

Notations and Definitions

The architectures used in the experiments of this thesis are based on the Transformer [1], and vision Transformer [2] architecture. Therefore, both image and text are represented as sequences of embeddings, which are processed by the Transformer blocks.

Image Representation

We define an image as a 3-dimensional tensor $\mathbf{v} \in \mathbb{R}^{C \times H \times W}$. Because we will use the base variant of the vision Transformer, ViT-B/16 [2], the image is patchified into 14x14 patches, each being a square of size 16x16 pixels. Each image patch represents one timestep in the sequence, and the number of patches N is given by $N = H \times \frac{W}{P^2}$, with P being the number of patches per dimension, and $P = 14$. Since we use an image size of 224x224 pixels, so $\mathbf{v} \in \mathbb{R}^{3 \times 224 \times 224}$, we will have $N = 224 \times \frac{224}{14^2} = 196$ patches, or timesteps respectively. Each patch is flattened into a 256-dimensional vector, and then projected into a 768 dimensions $\mathbf{e}_i^v \in \mathbb{R}^{768}$, using a fully connected layer. The image sequence is prepended with a special learnable $[\mathbf{I_CLS}] \in \mathbb{R}^{768}$ token, which is used to aggregate the global information/content of the image, and following [2]. The result is a sequence of patch embeddings, which we define as \mathbf{E}_v , where v indicates an image:

$$\mathbf{E}_v = [\mathbf{e}_{[\mathbf{I_CLS}]}^v, \mathbf{e}_1^v, \mathbf{e}_2^v, \dots, \mathbf{e}_N^v] \quad (1)$$

To give the Transformer a sense of order in the image patches/timestep, a unique positional encoding is added to each patch embedding. This can either be learned or fixed, with the latter being for example a sinusoidal positional encoding [1]. This positional encoding is also represented as a sequence of 768-dimensional vectors:

$$\mathbf{T}_v^{\text{pos}} = [0, \mathbf{t}_{\text{pos}_1}^v, \mathbf{t}_{\text{pos}_2}^v, \dots, \mathbf{t}_{\text{pos}_N}^v] \quad (2)$$

Since the $[\mathbf{I_CLS}]$ token is not part of the image, the positional encoding for the $[\mathbf{I_CLS}]$ token is set to zero, so nothing is added to it. An image representation is defined as:

$$\mathbf{H}_{v,l}^s = [\mathbf{h}_{v,l,[\mathbf{I_CLS}]}^s, \mathbf{h}_{v,l,1}^s, \dots, \mathbf{h}_{v,l,N}^s] \quad (3)$$

In Equation 3, l denotes the layer of the Transformer block that returned the image representation, and v indicates that the representation is an image. Since we use Knowledge Distillation (KD) in some parts of this thesis, representations will be, if necessary, superscripted with s or t , for a student and teacher representation, respectively.

We define $l = 0$ as the input to the Transformer, and $l = L$ as the output of the Transformer, where L is the number of layers in the Transformer. Consequently, the image input to the Transformer is defined as:

$$\mathbf{H}_{v,0}^s = [\mathbf{h}_{v,0,[\mathbf{I_CLS}]}^s, \mathbf{h}_{v,0,1}^s, \dots, \mathbf{h}_{v,0,N}^s] = \mathbf{E}_v + \mathbf{T}_v^{\text{pos}} \quad (4)$$

The output of the Transformer is defined as:

$$\mathbf{H}_{v,L}^s = [\mathbf{h}_{v,L,[\mathbf{I_CLS}]}^s, \mathbf{h}_{v,L,1}^s, \dots, \mathbf{h}_{v,L,N}^s] \quad (5)$$

Text Representation

We define a text as a sequence of discrete tokens, which are, similar to image patches, embedded into 768-dimensional vectors using an embedding matrix. A single token i is represented as $\mathbf{e}_i^t \in \mathbb{R}^{768}$, and the sequence of tokens, representing the text, is prepended with a start-of-sequence token $[\mathbf{T_CLS}] \in \mathbb{R}^{768}$, and appended with an end-of-sequence token $[\mathbf{T_SEP}] \in \mathbb{R}^{768}$. The purpose of the $[\mathbf{T_CLS}]$ token is, as with $[\mathbf{I_CLS}]$, to aggregate the global information/content of the text. The $[\mathbf{T_SEP}]$

token is used to indicate the end of the text sequence. A text sequence consists of M tokens, and we use w to denote a text sequence:

$$E_w = [e_{[T_CLS]}^w, e_1^w, e_2^w, \dots, e_M^w, e_{[T_SEP]}^w] \quad (6)$$

The maximum text sequence length M is not fixed, and will be defined when necessary in the experimental part of this work.

As mentioned (TODO: cite data preparation), to obtain discrete tokens, a sentence is tokenized into subwords using the GPT-2 byte-pair encoder, so one token does not necessarily represent a whole word.

A positional encoding is also added to the text embeddings, to give the Transformer a sense of order in the text sequence. Since the special token $[T_SEP]$ denotes the end of the text sequence, it is part of the sequence, and therefore has a positional encoding. The latter does not hold for the $[T_CLS]$ token, as it is used to aggregate the global information/content of the text.

$$T_w^{\text{pos}} = [0, t_{\text{pos}_1}^w, t_{\text{pos}_2}^w, \dots, t_{\text{pos}_M}^w, t_{\text{pos}_{[T_SEP]}}^w] \quad (7)$$

A text representation is defined as:

$$H_{w,l}^s = [h_{w,l,[T_CLS]}^s, h_{w,l,1}^s, \dots, h_{w,l,M}^s, h_{w,l,[T_SEP]}^s] \quad (8)$$

Equation 8 denotes the representation denoted by a student model s , but it can also be a teacher representation t .

The input to the Transformer for text is Equation 8 with $l = 0$, and the output of the Transformer is Equation 8 with $l = L$.

Transformer Block

Unless we use pretrained architectures that follow a different architecture, which we will then specify, we follow the Pre-LayerNorm definition of the Transformer block as given in [3]. As the name suggests, it applies LayerNorm before the Multi-Head Attention, instead of after.

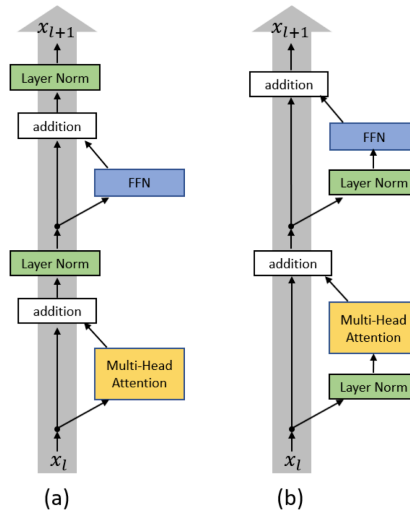


Figure 1: Comparison of a Post-Norm Transformer block/layer (a), and a Pre-Norm Transformer block/layer (b). (a) is the architecture as defined in the original “Attention is all you need” paper [1].

We follow the Pre-Norm architecture [3].

One Transformer block performs the following operations:

$$\mathbf{H}'_l = \text{MHA}(\text{LN}(\mathbf{H}_{l-1})) + \mathbf{H}_{l-1} \quad (9)$$

$$\mathbf{H}_l = \text{FFN}(\mathbf{H}_{l-1}) + \mathbf{H}'_l \quad (10)$$

We denote LN as LayerNorm, MHA as Multi-Head Attention, and FFN as a 2 layer MLP, all following the original Transformer of [1]. As previously mentioned, the only difference is the order of operations [3]. $\mathbf{H}_{v,l}^s$ and $\mathbf{H}_{w,l}^s$ can be used as a drop-in replacement for image and text, respectively. Both equations are inspired by VLMo [4].

We define a Transformer as multiple Transformer blocks stacked on top of each other.

Bibliography

- [1] A. Vaswani *et al.*, “Attention is all you need,” in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, in NIPS'17. Long Beach, California, USA: Curran Associates Inc., 2017, pp. 6000–6010.
- [2] A. Dosovitskiy *et al.*, “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale,” *ICLR*, 2021.
- [3] R. Xiong *et al.*, “On layer normalization in the transformer architecture,” in *Proceedings of the 37th International Conference on Machine Learning*, in ICML'20. JMLR.org, 2020.
- [4] H. Bao *et al.*, “VLMo: Unified Vision-Language Pre-Training with Mixture-of-Modality-Experts,” in *Advances in Neural Information Processing Systems*, 2022. [Online]. Available: <https://openreview.net/forum?id=bydKs84JEyw>