

Attention all you need paper

Zaina Almarede 202010961

November 2023

0.1 Summary

The paper titled "Attention Is All You Need" introduces the Transformer, a novel neural network architecture for sequence transduction tasks, such as machine translation. Unlike traditional models that rely on recurrent or convolutional layers, the Transformer is built solely on attention mechanisms. The authors demonstrate the superiority of the Transformer through experiments on machine translation tasks, achieving improved quality, enhanced parallelizability, and significantly reduced training time.

The Transformer's key innovation is the use of self-attention, allowing the model to capture global dependencies between input and output sequences without the need for sequential computation. The architecture consists of encoder and decoder stacks, each comprising multiple layers with attention and feed-forward sub-layers. Multi-head attention is employed to enable the model to attend to different representation subspaces simultaneously.

The authors compare the Transformer with other models, emphasizing its computational efficiency and ability to handle long-range dependencies. They also introduce positional encodings to enable the model to consider the order of sequence elements. The training process involves the WMT 2014 English-German data-set, and the authors achieve state-of-the-art results in both English-German and English-French translation tasks.

In conclusion, the Transformer architecture demonstrates remarkable performance improvements in sequence transduction tasks, setting new standards in translation quality while being more parallelizable and efficient in training compared to existing models. The paper presents a significant advancement in neural network architectures for sequence-to-sequence tasks.