

# AI Applications Lecture 15

## Image Generation AI 5: Convolutional Neural Networks for Image Generation

SUZUKI, Atsushi  
Jing WANG

### Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Roadmap Recap . . . . .	3
1.2	Learning Outcomes . . . . .	3
<b>2</b>	<b>Preparation: Mathematical Notations</b>	<b>3</b>
<b>3</b>	<b>General Theory: Reconfirming Architectural Freedom</b>	<b>5</b>
3.1	Training and Inference of Noise Estimators are Architecture-Independent . .	5
3.2	VAE Training and Inference are Likewise Architecture-Independent . . . . .	5
3.3	Consequence of the General Theory . . . . .	6
<b>4</b>	<b>Necessity of Variable I/O Sizes and Convolutional Layers</b>	<b>6</b>
4.1	Necessity of Variable Input/Output Sizes . . . . .	6
4.2	Layer Achieving Variable I/O Sizes: Convolutional Layer . . . . .	6
<b>5</b>	<b>Overview of Stable Diffusion 1.5 U-Net and VAE with Diagrams</b>	<b>7</b>
5.1	Influence of Different Motivations for U-Net and VAE on Architectural Differences	7
5.2	Understanding the Implemented Architecture . . . . .	7
5.3	Block Diagram of U-Net (Conditional) . . . . .	7
5.4	Block Diagram of VAE (Decoder) . . . . .	8
5.5	Confirmation of Variable Input/Output Sizes . . . . .	8

<b>6</b>	<b>Formal Definitions of Layers (Using only cross-correlation, matrix operations, and elementary functions)</b>	<b>8</b>
6.1	Cross-Correlation (2D "Convolution") . . . . .	8
6.2	Linear Layer (Fully-Connected) . . . . .	9
6.3	Activation Functions and Softmax . . . . .	9
6.4	Normalization . . . . .	10
6.5	Up/Downsample . . . . .	10
6.5.1	Definition of Downsample by Strided Cross-Correlation . . . . .	10
6.5.2	Explicit Definition based on Element-wise Formulas (Nearest/Average Pooling) . . . . .	10
6.5.3	Upsample by Nearest-Neighbor Interpolation and Cross-Correlation . . . . .	11
6.6	Reshape, Concatenate . . . . .	11
6.7	Timestep Embedding (Noise Level) . . . . .	11
6.8	Scaled Dot-Product Attention (Self/Cross Attention) . . . . .	12
6.8.1	Programming Intuition of a Dictionary (Key-Value Map) . . . . .	12
6.9	Time-Conditioned Residual Block . . . . .	15
6.10	U-Net Construction Blocks (Down/Up/Mid) . . . . .	16
6.11	VAE Decoder's Terminal Mapping (RGB Output; $1 \times 1$ conv) . . . . .	18
<b>7</b>	<b>Why U-Net Reaches the Entire Area "Shallowly": Quantitative Comparison of Receptive Fields</b>	<b>18</b>
7.1	Receptive Field of a Pure CNN (Conv2d only) . . . . .	18
7.2	Receptive Field of U-Net (with staged Downsample2D) . . . . .	18
<b>8</b>	<b>Full Definition of U-Net and VAE Decoder "as Functions"</b>	<b>19</b>
8.1	The U-Net (Conditional) Overall Function . . . . .	19
8.2	The VAE Decoder Overall Function . . . . .	20
<b>9</b>	<b>Summary (Correspondence to Learning Outcomes)</b>	<b>20</b>
<b>10</b>	<b>Next Lecture Preview</b>	<b>21</b>

# 1 Introduction

## 1.1 Roadmap Recap

We will review the content learned so far. Three lectures ago, we learned about the **Variational Autoencoder (VAE)** as a natural image decoder [3]. Two lectures ago, we learned about the **reverse diffusion process** that generates low-resolution latent images, namely the **denoising scheduler** [1]. In the previous lecture, we mathematically understood the sense in which the reverse diffusion process performs **distribution learning**, using the continuity equation, score, and KL divergence (see [2, 6]).

In this lecture, we will focus our discussion on **practical image generation AI** and look in detail at the **neural network architectures** used in denoising schedulers and VAEs. In particular, we will focus on the differences, as the **architecture in the original paper** and the **architecture used in actual implementations** often differ (e.g., Latent Diffusion/Stable Diffusion [4]).

## 1.2 Learning Outcomes

By the end of this lecture, students should be able to:

- **Mathematically describe** the **neural network architectures** used in **practical image generation AI**.
- **Explain** how practical image generation AI achieves support for **variable input/output sizes** by using specific **layers**.
- **Explain** how the neural network architectures used in practical image generation AI have **changed from their original proposals**.

# 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. E.g.,  $\mathcal{A}$ .
  - $x \in \mathcal{A}$ : Indicates that element  $x$  belongs to set  $\mathcal{A}$ .
  - $\{\}$ : The empty set.
  - $\{a, b, c\}$ : The set consisting of elements  $a, b, c$  (set-builder notation by extension).

- $\{x \in \mathcal{A} \mid P(x)\}$ : The set of elements in set  $\mathcal{A}$  for which the proposition  $P(x)$  is true (set-builder notation by intension).
- $|\mathcal{A}|$ : The number of elements in set  $\mathcal{A}$  (in this lecture, used only for finite sets).
- $\mathbb{R}$ : The set of all real numbers. Similarly for  $\mathbb{R}_{>0}$ ,  $\mathbb{R}_{\geq 0}$ , etc.
- $\mathbb{Z}$ : The set of all integers. Similarly for  $\mathbb{Z}_{>0}$ ,  $\mathbb{Z}_{\geq 0}$ , etc.
- $[1, k]_{\mathbb{Z}}$ : 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 a mapping.
- $y = f(x)$  denotes the output  $y \in \mathcal{Y}$  for the input  $x \in \mathcal{X}$ .

• **Vector:**

- Vectors are denoted by bold italic lowercase letters. E.g.,  $\mathbf{v}$ .  $\mathbf{v} \in \mathbb{R}^n$ .
- 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:**

- We call  $\mathbf{a} : [1, n]_{\mathbb{Z}} \rightarrow \mathcal{A}$  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. E.g.,  $\mathbf{A} \in \mathbb{R}^{m,n}$ .
- 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 General Theory: Reconfirming Architectural Freedom

#### 3.1 Training and Inference of Noise Estimators are Architecture-Independent

The training of the **noise estimator** used in the denoising scheduler was given by the optimization problem of minimizing the squared error of the noise estimation. The objective function is identical to the previous lecture:

$$\min_{\theta} \sum_{i=1}^m \left\| \epsilon^{(i)} - \hat{\epsilon}_{\theta}(\zeta^{(i)}, \mathbf{c}^{(i)}, t^{(i)}) \right\|_2^2. \quad (6)$$

(6) is defined **independently of the neural network's architecture**. Inference also just involves inserting the trained  $\hat{\epsilon}_{\theta}$  into a sequential algorithm, which is also **architecture-independent** (e.g., DDPM/DDIM steps [1]).

#### 3.2 VAE Training and Inference are Likewise Architecture-Independent

VAE training is regularized reconstruction error minimization [3]. For each input image  $\underline{\mathbf{X}} \in \mathcal{I}$ , we define an encoder that outputs a mean vector

$$\text{MeanEnc}_{\eta_{\text{mean}}} : \mathcal{I} \rightarrow \mathbb{R}^d, \quad \mu(\underline{\mathbf{X}}) := \text{MeanEnc}_{\eta_{\text{mean}}}(\underline{\mathbf{X}}), \quad (7)$$

and an encoder that outputs the standard deviation for each component

$$\text{SDEnc}_{\eta_{\text{SD}}} : \mathcal{I} \rightarrow \mathbb{R}_{>0}^d, \quad \sigma(\underline{\mathbf{X}}) := \text{SDEnc}_{\eta_{\text{SD}}}(\underline{\mathbf{X}}) \quad (8)$$

( $\sigma$  is a positive real number for each element). The entire set of parameters for the VAE encoder is then

$$\boldsymbol{\eta} := (\eta_{\text{mean}}, \eta_{\text{SD}}) \quad (9)$$

. We **sample** a standard normal random vector  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and generate the latent vector by

$$\mathbf{z}_\epsilon := \boldsymbol{\mu}(\underline{\mathbf{X}}) + \boldsymbol{\sigma}(\underline{\mathbf{X}}) \odot \epsilon \quad (10)$$

( $\odot$  is the element-wise product). The decoder provides

$$\hat{\mathbf{X}}_\epsilon := \text{Dec}_\gamma(\mathbf{z}_\epsilon) \quad (11)$$

. Using a reconstruction loss function  $\ell : \mathcal{I} \times \mathcal{I} \rightarrow \mathbb{R}_{\geq 0}$ , the objective function is

$$\mathcal{L}(\boldsymbol{\eta}, \boldsymbol{\gamma}) := \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \underbrace{\ell(\underline{\mathbf{X}}, \hat{\mathbf{X}}_\epsilon)}_{\text{Reconstruction term}} \right] + \beta \underbrace{\sum_{i=1}^d \left( \mu_i(\underline{\mathbf{X}})^2 + \sigma_i(\underline{\mathbf{X}})^2 - \log \sigma_i(\underline{\mathbf{X}})^2 - 1 \right)}_{\text{Regularization (concentration to origin)}}. \quad (12)$$

Each term in (12) is also **architecture-independent**. Once the decoder is trained, inference is **only the application of the decoder**.

### 3.3 Consequence of the General Theory

From the above, it is clear that, outside the context of image generation, both the noise estimator and the VAE decoder can adopt **any architecture**.

## 4 Necessity of Variable I/O Sizes and Convolutional Layers

### 4.1 Necessity of Variable Input/Output Sizes

In practical image generation, the **output resolution (dimensions)** depends on user requirements and cannot be fixed at training time. Therefore, an architecture with **variable input/output sizes**, which allows selecting the output size by choosing the input size (latent resolution), is necessary.

### 4.2 Layer Achieving Variable I/O Sizes: Convolutional Layer

To construct variable input/output sizes, each layer only needs to be a **parametric function compatible with variable sizes**. A typical example is the **convolution layer (implemented as cross-correlation)**. The noise estimator in Stable Diffusion 1.5 is a **U-Net** [5] family (conditional, with attention), and the VAE is also constructed with convolutional systems [4].

## 5 Overview of Stable Diffusion 1.5 U-Net and VAE with Diagrams

### 5.1 Influence of Different Motivations for U-Net and VAE on Architectural Differences

U-Net and VAE have very different expected functionalities.

- U-Net's purpose is to generate a **globally coherent** low-resolution latent image from noise (which has no information), using information from a text encoder.
- VAE's purpose is to convert a low-resolution latent image, which already has some global coherence, into a high-resolution natural image by defining the details.

This is evident even when observing the intermediate states of image generation.

Corresponding to these differences in motivation, the architectures actually used for U-Net and VAE encoders also differ.

- U-Net, to achieve global coherence, uses downsampling, giving it a structure that efficiently allows pixels at one edge to influence pixels at the opposite edge with a relatively small number of layers.
- The VAE encoder is composed of pure convolutional layers and upsampling layers, adopting a structure that restricts the use of parameters to local value transformations.

Let's confirm these differences by looking at the actual architectures.

### 5.2 Understanding the Implemented Architecture

The actual architecture may differ in details from the paper's description. Using tools like **torchview**, one can visualize the **computation graph** from the implemented model, making it easier to grasp implementation differences<sup>1</sup>.

### 5.3 Block Diagram of U-Net (Conditional)

Figure 1 shows a schematic of the conditional U-Net in Stable Diffusion 1.5 (including down-sample, bottleneck, upsample, skip connections, and cross-attention). This figure is obtained by applying torchview with `depth=1` to the U-Net part of Stable Diffusion 1.5 uploaded to Hugging Face, reflecting the actual implementation.

**Remark 5.1. Difference from the original U-Net:** Ronneberger et al.'s U-Net [5] is an **image-to-image** mapping for **medical image segmentation**, where both input and output

---

<sup>1</sup>torchview: <https://github.com/mert-kurttutan/torchview>

are images. In image generation AI, it needs to accept **time (noise level)** and **text conditions** as input, and perform **noise estimation (or  $v$ -prediction)**. Therefore, the significant differences are the addition of a **time encoder** and **cross-attention layers (for text)** [4].

## 5.4 Block Diagram of VAE (Decoder)

Figure 2 shows the architecture of the VAE decoder in Stable Diffusion 1.5 (Conv/ResBlock/Upsample). This figure is obtained by applying torchview with `depth=1` to the VAE part of Stable Diffusion 1.5 uploaded to Hugging Face, reflecting the actual implementation.

## 5.5 Confirmation of Variable Input/Output Sizes

The fact that U-Net and VAE have variable input/output sizes follows from each component layer (convolution, normalization, attention, up/down-sample) being defined **convolutionally (translationally equivariant under isomorphism)** with respect to the **spatial resolution**  $H \times W$ .

**Remark 5.2.** The **variable input/output sizes** mentioned here refer to the **variability in the image width  $W$  and height  $H$** ; the **number of channels  $C$  is fixed** (although the channel width may change in steps inside the U-Net, the number of channels at the input/output interface is specified).

# 6 Formal Definitions of Layers (Using only cross-correlation, matrix operations, and elementary functions)

In the following, all layers (functions) used in the U-Net (conditional) and VAE decoder will be **formally** defined using only **cross-correlation (convolution in practice)**, **matrix operations**, **elementary functions**, and **well-known special functions**. In each definition, the **function name (matching the library class/function name)** is given by  $\cdot$ , **non-learnable hyperparameters** (e.g., stride, padding, groups) are indicated in superscript parentheses  $(\cdot)$ , and the **entire set of learnable parameters** is specified in the subscript. The URL of the corresponding implementation is provided in a footnote.

## 6.1 Cross-Correlation (2D "Convolution")

**Definition 6.1** (Conv2d (Convolution in practice is cross-correlation)). For input  $\underline{X} \in \mathbb{R}^{C_{\text{in}} \times H \times W}$  and output  $\underline{Y} \in \mathbb{R}^{C_{\text{out}} \times H' \times W'}$ , we fix **hyperparameters** kernel size  $(k_h, k_w)$ , stride  $(s_h, s_w)$ , and padding  $(p_h, p_w)$ . The set of **learnable parameters** (weights and biases) is

$$\Theta_{\text{Conv2d}} = \left\{ \mathbf{W}^{(o)} \in \mathbb{R}^{C_{\text{in}} \times k_h \times k_w}, b_o \in \mathbb{R} \right\}_{o=1}^{C_{\text{out}}} \quad (13)$$



At this time,

$$(\text{Conv2d}_{\Theta_{\text{Conv2d}}}^{(k_h, k_w; s_h, s_w; p_h, p_w)}(\underline{X}))_{o,i,j} = b_o + \sum_{c=1}^{C_{\text{in}}} \sum_{u=1}^{k_h} \sum_{v=1}^{k_w} W_{c,u,v}^{(o)} X_{c, i \cdot s_h + u - p_h, j \cdot s_w + v - p_w}. \quad (14)$$

The output spatial size is  $H' = \left\lfloor \frac{H - k_h + 2p_h}{s_h} \right\rfloor + 1$ ,  $W' = \left\lfloor \frac{W - k_w + 2p_w}{s_w} \right\rfloor + 1$ .<sup>a</sup>

<sup>a</sup>`torch.nn.Conv2d`: <https://pytorch.org/docs/stable/generated/torch.nn.Conv2d.html>,  
`torch.nn.functional.conv2d`: <https://pytorch.org/docs/stable/generated/torch.nn.functional.conv2d.html>.

**Remark 6.1.** "Convolution" in implementations is cross-correlation (**does not flip the kernel**) and matches (14).

## 6.2 Linear Layer (Fully-Connected)

**Definition 6.2** (Linear). For input  $\mathbf{x} \in \mathbb{R}^{d_{\text{in}}}$ , output  $\mathbf{y} \in \mathbb{R}^{d_{\text{out}}}$ , and learnable parameters  $\Theta_{\text{Linear}} = \{\mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}, \mathbf{b} \in \mathbb{R}^{d_{\text{out}}}\}$ ,

$$\text{Linear}_{\Theta_{\text{Linear}}}(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}. \quad (15)$$

<sup>a</sup>

<sup>a</sup>`torch.nn.Linear`: <https://pytorch.org/docs/stable/generated/torch.nn.Linear.html>,  
`torch.nn.functional.linear`: <https://pytorch.org/docs/stable/generated/torch.nn.functional.linear.html>.

## 6.3 Activation Functions and Softmax

**Definition 6.3** (SiLU (Swish)). For a component  $u$  of an arbitrary-dimensional tensor,

$$\text{SiLU}(u) = u \sigma(u), \quad \sigma(u) = \frac{1}{1 + e^{-u}}. \quad (16)$$

<sup>a</sup>

<sup>a</sup>`torch.nn.SiLU`: <https://pytorch.org/docs/stable/generated/torch.nn.SiLU.html>,  
`torch.nn.functional.silu`: <https://pytorch.org/docs/stable/generated/torch.nn.functional.silu.html>.

**Definition 6.4** (Softmax). For  $\mathbf{x} \in \mathbb{R}^n$ , fixing the temperature  $\tau > 0$ ,

$$(\text{Softmax}^{(\tau)}(\mathbf{x}))_i = \frac{\exp(x_i/\tau)}{\sum_{j=1}^n \exp(x_j/\tau)}. \quad (17)$$

## 6.4 Normalization

**Definition 6.5** (GroupNorm). For input  $\underline{X} \in \mathbb{R}^{C \times H \times W}$ , **hyperparameter** number of groups  $G \mid C$ , and **learnable parameters**  $\Theta_{\text{GroupNorm}} = \{\gamma \in \mathbb{R}^C, \beta \in \mathbb{R}^C\}$ , the mean and variance for each group  $g$  are

$$\mu_g = \frac{1}{|S_g|} \sum_{(c,i,j) \in S_g} X_{c,i,j}, \quad \sigma_g^2 = \frac{1}{|S_g|} \sum_{(c,i,j) \in S_g} (X_{c,i,j} - \mu_g)^2, \quad (18)$$

The output is

$$(\text{GroupNorm}_{\Theta_{\text{GroupNorm}}}^{(G)}(\underline{X}))_{c,i,j} = \gamma_c \frac{X_{c,i,j} - \mu_{g(c)}}{\sqrt{\sigma_{g(c)}^2 + \varepsilon}} + \beta_c. \quad (19)$$

<sup>a</sup>

<sup>a</sup>torch.nn.GroupNorm: <https://pytorch.org/docs/stable/generated/torch.nn.GroupNorm.html>.

## 6.5 Up/Downsample

### 6.5.1 Definition of Downsample by Strided Cross-Correlation

Using Conv2d with stride 2:

**Definition 6.6** (Downsample2D (by Conv2d)).

$$\text{Downsample2D}_{\Theta_{\text{Down}}}^{(2)}(\underline{X}) := \text{Conv2d}_{\Theta_{\text{Down}}}^{(k_h, k_w; 2, 2; p_h, p_w)}(\underline{X}). \quad (20)$$

( $\Theta_{\text{Down}}$  is the set of learnable parameters for Conv2d).

### 6.5.2 Explicit Definition based on Element-wise Formulas (Nearest/Average Pooling)

**Definition 6.7** (Downsample2DNearest (nearest pooling; element-wise)). For  $\underline{X} \in \mathbb{R}^{C \times H \times W}$ ,  $\underline{Y} \in \mathbb{R}^{C \times \lfloor H/2 \rfloor \times \lfloor W/2 \rfloor}$  is defined as

$$Y_{c,i,j} = X_{c,2i,2j} \quad (0 \leq i < \lfloor H/2 \rfloor, 0 \leq j < \lfloor W/2 \rfloor) \quad (21)$$

(No learnable parameters)

**Definition 6.8** (Downsample2DAverage (average pooling; element-wise)). For  $\underline{X} \in \mathbb{R}^{C \times H \times W}$ ,

$$Y_{c,i,j} = \frac{1}{4} \sum_{u=0}^1 \sum_{v=0}^1 X_{c,2i+u,2j+v}. \quad (22)$$

(No learnable parameters)

### 6.5.3 Upsample by Nearest-Neighbor Interpolation and Cross-Correlation

**Definition 6.9** (Interpolate (nearest; element-wise)). For  $\underline{X} \in \mathbb{R}^{C \times H \times W}$ ,  $\underline{Z} = \text{Interpolate}^{(\times 2, \text{nearest})}(\underline{X}) \in \mathbb{R}^{C \times 2H \times 2W}$  is defined as

$$Z_{c, 2i+u, 2j+v} = X_{c, i, j} \quad (u, v \in \{0, 1\}) \quad (23)$$

<sup>a</sup>

<sup>a</sup>torch.nn.functional.interpolate: <https://pytorch.org/docs/stable/generated/torch.nn.functional.interpolate.html>.

**Definition 6.10** (Upsample2D (nearest+conv)).

$$\text{Upsample2D}_{\Theta_{\text{Up}}}^{(2)}(\underline{X}) := \text{Conv2d}_{\Theta_{\text{Up}}}^{(k_h, k_w; 1, 1; p_h, p_w)}(\text{Interpolate}^{(\times 2, \text{nearest})}(\underline{X})). \quad (24)$$

## 6.6 Reshape, Concatenate

**Definition 6.11** (flatten, unflatten, Concat). For  $\underline{X} \in \mathbb{R}^{C \times H \times W}$ ,

$$\text{flatten}_{(H, W)}(\underline{X}) \in \mathbb{R}^{(HW) \times C}, \quad (\text{flatten}_{(H, W)}(\underline{X}))_{(i-1)W+j, c} = X_{c, i, j}, \quad (25)$$

$$\text{unflatten}_{(H, W)}(\underline{Y}) \in \mathbb{R}^{C \times H \times W}, \quad (\text{unflatten}_{(H, W)}(\underline{Y}))_{c, i, j} = Y_{(i-1)W+j, c}, \quad (26)$$

$$\text{Concat}(\underline{A}, \underline{B}) = \underline{A} \oplus \underline{B} \text{ (concatenation along the channel dimension)}. \quad (27)$$

## 6.7 Timestep Embedding (Noise Level)

**Definition 6.12** (TimestepEmbedding (Sinusoidal + MLP)). For scalar  $t \in \mathbb{R}$  and frequency sequence  $\omega_r = \omega_0 \beta^{r-1}$  ( $r = 1, \dots, R$ ),

$$\mathbf{e}(t) = [\cos(\omega_1 t), \sin(\omega_1 t), \dots, \cos(\omega_R t), \sin(\omega_R t)]^\top \in \mathbb{R}^{2R}. \quad (28)$$

For learnable parameters  $\Theta_{\text{TE}} = \{U_1 \in \mathbb{R}^{d_h \times 2R}, \mathbf{b}_1 \in \mathbb{R}^{d_h}, U_2 \in \mathbb{R}^{d_t \times d_h}, \mathbf{b}_2 \in \mathbb{R}^{d_t}\}$ ,

$$\text{TimestepEmbedding}_{\Theta_{\text{TE}}}(t) = U_2 \text{SiLU}(U_1 \mathbf{e}(t) + \mathbf{b}_1) + \mathbf{b}_2 \in \mathbb{R}^{d_t}. \quad (29)$$

<sup>a</sup>

<sup>a</sup>Diffusers implementation (embeddings.py / unet\_2d\_condition.py): <https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/embeddings.py>, [https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/unet\\_2d\\_condition.py](https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/unet_2d_condition.py).

**Remark 6.2.** Components with small  $\omega$  represent the **coarse position (low frequency, long period)** of  $t$ , while components with large  $\omega$  represent the **fine position (high frequency, short period)**. This is analogous to the **positional numeral system**, where **upper digits**

are useful for approximate estimation, and **lower digits** independently provide information about divisors or remainders.

## 6.8 Scaled Dot-Product Attention (Self/Cross Attention)

**Definition 6.13** (ScaledDotProductAttention). For query  $Q \in \mathbb{R}^{N \times d}$ , key  $K \in \mathbb{R}^{M \times d}$ , and value  $V \in \mathbb{R}^{M \times d_v}$ ,

$$\text{ScaledDotProductAttention}(Q, K, V) = \text{Softmax}^{(\sqrt{d})^{-1}} \left( \frac{QK^T}{\sqrt{d}} \right) V. \quad (30)$$

<sup>a</sup>

<sup>a</sup>torch.nn.functional.scaled\_dot\_product\_attention: [https://pytorch.org/docs/stable/generated/torch.nn.functional.scaled\\_dot\\_product\\_attention.html](https://pytorch.org/docs/stable/generated/torch.nn.functional.scaled_dot_product_attention.html).

**Definition 6.14** (MultiheadAttention (SDPA with Projections)). For input sequence  $X \in \mathbb{R}^{N \times d_{\text{in}}}$  and context sequence  $C \in \mathbb{R}^{M \times d_{\text{ctx}}}$ , using learnable parameters

$$\Theta_{\text{MHA}} = \{W_Q \in \mathbb{R}^{d_{\text{in}} \times d}, W_K \in \mathbb{R}^{d_{\text{ctx}} \times d}, W_V \in \mathbb{R}^{d_{\text{ctx}} \times d_v}, W_O \in \mathbb{R}^{d_v \times d_{\text{out}}}\} \quad (31)$$

$$Q = XW_Q, \quad K = CW_K, \quad V = CW_V, \quad (32)$$

$$\text{MultiheadAttention}_{\Theta_{\text{MHA}}}(X, C) = \text{ScaledDotProductAttention}(Q, K, V) W_O. \quad (33)$$

<sup>a</sup>

<sup>a</sup>torch.nn.MultiheadAttention: <https://pytorch.org/docs/stable/generated/torch.nn.MultiheadAttention.html>. Diffusers attention mechanism (attn\_processor.py): [https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/attn\\_processor.py](https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/attn_processor.py).

### 6.8.1 Programming Intuition of a Dictionary (Key-Value Map)

A **dictionary (map)** in programming is a correspondence of {key : value}, a structure that retrieves the corresponding value when a key is given.

**Example 6.1** (Python Dictionary). For  $D = \{\text{"cat"} : 1, \text{"dog"} : 2\}$ ,  $D[\text{"dog"}] = 2$ . The ScaledDotProductAttention in attention implements a **soft dictionary** that "retrieves a weighted sum of values closest to the key" in a continuous vector space.

**Proposition 6.1** (Attention as a Soft Dictionary (Symbol correspondence for ScaledDotProductAttention)). Let  $K = [k_1^T; \dots; k_M^T] \in \mathbb{R}^{M \times d}$ ,  $V = [v_1^T; \dots; v_M^T] \in \mathbb{R}^{M \times d_v}$ , and

consider a single query  $\mathbf{q} \in \mathbb{R}^d$  with  $\mathbf{Q} = [\mathbf{q}^\top]$ . Then

$$\text{ScaledDotProductAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \left[ \sum_{m=1}^M \pi_m(\mathbf{q}) \mathbf{v}_m \right], \quad (34)$$

$$\pi_m(\mathbf{q}) = \frac{\exp(\langle \mathbf{q}, \mathbf{k}_m \rangle / \sqrt{d})}{\sum_{j=1}^M \exp(\langle \mathbf{q}, \mathbf{k}_j \rangle / \sqrt{d})}. \quad (35)$$

Here,  $\pi(\mathbf{q})$  is the first row of  $\text{Softmax}((\mathbf{Q}\mathbf{K}^\top)/\sqrt{d})$  and corresponds perfectly to the Softmax in (30).

*Proof.* The  $m$ -th component of  $\mathbf{Q}\mathbf{K}^\top/\sqrt{d} \in \mathbb{R}^{1 \times M}$  is  $\langle \mathbf{q}, \mathbf{k}_m \rangle / \sqrt{d}$ , so applying Softmax yields the weight  $\pi_m(\mathbf{q})$ . Multiplying by  $\mathbf{V}$  from the right yields  $\sum_m \pi_m(\mathbf{q}) \mathbf{v}_m$ . This can be shown by expanding equation (30) component-wise.  $\square$

**Remark 6.3.** If  $\pi(\mathbf{q})$  becomes a one-hot vector (1 for some  $m^*$ , 0 otherwise), then  $\sum_m \pi_m(\mathbf{q}) \mathbf{v}_m = \mathbf{v}_{m^*}$ , which matches the exact retrieval from a dictionary.

**Example 6.2** (Numerical Calculation of ScaledDotProductAttention). Let  $d = 2$ ,  $\mathbf{q} = (1, 0)^\top$ ,  $\mathbf{k}_1 = (1, 0)^\top$ ,  $\mathbf{k}_2 = (0, 1)^\top$ ,  $\mathbf{v}_1 = (2, 0)^\top$ ,  $\mathbf{v}_2 = (0, 3)^\top$ . At this time, the inner products are

$$\langle \mathbf{q}, \mathbf{k}_1 \rangle = 1 \cdot 1 + 0 \cdot 0 = 1, \quad \langle \mathbf{q}, \mathbf{k}_2 \rangle = 1 \cdot 0 + 0 \cdot 1 = 0 \quad (36)$$

As the exponentials of the scores divided by the square root of  $d$ ,

$$\sqrt{d} = \sqrt{2}, \quad (37)$$

$$s_1 = \exp(\langle \mathbf{q}, \mathbf{k}_1 \rangle / \sqrt{2}) = \exp(1/\sqrt{2}), \quad s_2 = \exp(\langle \mathbf{q}, \mathbf{k}_2 \rangle / \sqrt{2}) = \exp(0) = 1 \quad (38)$$

we obtain. Numerically, from  $1/\sqrt{2} \approx 0.70710678$ ,

$$s_1 \approx e^{0.70710678} \approx 2.02811498, \quad s_2 = 1. \quad (39)$$

The softmax weights are

$$\pi_1 = \frac{s_1}{s_1 + s_2} = \frac{e^{1/\sqrt{2}}}{e^{1/\sqrt{2}} + 1}, \quad \pi_2 = \frac{s_2}{s_1 + s_2} = \frac{1}{e^{1/\sqrt{2}} + 1} \quad (40)$$

and numerically,

$$\pi_1 \approx \frac{2.02811498}{2.02811498 + 1} \approx 0.66976155, \quad \pi_2 \approx \frac{1}{2.02811498 + 1} \approx 0.33023845. \quad (41)$$

The output is

$$\begin{aligned}\mathbf{o} &= \pi_1 \mathbf{v}_1 + \pi_2 \mathbf{v}_2 = \pi_1 \begin{bmatrix} 2 \\ 0 \end{bmatrix} + \pi_2 \begin{bmatrix} 0 \\ 3 \end{bmatrix} = \begin{bmatrix} 2\pi_1 \\ 3\pi_2 \end{bmatrix} \\ &= \left( \frac{2s_1}{s_1 + 1}, \frac{3}{s_1 + 1} \right)^\top \quad (\text{Substituting Eq. (38)}),\end{aligned}\tag{42}$$

and numerically,

$$\mathbf{o} \approx \begin{bmatrix} 2 \times 0.66976155 \\ 3 \times 0.33023845 \end{bmatrix} = \begin{bmatrix} 1.33952310 \\ 0.99071535 \end{bmatrix}.\tag{43}$$

As shown above, we can follow step-by-step in the order of exponential calculation  $\rightarrow$  normalization  $\rightarrow$  weighted sum of values.

**Exercise 6.1** (Numerical Example with 2D Vectors). Let  $d = 2$ ,  $\mathbf{q} = (2, 1)^\top$ ,  $\mathbf{k}_1 = (1, 0)^\top$ ,  $\mathbf{k}_2 = (0, 1)^\top$ ,  $\mathbf{v}_1 = (1, 2)^\top$ ,  $\mathbf{v}_2 = (4, -1)^\top$ . Calculate the output vector  $\mathbf{o}$  of ScaledDotProductAttention both as an exact expression and numerically ( $\sqrt{d} = \sqrt{2}$ ).

**Answer.** First, calculate the inner products:

$$\langle \mathbf{q}, \mathbf{k}_1 \rangle = 2 \cdot 1 + 1 \cdot 0 = 2, \quad \langle \mathbf{q}, \mathbf{k}_2 \rangle = 2 \cdot 0 + 1 \cdot 1 = 1.\tag{44}$$

Thus, the scores are

$$\begin{aligned}s_1 &= \exp(\langle \mathbf{q}, \mathbf{k}_1 \rangle / \sqrt{2}) = \exp(2 / \sqrt{2}) = \exp(\sqrt{2}), \\ s_2 &= \exp(\langle \mathbf{q}, \mathbf{k}_2 \rangle / \sqrt{2}) = \exp(1 / \sqrt{2})\end{aligned}\tag{45}$$

Numerically, from  $\sqrt{2} \approx 1.41421356$  and  $1/\sqrt{2} \approx 0.70710678$ ,

$$s_1 \approx e^{1.41421356} \approx 4.11325038, \quad s_2 \approx e^{0.70710678} \approx 2.02811498.\tag{46}$$

Therefore, the softmax weights are

$$\pi_1 = \frac{s_1}{s_1 + s_2} = \frac{e^{\sqrt{2}}}{e^{\sqrt{2}} + e^{1/\sqrt{2}}}, \quad \pi_2 = \frac{s_2}{s_1 + s_2} = \frac{e^{1/\sqrt{2}}}{e^{\sqrt{2}} + e^{1/\sqrt{2}}},\tag{47}$$

Numerically,

$$\pi_1 \approx \frac{4.11325038}{4.11325038 + 2.02811498} \approx 0.66976155, \quad \pi_2 \approx \frac{2.02811498}{4.11325038 + 2.02811498} \approx 0.33023845.\tag{48}$$

The output is

$$\begin{aligned} \mathbf{o} &= \pi_1 \mathbf{v}_1 + \pi_2 \mathbf{v}_2 = \pi_1 \begin{bmatrix} 1 \\ 2 \end{bmatrix} + \pi_2 \begin{bmatrix} 4 \\ -1 \end{bmatrix} = \begin{bmatrix} \pi_1 + 4\pi_2 \\ 2\pi_1 - \pi_2 \end{bmatrix} \\ &= \begin{bmatrix} 4 - 3\pi_1 \\ 3\pi_1 - 1 \end{bmatrix} \quad (\text{Simplified using } \pi_2 = 1 - \pi_1). \end{aligned} \quad (49)$$

Numerically,

$$\mathbf{o} \approx \begin{bmatrix} 4 - 3 \times 0.66976155 \\ 3 \times 0.66976155 - 1 \end{bmatrix} = \begin{bmatrix} 1.99071535 \\ 1.00928465 \end{bmatrix}. \quad (50)$$

This completes the derivation of the exact expression and the numerical calculation.

**Proposition 6.2** (Limit to a Hard Dictionary (ScaledDotProductAttention)). Suppose for some  $m^\star$ ,  $\mathbf{k}_{m^\star} \parallel \mathbf{q}$  and  $\mathbf{k}_m \perp \mathbf{q}$  ( $m \neq m^\star$ ). Then, for any  $\alpha > 0$ , let  $\mathbf{q}_\alpha = \alpha \mathbf{q}$ ,

$$\lim_{\alpha \rightarrow +\infty} \pi_m(\mathbf{q}_\alpha) = \begin{cases} 1, & m = m^\star, \\ 0, & m \neq m^\star. \end{cases} \quad (51)$$

*Proof.*  $\langle \mathbf{q}_\alpha, \mathbf{k}_{m^\star} \rangle = \alpha \|\mathbf{q}\| \|\mathbf{k}_{m^\star}\|$ , and others are 0. Therefore, only the  $m^\star$ -th component of  $\mathbf{QK}^\top / \sqrt{d}$  increases linearly with  $\alpha$ . From the definition of Softmax (17), the exponential of the largest component dominates the others, and (51) follows.  $\square$

**Remark 6.4.**  $\mathbf{k}_{m^\star} \parallel \mathbf{q}$  means "query and key are in the same direction (identical content)", which causes the corresponding  $\mathbf{v}_{m^\star}$  to be selected. That is, it matches the "retrieval of the value corresponding to the key".

## 6.9 Time-Conditioned Residual Block

**Definition 6.15** (ResnetBlock2D (Affine modulation by time embedding; FiLM)). For input  $\underline{\mathbf{X}} \in \mathbb{R}^{C_{\text{in}} \times H \times W}$  and time embedding  $\mathbf{h} \in \mathbb{R}^{d_t}$ , using learnable parameters

$$\Theta_{\text{ResnetBlock2D}} = \left( \Theta_{\text{Conv2d}}^{(1)}, \Theta_{\text{Conv2d}}^{(2)}, \Theta_{\text{Conv2d}}^{(s)}, \Theta_{\text{GN}}^{(1)}, \Theta_{\text{GN}}^{(2)}, \Theta_{\text{Linear}}^{(\gamma)}, \Theta_{\text{Linear}}^{(\beta)} \right),$$

$$\underline{U}_1 = \text{GroupNorm}_{\Theta_{\text{GN}}^{(1)}}^{(G)}(\underline{X}), \quad \underline{V}_1 = \text{SiLU}(\underline{U}_1), \quad \underline{W}_1 = \text{Conv2d}_{\Theta_{\text{Conv2d}}^{(1)}}^{(k,k; 1,1; p,p)}(\underline{V}_1) \quad (52)$$

$$\gamma(\mathbf{h}) = \text{Linear}_{\Theta_{\text{Linear}}^{(\gamma)}}(\mathbf{h}), \quad \beta(\mathbf{h}) = \text{Linear}_{\Theta_{\text{Linear}}^{(\beta)}}(\mathbf{h}), \quad (53)$$

$$\underline{U}_2 = \text{GroupNorm}_{\Theta_{\text{GN}}^{(2)}}^{(G)}(\underline{W}_1), \quad \widehat{\underline{U}}_2 = \gamma(\mathbf{h}) \odot \underline{U}_2 + \beta(\mathbf{h}), \quad (54)$$

$$\underline{V}_2 = \text{SiLU}(\widehat{\underline{U}}_2), \quad \underline{W}_2 = \text{Conv2d}_{\Theta_{\text{Conv2d}}^{(2)}}^{(k,k; 1,1; p,p)}(\underline{V}_2), \quad (55)$$

$$\underline{S} = \text{Conv2d}_{\Theta_{\text{Conv2d}}^{(s)}}^{(1,1; 1,1; 0,0)}(\underline{X}) \quad (\text{channel matching}), \quad (56)$$

$$\text{ResnetBlock2D}_{\Theta_{\text{ResnetBlock2D}}}(\underline{X}, \mathbf{h}) = \underline{S} + \underline{W}_2. \quad (57)$$

a

<sup>a</sup>Diffusers ResnetBlock2D implementation: <https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/resnet.py>.

## 6.10 U-Net Construction Blocks (Down/Up/Mid)

**Definition 6.16** (DownBlock2D). Taking the number of residual layers within the level  $n \in \mathbb{Z}_{>0}$  as a hyperparameter, and for learnable parameters  $\Theta_{\text{DownBlock2D}} = (\{\Theta_{\text{Res}}^{(r)}\}_{r=1}^n, \Theta_{\text{Down}})$ ,

$$\underline{H}_0 = \underline{X}, \quad \underline{H}_r = \text{ResnetBlock2D}_{\Theta_{\text{Res}}^{(r)}}(\underline{H}_{r-1}, \mathbf{h}) \quad (r = 1, \dots, n), \quad (58)$$

$$\text{DownBlock2D}_{\Theta_{\text{DownBlock2D}}}(\underline{X}, \mathbf{h}) = \text{Downsample2D}_{\Theta_{\text{Down}}}^{(2)}(\underline{H}_n). \quad (59)$$

a

<sup>a</sup>Diffusers block implementation: [https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/unet\\_2d\\_blocks.py](https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/unet_2d_blocks.py).

**Definition 6.17** (UpBlock2D). Combining the skip connection  $\underline{S}$  and the input from the bottom  $\underline{X}$  with Concat, and for learnable parameters  $\Theta_{\text{UpBlock2D}} = (\{\Theta_{\text{Res}}^{(r)}\}_{r=1}^n, \Theta_{\text{Up}})$ ,

$$\underline{Y}_0 = \text{Concat}\left(\text{Upsample2D}_{\Theta_{\text{Up}}}^{(2)}(\underline{X}), \underline{S}\right), \quad (60)$$

$$\underline{Y}_r = \text{ResnetBlock2D}_{\Theta_{\text{Res}}^{(r)}}(\underline{Y}_{r-1}, \mathbf{h}) \quad (r = 1, \dots, n), \quad (61)$$

$$\text{UpBlock2D}_{\Theta_{\text{UpBlock2D}}}(\underline{X}, \underline{S}, \mathbf{h}) = \underline{Y}_n. \quad (62)$$

a

<sup>a</sup>Implementation reference: [https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/unet\\_2d\\_blocks.py](https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/unet_2d_blocks.py).

**Definition 6.18** (CrossAttentionMidBlock2D (Using Self/Cross Attention)). For learnable pa-



parameters  $\Theta_{\text{MidBlock2D}} = (\Theta_{\text{Res}}^{(1)}, \Theta_{\text{MHA}}^{\text{self}}, \Theta_{\text{MHA}}^{\text{cross}}, \Theta_{\text{Res}}^{(2)})$  and text context  $\mathbf{C} \in \mathbb{R}^{M \times d_{\text{ctx}}}$ ,

$$\underline{\mathbf{A}}_0 = \text{ResnetBlock2D}_{\Theta_{\text{Res}}^{(1)}}(\underline{\mathbf{X}}, \mathbf{h}), \quad (63)$$

$$\mathbf{X}_{\text{flat}} = \text{flatten}_{(H,W)}(\underline{\mathbf{A}}_0) \in \mathbb{R}^{(HW) \times d_{\text{in}}}, \quad (64)$$

$$\mathbf{B}_1 = \text{MultiheadAttention}_{\Theta_{\text{MHA}}^{\text{self}}}(\mathbf{X}_{\text{flat}}, \mathbf{X}_{\text{flat}}), \quad (65)$$

$$\mathbf{B}_2 = \text{MultiheadAttention}_{\Theta_{\text{MHA}}^{\text{cross}}}(\mathbf{B}_1, \mathbf{C}), \quad (66)$$

$$\underline{\mathbf{A}}_1 = \text{unflatten}_{(H,W)}(\mathbf{B}_2), \quad \text{CrossAttentionMidBlock2D}_{\Theta_{\text{MidBlock2D}}}(\underline{\mathbf{X}}, \mathbf{h}, \mathbf{C}) = \text{ResnetBlock2D}_{\Theta_{\text{Res}}^{(2)}}(\underline{\mathbf{A}}_1, \mathbf{h}), \quad (67)$$

a

<sup>a</sup>Processor for self/cross attention: [https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/attn\\_processor.py](https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/attn_processor.py).

**Definition 6.19** (CrossAttentionDownBlock2D). Let  $n \in \mathbb{Z}_{>0}$  be the number of residual layers in the level. Self-attention and cross-attention use the mapping defined in §6.8 for the ‘sequence of tokens per spatial position’. Let the learnable parameters be

$$\Theta_{\text{XDown}} = \left( \{\Theta_{\text{Res}}^{(r)}\}_{r=1}^n, \{\Theta_{\text{MHA},\text{self}}^{(r)}\}_{r=1}^n, \{\Theta_{\text{MHA},\text{cross}}^{(r)}\}_{r=1}^n, \Theta_{\text{Down}} \right)$$

For an input  $\underline{\mathbf{X}} \in \mathbb{R}^{C \times H \times W}$ , time embedding  $\mathbf{h} \in \mathbb{R}^{d_t}$ , and text context  $\mathbf{C} \in \mathbb{R}^{M \times d_{\text{ctx}}}$ ,

$$\underline{\mathbf{H}}_0 = \underline{\mathbf{X}},$$

$$\underline{\mathbf{H}}_r = \text{ResnetBlock2D}_{\Theta_{\text{Res}}^{(r)}}(\underline{\mathbf{H}}_{r-1}, \mathbf{h}), \quad (68)$$

$$\mathbf{T}_r = \text{flatten}_{(H,W)}(\underline{\mathbf{H}}_r), \quad (69)$$

$$\mathbf{U}_r = \text{MultiheadAttention}_{\Theta_{\text{MHA},\text{self}}^{(r)}}(\mathbf{T}_r, \mathbf{T}_r), \quad (70)$$

$$\mathbf{V}_r = \text{MultiheadAttention}_{\Theta_{\text{MHA},\text{cross}}^{(r)}}(\mathbf{U}_r, \mathbf{C}), \quad (71)$$

$$\underline{\mathbf{H}}_r^* = \text{unflatten}_{(H,W)}(\mathbf{V}_r), \quad r = 1, \dots, n, \quad (72)$$

$$\text{CrossAttentionDownBlock2D}_{\Theta_{\text{XDown}}}(\underline{\mathbf{X}}, \mathbf{h}, \mathbf{C}) = \text{Downsample2D}_{\Theta_{\text{Down}}^{(2)}}(\underline{\mathbf{H}}_n^*). \quad (73)$$

**Definition 6.20** (CrossAttentionUpBlock2D). For learnable parameters

$$\Theta_{\text{XUp}} = \left( \{\Theta_{\text{Res}}^{(r)}\}_{r=1}^n, \{\Theta_{\text{MHA},\text{self}}^{(r)}\}_{r=1}^n, \{\Theta_{\text{MHA},\text{cross}}^{(r)}\}_{r=1}^n, \Theta_{\text{Up}} \right)$$

and for an input  $\underline{\mathbf{X}}$  from the bottom, a skip connection  $\underline{\mathbf{S}}$  with the same resolution, a time

embedding  $\mathbf{h}$ , and context  $\mathbf{C}$ :

$$\underline{\mathbf{Y}}_0 = \text{Concat}\left(\text{Upsample2D}_{\Theta_{\text{Up}}}^{(2)}(\underline{\mathbf{X}}), \underline{\mathbf{S}}\right), \quad (74)$$

$$\underline{\mathbf{Y}}_r = \text{ResnetBlock2D}_{\Theta_{\text{Res}}^{(r)}}(\underline{\mathbf{Y}}_{r-1}, \mathbf{h}), \quad r = 1, \dots, n, \quad (75)$$

$$\mathbf{T}_r = \text{flatten}_{(H,W)}(\underline{\mathbf{Y}}_r), \quad \mathbf{U}_r = \text{MultiheadAttention}_{\Theta_{\text{MHA,self}}^{(r)}}(\mathbf{T}_r, \mathbf{T}_r), \quad (76)$$

$$\mathbf{V}_r = \text{MultiheadAttention}_{\Theta_{\text{MHA,cross}}^{(r)}}(\mathbf{U}_r, \mathbf{C}), \quad \mathbf{Y}_r^* = \text{unflatten}_{(H,W)}(\mathbf{V}_r), \quad (77)$$

$$\text{CrossAttentionUpBlock2D}_{\Theta_{\text{XUp}}}(\underline{\mathbf{X}}, \underline{\mathbf{S}}, \mathbf{h}, \mathbf{C}) = \underline{\mathbf{Y}}_n^*. \quad (78)$$

## 6.11 VAE Decoder's Terminal Mapping (RGB Output; $1 \times 1$ conv)

**Definition 6.21** (Conv1x1 (Final Projection)). We define  $\text{Conv1x1}_{\Theta_{\text{out}}} := \text{Conv2d}_{\Theta_{\text{out}}}^{(1,1; 1,1; 0,0)}$  and use it for the mapping to RGB output  $\mathbb{R}^{3 \times H \times W}$ .

# 7 Why U-Net Reaches the Entire Area "Shallowly": Quantitative Comparison of Receptive Fields

## 7.1 Receptive Field of a Pure CNN (Conv2d only)

When  $L$  layers of Conv2d with kernel size  $k = 3$ , stride 1, and padding 1 are stacked, the **receptive field** in one dimension is

$$R_{\text{pure}}(L) = 1 + (k - 1)L = 1 + 2L. \quad (79)$$

The condition to reach the entire width  $W$  is  $R_{\text{pure}}(L) \geq W$ , i.e.,

$$L \geq \frac{W - 1}{2}. \quad (80)$$

## 7.2 Receptive Field of U-Net (with staged Downsample2D)

Performing  $\text{Downsample2D}^{(2)}$  with stride 2  $L$  times at each level, and performing  $n_\ell$   $3 \times 3$  Conv2d (stride 1) at each resolution, one step at the final (coarsest) level corresponds to  $2^L$  pixels in the original resolution. Therefore, the receptive field converted to the original resolution is

$$R_{\text{unet}} = 1 + \sum_{\ell=0}^L (2^\ell) \cdot (k - 1) n_\ell = 1 + 2 \sum_{\ell=0}^L 2^\ell n_\ell. \quad (81)$$

If we uniformly set  $n_\ell = n$ ,

$$R_{\text{unet}} = 1 + 2n(2^{L+1} - 1). \quad (82)$$

**Theorem 7.1** (U-Net reaches the entire area with  $O(\log W)$  depth). Assuming  $k = 3$  and  $n \geq 1$  layers at each level, the sufficient condition  $R_{\text{unet}} \geq W$  to reach the entire width  $W$  is

$$L \geq \log_2 \left( \frac{W-1}{2n} + 1 \right) - 1. \quad (83)$$

Therefore, the required number of levels  $L$  is  $O(\log W)$ , which is **significantly fewer layers** to express dependencies from end to end compared to the  $O(W)$  of a pure CNN (Conv2d only) in (80).

*Proof.* Substituting (82) into  $R_{\text{unet}} \geq W$ , we get  $1 + 2n(2^{L+1} - 1) \geq W \iff 2^{L+1} \geq \frac{W-1}{2n} + 1$ . Taking  $\log_2$ , we get  $L + 1 \geq \log_2 \left( \frac{W-1}{2n} + 1 \right)$ , which yields (83).  $\square$

## 8 Full Definition of U-Net and VAE Decoder "as Functions"

### 8.1 The U-Net (Conditional) Overall Function

**Definition 8.1** (Parametric Function of UNet2DConditionModel). The inputs are latent  $\underline{Z} \in \mathbb{R}^{C \times H \times W}$ , time  $t \in \mathbb{R}$ , and text embedding sequence  $\underline{C} \in \mathbb{R}^{M \times d_{\text{ctx}}}$ . The learnable parameter vector is

$$\Theta_{\text{U}} = \left( \Theta_{\text{TE}}, \{\Theta_\ell^\downarrow\}_{\ell=1}^L, \Theta^{\text{mid}}, \{\Theta_\ell^\uparrow\}_{\ell=1}^L, \Theta^{\text{out}} \right) \quad (84)$$

(where each  $\Theta$  is the concatenation of all coefficients from §6 for Conv2d, GroupNorm, SiLU, Linear, ScaledDotProductAttention, MultiheadAttention, TimestepEmbedding, ResnetBlock2D, Downsample2D, Upsample2D, DownBlock2D, UpBlock2D, CrossAttentionMidBlock2D, Conv1x1). Injecting the time embedding  $\underline{h} = \text{TimestepEmbedding}_{\Theta_{\text{TE}}}(t)$  into each residual block,

$$\underline{D}_0 = \underline{Z}, \quad (85)$$

$$\underline{D}_\ell = (\text{CrossAttention})\text{DownBlock2D}_{\Theta_\ell^\downarrow}(\underline{D}_{\ell-1}, \underline{h}, \underline{C}), \quad \ell = 1, \dots, L, \quad (86)$$

$$\underline{B} = \text{CrossAttentionMidBlock2D}_{\Theta^{\text{mid}}}(\underline{D}_L, \underline{h}, \underline{C}), \quad (87)$$

$$\underline{U}_L = \text{UpBlock2D}_{\Theta_L^\uparrow}(\underline{B}, \underline{D}_L, \underline{h}), \quad (88)$$

$$\underline{U}_{\ell-1} = (\text{CrossAttention})\text{UpBlock2D}_{\Theta_{\ell-1}^\uparrow}(\underline{U}_\ell, \underline{D}_{\ell-1}, \underline{h}, \underline{C}), \quad \ell = L, \dots, 1, \quad (89)$$

$$\underline{\hat{E}} = \text{Conv1x1}_{\Theta^{\text{out}}}(\underline{U}_0) \in \mathbb{R}^{C \times H \times W}, \quad (90)$$

$$\underline{\hat{E}} = \text{UNet2DConditionModel}_{\Theta_{\text{U}}}(\underline{Z}, t, \underline{C}). \quad (91)$$

<sup>a</sup>

<sup>a</sup>Diffusers UNet2DConditionModel (U-Net Conditional): <https://huggingface.co/docs/diffusers/api/>

**Remark 8.1.** Each block is a composition of the basic operations from §6, and the output spatial size continuously follows  $H \times W$  (**variable input/output sizes**).

## 8.2 The VAE Decoder Overall Function

Stable Diffusion series VAEs have a decoder structure (Decoder) included in AutoencoderKL. Here, we formulate the **deterministic architecture** of the decoder as a composition of ResnetBlock2D and Upsample2D.

**Definition 8.2** (Decoder (VAE Decoder)). For input latent  $\underline{Z} \in \mathbb{R}^{C_z \times H_z \times W_z}$ , using learnable parameters

$$\Theta_{\text{Dec}} = \left( \Theta^{\text{in}}, \{\Theta_\ell^\uparrow\}_{\ell=1}^{L_d}, \Theta^{\text{out}} \right) \quad (92)$$

$$\underline{H}_0 = \text{ResnetBlock2D}_{\Theta^{\text{in}}}(\underline{Z}, \mathbf{0}) \quad (\text{no time dependence, so } \mathbf{h} = \mathbf{0}), \quad (93)$$

$$\underline{H}_\ell = \text{ResnetBlock2D}_{\Theta_\ell^\uparrow}(\text{Upsample2D}_{\Theta_{\text{Up}}^{(\ell)}}^{(2)}(\underline{H}_{\ell-1}), \mathbf{0}), \quad \ell = 1, \dots, L_d, \quad (94)$$

$$\hat{\underline{X}} = \text{Conv1x1}_{\Theta^{\text{out}}}(\underline{H}_{L_d}) \in \mathbb{R}^{3 \times H \times W}, \quad (H = 2^{L_d} H_z, W = 2^{L_d} W_z). \quad (95)$$

**Definition 8.3** (AutoencoderKL (Decoder part)). We define the decoder mapping  $\mathcal{D}$  of AutoencoderKL as

$$\mathcal{D}_{\Theta_{\text{Dec}}} : \mathbb{R}^{C_z \times H_z \times W_z} \rightarrow \mathbb{R}^{3 \times (2^{L_d} H_z) \times (2^{L_d} W_z)}, \quad \mathcal{D}_{\Theta_{\text{Dec}}}(\underline{Z}) = \text{Decoder}_{\Theta_{\text{Dec}}}(\underline{Z}) \quad (96)$$

<sup>a</sup>

<sup>a</sup>Diffusers AutoencoderKL implementation (includes Decoder): [https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/autoencoder\\_kl.py](https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/autoencoder_kl.py), API: <https://huggingface.co/docs/diffusers/api/models/autoencoderkl>.

## 9 Summary (Correspondence to Learning Outcomes)

### Correspondence with Learning Outcomes

- **Mathematical description of architectures:** We formally defined as functions the U-Net and VAE decoder.
- **Explanation of variable I/O sizes:** We confirmed that U-Net/VAE satisfy variable input/output sizes because each layer is defined in a **form independent of spatial size**.
- **Explanation of differences from the proposal:** We clarified the configuration of the U-Net in image generation AIs is different from the originally proposed form.

## 10 Next Lecture Preview

Next time, we will explain the **Text Encoder**.

## References

- [1] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Advances in Neural Information Processing Systems (NeurIPS), 2020.
- [2] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In Advances in Neural Information Processing Systems (NeurIPS), 2022.
- [3] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In International Conference on Learning Representations (ICLR), 2014.
- [4] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022.
- [5] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In Medical Image Computing and Computer-Assisted Intervention (MICCAI), 2015.
- [6] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In International Conference on Learning Representations (ICLR), 2021.

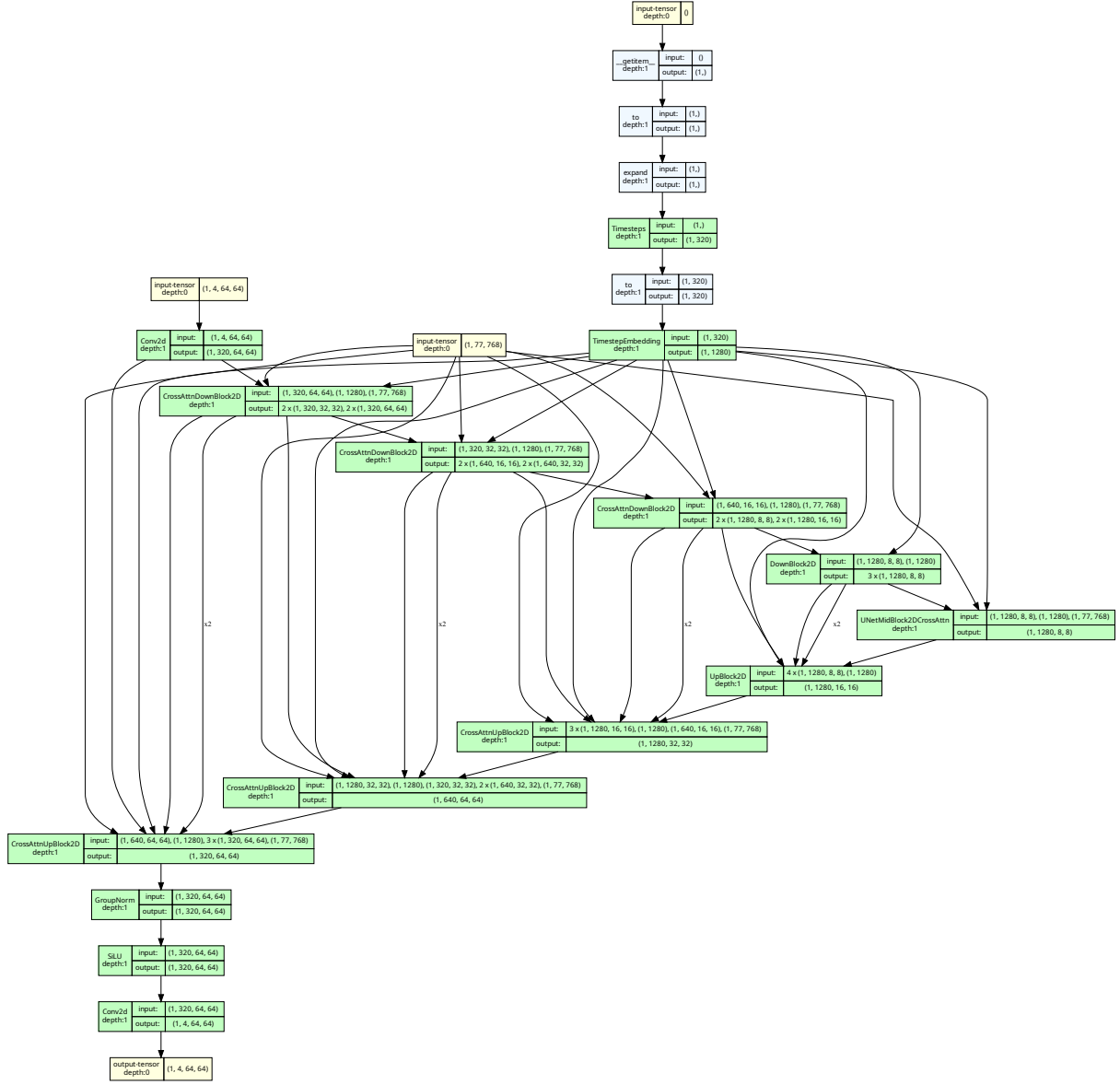


Figure 1: The neural network architecture of the U-Net in Stable Diffusion 1.5. The tensor size corresponds to the input with: batch size = 1, the number of channels = 4, width = 64, height = 64.

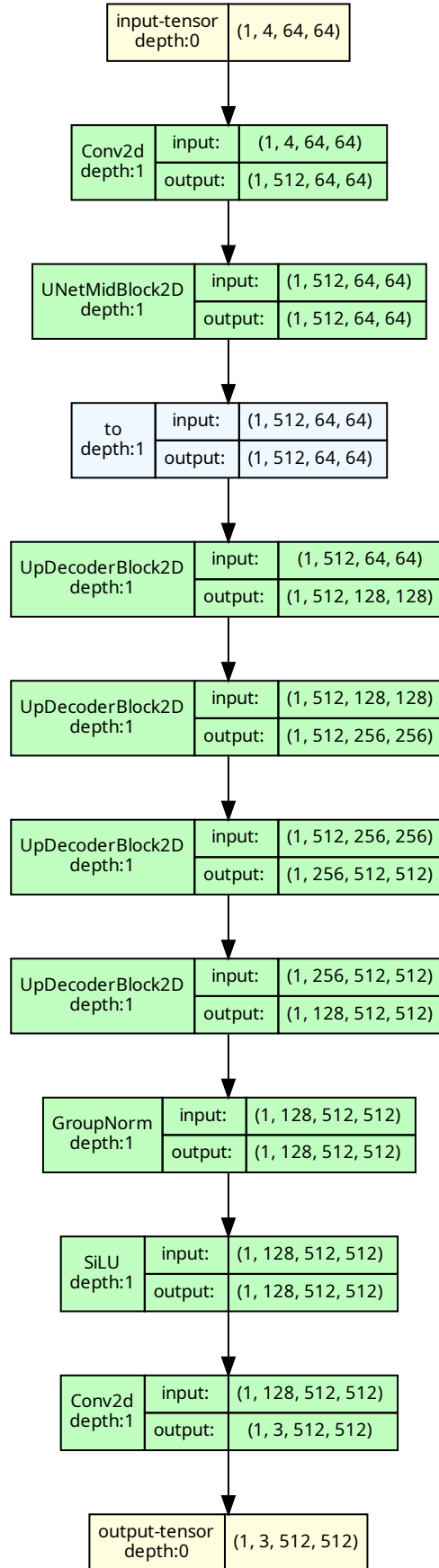


Figure 2: The neural network architecture of the VAE Decoder in Stable Diffusion 1.5. The tensor size corresponds to the input with: batch size = 1, the number of channels = 4, width = 64, height = 64.

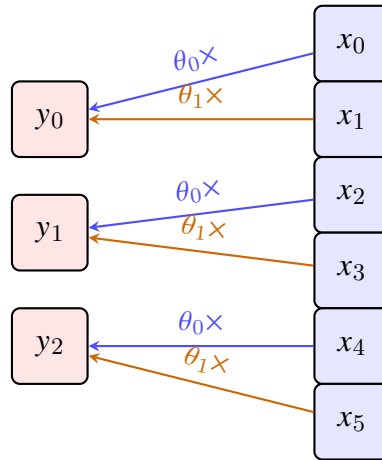


Figure 3: 1D **Downsample1D** (Example: input length 6, output length 3, filter width 2, stride 2). Identical weights  $\theta_0, \theta_1$  are shown by edges of the same color.

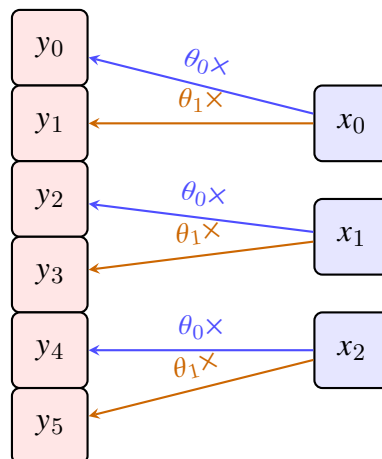


Figure 4: 1D **Upsample1D** (Example: input length 3, output length 6, filter width 2, stride 2). Identical weights  $\theta_0, \theta_1$  are shown by edges of the same color.