

AI Applications Lecture 16

Image Generation AI 6: Text Encoders for Image Generation

SUZUKI, Atsushi
Jing WANG

Contents

1	Introduction	3
1.1	Roadmap Recap	3
1.2	Learning Outcomes	3
2	Preliminaries: Mathematical Notations	3
3	What Does the Text Encoder Receive, and What Information Does It Pass On?	5
3.1	General Theory	5
3.2	In the Case of Stable Diffusion 1.5	6
4	Controlling the Final Output Image Using Prompt Weighting	7
4.1	Basic Idea	7
4.2	Text Encoder Class and Input/Output Correspondence	8
4.3	General Definition of Prompt Weighting	8
4.4	Functionalization of Automatic1111 Syntax and Compel Implementation Correspondence	9
5	Text Encoders in Practical Applications (Stable Diffusion 1.5's CLIP ViT-L/14 Text Encoder)	10
5.1	Overview via Diagram	10
5.2	Rigorous Sub-layer Definitions	10
5.3	CLIP ViT-L/14 Text Encoder Whole Model (Mask-Compatible)	14
6	Training the CLIP Text Encoder	14
6.1	CLIP Architecture and Objective Function	15

1 Introduction

1.1 Roadmap Recap

Let's review what we have learned so far. The core of the image generation AI pipeline is the **reverse diffusion process** managed by the **denoising scheduler**. The reverse diffusion process receives information from the **text encoder** and passes the low-resolution latent image to the **natural image decoder**, which is composed of a VAE. Then, the natural image decoder (composed of a VAE) converts the low-resolution latent image into a natural image. This time, we will explain the **text encoder**.

1.2 Learning Outcomes

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

- **Control the output image** by **prompt weighting**, which processes the **output of the text encoder**.
- Explain the **objective function used to train the text encoders** employed in practical image generation AI.
- **Mathematically describe the architecture of the text encoder** used in practical image generation AI and explain its characteristics.

2 Preliminaries: Mathematical Notations

- **Definitions:**
 - (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 extension notation).
 - $\{x \in \mathcal{A} \mid P(x)\}$: The set of elements in \mathcal{A} for which the proposition $P(x)$ is true (set-builder notation).
 - $|\mathcal{A}|$: The number of elements in set \mathcal{A} (used only for finite sets in this lecture).
 - \mathbb{R} : The set of all real numbers. $\mathbb{R}_{>0}$, $\mathbb{R}_{\geq 0}$ are defined similarly.

- \mathbb{Z} : The set of all integers. $\mathbb{Z}_{>0}$, $\mathbb{Z}_{\geq 0}$ 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 a mapping.
- $y = f(x)$ denotes the output $y \in \mathcal{Y}$ for an 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: $\mathbf{A} \in \mathbb{R}^{m,n}$.
- Elements are denoted as $a_{i,j}$:

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

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

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

- Transpose of a vector:

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

- **Tensor:**

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

3 What Does the Text Encoder Receive, and What Information Does It Pass On?

3.1 General Theory

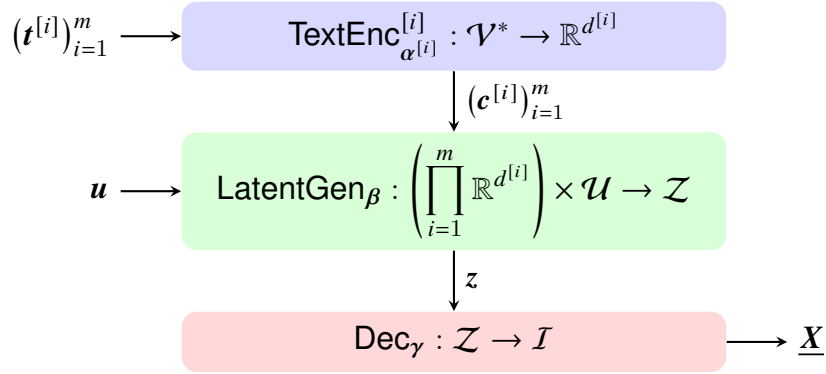


Figure 1: Recap of the three-component pipeline.

Figure 1 shows the pipeline of the text-to-image AI.

First, let's clarify the input and output of a text encoder in general terms. The input to the text encoder is originally a single byte sequence corresponding to the user's input prompt. Let's denote this as b . In reality, a **tokenizer** (represented as the function `Tokenizer` below), often based on **Byte Pair Encoding (BPE)**, is used to convert each byte sequence into a token sequence t . This is then fed into the neural network that constitutes the text encoder (hereafter $\text{TextEncoderNN}_{\theta_{\text{TE}}}$). The output $C = \text{TextEncoderNN}_{\theta_{\text{TE}}}(t)$ is a tensor that transforms the user's input prompt (or its corresponding token sequence) into a shape suitable for image generation, and this is passed to the **noise estimator** that constitutes the **reverse diffusion process**. Therefore, the size of the output tensor C must be acceptable as input by the noise estimator. If the noise estimator does not accept variable-sized inputs, the output size must be adjusted accordingly.

In practical applications, multiple prompts are input. For example, if there are a positive prompt b^+ and a negative prompt b^- , they are tokenized into t^+ and t^- respectively, and both $C^+ = \text{TextEncoderNN}_{\theta_{\text{TE}}}(t^+)$ and $C^- = \text{TextEncoderNN}_{\theta_{\text{TE}}}(t^-)$ are passed to the noise estimator of the reverse diffusion process.

3.2 In the Case of Stable Diffusion 1.5

Example 3.1 (Input/Output Dimensions in Stable Diffusion 1.5). In text, let L_{tok} be the input token length and d_{ctx} be the embedding dimension. In the CLIP ViT-L/14 Text Encoder of Stable Diffusion 1.5, the **input length is fixed** at $L_{\text{tok}} = 77$, and the **embedding dimension is also fixed** at $d_{\text{ctx}} = 768$. Therefore,

$$\text{Input: } \mathbf{t} \in \{1, 2, \dots, V\}^{77}, \quad \text{Output: } \mathbf{C} \in \mathbb{R}^{77 \times 768}. \quad (6)$$

Here, $V = 49408$ is the vocabulary size. The specific values in the Hugging Face CLIPTextModel implementation are Embedding(49408, 768), position_embedding(77, 768), 12 Transformer layers, and the shape of the final output last_hidden_state is (batch, 77, 768)^a.

^aThe Diffusers/Transformers CLIPTextModel implementation conforms to <https://huggingface.co/openai/clip-vit-large-patch14> and returns last_hidden_state and pooler_output. See the CLIP documentation in Transformers at https://huggingface.co/docs/transformers/en/model_doc/clip.

Remark 3.1 (Overview of Fixed-Length Conversion in SD1.5). In practical Stable Diffusion (Diffusers), **padding/truncation to a fixed length $L = 77$ is performed during tokenization**. In the CLIP tokenizer settings, the **PAD (padding) token is set to be identical to the EOS token**. Furthermore, the **attention mask** is a binary sequence with 1 for non-PAD tokens and 0 for PAD tokens. Below, we rigorously define the fixed-length normalization function LengthNormalization⁽⁷⁷⁾ and the attention mask generation function MaskPAD that masks PADs.

Definition 3.1 (Length Normalization Function LengthNormalization⁽⁷⁷⁾). Let the **content tokens** (excluding special tokens) obtained from the input string be $s = (s_1, \dots, s_m)$, and abbreviate BOS (begin-of-sentence), EOS (end-of-sentence), and PAD (padding) as BOS, EOS, PAD respectively. Define the **preprocessed sequence** $\tilde{\mathbf{t}}$ as

$$\tilde{\mathbf{t}} := (\text{BOS}, s_1, \dots, s_m, \text{EOS}). \quad (7)$$

(addition of BOS/EOS). At this time, the function for **length normalization**

$$\text{LengthNormalization}^{(77)} : \{1, \dots, V\}^* \rightarrow \{1, \dots, V\}^{77}, \quad (8)$$

$$\mathbf{t} = \text{LengthNormalization}^{(77)}(\tilde{\mathbf{t}}) \quad (9)$$

is defined as follows:

$$\mathbf{t} = \begin{cases} (\tilde{t}_1, \dots, \tilde{t}_\ell, \underbrace{\text{PAD}, \dots, \text{PAD}}_{77-\ell \text{ items}}) & \text{if } \ell := |\tilde{\mathbf{t}}| \leq 77, \\ (\tilde{t}_1, \dots, \tilde{t}_{76}, \text{EOS}) & \text{if } |\tilde{\mathbf{t}}| > 77. \end{cases} \quad (10)$$

That is, if the length is insufficient, it is right-padded with PAD (=EOS); if the length is excessive, the first 76 items are retained, and the last token is forcibly truncated to EOS.

Definition 3.2 (PAD Mask Generation Function MaskPAD (attention mask for PAD)). For the normalized sequence $\mathbf{t} = (t_1, \dots, t_{77})$ from equation (9), the function

$$\text{MaskPAD} : \{1, \dots, V\}^{77} \rightarrow \{0, 1\}^{77}, \quad (11)$$

$$\mathbf{a} = \text{MaskPAD}(\mathbf{t}) \quad (12)$$

that generates the **attention mask** $\mathbf{a} \in \{0, 1\}^{77}$ is defined as

$$a_i := \begin{cases} 1 & \text{if } t_i \neq \text{PAD}, \\ 0 & \text{if } t_i = \text{PAD}, \end{cases} \quad i \in [1, 77]_{\mathbb{Z}}, \quad (13)$$

\mathbf{a} is used for **masking the key side** in the **self-attention** mechanism used downstream.

4 Controlling the Final Output Image Using Prompt Weighting

4.1 Basic Idea

The text encoder outputs a tensor, which is passed to the noise estimator that constitutes the reverse diffusion process. **Prompt weighting** is a technique that processes this tensor between the text encoder and the noise estimator, reflecting user intentions that cannot be fully expressed in the natural language input to the text encoder.

Example 4.1 (Prompt weighting Syntax in Automatic1111). **Intuition for beginners: Emphasizing** a part of the text makes the visual elements related to that word **more likely to appear strongly**. Conversely, de-emphasizing makes them **less likely to appear**. The typical syntax in Automatic1111 is as follows^a.

- Parentheses: (word) emphasizes by 1.1 times, multiplication is applied for nesting (e.g., ((word)) is $\approx 1.1^2$).
- Brackets: [word] de-emphasizes by 0.9 times (can be nested).
- Explicit weight: (word:1.5) specifies 1.5 times.

Expected effect: For example, in (cinematic lighting:1.4), soft focus, features related to cinematic lighting are supplied **more strongly** to the cross-attention, making the lighting expression more likely to be emphasized.

^aAutomatic1111 WebUI: <https://github.com/AUTOMATIC1111/stable-diffusion-webui>. Much of the syntax is convention appearing in the WebUI documentation or implementation (Python), and the precise mathematical specification conforms to the implementation.

Remark 4.1. In **diffusers**, the above Prompt weighting syntax can also be used by employing the **Compel** library^a.

^aCompel: <https://github.com/damian0815/compel>. Diffusers: <https://github.com/huggingface/diffusers>.

4.2 Text Encoder Class and Input/Output Correspondence

Definition 4.1 (Text Encoder Class (Input/Output Size Correspondence)). A text encoder that takes a **token sequence** $t \in \{1, 2, \dots, V\}^{L_{\text{tok}}}$ as input and outputs a **context sequence** $C \in \mathbb{R}^{L_{\text{ctx}} \times d_{\text{ctx}}}$ is defined as

$$\text{TextEncoder}_{\theta_{\text{TE}}}^{(V, L_{\text{tok}}, L_{\text{ctx}}, d_{\text{ctx}})} : \{1, 2, \dots, V\}^{L_{\text{tok}}} \rightarrow \mathbb{R}^{L_{\text{ctx}} \times d_{\text{ctx}}} \quad (14)$$

The CLIP ViT-L/14 of Stable Diffusion 1.5 has $L_{\text{tok}} = L_{\text{ctx}} = 77$, $d_{\text{ctx}} = 768$, $V = 49408$ (see Example 3.1).

Example 4.2 (SD1.5 Text Encoder is Included in Definition 4.1). In Definition 4.1, setting $(V, L_{\text{tok}}, L_{\text{ctx}}, d_{\text{ctx}}) = (49408, 77, 77, 768)$ and letting θ_{TE} be all learnable parameters of CLIP ViT-L/14, the `last_hidden_state` of the `CLIPTextModel` in implementation provides $C \in \mathbb{R}^{77 \times 768}$.

4.3 General Definition of Prompt Weighting

Definition 4.2 (General Prompt Weighting (Action on Token Subsets)). Let the text encoder output be $C = [c_1^\top; \dots; c_{L_{\text{ctx}}}^\top] \in \mathbb{R}^{L_{\text{ctx}} \times d_{\text{ctx}}}$. The index set $[1, L_{\text{ctx}}]_{\mathbb{Z}}$ is partitioned into a family of **disjoint subsets** $\mathcal{S} = \{S_1, \dots, S_K\}$ (i.e., $S_k \subseteq [1, L_{\text{ctx}}]_{\mathbb{Z}}$, $S_i \cap S_j = \emptyset$, $\bigcup_{k=1}^K S_k = [1, L_{\text{ctx}}]_{\mathbb{Z}}$). For each subset S_k , a mapping $\psi_k : \mathbb{R}^{d_{\text{ctx}}} \rightarrow \mathbb{R}^{d_{\text{ctx}}}$ is given (classes such as continuous maps or linear maps are assumed, but not required). At this time, the **weighting operation** $\text{PW}_{\{\psi_k\}, \mathcal{S}}$ is defined as

$$\widehat{C} = \text{PW}_{\{\psi_k\}, \mathcal{S}}(C) \in \mathbb{R}^{L_{\text{ctx}} \times d_{\text{ctx}}}, \quad (15)$$

$$\widehat{c}_i = \psi_k(c_i) \quad \text{if } i \in S_k \quad (16)$$

Equations (15)–(16) mean that for each token position i , the mapping ψ_k associated with the subset S_k it belongs to is applied to the component vector c_i .

Remark 4.2. In practice, ψ_k is often **scalar multiplication** $\psi_k(x) = r_k x$ ($r_k > 0$), representing up-weighting. On the other hand, for down-weighting, implementations exist that are used in combination with **attention masks** or **linear blending**, and in these cases, it goes beyond simple output processing, becoming more complex by using a different context sequence based on input masked at the corresponding locations^a.

^aIn the Compel implementation, token weights are handled by scalar application to the embedding

tensor and, if necessary, masking/blending. See code: `get_embeddings_for_weighted_prompt_fragments` and `get_token_ids_and_expand_weights` in `embeddings_provider.py` (https://github.com/damian0815/compel/blob/main/src/compel/embeddings_provider.py). Furthermore, the policy of handling down-weighting **not by removal but by masking** is specified in the README/Changelog (see Changelog "1.0.0 - new downweighting algorithm" at <https://pypi.org/project/compel/>).

4.4 Functionalization of Automatic1111 Syntax and Compel Implementation Correspondence

Example 4.3 (Mapping ψ_k for Parentheses, Brackets, and Explicit Weights). Suppose a parser extracts a family of subsets $\mathcal{S} = \{S_1, \dots, S_K\}$ for text fragments (spans) and a scalar $r_k > 0$ for each S_k from the user prompt. At this time, if we define the **mapping** $\psi_k : \mathbb{R}^{d_{\text{ctx}}} \rightarrow \mathbb{R}^{d_{\text{ctx}}}$ as

$$\psi_k(\mathbf{x}) = r_k \mathbf{x} \quad (17)$$

then, according to Definition 4.2, a weight of r_k times is applied to the embedding \mathbf{c}_i at the corresponding token position $i \in S_k$. The correspondence between Automatic1111-style syntax and r_k is as follows:

$$(\text{span}) \rightsquigarrow r_k = \lambda \quad (\lambda = 1.1 \text{ default; } \lambda^{\#_{\text{nest}}} \text{ for nesting count}), \quad (18)$$

$$(\text{span: } r) \rightsquigarrow r_k = r \quad (r > 1). \quad (19)$$

In Compel, the syntax is as follows.

$$\text{span+} \rightsquigarrow r_k = \lambda \quad (\lambda = 1.1 \text{ default; } \lambda^{\#_{\text{nest}}} \text{ for nesting count}), \quad (20)$$

$$(\text{span})r \rightsquigarrow r_k = r \quad (r > 1). \quad (21)$$

In Compel, this r_k is expanded as **per-token weights** and applied to the embedding tensor in the form of equation (17)^a.

^aOne example of the application path: **per_token_weights** are constructed in `get_token_ids_and_expand_weights` of `embeddings_provider.py`, and reflected in the embeddings in `get_embeddings_for_weighted_prompt_fragments` (https://github.com/damian0815/compel/blob/main/src/compel/embeddings_provider.py).

Remark 4.3. In Compel's current implementation, **down-weighting** can be achieved with the following syntax:

$$\text{span-} \quad (22)$$

$$(\text{span})r \quad (0 < r < 1). \quad (23)$$

Here, `(some phrase: 0.9)` is equivalent to `[some phrase]`. Interestingly, unlike up-weighting, the implementation of down-weighting is not composed of scalar multiplication of parts of the output, but is designed to use linear combinations with **another context**

output using attention masks, and in this case, it does not stay within the scope of just processing the output (as it depends on things other than the original context output)^a.

^aChangelog "1.0.0 - new downweighting algorithm": <https://pypi.org/project/compel/>.

5 Text Encoders in Practical Applications (Stable Diffusion 1.5's CLIP ViT-L/14 Text Encoder)

This chapter describes the architecture of the **CLIP ViT-L/14 Text Encoder** used by Stable Diffusion 1.5 as the Text Encoder in Text-to-image (text-conditioned image generation). The text encoder of Stable Diffusion 1.5 has a fixed input token length of $L = 77$, and we will confirm which component this property originates from.

Remark 5.1. CLIP (Contrastive Language-Image Pre-training) is a **learning method**, and **ViT-L/14** is an **architecture**. CLIP is a learning paradigm that encompasses multiple architectures such as ResNet/ViT [1, 2].

5.1 Overview via Diagram

Figure 2 shows the architecture of the CLIP ViT-L/14 Text Encoder.

5.2 Rigorous Sub-layer Definitions

Below, we use function names corresponding to library class names, denote **learnable parameters** in **parentheses**, and explicitly indicate **hyperparameters** in **superscript parentheses**. We will also state the sizes of all vectors, matrices, and tensors, including inputs, outputs, and intermediate variables. **Attention-related functions are defined to accept attention masks** (The Transformers implementation applies `attention_mask` and `causal_attention_mask` sequentially)¹.

Definition 5.1 (Embedding). Hyperparameters are vocabulary size $V \in \mathbb{Z}_{>0}$ and embedding dimension $d \in \mathbb{Z}_{>0}$. The learnable parameters are

$$\Theta_{\text{Emb}} = (E \in \mathbb{R}^{V \times d}). \quad (24)$$

Taking a token sequence $\mathbf{t} \in \{1, \dots, V\}^L$ of length $L \in \mathbb{Z}_{>0}$ as input,

$$\text{Embedding}_{\Theta_{\text{Emb}}}^{(V,d)} : \{1, \dots, V\}^L \rightarrow \mathbb{R}^{L \times d}, \quad \text{Embedding}_{\Theta_{\text{Emb}}}^{(V,d)}(\mathbf{t}) = \mathbf{X} \in \mathbb{R}^{L \times d}, \quad (25)$$

¹In the Hugging Face Transformers CLIP implementation (`modeling_clip.py`), `causal_attention_mask` and `attention_mask` are combined and used in the text-side attention. Reference: https://huggingface.co/transformers/v4.8.0/_modules/transformers/models/clip/modeling_clip.html, latest documentation: https://huggingface.co/docs/transformers/en/model_doc/clip. For an application example using PyTorch's SDPA (scaled dot product attention), see https://huggingface.co/microsoft/LLM2CLIP-Openai-B-16/blob/main/modeling_clip.py.

$$\text{where } X_{i,:} = E_{t_i,:} \quad (1 \leq i \leq L). \quad (26)$$

(This function does not accept an attention mask.)

Definition 5.2 (LayerNorm (Layer Normalization)). For an input $X \in \mathbb{R}^{L \times d}$, the learnable parameters are

$$\Theta_{\text{LN}} = (\gamma \in \mathbb{R}^d, \beta \in \mathbb{R}^d). \quad (27)$$

The mean and variance at each position $i \in \{1, \dots, L\}$ are

$$\mu_i = \frac{1}{d} \sum_{j=1}^d X_{i,j} \in \mathbb{R}, \quad (28)$$

$$\sigma_i^2 = \frac{1}{d} \sum_{j=1}^d (X_{i,j} - \mu_i)^2 \in \mathbb{R}, \quad (29)$$

and the output $Y = \text{LayerNorm}_{\Theta_{\text{LN}}}(X) \in \mathbb{R}^{L \times d}$ is defined as

$$Y_{i,j} = \gamma_j \frac{X_{i,j} - \mu_i}{\sqrt{\sigma_i^2 + \varepsilon}} + \beta_j \quad (1 \leq i \leq L, 1 \leq j \leq d) \quad (30)$$

(This function does not accept an attention mask.)

Definition 5.3 (QuickGELUActivation (QuickGELU Activation)). For an element $u \in \mathbb{R}$,

$$\text{QuickGELU}(u) = u \sigma(1.702 u), \quad \sigma(x) = \frac{1}{1 + e^{-x}}. \quad (31)$$

(This function does not accept an attention mask.)

Definition 5.4 ($\text{Softmax}^{(i)}$ (Softmax; Axis Specification)). For a tensor (matrix) $A \in \mathbb{R}^{m \times n}$, specifying axis $i \in \{1, 2\}$,

$$(\text{Softmax}^{(1)}(A))_{p,q} = \frac{\exp(A_{p,q})}{\sum_{p'=1}^m \exp(A_{p',q})} \in \mathbb{R}, \quad (32)$$

$$(\text{Softmax}^{(2)}(A))_{p,q} = \frac{\exp(A_{p,q})}{\sum_{q'=1}^n \exp(A_{p,q'})} \in \mathbb{R} \quad (33)$$

(This function does not accept an attention mask.)

Definition 5.5 (CLIPAttention (Self-Attention; Mask-Compatible)). Takes input $X \in \mathbb{R}^{L \times d}$ and **attention mask** $M \in \{0, 1\}^{L \times L}$. Hyperparameters are **number of heads** $h \in \mathbb{Z}_{>0}$ and **dimension per head** $d_h \in \mathbb{Z}_{>0}$ ($d = h d_h$). The learnable parameters are

$$\Theta_{\text{Attn}} = \left(\{W_Q^{(\ell)} \in \mathbb{R}^{d \times d_h}, W_K^{(\ell)} \in \mathbb{R}^{d \times d_h}, W_V^{(\ell)} \in \mathbb{R}^{d \times d_h}\}_{\ell=1}^h, W_O \in \mathbb{R}^{(h d_h) \times d} \right). \quad (34)$$

For each head $\ell \in \{1, \dots, h\}$,

$$\mathbf{Q}^{(\ell)} = \mathbf{X} \mathbf{W}_Q^{(\ell)} \in \mathbb{R}^{L \times d_h}, \quad \mathbf{K}^{(\ell)} = \mathbf{X} \mathbf{W}_K^{(\ell)} \in \mathbb{R}^{L \times d_h}, \quad \mathbf{V}^{(\ell)} = \mathbf{X} \mathbf{W}_V^{(\ell)} \in \mathbb{R}^{L \times d_h}, \quad (35)$$

$$\mathbf{S}^{(\ell)} = \frac{\mathbf{Q}^{(\ell)} (\mathbf{K}^{(\ell)})^\top}{\sqrt{d_h}} \in \mathbb{R}^{L \times L}. \quad (36)$$

Define **mask application** as

$$(\tilde{\mathbf{S}}^{(\ell)})_{i,j} = \begin{cases} S_{i,j}^{(\ell)} & \text{if } M_{i,j} = 1, \\ -\infty & \text{if } M_{i,j} = 0, \end{cases} \quad (37)$$

(where $-\infty$ implies the limit operation in softmax), and

$$\mathbf{A}^{(\ell)} = \text{Softmax}^{(2)}(\tilde{\mathbf{S}}^{(\ell)}) \mathbf{V}^{(\ell)} \in \mathbb{R}^{L \times d_h}, \quad (38)$$

$$\text{CLIPAttention}_{\Theta_{\text{Attn}}}^{(h, d_h)}(\mathbf{X}, \mathbf{M}) = [\mathbf{A}^{(1)} \mid \dots \mid \mathbf{A}^{(h)}] \mathbf{W}_O \in \mathbb{R}^{L \times d}. \quad (39)$$

(In the Transformers implementation, \mathbf{M} is passed as a combination of `causal_attention_mask` and `attention_mask`, and applied internally with SDPA or an equivalent additive mask^a.)

^aAn example of mask composition in CLIP attention: https://huggingface.co/microsoft/LLM2CLIP-Openai-B-16/blob/main/modeling_clip.py.

Definition 5.6 (CLIPMLP (Per-Position MLP)). For input $\mathbf{X} \in \mathbb{R}^{L \times d}$, the intermediate dimension $d_{\text{mlp}} \in \mathbb{Z}_{>0}$ is a hyperparameter, and the learnable parameters are

$$\Theta_{\text{MLP}} = (\mathbf{W}_1 \in \mathbb{R}^{d \times d_{\text{mlp}}}, \mathbf{b}_1 \in \mathbb{R}^{d_{\text{mlp}}}, \mathbf{W}_2 \in \mathbb{R}^{d_{\text{mlp}} \times d}, \mathbf{b}_2 \in \mathbb{R}^d). \quad (40)$$

$$\mathbf{H}_1 = \mathbf{X} \mathbf{W}_1 + \mathbf{1}_L \mathbf{b}_1^\top \in \mathbb{R}^{L \times d_{\text{mlp}}}, \quad (41)$$

$$\mathbf{H}_2 = \text{QuickGELU}(\mathbf{H}_1) \in \mathbb{R}^{L \times d_{\text{mlp}}}, \quad (42)$$

$$\text{CLIPMLP}_{\Theta_{\text{MLP}}}^{(d_{\text{mlp}})}(\mathbf{X}) = \mathbf{H}_2 \mathbf{W}_2 + \mathbf{1}_L \mathbf{b}_2^\top \in \mathbb{R}^{L \times d}. \quad (43)$$

(This function does not accept an attention mask.)

Definition 5.7 (CLIPEncoderLayer (Encoder Layer; Mask-Compatible)). The learnable parameters are

$$\Theta_{\text{EL}} = (\Theta_{\text{LN1}}, \Theta_{\text{Attn}}, \Theta_{\text{LN2}}, \Theta_{\text{MLP}}). \quad (44)$$

For input $X \in \mathbb{R}^{L \times d}$ and attention mask $M \in \{0, 1\}^{L \times L}$,

$$U = \text{LayerNorm}_{\Theta_{\text{LN1}}}(X) \in \mathbb{R}^{L \times d}, \quad (45)$$

$$H = \text{CLIPAttention}_{\Theta_{\text{Attn}}}^{(h, d_h)}(U, M) \in \mathbb{R}^{L \times d}, \quad (46)$$

$$Y = X + H \in \mathbb{R}^{L \times d}, \quad V = \text{LayerNorm}_{\Theta_{\text{LN2}}}(Y) \in \mathbb{R}^{L \times d}, \quad (47)$$

$$Z = \text{CLIPMLP}_{\Theta_{\text{MLP}}}^{(d_{\text{mlp}})}(V) \in \mathbb{R}^{L \times d}, \quad \text{CLIPEncoderLayer}_{\Theta_{\text{EL}}}(X, M) = Y + Z \in \mathbb{R}^{L \times d}. \quad (48)$$

(In Transformers, the corresponding layer `CLIPEncoderLayer` accepts `attention_mask` and propagates it internally to the attention^a.)

^aAn example implementation (documentation-generated version): https://huggingface.co/transformers/v4.11.3/_modules/transformers/models/clip/modeling_clip.html.

Definition 5.8 (CLIPEncoder (Encoder Stack; Mask-Compatible)). The number of layers $N \in \mathbb{Z}_{>0}$ is a hyperparameter, and the learnable parameters are

$$\Theta_{\text{Enc}} = (\Theta_{\text{EL}}^{(1)}, \dots, \Theta_{\text{EL}}^{(N)}, \Theta_{\text{LNf}}). \quad (49)$$

For input $X \in \mathbb{R}^{L \times d}$ and attention mask $M \in \{0, 1\}^{L \times L}$,

$$H_0 = X \in \mathbb{R}^{L \times d}, \quad (50)$$

$$H_\ell = \text{CLIPEncoderLayer}_{\Theta_{\text{EL}}^{(\ell)}}(H_{\ell-1}, M) \in \mathbb{R}^{L \times d} \quad (\ell = 1, \dots, N), \quad (51)$$

$$\text{CLIPEncoder}_{\Theta_{\text{Enc}}}^{(N)}(X, M) = \text{LayerNorm}_{\Theta_{\text{LNf}}}(H_N) \in \mathbb{R}^{L \times d}. \quad (52)$$

(CLIPTextTransformer passes `attention_mask` to the encoder^a.)

^aCLIP text transformer forward arguments: https://huggingface.co/docs/transformers/en/model_doc/clip.

Definition 5.9 (CLIPTextEmbeddings (Token + Position Embedding Sum)). Uses learnable parameters for token embedding $\Theta_{\text{Tok}} = (E_{\text{tok}} \in \mathbb{R}^{V \times d})$ and learnable parameters for position embedding $\Theta_{\text{Pos}} = (E_{\text{pos}} \in \mathbb{R}^{L \times d})$.

$$T = \text{Embedding}_{\Theta_{\text{Tok}}}^{(V, d)}(t) \in \mathbb{R}^{L \times d}, \quad (53)$$

$$P = \text{Embedding}_{\Theta_{\text{Pos}}}^{(L, d)}([0, 1, \dots, L-1]) \in \mathbb{R}^{L \times d}, \quad (54)$$

$$\text{CLIPTextEmbeddings}_{\Theta_{\text{Tok}}, \Theta_{\text{Pos}}}^{(V, L, d)}(t) = T + P \in \mathbb{R}^{L \times d}. \quad (55)$$

(This function does not accept an attention mask.)

5.3 CLIP ViT-L/14 Text Encoder Whole Model (Mask-Compatible)

Definition 5.10 (CLIP ViT-L/14 Text Encoder). The hyperparameters are

$$(V, L, d, h, d_h, N, d_{\text{mlp}}) = (49408, 77, 768, 12, 64, 12, 3072) \quad (\text{where } d = h d_h) \quad (56)$$

The set of all learnable parameters is

$$\theta_{\text{TE}} = (\Theta_{\text{Tok}}, \Theta_{\text{Pos}}, \Theta_{\text{Enc}}) \quad (57)$$

For input $t \in \{1, \dots, V\}^L$ and PAD mask sequence $a = \text{MaskPAD}(t)$ (Definition 3.2),

$$X_0 = \text{CLIPTextEmbeddings}_{\Theta_{\text{Tok}}, \Theta_{\text{Pos}}}^{(V, L, d)}(t) \in \mathbb{R}^{L \times d}, \quad (58)$$

$$C_{\text{causal}} = [\mathbf{1}_{\{j \leq i\}}]_{i,j=1}^L \in \{0, 1\}^{L \times L} \quad (\text{Lower-triangular causal mask}), \quad (59)$$

$$M_{\text{pad}} = [a_j]_{i,j=1}^L \in \{0, 1\}^{L \times L} \quad (\text{Column replication of PAD mask to key direction}), \quad (60)$$

$$M = C_{\text{causal}} \odot M_{\text{pad}} \in \{0, 1\}^{L \times L}, \quad (61)$$

$$X_N = \text{CLIPEncoder}_{\Theta_{\text{Enc}}}^{(N)}(X_0, M) \in \mathbb{R}^{L \times d}, \quad (62)$$

$$\text{CLIPTextModel}_{\theta_{\text{TE}}}^{(V, L, d, h, d_h, N, d_{\text{mlp}})}(t) = X_N \in \mathbb{R}^{L \times d}. \quad (63)$$

That is, M , which combines the **causal mask** and the **PAD mask (attention mask for padding)**, is passed to the self-attention in each layer. The Hugging Face implementation `CLIPTextModel` also accepts an `attention_mask` argument and applies it internally by combining it with the `causal_attention_mask`^a. In Stable Diffusion, $C = X_N$ is supplied to the U-Net's **cross-attention**.

^aCLIP text model API: https://huggingface.co/docs/transformers/en/model_doc/clip. Implementation example: https://huggingface.co/transformers/v4.8.0/_modules/transformers/models/clip/modeling_clip.html.

Remark 5.2. The fixed length $L = 77$ is stipulated by the position embedding $\text{Embedding}^{(L, d)}$ (Equation (54)) and the tensor shaping based on it. This design is different from relative position embeddings or rotational position embeddings (RoPE), and this is the direct origin of the **fixed token length** in SD1.5 [2].

6 Training the CLIP Text Encoder

What output was the Text Encoder targeted to produce and trained for? In other words, under what objective function was it trained? The Text Encoder in the Text-to-image pipeline is expected to take a tensor mechanically obtained from the input prompt and output a tensor

that contains that information and is useful for image generation. That is, we want the output to reflect not only the abstract meaning of the text but also the similarity when represented as an image. The CLIP ViT-L/14 adopted by Stable Diffusion 1.5 is trained independently of image generation, but it is trained by **CLIP (Contrastive Language-Image Pre-training)**, which reflects the similarity structures in both images and text [2]. This training, which considers both images and text, is thought to be one reason why the image generation pipeline that includes it can show high performance even when it is used frozen as a text-to-tensor converter without being trained for image generation. This chapter describes CLIP, especially its objective function.

Remark 6.1. Since the **actual training source code** used for the CLIP ViT-L/14 in Stable Diffusion 1.5 does not seem to be publicly available, we will explain based on `open_clip` (`mlfoundations/open_clip`) here^a.

^aOpenCLIP repository: https://github.com/mlfoundations/open_clip.

6.1 CLIP Architecture and Objective Function

Definition 6.1 (CLIP Components and Contrastive Loss). Consider text sequences and image sequences of batch size N , $(t_i, \underline{I}_i)_{i=1}^N$.

- **Text Encoder** $\text{TextEncoder}_{\theta_{\text{TE}}} : \{1, 2, \dots, V\}^L \rightarrow \mathbb{R}^d$ (corresponding to Definition ??).
- **Image Encoder** $\text{ImageEncoder}_{\theta_{\text{IE}}} : \mathcal{X}_{\text{img}} \rightarrow \mathbb{R}^d$ (ViT/ResNet etc. [1]).
- **Temperature** $\tau > 0$ (often made learnable in implementations).

For each sample pair (t_i, \underline{I}_i) ,

$$\mathbf{u}_i = \frac{\text{TextEncoder}_{\theta_{\text{TE}}}(t_i)}{\|\text{TextEncoder}_{\theta_{\text{TE}}}(t_i)\|_2}, \quad \mathbf{v}_i = \frac{\text{ImageEncoder}_{\theta_{\text{IE}}}(\underline{I}_i)}{\|\text{ImageEncoder}_{\theta_{\text{IE}}}(\underline{I}_i)\|_2} \in \mathbb{R}^d. \quad (64)$$

The similarity matrix $\mathbf{S} \in \mathbb{R}^{N \times N}$ is

$$S_{ij} = \frac{\langle \mathbf{u}_i, \mathbf{v}_j \rangle}{\tau}. \quad (65)$$

The **cross-entropy losses** for text→image and image→text are

$$\mathcal{L}_{\text{Text} \rightarrow \text{Image}} = \frac{1}{N} \sum_{i=1}^N \left(-\log \frac{\exp(S_{ii})}{\sum_{j=1}^N \exp(S_{ij})} \right), \quad (66)$$

$$\mathcal{L}_{\text{Image} \rightarrow \text{Text}} = \frac{1}{N} \sum_{i=1}^N \left(-\log \frac{\exp(S_{ii})}{\sum_{j=1}^N \exp(S_{ji})} \right), \quad (67)$$

and the **symmetric loss**

$$\mathcal{L}_{\text{CLIP}}(\theta_{\text{TE}}, \theta_{\text{IE}}, \tau) = \frac{1}{2} (\mathcal{L}_{\text{Text} \rightarrow \text{Image}} + \mathcal{L}_{\text{Image} \rightarrow \text{Text}}) \quad (68)$$

is minimized.

Remark 6.2. Equation (68) applies pressure to make the inner product (cosine similarity) of the corresponding pair (i, i) **larger than others**. The temperature τ controls the sharpness of the softmax, and as $\tau \downarrow 0$, it approaches harder discrimination [2]. The OpenCLIP implementation conforms to this formulation^a.

^aThe OpenCLIP loss implementation is included in the training code at https://github.com/mlfoundations/open_clip.

7 Summary and Next Time

Correspondence with Learning Outcomes

- **Control via prompt weighting:** We formalized the weight action on token groups with Definition 4.2, and explained the correspondence with practical Automatic111/Compel syntax in Examples 4.1–4.3 and Remark 4.3.
- **Explanation of objective function:** We specified CLIP’s **symmetric contrastive learning loss** in Definition 6.1 and supplemented the intuition in Remark 6.2 (see [2]).
- **Mathematical description of architecture:** We **rigorously defined as functions** all layers of the **CLIP ViT-L/14 Text Encoder** in Definitions ??–??, and stated the origin of the fixed length in Remark 5.2 (see [1, 2]).

Next Time

From next time, we will describe the **learning (parameter update) of neural networks**.

References

- [1] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In International Conference on Learning Representations (ICLR), 2021.
- [2] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In Proceedings of the 38th International Conference on Machine Learning (ICML), 2021. CLIP.

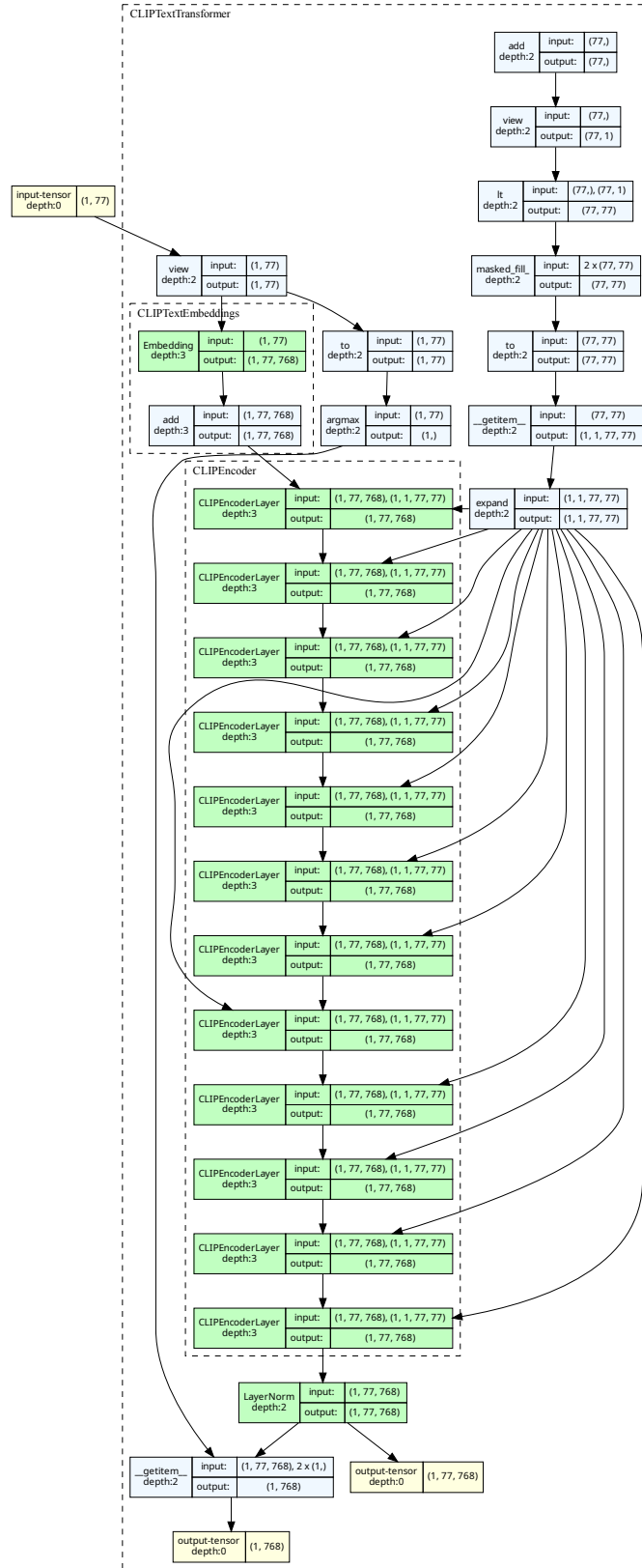


Figure 2: The neural network architecture of the CLIP ViT-L/14 text encoder in Stable Diffusion 1.5. The tensor size corresponds to the input with: batch size = 1. The left most tensor is the input token sequence, and the right-hand side corresponds to the attention mask. The left-hand side main stream is the core part with CLIPEncoderLayers. The output tensor with size $(1, 768)$ refers to the vector corresponding to the EOS token, which is also ignored in the image generation task.

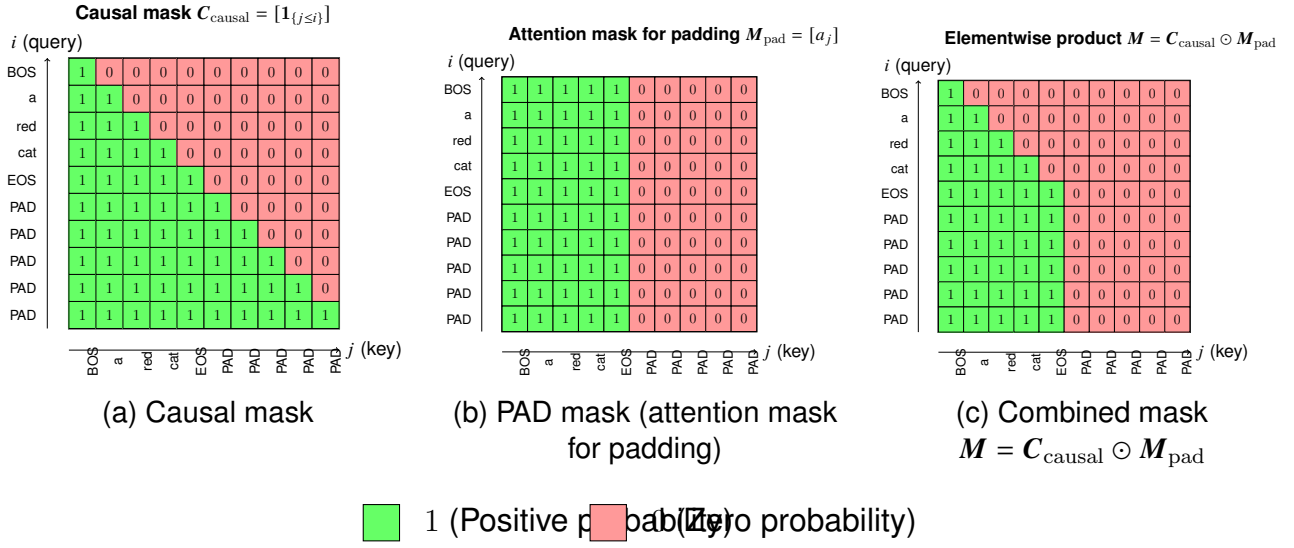


Figure 3: Visualization of the causal mask, PAD mask (attention mask for padding), and the elementwise product $M = C_{\text{causal}} \odot M_{\text{pad}}$. The token sequence $\{\text{BOS}, \text{a}, \text{red}, \text{cat}, \text{EOS}, \text{PAD}, \dots\}$ is shown as an example on the vertical and horizontal axes. The matrix and each element are drawn with a 1 : 1 aspect ratio.

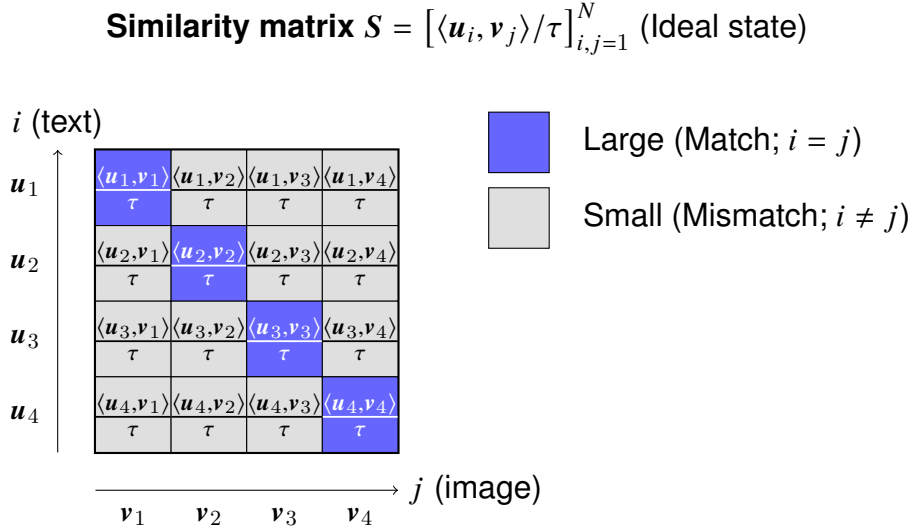


Figure 4: Visualization of the ideal similarity matrix. The color indicates the state where diagonal components $\langle u_i, v_i \rangle / \tau$ are large and off-diagonal components $\langle u_i, v_j \rangle / \tau$ ($i \neq j$) are small. The corresponding vectors u_i (text) and v_j (image) are displayed outside each row and column. The matrix and each element are drawn with a 1 : 1 aspect ratio.