# Implementation of a Transformer-Based Architecture for Semantic Paraphrasing on Resource-Constrained Hardware

Fegyó Benedek

*Department of Artificial Intelligence*

November 29, 2025

## Abstract

This paper presents an implementation of the Transformer architecture from scratch, designed to perform a semantic paraphrasing tasks. While modern Natural Language Processing (NLP) relies heavily on pre-trained large language models, this work constructs the architecture from first principles using PyTorch to verify the efficacy of the attention mechanism. The model was trained on the MS COCO captioning dataset, adapting a computer vision benchmark into a text-to-text paraphrasing task. Despite hardware constraints (RTX 2060 Mobile), we achieved stable convergence and generated semantically viable paraphrases by optimizing hyperparameters and utilizing Mixed Precision training.

## 1   Introduction

Sequence-to-sequence learning has revolutionized NLP, moving from Recurrent Neural Networks (RNNs) to the Transformer architecture, which relies entirely on self-attention mechanisms to draw global dependencies between input and output.

The objective of this project was two-fold:

1. To validate the mathematical foundations of the "Attention Is All You Need" paper by implementing the architecture from scratch without high-level wrappers (e.g., HuggingFace).

2. To deploy this model on consumer-grade hardware (6GB VRAM) for the task of sentence paraphrasing.

We demonstrate that a scaled-down Transformer, when optimized with specific memory-saving techniques, can learn syntactic structures and semantic relationships from a relatively small corpus of image captions.

## 2   Related Work

The foundational architecture is based on Vaswani et al. (2017) [1], which introduced the Multi-Head Attention mechanism. Unlike prior LSTM-based approaches (Sutskever et al., 2014), the Transformer allows for parallelization during training, significantly reducing training time.

For the dataset, we utilized the Microsoft COCO dataset (Lin et al., 2014) [2]. While primarily a computer vision benchmark, COCO contains five human-annotated captions per image. We exploited this by generating pairs of captions describing the same image, effectively creating a high-quality, ground-truth paraphrasing dataset.

## 3   Architecture Description

The model follows the standard Encoder-Decoder structure composed of stacked layers.
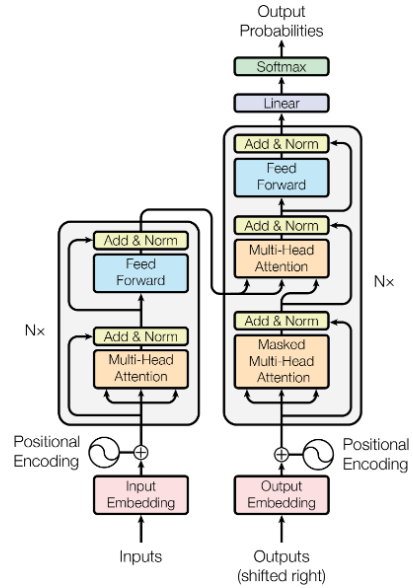


Figure 1: The transformer model architecture. [1]

### 3.1   The Encoder-Decoder Stack

The encoder consists of $N = 3$ identical layers. Each layer has two sub-layers: a Multi-Head Self-Attention mechanism and a position-wise fully connected feed-forward network. We employ a residual connection around each of the two sub-layers, followed by layer normalization.

The decoder shares a similar structure but includes a third sub-layer: Multi-Head Cross-Attention, which performs attention over the output of the encoder stack.

## 3.2 Scaled Dot-Product Attention

The core computational unit is defined as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \qquad (1)$$

Where inputs consisting of queries $Q$ and keys $K$ of dimension $d_k$ are scaled to prevent vanishing gradients in the softmax function. It is needed as otherwise the variance of the distribution would scale with $d_k$ in each iteration.

## 3.3 Positional Encoding

Since the model contains no recurrence, we inject information about the relative or absolute position of the tokens using sinusoidal functions:

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}}) \qquad (2)$$

## 3.4 Hardware Optimizations

To accommodate the training on an NVIDIA RTX 2060 (6GB VRAM), we deviated from the base paper's hyperparameters. We reduced $d_{model}$ to 256 and utilized Automatic Mixed Precision (AMP) to perform forward passes in FP16, reducing memory footprint by approximately 40%.

## 4 Results

### 4.1 Experimental Setup

The model was trained using the Adam optimizer with $\beta_1 = 0.9, \beta_2 = 0.98$ and a dynamic learning rate scheduler. We used a Byte-Pair Encoding (BPE) tokenizer with a shared vocabulary of 30,000 tokens.

Table 1: Model Hyperparameters

| Parameter | Value |
|---|---|
| Encoder/Decoder Layers ($N$) | 3 |
| Model Dimension ($d_{model}$) | 256 |
| Feed-Forward Dimension ($d_{ff}$) | 1024 |
| Attention Heads ($h$) | 8 |
| Effective Batch Size | 64 |

### 4.2 Qualitative Analysis

After 4 epochs of training, the model demonstrates the ability to generalize syntax and identify semantic categories, though it occasionally hallucinates specific object details.

Table 2 shows sample outputs generated using Greedy Decoding.

### 4.3 Discussion

The results indicate that the model successfully learned the structure of the English language and the concept of paraphrasing. For instance, "dig" (likely a typo for

Table 2: Paraphrasing Output Samples

| Input Sentence | Generated Paraphrase |
|---|---|
| *A dig sitting on the grass* | A bench sitting on a lush green field. |
| *A gir and a man playing on the filed* | A man in a suit and tie playing a game of frisbee. |
| *A woman interacting with an apple* | A woman holding a banana in her hand. |

dog) was interpreted as a stationary object ("bench"), and "field" triggered the context of "lush green."

The transformation of "apple" to "banana" highlights that the model has learned categorical embeddings (fruit) but lacks the precision required to distinguish specific instances without further training time or a larger model dimension.

## 5 Conclusion

We successfully implemented and validated a Transformer architecture from scratch. The results confirm that even with reduced depth and dimensionality, the attention mechanism is powerful enough to capture semantic meaning. Future work will involve implementing Beam Search to improve generation quality and scaling the model parameters as hardware permits.

## References

[1] Vaswani, A., et al. (2017). Attention Is All You Need. *Advances in Neural Information Processing Systems*, 30.

[2] Lin, T. Y., et al. (2014). Microsoft COCO: Common Objects in Context. *European Conference on Computer Vision*.