

## MT Summer Term 2021 Ex14: Attention is all you need

1. In your own words, please explain the following pictorial representation of the transformer NMT system (source: Vaswani et al. 2017):

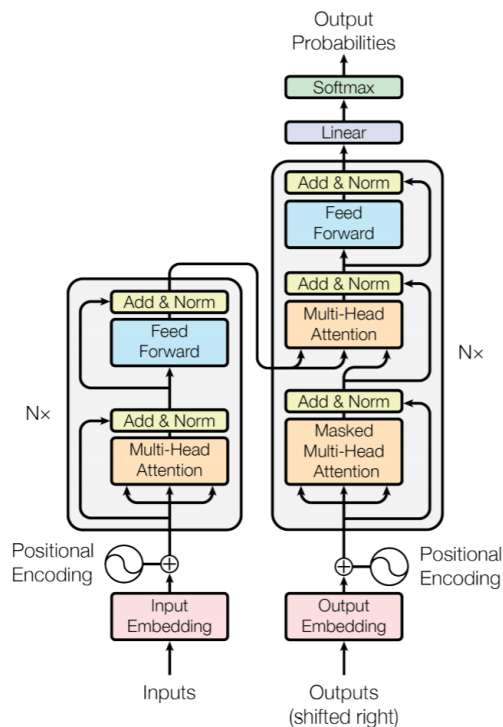


Figure 1: The Transformer - model architecture.

In your discussion, please cover

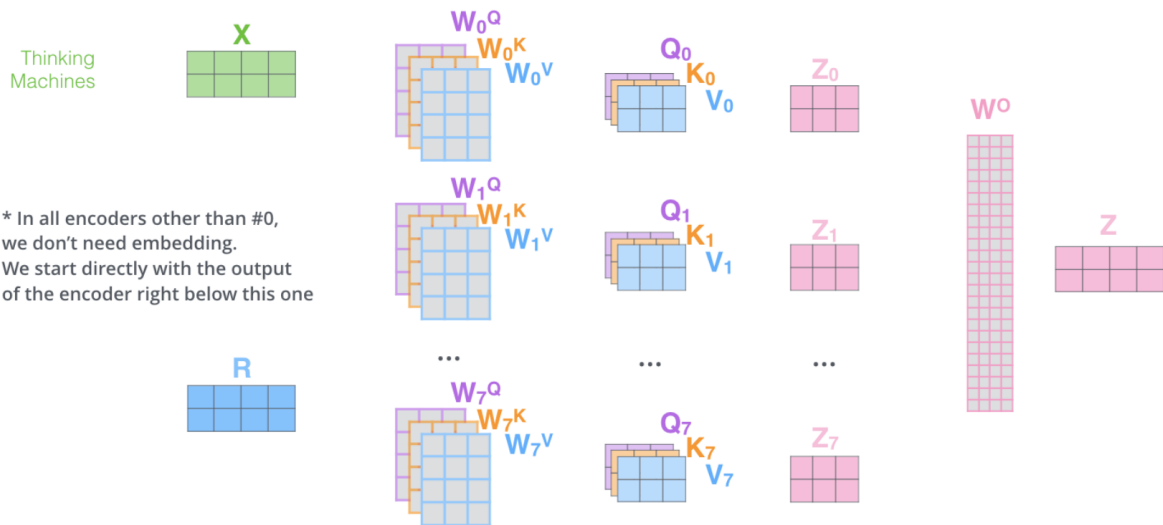
- Encoder-decoder stacks
- Self-attention
- Cross-attention
- Masked-attention
- Multi-head attention
- Query, key and values in attention computation
- Positional encoding
- Residual connections and layer normalisation
- Static embeddings and contextualised embeddings
- Autoregressive decoding

2. Are parameters shared between encoder blocks in the Transformer? Are parameters shared between decoder blocks in the Transformer? How does this compare to RNNs?

3. Comparing Transformer with simple (!) RNN based encoder-decoder systems, what are their computational complexities?

4. Why do you need the query-, key- and value-projection matrices in the picture below:

- 1) This is our input sentence\*
- 2) We embed each word\*
- 3) Split into 8 heads. We multiply  $X$  or  $R$  with weight matrices
- 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



5. In your own words please explain:

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

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

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

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) V$$

$Z$

The self-attention calculation in matrix form

6. Which parts of the Transformer make up BERT? Which parts of the Transformer make up GPT-X?

7. Please read the "Attention is all you need" paper.