

AI Applications Lecture 13

Image Generation AI

Reverse Diffusion Process

SUZUKI, Atsushi

Jing WANG

October 31, 2025

Contents

1	Introduction	3
1.1	Recap of Previous Lecture	3
1.2	Learning Outcomes	3
2	Preparation: Mathematical Notations	3
3	Revisiting the Goal of the Low-Resolution Latent Image Generator	5
4	Motivation for Introducing the Reverse Diffusion Process	5
4.1	Problem Formulation	5
4.2	Framework of Pushforward Measure	6
4.3	Distribution Approximation is Not Input-Output Correspondence Learning .	7
4.4	Difficulty of Naive Likelihood Maximization and Flow-based Models	7
4.5	Positioning of the Reverse Diffusion Process	8
5	Reverse Diffusion Process: Neural Network Inputs, Outputs, and Conditions	8
6	Denoising Schedulers (DDPM, Euler a, DPM 2+ Karras)	10
6.1	DDPM (Variance-Preserving Discrete Process)	10
6.2	Euler a (Euler Ancestral)	12
6.3	DPM 2+ (Karras Sigma Sequence)	14
7	Advantages of Each Scheduler (Practical Perspective)	17

8 Summary and Next Lecture Preview **17**

8.1 Summary Corresponding to Learning Outcomes 17

8.2 Next Lecture Preview 18

1 Introduction

1.1 Recap of Previous Lecture

In the previous lecture, we learned that a **natural image decoder**, such as those realized by **Variational Autoencoders (VAE)** [6, 8], can construct high-resolution **natural images** from low-resolution **latent images**. Specifically, we reconfirmed that the VAE decoder $\text{Dec}_\gamma : \mathcal{Z} \rightarrow \mathcal{I}$ can output high-resolution images in the pixel space $\mathcal{I} \subset \mathbb{R}^{H \times W \times C}$ given inputs from the latent space \mathcal{Z} .

1.2 Learning Outcomes

Upon completion of this lecture, students should be able to explain and execute the following:

- Explain **what the inputs and outputs** of the **reverse diffusion process** are in the text-to-image pipeline.
- Explain the **conditions that the neural network** constituting the reverse diffusion process must **satisfy**.
- Explain how the **output of the text encoder** is handled in the reverse diffusion process and how it **controls generation**.
- Explain the differences between **denoising schedulers** and be able to appropriately select and **execute** the reverse diffusion process according to the situation.

2 Preparation: Mathematical Notations

- **Definition:**

- (LHS) := (RHS): Indicates that the left-hand side is defined by the right-hand side. For example, $a := b$ indicates that a is defined as b .

- **Set:**

- Sets are often denoted by uppercase calligraphic letters. Example: \mathcal{A} .
- $x \in \mathcal{A}$: Indicates that the element x belongs to the set \mathcal{A} .
- $\{\}$: The empty set.
- $\{a, b, c\}$: The set consisting of elements a, b, c (set-builder notation, extension).
- $\{x \in \mathcal{A} \mid P(x)\}$: The set of elements in \mathcal{A} for which the proposition $P(x)$ is true (set-builder notation, intension).
- $|\mathcal{A}|$: The number of elements in the set \mathcal{A} (used only for finite sets in this lecture).

- \mathbb{R} : The set of all real numbers. $\mathbb{R}_{>0}$, $\mathbb{R}_{\geq 0}$, etc., are defined similarly.
- \mathbb{Z} : The set of all integers. $\mathbb{Z}_{>0}$, $\mathbb{Z}_{\geq 0}$, etc., are defined similarly.
- $[1, k]_{\mathbb{Z}}$: For $k \in \mathbb{Z}_{>0} \cup \{+\infty\}$, if $k < +\infty$ then $\{1, \dots, k\}$, if $k = +\infty$ then $\mathbb{Z}_{>0}$.

• **Function:**

- $f : \mathcal{X} \rightarrow \mathcal{Y}$ denotes that f is a map.
- $y = f(x)$ denotes the output $y \in \mathcal{Y}$ for the input $x \in \mathcal{X}$.

• **Vector:**

- Vectors are denoted by bold italic lowercase letters. Example: \mathbf{v} . $\mathbf{v} \in \mathbb{R}^n$.
- The i -th component is written as v_i :

$$\mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}. \quad (1)$$

- Standard inner product:

$$\langle \mathbf{u}, \mathbf{v} \rangle := \sum_{i=1}^{d_{\text{emb}}} u_i v_i. \quad (2)$$

• **Sequence:**

- $\mathbf{a} : [1, n]_{\mathbb{Z}} \rightarrow \mathcal{A}$ is called a sequence of length n . If $n < +\infty$, $\mathbf{a} = (a_1, \dots, a_n)$; if $n = +\infty$, $\mathbf{a} = (a_1, a_2, \dots)$.
- The length is written as $|\mathbf{a}|$.

• **Matrix:**

- Matrices are denoted by bold italic uppercase letters. Example: $\mathbf{A} \in \mathbb{R}^{m,n}$.
- The elements are $a_{i,j}$, and

$$\mathbf{A} = \begin{bmatrix} a_{1,1} & \cdots & a_{1,n} \\ \vdots & \ddots & \vdots \\ a_{m,1} & \cdots & a_{m,n} \end{bmatrix}. \quad (3)$$

- Transpose $\mathbf{A}^{\top} \in \mathbb{R}^{n,m}$:

$$\mathbf{A}^{\top} = \begin{bmatrix} a_{1,1} & \cdots & a_{m,1} \\ \vdots & \ddots & \vdots \\ a_{1,n} & \cdots & a_{m,n} \end{bmatrix}. \quad (4)$$

- The transpose of a vector is

$$\mathbf{v}^\top = \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix}. \quad (5)$$

- **Tensor:**

- A tensor as a multi-dimensional array is denoted by an underlined bold italic uppercase letter $\underline{\mathbf{A}}$.

3 Revisiting the Goal of the Low-Resolution Latent Image Generator

The goal of the latent generator LatentGen_β is as follows:

- To receive the text encoder output $(\mathbf{c}^{[i]})_{i=1}^m$ and generate a **low-”resolution” latent image** $\mathbf{x} \in \mathcal{X}$ to be passed to the **natural image decoder**. The decoder Dec_γ can then convert this \mathbf{x} into a high-resolution natural image.
- Even for the same $(\mathbf{c}^{[i]})_{i=1}^m$ obtained from the same prompt, to be able to output a **diverse** set of candidate images by changing the **random seed** (pseudo-random number). Since text-to-image tasks tend to lack sufficient input information, **diversity** in the output is essential.

4 Motivation for Introducing the Reverse Diffusion Process

4.1 Problem Formulation

Given a text condition $(\mathbf{c}^{[i]})_{i=1}^m$, let $P_{(\mathbf{c}^{[i]})_{i=1}^m}$ denote the distribution that covers the corresponding **set of images** with high probability. The goal is to obtain an algorithm that can pseudo-**sample/simulate** \mathbf{x} such that

$$\mathbf{x} \sim P_{(\mathbf{c}^{[i]})_{i=1}^m}. \quad (6)$$

What is important is **data generation** itself, not explicitly obtaining the **formula of the distribution** (probability mass function/probability density function).

Here, it should be noted that the goal, although often confused, is not to obtain a specific representation of $P_{(\mathbf{c}^{[i]})_{i=1}^m}$ by its probability density function or probability mass function. Even if we obtain the representation of the probability mass function or probability density function, it does not mean we can immediately sample data according to that distribution. We need to re-recognize that our goal is data generation, not the representation of the distribution. Also, representing the probability mass function or probability density function is difficult.

For example, if we try to represent the support as a discrete object using a probability mass function, there are $256^{C \times W \times H}$ possible values, making it practically impossible to specify them.

If we try to represent the support as a continuous object using a probability density function, it is impossible naively due to infinite degrees of freedom. If we try to represent it directly with a parametric function, practically the only useful parametric probability model on a multidimensional vector space is the multivariate normal distribution.

It is also difficult to construct a probability mass function or probability density function using a neural network. This is also often confused; neural networks are suitable for representing complex functions, but that is when there are no constraints on the function. Representing a function that satisfies the condition of being a probability, i.e., summing to 1, is not easy, regardless of whether it's a neural network. In natural language processing, softmax was used to construct a probability mass function, but this was possible because the number of elements in the support was not that large. In the case of a probability density function, construction by softmax is itself impossible.

4.2 Framework of Pushforward Measure

For the reasons stated in the previous section, instead of directly constructing the probability mass function or probability density function, we take the method of transforming random numbers using a function, thereby implicitly sampling from a **pushforward measure**.

Let \mathcal{S} be the support, which has the same dimension as the output, and let λ be the base distribution on it (in practice, the standard multivariate normal distribution). We find a **measurable function**

$$f(\cdot) : \left(\prod_{i=1}^m \mathbb{R}^{d^{(i)}} \right) \times \mathcal{S} \rightarrow \mathcal{I} \quad (7)$$

and define the **pushforward** as

$$f\left((\mathbf{c}^{[i]})_{i=1}^m, \cdot\right) \# \lambda(A) := \lambda\left(\left\{ \epsilon \in \mathcal{S} \mid f\left((\mathbf{c}^{[i]})_{i=1}^m, \epsilon\right) \in A \right\}\right), \quad A \subset \mathcal{I}. \quad (8)$$

The target condition is

$$f\left((\mathbf{c}^{[i]})_{i=1}^m, \cdot\right) \# \lambda \approx P_{(\mathbf{c}^{[i]})_{i=1}^m}. \quad (9)$$

If this is achieved, we can realize (6) simply by generating $\mathbf{z} \sim \lambda$ with pseudo-random numbers and setting

$$\mathbf{x} = f\left((\mathbf{c}^{[i]})_{i=1}^m, \mathbf{z}\right). \quad (10)$$

4.3 Distribution Approximation is Not Input-Output Correspondence Learning

What is given is the **distribution on the output side**, and the **correct answer** for the input z is not uniquely determined. For example, with Gaussian input, composing a rotation centered at the origin in the input space leaves the output distribution invariant. To put it more formally, suppose the pushforward distribution $f\#\lambda_{\text{normal}}$ of the standard Gaussian distribution λ_{normal} by some bijective map $f : S \rightarrow S$ is ideal. In this case, an output $x \in S$ corresponds to $z = f^{-1}(x)$. However, if we consider a rotation map Rot centered at the origin, $(f \circ \text{Rot})\#\lambda_{\text{normal}}$ is equal to $f\#\lambda_{\text{normal}}$, so it is also ideal. In this case, x corresponds to $z' = \text{Rot}^{-1} \circ f^{-1}(x)$. We cannot say that z is the correct input, nor that z' is the correct input. Therefore, it is not a simple "input-output correspondence regression problem".

Therefore, if we want to formulate it as learning an input-output relationship, we need to determine the "input" $z^{(i)}$ for each output $x^{(i)}$, which is not uniquely determined in principle, and then learn the input-output relationship. Although it might seem mathematically simpler due to the freedom in determining the input, if the determination of the "input" is poor, the input-output relationship becomes complicated, and a slight change in the input (inputting values not used during training is essential for achieving diverse outputs, so this case must always be considered) can lead to a large change in the output, resulting in an unstable distribution unsuitable for practical use. Therefore, it becomes a difficult problem of not just learning the input-output relationship, but also determining the "input" appropriately in some sense. In fact, many distribution generation problems in neural networks determine the "input" in some sense. Autoencoders, including the VAE explained last time, are also based on this idea.

4.4 Difficulty of Naive Likelihood Maximization and Flow-based Models

Now, since we want to select the optimal function on a computer, we use a parametric function $f_{(\cdot)}$. If we try to achieve (9) using this parametric function and optimize the log-likelihood directly with a data sequence $\{\mathbf{x}^{(n)}\}_{n=1}^N$, we need the **inverse map** f_{θ}^{-1} and the **Jacobian determinant** (change of variables rule for invertible transformations).

Here, suppose the observed data $\mathbf{x}^{(n)} \in \mathbb{R}^d$ are independent and identically distributed, and are mapped to the base distribution p_Z by $z = f_{\theta}^{-1}(x)$. Then the likelihood of each sample is

$$p_X(\mathbf{x}^{(n)} \mid \theta) = p_Z(f_{\theta}^{-1}(\mathbf{x}^{(n)})) \left| \det \frac{\partial f_{\theta}^{-1}(\mathbf{x}^{(n)})}{\partial \mathbf{x}^{\top}} \right|, \quad (11)$$

so the log-likelihood of all data is

$$\mathcal{L}(\theta) = \sum_{n=1}^N \log p_X(\mathbf{x}^{(n)} | \theta) = \sum_{n=1}^N \left\{ \log p_Z(f_{\theta}^{-1}(\mathbf{x}^{(n)})) + \log \left| \det \frac{\partial f_{\theta}^{-1}(\mathbf{x}^{(n)})}{\partial \mathbf{x}^{\top}} \right| \right\}. \quad (12)$$

In this case, the gradient (first derivative with respect to the parameters) is given by

$$\nabla_{\theta} \mathcal{L}(\theta) = \sum_{n=1}^N \left[\nabla_{\theta} \log p_Z(f_{\theta}^{-1}(\mathbf{x}^{(n)})) + \nabla_{\theta} \log \left| \det \frac{\partial f_{\theta}^{-1}(\mathbf{x}^{(n)})}{\partial \mathbf{x}^{\top}} \right| \right], \quad (13)$$

Here, the first term is the derivative of the log-density of the base distribution p_Z through f_{θ}^{-1} , and the second term originates from the Jacobian term of the variable transformation. Therefore, calculating this gradient generally requires the explicit calculation of f_{θ}^{-1} and its Jacobian determinant.

However, calculating f_{θ}^{-1} for a general neural network is difficult. To avoid this, there are **flow-based** models (RealNVP [1], Glow [5]) that impose constraints of **invertibility and efficiently computable Jacobian**, but this tends to sacrifice **freedom in architecture design**. One of the advantages of neural networks is the flexibility in design, being able to freely design the graph structure according to the application field, so losing that is a pain in practice.

4.5 Positioning of the Reverse Diffusion Process

From the above, distribution approximation requires **non-trivial methods**. One of them is the **diffusion model**. Here, we will not go into the details of the **training stage**, but formulate the **reverse diffusion process** as the **inference stage (data generation stage)**. It is called "reverse" because during training, the meaning of the **input** is established through **forward diffusion (diffusion process)**, and during generation, the **reverse** of that is traced [2, 4]. The **denoising scheduler** provides the specific **implementation procedure** for the reverse diffusion process.

5 Reverse Diffusion Process: Neural Network Inputs, Outputs, and Conditions

The reverse diffusion process is an algorithm that achieves the target distribution by repeatedly applying a function that is close to, but slightly different from, the identity function (meaning the input and output do not change much). Each time the function is applied, the input moves slightly. If we view this as a process where the distribution is deformed by the composite function, the initial input Gaussian distribution is slightly deformed by the pushforward each time the function is applied. It is expected that this will eventually approach the target distribution.

To explain the construction using a neural network more concretely, it is as follows. There are $k = 0, 1, \dots, N - 1$ and a corresponding monotonically decreasing sequence called timesteps $\tau_0 > \tau_1 > \dots > \tau_{N-1}$, and

- At each step, the neural network receives the output $x^{(k)}$ from the previous step, the input $(c^{[i]})_{i=1}^m$ from the text encoder, and τ_k , and determines the negative direction $\widehat{\epsilon}^{(k)}$ in which $x^{(k)}$ moves.
- Roughly, the point is moved as $x^{(k+1)} \leftarrow x^{(k)} - \gamma_t \widehat{\epsilon}^{(k)}$. How γ_t is determined and the detailed behavior vary depending on the algorithm called the **denoising scheduler**. In some cases, Gaussian noise is further added.

Definition 5.1 (Input/Output Types of the Network Used in the Reverse Diffusion Process). Let $d_{\text{latent}} \in \mathbb{Z}_{>0}$ be the desired latent dimension, $d_{\text{context}} \in \mathbb{Z}_{>0}$ be the text condition dimension, and \mathcal{T} be the set of timesteps. Here \mathcal{T} is the set of **discrete timesteps** generated by the scheduler before inference, given by (21), (38), (54) described later. The **noise predictor** (or **data predictor**) f_θ in the reverse diffusion process is

$$f_\theta : \mathbb{R}^{d_{\text{latent}}} \times \mathbb{R}^{d_{\text{context}}} \times \mathcal{T} \rightarrow \mathbb{R}^{d_{\text{latent}}} \quad (14)$$

and the input is (i) the current sample $x_t \in \mathbb{R}^{d_{\text{latent}}}$, (ii) the text encoder output (including both positive and negative) $c \in \mathbb{R}^{d_{\text{context}}}$, and (iii) the timestep index $t \in \mathcal{T}$. The output is a tensor of the **same dimension** (e.g., **predicted noise** $\widehat{\epsilon}$).

Definition 5.2 (Sequential Application of Classifier-Free Guidance). Given $x_t \in \mathbb{R}^{d_{\text{latent}}}$, positive embedding $c_+ \in \mathbb{R}^{d_{\text{context}}}$, negative (unconditional) embedding $c_- \in \mathbb{R}^{d_{\text{context}}}$, and timestep $t \in \mathcal{T}$. The composition of model outputs in CFG is defined as

$$\widehat{\epsilon}_{\text{cond}} := f_\theta(x_t, c_+, t), \quad (15)$$

$$\widehat{\epsilon}_{\text{uncond}} := f_\theta(x_t, c_-, t), \quad (16)$$

$$\widehat{\epsilon}_{\text{CFG}} := \widehat{\epsilon}_{\text{uncond}} + s(\widehat{\epsilon}_{\text{cond}} - \widehat{\epsilon}_{\text{uncond}}), \quad (17)$$

where $s \in \mathbb{R}_{\geq 1}$ is the **guidance scale**. This definition allows CFG to be applied without breaking the function type (14) of f_θ .

Remark 5.1. In practice, implementations often batch x_t and c_\pm together to compute (15)–(16) simultaneously in a single forward pass, but the mathematical type is always (14). In this document, we will henceforth describe it in the form of applying twice and then composing, as in (15)–(17). The timestep t called at this time is the discrete timestep given by the scheduler.

6 Denoising Schedulers (DDPM, Euler a, DPM 2+ Karras)

Hereafter, focusing only on the inference stage, we will strictly provide the step update equations of representative schedulers in a one-to-one correspondence with the relevant code in the implementation (Hugging Face Diffusers). The original papers are also cited [2, 4, 7].

6.1 DDPM (Variance-Preserving Discrete Process)

Definition 6.1 (Reverse Diffusion Algorithm by DDPM). The following data are given.

- **Quantities determined during training:**

- Number of training timesteps $T \in \mathbb{Z}_{>0}$.
- Beta sequence $\beta = (\beta_0, \dots, \beta_{T-1}) \in (0, 1)^T$ and derived from it

$$\alpha_t := (1 - \beta_t), \quad \bar{\alpha}_t := \prod_{s=0}^t (1 - \beta_s), \quad t = 0, 1, \dots, T-1. \quad (18)$$

For convenience, let $\bar{\alpha}_{-1} := 1$.

- Noise prediction network trained with epsilon prediction type

$$f_\theta : \mathbb{R}^{d_{\text{latent}}} \times \mathbb{R}^{d_{\text{context}}} \times \{0, 1, \dots, T-1\} \rightarrow \mathbb{R}^{d_{\text{latent}}}. \quad (19)$$

- **Quantities determined by other networks:**

- Positive embedding $c_+ \in \mathbb{R}^{d_{\text{context}}}$ obtained from the text encoder.
- Negative (unconditional) embedding $c_- \in \mathbb{R}^{d_{\text{context}}}$ obtained from the text encoder.

- **Quantities specified by the user at inference time:**

- Number of inference timesteps $N \in \mathbb{Z}_{>0}$ satisfying $N \leq T$.
- Guidance scale $s \in \mathbb{R}_{\geq 1}$.

At this time, the reverse diffusion algorithm by DDPM is as follows.

(0) Initialization. Set the step ratio as

$$r := \left\lfloor \frac{T}{N} \right\rfloor \in \mathbb{Z}_{>0} \quad (20)$$

and define the discrete timestep sequence as

$$t_k := r(N - 1 - k), \quad k = 0, 1, \dots, N-1. \quad (21)$$

The initial sample is generated by

$$\mathbf{x}_{t_0} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_{\text{latent}}}). \quad (22)$$

(1) Main loop. For $k = 0, 1, \dots, N - 2$, perform the following.

- Let the current time be t_k and the next time be t_{k+1} . In this case, $t_{k+1} < t_k$.
- Let the current sample be $\mathbf{x}_{t_k} \in \mathbb{R}^{d_{\text{latent}}}$.
- Composite the noise prediction by CFG:

$$\hat{\boldsymbol{\epsilon}}_{\text{cond}} := f_{\boldsymbol{\theta}}(\mathbf{x}_{t_k}, \mathbf{c}_+, t_k), \quad (23)$$

$$\hat{\boldsymbol{\epsilon}}_{\text{uncond}} := f_{\boldsymbol{\theta}}(\mathbf{x}_{t_k}, \mathbf{c}_-, t_k), \quad (24)$$

$$\hat{\boldsymbol{\epsilon}}_{t_k} := \hat{\boldsymbol{\epsilon}}_{\text{uncond}} + s(\hat{\boldsymbol{\epsilon}}_{\text{cond}} - \hat{\boldsymbol{\epsilon}}_{\text{uncond}}). \quad (25)$$

- Set the variance as

$$\sigma_{t_k}^2 := \frac{1 - \bar{\alpha}_{t_{k+1}}}{1 - \bar{\alpha}_{t_k}} \beta_k \quad (26)$$

and using Gaussian noise

$$\mathbf{z}_{t_k} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_{\text{latent}}}) \quad (27)$$

update as

$$\mathbf{x}_{t_{k+1}} := \frac{1}{\sqrt{\bar{\alpha}_{t_k}}} \mathbf{x}_{t_k} - \frac{\beta_{t_k}}{\sqrt{\bar{\alpha}_{t_k}} \sqrt{1 - \bar{\alpha}_{t_k}}} \hat{\boldsymbol{\epsilon}}_{t_k} + \sigma_{t_k} \mathbf{z}_{t_k}. \quad (28)$$

Remark 6.1 (β sequence in Stable Diffusion 1.5). In the public config file for Stable Diffusion 1.5^a, the β sequence is defined by the **scaled linear** method as

$$\beta_t = \left(\text{linspace}(\sqrt{0.00085}, \sqrt{0.012}, 1000)_t \right)^2, \quad t = 0, 1, \dots, 999. \quad (29)$$

Specifically

$$\beta_0 = 0.00085, \quad (30)$$

$$\beta_1 \approx 0.0008546986, \quad (31)$$

$$\beta_2 \approx 0.0008594103, \quad (32)$$

$$\beta_{10} \approx 0.0008975693, \quad (33)$$

$$\beta_{100} \approx 0.0013839726, \quad (34)$$

$$\beta_{999} = 0.012, \quad (35)$$

and $\bar{\alpha}_t$ in (18) is uniquely determined from this sequence.

Very roughly speaking, in the update step, the coefficient of \mathbf{x}_{t_k} is relatively close to 1, and the

values of $-\widehat{\epsilon}_{t_k}$ and z_{t_k} are small compared to 1. That is, roughly speaking, the DDPM update has a structure of obtaining the step direction with a neural network, scaling it, subtracting it, and adding small Gaussian noise close to standard deviation $\sqrt{\beta_t}$ to it.

^a<https://github.com/runwayml/stable-diffusion/blob/main/configs/stable-diffusion/v1-inference.yaml>

Remark 6.2. (6.1) corresponds one-to-one with the default settings (epsilon prediction, fixed small variance) of `DDPMScheduler.step`, and applying CFG twice does not require changes to the scheduler’s API.

6.2 Euler a (Euler Ancestral)

Definition 6.2 (Reverse Diffusion Algorithm by Euler a). The following data are given.

- **Quantities determined during training:**

- Number of training timesteps $T \in \mathbb{Z}_{>0}$.
- Cumulative product sequence on training timesteps

$$\bar{\alpha}_t \in (0, 1], \quad t = 0, 1, \dots, T - 1. \quad (36)$$

- U-Net (noise predictor) trained with epsilon prediction

$$f_{\theta} : \mathbb{R}^{d_{\text{latent}}} \times \mathbb{R}^{d_{\text{context}}} \times [0, T - 1] \rightarrow \mathbb{R}^{d_{\text{latent}}}. \quad (37)$$

- **Quantities determined by other networks:**

- Positive embedding $c_+ \in \mathbb{R}^{d_{\text{context}}}$.
- Negative embedding $c_- \in \mathbb{R}^{d_{\text{context}}}$.

- **Quantities specified by the user at inference time:**

- Number of inference timesteps $N \in \mathbb{Z}_{>0}$.
- Guidance scale $s \in \mathbb{R}_{\geq 1}$.

At this time, Euler a operates by the following procedure.

(0) Timestep sequence generation. For $k = 0, 1, \dots, N - 1$, define

$$\tau_k := \frac{T - 1}{N - 1}(N - 1 - k). \quad (38)$$

(1) Sigma sequence generation. Corresponding to the training timesteps, set

$$\sigma_t^{\text{train}} := \sqrt{\frac{1 - \bar{\alpha}_t}{\bar{\alpha}_t}}, \quad t = 0, 1, \dots, T - 1 \quad (39)$$

and by linear interpolation

$$\text{LinInterp}(\{\sigma_t^{\text{train}}\}_{t=0}^{T-1}; \tau) := \sigma_{\lfloor \tau \rfloor}^{\text{train}} + (\tau - \lfloor \tau \rfloor)(\sigma_{\lceil \tau \rceil}^{\text{train}} - \sigma_{\lfloor \tau \rfloor}^{\text{train}}) \quad (40)$$

let

$$\sigma_k := \text{LinInterp}(\{\sigma_t^{\text{train}}\}; \tau_k), \quad k = 0, 1, \dots, N-1. \quad (41)$$

Furthermore, add

$$\sigma_N := 0. \quad (42)$$

(2) Initialization. Let

$$\mathbf{x}^{(0)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_{\text{latent}}}). \quad (43)$$

(3) Main loop. For $k = 0, 1, \dots, N-1$, perform the following.

- Let the current sample be $\mathbf{x}^{(k)}$, current time τ_k , current sigma σ_k , and next sigma σ_{k+1} .
- Before inputting to the U-Net, normalize as

$$\tilde{\mathbf{x}}^{(k)} := \frac{\mathbf{x}^{(k)}}{\sqrt{1 + \sigma_k^2}}. \quad (44)$$

- Composite the noise prediction by CFG:

$$\hat{\boldsymbol{\epsilon}}_{\text{cond}}^{(k)} := f_{\theta}(\tilde{\mathbf{x}}^{(k)}, \mathbf{c}_+, \tau_k), \quad (45)$$

$$\hat{\boldsymbol{\epsilon}}_{\text{uncond}}^{(k)} := f_{\theta}(\tilde{\mathbf{x}}^{(k)}, \mathbf{c}_-, \tau_k), \quad (46)$$

$$\hat{\boldsymbol{\epsilon}}^{(k)} := \hat{\boldsymbol{\epsilon}}_{\text{uncond}}^{(k)} + s(\hat{\boldsymbol{\epsilon}}_{\text{cond}}^{(k)} - \hat{\boldsymbol{\epsilon}}_{\text{uncond}}^{(k)}). \quad (47)$$

- Let the Euler coefficients be

$$\sigma_k^{\uparrow} := \sqrt{\frac{\sigma_{k+1}^2(\sigma_k^2 - \sigma_{k+1}^2)}{\sigma_k^2}}, \quad (48)$$

$$\sigma_k^{\downarrow} := \sqrt{\sigma_{k+1}^2 - (\sigma_k^{\uparrow})^2}. \quad (49)$$

Note that $(\sigma_k^{\uparrow})^2 + (\sigma_k^{\downarrow})^2 = \sigma_{k+1}^2$.

- Let the Gaussian noise be

$$\mathbf{z}^{(k)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_{\text{latent}}}) \quad (50)$$

and define the next sample as

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - (\sigma_k - \sigma_k^{\downarrow})\hat{\boldsymbol{\epsilon}}^{(k)} + \sigma_k^{\uparrow}\mathbf{z}^{(k)}. \quad (51)$$

(4) Output. The final sample $\mathbf{x}^{(N)}$ corresponds to $\sigma_N = 0$ and is output.

Remark 6.3. Although Euler a is a different algorithm from DDPM, it shares with DDPM the structure of obtaining the step direction with a neural network, scaling it, subtracting it, and adding small Gaussian noise close to standard deviation $\sqrt{\beta_t}$ to it.

6.3 DPM 2+ (Karras Sigma Sequence)

Definition 6.3 (Reverse Diffusion Algorithm by DPM 2+ (Karras)). The following data are given.

- **Quantities determined during training:**

- Number of training timesteps $T \in \mathbb{Z}_{>0}$.
- Cumulative product sequence on training timesteps

$$\bar{\alpha}_t \in (0, 1], \quad t = 0, 1, \dots, T - 1. \quad (52)$$

- Noise predictor trained with epsilon prediction

$$f_\theta : \mathbb{R}^{d_{\text{latent}}} \times \mathbb{R}^{d_{\text{context}}} \times [0, T - 1] \rightarrow \mathbb{R}^{d_{\text{latent}}}. \quad (53)$$

- **Quantities determined by other networks:**

- Positive embedding $\mathbf{c}_+ \in \mathbb{R}^{d_{\text{context}}}$.
- Negative embedding $\mathbf{c}_- \in \mathbb{R}^{d_{\text{context}}}$.

- **Quantities specified by the user at inference time:**

- Number of inference timesteps $N \in \mathbb{Z}_{>0}$.
- Guidance scale $s \in \mathbb{R}_{\geq 1}$.

At this time, the algorithm is as follows.

(0) Base timestep sequence generation. For $k = 0, 1, \dots, N - 1$, define

$$\tau_k := \frac{T - 1}{N - 1}(N - 1 - k). \quad (54)$$

(1) Base sigma sequence generation. On the training timesteps, define

$$\sigma_t^{\text{train}} := \sqrt{\frac{1 - \bar{\alpha}_t}{\bar{\alpha}_t}}, \quad t = 0, 1, \dots, T - 1 \quad (55)$$

and by linear interpolation

$$\text{LinInterp}(\{\sigma_t^{\text{train}}\}; \tau) := \sigma_{\lfloor \tau \rfloor}^{\text{train}} + (\tau - \lfloor \tau \rfloor)(\sigma_{\lceil \tau \rceil}^{\text{train}} - \sigma_{\lfloor \tau \rfloor}^{\text{train}}) \quad (56)$$

let

$$\sigma_k := \text{LinInterp}(\{\sigma_t^{\text{train}}\}; \tau_k), \quad k = 0, 1, \dots, N-1 \quad (57)$$

and furthermore, add

$$\sigma_N := 0. \quad (58)$$

(2) Log-sigma sequence and midpoint sigma sequence. For $k = 0, 1, \dots, N$, define

$$\lambda_k := \log \sigma_k \quad (59)$$

and give the midpoint sigma sequence by

$$\sigma_0^{\text{mid}} := \sigma_0, \quad (60)$$

$$\sigma_k^{\text{mid}} := \exp\left(\frac{\log \sigma_k + \log \sigma_{k-1}}{2}\right), \quad k = 1, 2, \dots, N. \quad (61)$$

(3) Construction of the long sigma sequence. Define a sequence of length $2N + 2$, $(\tilde{\sigma}_m)_{m=0}^{2N+1}$, by

$$\tilde{\sigma}_0 := \sigma_0, \quad \tilde{\sigma}_{2k-1} := \sigma_k, \quad \tilde{\sigma}_{2k} := \sigma_k, \quad \tilde{\sigma}_{2N+1} := \sigma_N, \quad k = 1, 2, \dots, N. \quad (62)$$

Similarly, define a long midpoint sigma sequence of length $2N + 2$, $(\tilde{\sigma}_m^{\text{mid}})_{m=0}^{2N+1}$, by

$$\tilde{\sigma}_0^{\text{mid}} := \sigma_0^{\text{mid}}, \quad \tilde{\sigma}_{2k-1}^{\text{mid}} := \sigma_k^{\text{mid}}, \quad \tilde{\sigma}_{2k}^{\text{mid}} := \sigma_k^{\text{mid}}, \quad \tilde{\sigma}_{2N+1}^{\text{mid}} := \sigma_N^{\text{mid}}, \quad k = 1, 2, \dots, N. \quad (63)$$

(4) Conversion to pseudo-timesteps ($\sigma \mapsto t$). Let the log-sigmas on the training timesteps be

$$L_t := \log \sigma_t^{\text{train}}, \quad t = 0, 1, \dots, T-1 \quad (64)$$

and for $\sigma > 0$, define

$$j(\sigma) := \max\{i \in \{0, 1, \dots, T-2\} \mid L_i \leq \log \sigma\}, \quad (65)$$

$$j^+(\sigma) := j(\sigma) + 1, \quad (66)$$

$$w(\sigma) := \frac{L_{j(\sigma)} - \log \sigma}{L_{j(\sigma)} - L_{j^+(\sigma)}} \in [0, 1], \quad (67)$$

$$\text{SigmaToT}(\sigma) := (1 - w(\sigma)) j(\sigma) + w(\sigma) j^+(\sigma). \quad (68)$$

(5) Construction of the long timestep sequence. For the base timestep sequence $(\tau_k)_{k=0}^{N-1}$, define a sequence of length $2N - 1$, $(\theta_m)_{m=0}^{2N-2}$, by

$$\theta_0 := \tau_0, \quad \theta_{2k-1} := \text{SigmaToT}(\sigma_k^{\text{mid}}), \quad \theta_{2k} := \tau_k, \quad k = 1, 2, \dots, N-1. \quad (69)$$

(6) Initialization. Let

$$\mathbf{x}^{(0)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_{\text{latent}}}). \quad (70)$$

(7) Main loop (Heun-type two-stage method). For $m = 0, 1, \dots, 2N - 2$, perform the following.

- Let the current pseudo-timestep be θ_m , current sigma $\tilde{\sigma}_m$, midpoint sigma $\tilde{\sigma}_{m+1}^{\text{mid}}$, and next sigma $\tilde{\sigma}_{m+1}$.
- Let the current sample be $\mathbf{x}^{(m)}$ and compute the following:

$$\hat{\boldsymbol{\epsilon}}_{\text{cond}}^{(m)} := f_{\theta}(\mathbf{x}^{(m)}, \mathbf{c}_+, \theta_m), \quad (71)$$

$$\hat{\boldsymbol{\epsilon}}_{\text{uncond}}^{(m)} := f_{\theta}(\mathbf{x}^{(m)}, \mathbf{c}_-, \theta_m), \quad (72)$$

$$\hat{\boldsymbol{\epsilon}}^{(m)} := \hat{\boldsymbol{\epsilon}}_{\text{uncond}}^{(m)} + s(\hat{\boldsymbol{\epsilon}}_{\text{cond}}^{(m)} - \hat{\boldsymbol{\epsilon}}_{\text{uncond}}^{(m)}). \quad (73)$$

- When m is even (1st stage), set the input sigma $\sigma_{\text{in}}^{(m)} := \tilde{\sigma}_m$, and find the original image prediction by

$$\hat{\mathbf{x}}_0^{(m)} := \mathbf{x}^{(m)} - \sigma_{\text{in}}^{(m)} \hat{\boldsymbol{\epsilon}}^{(m)}. \quad (74)$$

Let the derivative be

$$\mathbf{g}^{(m)} := \frac{\mathbf{x}^{(m)} - \hat{\mathbf{x}}_0^{(m)}}{\sigma_{\text{in}}^{(m)}} \quad (75)$$

and the step size be

$$\Delta\sigma^{(m)} := \tilde{\sigma}_{m+1}^{\text{mid}} - \tilde{\sigma}_m \quad (76)$$

and update as

$$\mathbf{x}^{(m+1)} := \mathbf{x}^{(m)} + \Delta\sigma^{(m)} \mathbf{g}^{(m)} \quad (77)$$

and temporarily store $\mathbf{x}^{(m)}$ and $\mathbf{g}^{(m)}$.

- When m is odd (2nd stage), set the input sigma $\sigma_{\text{in}}^{(m)} := \tilde{\sigma}_m^{\text{mid}}$, and from

$$\hat{\mathbf{x}}_0^{(m)} := \mathbf{x}^{(m)} - \sigma_{\text{in}}^{(m)} \hat{\boldsymbol{\epsilon}}^{(m)} \quad (78)$$

obtain

$$\mathbf{g}^{(m)} := \frac{\mathbf{x}^{(m)} - \hat{\mathbf{x}}_0^{(m)}}{\sigma_{\text{in}}^{(m)}}. \quad (79)$$

Using $\mathbf{x}^{(m-1)}$ and $\mathbf{g}^{(m-1)}$ stored in the 1st stage, update as

$$\mathbf{x}^{(m+1)} := \mathbf{x}^{(m-1)} + \frac{\tilde{\sigma}_{m+1} - \tilde{\sigma}_{m-1}}{2} (\mathbf{g}^{(m-1)} + \mathbf{g}^{(m)}). \quad (80)$$

(8) Output. Output $\mathbf{x}^{(2N-1)}$ at $m = 2N - 2$ as the final latent.

Remark 6.4. Definition (6.3) is very complicated, but there are two major differences from DDPM and Euler a.

- It uses the neural network twice per step. Therefore, even if the same number of inference steps N is specified, the inference time tends to be longer.
- It does not use pseudo-random numbers in steps other than initialization.

7 Advantages of Each Scheduler (Practical Perspective)

Remark 7.1 (DDPM). Theoretically clear and consistent with the training formulation, but slow if the number of inference steps is large. In practice, faster schedulers like DDIM, PNDM, Euler, DPM-based are often selected^a.

^aDiffusers guide: https://huggingface.co/docs/diffusers/en/using-diffusers/write_own_pipeline

Remark 7.2 (Euler / Euler a). High quality in few steps (20–30 steps target), easy to implement/tune, and highly robust^a.

^aDiffusers Euler-based API: <https://huggingface.co/docs/diffusers/en/api/schedulers/euler>, https://huggingface.co/docs/diffusers/en/api/schedulers/euler_ancestral

Remark 7.3 (DPM 2+ (Karras)). A second-order multi-step method, easy to obtain a high quality/speed trade-off. Implementations combining with the Karras sigma sequence for high quality in few steps are widely used [7].

8 Summary and Next Lecture Preview

8.1 Summary Corresponding to Learning Outcomes

- Input and Output: The reverse diffusion process takes the current sample, text embedding (pos/neg), and time (DDPM's integer time, Euler a's float time, or KDPM2's pseudo-time) as input, and outputs the next sample.
- Network Conditions: The output dimension matches the input sample dimension ((14)), and it can accept conditional input and time.
- Text Handling/Control: The quality/diversity trade-off can be adjusted by CFG (17) [3].
- Scheduler Selection: DDPM is fundamental, Euler/Euler a are for fast general-purpose use, and DPM 2+ (Karras) has strengths in high quality in few steps.

8.2 Next Lecture Preview

Next time, we will move away from inference for a moment and outline the diffusion process used for training the reverse diffusion process.

References

- [1] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. In International Conference on Learning Representations (ICLR), 2017.
- [2] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems, 33:6840–6851, 2020.
- [3] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598, 2022.
- [4] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Elucidating the design space of diffusion-based generative models. Advances in Neural Information Processing Systems, 35:26565–26577, 2022.
- [5] Diederik P. Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In Advances in Neural Information Processing Systems (NeurIPS), volume 31, 2018.
- [6] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114, 2013.
- [7] Chunqi Lu, Junxian He, Chenyang Xu, and Yang Song. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. arXiv preprint arXiv:2211.01095, 2022.
- [8] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 10684–10695, 2022.