

Neural Networks for Data Science Applications

Master's Degree in Data Science

Lecture 9: Transformer (attention-based) models

Lecturer: S. Scardapane



SAPIENZA
UNIVERSITÀ DI ROMA

Designing the transformer

Moving beyond convolutional layers

The assumption of **locality** embedded in convolutional layers is not always optimal: in a text, for example, a subject can depend on an object quite far from its position. In text, audio, graphs, etc., dependencies can be **sparse**, **long-range**, and possibly **dynamic**.

For example, in *'the **cat** is on the **table**'* and *'the **cat**, which belonged to my mother, is on the **table**'*, the relation between the two words is similar, but their relative positioning is quite different.

In 2017 a new architecture, called **transformer**, was proposed for text processing and then extended to most other modalities. The original model ([Vaswani et al., 2017](#)), was an encoder-decoder model for NLP tasks. Today, similar models are widespread in computer vision, audio, biology, etc.

The core of the transformer is a new layer called **multi-head attention** (MHA). It replaces the assumption of locality with a more general notion of (soft) **sparsity** of interactions.

Transformers have better scaling laws

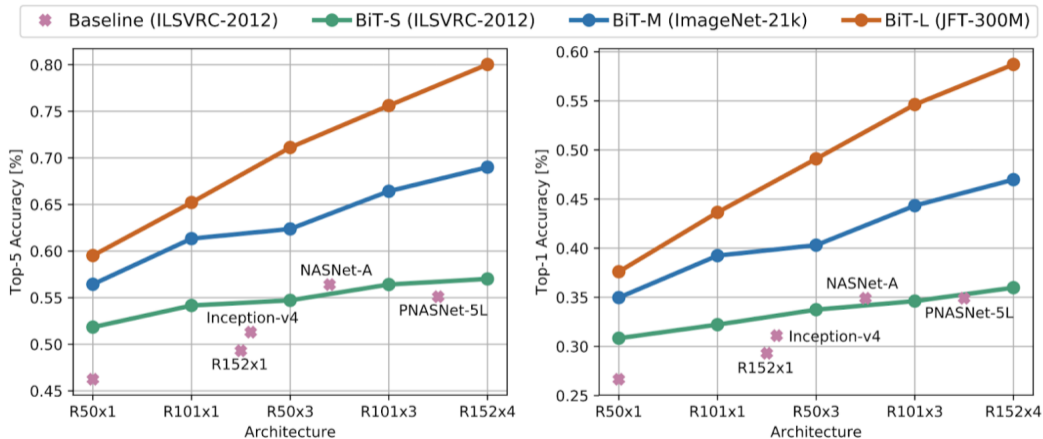


Figure 1: Open-Sourcing BiT: Exploring Large-Scale Pre-training for Computer Vision (Google AI Blog).

Designing the transformer

Self-attention

Consider a 1D sequence $\mathbf{x}_1, \dots, \mathbf{x}_n$, where $\mathbf{x}_i \in \mathbb{R}^d$. Because transformers originate from NLP, we call each element of the sequence a **token** and d the **embedding dimension**.

We can write a 1D convolutional layer (ignoring padding) of kernel size k as:

$$\mathbf{h}_i = \sum_{j=-k}^k \mathcal{W}_j \mathbf{x}_{i+j}, \quad (1)$$

where \mathcal{W} is the kernel tensor. We want to remove the assumption of locality while maintaining parameter efficiency.

Increasing k increases the number of parameters linearly. As an alternative, we can *learn* the parameters of the kernel for each possible position i via a trainable block $g(i) : \mathbb{N} \rightarrow \mathbb{R}^{d \times d}$ taking as input the shift:

$$\mathbf{h}_i = \sum_{j=-i+1}^{n-i} g(i+j)\mathbf{x}_j. \quad (2)$$

These are called **continuous convolutions**, and they work well with image-like data with, e.g., variable sizes and resolutions.

Romero, D.W. et al., 2022. Towards a General Purpose CNN for Long Range Dependencies in ND. *arXiv preprint arXiv:2206.03398*.

The previous model works well at handling non-locality, but it still assumes that dependencies are *regular*, i.e., they only depend on j . We can make it more general by letting them depend on the *values* of tokens instead:

$$\mathbf{h}_i = \sum_{j=1}^n \alpha(\mathbf{x}_i, \mathbf{x}_j) \mathbf{x}_j. \quad (3)$$

This is an example of a **non-local neural network** model. By a proper choice of the weighting function $\alpha(\cdot, \cdot)$ we can obtain the MHA layer.

Wang, X., Girshick, R., Gupta, A. and He, K., 2018. **Non-local neural networks**. In *IEEE/CVF CVPR*.

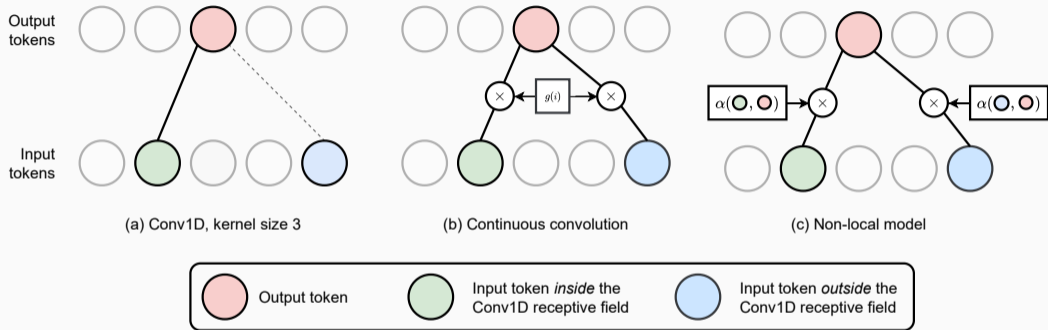


Figure 2: Example of short-term (ST, blue) and long-term (LT, green) interaction. (a) Conv1D model: ST has a trainable weight, LT is removed; (b) both connections have weights given by $g(-1)$ and $g(2)$; both connections have weight that depend on the tokens' similarities.

We make a few assumptions to simplify the layer:

- ▶ The output of α is a scalar (not a matrix). We call α the **attention scoring function** (or attention function), and its outputs the **attention scores** for token i .
- ▶ For each token, its attention scores are normalized in the simplex (they are positive and they sum to one). With this formulation, each token will have a ‘budget’ of attention to allocate, i.e., increasing an attention score necessarily decreases the attention over the remaining tokens.

There are many choices for the attention function; commonly, we use the normalized dot product $\alpha(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{\sqrt{d}} \mathbf{x}_i^\top \mathbf{x}_j$ because it is fast and efficient to parallelize.

Putting everything together we obtain (always for a single token):

$$\mathbf{h}_i = \sum_{j=1}^n \text{softmax}_j \left(\frac{1}{\sqrt{d}} \mathbf{x}_i^\top \mathbf{x}_j \right) \mathbf{x}_j. \quad (4)$$

Why the extra factor \sqrt{d} ? Suppose the elements of \mathbf{x}_i are sampled according to $\mathcal{N}(0, \sigma^2)$. The variance of $\mathbf{x}_i^\top \mathbf{x}_j$ will be σ^4 (check!), which can easily saturate the softmax with a single large (positive or negative) value.

In its current formulation, the layer lacks trainable parameters. To this end, we first reproject the input three times using three trainable matrices:

$$\mathbf{q}_i = \mathbf{W}_q^\top \mathbf{x}_i, \quad \mathbf{k}_i = \mathbf{W}_k^\top \mathbf{x}_i, \quad \mathbf{v}_i = \mathbf{W}_v^\top \mathbf{x}_i.$$

We call these the **query**, **key**, and **value** (for reasons to be explained in detail later on). The **self-attention** (SA) layer can now be written as:

$$\mathbf{h}_i = \sum_j \text{softmax} \left(\frac{1}{\sqrt{d}} \mathbf{q}_i^\top \mathbf{k}_j \right) \mathbf{v}_j.$$

Let us write the previous equation in a vectorized form, by stacking the n input vectors $\{\mathbf{x}_i\}$ into a matrix \mathbf{X} . SA can be rewritten as:

$$\underset{(n,q)}{\mathbf{Q}} = \mathbf{X} \underset{(n,d)}{\mathbf{W}}_q, \quad \underset{(n,q)}{\mathbf{K}} = \mathbf{X} \underset{(n,d)}{\mathbf{W}}_k, \quad \underset{(n,v)}{\mathbf{V}} = \mathbf{X} \underset{(n,d)}{\mathbf{W}}_v$$

$$\underset{(n,v)}{\mathbf{H}} = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{q}} \right) \mathbf{V}.$$

where the hyper-parameters are q and v . When the layer is applied to a batch of elements (e.g., sentences), it computes the attention function independently for every element of the batch (i.e., each token can *attend* only to tokens in the same sentence).

Visualizing the attention operation

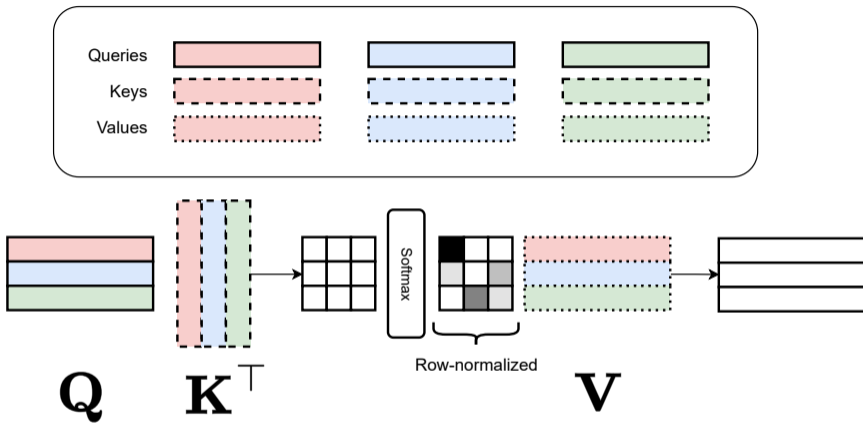


Figure 3: Visualization of the attention operation (ignoring the initial projections).

To understand the terminology, consider a Python dictionary $d = \text{dict}(\dots)$. It is a collection of key/values (k, v) pairs, such that for a given query $d[q] = v$ if k is stored inside. If the key does not exist, an error or default value is returned.

We can consider instead a 'soft' variant that always returns a value, by considering the value associated to the most similar key, even if a perfect match does not occur. If the keys, queries, and values are vectors and the distance is the dot product, this is equivalent to SA when replacing the softmax with an argmax over rows!

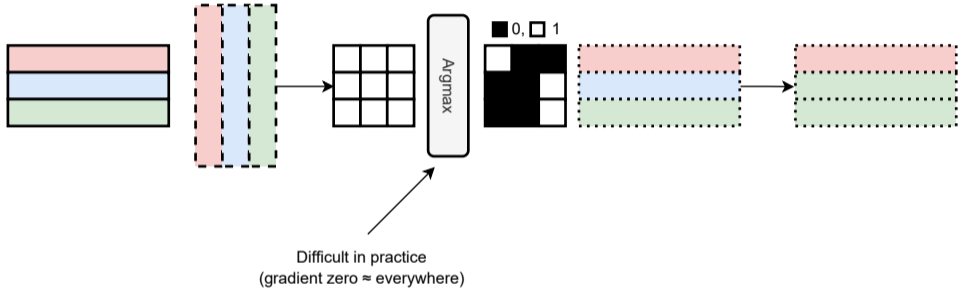


Figure 4: Hard attention is fundamentally equivalent to a dictionary with associative recall.

The SA layer can handle quasi-sparse dependencies (because of the softmax), and also dynamic ones (because of the attention function). However, what happens when the token can depend on multiple subsets of tokens?

A common generalization in this case is **multi-head attention** (MHA). It works by computing $i = 1, \dots, h$ separate sets of keys, queries, and values:

$$\mathbf{Q}_t = \mathbf{XW}_{q,t}, \quad \mathbf{K}_t = \mathbf{XW}_{k,t}, \quad \mathbf{V}_t = \mathbf{XW}_{v,t}$$

$$\mathbf{H}_t = \text{softmax} \left(\frac{\mathbf{Q}_t \mathbf{K}_t^\top}{\sqrt{q}} \right) \mathbf{V}_t.$$

We now have $3h$ trainable matrices, or a $3 \times h \times q$ tensor assuming $q = v$.

The previous operation has h separate outputs; we combine them by concatenation over the embedding dimension, and a final reprojection with a trainable output matrix \mathbf{W}_o :

$$\mathbf{H} = \left[\mathbf{H}_1 \quad \dots \quad \mathbf{H}_h \right] \mathbf{W}_o.$$

As hyperparameters, we typically choose an embedding dimension m , an output size o , and a number of heads h , and we set $q = v = m//h$ for all heads.

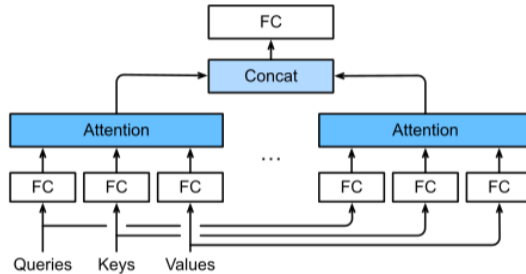


Figure 5: Visualization of the multi-head attention operation (D2L, Chapter 11.5).

Designing the transformer

The Transformer block

In transformers, the MHA layer is always used inside a more complex block, called the **transformer** block. Originally, this was composed of a MHA layer, two layer normalization operations, two residual connections, and a so-called **position-wise** network as follows:

1. Start with a MHA layer: $\mathbf{H} = \text{MHA}(\mathbf{X})$.
2. Add a residual connection and a layer normalization operation:
 $\mathbf{H} = \text{LayerNorm}(\mathbf{H} + \mathbf{X})$.
3. Apply a fully-connected model $g(\cdot)$ on each row: $\mathbf{F} = g(\mathbf{H})$.
4. Do again step 2: $\mathbf{H} = \text{LayerNorm}(\mathbf{F} + \mathbf{H})$.

Vaswani, A. et al., 2017. Attention is all you need. *NeurIPS*.

The design of the block was mostly based on empirical considerations. Roughly speaking, steps (1)-(2) correspond to a **token mixing** operation, while steps (3)-(4) are a per-token update which is akin to a 1x1 convolution. The block is similar in spirit to the depthwise separable convolution model.

The intermediate MLP is typically designed as a 2-layer MLP, with hidden dimension an integer multiple of the input dimension (e.g., 3x, 4x), and no biases.

Vaswani, A. et al., 2017. **Attention is all you need**. *NeurIPS*.

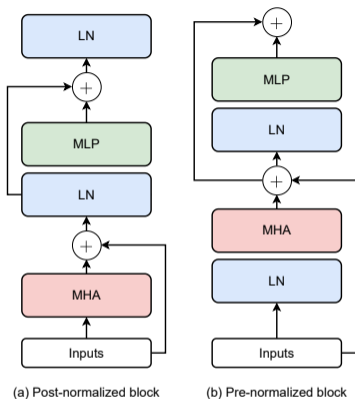


Figure 6: The original block is called **post-normalized**. A **pre-normalized** variant is also common due to it being simpler to train in most cases.¹

¹Xiong, R. et al., 2020. On layer normalization in the transformer architecture. *ICML*.

Many other variants are now common, depending on the application and computational considerations, e.g.:

- ▶ **Parallel** variants perform the MHA and MLP operations in parallel, i.e., $\mathbf{H} = \mathbf{H} + \text{MLP}(\mathbf{H}) + \text{MHA}(\mathbf{H})$. In this way, the initial and final projections of the MLP and MHA layers can be fused.²
- ▶ **Q/K normalized** variants add additional LN operations over the keys and queries (ibidem).
- ▶ **Multi-query** variants share the same keys and values over different heads to save computations.³

²Dehghani et al., 2023. *Scaling vision transformers to 22 billion parameters*. *ICML*.

³Shazeer, N., 2019. *Fast transformer decoding: One write-head is all you need*. *arXiv preprint arXiv:1911.02150*.

Designing the transformer

Positional embeddings

Consider the 3×3 matrix defined as:

$$\mathbf{P} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} .$$

It is easy to check that:

$$\mathbf{P} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} x_1 \\ x_3 \\ x_2 \end{pmatrix} .$$

These are called **permutation matrices**.

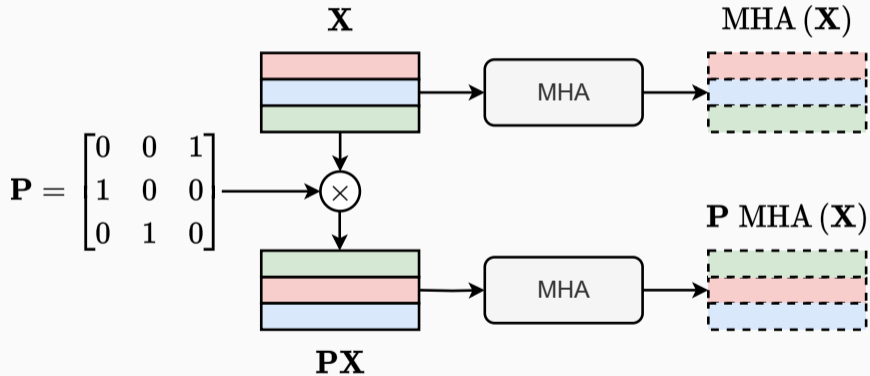
In audio and text, the i th row of \mathbf{X} represents a single time-step or a single text token (e.g., a word). In a MHA layer, their ordering is lost, because the layer is *equivariant* to the ordering (similar to the GAT layer for graphs).

If we multiply \mathbf{X} by a permutation matrix \mathbf{P} , then (the same holds trivially for the entire block):

$$\text{MHA}(\mathbf{PX}) = \mathbf{P} \cdot \text{MHA}(\mathbf{X}) .$$

This is not a good property to have for sequences.

This is easy to show, since $(\mathbf{PK})^\top = \mathbf{K}^\top \mathbf{P}^\top$ and $\text{softmax}(\mathbf{PK}^\top \mathbf{P}^\top) \mathbf{PV} = \mathbf{P} \text{softmax}(\mathbf{K}^\top) \mathbf{V}$.



Before the first MHA layer, we concatenate or sum to the input X a matrix of **positional embeddings** E :

(n,e)

$$X' = [X \parallel E] \quad \text{or} \quad X' = X + E,$$

where each row $[E]_i$ should *uniquely* encode the position of every element of the sequence.

Using this strategy, we ‘break’ the equivariance:

$$\text{MHA}(PX \parallel E) \neq P \cdot \text{MHA}(X \parallel E).$$

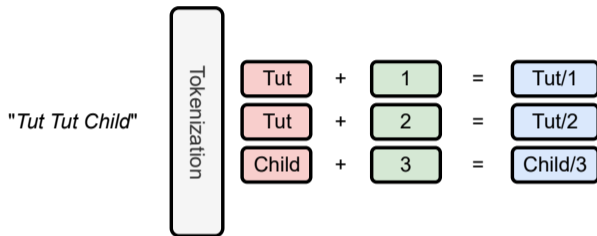


Figure 7: The same token at two different positions is now represented by different vectors after adding the positional embeddings.

We can encode the position for a sequence of maximum length p with a one-hot vector of dimension p , e.g.:

$$\mathbf{E}_0 = [1, 0, 0, \dots] , \quad \mathbf{E}_1 = [0, 1, 0, \dots] , \quad \mathbf{E}_2 = [0, 0, 1, \dots] , \quad \dots$$

Or with a single increasing scalar:

$$\mathbf{E}_0 = [0/p] , \quad \mathbf{E}_1 = [1/p] , \quad \mathbf{E}_2 = [2/p] , \quad \dots$$

Both strategies are not particularly good empirically.

We can *learn* the positional embeddings using the `tf.keras.layers.Embedding` layer:

- ▶ To each position i we associate an embedding vector of fixed dimension.
- ▶ The embeddings are trained with the rest of the network.

Note that we need to fix the maximum length of the sentence. For longer sentences, we need to linearly interpolate the set of vectors up to a larger dimension (this is the strategy used in BERT and the Vision Transformer described below).

Consider a single sinusoidal function of frequency ω :

$$\mathbf{E}_i = [\sin(i\omega)] .$$

We can interpret this as a clock with frequency ω : for two points inside a single rotation, it will give us their relative distance. For other points, the distance will be precise modulo the frequency.

To uniquely identify any possible position, we can consider multiple sinusoids, each with a frequency $\omega_j, j = 1, \dots, e$:

$$\mathbf{E}_i = [\sin(i\omega_0), \sin(i\omega_1), \dots, \sin(i\omega_e)] .$$

You can think of this as a clock with e different hands, each rotating at its own frequency. This is a nice representation because it can possibly generalize to any length, without the need to impose a maximum length *a priori*.

An empirically good choice for the frequencies (popularized by (Vaswani et al., 2017)) is:

$$\omega_j = \frac{1}{10000^{j/e}}.$$

For $j = 0$, this has frequency 2π . For $j = e$, this has frequency $10000 \cdot 2\pi$. In the middle, the frequency are increasing at a geometric progression.

To reduce the number of parameters, it is also common to *sum* the positional encodings instead of concatenating (in which case the dimension e is equal to d):

$$\mathbf{X}' = \mathbf{X} + \mathbf{E}.$$

A popular extension is to alternate sines and cosines of the same frequency:

$$[\mathbf{E}]_{i,2j} = \sin \left(\frac{i}{10000^{2j/e}} \right), \quad (5)$$

$$[\mathbf{E}]_{i,2j+1} = \cos \left(\frac{i}{10000^{2j/e}} \right). \quad (6)$$

One important property of this encoding is that it is possible to translate an encoding via matrix multiplication:

$$[\mathbf{E}]_{i+p} = [\mathbf{E}]_i \mathbf{T}(p) \quad \text{for some } \mathbf{T}(p).$$

See https://kazemnejad.com/blog/transformer_architecture_positional_encoding/ and references therein.

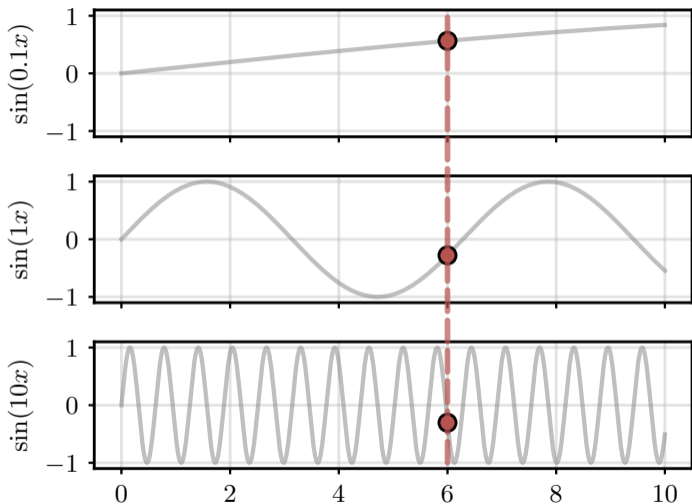


Figure 8: Visualization of the sinusoidal positional encodings (book, Chapter 10.6).

Another possibility is using **relative positional embeddings**. In this case, we modify the attention function to make it depend on the relative distance $i - j$ between tokens.

For example, **attention with linear biases**⁴ (ALiBi) adds trainable biases b_{ij} :

$$\alpha(\mathbf{x}_i, \mathbf{x}_j, i - j) = \mathbf{x}_i^\top \mathbf{x}_j + b_{ij}. \quad (7)$$

Another common option are **rotary position embeddings**⁵ (RoPE).

⁴Press, O., Smith, N.A. and Lewis, M., 2021. Train short, test long: Attention with linear biases enables input length extrapolation. *arXiv preprint arXiv:2108.12409*.

⁵Su, J. et al., 2021. Roformer: Enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*.

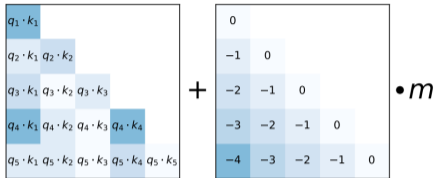


Figure 3: When computing attention scores for each head, our linearly biased attention method, ALiBi, adds a constant bias (right) to each attention score ($q_i \cdot k_j$, left). As in the unmodified attention sublayer, the softmax function is then applied to these scores, and the rest of the computation is unmodified. **m is a head-specific scalar** that is set and not learned throughout training. We show that our method for setting m values generalizes to multiple text domains, models and training compute budgets. When using ALiBi, we do *not* add positional embeddings at the bottom of the network.

Figure 9: Linear biases for attention (reproduced from Press, Smith, Lewis, 2021.).

Designing the transformer

The complete transformer model

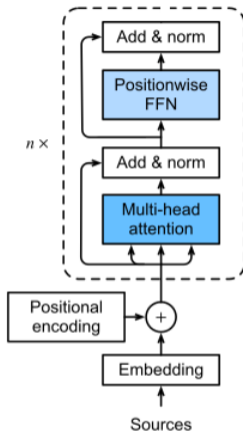


Figure 10: The final model is built with positional encodings and a stack of n transformer blocks (adapted from Chapter 10.7 of the book).

To perform classification or regression, we can apply a final global pooling on the n tokens and one or more fully-connected layers.

An alternative that is empirically found to work well is the **class token**, which is an additional *trainable* token \mathbf{c} added to the input matrix:

$$\mathbf{X}'_{(n+1,d)} = \begin{bmatrix} \mathbf{X} \\ \mathbf{c}^\top \end{bmatrix} .$$

The transformer model is applied to the matrix \mathbf{X}' as input ($\mathbf{H} = \text{Transformer}(\mathbf{X}')$), and classification is performed on its last row:

$$\mathbf{y} = \text{softmax}(\mathbf{W}^\top [\mathbf{H}]_{n+1}) .$$

Designing the transformer

Causal models and encoder-decoder models

The original transformer model was a more general model defined for **sequence to sequence** (seq2seq) tasks, such as machine translation (variable number of tokens in inputs *and* in output).

It performed an **encoding** of the input sequence, which was then **decoded** by a second, masked transformer to generate the output sequence. In order to understand it, we need to introduce two further mechanisms: **masked attention** and **cross-attention**.

For this reason, the model described before is sometimes called an **encoder-only** Transformer.

Vaswani, A., et al., 2017. **Attention is all you need**. In *Advances in neural information processing systems* (pp. 5998-6008).

In order to build a causal transformer variant, we can replace the SA layer with a masked variant:

$$\mathbf{H} = \text{softmax} \left(\frac{\mathbf{QK}^T \odot \mathbf{M}}{\sqrt{q}} \right) \mathbf{V}.$$

where \mathbf{M} has a triangular structure:

$$M_{ij} = \begin{cases} 1 & \text{if } j \leq i \\ -\infty & \text{otherwise} \end{cases}. \quad (8)$$

In practice we can use very small numbers, e.g., 10^{-10} .

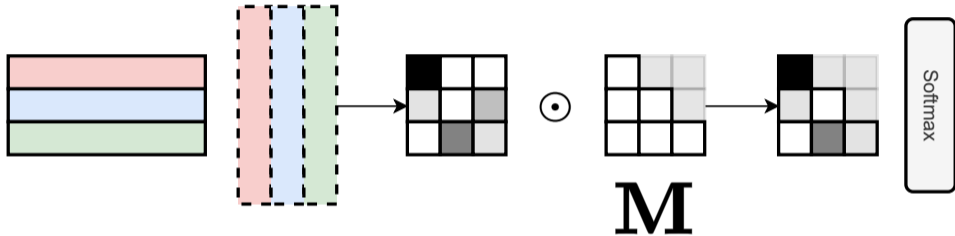


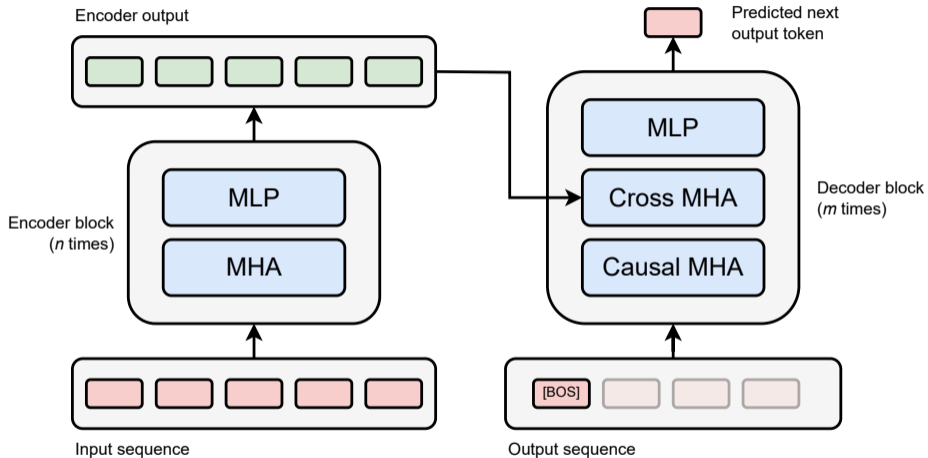
Figure 11: Note that masking with 0 is invalid because $\exp(0) = 1$, and masking after the softmax is invalid because of its denominator.

Given two sets \mathbf{X} and \mathbf{Z} , **cross attention** is defined as:

$$\text{CA}(\mathbf{X}, \mathbf{Z}) = \text{softmax} \left((\mathbf{Z}\mathbf{W}_q) (\mathbf{X}\mathbf{W}_k)^\top \right) \mathbf{X}\mathbf{W}_v. \quad (9)$$

This is a useful operation that can combine information coming from different streams of information (e.g., audio-visual datasets).

Full encoder-decoder architecture



Designing the transformer

Computational considerations

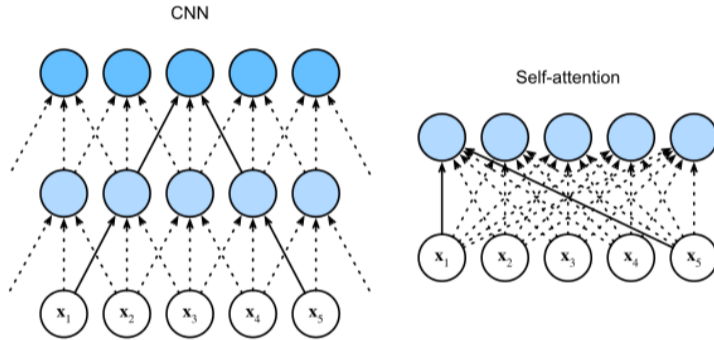


Figure 12: Adapted from Chapter 10.6 of Dive into Deep Learning.

Consider a 1D convolutional operation $\mathbf{H} = \text{Conv1D}(\mathbf{X})$ with a filter size of k . Computing the output requires $\mathcal{O}(nkd^2)$ operations.

Self-attention (with one head) requires $\mathcal{O}(nd^2)$ time for computing keys, queries, and values, and $\mathcal{O}(n^2d)$ time for computing the output. The n^2 term limited the application to large sequences, although very efficient implementations and hardware are available nowadays (e.g., FlashAttention⁶).

A single layer of MHA has a receptive field of n , while the convolutive layer has a receptive field of k .

⁶<https://github.com/Dao-AILab/flash-attention>

Materializing the \mathbf{QK}^T matrix also requires $\mathcal{O}(n^2)$ memory. Modern implementations (e.g., **FlashAttention**) avoid this by processing tokens in multiple chunks (see book). This can be done by storing intermediate values on the denominator of the softmax, and only applying the normalization at the end (**lazy softmax**).

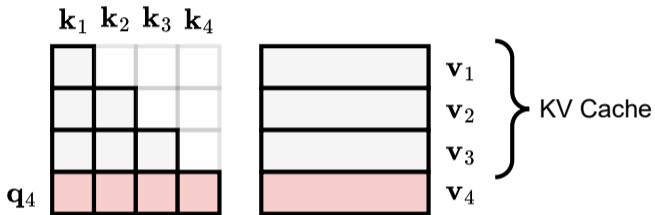


Figure 13: The **KV Cache** is an essential component of autoregressive implementations (see book).

Practical transformer models

Text Transformers

The majority of pre-trained word embedding models are standard transformer models trained on the sequence of text tokens.

- ▶ **BERT-like** models are pre-trained by masking one word in a sentence, and reconstructing the full sentence in output.
- ▶ **GPT-like** models are (causal) variants pre-trained to generate the sequence auto-regressively.

These models are called **contextual** embeddings because the same word in different sentences can be encoded to different vectorial representations.

Qiu, X., et al., 2020. Pre-trained models for natural language processing: A survey. *Science China Technological Sciences*, pp. 1-26.

Because these models are trained from the raw text alone (no specific targets) they are called **self-supervised** models (we will cover this more in-depth later).

Their strength is that scaling laws for transformers are empirically better than for other models (i.e., they benefit more from increasing the dataset by order of magnitude).

In natural language processing, this is also shown by the emergence of paradigms like **text-prompting** and **zero-shot learning**.

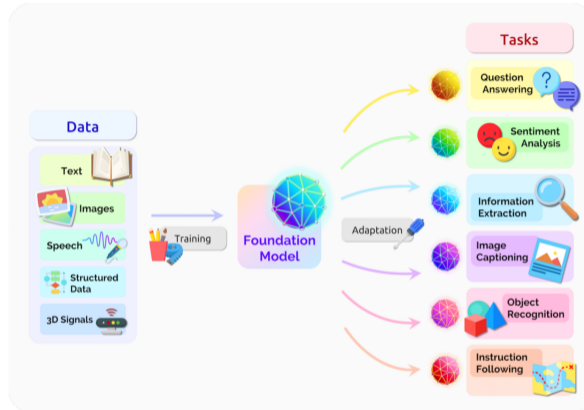


Fig. 2. A foundation model can centralize the information from all the data from various modalities. This one model can then be adapted to a wide range of downstream tasks.

Figure 14: An emerging name for these huge, pre-trained models is **foundation models**.

Practical transformer models

Vision, Audio, & Graph Transformers

One important realization of the last two years is that transformers can also benefit computer vision, especially when trained on huge datasets (e.g., ImageNet21k).

However, this requires to convert the original image (a 2D grid) into a 1D sequence (actually, a set together with the positional embeddings). Because this would scale quadratically in the number of pixels, a common solution is to work on **patches** of the original image.

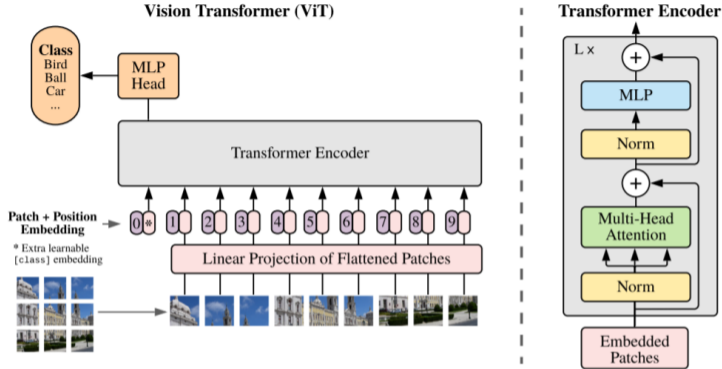


Figure 15: The Vision Transformer (ViT) is a standard transformer applied on top of image patches.

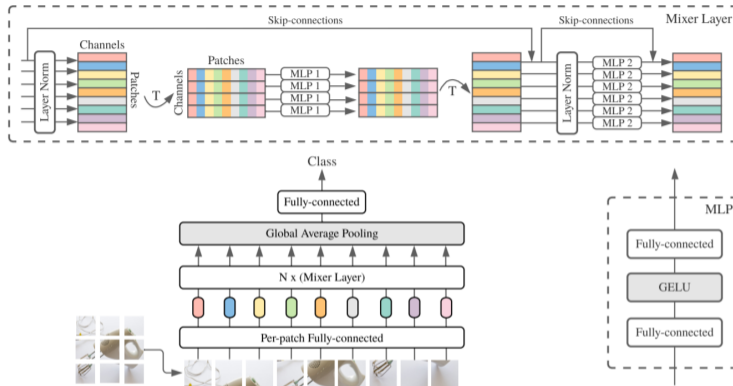


Figure 16: Mixer models are variants of the ViT, where the MHA is replaced by fully-connected layers.

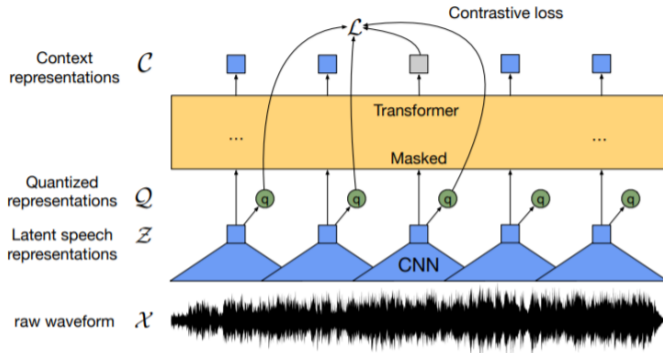


Figure 17: Architectures like **Wav2Vec 2.0** are pre-trained audio models exploiting transformers. However, this is harder because of the nature of the audio signal.

Baevski, A., Zhou, H., Mohamed, A. and Auli, M., 2020. **wav2vec 2.0: A framework for self-supervised learning of speech representations.** *arXiv preprint arXiv:2006.11477*.

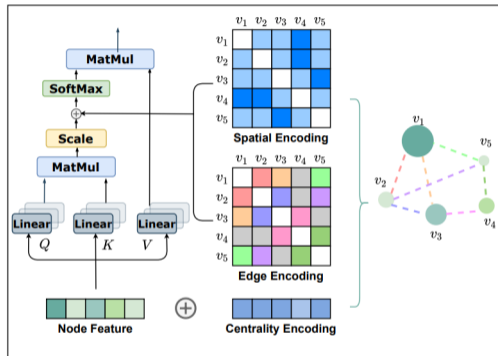


Figure 1: An illustration of our proposed centrality encoding, spatial encoding, and edge encoding in Graphormer.

Figure 18: We can also design **graph transformers**, where nodes become tokens and the connectivity is embedded inside the positional embeddings.

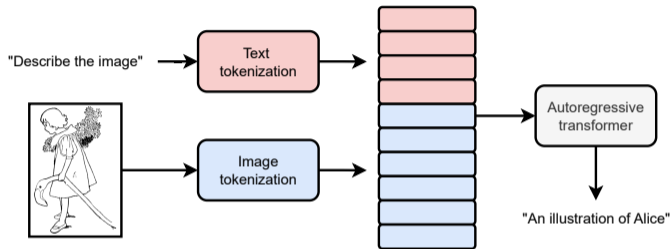


Figure 19: Transformers also simplify **multimodal** models by projecting tokens of different modalities in a single embedding space.

Ying, C., et al., 2021. **Do transformers really perform badly for graph representation?**. *Advances in Neural Information Processing Systems*, 34, pp. 28877-28888.

- ▶ **Chapters 10 and 11** from the book.
- ▶ <https://jalammr.github.io/illustrated-transformer/>.
- ▶ <https://srush.github.io/raspy/> for intuitions about how transformers can work.