

---

## Paper:2

1. Title: Attention Is All You Need
2. Authors: Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, Illia Polosukhin
3. Affiliation: Google Brain (Ashish Vaswani, Noam Shazeer, Jakob Uszkoreit, Llion Jones, Łukasz Kaiser), Google Research (Niki Parmar, Illia Polosukhin), University of Toronto (Aidan N. Gomez)
4. Keywords: attention mechanism, machine translation, neural networks, sequence transduction, Transformer architecture
5. Urls: Arxiv link: <https://arxiv.org/abs/1706.03762>
6. Github: None
7. Summary:
  - (1): The research background of this article is on sequence transduction models based on recurrent or convolutional neural networks that use an encoder and decoder, and with the best performing models using an attention mechanism to connect them.
  - (2): The past methods involved complex and computationally expensive models with recurrent or convolutional networks, which have limitations in scaling and parallelizing. The approach of using solely attention mechanisms is well motivated due to its success in past applications and simpler implementation.
  - (3): The research methodology proposed in this paper is a new network architecture called the Transformer, which only uses attention mechanisms without recurrence or convolution. The architecture has a vastly reduced training time and is more parallelizable, and generalizes well to other tasks.
  - (4): The methods in this paper achieve state-of-the-art performance on two machine translation tasks, with BLEU scores of 28.4 and 41.8, surpassing previous results, including using ensembles, while requiring significantly less training time. The performance supports the goals of providing a more efficient and scalable solution to sequence transduction models.
8. Methods:
  - (1): The proposed method in this paper introduces a new network architecture called the Transformer, which utilizes attention mechanisms solely without recurrence or convolution, making it more efficient and parallelizable than previous models.
  - (2): The Transformer employs a technique known as self-attention, allowing it to capture dependencies across an input sequence, and uses multi-head attention to improve performance and allow for scaling to larger datasets.

- 
- (3): The model is trained using standard supervised learning with a maximum likelihood objective and incorporates techniques such as label smoothing to improve generalization.
  - (4): The effectiveness of the method is demonstrated through its state-of-the-art performance on two machine translation tasks, with significantly less training time and improved scalability compared to previous models.

#### 9. Conclusion:

- (1): The significance of this piece of work is that it introduces a new network architecture called the Transformer, which utilizes attention mechanisms solely without recurrence or convolution. This makes it more efficient and parallelizable than previous models in the field of sequence transduction.
- (2): Innovation point: The innovation point of this article lies in the implementation of attention mechanisms solely without recurrence or convolution in the Transformer architecture, which significantly reduces its training time and improves scalability. Performance: This method achieves state-of-the-art performance on two machine translation tasks, with BLEU scores of 28.4 and 41.8, surpassing previous models and even ensembles while requiring significantly less training time. Workload: The workload in terms of training time and scalability is improved greatly with this approach. However, the text-only nature of the model could be a weakness for addressing tasks involving other modalities such as image, audio, or video data.