

$$\mathbf{s}_{\theta}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p(\mathbf{x})$$

- 1) [Song & Ermon, Generative Modeling by Estimating Gradients of the Data Distribution](#)
- 2) [Vincent, denoising score matching](#)
- 3) [Song et al., Sliced score matching](#)

- ❖ **Stochastic gradient Langevin dynamics** can produce samples from a probability density $p(\mathbf{x})$ using only the gradients $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ in a Markov chain of updates:

$$\mathbf{x}_t = \mathbf{x}_{t-1} + \frac{\delta}{2} \nabla_{\mathbf{x}} \log p(\mathbf{x}_{t-1}) + \sqrt{\delta} \boldsymbol{\epsilon}_t, \quad \text{where } \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

What are Diffusion Models?

Connection with noise-conditioned score networks (NCSN)

- However, according to the manifold hypothesis, most of the data is expected to concentrate in a **low dimensional manifold**, even though the observed data might look only arbitrarily high-dimensional. It brings a negative effect on score estimation since the data points cannot cover the whole space.
- In regions where data density is low, the score estimation is less reliable. After adding a small Gaussian noise to make the perturbed data distribution cover the full space \mathbb{R}^D , the training of the score estimator network becomes more stable.

➤ [Song & Ermon \(2019\)](#) improved it by **perturbing the data with the noise of different levels** and **train a noise-conditioned score network to jointly estimate the scores of all the perturbed data at different noise levels.**

- The schedule of increasing *noise levels* resembles the forward diffusion process. If we use the diffusion process annotation, **the score approximates $\mathbf{s}_\theta(\mathbf{x}_t, t) \approx \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)$** . Given a Gaussian distribution $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$, we can write the derivative of the logarithm of its density function as $\nabla_{\mathbf{x}} \log p(\mathbf{x}) = \nabla_{\mathbf{x}} \left(-\frac{1}{2\sigma^2} (\mathbf{x} - \boldsymbol{\mu})^2 \right) = -\frac{\mathbf{x} - \boldsymbol{\mu}}{\sigma^2} = -\frac{\boldsymbol{\epsilon}}{\sigma}$, where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

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

Therefore,

$$\begin{aligned} \mathbf{s}_\theta(\mathbf{x}_t, t) &\approx \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) \\ &= \mathbb{E}_{q(\mathbf{x}_0)} [\nabla_{\mathbf{x}_t} q(\mathbf{x}_t | \mathbf{x}_0)] = \mathbb{E}_{q(\mathbf{x}_0)} \left[-\frac{\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\sqrt{1 - \bar{\alpha}_t}} \right] = -\frac{\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\sqrt{1 - \bar{\alpha}_t}} \end{aligned}$$

What are Diffusion Models?

4. Parameterization of β_t

- The **forward variances** are set to be a sequence of *linearly* increasing constants in [Ho et al. \(2020\)](#), from $\beta_t = 10^{-4}$ to $\beta_T = 0.02$. They are relatively small compared to the normalized image pixel values between $[-1, 1]$. Diffusion models in their experiments showed *high-quality samples* but still could *not achieve* competitive model *log-likelihood* as other generative models.
- [Nichol & Dhariwal \(2021\)](#) proposed several improvement techniques to help diffusion models to obtain lower NLL(Negative Log Likelihood). One of the improvements is to use a *cosine-based variance schedule*. The choice of the scheduling function can be arbitrary, as long as it provides *a near-linear drop in the middle of the training process and subtle changes around $t = 0$ and $t = T$* .

$$\beta_t = \text{clip}\left(1 - \frac{\bar{\alpha}_t}{\bar{\alpha}_{t-1}}, 0.999\right) \quad \bar{\alpha}_t = \frac{f(t)}{f(0)} \quad \text{where } f(t) = \cos\left(\frac{t/T + s}{1 + s} \cdot \frac{\pi}{2}\right)^2$$

where the small offset s is to prevent β_t from being too small when close to t .

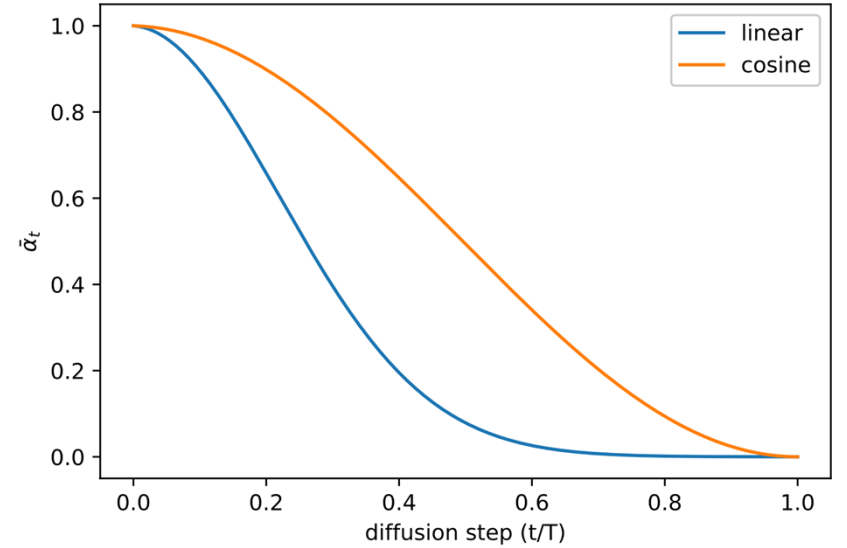


Fig. 5. Comparison of linear and cosine-based scheduling of β_t during training. (Image source: Nichol & Dhariwal, 2021)

What are Diffusion Models?

5. Parameterization of reverse process variance Σ_θ

- [Ho et al. \(2020\)](#) (DDPM) chose to fix β_t as constants instead of making them learnable and set $\Sigma_\theta(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}$, where σ_t is not learned but set to β_t or $\tilde{\beta}_t = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}$. Because they found that learning a diagonal variance Σ_θ leads to unstable training and poorer sample quality.
- [Nichol & Dhariwal \(2021\)](#) (Improved DDPM) proposed to learn $\Sigma_\theta(\mathbf{x}_t, t)$ as an interpolation between β_t and $\tilde{\beta}_t$ by model predicting a mixing vector \mathbf{v} :

$$\Sigma_\theta(\mathbf{x}_t, t) = \exp(\mathbf{v} \log \beta_t + (1 - \mathbf{v}) \log \tilde{\beta}_t)$$

- However, the simple objective L_{simple} does not depend on Σ_θ . To add the dependency, they constructed a hybrid objective

$$L_{\text{hybrid}} = L_{\text{simple}} + \lambda L_{\text{VLB}}$$

where $\lambda = 0.001$ is small and stop gradient on μ_θ in the L_{VLB} term such that L_{VLB} only guides the learning of Σ_θ .

- Empirically they observed that L_{VLB} is pretty challenging to optimize likely due to noisy gradients, so they proposed to use a time-averaging smoothed version of L_{VLB} with importance sampling.

Model	ImageNet	CIFAR
Glow (Kingma & Dhariwal, 2018)	3.81	3.35
Flow++ (Ho et al., 2019)	3.69	3.08
PixelCNN (van den Oord et al., 2016c)	3.57	3.14
SPN (Menick & Kalchbrenner, 2018)	3.52	-
NVAE (Vahdat & Kautz, 2020)	-	2.91
Very Deep VAE (Child, 2020)	3.52	2.87
PixelSNAIL (Chen et al., 2018)	3.52	2.85
Image Transformer (Parmar et al., 2018)	3.48	2.90
Sparse Transformer (Child et al., 2019)	3.44	2.80
Routing Transformer (Roy et al., 2020)	3.43	-
DDPM (Ho et al., 2020)	3.77	3.70
DDPM (cont flow) (Song et al., 2020b)	-	2.99
Improved DDPM (ours)	3.53	2.94

Fig. 6. Comparison of negative log-likelihood (NLL) of improved DDPM with other likelihood-based generative models. NLL is reported in the unit of bits/dim. (Image source: Nichol & Dhariwal, 2021)

Conditioned Generation

While training generative models on images with conditioning information such as ImageNet dataset, it is common to **generate samples conditioned on class labels or a piece of descriptive text**.

Classifier Guided Diffusion

To explicitly incorporate class information into the diffusion process, [Dhariwal & Nichol \(2021\)](#) trained a classifier $f_\phi(y|\mathbf{x}_t, t)$ on noisy image \mathbf{x}_t and use gradients $\nabla_{\mathbf{x}} \log f_\phi(y|\mathbf{x}_t)$ to guide the diffusion sampling process toward the conditioning information y (e.g. a target class label) by altering the noise prediction.

Recall that $\nabla_{\mathbf{x}_t} \log p_\theta(\mathbf{x}_t) = -\frac{1}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t)$ and we can write the **score function for the joint distribution $q(\mathbf{x}_t, y)$** as following,

$$\begin{aligned} \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t, y) &= \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log q(y|\mathbf{x}_t) \\ &\approx -\frac{1}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) + \nabla_{\mathbf{x}_t} \log f_\phi(y|\mathbf{x}_t) \\ &= -\frac{1}{\sqrt{1-\bar{\alpha}_t}} (\epsilon_\theta(\mathbf{x}_t, t) - \sqrt{1-\bar{\alpha}_t} \nabla_{\mathbf{x}_t} \log f_\phi(y|\mathbf{x}_t)) \\ &\quad \bar{\epsilon}_\theta(\mathbf{x}_t, t) \end{aligned}$$

[Dhariwal & Nichol \(2021\), Diffusion Models Beat GANs on Image Synthesis](#)

Thus, a new classifier-guided predictor $\bar{\epsilon}_\theta$ would take the form as following,

$$\bar{\epsilon}_\theta(\mathbf{x}_t, t) = \epsilon_\theta(\mathbf{x}_t, t) - \sqrt{1-\bar{\alpha}_t} \nabla_{\mathbf{x}_t} \log f_\phi(y|\mathbf{x}_t)$$

To control the strength of the classifier guidance, we can add a weight w to the delta part,

$$\bar{\epsilon}_\theta(\mathbf{x}_t, t) = \epsilon_\theta(\mathbf{x}_t, t) - \sqrt{1-\bar{\alpha}_t} w \nabla_{\mathbf{x}_t} \log f_\phi(y|\mathbf{x}_t)$$

The resulting *ablated diffusion model (ADM)* and the one with additional classifier guidance (**ADM-G**) are able to achieve better results than SOTA generative models (e.g. BigGAN).

$$\nabla_{x_t} \log p_\theta(x_t) = -\frac{1}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(x_t) \quad \text{Original Paper notation}$$

$$\begin{aligned} \nabla_{x_t} \log(p_\theta(x_t)p_\phi(y|x_t)) &= \nabla_{x_t} \log p_\theta(x_t) + \nabla_{x_t} \log p_\phi(y|x_t) \\ &= -\frac{1}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(x_t) + \nabla_{x_t} \log p_\phi(y|x_t) \end{aligned}$$

Conditioned Generation

Classifier Guided Diffusion

Algorithm 1 Classifier guided diffusion sampling, given a diffusion model $(\mu_\theta(x_t), \Sigma_\theta(x_t))$, classifier $f_\phi(y|x_t)$, and gradient scale s .

```
Input: class label  $y$ , gradient scale  $s$ 
 $x_T \leftarrow$  sample from  $\mathcal{N}(0, \mathbf{I})$ 
for all  $t$  from  $T$  to 1 do
   $\mu, \Sigma \leftarrow \mu_\theta(x_t), \Sigma_\theta(x_t)$ 
   $x_{t-1} \leftarrow$  sample from  $\mathcal{N}(\mu + s\Sigma \nabla_{x_t} \log f_\phi(y|x_t), \Sigma)$ 
end for
return  $x_0$ 
```

Algorithm 2 Classifier guided DDIM sampling, given a diffusion model $\epsilon_\theta(x_t)$, classifier $f_\phi(y|x_t)$, and gradient scale s .

```
Input: class label  $y$ , gradient scale  $s$ 
 $x_T \leftarrow$  sample from  $\mathcal{N}(0, \mathbf{I})$ 
for all  $t$  from  $T$  to 1 do
   $\hat{\epsilon} \leftarrow \epsilon_\theta(x_t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{x_t} \log f_\phi(y|x_t)$ 
   $x_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \left( \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}$ 
end for
return  $x_0$ 
```

$$\frac{1}{\sqrt{\bar{\alpha}_t}} \mathbf{x}_t - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} \hat{\epsilon} + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon} \quad ?$$

Additionally with some modifications on the **U-Net architecture**, Dhariwal & Nichol (2021) showed performance better than GAN with diffusion models.

The architecture modifications include larger model depth/width, more attention heads, multi-resolution attention, BigGAN residual blocks for up/downsampling, residual connection rescale by $1/\sqrt{2}$ and adaptive group normalization (AdaGN).

Fig. 7. The algorithms use guidance from a classifier to run conditioned generation with DDPM and DDIM. (Image source: [Dhariwal & Nichol, 2021](#))]

Conditioned Generation

Classifier-Free Guidance

Without an independent classifier f_ϕ , it is still possible to run conditional diffusion steps by incorporating the scores from a conditional and an unconditional diffusion model (Ho & Salimans, 2021).

Let unconditional denoising diffusion model $p_\theta(\mathbf{x})$ parameterized through a score estimator $\epsilon_\theta(\mathbf{x}_t, t)$ and the conditional model $p_\theta(\mathbf{x}|y)$ parameterized through $\epsilon_\theta(\mathbf{x}_t, t, y)$.

These two models can be learned via a single neural network. Precisely, a conditional diffusion model $p_\theta(\mathbf{x}|y)$ is trained on paired data (\mathbf{x}, y) , where the conditioning information y gets discarded periodically at random such that the model knows how to generate images unconditionally as well, i.e.

$$\epsilon_\theta(\mathbf{x}_t, t) = \epsilon_\theta(\mathbf{x}_t, t, y = \emptyset).$$

[Ho & Salimans, \(2021\), Classifier-Free Diffusion Guidance](#)

[Nichol, Dhariwal & Ramesh, et al. \(2022\), GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models](#)

The gradient of an implicit classifier can be represented with conditional and unconditional score estimators. Once plugged into the classifier-guided modified score, the score contains no dependency on a separate classifier.

$$\begin{aligned}\nabla_{\mathbf{x}_t} \log p(y|\mathbf{x}_t) &= \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|y) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \\ &= -\frac{1}{\sqrt{1-\bar{\alpha}_t}} \left(\epsilon_\theta(\mathbf{x}_t, t, y) - \epsilon_\theta(\mathbf{x}_t, t) \right) \\ \rightarrow \bar{\epsilon}_\theta(\mathbf{x}_t, t, y) &= \epsilon_\theta(\mathbf{x}_t, t, y) - \sqrt{1-\bar{\alpha}_t} w \nabla_{\mathbf{x}_t} \log p(y|\mathbf{x}_t) \\ &= \epsilon_\theta(\mathbf{x}_t, t, y) + w(\epsilon_\theta(\mathbf{x}_t, t, y) - \epsilon_\theta(\mathbf{x}_t, t)) \\ &= (w+1)\epsilon_\theta(\mathbf{x}_t, t, y) - w\epsilon_\theta(\mathbf{x}_t, t)\end{aligned}$$

Their experiments showed that classifier-free guidance can achieve a good balance between FID (distinguish between synthetic and generated images) and IS (quality and diversity).

The **guided diffusion model, GLIDE** (Nichol, Dhariwal & Ramesh, et al. 2022), explored both **guiding strategies, CLIP guidance and classifier-free guidance**, and found that **the latter is more preferred**. They hypothesized that it is because CLIP guidance exploits the model with adversarial examples towards the CLIP model, rather than optimize the better matched images generation.

Conditioned Generation

Classifier-Free Guidance

Algorithm 1 Joint training a diffusion model with classifier-free guidance

Require: p_{uncond} : probability of unconditional training

- 1: **repeat**
 - 2: $(\mathbf{x}, \mathbf{c}) \sim p(\mathbf{x}, \mathbf{c})$ ▷ Sample data with conditioning from the dataset \mathbf{c}
 - 3: $\mathbf{c} \leftarrow \emptyset$ with probability p_{uncond} ▷ Randomly discard conditioning to train unconditionally
 - 4: $\lambda \sim p(\lambda)$ ▷ Sample log SNR value
 - 5: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 6: $\mathbf{z}_\lambda = \alpha_\lambda \mathbf{x} + \sigma_\lambda \epsilon$ ▷ Corrupt data to the sampled log SNR value
 - 7: Take gradient step on $\nabla_\theta \|\epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c}) - \epsilon\|^2$ ▷ Optimization of denoising model
 - 8: **until** converged
-

Classifier-free guidance는 $\epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c}) (= \epsilon_\theta(\mathbf{x}_t, t, y))$ 를 수정하여 classifier 없이 classifier guidance와 같은 효과를 얻는 방법

Algorithm 2 Conditional sampling with classifier-free guidance

Require: w : guidance strength

Require: \mathbf{c} : conditioning information for conditional sampling

Require: $\lambda_1, \dots, \lambda_T$: increasing log SNR sequence with $\lambda_1 = \lambda_{\min}$, $\lambda_T = \lambda_{\max}$

- 1: $\mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 2: **for** $t = 1, \dots, T$ **do**
 - ▷ Form the classifier-free guided score at log SNR λ_t
 - 3: $\tilde{\epsilon}_t = (1 + w)\epsilon_\theta(\mathbf{z}_t, \mathbf{c}) - w\epsilon_\theta(\mathbf{z}_t)$
 - ▷ Sampling step (could be replaced by another sampler, e.g. DDIM)
 - 4: $\tilde{\mathbf{x}}_t = (\mathbf{z}_t - \sigma_{\lambda_t} \tilde{\epsilon}_t) / \alpha_{\lambda_t}$
 - 5: $\mathbf{z}_{t+1} \sim \mathcal{N}(\tilde{\mu}_{\lambda_{t+1}|\lambda_t}(\mathbf{z}_t, \tilde{\mathbf{x}}_t), (\tilde{\sigma}_{\lambda_{t+1}|\lambda_t}^2)^{1-v} (\sigma_{\lambda_t|\lambda_{t+1}}^2)^v)$ if $t < T$ else $\mathbf{z}_{t+1} = \tilde{\mathbf{x}}_t$
 - 6: **end for**
 - 7: **return** \mathbf{z}_{T+1}
-

$$\begin{aligned}
 \nabla_{\mathbf{x}_t} \log p(y|\mathbf{x}_t) &= \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|y) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \\
 &= -\frac{1}{\sqrt{1 - \bar{\alpha}_t}} \left(\epsilon_\theta(\mathbf{x}_t, t, y) - \epsilon_\theta(\mathbf{x}_t, t) \right) \\
 \bar{\epsilon}_\theta(\mathbf{x}_t, t, y) &= \epsilon_\theta(\mathbf{x}_t, t, y) - \sqrt{1 - \bar{\alpha}_t} w \nabla_{\mathbf{x}_t} \log p(y|\mathbf{x}_t) \\
 &= \epsilon_\theta(\mathbf{x}_t, t, y) + w(\epsilon_\theta(\mathbf{x}_t, t, y) - \epsilon_\theta(\mathbf{x}_t, t)) \\
 &= (w + 1)\epsilon_\theta(\mathbf{x}_t, t, y) - w\epsilon_\theta(\mathbf{x}_t, t)
 \end{aligned}$$

Speed up Diffusion Model Sampling

It is **very slow** to generate a sample from **DDPM** by following [the Markov chain of the reverse diffusion process](#), as T can be up to **one or a few thousand steps**. One data point from [Song et al. 2020](#): “For example, it takes around 20 hours to sample 50k images of size 32×32 from a DDPM, but less than a minute to do so from a GAN on an Nvidia 2080 Ti GPU.”

Fewer Sampling Steps (DDIM, IDDPM)

- One simple way is to run a **strided sampling schedule** ([Nichol & Dhariwal, 2021](#)) by **taking the sampling update every $[T/S]$ steps to reduce the process from T to S steps**. The new sampling schedule for generation is $\{\tau_1, \tau_2, \dots, \tau_S\}$ where $\tau_1 < \tau_2 < \dots < \tau_S \in [1, T]$ and $S < T$.
- For another approach, let's rewrite $q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ to be **parameterized by a desired standard deviation σ_t** according to the nice property:

$$\begin{aligned}\mathbf{x}_{t-1} &= \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1}}\epsilon_{t-1} \\ &= \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2}\epsilon_t + \sigma_t\epsilon \\ &= \sqrt{\bar{\alpha}_{t-1}}\left(\frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_\theta^{(t)}(\mathbf{x}_t)}{\sqrt{\bar{\alpha}_t}}\right) + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2}\epsilon_\theta^{(t)}(\mathbf{x}_t) + \sigma_t\epsilon \\ \mathbf{x}_0 &= \frac{1}{\sqrt{\bar{\alpha}_t}}(\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_t)\end{aligned}$$

[Song et al, Denoising Diffusion Implicit Models, 2021](#)

[Nichol & Dhariwal, Improved Denoising Diffusion Probabilistic Models, 2021](#)

$$q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\left(\frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_\theta^{(t)}(\mathbf{x}_t)}{\sqrt{\bar{\alpha}_t}}\right) + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2}\epsilon_\theta^{(t)}(\mathbf{x}_t), \sigma_t^2\mathbf{I})$$

where the model $\epsilon_\theta^{(t)}(\cdot)$ predicts the ϵ_t from \mathbf{x}_t .

Recall that in $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t\mathbf{I})$, therefore we have:

$$\tilde{\beta}_t = \sigma_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \cdot \beta_t \quad \text{Page 16}$$

Let $\sigma_t^2 = \eta \cdot \tilde{\beta}_t$ such that we can adjust $\eta \in \mathbb{R}^+$ as a hyperparameter to control the sampling stochasticity. The special case of $\eta = 0$ makes the sampling process deterministic. Such a model is named the **denoising diffusion implicit model (DDIM)**; Song et al., 2020). DDIM has the same marginal noise distribution but deterministically maps noise back to the original data samples.

Speed up Diffusion Model Sampling

Fewer Sampling Steps (DDIM, IDDPM)

- During generation, we don't have to follow the whole chain $t = 1, \dots, T$, but rather a subset of steps. Let's denote $s < t$ as two steps in this accelerated trajectory. The **DDIM update step** is:

$$q_{\sigma, s < t}(\mathbf{x}_s | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_s; \sqrt{\bar{\alpha}_s} \left(\frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_{\theta}^{(t)}(\mathbf{x}_t)}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_s - \sigma_t^2 \epsilon_{\theta}^{(t)}(\mathbf{x}_t)}, \sigma_t^2 \mathbf{I})$$

- While all the models are trained with $T = 1000$ diffusion steps in the experiments, they observed that DDIM ($\eta = 0$) can produce the best quality samples when S is small, while DDPM ($\eta = 1$) performs much worse on small S .
- DDPM does perform better when we can afford to run the full reverse Markov diffusion steps ($S = T = 1000$). With **DDIM**, it is possible to train the diffusion model up to any arbitrary number of forward steps but only sample from a subset of steps in the generative process.

S	CIFAR10 (32×32)					CelebA (64×64)					
	10	20	50	100	1000	10	20	50	100	1000	
η	0.0	13.36	6.84	4.67	4.16	4.04	17.33	13.73	9.17	6.53	3.51
	0.2	14.04	7.11	4.77	4.25	4.09	17.66	14.11	9.51	6.79	3.64
	0.5	16.66	8.35	5.25	4.46	4.29	19.86	16.06	11.01	8.09	4.28
	1.0	41.07	18.36	8.01	5.78	4.73	33.12	26.03	18.48	13.93	5.98
$\hat{\sigma}$	367.43	133.37	32.72	9.99	3.17	299.71	183.83	71.71	45.20	3.26	

Fig. 8. FID scores on CIFAR10 and CelebA datasets by diffusion models of different settings, including **DDIM**($\eta = 0$) and **DDPM**($\hat{\sigma}$). (Image source: [Song et al., 2020](#))

Compared to DDPM, **DDIM** is able to:

1. Generate higher-quality samples using **a much fewer number of steps**.
2. Have “**consistency**” property since the generative process is deterministic, meaning that **multiple samples conditioned on the same latent variable should have similar high-level features**.
3. Because of the consistency, DDIM can do **semantically meaningful interpolation in the latent variable**.

Speed up Diffusion Model Sampling

Progressive Distillation

Progressive Distillation (Salimans & Ho, 2022) is a method for distilling trained deterministic samplers into new models of halved sampling steps.

- The student model is initialized from the teacher model and denoises towards a target where **one student DDIM step matches 2 teacher steps**, instead of using the original sample \mathbf{x}_0 as the denoise target.
- In every progressive distillation iteration, we can half the sampling steps.

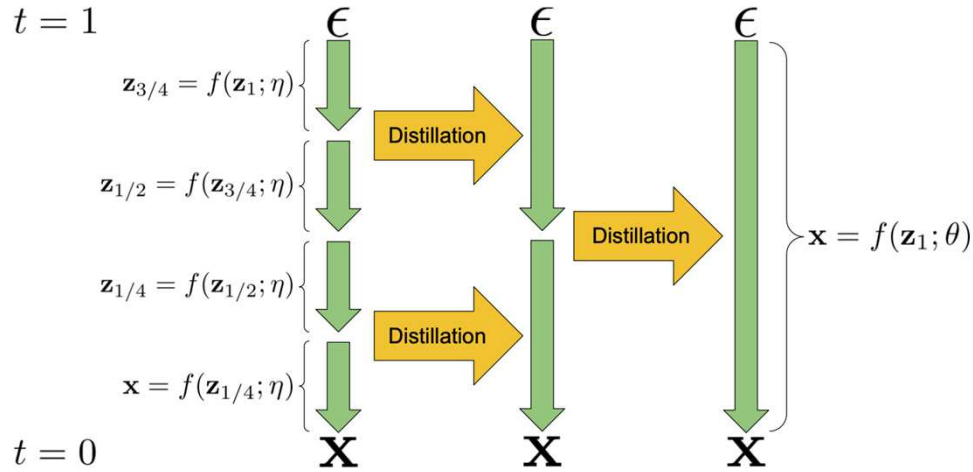


Fig. 9. Progressive distillation can reduce the diffusion sampling steps by half in each iteration. (Image source: Salimans & Ho, 2022)

Tim Salimans, Jonathan Ho, "Progressive Distillation for Fast Sampling of Diffusion Models," 2022.

Algorithm 1 Standard diffusion training

Require: Model $\hat{\mathbf{x}}_\theta(\mathbf{z}_t)$ to be trained
Require: Data set \mathcal{D}
Require: Loss weight function $w()$

while not converged **do**

$\mathbf{x} \sim \mathcal{D}$ \triangleright Sample data
 $t \sim U[0, 1]$ \triangleright Sample time
 $\epsilon \sim N(0, I)$ \triangleright Sample noise
 $\mathbf{z}_t = \alpha_t \mathbf{x} + \sigma_t \epsilon$ \triangleright Add noise to data

$\tilde{\mathbf{x}} = \mathbf{x}$ \triangleright Clean data is target for $\hat{\mathbf{x}}$
 $\lambda_t = \log[\alpha_t^2 / \sigma_t^2]$ \triangleright log-SNR
 $L_\theta = w(\lambda_t) \|\tilde{\mathbf{x}} - \hat{\mathbf{x}}_\theta(\mathbf{z}_t)\|_2^2$ \triangleright Loss
 $\theta \leftarrow \theta - \gamma \nabla_\theta L_\theta$ \triangleright Optimization

end while

Algorithm 2 Progressive distillation

Require: Trained teacher model $\hat{\mathbf{x}}_\eta(\mathbf{z}_t)$
Require: Data set \mathcal{D}
Require: Loss weight function $w()$
Require: Student sampling steps N

for K iterations **do**

$\theta \leftarrow \eta$

\triangleright Init student from teacher

while not converged **do**

$\mathbf{x} \sim \mathcal{D}$

$t = i/N, i \sim \text{Cat}[1, 2, \dots, N]$

$\epsilon \sim N(0, I)$

$\mathbf{z}_t = \alpha_t \mathbf{x} + \sigma_t \epsilon$

2 steps of DDIM with teacher

$t' = t - 0.5/N, t'' = t - 1/N$

$\mathbf{z}_{t'} = \alpha_{t'} \hat{\mathbf{x}}_\eta(\mathbf{z}_t) + \frac{\sigma_{t'}}{\sigma_t} (\mathbf{z}_t - \alpha_t \hat{\mathbf{x}}_\eta(\mathbf{z}_t))$

$\mathbf{z}_{t''} = \alpha_{t''} \hat{\mathbf{x}}_\eta(\mathbf{z}_{t'}) + \frac{\sigma_{t''}}{\sigma_{t'}} (\mathbf{z}_{t'} - \alpha_{t'} \hat{\mathbf{x}}_\eta(\mathbf{z}_{t'}))$

$\tilde{\mathbf{x}} = \frac{\mathbf{z}_{t''} - (\sigma_{t''}/\sigma_t) \mathbf{z}_t}{\alpha_{t''} - (\sigma_{t''}/\sigma_t) \alpha_t}$ \triangleright Teacher $\hat{\mathbf{x}}$ target

$\lambda_t = \log[\alpha_t^2 / \sigma_t^2]$

$L_\theta = w(\lambda_t) \|\tilde{\mathbf{x}} - \hat{\mathbf{x}}_\theta(\mathbf{z}_t)\|_2^2$

$\theta \leftarrow \theta - \gamma \nabla_\theta L_\theta$ **Student model**

end while

$\eta \leftarrow \theta$

\triangleright Student becomes next teacher

$N \leftarrow N/2$ \triangleright Halve number of sampling steps

end for

Fig. 10. Comparison of Algorithm 1 (diffusion model training) and Algorithm 2 (progressive distillation) side-by-side, where the relative changes in progressive distillation are highlighted in green. (Image source: Salimans & Ho, 2022)

Speed up Diffusion Model Sampling

Consistency Models

Consistency Models (Song et al. 2023) learns to map any intermediate noisy data points \mathbf{x}_t , $t > 0$ on the diffusion sampling trajectory back to its origin \mathbf{x}_0 directly.

It is named as **consistency model** because of its **self-consistency** property as any data points on the same trajectory is mapped to the same origin.

- Given a trajectory $\{\mathbf{x}_t | t \in [\epsilon, T]\}$, the consistency function f is defined as $f: (\mathbf{x}_t, t) \mapsto \mathbf{x}_\epsilon$ and the equation $f(\mathbf{x}_t, t) = f(\mathbf{x}_{t'}, t') = \mathbf{x}_\epsilon$ holds true for all $t, t' \in [\epsilon, T]$. When $t = \epsilon$, f is an identify function.
- The model can be parameterized as follows, where $c_{\text{skip}}(t)$ and $c_{\text{out}}(t)$ functions are designed in a way that $c_{\text{skip}}(\epsilon) = 1$, $c_{\text{out}}(\epsilon) = 0$:

$$f_\theta(\mathbf{x}, t) = c_{\text{skip}}(t)\mathbf{x} + c_{\text{out}}(t)F_\theta(\mathbf{x}, t)$$

- ✓ $F_\theta(\mathbf{x}, t)$: a free-form deep neural network, Output has the same dimensionality as \mathbf{x} .

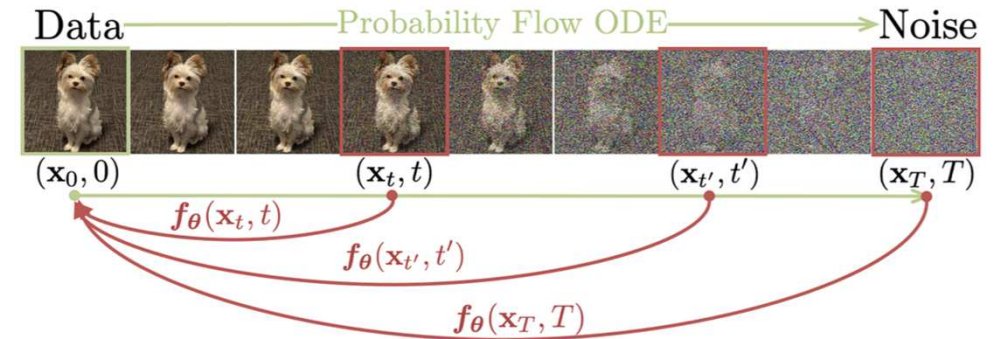


Fig. 11. Consistency models learn to map any data point on the trajectory back to its origin. (Image source: Song et al., 2023)

- It is possible for the consistency model to generate **samples in a single step**, while still maintaining the flexibility of trading computation for better quality following a **multi-step sampling process**.

Speed up Diffusion Model Sampling

Consistency Models

The paper introduced two ways to train consistency models:

1) Consistency Distillation (CD)

- Distill a diffusion model into a consistency model by minimizing the difference between model outputs for pairs generated out of the same trajectory. This enables a much cheaper sampling evaluation.
- The consistency distillation (CD) loss is:

$$\mathcal{L}_{\text{CD}}^N(\theta, \theta^-; \phi) = \mathbb{E}[\lambda(t_n) d(f_\theta(\mathbf{x}_{t_{n+1}}, t_{n+1}), f_{\theta^-}(\hat{\mathbf{x}}_{t_n}^\phi, t_n))]$$
$$\hat{\mathbf{x}}_{t_n}^\phi = \mathbf{x}_{t_{n+1}} - (t_n - t_{n+1}) \Phi(\mathbf{x}_{t_{n+1}}, t_{n+1}; \phi)$$

- $\Phi(\cdot; \phi)$ is the update function of a one-step [ODE](#) (Ordinary differential equation) solver;
- $n \sim \mathcal{U}[1, N - 1]$ has a uniform distribution over $1, \dots, N - 1$;
- The network parameters θ^- is EMA(Exponential Moving Average) version of θ which greatly stabilizes the training (just like in [DQN](#) or [momentum contrastive](#) learning);
- $d(\cdot, \cdot)$ is a positive distance metric function that satisfies $\forall \mathbf{x}, \mathbf{y}: d(\mathbf{x}, \mathbf{y}) \geq 0$ and $d(\mathbf{x}, \mathbf{y}) = 0$ if and only if $\mathbf{x} = \mathbf{y}$ such as l_2 , l_1 or [LPIPS \(learned perceptual image patch similarity\)](#) distance;
- $\lambda(\cdot) \in \mathbb{R}^+$ is a positive weighting function and the paper sets $\lambda(t_n) = 1$.

2) Consistency Training (CT)

- The other option is to train a consistency model independently. Note that in CD, a pre-trained score model $s_\phi(\mathbf{x}, t)$ is used to approximate the ground truth score $\nabla \log p_t(\mathbf{x})$ but in CT we need a way to estimate this score function and it turns out an unbiased estimator of $\nabla \log p_t(\mathbf{x})$ exists as $-\frac{\mathbf{x}_t - \mathbf{x}}{t^2}$.
- The consistency training (CT) loss is defined as follows:

$$\mathcal{L}_{\text{CT}}^N(\theta, \theta^-; \phi) = \mathbb{E}[\lambda(t_n) d(f_\theta(\mathbf{x} + t_{n+1} \mathbf{z}, t_{n+1}), f_{\theta^-}(\mathbf{x} + t_n \mathbf{z}, t_n))]$$

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

Speed up Diffusion Model Sampling

Consistency Models

According to the experiments in the paper, they found,

- **Heun ODE solver** works better than **Euler's first-order solver**, since **higher order ODE solvers** have smaller estimation errors with the same N .
- Among different options of the distance metric function $d(\cdot)$, the **LPIPS** metric works better than l_1 and l_2 distance.
- **Smaller N leads to faster convergence** but **worse samples**, whereas **larger N leads to slower convergence** but **better samples** upon convergence.

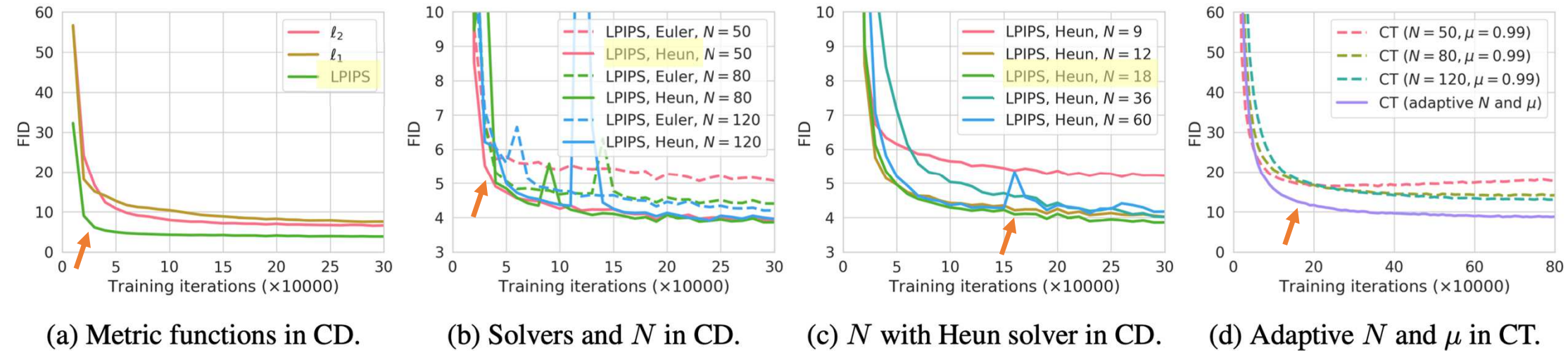


Fig. 12. Comparison of consistency models' performance under different configurations.

The **best configuration for CD** is **LPIPS distance metric, Heun ODE solver, and $N = 18$** . (Image source: Song et al., 2023)

Latent Variable Space

Latent diffusion model (LDM)

Latent diffusion model (LDM; Rombach & Blattmann, et al. 2022) runs the **diffusion process in the latent space** instead of pixel space, making training cost lower and inference speed faster.

- It is motivated by the observation that most bits of an image contribute to perceptual details and the semantic and conceptual composition still remains after aggressive compression.
- LDM loosely **decomposes the perceptual compression and semantic compression** with generative modeling learning by first **trimming off pixel-level redundancy with autoencoder** and then **manipulate/generate semantic concepts with diffusion process on learned latent**.

- **The perceptual compression process** relies on an **autoencoder** model. An encoder \mathcal{E} is used to compress the input image $\mathbf{x} \in \mathbb{R}^{H \times W \times 3}$ to a smaller 2D latent vector $\mathbf{z} = \mathcal{E}(\mathbf{x}) \in \mathbb{R}^{h \times w \times c}$, where the downsampling rate $f = H/h = W/w = 2^m$, $m \in \mathbb{N}$. Then a decoder \mathcal{D} reconstructs the images from the latent vector, $\tilde{\mathbf{x}} = \mathcal{D}(\mathbf{z})$.
- The paper explored **two types of regularization** in **autoencoder training to avoid arbitrarily high-variance in the latent spaces**.
 - KL-reg**: A small KL penalty towards a standard normal distribution over the learned latent, similar to VAE.
 - VQ-reg**: Uses a vector quantization layer within the decoder, like VQVAE but the quantization layer is absorbed by the decoder.

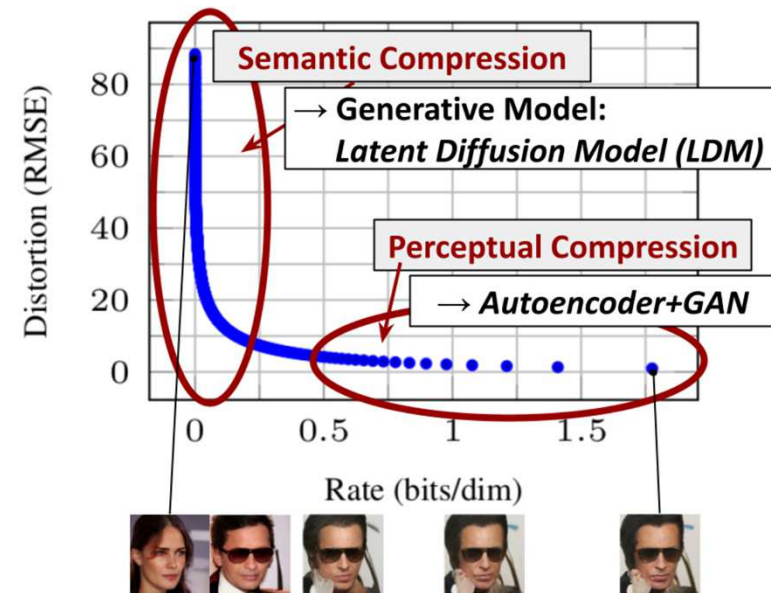


Fig. 8. The plot for tradeoff between compression rate and distortion, illustrating **two-stage compressions - perceptual and semantic compression**. (Image source: Rombach & Blattmann, et al. 2022)

VAE : <https://lilianweng.github.io/posts/2018-08-12-vae/>

VQVAE : <https://lilianweng.github.io/posts/2018-08-12-vae/#vq-vae-and-vq-vae-2>

Latent Variable Space

Latent diffusion model (LDM)

- The diffusion and denoising processes happen on the latent vector z .
- The denoising model is a **time-conditioned U-Net**, augmented with the **cross-attention mechanism** to handle flexible conditioning information for image generation (e.g. class labels, semantic maps, blurred variants of an image).
- The design is equivalent to **fuse representation of different modality into the model with cross-attention mechanism**.
- Each type of conditioning information is paired with a **domain-specific encoder** τ_θ to project the conditioning input y to an intermediate representation that can be mapped into cross-attention component, $\tau_\theta(y) \in \mathbb{R}^{M \times d_\tau}$:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right) \cdot \mathbf{V}$$

where $\mathbf{Q} = \mathbf{W}_Q^{(i)} \cdot \varphi_i(\mathbf{z}_i)$, $\mathbf{K} = \mathbf{W}_K^{(i)} \cdot \tau_\theta(y)$, $\mathbf{V} = \mathbf{W}_V^{(i)} \cdot \tau_\theta(y)$
 and $\mathbf{W}_Q^{(i)} \in \mathbb{R}^{d \times d_\epsilon^i}$, $\mathbf{W}_K^{(i)}, \mathbf{W}_V^{(i)} \in \mathbb{R}^{d \times d_\tau}$, $\varphi_i(\mathbf{z}_i) \in \mathbb{R}^{N \times d_\epsilon^i}$, $\tau_\theta(y) \in \mathbb{R}^{M \times d_\tau}$

$\varphi_i(\mathbf{z}_t) \in \mathbb{R}^{M \times d_\epsilon^i}$: a (flattened) intermediate representation of the UNet implementing ϵ_θ

[Robin Rombach, Andreas Blattmann, et al., "High-Resolution Image Synthesis with Latent Diffusion Models," CVPR 2022.](#)

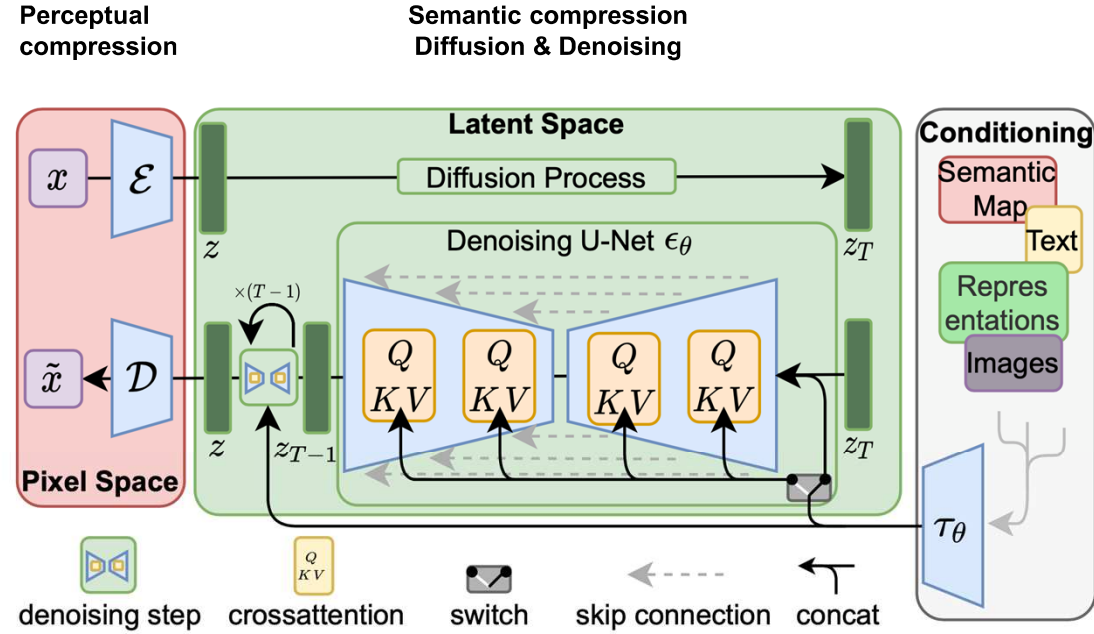


Fig. 14. The architecture of latent diffusion model.

- Based on image-conditioning pairs, we then learn the conditional LDM via

$$L_{LDM} := \mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_\theta(z_t, t, \tau_\theta(y))\|_2^2 \right]$$

where both τ_θ and ϵ_θ are jointly optimized.

Scale up Generation Resolution and Quality

Cascaded Diffusion Model – Noise Conditioning Augmentation

To generate high-quality images at high resolution, Ho et al. (2021) proposed to use a pipeline of multiple diffusion models at increasing resolutions.

- **Noise conditioning augmentation** between pipeline models is crucial to the final image quality, which is to apply strong data augmentation to the conditioning input \mathbf{z} of each super-resolution model $p_\theta(\mathbf{x}|\mathbf{z})$. The conditioning noise helps reduce compounding error in the pipeline setup.
- **U-net** is a common choice of model architecture in diffusion modeling for high-resolution image generation.

[Ho et al. \(2021\), Cascaded Diffusion Models for High Fidelity Image Generation](#)

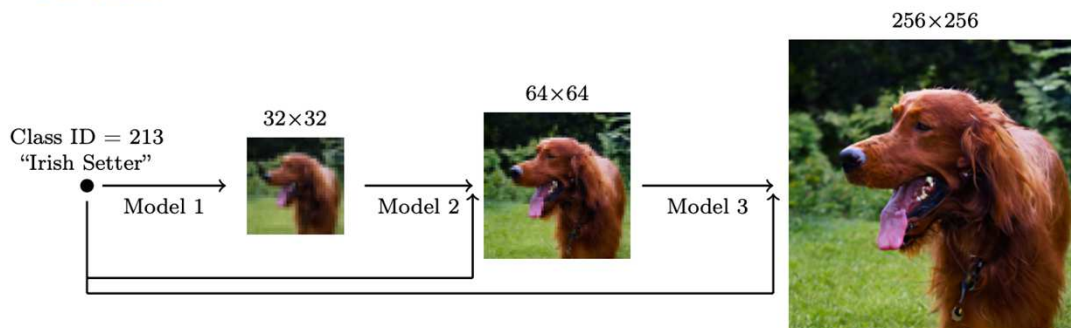


Fig. 15. A cascaded pipeline of multiple diffusion models at increasing resolutions. (Image source: Ho et al. 2021)]

- They found the most effective noise is to apply Gaussian noise at low resolution and Gaussian blur at high resolution.
- In addition, they also explored **two forms of conditioning augmentation** that require small modification to the training process. Note that conditioning noise is only applied to training but not at inference.
 - 1) **Truncated conditioning augmentation** stops the diffusion process early at step $t > 0$ for low resolution.
 - 2) **Non-truncated conditioning augmentation** runs the full low resolution reverse process until step 0 but then corrupt it by $\mathbf{z}_t \sim q(\mathbf{x}_t|\mathbf{x}_0)$ and then feeds the corrupted \mathbf{z}_t s into the super-resolution model.

Scale up Generation Resolution and Quality

Cascaded Diffusion Model – Noise Conditioning Augmentation

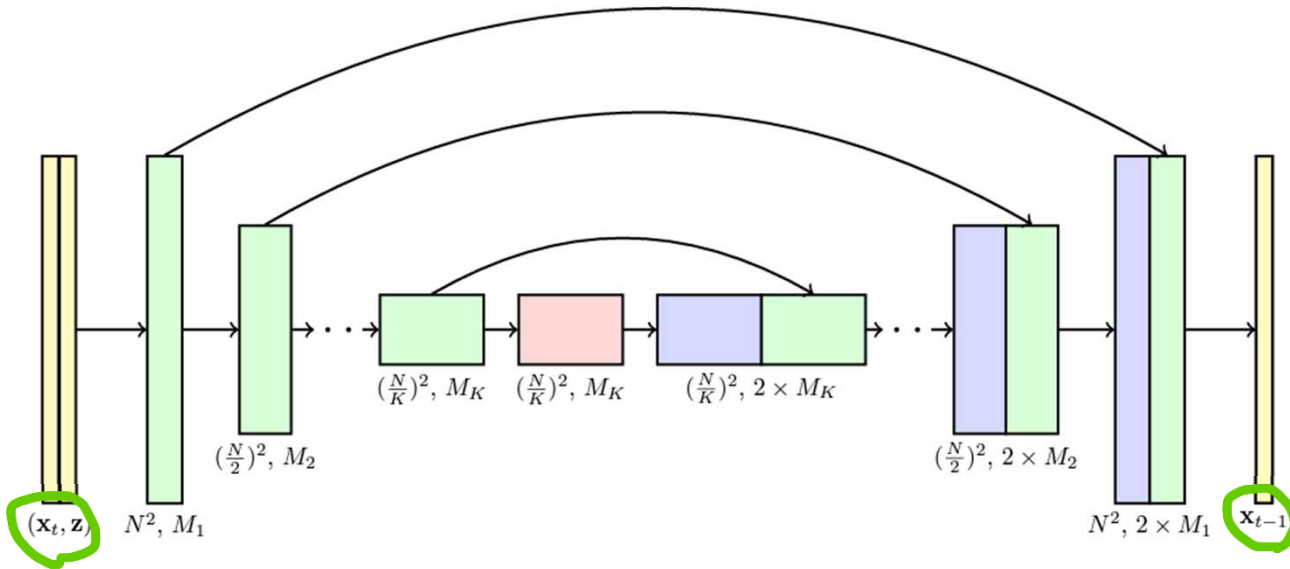


Fig 3. The **U-Net architecture** used in each model of a CDM pipeline.

The first model is a **class-conditional diffusion model** that receives the noisy image \mathbf{x}_t and the class label y and as input. (The class label y and timestep t are injected into each block as an embedding, not depicted here).

The remaining models in the pipeline are **class-conditional super-resolution models** that receive \mathbf{x}_t, y , and an additional upsampled low-resolution image \mathbf{z} as input. The downsampling/upsampling blocks adjust the image input resolution $N \times N$ by a factor of 2 through each of the K blocks. The channel count at each block is specified using channel multipliers M_1, M_2, \dots, M_K , and the upsampling pass has concatenation skip connections to the downsampling pass.

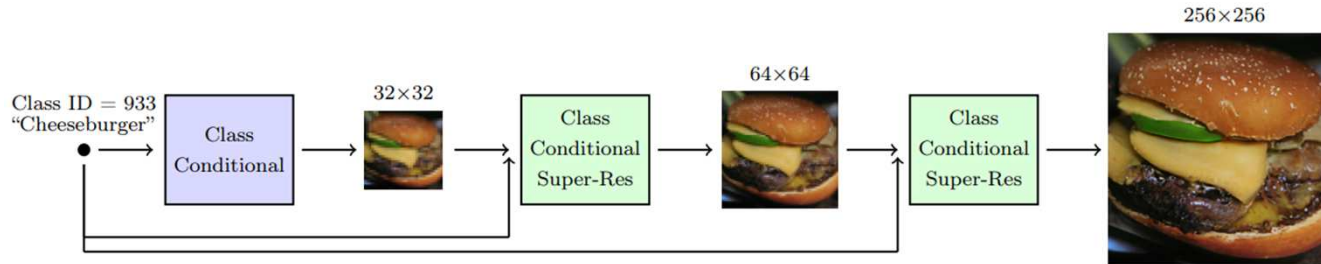


Fig 4. Detailed CDM pipeline for generation of class conditional 256×256 images. The first model is a class-conditional diffusion model, and it is followed by a sequence of two class-conditional super-resolution diffusion models. Each model has a U-Net architecture as depicted in Fig. 3.

Scale up Generation Resolution and Quality

unCLIP

The two-stage diffusion model **unCLIP** (Ramesh et al. 2022) heavily utilizes the CLIP text encoder to produce text-guided images at high quality.

- Given a pretrained CLIP model \mathbf{c} and paired training data for the diffusion model, (\mathbf{x}, \mathbf{y}) , where \mathbf{x} is an image and \mathbf{y} is the corresponding caption, we can compute the CLIP text and image embedding, $\mathbf{c}^i(\mathbf{y})$ and $\mathbf{c}^i(\mathbf{x})$, respectively. The unCLIP learns two models in parallel:

- ✓ A prior model $P(\mathbf{c}^i|\mathbf{y})$: outputs CLIP image embedding \mathbf{c}^i given the text \mathbf{y} .
- ✓ A decoder $P(\mathbf{x}|\mathbf{c}^i, [\mathbf{y}])$: generates the image \mathbf{x} given CLIP image embedding \mathbf{c}^i and optionally the original text \mathbf{y} .

These two models enable conditional generation, because

$$\underbrace{P(\mathbf{x}|\mathbf{y}) = P(\mathbf{x}, \mathbf{c}^i|\mathbf{y})}_{\mathbf{c}^i \text{ is deterministic given } \mathbf{x}} = P(\mathbf{x}|\mathbf{c}^i, \mathbf{y})P(\mathbf{c}^i|\mathbf{y})$$

[Ramesh et al. 2022, Hierarchical Text-Conditional Image Generation with CLIP Latents](#)

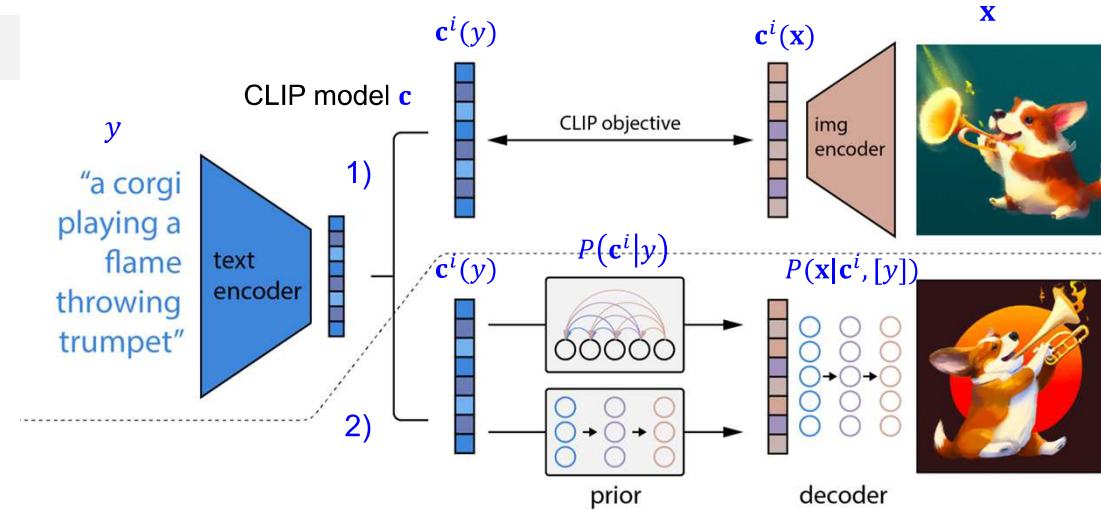


Fig. 12. The architecture of unCLIP. (Image source: Ramesh et al. 2022)]

unCLIP follows a two-stage image generation process:

- 1) Given a text \mathbf{y} , a CLIP model \mathbf{c} is first used to generate a text embedding $\mathbf{c}^i(\mathbf{y})$. Using CLIP latent space enables zero-shot image manipulation via text.
- 2) A diffusion or autoregressive prior $P(\mathbf{c}^i|\mathbf{y})$ processes this CLIP text embedding to construct an image prior and then a diffusion decoder $P(\mathbf{x}|\mathbf{c}^i, [\mathbf{y}])$ generates an image, conditioned on the prior. This decoder can also generate image variations conditioned on an image input, preserving its style and semantics.

Scale up Generation Resolution and Quality

Imagen

Instead of CLIP model, [Imagen](#) (Saharia et al. 2022) uses a pre-trained large LM (i.e. a frozen **T5-XXL** text encoder) to encode text for image generation. There is a general trend that larger model size can lead to better image quality and text-image alignment. They found that **T5-XXL** and **CLIP** text encoder achieve similar performance on MS-COCO, but human evaluation prefers T5-XXL on DrawBench (a collection of prompts covering 11 categories).

When applying classifier-free guidance, increasing w may lead to better image-text alignment but worse image fidelity. They found that it is due to train-test mismatch, that is saying, because training data x stays within the range $[-1,1]$, the test data should be so too. Two thresholding strategies are introduced:

- **Static thresholding**: clip x prediction to $[-1,1]$
- **Dynamic thresholding**: at each sampling step, compute s as a certain percentile absolute pixel value; if $s > 1$, clip the prediction to $[-s, s]$ and divide by s .

[Saharia et al. 2022, Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding](#)

Imagen modifies several designs in U-net to make it **efficient U-Net**.

- Shift model parameters from high resolution blocks to low resolution by adding more residual locks for the lower resolutions;
- Scale the skip connections by $1/\sqrt{2}$
- Reverse the order of downsampling (move it before convolutions) and upsampling operations (move it after convolution) in order to improve the speed of forward pass.

They found that **noise conditioning augmentation**, **dynamic thresholding** and **efficient U-Net** are critical for image quality, but **scaling text encoder size** is more important than U-Net size.

Scale up Generation Resolution and Quality

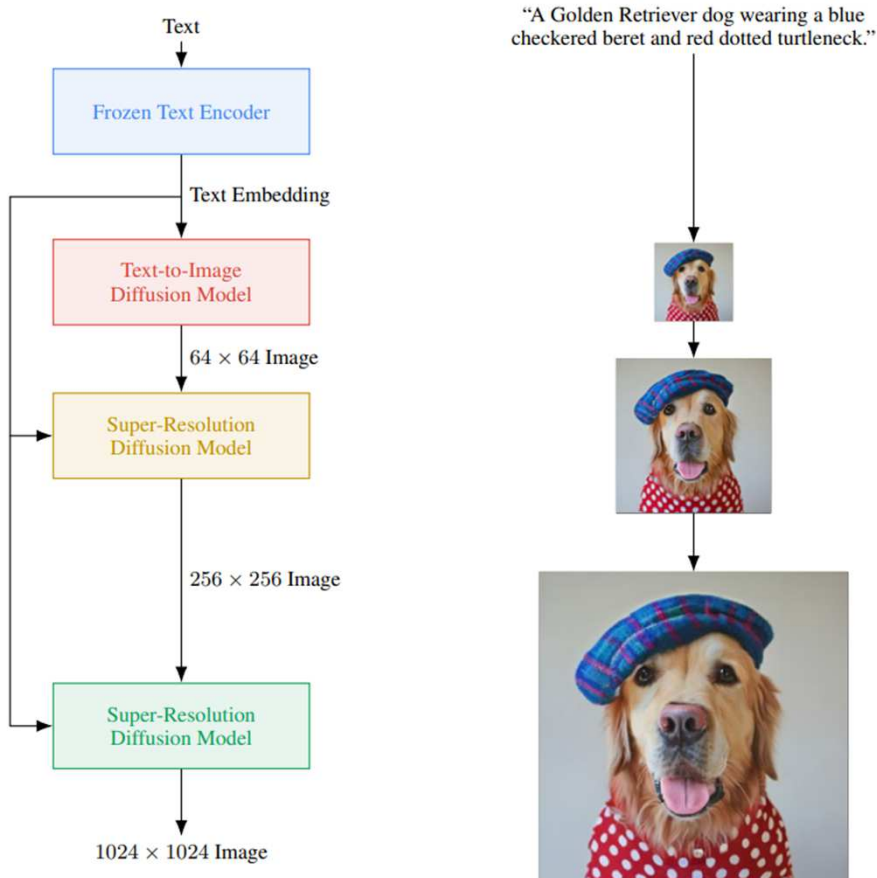
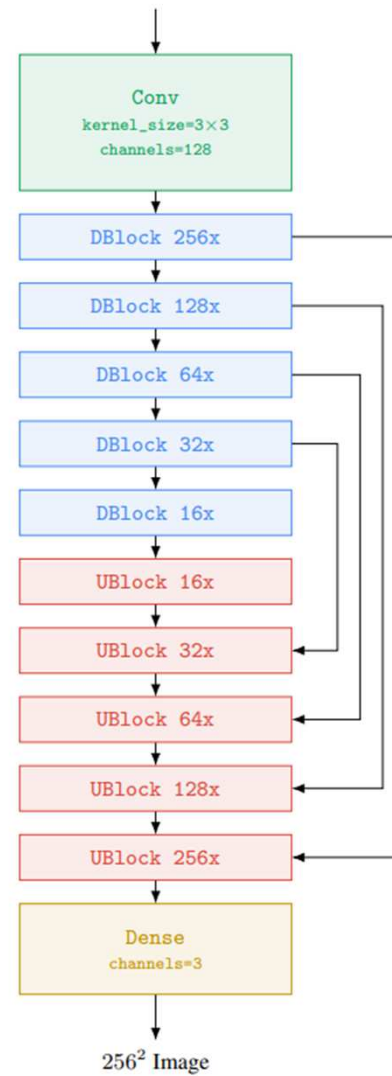
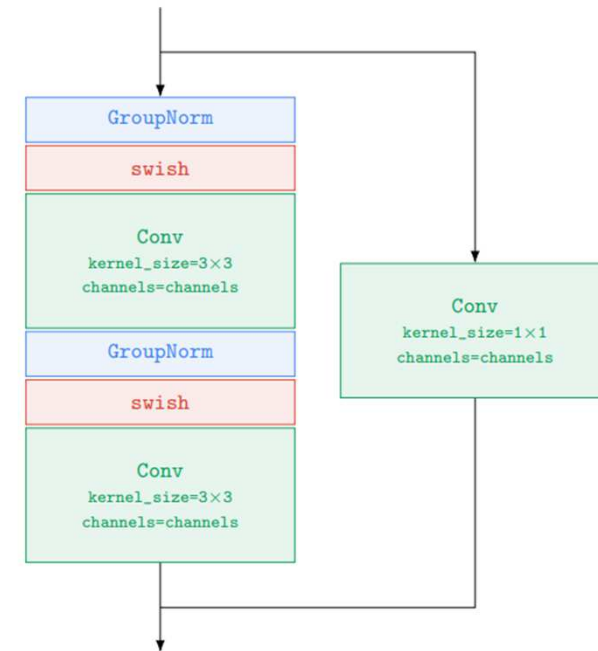


Figure A.4: Visualization of Imagen. Imagen uses a **frozen text encoder** to encode the input text into text embeddings. A **conditional diffusion model** maps the text embedding into a 64×64 image. Imagen further utilizes **text-conditional super-resolution diffusion models** to upsample the image, first $64 \times 64 \rightarrow 256 \times 256$, and then $256 \times 256 \rightarrow 1024 \times 1024$.



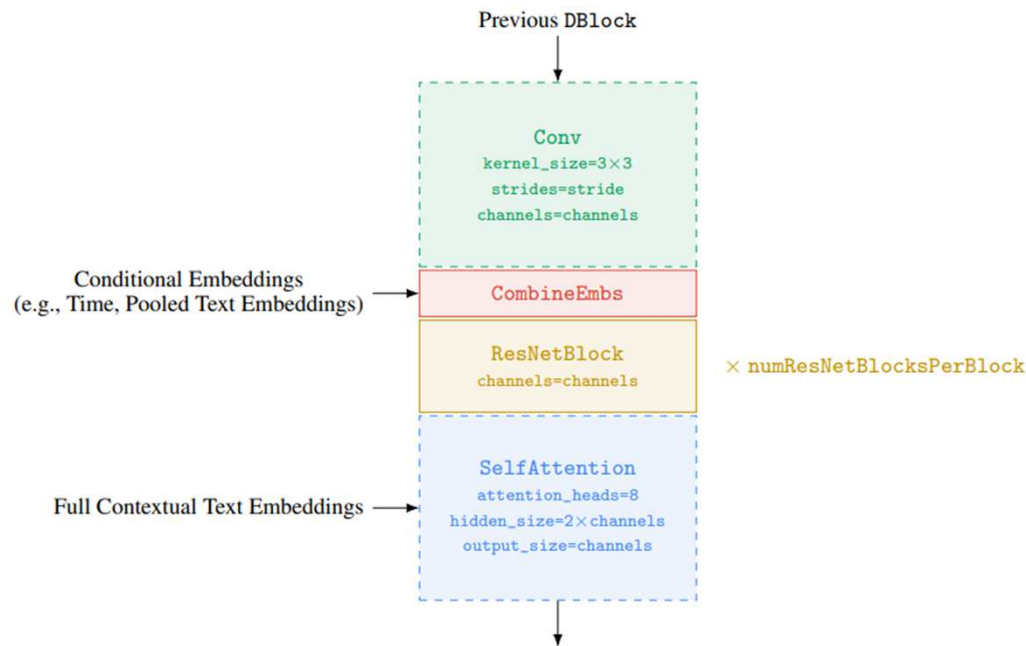
Efficient U-Net architecture for $64^2 \rightarrow 256^2$



Efficient U-Net *ResNetBlock*.

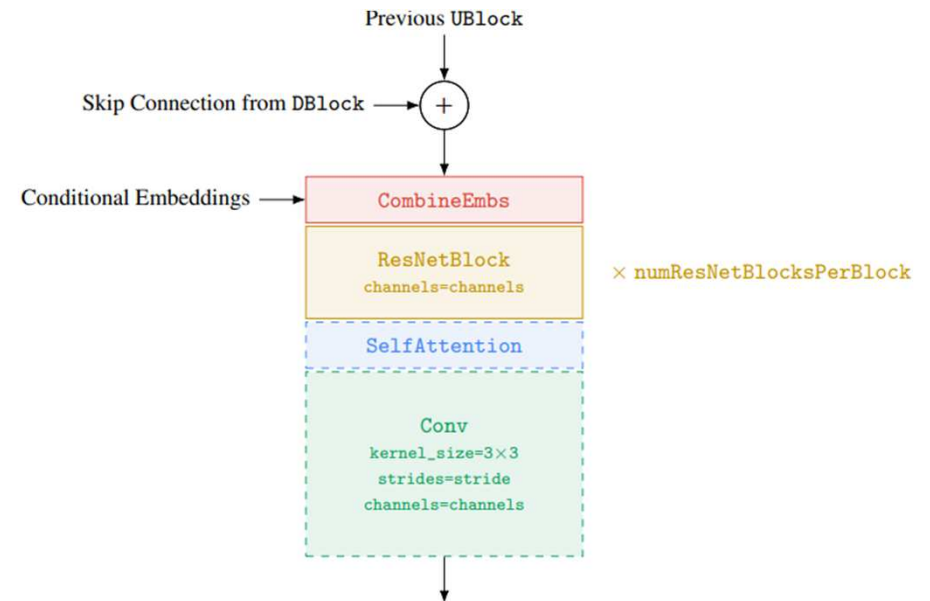
The *ResNetBlock* is used both by the *DBlock* and *UBlock*. Hyperparameter of the *ResNetBlock* is the number of channels: int.

Scale up Generation Resolution and Quality



Efficient U-Net *DBlock*.

Hyperparameters of *DBlock* are: the stride of the block if there is downsampling stride: `Optional[Tuple[int, int]]`, number of *ResNetBlock* per *DBlock* `numResNetBlocksPerBlock`: `int`, and number of channels `channels`: `int`. The dashed lined blocks are optional, e.g., not every *DBlock* needs to downsample or needs self-attention



Efficient U-Net *UBlock*.

Hyperparameters of *UBlock* are: the stride of the block if there is upsampling stride: `Optional[Tuple[int, int]]`, number of *ResNetBlock* per *DBlock* `numResNetBlocksPerBlock`: `int`, and number of channels `channels`: `int`. The dashed lined blocks are optional, e.g., not every *UBlock* needs to upsample or needs self-attention.

Model Architecture

There are two common backbone architecture choices for diffusion models: U-Net and Transformer.

U-Net

[U-Net \(Ronneberger, et al. 2015\)](#) consists of a downsampling stack and an upsampling stack.

- *Downsampling*: Each step consists of the repeated application of two 3x3 convolutions (unpadded convolutions), each followed by a ReLU and a 2x2 max pooling with stride 2. At each downsampling step, the number of feature channels is doubled.
- *Upsampling*: Each step consists of an upsampling of the feature map followed by a 2x2 convolution and each halves the number of feature channels.
- *Shortcuts*: Shortcut connections result in a concatenation with the corresponding layers of the downsampling stack and provide the essential high-resolution features to the upsampling process.

[Olaf Ronneberger et al., "U-Net: Convolutional Networks for Biomedical Image Segmentation," , 2015](#)

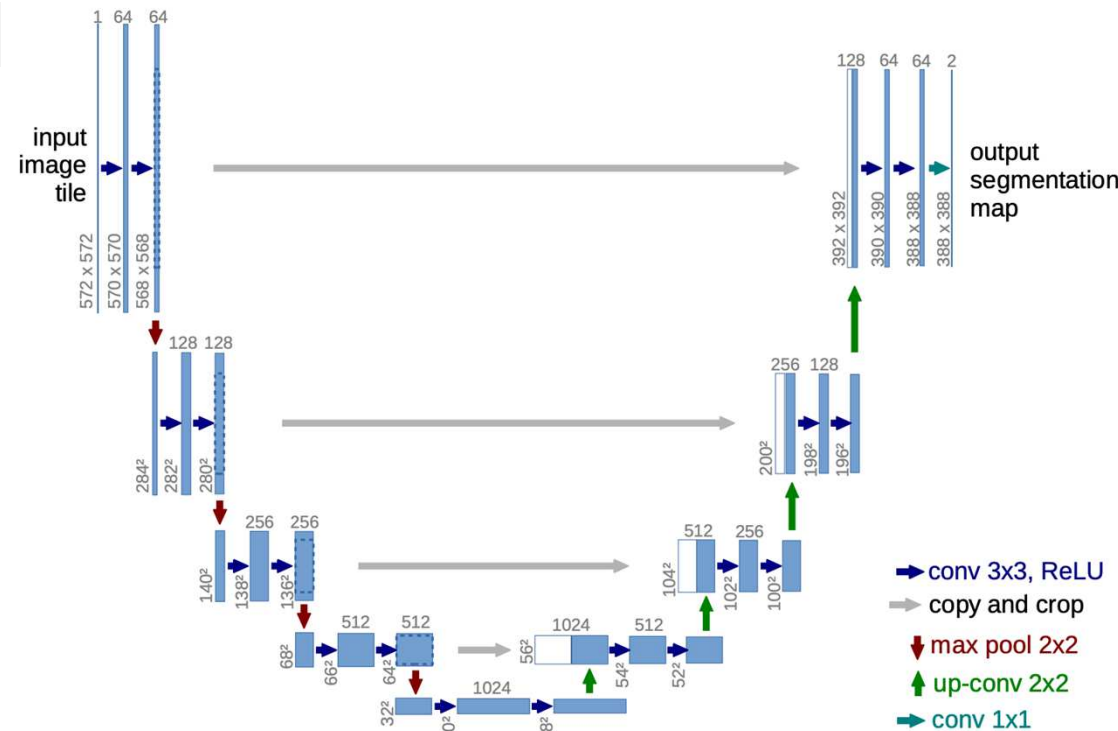


Fig. 17. The U-net architecture. Each blue square is a feature map with the number of channels labeled on top and the height x width dimension labeled on the left bottom side. The gray arrows mark the shortcut connections. (Image source: Ronneberger, 2015)

Model Architecture

ControlNet

To enable image generation conditioned on additional images for **composition info** like **Canny edges, Hough lines, user scribbles, human post skeletons, segmentation maps, depths and normals**, [ControlNet \(Zhang et al. 2023\)](#) introduces architectural changes via adding a “sandwiched” zero convolution layers of a trainable copy of the original model weights into each encoder layer of the U-Net.

Precisely, given a neural network block $\mathcal{F}_\theta(\cdot)$, ControlNet does the following:

- First, freeze the original parameters θ of the original block
- Clone it to be a copy with trainable parameters θ_c and an additional conditioning vector \mathbf{c} .
- Use two zero convolution layers, denoted as $\mathcal{Z}_{\theta_{z1}}(\cdot; \cdot)$ and $\mathcal{Z}_{\theta_{z2}}(\cdot; \cdot)$, which is 1x1 conv layers with both weights and biases initialized to be zeros, to connect these two blocks. *Zero convolutions protect this back-bone by eliminating random noise as gradients in the initial training steps.*
- The final output is:

$$\mathbf{y}_c = \mathcal{F}_\theta(\mathbf{x}) + \mathcal{Z}_{\theta_{z2}}(\mathcal{F}_{\theta_c}(\mathbf{x} + \mathcal{Z}_{\theta_{z1}}(\mathbf{c})))$$

[Lvmin Zhang et al., “Adding Conditional Control to Text-to-Image Diffusion Models,” 2023](#)

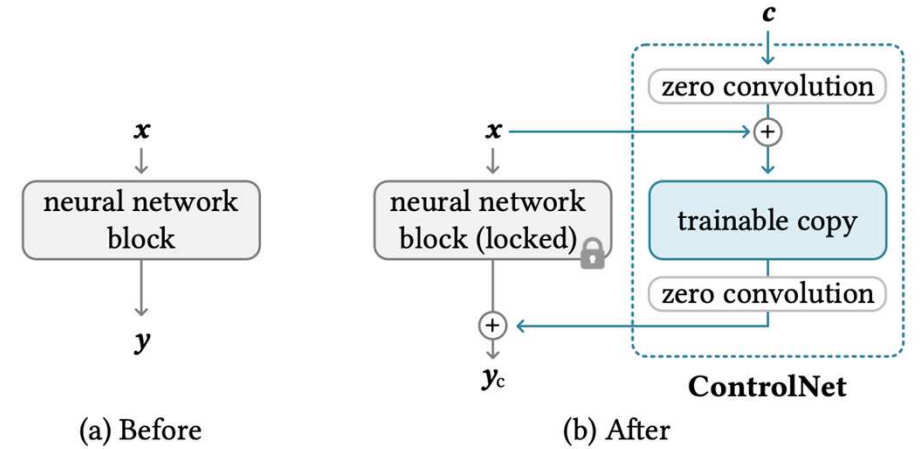
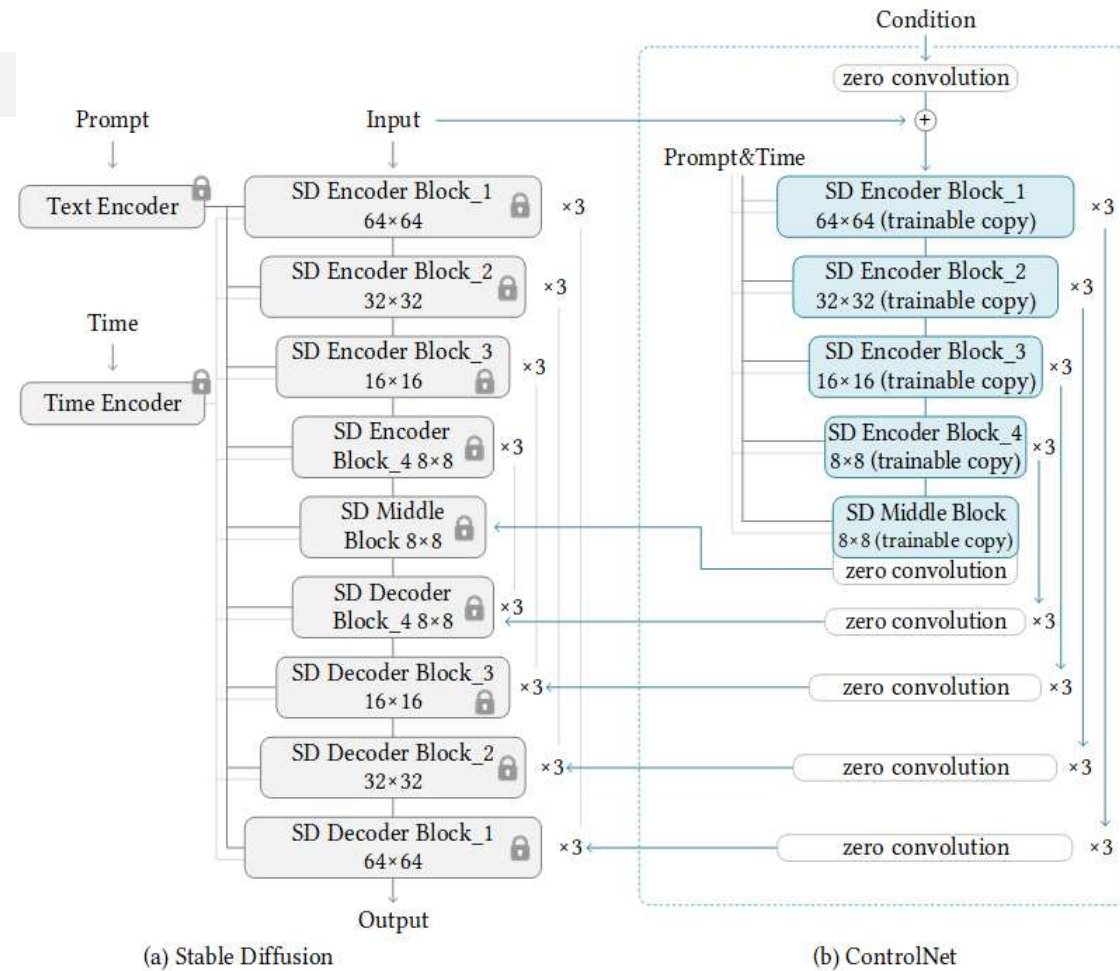


Fig. 18. The ControlNet architecture. (Image source: Zhang et al. 2023)

Model Architecture

ControlNet

- (a) Stable Diffusion 사전 훈련된 Diffusion 모델(예: Stable Diffusion의 latent UNet)의 매개변수를 복제하여 “학습 가능한 복사본”과 “잠긴 복사본(자물쇠)”으로 나눈다.
- 잠긴 복사본은 대규모 데이터셋에서 학습한 방대한 지식을 보존하며, 학습 가능한 복사본은 과제별 측면을 학습하는 데 사용됨
- 학습 가능한 복사본과 잠긴 복사본의 매개변수는 “zero convolution” 레이어로 연결됨. 이것은 ControlNet 프레임워크의 일부로 최적화되며, 새로운 조건이 학습될 때 이미 학습된 의미를 보존하는 학습 기법임



모두의 연구소, ControlNet 으로 이미지생성 (feat. Stable Diffusion)

<https://modulabs.co.kr/blog/controlnet-stable-diffusion/>

출처 : <https://huggingface.co/blog/controlnet>

Model Architecture

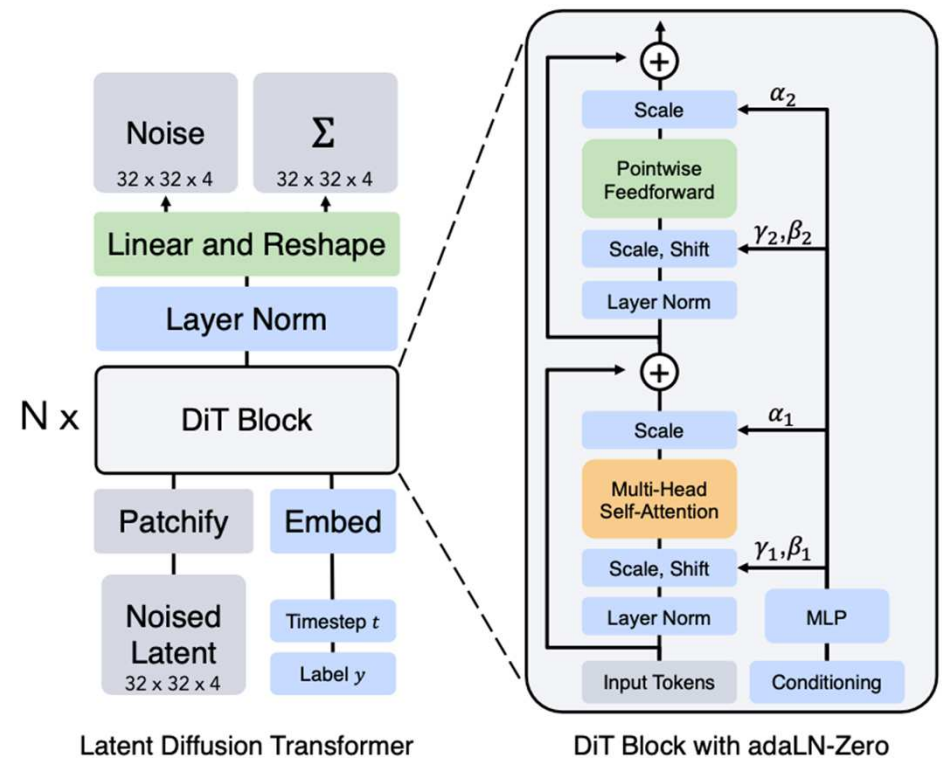
Diffusion Transformer (DiT)

[Diffusion Transformer \(DiT; Peebles & Xie, 2023\)](#) for diffusion modeling operates on [latent patches](#), using [the same design space of LDM](#) (Latent Diffusion Model)]. DiT has the following setup:

- 1) Take the latent representation of an input z as input to DiT.
- 2) “Patchify” the noise latent of size $I \times I \times C$ into patches of size p and convert it into a sequence of patches of size $(I/p)^2$.
- 3) Then this sequence of tokens go through Transformer blocks. They are exploring three different designs for how to do generation conditioned on contextual information like timestep t or class label c . Among three designs, *adaLN* (*Adaptive layer norm*)-Zero works out the best, better than in-context conditioning and cross-attention block. The scale and shift parameters, γ and β , are regressed from the sum of the embedding vectors of t and c . The dimension-wise scaling parameters α is also regressed and applied immediately prior to any residual connections within the DiT block.
- 4) The transformer decoder outputs noise predictions and an output diagonal covariance prediction.

[William Peebles and Saining Xie, “Scalable Diffusion Models with Transformers,” 2023](#)

Transformer architecture can be easily scaled up and it is well known for that. This is one of the biggest benefits of DiT as its performance scales up with more compute and larger DiT models are more compute efficient according to the experiments.



Quick Summary

Pros

Tractability and **flexibility** are two conflicting objectives in generative modeling. **Tractable models** can be analytically evaluated and cheaply fit data (e.g. via a Gaussian or Laplace), but they cannot easily describe the structure in rich datasets. **Flexible models** can fit arbitrary structures in data, but evaluating, training, or sampling from these models is usually expensive.

Diffusion models are both *analytically tractable* and *flexible*

Cons

Diffusion models *rely on a long Markov chain of diffusion steps to generate samples*, so it can be quite expensive in terms of time and compute.

New methods have been proposed to make the process much faster, but the sampling is still slower than GAN.