

Transformers

CNN, RNN, LSTM are sequential architectures has limited the possibilities to train on the same scale of data that showed value for computer vision

• Transformers processes sequential data with much more parallelization.

* Types of modifications to the transformers block,

| Modification | Transformer |
|---------------------------|---|
| Lightweight transformers | Lite Transformer [274] Funnel Transformer [66] DeLighT [180] |
| Cross-block connectivity | Reformer [107] Transparent Attention [19] |
| Adaptive computation time | Universal Transformer [69] Conditional Computation Transformer [18] DeeBERT [276] |
| Recurrent | Transformer-XL [67] Compressive Transformer [204] Memformer [287] |
| Hierarchical | HIBERT [296] Hi-Transformer [270] |
| Different architectures | Macaron Transformer [174] Sandwich Transformer [201] Differentiable Architecture Search [299] |

* Transformer block changes..

- Decreasing memory footprint and compute.
- Adding connections between transformers blocks.
- Adaptive computation time (allow early stopping during training).
- Recurrence or hierarchical structure.
- changing the architecture more drastically.

* Parts of Transformer block:-

- Positional encodings.
- multi-head attention.
- residual connections with layer normalization.
- Position-wise feedforward network.

TABLE 1.2 Types of modifications to the multi-head attention module

| Modification | Transformer |
|-------------------------------|--|
| Low-rank | Performer [53] Nystromformer [277] Synthesizer [241] |
| Attention with prior | Gaussian Transformer [102] Realformer [107] Synthesizer [241] Longformer [25] |
| Improved multi-head attention | Talking-heads Attention [227] Multi-Scale Transformer [234] |
| Complexity reduction | |

| Modification | Transformer |
|----------------------------|--|
| | Longformer [25] Reformer [142] Big Bird [292] Performer [53] Routing Transformer [214] |
| Prototype queries | Clustered Attention [256] Informer [302] |
| Clustered key-value memory | Set Transformer [151] Memory Compressed Transformer [167] Linformer [259] |

Kinds of Positional encoding:-

- absolute positional encodings (standard transformer).
- relative \Rightarrow (Transformer-XL).
- hybrid encodings have absolute and relative.
- implicit encodings that provide information about sequence.

* Pre-training Methods:-

- encoder only : BERT
- decoder only : GPT-3
- encoder-decoder : T5 and BERT

TABLE 1.4 Libraries and Tools

| Organization | Language and Framework | API | Pre-trained |
|--------------|------------------------------|-----|-------------|
| AllenNLP | Python and PyTorch | Yes | Yes |
| HuggingFace | Jax, PyTorch, and TensorFlow | Yes | Yes |
| Google Brain | TensorFlow | Yes | |
| GluonNLP | MXNet | Yes | Yes |

* Encoder-Decoder Architecture:-

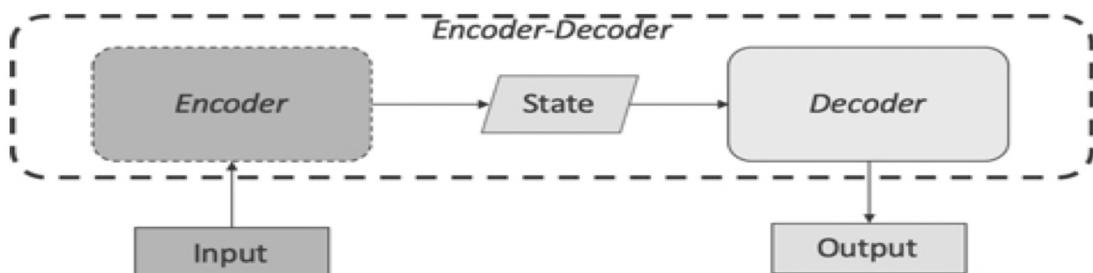
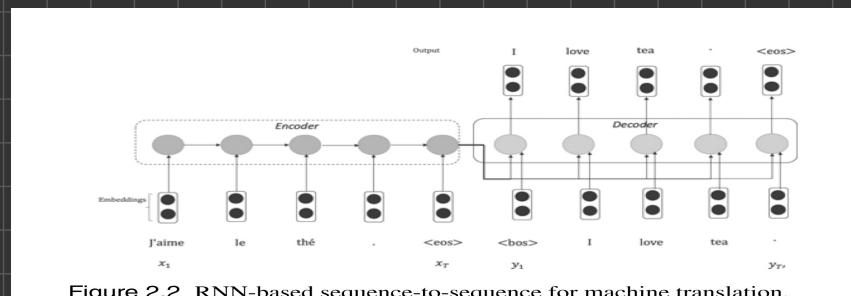


Figure 2.1 Encoder-Decoder architecture.

- The encoder component takes a variable-length sequence and converts it into a fixed-length output-state.
- The decoder component takes a fixed-length state and converts it back into a variable-length output.

* Sequence-to-Sequence



Encoder :-

Sentence is tokenized into words and words are mapped into feature vectors.

$$\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t)$$

Input: generates a new hidden state.
hidden state
previous hidden state

$$\mathbf{c} = m(\mathbf{h}_1, \dots, \mathbf{h}_T)$$

mapping function.
context vector
encodes the information of the entire input sequence.

* Embedding mappings transform the input

$\mathbf{x}_1, \dots, \mathbf{x}_T$ to $\mathbf{x}_1, \dots, \mathbf{x}_T$ vectors.

* Decoder:-

$$P(\mathbf{y}_{t'} | \mathbf{y}_{t'-1}, \dots, \mathbf{y}_1, \mathbf{c}) = \text{softmax}(\mathbf{s}_{t-1}, \mathbf{y}_{t'-1}, \mathbf{c})$$

Decoder predicts a probability distribution for output tokens at each time step, and softmax gives the distribution over the words.

- Loss function:- cross-entropy is used for optimization.

$$\max_{\theta} \frac{1}{N} \sum_{n=1}^N \log p_{\theta}(\mathbf{y}^{(n)} | \mathbf{x}^{(n)})$$

- RNNs have issues with vanishing and explosive gradients.

* ATTention Mechanism :-

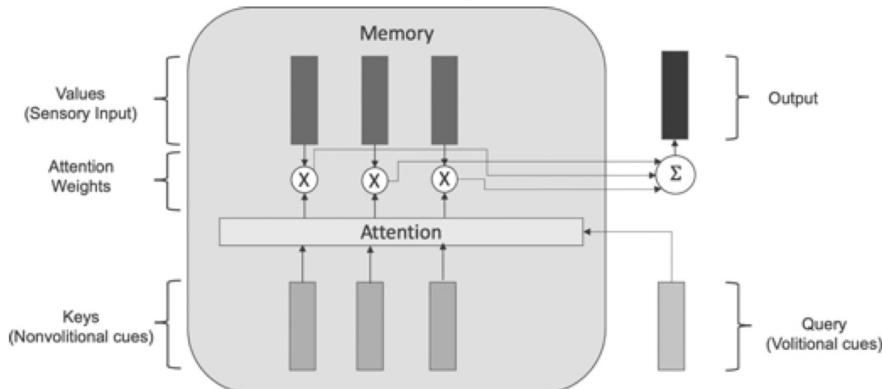
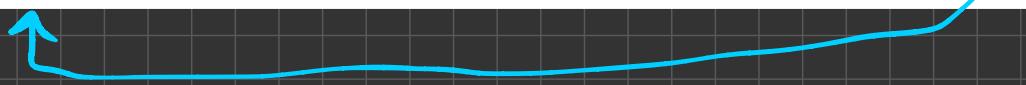


Figure 2.3 Attention mechanism showing query, keys, values, and output vector interactions.

- Attention mechanism can be considered as a memory with keys and values.

To formalize, let us consider the memory unit consisting of n key-value pairs $(\mathbf{k}_1, \mathbf{v}_1), \dots, (\mathbf{k}_n, \mathbf{v}_n)$ with $\mathbf{k}_i \in \mathbb{R}^{d_k}$ and $\mathbf{v}_i \in \mathbb{R}^{d_v}$. The attention layer receives an input as query $\mathbf{q} \in \mathbb{R}^{d_q}$ and returns an output $\mathbf{o} \in \mathbb{R}^{d_v}$ with same shape as the value \mathbf{v} .



- The attention layer measures the similarity between the query and the key using a score function α which returns a_1, \dots, a_n and k_1, \dots, k_n .

$$a_i = \alpha(\mathbf{q}, \mathbf{k}_i)$$

score

score
function

- Attention weights are computed as softmax of the score :-

$$\mathbf{b} = \text{softmax}(\mathbf{a})$$

Each element of \mathbf{b} is

$$b_i = \frac{\exp(a_i)}{\sum_j \exp(a_j)}$$

- Output is the sum of the attention weights and values.

$$\mathbf{o} = \sum_{i=1}^n b_i \mathbf{v}_i$$

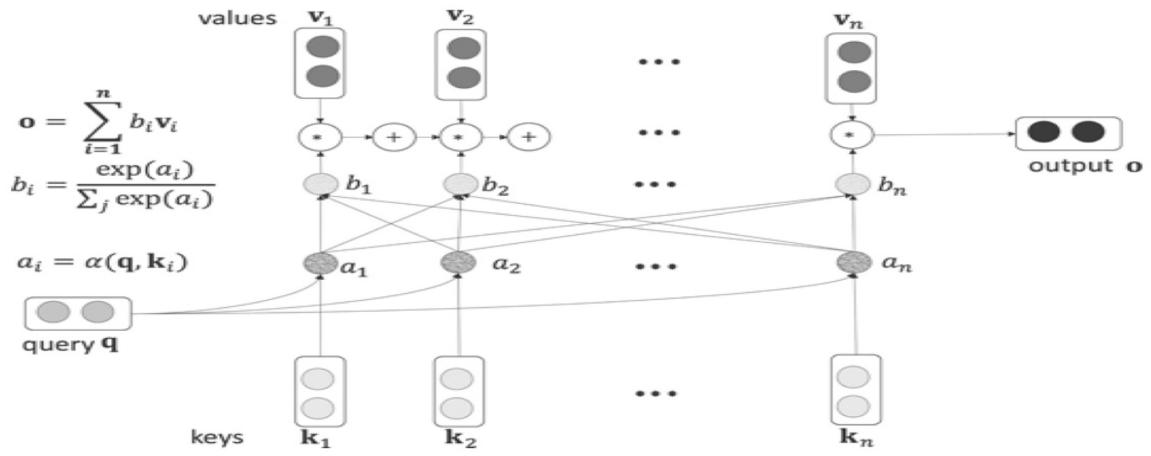


Figure 2.4 Attention mechanism showing query, keys, values, and output vector interactions.

* Types of Score-based Attention:

- Dot-product:

$$\alpha(\mathbf{q}, \mathbf{k}) = \mathbf{q} \cdot \mathbf{k}$$

• has no parameters to tune.

- Scaled dot product:

divides the dot product by $\sqrt{d_k}$.

To remove the influence of dimension d_k dim of k

$$\alpha(\mathbf{q}, \mathbf{k}) = \frac{\mathbf{q} \cdot \mathbf{k}}{\sqrt{d_k}}$$

* Attention-Based Sequence-to-Sequence:-

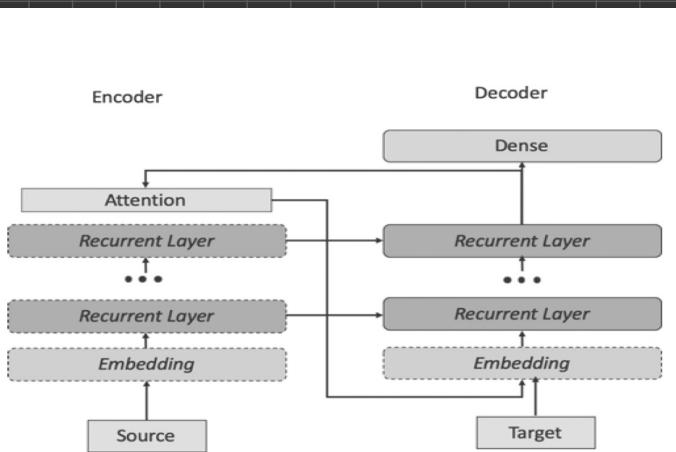


Figure 2.5 Encoder-decoder with attention layer.

The output of last encoder state are used as keys and values.

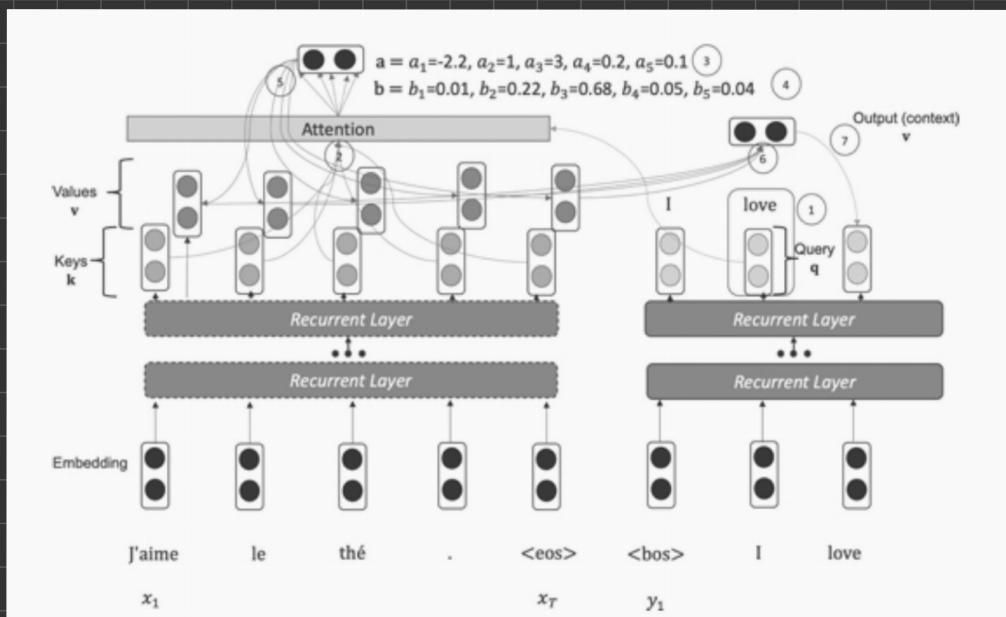
$\rightarrow \hat{z} = \text{softmax}(a)$ - decoder
at time $t-1$ is used as query

• Output from Attention layer \hat{z}
the context variable used for the next decoder state.

example:-

French Translation.

- Encoder state for tokens $\{\text{J'aime, le, thé...}\}$ are keys and values.
- Decoder at time $t-1$ generated the tokens I, Love
- Decoder output at time $t-1$, love, flows into attention layer as query.



* **Transformers**:- based on encoder-decoder architecture.
 Combines the advantages of CNN to parallelize the computations.
 and RNN to capture long-range, sequential information.

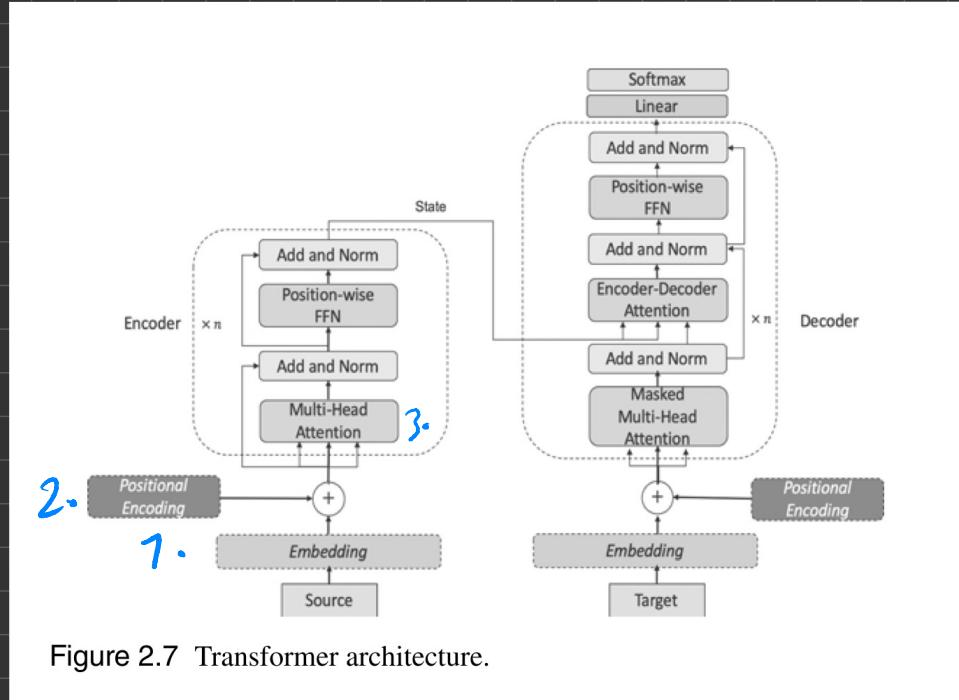


Figure 2.7 Transformer architecture.

1. Word Embedding:-

Convert Sentence of length ℓ to Matrix (vector) W of dim. (ℓ, d) .

2. Positional Encoding:- (word order information)

$$\mathbf{P}_{i,2j} = \sin(i/1000^{2j/d}) \rightarrow \text{even}$$

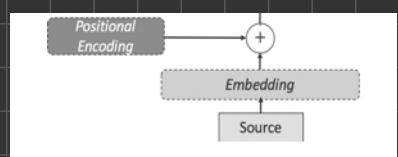
$$\mathbf{P}_{i,2j+1} = \cos(i/1000^{2j/d}) \rightarrow \text{odd}$$

⇒ the distance between tow time-steps across sentences is various lengths.

also

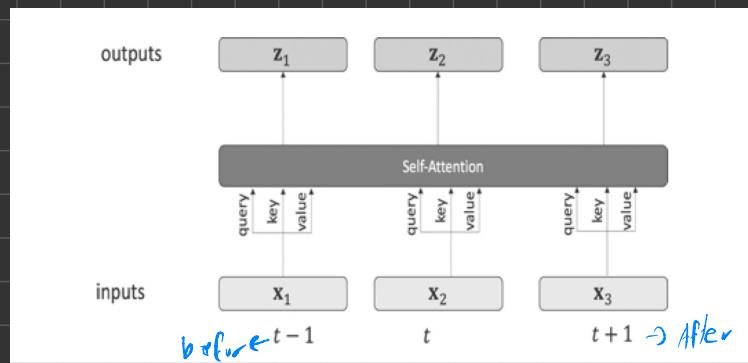
$$\text{Input representation } \mathbf{X} = \mathbf{W} + \mathbf{P} \in \mathbb{R}^{\ell \times d}.$$

Word embeddings Positional encoding



3. Self-attention:

- Input vectors x_i converted to the vectors z_i , through the self-attention layer.



- Each input vector (x_i) has 3 different vectors (Q_i, K_i, V_i).
 - Q, K, V are obtained from the input vector (x_i) at time i on the learnable weight matrices. W_Q, W_K and W_V .
 - Weight matrices are randomly initialized and weights learning from the training.
 - For the first attention layer of en. and de., the inputs are the summation of word embeddings (w) and positional encodings (p).
- Roles:
 - query vector q_i : to combine with every other key vectors.
 $\sum_{j=0}^P q_i K_j^T$ to influence the weights for its own output z_i
 - key vector k_i : to be matched with every other query vectors to get similarity with query and to influence the output through query-key product scoring.
 - value vector v_i : is extracting information by combining with output of the query-key scores to get the output vector z_i .

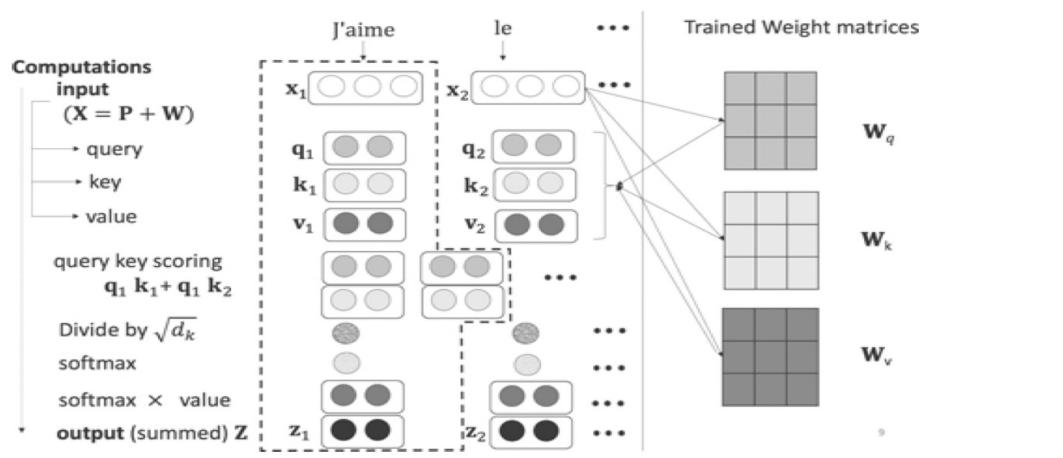


Figure 2.11 The dotted lines show the complete flow of computation for one input through a self-attention layer.

Instead of a vector computation for each token i , input matrix $\mathbf{X} \in \mathbb{R}^{l \times d}$ where l is the maximum length of the sentence and d is the dimension of the inputs, combines with each of the query, key, and value matrices as a single computation given by

$$\text{attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

* Multi-head Attention:

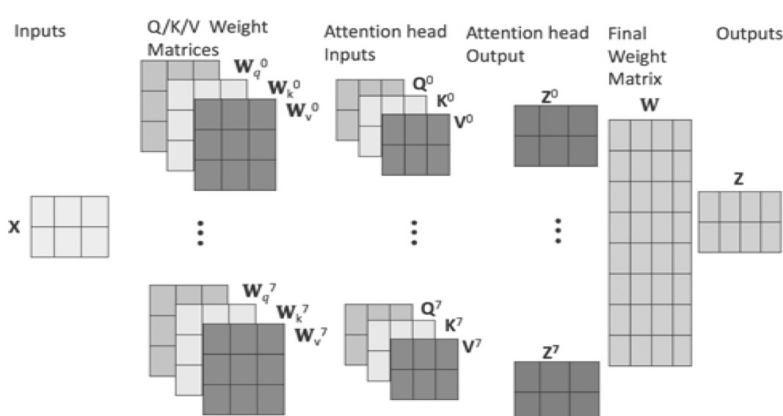


Figure 2.12 Multi-head attention.

$$\text{head}_i = \text{attention}(\mathbf{W}_q^i \mathbf{Q}, \mathbf{W}_k^i \mathbf{K}, \mathbf{W}_v^i \mathbf{V})$$

$$\text{multihead}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathbf{W}_O \text{concat}(\text{head}_1, \dots, \text{head}_h)$$

, each head can learn

something different.

• multi head Attention have multiple sets of q/k/v weight metrics.

• generating output matrices Z_i .

• Output matrices from each head are concatenated and multiplied with a Additional weight matrix \mathbf{W}_O to get a single final matrix \mathbf{Z} with vectors Z_i as output for each input X_i .

* Masked multi-head attention:

- Only previous target tokens need to be present and others to be masked.
- This is implemented by having a masking weight matrix M that has $-\infty$ for future tokens and 0 for previous tokens.

This computation is inserted after the scaling of the multiplication of \mathbf{Q} and \mathbf{K}^T and before the softmax so that the softmax results in the actual scaled values for previous tokens and the value 0 for future tokens.

$$\text{masked Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T + \mathbf{M}}{\sqrt{d_k}}\right)\mathbf{V}$$

* Residuals and layer Normalization:

- The input X are circuit to the Output Z and both are added and passed through layer Normalization addAndNorm ($X + Z$)
- Layer normalization ensures each layer to have 0 mean and a unit(1) variance.
- For each hidden unit, h_i we can compute:-

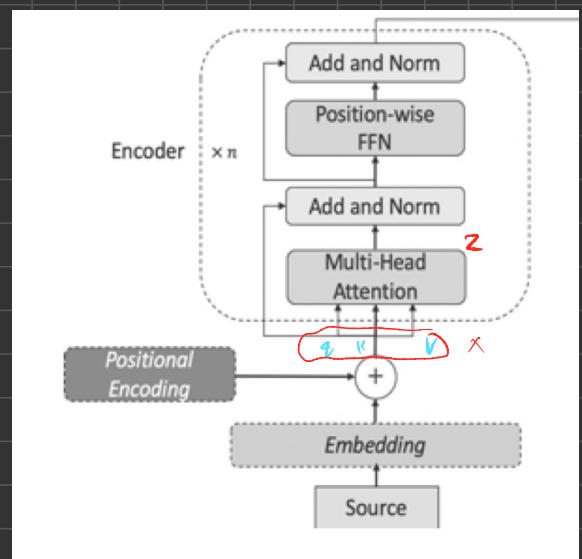
$$h_i = \frac{g}{\sigma}(h_i - \mu)$$

where g is the gain variable (can be set to 1), μ is the mean given by $\frac{1}{H} \sum_{i=1}^H h_i$ and σ is the standard deviation given by $\sqrt{\frac{1}{H} (h_i - \mu)^2}$.

* Position wise Feed-forward Networks:

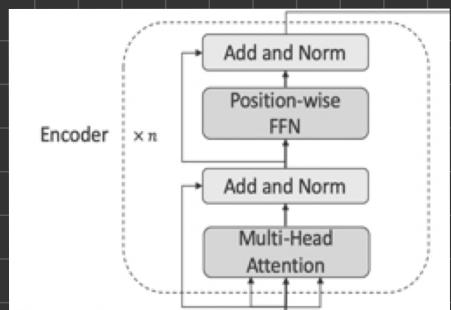
- For each position, similar linear transformations with a ReLU activation in between is performed.

$$FFN(\mathbf{x}) = \max(0, \mathbf{x}\mathbf{W}_1 + b_1)\mathbf{W}_2 + b_2$$



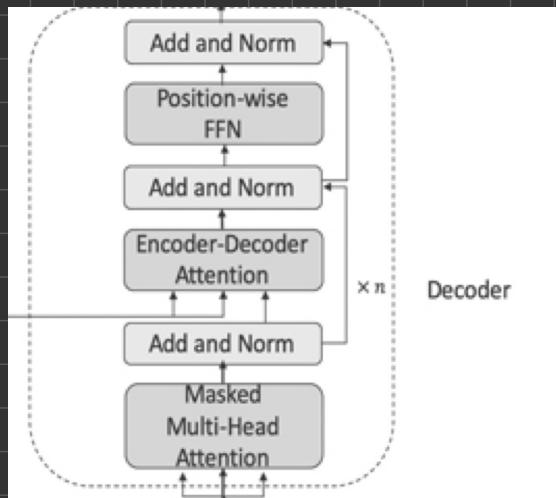
* Encoder:

- encoder block consists of n blocks of {Multi-Head Attention, addAndNorm, FFN, addandNorm}.



* Decoder:

- multi-head-attention in decoder attends to the target attention between masked outputs with themselves.
- The encoder-decoder attention layer creates attention between the source and the target.



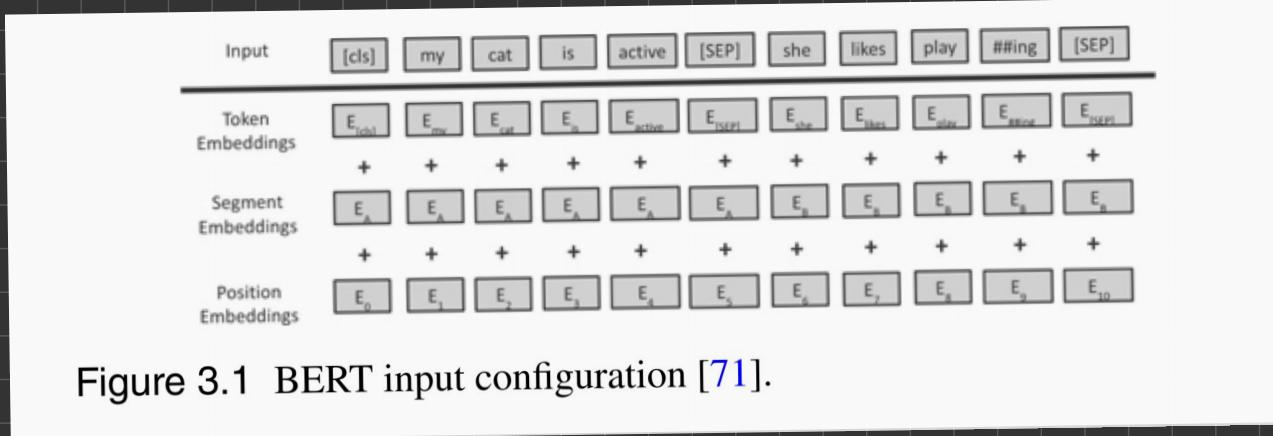


Figure 3.1 BERT input configuration [71].

Masked Language Model (MLM):

The idea is to randomly mask out a percentage of the input sequence tokens, replacing them with the special [MASK] token.

- modified input sequence runs through BERT.

- output representations of masked tokens are fed into a softmax layer over the word piece vocabulary.

Next Sentence Prediction (NSP):

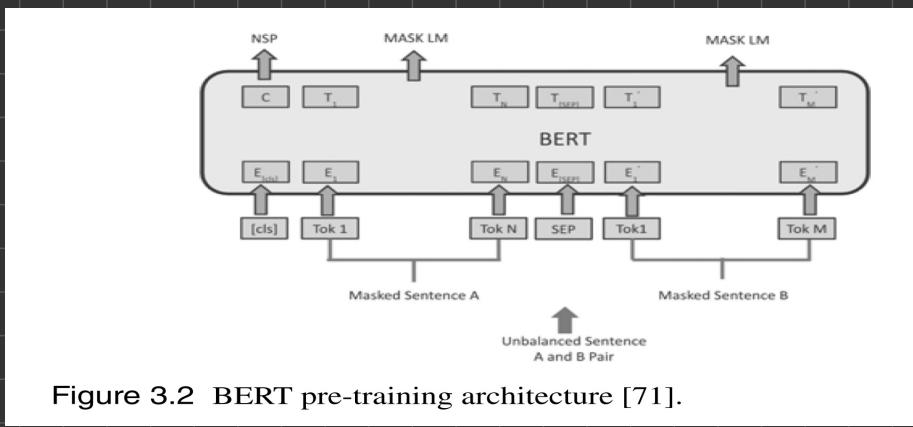


Figure 3.2 BERT pre-training architecture [71].

- the first sentence is prefixed with [CLS] token, then the two sentences are delimited by special token [SEP].

To summarize, using the same basic architecture as BERT, RoBERTa has shown that the following design changes significantly improves RoBERTa performance over BERT.

1. Pre-training the model for longer time
2. Using bigger batches
3. Using more training data
4. Removing the NSP pre-training task
5. Training on longer sequences
6. Dynamically changing the masking pattern applied to the training data.