



FACULTY DIALOG SYSTEMS AND MACHINE LEARNING

IMPLEMENTING TRANSFORMERS

Project Report

Author Nils Reck 2898155

Supervisors Dr. Carel van Niekerk Dr. Hsien-Chin Lin

29. 01. 2025

1 Introduction

The Transformer model [Vas17] established the foundation for more performant and context-aware sequence transduction by removing the recurrent or convolutional means of previous state-of-the-art models. This is mainly due to the Transformer's ability to process multiple sequences at once. Up until today, the Transformer model remains the architecture of choice for top-performing large language models, like GPT-4o.

This report showcases my attempt to implement a custom Transformer model for a German-English translation task and aligns with the practicals conducted during the course. In particular, it focuses on providing the reader with detailed knowledge about its components and their interplay, as well as my personal insights and struggles during development and training. Thus, the report commences with a methodology section, explaining the modular components of a Transformer.

2 Methodology

However, before the model can process the input data, it has to be encoded in a numerical representation that the model can interpret. For this, a shared tokenizer is trained over the source and target sequences, which maps a token (a word or subword) of a sequence to a number and vice versa.

add info about alignment of sequences (maybe add to training section)

2.1 Embedding Layers

The embedding layer creates a d_{model} -dimensional vector representation for each encoded token of the input and target sequence. Consistent with the original Transformer architecture, we apply parameter sharing by using the same set of weights for both embedding layers and the pre-softmax linear transformation, which maps the embeddings back to their respective token index. Sharing parameters between the encoder and decoder embedding layers offers several advantages. First, it can significantly reduce the model size while maintaining model performance [PW17]. Second, parameter sharing reduces the degrees of freedom of the model, thus implicitly applying regularization by forcing different parts of the model to use the same parameters, preventing the model from overfitting. Additionally, the efficiency of the model improves because shared parameters allow for faster updates and fewer memory operations. Finally, by tying the input and output embeddings together, the model can enhance cross-lingual transfer learning, as aligned word representations across languages make it easier to generalize.

Unlike recurrent architectures, which process sequences step by step, the Transformer processes entire sequences in parallel. To compensate for the lack of sequence order awareness, the positional encoding layer enriches the representations with fixed positional information.

delete last point?

How does the model differentiate between embedding and position?

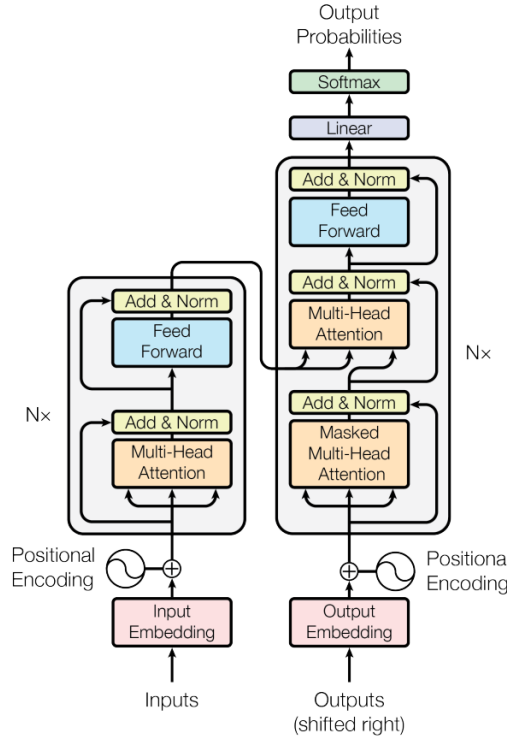


Figure 1: The original transformer architecture, adapted from Vaswani et al. [Vas17]

add formulas if space permits

2.2 Encoder Stack

The encoder consists of six identical layers, each designed to transform the input sequence into a context-rich representation. Each layer comprises two sub-layers: a multi-head self-attention mechanism and a position-wise feed-forward network (see Section 2.5), each followed by a residual connection [HZRS15] and layer normalization [BKH16] (see Section 2.6) to stabilize training and improve gradient flow. Residual connections, defined as $y = \mathbf{x} + f(\mathbf{x})$, preserve the original signal while adding important features from multi-head attention or feed-forward layers, alleviating the problem of vanishing gradients during backpropagation. If the transformation $f(\mathbf{x})$ collapses to zero (e.g. due to all weights and biases being pushed to zero), the output reduces to $y = \mathbf{x}$, ensuring that the original signal is preserved when the layer does not learn anything. Residual connections can also be described by the residual mapping

$g(\mathbf{x}) = f(\mathbf{x}) - \mathbf{x}$, emphasizing that the network only needs to learn a small transformation when $f(\mathbf{x})$ is close to the identity function. Learning a function close to the identity function, residual blocks slightly refine existing features instead of learning full (high variance) functions from scratch.

Additionally, by focussing on small residuals $g(\mathbf{x}) = f(\mathbf{x}) - \mathbf{x}$, the network increases the chance of generalizing better to unseen data, reducing overfitting.

why?

2.3 Decoder Stack

The decoder also consists of six identical layers. In addition to the two sub-layers of the encoder, it has a second multi-head attention mechanism over the outputs of the encoder. Consistent with the encoder, residual connections and layer normalization are employed after each sub-layer. In contrast to the multi-head self-attention layer in the encoder, the inputs to the attention mechanism in the decoder are masked such that the decoder cannot attend to future tokens. This prevents the decoder from cheating by attending to tokens it has not yet seen. Finally, the output of the decoder undergoes a linear transformation. After that, softmax is applied to convert the output into probabilities to predict the next token.

2.4 Attention

The attention function injects contextual information about related tokens into each token's representation. This process enables the model to capture dependencies between words, regardless of their position in the sequence.

The first step is to create the query(Q), key(K), and value(V) vectors from the encoder or decoder input vectors by multiplying them with three matrices that are learned during training. These matrices must be learned in a way that they reflect meaningful similarity relationships in terms of attention.

According to Equation (1), the attention function then first computes a score for each token in the sequence relative to every other token by taking the dot product of the query vector with the transposed key vector.

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

This computation, along with all other operations of the attention mechanism, is performed in parallel for all tokens in each sequence across the entire batch. Next, the result is scaled by $\sqrt{d_k}$ to avoid exploding gradients and improve stability. Then a softmax function is applied to maintain relevant words, subside words we can mostly ignore, and prepare the output to be summed up. Finally, by multiplying the softmax scores by V produces a new representation for each token. While it retains most of its original structure, it is enriched with contextual information from the most relevant tokens for our translation task.

In the attention mechanism, we employ two types of masks; first, an attention mask to prevent the attention function from attending to padding tokens.

Second, a causal mask in the Masked Multi-Head Attention of the decoder to prevent our model from cheating by attending to future tokens during inference.

2.5 Position-Wise Feed-Forward Networks

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2 \quad (2)$$

The FNN introduces a higher-dimensional space to explore combinations of features present in the token embeddings that it could not explore in the original embedding space. That happens by the first linear transformation $xW_1 + b_1$. Next, ReLU, $\max(0, xW_1 + b_1)$ introduces non-linearity (why does that help?) and helps prevent vanishing gradients (how?). The FFN has two linear layers of size $(d_{\text{model}}, d_{\text{ffn}})$ and $(d_{\text{ffn}}, d_{\text{model}})$, respectively. Finally, the non-linearly transformed representation is projected back into the original space, d_{model} , such that it (what is it?) is forced to focus on the most significant feature combinations (bring examples).

2.6 Normalization Layer

Layer normalization is applied after each self-attention and feed-forward sub-layer.

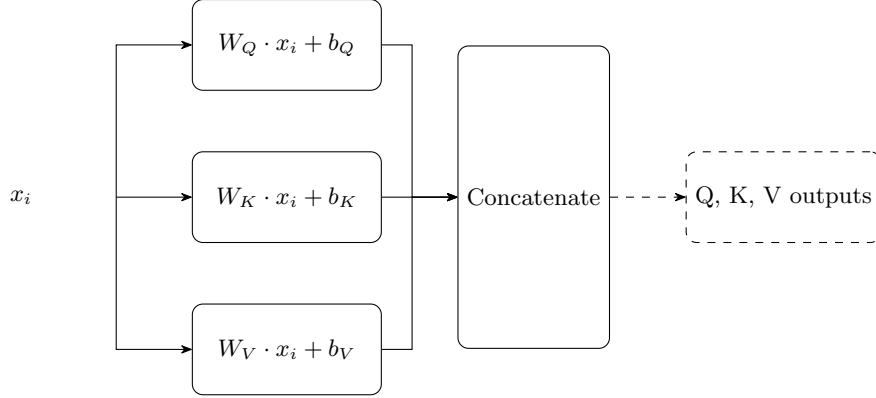
$$\text{LayerNorm}(x) = \frac{x - \mu}{\sigma} \cdot \gamma + \beta \quad (3)$$

Where x is the input vector, in our case the token embedding, μ the mean of x , calculated across the features, σ is the standard deviation, also calculated across the features:

$$\mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l \quad \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2} \quad (4)$$

H is the number of features for each token representation, γ and β are optional, learnable parameters to scale and shift the normalized values.

In a transformer architecture, the layer normalization layer serves different purposes: it stabilizes training by normalizing the distributions of the layer inputs, thus preventing exploding or vanishing gradients, which would also have adverse, covariate effects on the surrounding layers in the forward and backward passes. Additionally, contrary to batch normalization, layer normalization handles variations in sequence length better, since it computes the mean and variance along the features of the token and not across the individual features across the batch.



Three parallel linear transformations

3 Training

This section covers the training regime and a comparison between CPU and GPU training. The code is available on GitLab¹.

3.1 Data

We train on the WMT 17 German to English dataset², consisting of about 4.9 million sentence pairs after filtering out sequences exceeding 64 tokens in length. To not inflict unnecessary load on the GPU during training, we preprocess the datasets beforehand. The sequences are encoded using byte-pair encoding [BGLL17], which has a shared source-target vocabulary of 50000 tokens. We use only 5% of the training data for benchmarking CPU vs. GPU performance with a batch size of 32 on both. For the final model performance reported in Section 4, however, we train on the entire dataset with a batch size of 1024 sentence pairs.

3.2 Training and Schedule

We train our models on a single node with five processing cores and a single NVIDIA A100 GPU for XXX steps (30 epochs). Since GPUs are optimized for parallel computation, they are well-suited for the highly parallel nature of sequence processing in Transformer models. Unlike CPUs, which are designed for handling a wide range of sequential operations, GPUs distribute the workload across its many cores, consisting of 8192 FP32 CUDA cores and 432 Tensor

¹https://git.hhu.de/nirec101/transformer_project

²<https://www.statmt.org/wmt17/translation-task.html>

cores³, optimized for matrix operations. For benchmarking CPU vs. GPU performance under identical conditions, we use the exact same set of hyperparameters on both setups, with the CPU configuration consisting of five cores and 64GB of memory. Each step took about XXX seconds, utilizing mixed precision.

3.3 AdamW Optimizer

In all the experiments, we use the AdamW optimizer [LH19] with $\beta_1 = 0.9$, $\beta_2 = 0.99$ and $\epsilon = 10^{-8}$. This section elaborates the core differences between the Adam optimizer [KB17] used in the original Transformer architecture and AdamW.

Both, in Adam and AdamW, Equation (5) shows that the learning rate is adjusted for each parameter independently based on the history of gradients. The running averages, m_{t-1} and v_{t-1} , make it possible to include the history of the gradients in the calculation of the first and second moment:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t, \quad v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \quad (5)$$

The calculation of the first and second moment in this fashion ensures that parameters with larger gradient variances are updated more slowly than those with larger gradient variances to stabilize the optimization process.

The bias correction from Equation (6) is important because, without it, the first and second moments are biased toward zero at early time steps, because m_0 and v_0 are zero. Consequently, this results in overly careful parameter updates in the beginning, which hinder the performance and convergence of the training process.

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t} \quad (6)$$

In the original Adam, weight decay is added directly to the gradient. Consequently, this means that the weight decay term is included in the moment estimates (m_t and v_t). The AdamW optimizer circumvents this problem: The weight decay is applied directly to the weights after the adaptive gradient update, as shown in Equation (7):

$$\theta_t \leftarrow \theta_t - \eta \lambda \theta_t \quad (7)$$

Equation (8) shows the complete parameter update for the AdamW optimizer, where the weight decay is decoupled from the gradient calculation.

$$\theta_{t+1} = \theta_t - \eta \left(\frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \right) - \eta \lambda \theta_t \quad (8)$$

³<https://images.nvidia.com/aem-dam/en-zz/Solutions/data-center/nvidia-ampere-architecture-whitepaper.pdf>

Add figure from AdamW paper?

In the backward pass, if you have multiple rank-deficient matrices, your rank becomes even lower. because the composition of rank-deficient matrices leads to a further reduction in the rank, potentially causing the gradients to vanish or lose critical information needed for effective weight updates.

4 Results

This section provides the results of the CPU vs. GPU comparison, as well as the performance of the Transformer model on the translation task.

4.1 GPU versus CPU training

Table 1 illustrates the results of the performance comparison between CPU and GPU training. The GPU significantly accelerates the overall training process by a factor of 8.07. Accordingly, the GPU processes a single epoch, as well as a forward pass, faster by approximately the same margin. The GPU achieves the most significant speed-up (20x) during the backward pass. This highlights the GPU’s superior efficiency in handling computationally expensive tasks, such as gradient computation, which heavily rely on parallelized matrix calculations. Surprisingly, the GPU uses 38 times less memory per epoch compared to the CPU, which can mainly be attributed to the small batch size of 32 for both setups.

Due to approaching deadlines and long queues on the high performance cluster, we abstain from further optimizations of the GPU codebase. For instance, potential improvements like offloading BLEU score calculation to the CPU, and leveraging mixed precision (which is implemented in full model training), among other optimizations, could further improve performance.

why does that happen?

| Metric | CPU | GPU | CPU to GPU Ratio |
|-----------------------------------|---------|--------|------------------|
| Training time (hours) | 3.0017 | 0.3717 | 8.07 |
| Avg. Epoch times (seconds) | 2161.24 | 267.63 | 8.07 |
| Avg. Forward pass times (seconds) | 0.0770 | 0.0088 | 8.75 |
| Avg. Backward pass time (seconds) | 0.2020 | 0.0101 | 20 |
| Avg. Single step time (seconds) | 0.2790 | 0.0189 | 14.76 |
| Allocated memory per epoch (MB) | 32998.0 | 865.91 | 38.1 |

Table 1: Performance comparison of CPU vs. GPU

Additionally, the A100 GPU has a memory bandwidth of 1555 GB/sec, meaning it can transfer data from the Video Random Access Memory (VRAM) at an exceptionally high rate compared to a CPU.

4.2 Machine Translation

Where do we apply dropout?

What are learnable parameters in a transformer model?

Included questions from tests

References

- [BGLL17] Denny Britz, Anna Goldie, Minh-Thang Luong, and Quoc Le. Massive exploration of neural machine translation architectures, 2017.
- [BKH16] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016.
- [HZRS15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- [KB17] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- [LH19] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.
- [PW17] Ofir Press and Lior Wolf. Using the output embedding to improve language models, 2017.
- [Vas17] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.