

## 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-dimensional vector  $\mathbf{e}_i^v \in \mathbb{R}^{768}$ , using a fully connected layer. The image sequence is prepended with a special 768-dimensional learnable [I\_CLS]  $\in \mathbb{R}^{768}$  token, which is used to aggregate the global information/content of the image, and, following [2], also 768-dimensional. The result is a sequence of patch embeddings, which we define as  $\mathbf{E}_v$ , where  $v$  indicates an image:

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

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 [I\_CLS] token is not part of the image, the positional encoding for the [I\_CLS] token is set to zero, so nothing is added to it. An image representation is defined as:

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

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:

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

The output of the Transformer is defined as:

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

### 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 [T\_CLS]  $\in \mathbb{R}^{768}$ , and appended with an end-of-sequence token [T\_SEP]  $\in \mathbb{R}^{768}$ . The purpose of the [T\_CLS] token is, as with [I\_CLS], to aggregate the global information/content of the text. The [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:

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

As with the image, a positional encoding is 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:

For the Transformer blocks, we use the same structure for both image and text. As mentioned in (TODO: cite data preparation), text is tokenized into subwords using the GPT-2 byte-pair encoder also used in Data2Vec [3], [4]. Before being passed into the Transformer, a start-of-sequence token  $[T\_CLS]$  is added to the beginning of the sequence, and an end-of-sequence token  $[T\_SEP]$  is added to the end of the sequence. Then, the sequence is embedded into 768-dimensional vectors, and a positional encoding is added to the embeddings. In this thesis, we define a text sequence as follows:

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

Because we use KD in some parts, representations will be superscripted with  $s$  or  $t$ , for a student and teacher representation, respectively.

## 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] A. Baevski, W.-N. Hsu, Q. Xu, A. Babu, J. Gu, and M. Auli, “data2vec: A general framework for self-supervised learning in speech, vision and language,” *arXiv abs/2202.03555*, 2022.
- [4] A. Baevski, A. Babu, W.-N. Hsu, and M. Auli, “Efficient Self-supervised Learning with Contextualized Target Representations for Vision, Speech and Language.” 2022.