Transformer Network



Learning Objectives:

- Create positional encodings to capture sequential relationships in data
- Calculate scaled dot-product self-attention with word embeddings
- Implement masked multi-head attention
- Build and train a Transformer model
- Fine-tune a pre-trained transformer model for Named Entity Recognition
- Fine-tune a pre-trained transformer model for Question Answering
- Implement a QA model in TensorFlow and PyTorch
- Fine-tune a pre-trained transformer model to a custom dataset
- Perform extractive Question Answering

Transformer Network

- 1 Transformers Intuition
- 2 Self-Attention
- 3 Multi-Head Attention
- 4 Transformers





Transformer Network

Transformers Intuition

Transformers Motivation



- The development of Transformers was motivated by the limitations of traditional sequence-to-sequence models in capturing long-term dependencies in input sequences due to the vanishing gradient problem.
- Transformers use self-attention mechanisms to attend to different parts of the input sequence at different time steps, enabling them to capture longterm dependencies more effectively and leading to improved performance in natural language processing tasks.

Transformers Motivation



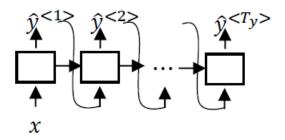
- The transformer architecture is very popular in the field of natural language processing (NLP).
- It is a complex architecture, but it allows for parallel processing of entire sequences, rather than ingesting input one word or token at a time. This is in contrast to RNNs, GRUs, and LSTMs, which are all sequential models that process input one unit at a time. As the complexity of sequence tasks increases, so does the complexity of the model.
- The major innovation of the transformer architecture is the combination of attention-based representations and a CNN-style of processing.

Transformers Motivation

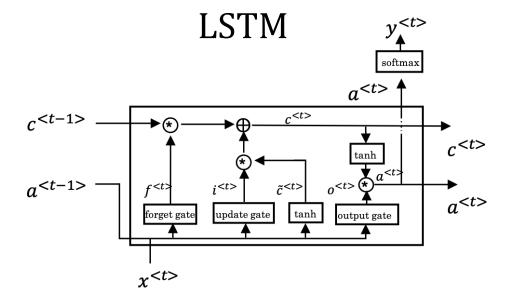




RNN



GRU



Transformers Intuition



- The attention network is a way of computing rich and useful representations of words in parallel using self-attention and multi-headed attention.
- The transformer architecture overcomes the limitations of traditional sequence-to-sequence models in capturing long-term dependencies.
- The next section will cover the self-attention and multi-headed attention processes in detail and put all the pieces together to give a complete understanding of how the transformer works.

Transformers Intuition



- Attention + CNN
 - Self-Attention
 - Multi-Head Attention



Sequence to sequence models

Self-Attention

Self-Attention Intuition



- The attention values are then computed using the dot product of the query and key, and scaled by the square root of the dimension of the key.
- The resulting attention weights are used to compute a weighted sum of the values, which gives the attention-based representation for each word.
- The self-attention mechanism allows for the parallel computation of these representations, making it more efficient than traditional sequence-to-sequence models.

Self-Attention Intuition



• A(q, K, V)= attention-based vector representation of a word

RNN Attention

$$\alpha^{< t, t'>} = \frac{\exp(e^{< t, t'>})}{\sum_{t=1}^{T_{,x}} \exp(e^{< t, t'>})}$$

Transformers Attention

$$A(q, K, V) = \sum_{i} \frac{\exp(e^{\langle q \cdot k < i \rangle})}{\sum_{j} \exp(e^{\langle q \cdot k} \langle j \rangle)} v^{\langle i \rangle}$$

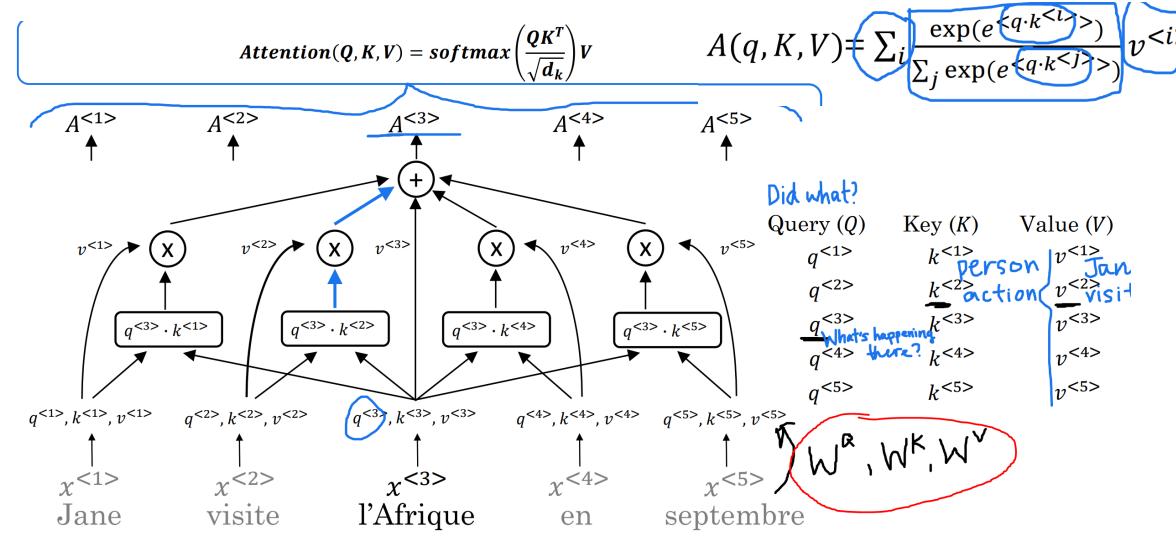
Self-Attention



- The value of each word in a sentence can be plugged into the attentionbased representation for the sentence, allowing for a richer and more nuanced representation of words than fixed word embeddings.
- This self-attention mechanism can be used to create attention-based representations for all the words in a sentence, which can be combined in a compressed or vectorized representation using the scaled dot-product attention mechanism.
- The multi-headed attention mechanism adds a loop over the self-attention mechanism, creating multiple attention-based representations for each word and allowing for even more complex and nuanced representations of input sequences.

Self-Attention





'aswani et al. 2017, Attention Is All You Need]

Andrew Ng



Sequence to sequence models

Multi-Head Attention

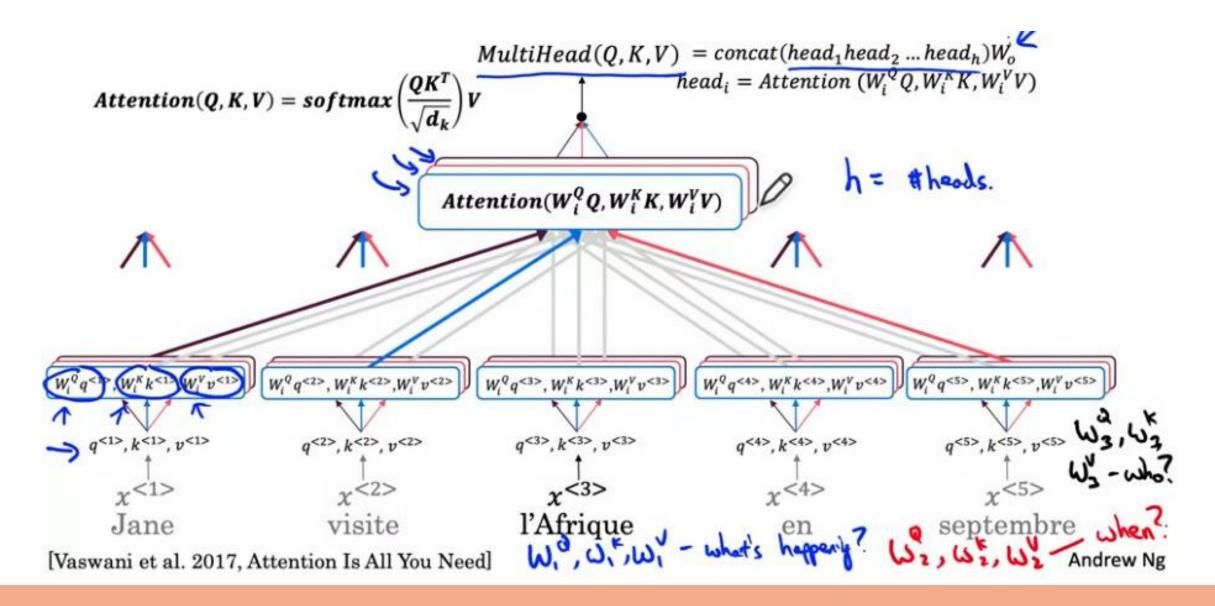
Multi-Head Attention



- In multi-head attention mechanism, the input is multiplied with weight matrices to obtain query, key, and value vectors.
- Self attention is computed multiple times, each time using a different set of weight matrices referred to as "heads". Each head represents a different feature and asks a different question about the input sequence.
- The output of each head is concatenated and multiplied by a final weight matrix to produce the final output. This computation can be done in parallel, making it computationally efficient and suitable for use in large-scale deep learning models such as the transformer network.

Multi-Head Attention







Sequence to sequence models

Transformers

Transformer Details



- The Transformer network is an architecture that uses self-attention and multi-headed attention mechanisms to perform sequence-tosequence tasks, such as machine translation.
- The network has an encoder and decoder block, each with multiple layers.
- The encoder produces a contextualized embedding for each word using multi-headed attention and feed-forward neural networks.
- The decoder generates the translated output using a similar multiheaded attention mechanism, generating the output one word at a time.
- Additional features such as positional encoding, residual connections, and masking during training are used to improve performance. Variations of the model, such as BERT and BERT distill, have also been proposed.

Transformer Details



<SOS> Jane visits Africa in September <EOS> Softmax Linear Encoder Decoder Add & Norm Feed Forward Add & Norm Neural Network Feed Forward i = 1**Neural Network** Add & Norm i=2 N_{x} Multi-Head Add & Norm 1=3 Attention Multi-Head 005 Attention Add & Norm Multi-Head Attention Positional Encoding $PE_{(pos,2i)} = sin(\frac{pos}{1000^{\frac{2i}{d}}})$ <SOS> $x^{<1>}$ $x^{<2>}$... $x^{<T_x-1>}$ $x^{<T_x>}$ <EOS> Jane visite l'Afrique en septembre $PE_{(pos,2i+1)} = cos(\frac{pos}{2i})$ $\langle SOS \rangle y^{<1}\rangle y^{<2}\rangle \dots y^{<T_y-1}\rangle y^{<T_y>}$ $\langle SOS \rangle$ Jane visits Africa in September

[Vaswani et al. 2017, Attention Is All You Need]

Andrew Ng

Summarization



- Transformers, introduced in "Attention is All You Need," use self-attention for sequential data.
- Self-attention allows dynamic weighing of words based on relationships.
 Multi-head attention enhances this by focusing on different parts simultaneously.
- Transformers consist of encoder and decoder layers with self-attention and feed-forward networks. They excel in NLP, surpassing RNNs and CNNs due to parallelization and long-range dependency handling.