

13.3 Natural language processing (NLP): Pretrained language models

13.3.1 Transformer

13.3.1.1 Overall structure

Transformer [2] is one of most successful architectures in tackling challenging seq2seq NLP tasks like translation, question-answering, etc.

Traditionally, seq2seq tasks heavily use RNN based encoder-decoder architectures, plus attention mechanisms, to transform one sequence into another one. Transformer, on the other hand, does not rely on any recurrent structure and is able to process all tokens in a sequence at the same time.

On a high level, transformer resembles the encoder-decoder architecture [Figure 13.3.1], where the encoder converts an input sequences into low-dimensional embeddings and the decoder outputs an output sequence probabilities. The decoder takes input sequence embeddings and previous output token as the input. The sequential information, original stored in recurrent network structure, is now stored in position encoding added at the entry point of the encoder and decoder modules.

The encoder module on the left [Figure 13.3.1] consists of blocks that are stacked on top of each other to obtain the embeddings that retain rich information in the input. Multi-head self-attentions are used in each block to enable the extraction of contextual information into the final embedding. Similarly, the decoder module on the right also consists of blocks that are stacked on top of each other to obtain the embeddings. Two types of multi-head attentions are used in each block, one is self-attention to capture contextual information among output sequence and one is to encoder-decoder attention to capture the dynamic information between input and output sequence.

In the following, we will present detailed analysis on each component of transformer architecture.

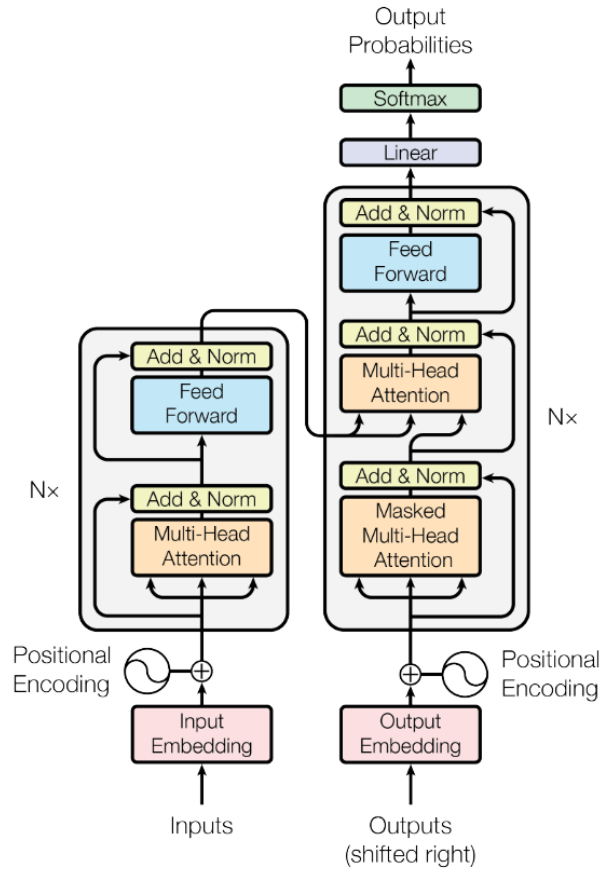


Figure 13.3.1: The transformer architecture [2].

13.3.1.2 Input output conventions

In the transformer architecture, we need to deal with different types of sequences

- Input sequence $s = (i_1, i_2, \dots, i_p, \dots, i_n), i_p \in \mathbb{N}$ and input position sequence $s^p = (1, 2, \dots, n)$.
- Output sequence $o = (o_1, o_2, \dots, o_p, \dots, o_m), o_p \in \mathbb{N}$ and output position sequence $o^p = (1, 2, \dots, m)$. Output sequence is the input to the decoder.
- Target sequence $t = (t_1, t_2, \dots, t_p, \dots, t_m), t_p \in \mathbb{N}$. Input sequence and target sequence form a pair in the training examples.

For example, consider a translational task with input and target sequence given by 'Ich möchte eine Flasche Wasser' and 'I want a bottle of water'. In the training, we can get a typical input sequence, target sequence, and output sequence in the following form

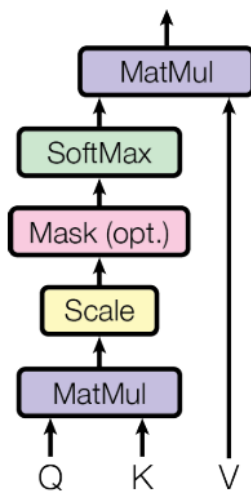
- $s = (\text{Ich}, \text{mchte}, \text{eine}, \text{Flasche}, \text{Wasser}, \text{PAD}, \text{PAD})$
- $t = (\text{I}, \text{want}, \text{a}, \text{bottle}, \text{of}, \text{water}, \text{EOS}, \text{PAD})$

- $o = (SOS, I, want, a, bottle, of, water, EOS)$

The output sequence is the right-shifted target sequence with a starting token. Because sequence data is fed into transformer via mini-batches, all input sequence in the same mini-batch will be pad to the maximum length of its batch.

13.3.1.3 Multihead attention with marks

Scaled Dot-Product Attention



Multi-Head Attention

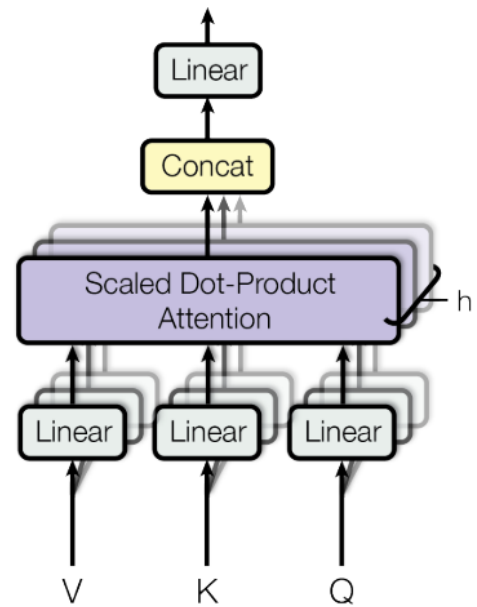


Figure 13.3.2: The multi-head attention architecture [2].

Multi-head attention mechanism plays important role in both encoder and decoder side; particularly, they are used in three places:

- In the encoder module, Multi-head attention is used to construct embedding for each token that depends of its context (i.e., self-attention of the input sequence).
- In the decoder module, Multi-head attention is used to construct embedding for each token that depends of its *seen* output context (i.e., masked self-attention of the output sequence).
- From the encoder module to the decoder module, Multi-head attention is used to construct embedding for each output token that depends of its *input* context (i.e., attention between input and output sequences).

Given a query matrix $Q \in \mathbb{R}^{n \times d_{\text{model}}}$, a key matrix $K \in \mathbb{R}^{m \times d_{\text{model}}}$, and a value matrix $V \in \mathbb{R}^{m \times d_{\text{model}}}$, the multi-head (h heads) attention associated with (Q, K, V) is given by

$$\text{MultiHeadAttention}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

where $\text{head}_i \in \mathbb{R}^{n \times d_v}$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

with $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$, $W^O \in \mathbb{R}^{h \times d_v \times d_{\text{model}}}$. In general, we require $d_k = d_v = d_{\text{model}}/h$ such that the output of $\text{MultiHeadAttention}(Q, K, V)$ has the dimensionality of $n \times d_{\text{model}}$.

The attention among (Q, K, V) is given by

$$\text{Attention}(QW_i^Q, KW_i^K, VW_i^V) = \text{softmax}\left(\frac{QW_i^Q (KW_i^K)^T}{\sqrt{d_k}}\right) VW_i^V.$$

Note that the softmax applies to each row such that the $\text{softmax}(\cdot)$ produces a weight matrix $w^{\text{att}} \in \mathbb{R}^{n \times m}$, with each row summing up to unit 1.

Explicitly, we have

$$\begin{bmatrix} w_{11}^{\text{att}} & w_{12}^{\text{att}} & \dots & w_{1m}^{\text{att}} \\ w_{21}^{\text{att}} & w_{22}^{\text{att}} & \dots & w_{2m}^{\text{att}} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n1}^{\text{att}} & w_{n2}^{\text{att}} & \dots & w_{nm}^{\text{att}} \end{bmatrix} \begin{bmatrix} [VW^V]_1 \\ [VW^V]_2 \\ \vdots \\ [VW^V]_m \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^m w_{1j}^{\text{att}} [VW^V]_j \\ \sum_{j=1}^m w_{2j}^{\text{att}} [VW^V]_j \\ \vdots \\ \sum_{j=1}^m w_{nj}^{\text{att}} [VW^V]_j \end{bmatrix}$$

where $[VW^V]_i$ is the i th row vector of value matrix VW^V .

Sometimes, for each query associated with a symbol, we like to only allow a subset of keys and values to be queried via a binary **mask** $\text{mask} \in \{0, 1\}^m$, where 1 indicates exclusion of the key. We can set the values in the intermediate matrix $QW_i^Q (KW_i^K)^T$ to be negative infinity if this column index is associated with mask value 1.

13.3.1.4 Encoder anatomy

Given an input sequence represented by integer sequence $s = (i_1, \dots, i_p, \dots, i_n)$, $i_p \in \mathbb{N}$, e.g., $s = (3, 5, 30, 2, \dots, 21)$, the goal of the encoder module [Figure 13.3.1] is to map s to n dense vectors that captures semantic, position, and contextual information. For the input

sequence s , we first project s into a dense word embedding space via $WE(s) \in \mathbb{R}^{n \times d_{model}}$. The word embedding $WE(s)$ does not encode positional information in the sequence, so we utilize a position encoding PE maps an integer index representing the position of the token in the sequence, $s^p = (1, 2, \dots, n)$ to n dense vector with same dimension d_{model} , i.e., $PE(s^p) \in \mathbb{R}^{n \times d_{model}}$. Notably, given the token position $i \in \{1, \dots, n\}$ $PE(i) \in \mathbb{R}^{d_{model}}$ is given by

$$[PE(i)]_j = \begin{cases} \sin\left(\frac{i}{10000^{2j/d_{model}}}\right) & \text{if } j \text{ is even} \\ \cos\left(\frac{i}{10000^{2j-1/d_{model}}}\right) & \text{if } j \text{ is odd} \end{cases}$$

Intuitively, each dimension of the positional encoding corresponds to a sinusoid of different wavelengths ranging from 2π to $10000 \cdot 2\pi$.

The whole computation in the encoder module can be summarized in the following.

Definition 13.3.1 (computation in encoder module). *Given an input sequence represented by integer sequence $s = (i_1, \dots, i_p, \dots, i_n)$ and its position $s^p = (1, \dots, p, \dots, n)$. The encoder module takes s, s^p as inputs and produce $e_N \in \mathbb{R}^{n \times d_{model}}$.*

$$\begin{aligned} e_0 &= WE(s) + PE(s^p) \\ e_1 &= \text{EncoderLayer}(e_0) \\ e_2 &= \text{EncoderLayer}(e_1) \\ &\dots \\ e_N &= \text{EncoderLayer}(e_{N-1}) \end{aligned}$$

where $e_i \in \mathbb{R}^{n \times d_{model}}$, $\text{EncoderLayer} : \mathbb{R}^{n \times d_{model}} \rightarrow \mathbb{R}^{n \times d_{model}}$ is an encoder sub-unit, N is the number of encoder sub-units. Specifically, this encoder sub-unit can be decomposed into following calculation procedures

$$\begin{aligned} e_{mid} &= \text{LayerNorm}(e_{in} + \text{MultiHeadAttention}(e_{in}, e_{in}, e_{in}, \text{padMask})) \\ e_{out} &= \text{LayerNorm}(e_{mid} + \text{FFN}(e_{mid})) \end{aligned}$$

where $e_{mid}, e_{out} \in \mathbb{R}^{n \times d_{model}}$,

$$\text{FFN}(e_{mid}) = \max(0, e_{mid}W_1 + b_1)W_2 + b_2,$$

with $W_1 \in \mathbb{R}^{d_{model} \times d_{ff}}$, $W_2 \in \mathbb{R}^{d_{ff} \times d_{model}}$, $b_1 \in \mathbb{R}^{d_{ff}}$, $b_2 \in \mathbb{R}^{d_{model}}$, and the padMask excludes padding symbols in s .

13.3.1.5 Decoder anatomy

In the decoder side, we are similarly given an output sequence represented by integer sequence $o = (o_1, \dots, o_p, \dots, o_m), o_p \in \mathbb{N}$, e.g., $o = (5, 10, 30, 2, \dots, 21)$. The decoder module [Figure 13.3.1] aims to convert an output sequence o , combining with resulting embedding e_N in the encoder module [Definition 13.3.1] to its corresponding probabilities over the vocabulary. The output probability can be used to compute categorical loss that drives the learning process of the encoder and the decoder.

The whole computation in the encoder module can be summarized in the following.

Definition 13.3.2 (computation in decoder module). Given an input sequence represented by integer sequence $o = (o_1, \dots, o_p, \dots, o_m)$ and its position $o^p = (1, \dots, p, \dots, m)$. The encoder module takes o, o^p as inputs, combines resulting embedding e_N from the encoder, and produce $d_N \in \mathbb{R}^{n \times d_{\text{model}}}$ and probabilities over the vocabulary.

$$\begin{aligned} d_0 &= \text{WE}(o) + \text{PE}(o^p) \\ d_1 &= \text{DecoderLayer}(d_0, e_N) \\ d_2 &= \text{DecoderLayer}(d_1, e_N) \\ &\dots \\ d_N &= \text{DecoderLayer}(d_{N-1}, e_N) \\ \text{output prob} &= \text{Softmax}(d_N W) \end{aligned}$$

where $d_0 \in \mathbb{R}^{m \times d_{\text{model}}}$, $\text{DecoderLayer} : \mathbb{R}^{m \times d_{\text{model}}} \rightarrow \mathbb{R}^{m \times d_{\text{model}}}$ is a decoder sub-unit, N is the number of decoder sub-units, $W \in \mathbb{R}^{d_{\text{model}} \times |V|}$. Specifically, this decoder sub-unit can be decomposed into following calculation procedures

$$\begin{aligned} d_{\text{mid1}} &= \text{LayerNorm}(d_{\text{in}} + \text{MaskedMultiHeadAttention}(d_{\text{in}}, d_{\text{in}}, d_{\text{in}}, \text{padMask})) \\ d_{\text{mid2}} &= \text{LayerNorm}(d_{\text{mid1}} + \text{MaskedMultiHeadAttention}(d_{\text{mid1}}, e_N, e_N, \text{padMask}|\text{seqMask})) \\ d_{\text{out}} &= \text{LayerNorm}(d_{\text{mid2}} + \text{FFN}(d_{\text{mid2}})) \end{aligned}$$

where $d_{\text{mid1}}, d_{\text{mid2}}, d_{\text{out}} \in \mathbb{R}^{m \times d_{\text{model}}}$,

$$\text{FFN}(d_{\text{out}}) = \max(0, d_{\text{mid}} W_1 + b_1) W_2 + b_2,$$

with $W_1 \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ff}}}$, $W_2 \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{model}}}$, $b_1 \in \mathbb{R}^{d_{\text{ff}}}$, $b_2 \in \mathbb{R}^{d_{\text{model}}}$.

Note that there are two type of attention in the decoder module, one is self-attention among output sequence and one is attention between encoder output and decoder output.

The encoder-decoder attention uses a mask that excludes padding symbol in the input sequence; the self-attention uses a mask that excludes padding symbol and future symbol, which can be computed via logical OR between *padMask* and *seqMask*. *seqMask* for a symbol at position i is a binary vector $seqMask_i \in \mathbb{R}^m$ whose value is 1 at position equal or greater than i . Therefore, for a sequence of length m , the complete *seqMask* would be a upper triangle matrix value 1.