

## Transformer: Neural machine translation without recurrence

<https://powcoder.com>

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

$$\mathbf{z}_m^{(i)} = \sum_{n=1}^{M^{(s)}} \alpha_{m \rightarrow n}^{(i)} (\mathbf{Q}_v \mathbf{h}_n^{(i-1)})$$

$$\mathbf{h}_n^{(i)} = \text{ReLU}(\mathbf{Q}_1 \mathbf{z}_m^{(i)} + \mathbf{b}_1) + \mathbf{b}_2$$

For each token  $m$  at level  $i$ , we compute self-attention over the entire source sentence. The keys, values, and queries 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://powcoder.com>

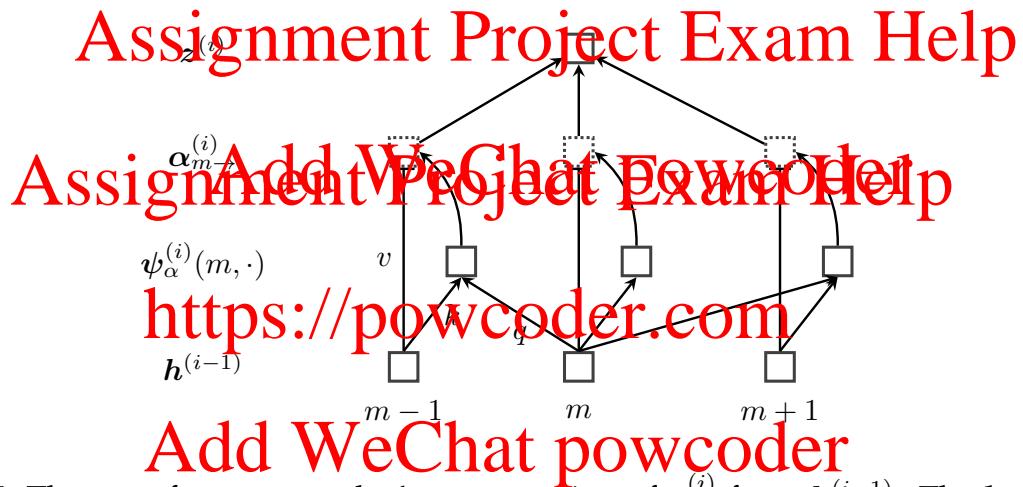


Figure 18.7: The transformer encoder’s computation of  $z_m^{(i)}$  from  $\mathbf{h}^{(i-1)}$ . The key, value, and query are shown for token  $m - 1$ . For example,  $\psi_\alpha^{(i)}(m, m - 1)$  is computed from the key  $\Theta_k \mathbf{h}_{m-1}^{(i-1)}$  and the query  $\Theta_q \mathbf{h}_m^{(i-1)}$ , and the gate  $\alpha_{m \rightarrow m-1}^{(i)}$  operates on the value  $\Theta_v \mathbf{h}_{m-1}^{(i-1)}$ . 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”

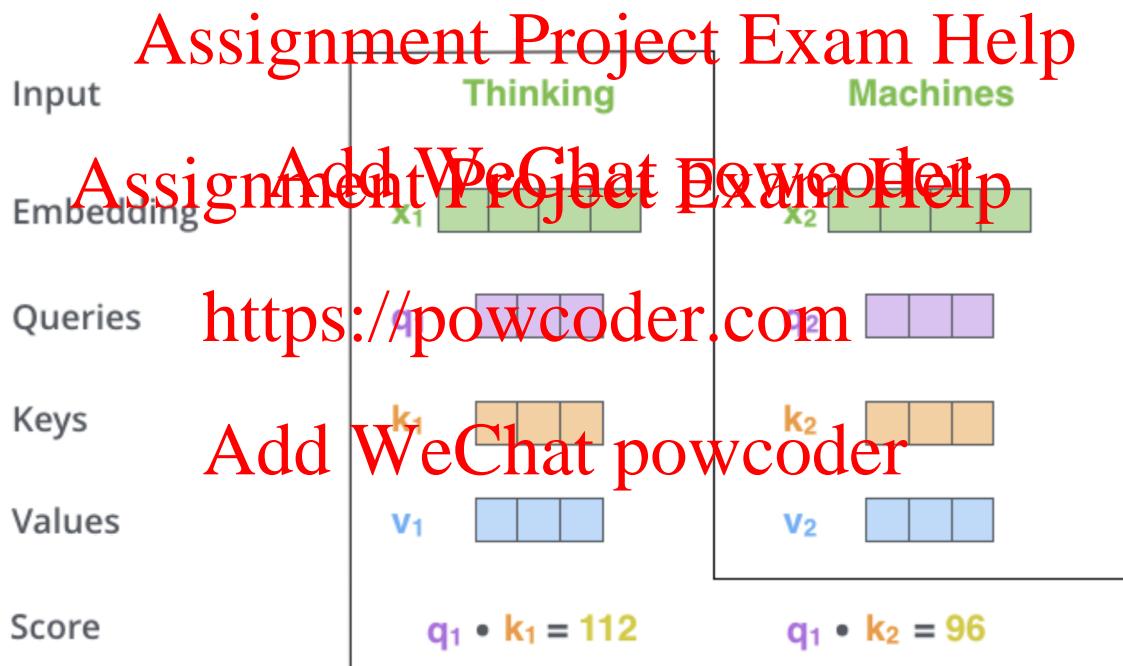
<https://powcoder.com>



Multiplying  $x_1$  by the  $W_Q$  weight matrix produces  $q_1$ , the "query" vector associated with that word. We end up creating a "query", a "key", and a "value" projection of each word in the input sentence.

“Self-attention”

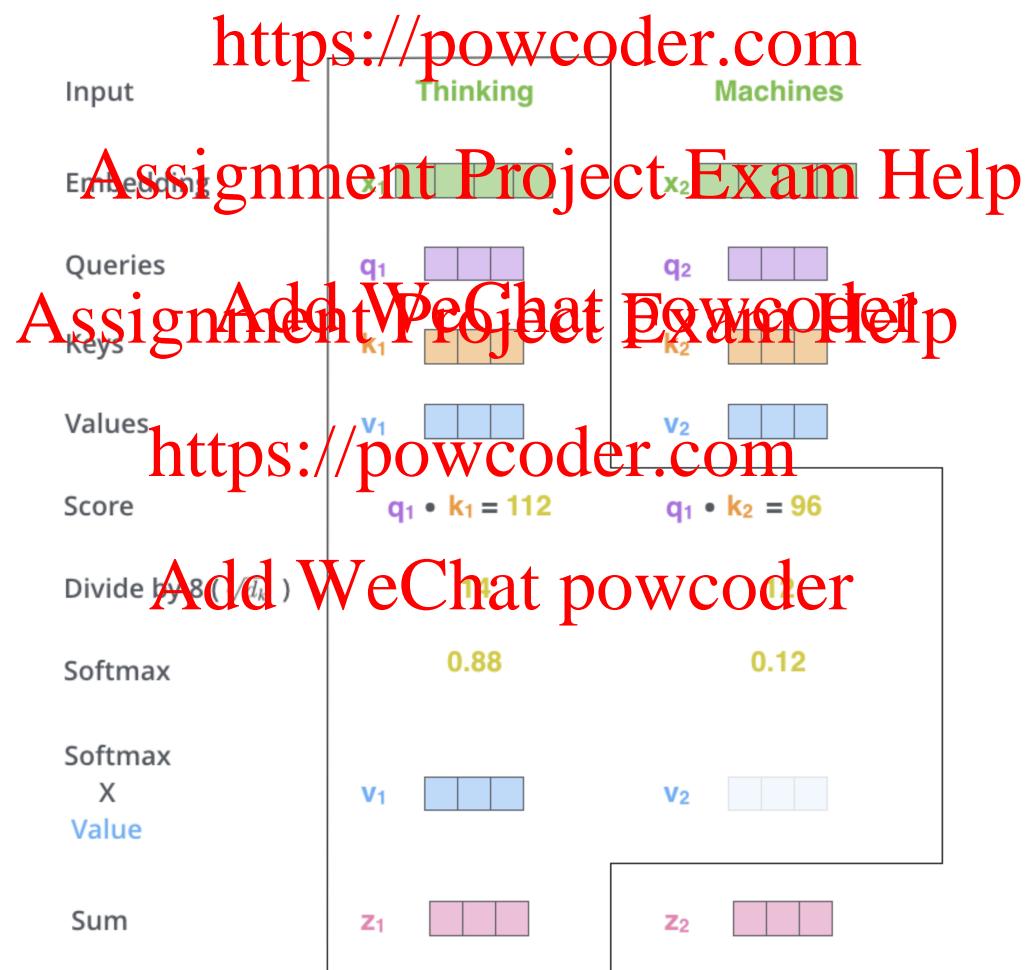
<https://powcoder.com>



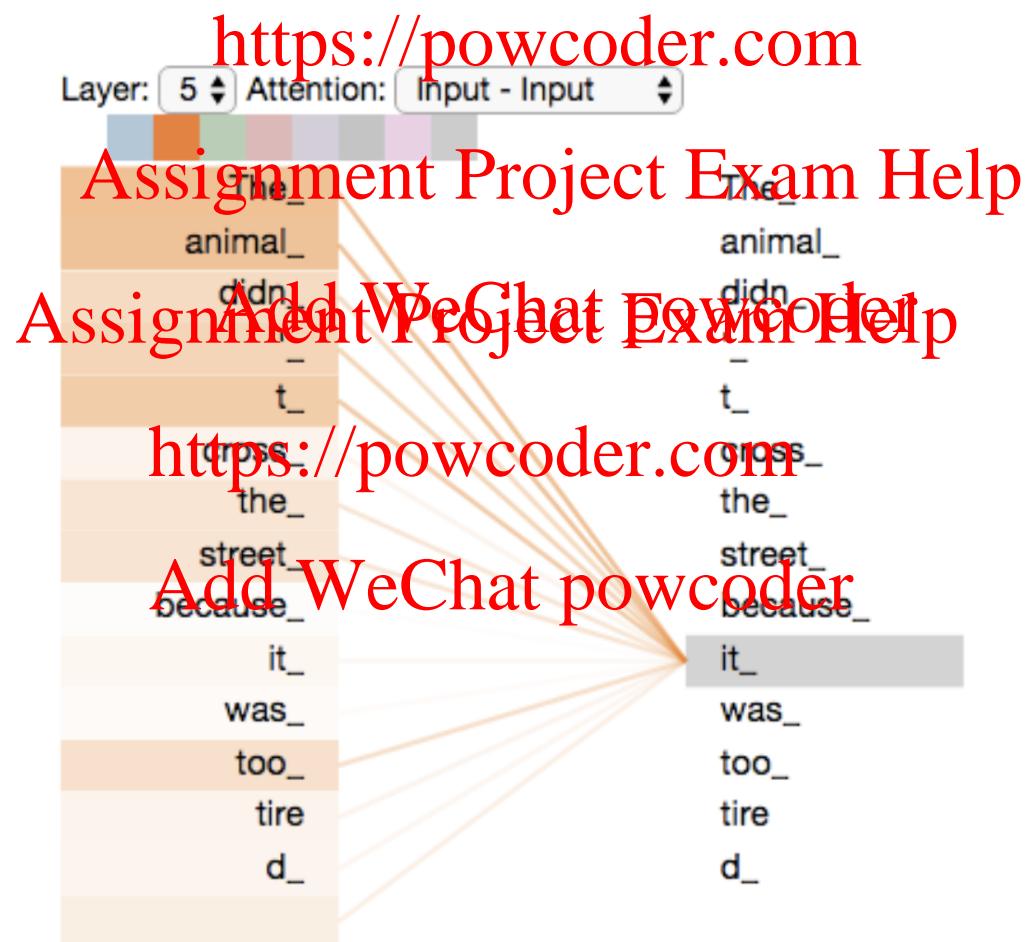
## “Self-attention”

|                              |  | <a href="https://powcoder.com">https://powcoder.com</a> |                      |
|------------------------------|--|---------------------------------------------------------|----------------------|
| Input                        |  | Thinking                                                | Machines             |
| Embedding                    |  | $x_1$                                                   | $x_2$                |
| Queries                      |  | $q_1$                                                   | $q_2$                |
| Keys                         |  | $k_1$                                                   | $k_2$                |
| Values                       |  | $v_1$                                                   | $v_2$                |
| Score                        |  | $q_1 \cdot k_1 = 112$                                   | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) |  | 14                                                      | 12                   |
| Softmax                      |  | 0.88                                                    | 0.12                 |

## “Self-attention”



“Self-attention”



“Self-attention”

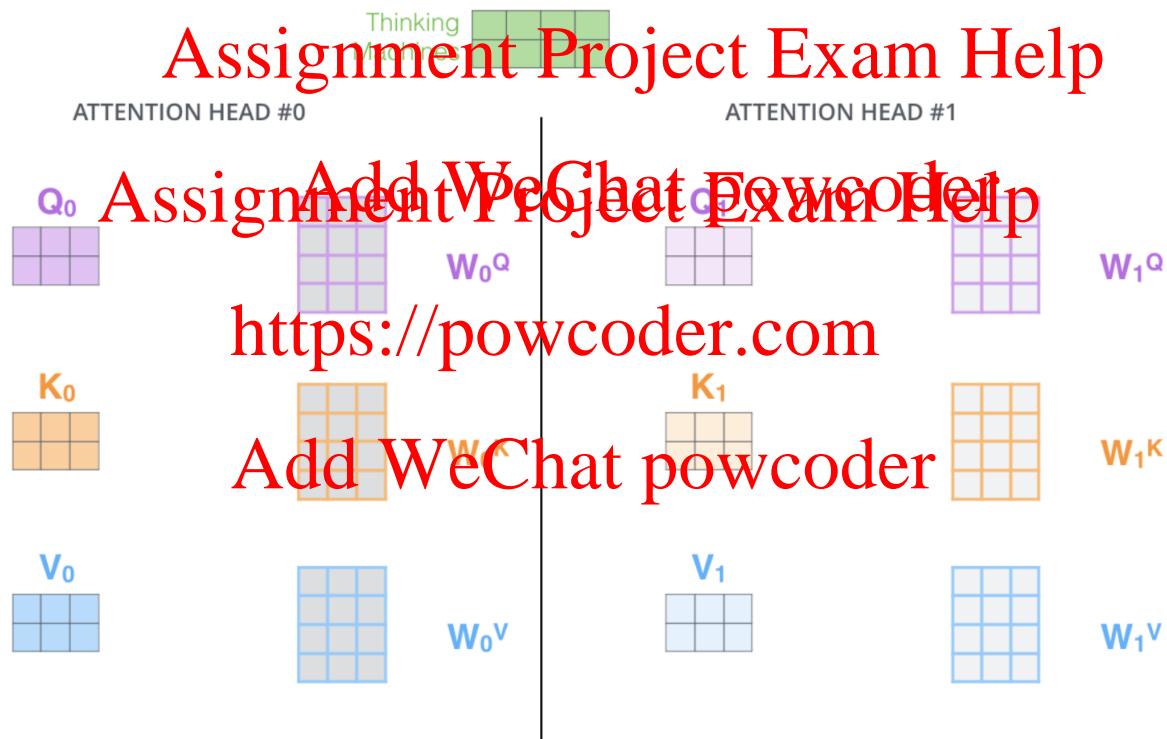
<https://powcoder.com>

## Assignment Project Exam Help

- ▶ Attention is an effective mechanism to model dependencies.
- ▶ By varying the keys, values, and queries, we get different variants of attention. We get self-attention by making the key, value and query to be projections of tokens in the same sentence. <https://powcoder.com>
- ▶ Intuitively self-attention contextualizes a word token in a sentence (or any sequence of word tokens), which is useful in interpreting the word token, which will in turns helps other tasks that build on the understanding of words.

## Multiple Attention “heads”

<https://powcoder.com>



With multi-headed attention, we maintain separate Q/K/V weight matrices for each head resulting in different Q/K/V matrices. As we did before, we multiply X by the WQ/WK/WV matrices to produce Q/K/V matrices.

## Multiple Attention “heads”

<https://powcoder.com>

If we do the same self-attention calculation we outlined above, just eight different times with different weight matrices, we end up with eight different  $Z$  matrices



## Multiple Attention “heads”

<https://powcoder.com>

How do we do that? We concat the matrices then multiple them by an additional weights matrix  $W_O$ .

## Assignment Project Exam Help

1) Concatenate all the attention heads



2) Multiply with a weight matrix  $W_O$  that was trained jointly with the model



<https://powcoder.com>

3) The result would be the  $Z$  matrix that captures information from all the attention heads. We can send this forward to the FFN!

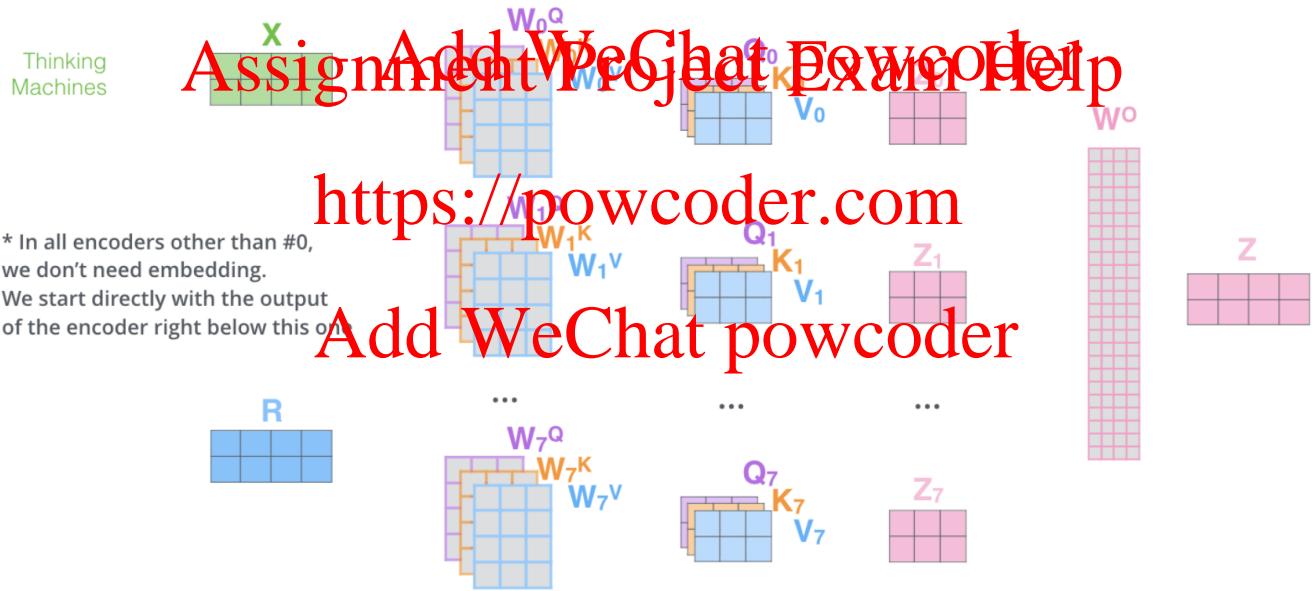
Add WeChat powcoder

$$= \begin{matrix} Z \\ \hline \end{matrix}$$

## Multiple Attention “heads”

<https://powcoder.com>

- 1) This is our input sentence\*
- 2) We embed each word\*\* with weight matrices
- 3) Split into 8 heads.
- 4) Calculate attention using the resulting Q/K/V matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix  $W^o$  to produce the output of the layer



What does self-attention do?

<https://powcoder.com>

## Assignment Project Exam Help

- ▶ Contextualization
  - ▶ <https://powcoder.com>
- ▶ why is contextualization a good thing?
  - ▶ “Work out the solution in your head.”
  - ▶ “Heat the solution to 75° Celsius.”
- ▶ Having the same embedding for both instances of “solution” doesn’t make sense

The Transformer architecture in (almost) its entirety  
<https://powcoder.com>

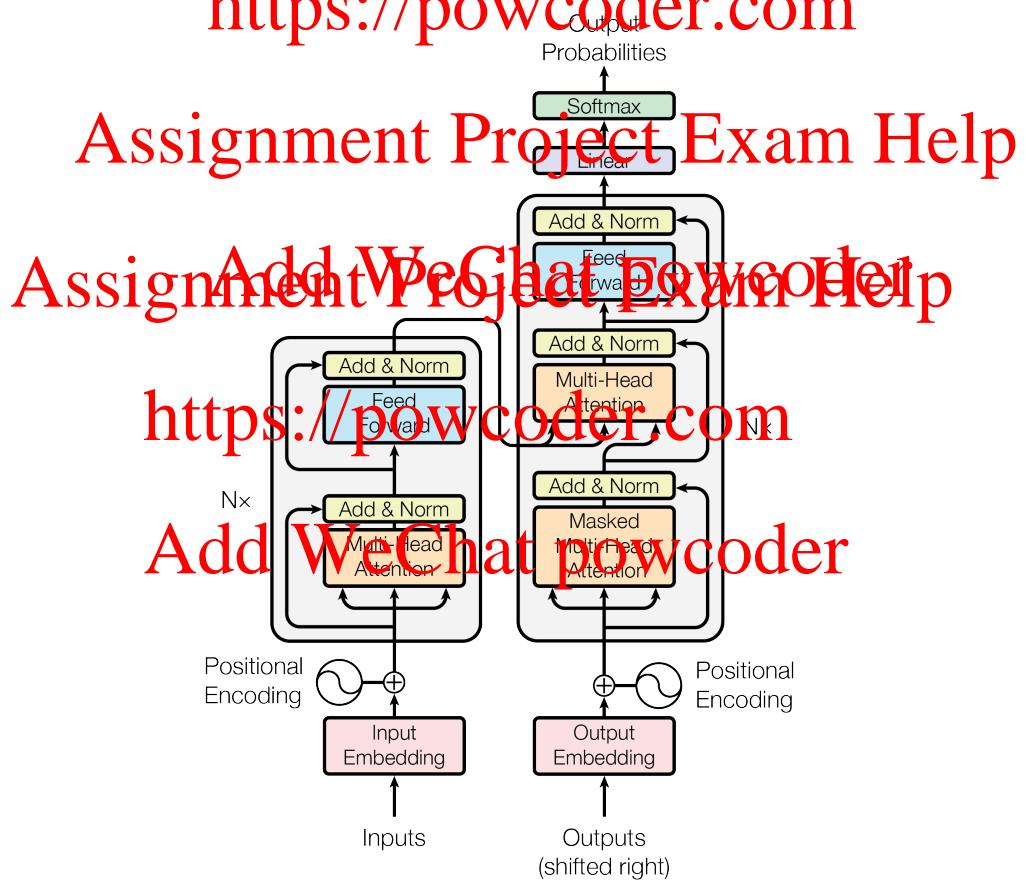


Figure 1: The Transformer - model architecture.

## Position Embedding

The original Transformer includes an absolute, precomputed 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 embedding  $e_w$ , and  $d_{model} = 4$ , the operation would be:

$$e_{w'} = e_w + \left[ \sin\left(\frac{pos}{10000^0}\right), \cos\left(\frac{pos}{10000^0}\right), \sin\left(\frac{pos}{10000^{2/4}}\right), \cos\left(\frac{pos}{10000^{2/4}}\right) \right]$$
$$= e_w + \left[ \sin(pos), \cos(pos), \sin\left(\frac{pos}{100}\right), \cos\left(\frac{pos}{100}\right) \right]$$

- ▶ 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://powcoder.com>

### Assignment Project Exam Help

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

$$FFN(x) = \Theta_2(\text{ReLU}(\Theta_1x + b_1)) + b_2$$

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

## Layer normalization

- ▶ Given a minibatch of input of size  $m$ ,  $B = \{x_1, x_2, \dots, x_m\}$ , after applying layer normalization to each sample in the batch, 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 steps:
  - ▶ Compute the mean and variance of each sample in the batch:

$$\mu_i = \frac{1}{K} \sum_{k=1}^K x_{i,k} \quad \sigma_i^2 = \sum_{k=1}^K (x_{i,k} - \mu_i)^2$$

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

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

- ▶ 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://powcoder.com>

BPE is a simple data compression technique

- ▶ Initialize the vocabulary with the character vocabulary, and segment each words into a sequence of characters, plus a end-of-word symbol.
- ▶ Iteratively count all symbol pairs, and replace each most frequent pair [A, B'] with a new symbol [AB']. Each merge produces a new symbol which represents a character ngram. 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://powcoder.com>

## Assignment Project Exam Help

In [13]: `import re, collections`

In [14]: `def get_stats(vocab):  
 pairs = collections.defaultdict(int)  
 for word, freq in vocab.items():  
 symbols = word.split()  
 for i in range(len(symbols)-1):  
 pairs[symbols[i],symbols[i+1]] += freq  
 return pairs`

In [15]: `def merge_vocab(pair, v_in):  
 v_out = {}  
 bigram = re.escape(' '.join(pair))  
 p = re.compile(r'(?<!\S)' + bigram + r'(?!\S)')  
 for word in v_in:  
 w_out = p.sub(''.join(pair), word)  
 v_out[w_out] = v_in[word]  
 return v_out`

# Python implementation

<https://powcoder.com>

## Assignment Project Exam Help

```
In [18]: vocab = {'l o w </w>':5, 'l o w e r </w>':2, 'n e w e s t </w>':6, 'w i d e s t </w>':3}
num_merges = 10
for i in range(num_merges):
    pairs = get_stats(vocab)
    best = max(pairs, key=pair_val)
    vocab = merge_vocab(best, vocab)
    print(best)
    print (vocab)
```

( 'e', 's')  
( 'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t </w>': 3}  
( 'es', 't')  
( 'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t </w>': 3}  
( 'est', '</w>')  
( 'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t </w>': 3}  
( 'l', 'o')  
( 'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t </w>': 3}  
( 'lo', 'w')  
( 'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t </w>': 3}  
( 'n', 'e')  
( 'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t </w>': 3}  
( 'ne', 'w')  
( 'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t </w>': 3}  
( 'new', 'est</w>')  
( 'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t </w>': 3}  
( 'low', '</w>')  
( 'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t </w>': 3}  
( 'w', 'i')  
( 'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t </w>': 3}

Add WeChat powcoder

<https://powcoder.com>

## Tutorials

<https://powcoder.com>

## Assignment Project Exam Help

- ▶ NMT with Tensorflow <https://github.com/tensorflow/nmt>
- ▶ Attention visualization:  
<https://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/>
- ▶ Illustrated transformer:  
<http://jalammar.github.io/illustrated-transformer/>

BERT: Bidirectional Encoder Representations from  
Transformers      <https://powcoder.com>

## Assignment Project Exam Help

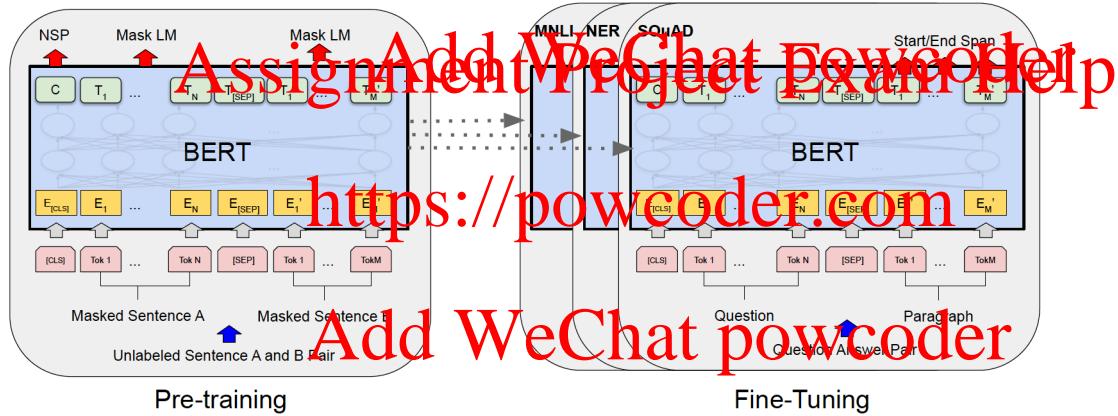
Input Representations:

- ▶ The input can be both a single sentence or a pair of sentences
- ▶ “Sentence” here just means a sequence of words.
- ▶ Sentences are broken down into WordPieces (token vocab size= 30K).
- ▶ The first token of the input sequence is a special symbol [CLS]
- ▶ Sentence pairs are packed together into one sentence and 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://powcoder.com>

Assignment Project Exam Help



## Pre-Training BERT

<https://powcoder.com>

## Assignment Project Exam Help

- ▶ The Transformer is trained with a Masked Language Model (MLM) task, and a next sentence prediction (NSP) task for sentence pairs.
- ▶ The MLM task (or the Cloze task) is trained by predicting 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://powcoder.com>

Use the output of the  
masked word's position  
to predict the masked word



# Assignment Project Exam Help

Add WeChat powcoder

<https://powcoder.com>

Add WeChat powcoder

Randomly mask  
15% of tokens

1 [CLS] 2 Let's 3 stick 4 to 5 [MASK] 6 in 7 this 8 skit ... 512

Input

[CLS] Let's stick to improvisation in this skit

BERT's clever language modeling task masks 15% of words in the input and asks the model to predict the missing word.

## BERT next sentence prediction

<https://powcoder.com>

Predict likelihood  
that sentence B  
belongs after  
sentence A

Assignment Project Exam Help

Assignment Project Exam Help

Tokenized  
Input

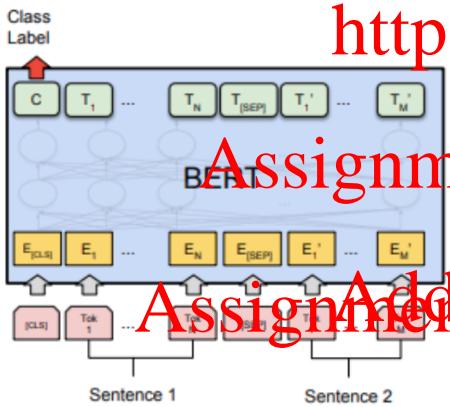
[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flightless birds [SEP]

Input

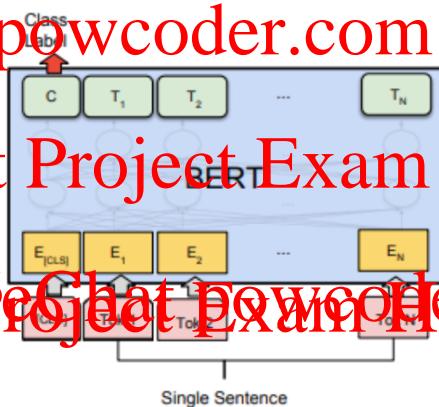


## BERT: fine tune to specific tasks

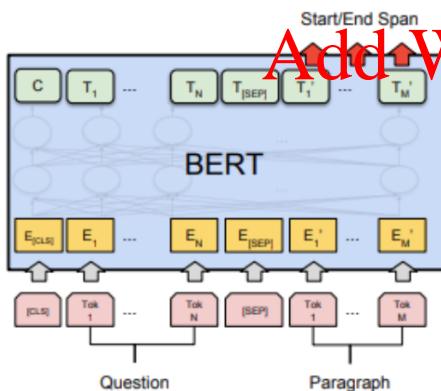
<https://powcoder.com>



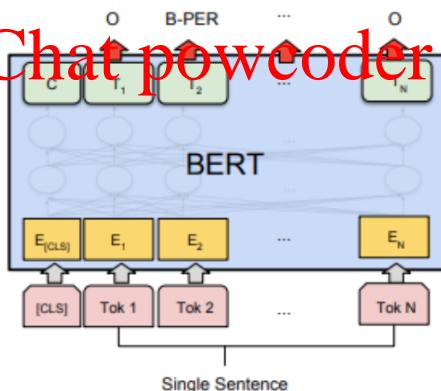
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



(b) Single Sentence Classification Tasks:  
SST-2, CoLA



(c) Question Answering Tasks:  
SQuAD v1.1



(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER