

Transformer: Neural machine translation without recurrence

<https://eduassistpro.github.io/>

- ▶ It is possible to replace recurrence with **self-attention** within the encoder and decoder, as in the transformer architecture

Assignment Project Exam Help  
Assign Add WeChat edu\_assist\_pro

<https://eduassistpro.github.io/>

For each token  $m$  at level  $i$ , we compute attention over the entire source sentence. The keys, values, and query are all projections of the vector  $\mathbf{h}^{(i-1)}$ .

- ▶ The attention scores  $\alpha_{m \rightarrow n}^{(i)}$  are computed using a scaled form of softmax attention,

$$\alpha_{m \rightarrow n} \propto \exp(\psi_\alpha(m, n)/M)$$

“Self-attention”

<https://eduassistpro.github.io/>

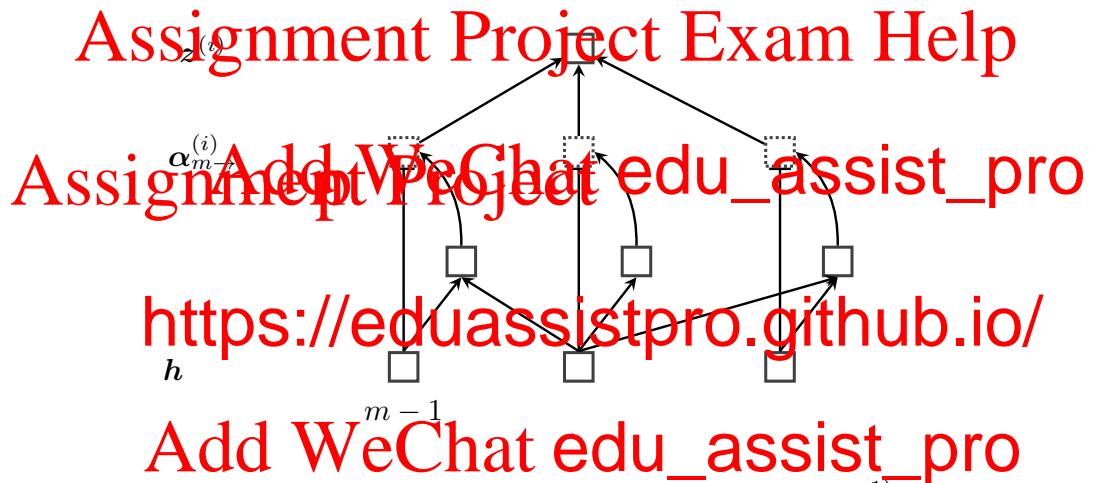


Figure 18.7: The transformer encoder’s computation of  $\psi_\alpha^{(i)}(m, m - 1)$ . The key, value, and query are shown for token  $m - 1$ . For example,  $\psi_\alpha^{(i)}(m, m - 1)$  is computed from the key  $\Theta_k h_{m-1}^{(i-1)}$  and the query  $\Theta_q h_m^{(i-1)}$ , and the gate  $\alpha_{m \rightarrow m-1}^{(i)}$  operates on the value  $\Theta_v h_{m-1}^{(i-1)}$ . The figure shows a minimal version of the architecture, with a single attention head. With multiple heads, it is possible to attend to different properties of multiple words.

“Self-attention”

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](#)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

“Self-attention”

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](#)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

“Self-attention”

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](#)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

“Self-attention”

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](#)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

“Self-attention”

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

“Self-attention”

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶ Attention is an effective mechanism to model dependencies.
- ▶ By varying the keys, values, and queries, we get variants of attention. We get self-attention value and query sentence. <https://eduassistpro.github.io/>
- ▶ Intuitively self-attention contextual sentence (or any sequence of word token) interpreting the word token, which will in turns helps other tasks that build on the understanding of words.

Multiple Attention “heads”

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

Multiple Attention “heads”

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](#)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

Multiple Attention “heads”

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

Multiple Attention “heads”

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

What does self-attention do?

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶ Contextualization <https://eduassistpro.github.io/> Add WeChat edu\_assist\_pro
- ▶ why is contextualization a good thing?
  - ▶ “Work o <https://eduassistpro.github.io/>
  - ▶ “Heat th
- ▶ Having the same embedding for both inst <https://eduassistpro.github.io/> Add WeChat edu\_assist\_pro doesn't make sense

The Transformer architecture in (almost) its entirety  
<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

Figure 1: The Transformer - model architecture.

## Position Embedding

The original Transformer uses position embedding.

- For word  $w$  at position  $pos \in [0, L-1]$  in the sequence  $\mathbf{w} = \{w_0, \dots, w_{L-1}\}$ , with a 4-dimensional word embedding  $e_w$ , and  $d_{model} = 4$ , the operation would be:

$$e_{w'} = e_w + \left[ \begin{array}{cccc} \sin\left(\frac{pos}{10000^{2/4}}\right) & \cos\left(\frac{pos}{10000^{2/4}}\right) \\ \vdots & \vdots \\ \sin\left(\frac{pos}{10000^{2/4}}\right) & \cos\left(\frac{pos}{10000^{2/4}}\right) \end{array} \right]$$

- The formula to calculate the position embedding is:

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

for  $i \in [0, d_{model}/2 - 1]$

## Position-wise feedforward neural network

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- ▶ A feed-forward network is applied to each position separately and identically, i.e., with shared weights across the sequence

<https://eduassistpro.github.io/>

### Add WeChat

- ▶ This is equivalent to applying two one-dimensional convolutional network with one kernel
- ▶ Question to think about: why not just use a one-dimensional convolution?

## Layer normalization

- ▶ Given a minibatch  $\{x_1, \dots, x_n\}$ , after applying layer normalization, we get a transformed minibatch  $B' = \{y_1, y_2, \dots, y_m\}$  where  $y_i = LN_{\gamma, \beta}(x_i)$
- ▶ Specifically, layer normalization needs the following:
  - ▶ Compute the mean and variance of each sample  $x_i$ :

<https://eduassistpro.github.io/>

- ▶ Normalize each sample such that the samples have zero mean and unit variance:

$$\hat{x}_{i,k} = \frac{x_{i,k} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

- ▶ Finally scaling and shifting with  $\gamma$  and  $\beta$ :

$$y_i = \gamma \odot \hat{x}_i + \beta \triangleq LN_{\gamma, \beta}(x_i)$$

## Reducing the output space with BPE

<https://eduassistpro.github.io/>

BPE is a simple data compression technique

- ▶ Initialize the vocabulary with the character vocabulary, and segment each words into a sequence of characters and end-of-word symbol.
- ▶ Iteratively co frequent pairs produces a new symbol which represents . Frequent character ngrams are eventually merged into a single symbol.
- ▶ The final vocabulary size equals to the initial vocabulary size plus the number of merge operations
- ▶ The number of merge operations is a hyperparameter

Python implementation

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

Python implementation

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](#)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

## Tutorials

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- ▶ NMT with Tensorflow <https://github.com/tensorflow/nmt>
- ▶ Attention visualization:  
<https://jamal-mathieu.github.io/machine-n-attention/>
- ▶ Illustrated transformer:  
<http://jalammar.github.io/illustrations/>

BERT: Bidirectional Encoder Representations from  
Transformers <https://eduassistpro.github.io/>

## Assignment Project Exam Help

Input Representations:

- ▶ The input can be both a single sentence or a pair of sentences.
- ▶ “Sentence” here just means a sequence of words.
- ▶ Sentences are packed together into one sequence (cab size = 30K).
- ▶ The first token of the input sequence is a special token [CLS].
- ▶ Sentence pairs are packed together into one sequence separated with a special token [SEP].
- ▶ Separate token embeddings for the first sentence and the second sentence in a sentence pair.

The pretrain-fine tune paradigm

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

## Pre-Training BERT

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶ The Transformer is trained with a Masked Language Model (MLM) task and a next sentence prediction (NSP) task.
- ▶ The MLM task randomly masks 15% of the tokens that are randomly masked.
- ▶ The final hidden vectors corresponding to the masked tokens are fed into a softmax over the token vocabulary.

BERT Pretraining with MLM

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](#)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

BERT next sentence prediction

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

BERT: fine tune to specific tasks

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)