

Details of the BERT Model

(**B**idirectional **E**ncoder **R**epresentations from
Transformers)

By: Hiva Mohammadzadeh

Structure of the model

- Pre-trained deep learning model that uses bidirectional transformers to create context-based word embeddings.
- An encoder only model - Reads the entire sequence at once allowing the model to learn the context of a word based on all of its surrounding words.
- Consists of 12 Transformer Encoder blocks where each block has two sub-layers: a multi-head self-attention mechanism and a position-wise fully connected feed-forward network.
- Two versions: BERT Base and BERT Large.
- We will use BERT base with 128 tokens and an embedding dimension of 768 and 3072 feed-forward filter size.

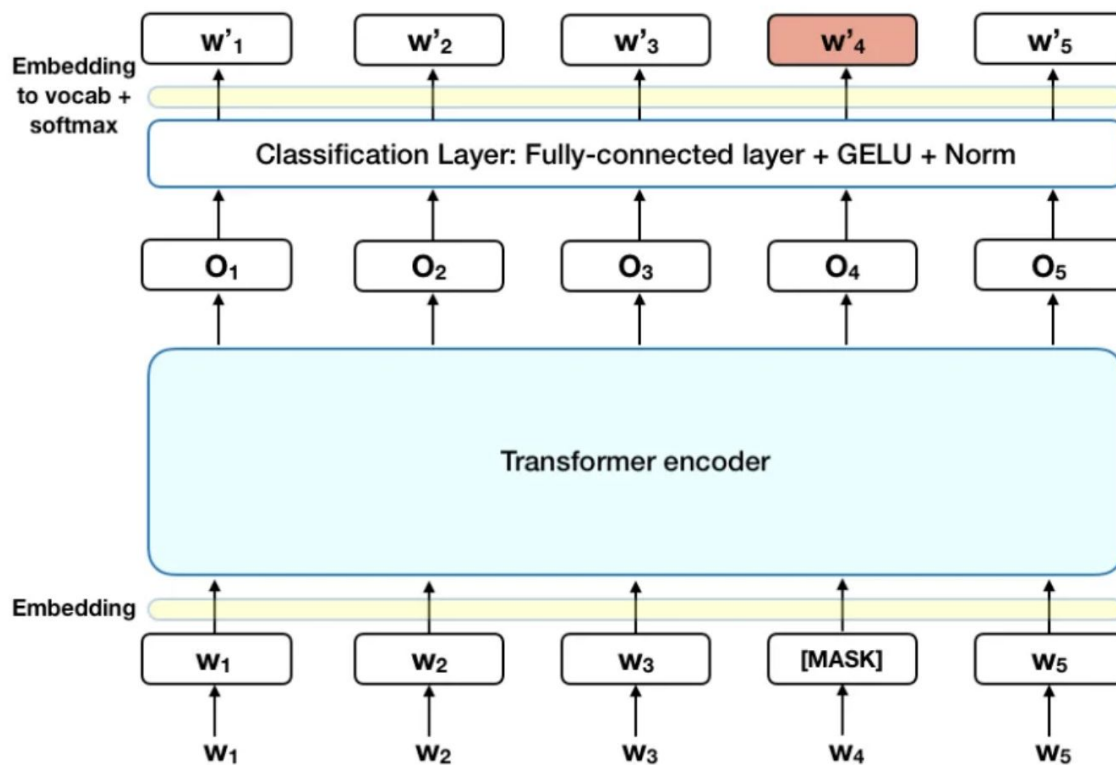
Training Scheme

- Pre-Training: Pre-trained in two stages
 - a. Masked Language Model Objective:
 - Randomly mask 15% of tokens in each input sequence.
 - Train the model to predict the masked tokens based on the surrounding context.
 - b. Next Sentence Prediction Objective:
 - Train the model to predict if two sentences are consecutive or not.
 - Improves the ability to handle tasks that involve multiple sentences.
- Fine-Tuning:
 - a. Fine-Tune the model on specific downstream tasks by adding a task-specific output layer.

BERT Architecture

Two key ingredients:

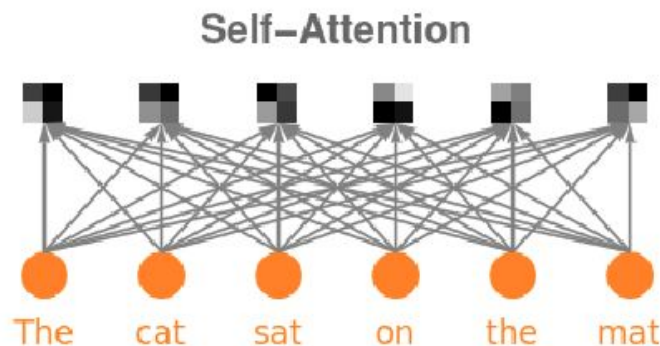
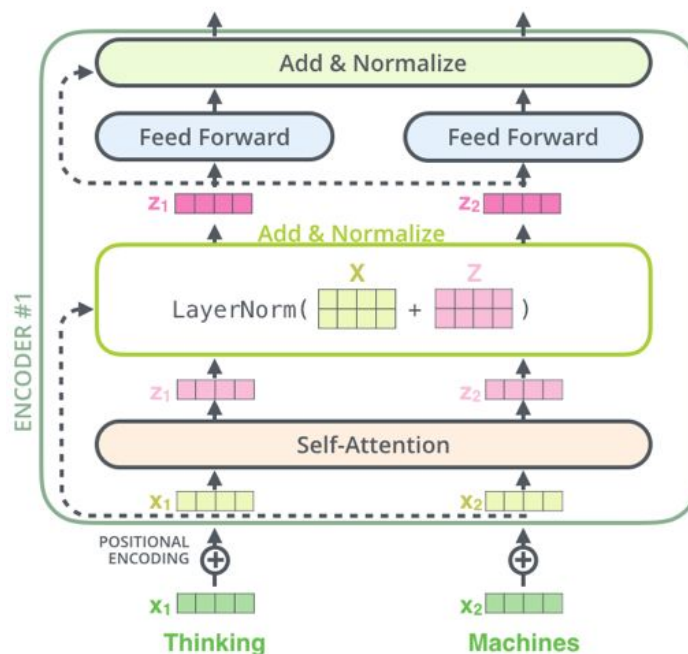
- Transformer Encoder
- Task specific output layer



Transformer Encoder

Two components: Attention and FFN

- Self-attention is used to capture the relationships between all the tokens in the input sequence.
- Feed Forward Network (FFN) module is used in every transformer block to process the output of the normalization layer in a way to better fit it to the next attention layer.
- Residual connections are used to avoid vanishing gradient problem.
- The Layer Normalizations are used to improve the model's convergence speed.
- GELU is the non-linearity used because it has been found to perform better than the other activation functions.



Dimensions of Weight Matrices

- The input embedding matrix has a dimension of $(\text{batch_size}, \text{sequence_length}, \text{embedding_size}) = (\text{batch_size}, 128, 768)$. But in order to do the attention parallel for different heads, we use $(\text{batch_size}, 128, 768/12 \text{ heads}) = (\text{batch_size}, 128, 64, 12)$
- The W_k , W_v , and W_q weight matrices used in the **self-attention mechanism** have dimensions of:
 - Query weight matrix: $(\text{batch_size}, \text{hidden_size}, \text{hidden_size}) = (\text{batch_size}, 768, 768)$
 - Key weight matrix: $(\text{batch_size}, \text{hidden_size}, \text{hidden_size}) = (\text{batch_size}, 768, 768)$
 - Value weight matrix: $(\text{batch_size}, \text{hidden_size}, \text{hidden_size}) = (\text{batch_size}, 768, 768)$
 - Output weight matrix: $(\text{batch_size}, \text{hidden_size}, \text{hidden_size}) = (\text{batch_size}, 768, 768)$
- The K , V , Q matrices also have dimensions of:
 - Query matrix: $(\text{batch_size}, \text{sequence_length}, \text{hidden_size}) = (\text{batch_size}, 128, 768)$ or $(\text{batch_size}, 128, 64, 12)$
 - Key matrix: $(\text{batch_size}, \text{sequence_length}, \text{hidden_size}) = (\text{batch_size}, 128, 768)$ or $(\text{batch_size}, 128, 64, 12)$
 - Value matrix: $(\text{batch_size}, \text{sequence_length}, \text{hidden_size}) = (\text{batch_size}, 128, 768)$ or $(\text{batch_size}, 128, 64, 12)$
- The weight matrices used in the **feed-forward neural network** have dimensions of:
 - First dense layer: $(\text{hidden_size}, 4*\text{hidden_size}) = (768, 3072)$.
 - Second dense layer: $(4*\text{hidden_size}, \text{hidden_size}) = (3072, 768)$.

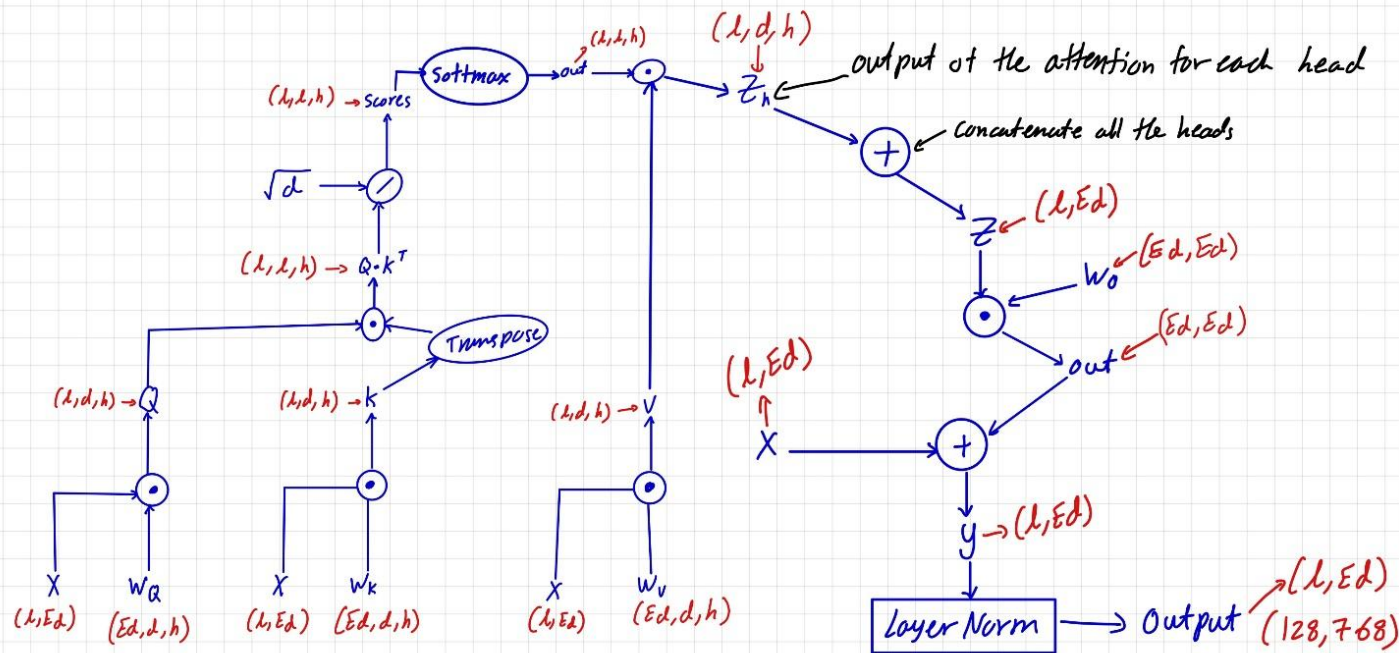
Computation Diagram of Attention Block

Multi-head self-Attention

Size of the embedding dimension = 768 $\rightarrow \epsilon d$

Length of the input tokens = 128 $\rightarrow l$

Dimension of keys, values, and queries per head = $\frac{\epsilon d}{H} = \frac{768}{12} = 64 \rightarrow d$ Number of Heads = 12 $\rightarrow H$

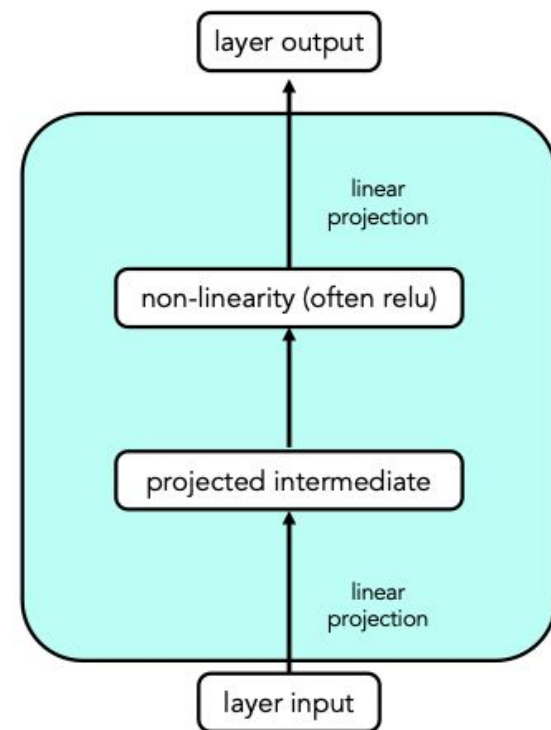
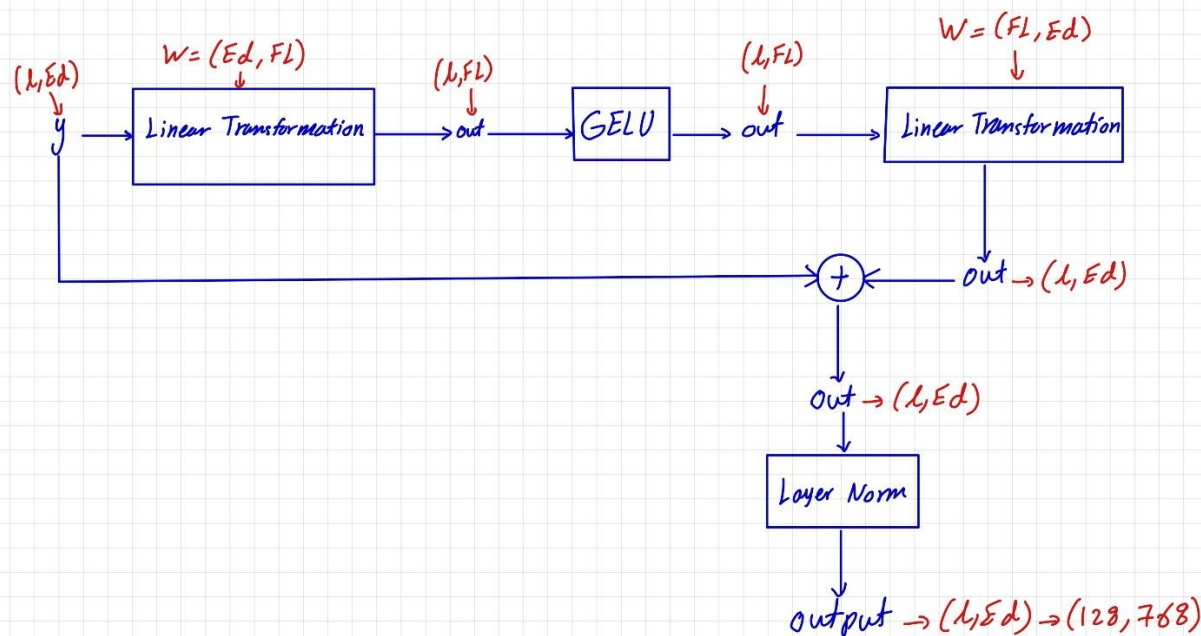


Computation Diagram of Feed Forward Layer

Position-wise Feed Forward Networks

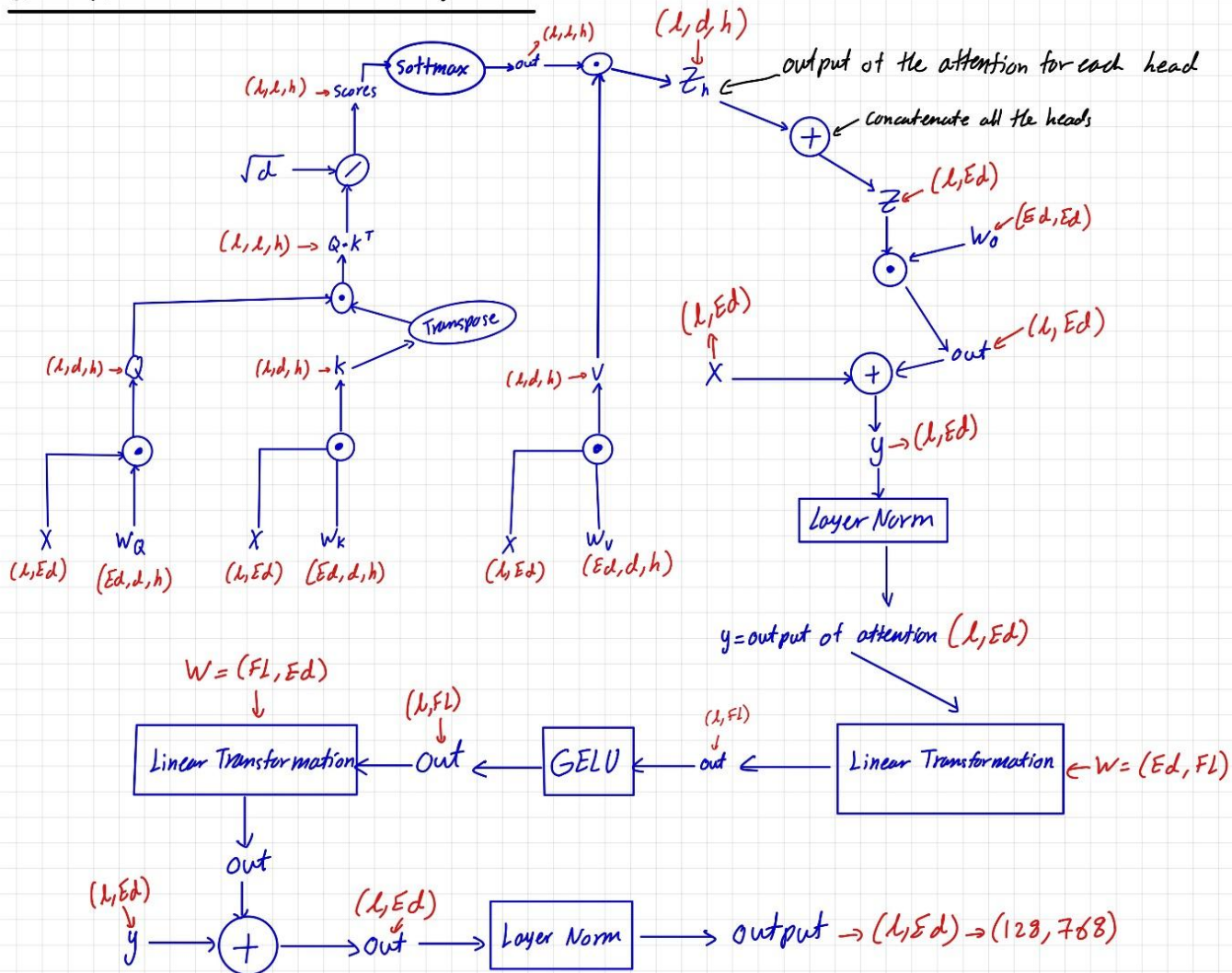
Input = output of the Attention. $\rightarrow y$ size of the embedding dimension = 768 $\rightarrow Ed$

Length of the input tokens = 128 $\rightarrow l$ fn. hidden size = 3072 $\rightarrow FL$



Computation Diagram of the whole BERT architecture

BERT Computation Diagram



FLOPS

Assuming that every dot product requires 1 multiplications and 1 addition, I calculated the following FLOPS:

- Flops for attention block:
 - Multiplication between embedding matrix and the weight matrices of K,Q and V : $3 \times 128 \text{ tokens} \times 768 \text{ embeddings} \times 768 \text{ embeddings} \times 2 \text{ FLOPs} = 452,984,832 \text{ FLOPs}$
 - Multi-Head Self-Attention for 12 layers: $2 \times 12 \times 128 \times 64 \times 128 \text{ FLOPs}$ (dot product of Q and K) + $2 \times 12 \times 128 \times 64 \times 128 \text{ FLOPs}$ (scaled dot-product attention between output of softmax and V) + $2 \times 128 \times 768 \times 768 \text{ FLOPs}$ (concatenation of heads and matrix multiplication of W_0) = 201,326,592
- FLOPS for Feed Forward Network Module:
 - Feed-Forward Network: $\text{FLOPs} = 128 \times 768 \times 3072 \times 2 \text{ FLOPs} \times 2 \text{ FLOPs} + 128 \text{ FLOPs} = 1,207,959,552$
- Total FLOPs for 12 transformer block: $\text{FLOPs} = 1,207,959,552 \text{ FLOPs (FFN)} + 201,326,592 \text{ FLOPs (MHA)} + 452,984,832 \text{ FLOPs (weight matrix calculations)} = 1,862,270,976 \text{ FLOPs}$