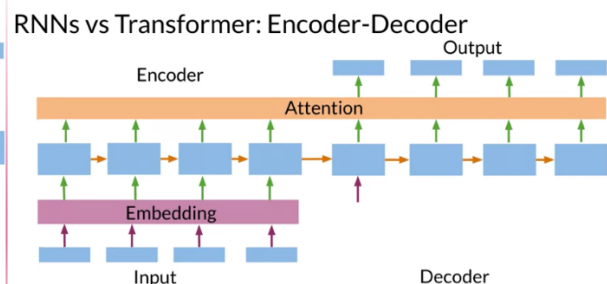
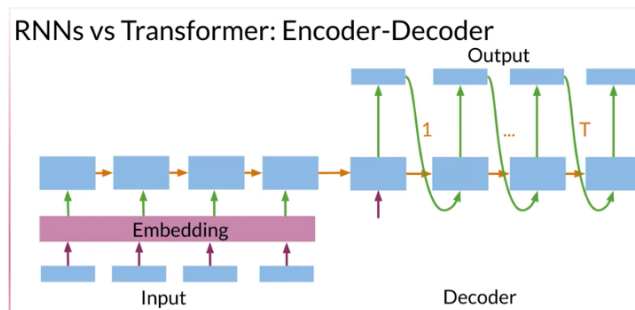
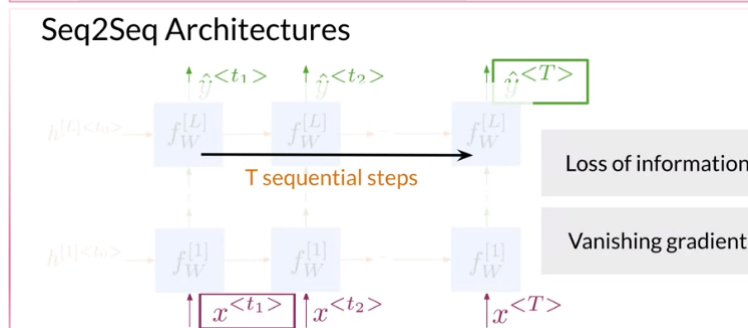
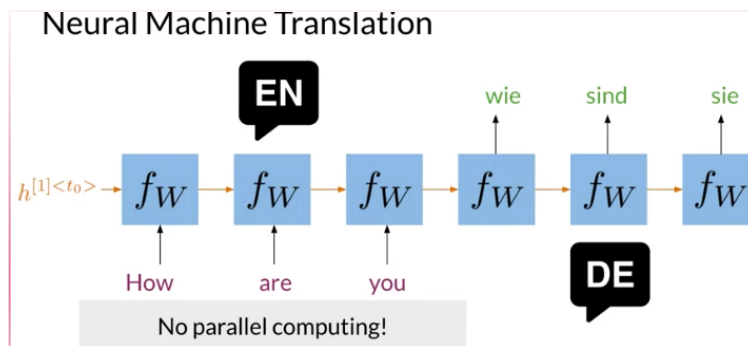


Text Summarization

Transformers vs RNNs

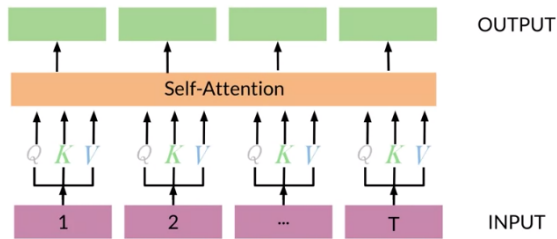
Issues with RNNs

- * No parallel computing
- * Loss of information
- * Vanishing gradients
- * Encoder-decoder based on RNNs is used to compute T sequential steps

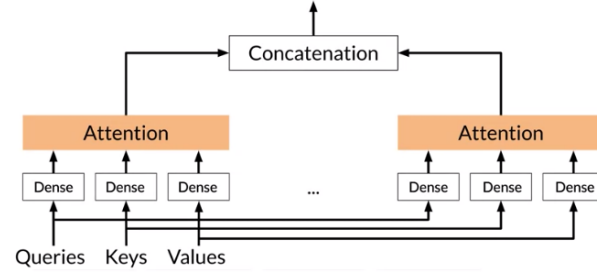


transformers are based on attention and don't require any sequential computation per layer, only one single step is needed. Additionally, the gradient steps that need to be taken from the last output to the first input in a transformer is just one. For RNNs, the number of steps is equal to T. Finally, transformers don't suffer from vanishing gradients problems that are related to the length of the sequences. Transformer differs from sequence to sequence by using multi-head attention layers instead of recurrent layers.

RNNs vs Transformer: Multi-headed attention

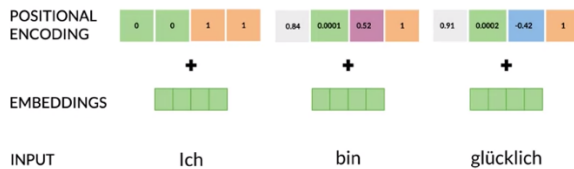


RNNs vs Transformer: Multi-headed attention



Transformers use a positional encoding to retain position info of input seq. Unlike the recurrent layer, the multi-head attention layer computes the outputs of each inputs in the sequence independently then it allows us to parallelize the computation. But it fails to model the sequential information for a given sequence. That is why you need to incorporate the positional encoding stage into the transformer model

RNNs vs Transformer: Positional Encoding



Transformer Applications

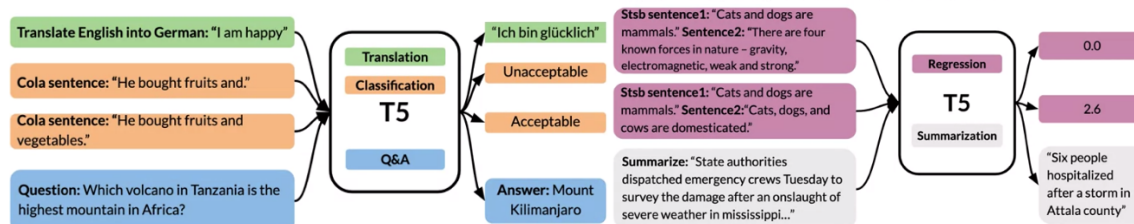
- * Text summarization
- * Auto complete
- * NER
- * Q&A
- * Translation
- * Chatbots
- * Sentiment Analysis
- * Market Intelligence
- * Text classification
- * Character recognition
- * Spell check

State of the Art Transformers

Radford, A., et al. (2018) Open AI	GPT-2: Generative Pre-training for Transformer
Devlin, J., et al. (2018) Google AI Language	BERT: Bidirectional Encoder Representations from Transformers
Colin, R., et al. (2019) Google	T5: Text-to-text transfer transformer

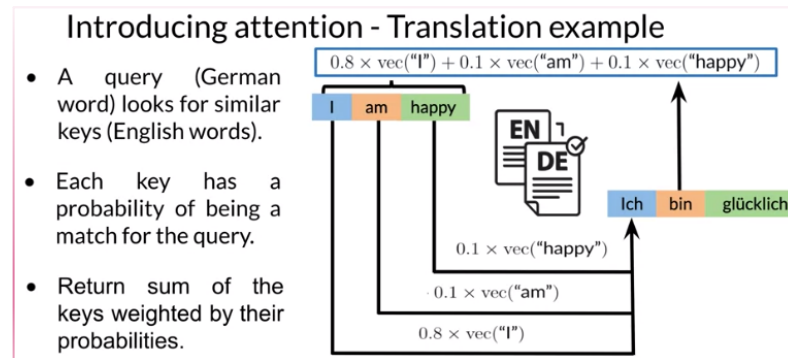
T5: Text-To-Text Transfer Transformer – all of the below tasks are performed by 1 model

T5: Text-To-Text Transfer Transformer

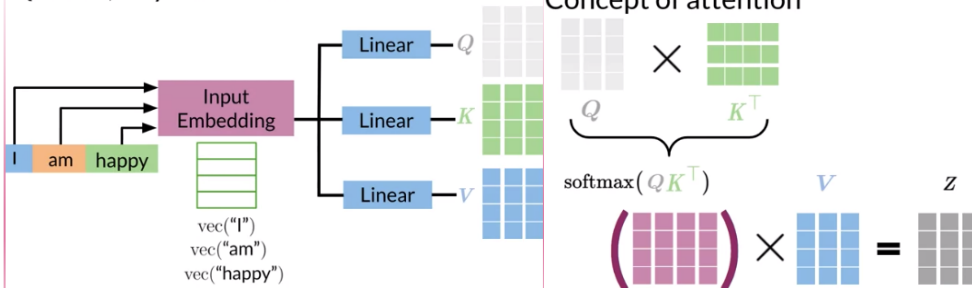


T5 Trivia demo

Dot Product Attention

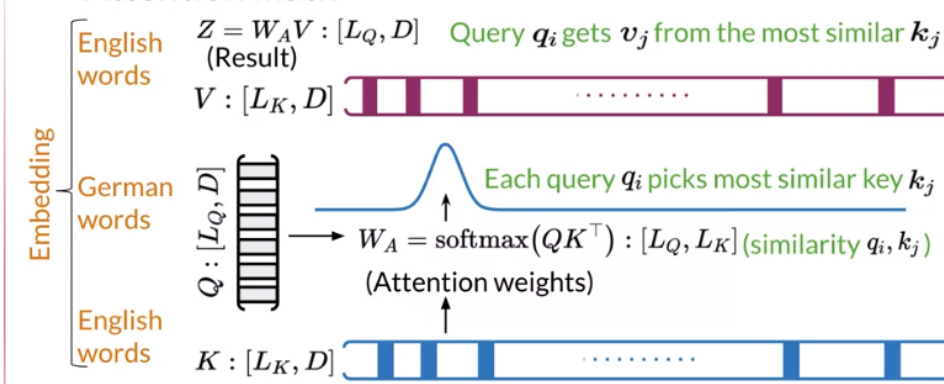


Queries, Keys and Values

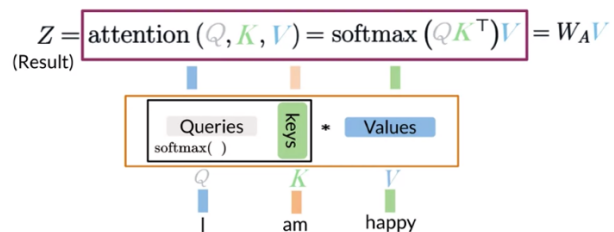


From both the capital Q matrix and the capital K matrix and attention model calculates weights or scores representing the relative importance of the keys for a specific query. These attention weights can be understood as alignments course as they come from a dot product. Additionally, to turn these weights into probabilities, a softmax function is required. Finally, multiplying these probabilities with the values, you will then get a weighted sequence, which is the attention results itself.

Attention math



Attention formula



Summary

- Dot-product Attention is essential for Transformer
- The input to Attention are queries, keys, and values
- A softmax function makes attention more focused on best keys
- GPUs and TPUs is advisable for matrix multiplications

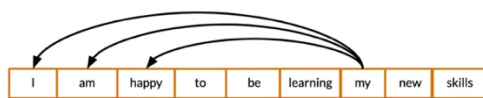
Causal Attention

Causal attention is the 2nd type of attention.

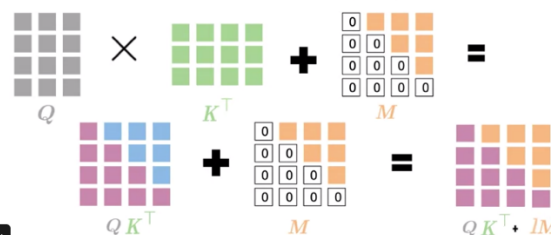
1. Encoder/Decoder attention – one sentence looks at another one.
2. Causal (self) attention – Words look at previous words (used for generation)
3. Bi directional self attention – Words look at prev and future words.
4. Queries and keys come from the same sentence and queries search among words before only.

Causal attention

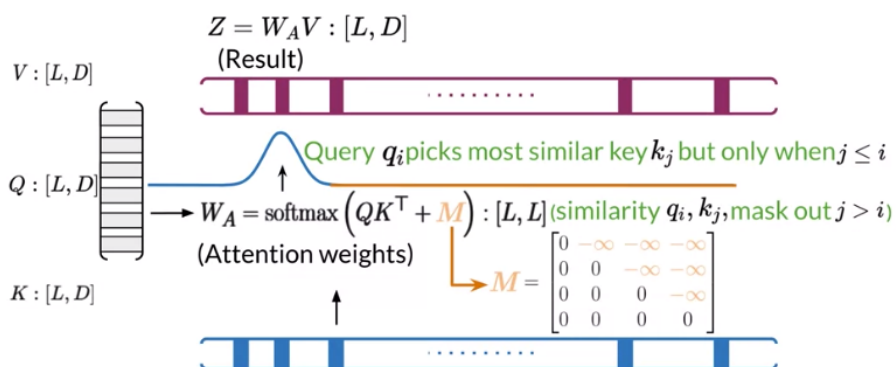
- Queries and keys are words from the same sentence
- Queries should only be allowed to look at words before



Causal attention math



Causal attention math

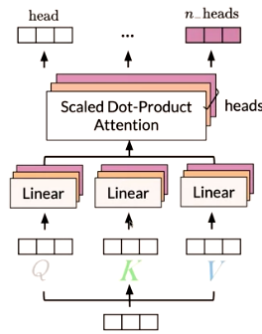


Multi-head Attention

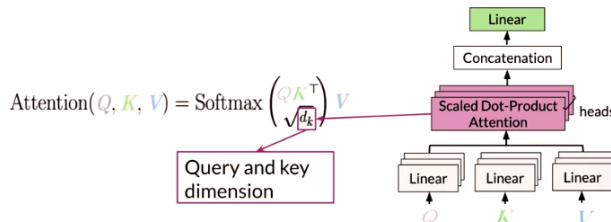
- * Different heads can learn diff relationship between words
- * Scaled dot product is adequate for Multi-Head attention
- * Multi headed models attend to info from diff representations at diff positions.

Multi-Head Attention

- Each head uses different linear transformations to represent words
- Different heads can learn different relationships between words

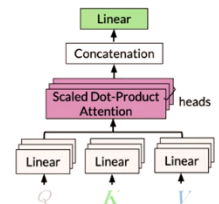


Multi-Head Attention - Scaled dot product



Multi-Head Attention - Concatenation

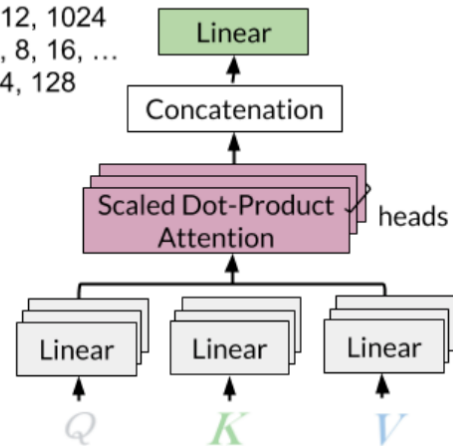
- Input(Q, K, V): [batch, length, d_model]
- Linear layer: [batch, length, n_heads * d_head]
- Transpose: [batch, n_heads, length, d_head]
- Apply attention treating n_heads like batch
- Result shape: [batch, n_heads, length, d_head]
- Transpose: [batch, length, n_heads * d_head]
- Linear layer into: [batch size, length, d_model]



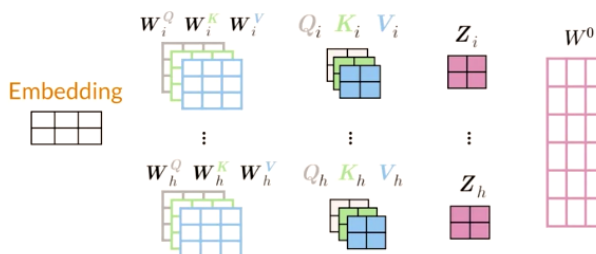
Multi-head causal attention

The layers and array dimensions involved in multi-head causal attention (which looks at previous words in the input text) are summarized in the figure below:

- Input(Q, K, V): [batch, length, d_model] ← 512, 1024
- Linear layer: [batch, length, n_heads * d_head] ← 4, 8, 16, ...
- Transpose: [batch, n_heads, length, d_head] ← 64, 128
- Apply attention treating n_heads like batch
- Result shape: [batch, n_heads, length, d_head]
- Transpose: [batch, length, n_heads * d_head]
- Linear layer into: [batch size, length, d_model]



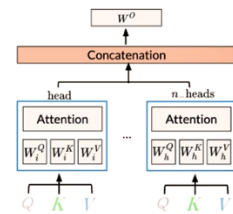
Multi-Head Attention math



Multi-Head Attention Formula

MultiHead(Q, K, V) = Concat(h₁, ..., h_n)W⁰
where h_i = Attention(QW_i^Q, KW_i^K, VW_i^V)

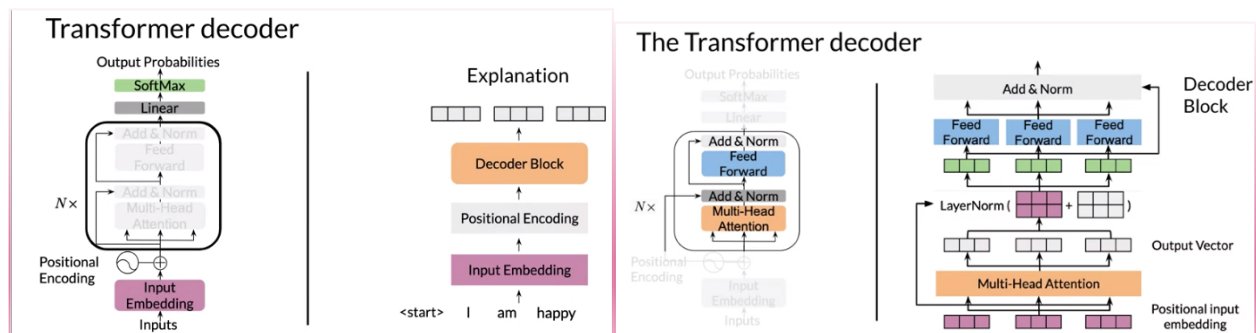
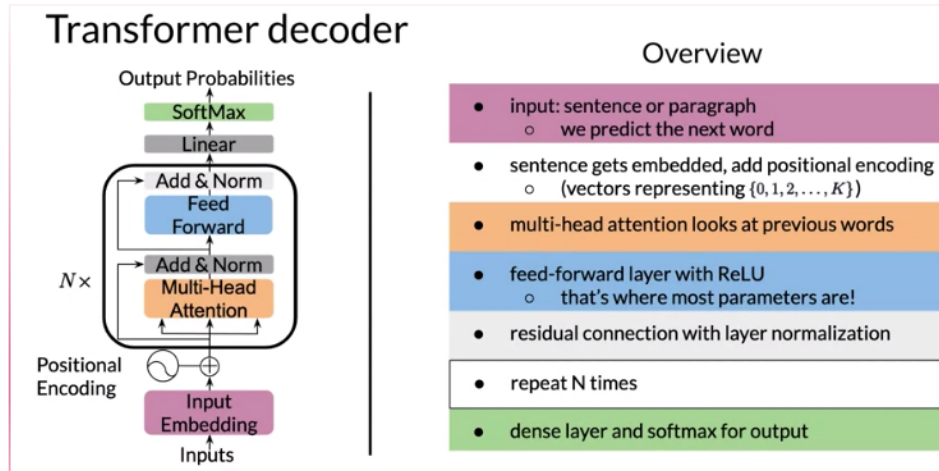
Each head h_i is the attention function of Query, Key and Value with trainable parameters (W_i^Q, W_i^K, W_i^V)



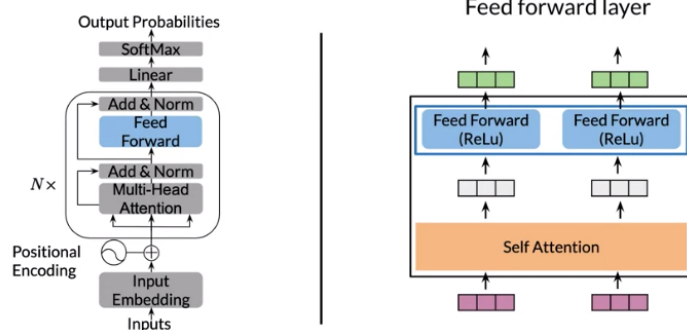
Transformer Decoder

* Transformer decoder consists of 3 layers

- * Decoder and feed forward blocks are the core of this model code
- * Also includes a module to calculate the cross-entropy loss.
- * For summarization, a weighted loss function is optimized.
- * Transformer Decoder summarizes predicting the next word using
- * Transformer uses tokenized versions of the input.

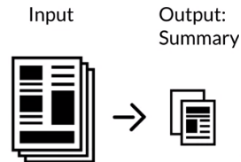
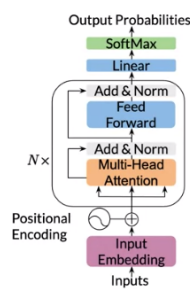


The Transformer decoder

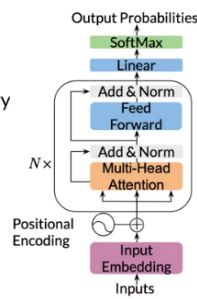


Transformer Summarizer

Transformer for summarization



Technical details for data processing



Model Input:

ARTICLE TEXT <EOS> SUMMARY <EOS> <pad> ...

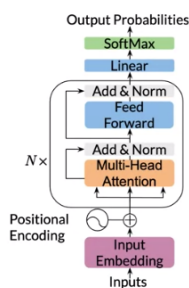
Tokenized version:

[2, 3, 5, 2, 1, 3, 4, 7, 8, 2, 5, 1, 2, 3, 6, 2, 1, 0, 0]

Loss weights: 0s until the first <EOS> and then 1 on the start of the summary.

* Use a weighted loss. Weight loss of words in articles with 0, and ones within summary with 1. Model only focuses on the summary, if summary has little data then use non-zero values for articles.

Cost function



Cross entropy loss

$$J = -\frac{1}{m} \sum_j^m \sum_i^K y_j^i \log \hat{y}_j^i$$

j : over summary

i : batch elements



Inference with a Language Model

Model input:

[Article] <EOS> [Summary] <EOS>

Inference:

- Provide: [Article] <EOS>
- Generate summary word-by-word
 - until the final <EOS>
- Pick the next word by random sampling
 - each time you get a different summary!

