# P1: Introduction to Transformers

Rajdeep Roshan Sahu
230215119
r.sahu@se23.qmul.ac.uk

## 1. Introduction

Transformer is a type of deep learning algorithm that has revolutionized machine learning, particularly natural language processing (NLP), with its unique attention mechanism, since it was first introduced in the paper,**" Attention is all you need."** This model surpasses traditional RNNs and CNNs by enabling parallel processing of sequences and superior understanding of long-range dependencies. Its innovative approach has also been successfully adapted to computer vision, as highlighted in **"An Image is Worth 16x16 Words,"** demonstrating its versatility across domains.

## 2. Transformers

**Architecture and Working:** Transformers use a self-attention mechanism to process data, with structures consisting of an encoder and a decoder, each with six layers. Encoders use multi-head self-attention and feed-forward networks to transform input data into context-rich representations, which are then improved further by normalization and residual connections. Decoders, which have an almost similar architecture, include an additional multi-head attention layer that focuses on critical input for output generation. Positional encoding maintains sequence order, which addresses the self-attention difficulty with sequential data. This approach allows for efficient, context-sensitive processing, particularly in activities like translation, by leveraging parallel processing and strong contextual knowledge. The attention function can be summarized by the equation:

**Attention(Q,K,V) = softmax[(Q* K ^T) / sqrt(dk)] * V**

Where, Q represents the query matrix, K represents the key matrix, V represents the value matrix, dk is the dimension of the key vectors that is used to scale the dot products. In contrast to RNNs and LSTMs, transformers handle complete data sequences in parallel. Multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions. This is achieved by performing the attention function in parallel for each "head," and then concatenating and linearly transforming the results.

**Limitations:** Despite their parallel processing capability, transformers have limitations including quadratic computational complexity with respect to input sequence length, making processing long sequences resource intensive. Additionally, storing entire sequence activations concurrently for parallel processing increases memory requirements. Transformers may struggle with generalizing to new, unseen data, particularly when training data lacks diversity. Moreover, they demand substantial data for effective training, limiting practicality in data-scarce tasks.

## 3. Vision Transformer (ViT)

**Architecture and Working:** The Vision Transformer (ViT) adapts the Transformer architecture, originally designed for natural language processing (NLP), to perform computer vision tasks. Its architecture is like the standard Transformer, with the primary adaptation being how the input data is prepared and fed into the model. Images are split into fixed-size patches (16x16 pixels), which are then flattened and linearly projected into embeddings. These embeddings, along with positional encodings to retain positional information, serve as the input to the Transformer encoder. The ViT encoder consists of multiple layers of multi-head self-attention and feed-forward networks, with normalization and residual connections. The intuition behind ViT is that the self-attention mechanism can learn to focus on relevant parts of an image, like how it identifies important words or phrases in a text sequence. This approach requires minimal changes to the Transformer encoder itself, demonstrating the applicability of the self-attention mechanism across various domains. The mathematics behind the ViT is same as in the standard transformer.

**Limitations:** Processing images with ViT, especially high-resolution images divided into many small patches, can be computationally expensive due to

the quadratic complexity of the self-attention mechanism. Moreover, ViT models tend to require large amounts of training data to achieve performance comparable to CNNs. Also, adapting ViTs to vision tasks other than classification (e.g., object detection, segmentation) requires additional architecture modifications.

# 4. Advancements in Transformers

**1) SparseMOE** (Sparse Mixture of Experts): is a Large Language Model (LLM) architecture by Google AI, optimizes efficiency by employing a mixture of smaller expert models rather than a singular large model. It utilizes a **'gating'** technique to identify and activate only the expert models deemed most relevant for each token in the input sequence, significantly cutting down on the required computations, especially for longer sequences. This selective activation approach allows SparseMOE to achieve a **lower computational cost and memory footprint**, while maintaining or enhancing accuracy across a variety of tasks, compared to conventional LLMs.

**2) GPT-3** (Generative Pre-trained Transformer 3): is a highly capable large language model renowned for producing human-like text, language translation, creative content creation, and delivering informative answers. It is pre-trained on extensive text and code datasets through unsupervised learning, laying a robust foundation for subsequent fine-tuning on specific tasks with limited data. This approach significantly **conserves time and resources by eliminating the need to train a model from scratch for every new task**. GPT-3's pre-trained knowledge enables it to adapt to a wide range of tasks, enhancing its versatility.

# 5. Advancements in Vision Transformers

**1) Swin Transformer** (Liu et al., 2021): The Swin transformer optimizes computational efficiency through shifted window-based attention. It employs hierarchical representation to reduce token count at higher levels, focusing on local windows to minimize attention computations. Window shifting enhances long-range dependency capture without full self-attention. Combined with standard Transformer features like residual connections and feedforward networks, Swin **achieves both high efficiency and effectiveness in processing**, particularly suited for images with its pyramid-like approach to **managing data complexity.**

**2) MoCo v3** (Chen et al., 2021): enhances Vision Transformer (ViT) training through contrastive learning and data augmentation, effectively **dealing with limited labeled data**. This method distinguishes between similar and dissimilar image pairs, encouraging the model to learn robust, invariant features through varied augmented images. Such augmentation increases data diversity, helping the model become more adaptable to unseen data variations. By applying this strategy, **MoCo v3 allows ViTs to achieve high performance, comparable to models trained on extensive datasets like ImageNet**, without the necessity for large volumes of labeled data.

# 6. Suggestions

Given the transformative role of transformers in natural language processing and computer vision, future research should focus on **developing innovative architectures that lower computational costs** without sacrificing accuracy and improve data efficiency to reduce dependency on large, labeled datasets. Additionally, **enhancing transfer learning** to boost model adaptability across different domains is vital for advancing transformer technology. Here are some of my suggestions aimed at leveraging various transformer architectures for various applications:
**1) Multimodal Learning**: Combining Transformers and ViTs with other modalities can be powerful. **2) Automated Medical Diagnosis:** ViTs, combined with other deep learning models, could analyze medical images for early disease detection and diagnosis. **3) Climate Change Mitigation:** Both Transformers and ViTs could be used for tasks like weather forecasting, optimizing renewable energy.

# 7. Conclusion

In conclusion, transformers have revolutionized machine learning, transforming natural language processing and computer vision. Their potential, however, extends far beyond these initial applications. By exploring new directions, addressing ethical concerns, and harnessing their power for good, transformers can make significant contributions to various fields and improve our lives in remarkable ways. As this technology continues to evolve, promoting responsible development and thoughtful implementation will be crucial to ensure a future where Transformers benefit all of humanity.

## 8. References

[1] Vaswani, Ashish, et al. "Attention is all you need." Advances in Neural Information Processing Systems, 2017, pp. 5998-6008.

[2] Dosovitskiy, Alexey, Lucas Beyer, Jan Khosla, Amaia Gordo, David Rodríguez, and Rob Fergus. 2020. "An image is worth 16x16 words: Transformers for image recognition at scale." In International Conference on Learning Representations (ICLR).

[3] Alistar, David, Armand Joulin, et al. 2023. "Sparsely MoE: Model-parallel methods for efficiently training multi-expert models." In Advances in Neural Information Processing Systems (NeurIPS) 36.

[4] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). Language Models are Few-Shot Learners.

[5] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., & Guo, B. (2021). Swin Transformer: Hierarchical Vision Transformer using Shifted Windows.

[6] Chen, X., Xie, S., & He, K. (2021). An Empirical Study of Training Self-Supervised Vision Transformers.