# Briefing Document: "Attention Is All You Need"

## Executive Summary

The paper "Attention Is All You Need" introduces the Transformer, a novel neural network architecture for sequence transduction tasks like machine translation. The key innovation is replacing traditional recurrent and convolutional layers with a mechanism called "attention," specifically "multi-head self-attention." This allows the model to process sequences in parallel, significantly improving training speed and achieving state-of-the-art performance on machine translation and constituency parsing tasks with significantly less computational cost. The Transformer architecture consists of stacked encoder and decoder layers, each utilizing self-attention and position-wise feed-forward networks, along with positional encodings to capture sequence order.

## Key Concepts and Themes

- **Dispensing with Recurrence and Convolutions:** The core idea is to build a sequence transduction model "based solely on attention mechanisms, dispensing with recurrence and convolutions entirely." This is a significant departure from the dominant architectures at the time.
- **Attention as the Primary Mechanism:** The Transformer relies entirely on attention mechanisms to "draw global dependencies between input and output." Attention functions map a query and a set of key-value pairs to an output, representing a weighted sum of the values based on the compatibility of the query with the keys.
- **Self-Attention:** Also known as intra-attention, self-attention is used to relate different positions of a single sequence to compute a representation of that sequence. In the Transformer, self-attention is applied within both the encoder and decoder.
- **Scaled Dot-Product Attention:** This is the specific attention function used in the Transformer. It computes dot products of the query with all keys, scales them by the square root of the key dimension, and applies a softmax function to get weights for the values.
- **Multi-Head Attention:** Instead of a single attention function, the model uses multiple "heads" running in parallel. This allows the model to "jointly attend to information from different representation subspaces at different positions."
- **Encoder-Decoder Architecture:** Like many sequence transduction models, the Transformer uses an encoder-decoder structure. The encoder processes the input sequence, and the decoder generates the output

sequence autoregressively, attending to both the encoder output and previously generated tokens.

- **Positional Encoding:** Since the Transformer lacks recurrence or convolutions, it needs a way to incorporate information about the order of tokens. This is achieved by adding "positional encodings" to the input embeddings. The paper uses sinusoidal functions of different frequencies for this purpose.
- **Parallelizability and Reduced Training Time:** A major advantage of the Transformer is its ability to perform computations in parallel, which "precludes parallelization within training examples" in recurrent models. This leads to "significantly more parallelization and can reach a new state of the art in translation quality after being trained for as little as twelve hours on eight P100 GPUs."
- **Improved Performance on Machine Translation:** The Transformer achieved state-of-the-art BLEU scores on the WMT 2014 English-to-German and English-to-French translation tasks, significantly outperforming previous models, including ensembles, with lower training costs.
- **Generalization to Other Tasks:** The paper demonstrates that the Transformer "generalizes well to other tasks by applying it successfully to English constituency parsing."

## Most Important Ideas and Facts

- **The Transformer Architecture:** The fundamental contribution is the introduction of the Transformer, a sequence transduction model "based solely on attention mechanisms."
- **Multi-Head Self-Attention as the Core Building Block:** The key novelty and source of the model's power is the use of "stacked self-attention and point-wise, fully connected layers" in both the encoder and decoder.
- **Superior Performance and Efficiency:** The paper empirically demonstrates that the Transformer achieves higher translation quality and requires significantly less training time compared to previous state-of-the-art models. "Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU."
- **Parallel Processing Capability:** The architecture's design allows for parallel computation, which is crucial for training efficiency, especially on

long sequences. "The Transformer allows for significantly more parallelization."

- **Positional Encoding is Necessary:** Explicitly adding positional information is critical for the model to handle sequence order in the absence of recurrence or convolutions.
- **Scaled Dot-Product Attention with Multi-Head Implementation:** The specific implementation of attention using scaled dot-product and multiple heads is a key technical detail contributing to the model's effectiveness.
- **Encoder and Decoder Structure:** The model follows the established encoder-decoder framework, but the internal layers are radically different.
- **Successful Application to Parsing:** The Transformer's success on English constituency parsing demonstrates its versatility beyond machine translation.
- **Open Source Code:** The authors released their code, making the model accessible for further research and development. "The code we used to train and evaluate our models is available at https://github.com/tensorflow/tensor2tensor."

## Supporting Details and Quotes

- **Problem Addressed:** Recurrent models' "inherently sequential nature precludes parallelization within training examples, which becomes critical at longer sequence lengths."
- **Attention as a Solution to Long-Range Dependencies:** Attention mechanisms "allowing modeling of dependencies without regard to their distance in the input or output sequences."
- **Transformer's Advantage in Path Length:** "In the Transformer this is reduced to a constant number of operations," compared to linear or logarithmic growth in convolutional models.
- **Scaled Dot-Product Attention Justification:** "While for small values of dk the two mechanisms perform similarly, additive attention outperforms dot product attention without scaling for larger values of dk [3]. We suspect that for large values of dk, the dot products grow large in magnitude, pushing the softmax function into regions where it has extremely small gradients."
- **Purpose of Multi-Head Attention:** "Multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions. With a single attention head, averaging inhibits this."

- **Applications of Attention in the Transformer:** The paper details the three ways multi-head attention is used: encoder-decoder attention, encoder self-attention, and decoder self-attention (with masking).
- **Model Architecture Details:** The encoder and decoder each have N=6 identical layers, with residual connections and layer normalization. The embedding dimension is dmodel = 512 (base model).
- **Training Details:** The paper provides specifics on training data, batching, hardware (8 NVIDIA P100 GPUs), training schedule, optimizer (Adam), and regularization (residual dropout, label smoothing).
- **Experimental Results:** Table 2 clearly shows the BLEU scores and training costs comparing the Transformer to previous models. "The Transformer (big) 28.4 41.8 2.3 · 1019" (demonstrating high performance and relatively low cost).
- **Interpretability:** The authors suggest that self-attention "could yield more interpretable models," noting that "individual attention heads clearly learn to perform different tasks, many appear to exhibit behavior related to the syntactic and semantic structure of the sentences." (See Figures 3, 4, and 5 in the appendix).
- **Future Work:** The authors outline plans to apply the Transformer to other modalities and investigate restricted attention mechanisms.

---

# What is the main innovation of the Transformer model?

The Transformer model represents a significant departure from previous dominant sequence transduction models, which were based on complex recurrent or convolutional neural networks. The core innovation of the Transformer is that it **eschews recurrence and convolutions entirely, relying solely on attention mechanisms** to model dependencies between input and output sequences. This fundamentally changes how the model processes sequential data.

# How does the Transformer achieve better parallelization compared to RNNs?

Recurrent neural networks (RNNs) process information sequentially, aligning computation with the positions in the input and output sequences. This inherent sequential nature limits parallelization within a single training example. The Transformer, by contrast, uses attention mechanisms to relate all positions in the sequence to each other simultaneously. This allows for **significantly more parallelization during training**, which is crucial for handling longer sequences and large datasets.

# What is the role of the encoder and decoder in the Transformer architecture?

Like most competitive neural sequence transduction models, the Transformer follows an **encoder-decoder structure**.

- The **encoder** takes an input sequence of symbol representations and maps it to a sequence of continuous representations. It is composed of a stack of identical layers, each containing a multi-head self-attention mechanism and a position-wise fully connected feed-forward network.
- The **decoder** takes the output of the encoder and generates an output sequence of symbols one element at a time. It is also composed of a stack of identical layers, which include the same sub-layers as the encoder layers, but with an additional multi-head attention layer that attends over the output of the encoder stack.

# How does the Transformer handle the order of the sequence without recurrence or convolutions?

Since the Transformer model does not use recurrence or convolution, which inherently process sequential order, it must explicitly inject information about the position of tokens. This is achieved through the use of **positional encodings**. These encodings are added to the input embeddings at the bottom of both the encoder and decoder stacks. The paper uses sinusoidal functions of different frequencies to create these positional encodings, allowing the model to learn to attend to relative positions.

# What is the difference between Scaled Dot-Product Attention and Multi-Head Attention?

- **Scaled Dot-Product Attention** is the specific attention function used within the Transformer. It takes a query, a set of keys, and a set of values as input. It computes the dot products of the query with all keys, scales the results, and applies a softmax function to get weights on the values. The output is a weighted sum of the values.
- **Multi-Head Attention** extends Scaled Dot-Product Attention. Instead of performing a single attention function, it linearly projects the queries, keys, and values multiple times (h times) with different learned linear projections. Attention is then performed in parallel on these projected versions, yielding h separate outputs. These outputs are concatenated and linearly projected again to produce the final output. This allows the model to attend to information from different representation subspaces at different positions.

# How is attention used in different parts of the Transformer model?

The Transformer utilizes multi-head attention in three distinct ways:

- **Encoder-Decoder Attention:** Used in the decoder, queries come from the previous decoder layer, and keys and values come from the output of the encoder. This allows the decoder to attend over the entire input sequence.
- **Encoder Self-Attention:** Used within the encoder layers, queries, keys, and values all come from the output of the previous layer in the encoder. This allows each position in the encoder to attend to all positions in the previous layer.
- **Decoder Self-Attention:** Used within the decoder layers, queries, keys, and values all come from the output of the previous layer in the decoder. A mask is applied to prevent positions from attending to subsequent positions, ensuring the auto-regressive property where predictions for a position only depend on known outputs at preceding positions.

# What are the key advantages of self-attention compared to recurrent and convolutional layers for sequence transduction?

The paper highlights several advantages:

- **Reduced Sequential Operations:** A self-attention layer connects all positions with a constant number of sequential operations ($O(1)$), while a recurrent layer requires $O(n)$ sequential operations, where n is the sequence length.
- **Shorter Maximum Path Length:** Self-attention provides a constant maximum path length between any two input and output positions ($O(1)$), facilitating the learning of long-range dependencies compared to RNNs ($O(n)$) or convolutional networks ($O(\log k(n))$ or $O(n/k)$).

- **Computational Efficiency (in certain cases):** Self-attention layers can be faster than recurrent layers when the sequence length is smaller than the representation dimensionality, which is common in modern machine translation models.

## How did the Transformer perform on machine translation and other tasks?

The Transformer demonstrated **superior performance on machine translation tasks**, achieving new state-of-the-art BLEU scores on both WMT 2014 English-to-German and English-to-French datasets. Notably, the "big" Transformer model surpassed existing best results, including ensembles, with significantly lower training costs. The paper also shows that the Transformer generalizes well to **English constituency parsing**, yielding competitive results even with limited task-specific tuning, outperforming some previous models trained on larger datasets.

Transformer Model Study Guide

Quiz

1. What is the fundamental difference between the Transformer architecture and dominant sequence transduction models that preceded it?
2. What are the three main components of the Transformer model architecture?
3. Describe the function of the encoder in the Transformer model.
4. What is the purpose of the masking applied to the self-attention sub-layer in the decoder?
5. How is the output from the Scaled Dot-Product Attention computed?
6. What is the primary benefit of using Multi-Head Attention?
7. Explain the purpose of positional encodings in the Transformer.
8. What are the three desiderata considered when comparing self-attention to recurrent and convolutional layers?
9. How did the Transformer perform on the WMT 2014 English-to-German and English-to-French translation tasks compared to previous models?
10. What other task, besides machine translation, did the authors successfully apply the Transformer to?

Answer Key

1. The Transformer is based solely on attention mechanisms, dispensing with recurrence and convolutions entirely, whereas previous dominant models relied on complex recurrent or convolutional neural networks.
2. The three main components are the encoder, the decoder, and the attention mechanisms (specifically multi-head self-attention, encoder-decoder attention, and position-wise feed-forward networks).
3. The encoder maps an input sequence of symbol representations to a sequence of continuous representations.
4. The masking prevents positions in the decoder from attending to subsequent positions, preserving the auto-regressive property where predictions for a position can only depend on known outputs at previous positions.
5. The output is computed as a weighted sum of the values, where the weights are obtained by taking the softmax of the dot product of the query with all keys, scaled by the square root of the key dimension.
6. Multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions, overcoming the averaging inhibition of a single attention head.
7. Positional encodings inject information about the relative or absolute position of tokens into the sequence since the model contains no recurrence or convolution to inherently capture this order.
8. The three desiderata are total computational complexity per layer, the amount of computation that can be parallelized (minimum number of sequential operations), and the path length between long-range dependencies.
9. The Transformer achieved a new state of the art on both tasks, outperforming previously reported models and ensembles with significantly lower training costs.
10. The authors successfully applied the Transformer to English constituency parsing.

Essay Questions

1. Discuss the advantages of the Transformer's attention-only architecture over traditional recurrent and convolutional models for sequence transduction tasks, focusing on computational efficiency, parallelization, and handling long-range dependencies.
2. Explain the different ways multi-head attention is utilized within the Transformer model architecture (encoder self-attention, decoder self-attention, and encoder-decoder attention) and the specific role each plays in the overall sequence transduction process.
3. Analyze the importance of scaled dot-product attention and multi-head attention to the performance of the Transformer, explaining the mechanics of each and why their combination is beneficial.
4. Detail the training regime described in the paper, including the data used, hardware and schedule, optimizer, and regularization techniques employed, and how these factors contributed to the model's success.
5. Evaluate the generalization capabilities of the Transformer architecture based on its performance on the English constituency parsing task, considering the specific challenges of this task and how the Transformer addresses them.

Glossary of Key Terms

- **Sequence Transduction:** Tasks that transform an input sequence into an output sequence, such as machine translation.
- **Recurrent Neural Networks (RNNs):** Neural networks that process sequences by maintaining a hidden state that is updated at each time step based on the current input and the previous hidden state.
- **Convolutional Neural Networks (CNNs):** Neural networks that use convolutional layers to process data, typically used for spatial hierarchies but adaptable to sequential data.
- **Attention Mechanism:** A mechanism that allows a model to focus on specific parts of the input sequence when processing or generating the output sequence, regardless of their distance.
- **Transformer:** A neural network architecture based entirely on attention mechanisms, eschewing recurrence and convolutions.
- **Encoder-Decoder Structure:** A common architecture for sequence transduction models where an encoder processes the input sequence and a decoder generates the output sequence.
- **Self-Attention (Intra-Attention):** An attention mechanism that relates different positions of a single sequence to compute a representation of that sequence.
- **Scaled Dot-Product Attention:** A specific attention function where the output is computed by taking the softmax of the scaled dot product of the query and keys, and then multiplying by the values.
- **Multi-Head Attention:** An attention mechanism that performs multiple attention functions in parallel with different linear projections of the queries, keys, and values, allowing the model to attend to information from different representation subspaces.
- **Position-wise Fully Connected Feed-Forward Network:** A feed-forward network applied independently and identically to each position in a sequence within the encoder and decoder layers.
- **Residual Connection:** A technique that adds the input of a layer to its output, helping to mitigate the vanishing gradient problem in deep networks.

- **Layer Normalization:** A normalization technique applied across the features within a layer for each training example.
- **Positional Encoding:** A method of injecting information about the position of tokens in a sequence into the model, typically by adding vectors to the input embeddings.
- **Auto-regressive:** A model property where the generation of an output at a given position depends on the previously generated outputs.
- **BLEU Score:** A metric used to evaluate the quality of machine translation output by comparing it to reference translations.
- **Byte-Pair Encoding (BPE):** A subword tokenization algorithm used to handle rare words and out-of-vocabulary tokens by breaking words into smaller units.
- **Word-Piece:** Another subword tokenization algorithm.
- **Dropout:** A regularization technique that randomly sets a fraction of the output units of a layer to zero during training to prevent overfitting.
- **Label Smoothing:** A regularization technique applied to the target probabilities during training, preventing the model from becoming overconfident in its predictions.
- **Beam Search:** A search algorithm used during inference for sequence generation that explores multiple promising output sequences simultaneously.
- **Constituency Parsing:** A task in natural language processing that aims to build a tree structure representing the syntactic constituents of a sentence.