

TRANSFORMER

Hallvard Høyland Lavik
hallvard.hoyland.lavik@nmbu.no



Contents

1	Motivation	5
2	Traditional sequence based deep learning	5
2.1	Basic structure	5
2.2	Long short-term memory	5
2.2.1	Γ_f Forget gate layer	6
2.2.2	Γ_u Update/input gate layer	6
2.2.3	Γ_o Output gate layer	6
3	Attention	7
3.1	Origin of attention in deep learning	7
3.2	Inputs	7
3.3	Decoder-only models	8
3.4	Encoder-decoder models	8
3.5	Query, key and value	8
3.6	Similarity score	9
3.6.1	Similarity functions	9
3.6.2	Masking	9
3.7	Python equivalent attention	10
3.8	Parallelization	10
3.9	Multi-head attention	10
4	Natural language processing	10
4.1	Tokenization	10
4.1.1	Byte-pair encoding	10
4.1.2	Regular expression	11
4.1.3	Training algorithm	12
4.2	Embedding	12
4.2.1	word2vec	13
5	Architecture	14
5.1	Positional encoding	15
5.2	Encoder	15
5.3	Decoder	15
5.3.1	Masked multi-head attention	16
5.3.2	Encoder-decoder attention	17
5.4	Attention	17
5.5	Multi-layer perceptron	18

6	Translation model implementation	18
7	Generation model implementation	19
8	GPT-2 finetuning implementation	19
9	Source code	22
10	Deep learning libraries	22
11	Methods	22

Figures

1	LSTM block.	6
2	Transformer architecture (from Huggingface [16]).	14

Listings

3.7	Masking of attention scores.	10
4.1.2	Regular expression used by GPT-4 for tokenization.	11
4.1.3	Merge logic of byte-pair encoding.	12
4.1.3	Training logic of byte-pair encoding.	12
5.1	Positional encoding layer.	15
5.3.1	Square mask helper function.	16
5.3.1	Masking of future tokens for multi-head attention.	16
5.4	Attention.	17
5.5	Multi-layer perceptron.	18
5.4	Transformer block.	18

Motivation

The transformer architecture, introduced in "Attention is All You Need" [1], aimed to overcome limitations of existing sequence-to-sequence models based on recurrent neural networks and convolutional neural networks. These models struggled with capturing long-range dependencies, processing long sequences, and efficiently utilizing parallel computation.

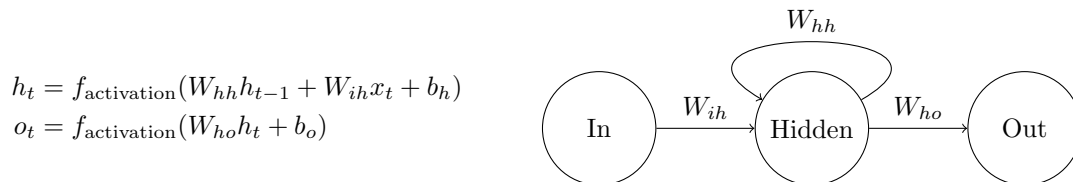
In order to overcome these challenges, the team [1] proposed a new architecture relying on attention mechanisms, to capture global dependencies between input and output elements. This design allows for more effectively attending long-range dependencies, and is the reason behind the transformers huge success in performing tasks related to natural language processing, among other things.

Traditional sequence based deep learning

Recurrent neural networks (RNNs) are a type of neural network designed for sequence-based data. Unlike traditional feedforward neural networks, RNNs have connections that loop back within the network, providing a way to maintain an internal state (*i.e.*, its memory).

BASIC STRUCTURE

A basic RNN has an input layer (In), a hidden recurrent layer (Hidden), and an output layer (Out). The recurrent layer processes sequences of data by looping its output back into its input, allowing it to learn from past information.



Where t represents the position in the sequence (*e.g.*, time), o the output, i the input and b the bias. h_{t-1} therefore represents the hidden output for the previous time-step, and h_t the current hidden output. $f_{\text{activation}}$ for the hidden and output layers may differ, and represent their activation functions.

LONG SHORT-TERM MEMORY

A Long Short-Term Memory (LSTM) model is a special type of RNN that can better learn long-term dependencies in the data compared to a simple RNN. It has a more complex internal structure involving gates that control the flow of information.

The LSTM structure consists of four gates, which combine or remove information. The operations done are linear, being

$$\begin{aligned} \oplus \text{ Element-wise addition. } & \begin{bmatrix} 0.8 \\ 0.8 \\ 0.8 \end{bmatrix} \oplus \begin{bmatrix} 1.0 \\ 0.5 \\ 0.0 \end{bmatrix} = \begin{bmatrix} 0.8 + 1.0 \\ 0.8 + 0.5 \\ 0.8 + 0.0 \end{bmatrix} = \begin{bmatrix} 1.8 \\ 1.3 \\ 0.8 \end{bmatrix} \\ \otimes \text{ Element-wise multiplication. } & \begin{bmatrix} 0.8 \\ 0.8 \\ 0.8 \end{bmatrix} \otimes \begin{bmatrix} 1.0 \\ 0.5 \\ 0.0 \end{bmatrix} = \begin{bmatrix} 0.8 \cdot 1.0 \\ 0.8 \cdot 0.5 \\ 0.8 \cdot 0.0 \end{bmatrix} = \begin{bmatrix} 0.8 \\ 0.4 \\ 0.0 \end{bmatrix} \end{aligned}$$

By inspecting these operations, we can see that gates using \otimes is able to either block or allow information to pass through (respectively through values of 0.0 or 1.0), or something in-between. This means, that the network can learn previous state values, and take these into account when filtering values of new inputs.

An LSTM network has a *cell state* C , which acts as the memory of the network. This state is being transferred across the time-steps, thus allowing for previous inputted information to be retained in future time-steps.

Γ_f FORGET GATE LAYER

The first step in an LSTM is the *forget gate* layer. This is a neural network which takes in the previous output along with current input. This layer has a sigmoid activation function, $\sigma(\cdot)$, where inputs resulting in 0's lead to variables being left out of the cell state C .

$$\Gamma_f = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \quad C_t^* = \Gamma_f \otimes C_{t-1}$$

Γ_u UPDATE/INPUT GATE LAYER

The next step is to decide what new information to store in the cell state.

$$\Gamma_u = \sigma(W_u \cdot [h_{t-1}, x_t] + b_u) \quad \tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

These two are then pairwise multiplied together, and added to the forgotten state.

$$\begin{aligned} \text{CELL STATE} \quad C_t &= C_{t-1}^* \oplus (\Gamma_u \otimes \tilde{C}_t) \\ &= (\Gamma_f \otimes C_{t-1}) \oplus (\Gamma_u \otimes \tilde{C}_t) \end{aligned}$$

Γ_o OUTPUT GATE LAYER

The output of the model is then calculated based on both the cell state and previous output as well as input.

$$\begin{aligned} \Gamma_o &= \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \\ \text{OUTPUT} \quad h_t &= \Gamma_o \otimes \tanh(C_t) \end{aligned}$$

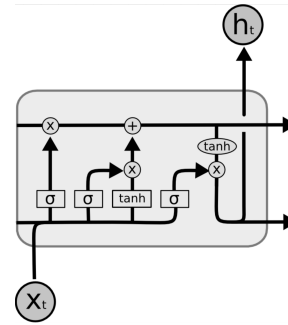


Figure 1: LSTM block.

Although models like the mentioned Recurrent Neural Network (RNN) and Long Short-Term Memory (LSTM) enable the retention of some prior information, they inherently struggle with handling long sequences. For example, when an LSTM cell state, C , is presented with an input sequence of length L , it is unable to maintain information from all previous steps when evaluating step l . This is because the cell state progressively becomes

more abstract as it processes each step in the sequence, making it challenging to retain information from the distant past.

Attention

Arguably, the most important core function of information processing in the brain is to selectively attend important impressions [2], and has also become the core feature of the transformer model [3]. For instance, when reading a document, much of its content is redundant and the human brain is able to attend to the important parts and their connection – thus disregarding most of the redundant information. Likewise, the attention mechanism of the transformer aims at selectively attending different parts of the input sequence when generating an output sequence. [1]

Attention mechanisms used in transformers can in simple terms be thought of as scalars that enhance or diminish some input, like the `forget gate` in the LSTM which use `element wise multiplication`. However, in order to capture more intricate connections in the input-sequence, a few tricks are applied.

When calculating the attention based on some input, it is important to note that the "input" consists of the full sequence. A practical example would be when predicting the next word based on the sentence;

At my dairy farm we always get our milk from ---,

where the model would need the full context in order to properly complete the sentence.

ORIGIN OF ATTENTION IN DEEP LEARNING

When trying to solve problems related to machine translation, Bahdanau *et al.* [4] began experimenting with how the context of a given sentence may be used when predicting an output. In the paper [4], they came up with a method where, for any given word i in the sequence, its context is composed of a weighed sum of the other words in the sequence.

While recent methods has optimized this approach by using matrices [1] instead of a neural network [4], the theory remains the same.

While it is possible to obtain the (additive) attention through a separate neural network, it is deemed more computational efficient to use (multiplicative attention through) optimized matrix multiplication algorithms [1].

INPUTS

When training a transformer model, the inputs are represented as a tensor \mathbf{X} of word embeddings of a sequence, hereafter thought of as a matrix of words to simplify the intuition;

$$\mathbf{X} = \begin{bmatrix} \text{sequence } 0 \\ \vdots \\ \text{sequence } N \end{bmatrix} = \begin{bmatrix} \text{word}_{0,0} & \dots & \text{word}_{0,M} \\ \vdots & \ddots & \vdots \\ \text{word}_{N,0} & \dots & \text{word}_{N,M} \end{bmatrix}. \quad (1)$$

Here, a sequence could represent a cutout from a text, *i.e.*, `sequence_0 = text[0:M]`. Thus, in a similar manner, we can also express the targets as the left-shifted version of \mathbf{X} (where a sequence therefore becomes `text[1:M+1]`);

$$\mathbf{Y} = \begin{bmatrix} \text{left-shifted sequence } 0 \\ \vdots \\ \text{left-shifted sequence } N \end{bmatrix} = \begin{bmatrix} \text{word}_{0,1} & \dots & \text{word}_{0,M+1} \\ \vdots & \ddots & \vdots \\ \text{word}_{N,1} & \dots & \text{word}_{N,M+1} \end{bmatrix}. \quad (2)$$

Where a practical example of a row of \mathbf{X} and \mathbf{Y} could be:

$$\begin{aligned} \mathbf{x} &= [\text{At, my, dairy, farm, ...}] \\ \mathbf{y} &= [\text{my, dairy, farm, we, ...}] \end{aligned} \quad (3)$$

where we see that the targets of $\mathbf{x}_{0 \rightarrow i}$ is \mathbf{y}_i .

DECODER-ONLY MODELS

For a generative model, such as GPT-2 [5], which consists of solely the decoder of the transformer, the \mathbf{X} and \mathbf{Y} seen in Equations (1) and (2) are used.

ENCODER-DECODER MODELS

When the goal of the transformer is to map some inputs to a specified output (*e.g.*, translation between languages), the previously explained procedure (of using the shifted inputs as targets) does not work. The matrices \mathbf{X} and \mathbf{Y} thus become;

$$\mathbf{X} = \begin{bmatrix} \text{source } 0 \\ \vdots \\ \text{source } N \end{bmatrix} \quad \text{and} \quad \mathbf{Y} = \begin{bmatrix} \text{target } 0 \\ \vdots \\ \text{target } N \end{bmatrix}, \quad (4)$$

where **source** and **target** represents the inputs and expected outputs (*e.g.*, a sentence in one language as the source and its translation as the target).

The transformer thus has to construct the target from scratch, whereas the **decoder-only** models build on the input by iteratively predicting the next word. This iterative nature is also present in the encoder-decoder models, but missing the initial starting-point.

QUERY, KEY AND VALUE

When calculating what to attend based on some input, \mathbf{X} , three new quantities, query, key and value, are introduced, respectively;

$$\begin{aligned} \mathbf{Q} &= \mathbf{XW}_Q \\ \mathbf{K} &= \mathbf{XW}_K, \\ \mathbf{V} &= \mathbf{XW}_V \end{aligned} \quad (5)$$

along with their respective sets of weights, \mathbf{W}_i .

The matrices containing the queries, keys and values can thus be seen as three different representations of the inputs which is used to enhance or diminish certain aspects of the context, when predicting the next element of the sequence.

While *Geometry of Deep Learning* tries to provide a biological analogy as to what the query and key represents, its actual representation in terms of the transformer is rather abstracted [6, 7]. It is however worth noting that all three of them is some form of embedding of the input, and is used to process the context of said input.

In practice, these three matrices are obtained by a single fully-connected layer which takes the inputs, \mathbf{X} , and outputs $\text{concat}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$.

SIMILARITY SCORE

The score of an arbitrary query vector, \mathbf{q}_i contained as a column in \mathbf{Q} , and the key vectors \mathbf{k}_j for $\mathbf{k}_j \in \mathbf{K}$, is found by taking the dot product between them

$$\text{score}_{i,j} = \mathbf{q}_i \cdot \mathbf{k}_j \quad \text{for } j = 1, \dots, M, \quad (6)$$

where M is the dimension of \mathbf{k} . Here, we get a score for each key element in the sequence \mathbf{k}_j and the chosen query element, \mathbf{q}_i . The intuition behind the dot product is to find the importance of the other elements contained in the sequence with respect to the current element.

These scores are then (typically) scaled by the root of their dimension, M , and normalized using the softmax function such that a probabilistic representation is obtained;

$$\text{weighting}_{i,j} = \text{softmax} \left(\frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{M}} \right) \quad \text{for } j = 1, \dots, M. \quad (7)$$

SIMILARITY FUNCTIONS

While the mentioned approach using the scaled dot product is the most common [6, 8, 7, 9], other functions for calculating the score may be used.

Other such functions include the non-scaled dot product, $\mathbf{q}_i \cdot \mathbf{k}_j$, and the cosine similarity, $\frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\|\mathbf{q}_i\| \|\mathbf{k}_j\|}$, as presented in *Geometry of Deep Learning* [6].

MASKING

Intuitively, (during training,) the attention scores need to be masked in order to prevent the influence of future words of preceding words. That is, when a model is generating a sequence, it does not know the future generated words before having generated them. Thus, when training, the scores of all words after the "current" word is set to `float('-inf')`, such that it is lost when `softmax` is performed;

$$\text{score} = \begin{bmatrix} \text{score}_{0,0} & \dots & -\infty \\ \vdots & \ddots & \vdots \\ \text{score}_{N,1} & \dots & \text{score}_{N,M} \end{bmatrix}. \quad (8)$$

This can be seen in practice through Listing 3.7.

When the softmax score has been calculated, it is multiplied with the value representation of the input, \mathbf{V} , much like the `forget gate` in the LSTM which enhance or diminish the focus on certain elements of the sequence. The output of the attention layer for element i is then the sum of all j products. That is;

$$\begin{aligned} \text{attention}_i &= \sum_{j=1}^N \text{weighting}_{i,j} \mathbf{v}_j \\ &= \sum_{j=1}^N \text{softmax} \left(\frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{M}} \right) \mathbf{v}_j \quad \text{for } j = 1, \dots, M. \end{aligned} \quad (9)$$

Which can be rewritten in terms of matrix notation [1]:

$$\text{attention} = \text{softmax} \left(\frac{\mathbf{QK}^T}{\sqrt{M}} \right) \mathbf{V}. \quad (10)$$

This allows the transformer model to selectively focus on different parts of the input sequence when generating each output element, improving its ability to capture long-range dependencies and context compared to just transferring the cell state C as is done in LSTMs.

PYTHON EQUIVALENT ATTENTION

The code implementation of Equation (10) here presented, with [masking](#) (Equation (8)):

```
1 score = (q @ k.transpose(-2, -1)) / math.sqrt(N)
2 mask = torch.tril(torch.ones(M, M))
3 score = score.masked_fill(mask[:T, :T] == 0, float('-inf'))
4 score = torch.nn.functional.softmax(score, dim=-1)
5 attention = score @ v
```

PARALLELIZATION

Unlike recurrent neural networks, which process input tokens one at a time, the attention mechanism calculates the weighted sum of value vectors for each input by considering all key and value vectors at once, enabling parallel computation, as seen in [Equation \(10\)](#).

MULTI-HEAD ATTENTION

Multi-head attention is an extension of attention that allows the model to attend each position of the sequence simultaneously. This is achieved by linearly projecting the input tokens into multiple attention heads, computing the attention scores for each head independently, and then concatenating the results. This process enables the model to capture a diverse range of contextual information and improve its overall understanding of the input sequence. [1]

Natural language processing

TOKENIZATION

Tokenization, in terms of language processing, involves breaking down text into small pieces called tokens. Said tokens vary between methods, and could include phrases, words, or even single characters. These tokens are then mapped to a numerical representation, as state-of-the-art artificial intelligence language models rely on numerical input. Tokenization methods are therefore used to preprocess the inputted text. [10] Thus, a tokenization model is in simple terms a dictionary of mappings between tokens and their corresponding numerical representations.

BYTE-PAIR ENCODING

Byte-pair encoding (BPE) is a type of tokenization that is commonly used in large language models [11]. It works by initially representing text as a sequence of characters, and then iteratively replacing the most frequent pair of bytes with a single, new byte. This process continues until the desired vocabulary size is

reached. The result is a set of tokens that represent common sequences of characters, which can be more efficient and effective for language modeling than word-based tokenization. [12, 10, 13]

BPE is commonly used because it strikes a balance between character-level and word-level tokenization. It can handle out-of-vocabulary words and rare words more effectively than word-level tokenization, while still capturing meaningful linguistic units. [12, 13]

REGULAR EXPRESSION

In order to simplify the vocabulary creation, regular expressions are useful. The regular expressions are increasingly complex [13], but help filter out redundant characters. The regular expression below is used for GPT-4 when tokenizing inputs (according to Karpathy [10]) [13]. This regular expression is extremely complex, and a description of its components are found in [Table 1](#).

```
'(?i:[sdmt]|ll|ve|re)|[^\r\n\p{L}\p{N}]?+\p{L}+|\p{N}{1,3}| ?[^\s\p{L}\p{N}
}]++[\r\n]*|s*[\r\n]|\s+(?!\S)|\s+
```

COMPONENT	DESCRIPTION
'(?i:[sdmt] ll ve re)	Uses a case-insensitive inline modifier (?i:...) to match common abbreviations and contractions in English text ('). Matches either a single character that is one of s, d, m, or t, or the two-letter sequences ll, ve, or re.
[^\r\n\p{L}\p{N}]?+\p{L}+	Uses a possessive quantifier ?+ to match one or more Unicode letters (\p{L}) that are preceded by zero or more characters that are not (^) Unicode letters, line breaks (\r\n), or Unicode numbers (\p{N}). Matches words that start with non-letter characters, such as hyphenated words.
\p{N}{1,3}	Matches up to three Unicode numbers (\p{N}). Matches numbers in the text.
._?[^\s\p{L}\p{N}]++[\r\n]*	Matches zero or one space character (.), followed by one or more characters that are not (^) whitespace (\s), Unicode letters (\p{L}), or Unicode numbers (\p{N}), followed by zero or more line breaks ([\r\n]*). The ? at the beginning of the pattern makes the preceding space optional. Matches punctuation and symbols that are not part of words or numbers.
\s*[\r\n]	Matches zero or more whitespace characters (\s*) followed by a line break ([\r\n]). Matches line breaks in the text.
\s+(?!\S)	Matches one or more whitespace characters (\s+) that are not followed by a non-whitespace character (?!\S). Matches whitespace at the end of lines.
\s+	Matches one or more whitespace characters, <i>i.e.</i> , spaces between words.

TRAINING ALGORITHM

When creating the BPE tokenizer, elements are continuously merged. The following code (from [10]) handles this merging; given a list of integers, `ids`, it replaces all consecutive occurrences of `pair` with the new token `idx`:

```
1 def merge(ids, pair, idx):
2     newids = []
3     i = 0
4     while i < len(ids):
5         if ids[i] == pair[0] and i < len(ids) - 1 and ids[i+1] == pair[1]:
6             newids.append(idx)
7             i += 2
8         else:
9             newids.append(ids[i])
10            i += 1
11    return newids
```

With the helper-function defined, the tokenizer (from [10]) can begin training on a given string, `text`, and a regular expression pattern, `pattern` (e.g., the pattern seen [above](#)):

```
1 merges = {}
2 vocabulary = {}
3
4 ids = [list(pt.encode("utf-8")) for pt in regex.findall(pattern, text)]
5
6 for i in range(vocabulary_size):
7     stats = {}
8     for chunk_ids in ids:
9         for pair in zip(ids, ids[1:]):
10            stats[pair] = counts.get(pair, 0) + 1
11
12    pair = max(stats, key=stats.get)
13
14    ids = [merge(chunk_ids, pair, i) for chunk_ids in ids]
15
16    merges[pair] = i
17    vocabulary[i] = vocabulary[pair[0]] + vocabulary[pair[1]]
```

The objective is to iteratively merge the most common pairs of characters or tokens in the given text to create new tokens until the desired `vocabulary_size` is achieved.

First, the `text` to be trained on is split into chunks (`pt`) based on the `pattern`, and each chunk is encoded into a list of bytes.

In each iteration, the frequency of consecutive pairs in the encoded chunks is calculated, and the pair with the highest frequency is merged into a new token. All occurrences of the pair in the encoded chunks are replaced with the new token, which is also added to the `vocabulary` and the `merges` dictionary. The updated `vocabulary` and `merges` dictionary are then stored, and is what is used when encoding or decoding new text or tokens, respectively.

EMBEDDING

Word (*i.e.*, token) embedding is a technique in natural language processing where words from the vocabulary are mapped to vectors of real numbers. These vectors capture the semantic properties of the words. For instance, words that are semantically similar are generally mapped to vectors that are close to each other in the vector space.

WORD2VEC

In order to achieve this vectorized word representation, word2vec is commonly used, as it is highly generalizable. The reason for this generalizability is that word2vec is a neural network-based method that comes in two variants: continuous bag of words (CBOW) (based on Cloze process from Taylor [14]) and Skip-Gram.

The CBOW model predicts a target word given its context words, while the Skip-Gram model predicts the context words given a target word. In both models, the input words are represented as one-hot vectors, which are then multiplied by a (learned) weight matrix to get the corresponding word embeddings. During training, the output words are also represented as one-hot vectors, which serve as the ground truth during training. The goal of the models is to learn a weight matrix that can map each word in the vocabulary to a dense, low-dimensional vector that captures its semantic and syntactic relationships with other words. [15]

For pure embedding models (*i.e.*, models which is used for embedding only), [word2vec](#) is a good approach. For large language models, however, the embedding is typically represented as a single dense layer (`torch.nn.Embedding` is in simple terms the same as `torch.nn.Linear`, except its handling of index representations instead of raw values), mapping token indices of dimensions `vocabulary_size` \mapsto `embedding_size`.

Architecture

The transformer architecture consists of two primary components: the encoder and the decoder, both of which contain multiple identical layers stacked upon each other. The sub-layers of the encoder and decoder can be seen in the left and right part of [Figure 2](#), respectively.

As the transformer architecture processes the entire input sequence simultaneously, its computational efficiency is far greater compared to traditional [RNNs](#) or [LSTMs](#), as explained in the original paper by Vaswani *et al.* [1].

While the architecture displayed in [Figure 2](#) is the originally presented transformer [1, 16], variations of the core architecture which aim at overcoming certain disadvantages can be found throughout the literature (like [17, 18, 19, 20, 21], among others).

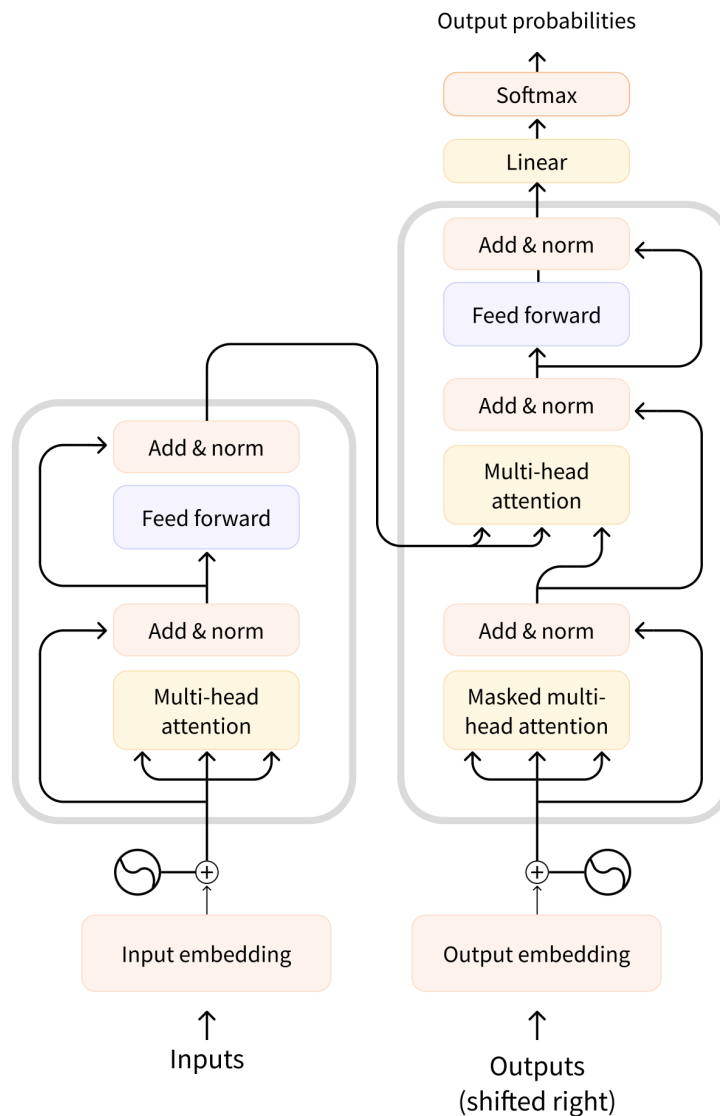


Figure 2: Transformer architecture (from Huggingface [16]).

POSITIONAL ENCODING

Since the transformer does not inherently consider the position or order of input tokens, positional encodings are added to provide this information. These encodings can be learned or fixed, and are added to the input embeddings before being passed through the model. [1] Positional encoding thus allow the model to maintain a sense of sequence order, which is important for understanding the context and relationships between tokens in the input data.

As seen in "Attention Is All You Need" [1], the positional encoding was done similarly to the code below, with modifications made by Karpathy [9].

```
1 class PositionalEncoding(torch.nn.Module):
2     def __init__(self, n_embd: int, dropout: float, maxlen: int = 5000):
3         super().__init__()
4         den = torch.exp(- torch.arange(0, n_embd, 2) * math.log(10000) /
5                             n_embd)
6         pos = torch.arange(0, maxlen).reshape(maxlen, 1)
7
8         pos_embedding = torch.zeros((maxlen, n_embd))
9         pos_embedding[:, 0::2] = torch.sin(pos * den)
10        pos_embedding[:, 1::2] = torch.cos(pos * den)
11
12        pos_embedding = pos_embedding.unsqueeze(-2)
13
14        self.dropout = torch.nn.Dropout(dropout)
15        self.register_buffer('pos_embedding', pos_embedding)
16
17    def forward(self, token_embedding: torch.Tensor):
18        return self.dropout(token_embedding + self.pos_embedding[:
19                                token_embedding.size(0), :])
```

While there are many ways to encode positions, the reasoning behind the chosen sinusoidal encoding is that "[t]he wavelengths form a geometric progression from 2π to $10000 \times 2\pi$ [...] because we hypothesized it would allow the model to easily learn to attend by relative positions, since for any fixed offset k , $[\text{pos_embedding}_{\text{pos}+k}]$ can be represented as a linear function of $[\text{pos_embedding}_{\text{pos}}]$." [1]

ENCODER

The encoder takes a sequence of (*e.g.*, tokenized word) positional embeddings as input, processes all elements of the sequence simultaneously using attention, and outputs a new sequence of vectors that represent the input sequence with the added context of the full sequence.

In order to achieve this, both a multi-head attention mechanism and a position-wise feed-forward network is used, the latter being a fully connected neural network applied independently to each position in the sequence [1], as seen in Figure 2.

DECODER

The decoder take the output of the encoder and generate an output sequence (*e.g.*, a translated sentence). The decoder does this by using several layers, each of which contains three sub-layers: a masked multi-head attention mechanism, an encoder-decoder attention mechanism, and a position-wise feed-forward network.

MASKED MULTI-HEAD ATTENTION

As explained when introducing attention, masking ensure that the predictions for position i are only dependent on positions before i . This is necessary during training to prevent the model from "cheating" by looking at the target sequence in its entirety. The mask ensures that the model only attends to previous positions when generating each token in the output sequence. [1] Intuitively, this masking has no effect when generating new text, but is important during training (where the output is known).

This masking may, for instance, be applied the in the following manner, (as opposed to the alternative approach in Listing 3.7,) based on [16]:

```
1 def square_mask(self, dim):
2     mask = (
3         torch.tril(torch.ones((dim, dim),
4                               device=self.config.device)) == 1
5     )
6     mask = mask.float().masked_fill(
7         mask == 0,
8         float('-inf'))
9     ).masked_fill(mask == 1, float(0.0))
10    return mask
```

Where the `square_mask` method creates a square mask of a given dimension. This mask is used to prevent future tokens from being used in the prediction of the current token during training. The method takes an integer `dim` as input, which represents the dimension of the square mask. It then creates a square matrix of ones with the same dimension, and applies the `torch.tril` function to it, which returns the lower triangular part of the matrix, and the values are replaced with `-inf` for zeros and `0.0` for ones, as per [1].

The `masking` method creates four different masks for the `src` (source) and `tgt` (target) data: the source mask, target mask, source padding mask, and target padding mask. The source and target masks are square masks created using the `square_mask` method above, with dimensions equal to the sequence lengths of the source and target data, respectively. The padding masks are created by comparing the source and target data to the padding token (here fetched from the tokenizer model), and transposing the resulting boolean tensors.

```
1 def masking(self, src, tgt):
2     src_seq_len = src.shape[0]
3     tgt_seq_len = tgt.shape[0]
4
5     tgt_m = self.square_mask(tgt_seq_len).to(self.config.device)
6     src_m = torch.zeros(
7         (src_seq_len, src_seq_len),
8         device=self.config.device
9     ).type(torch.bool)
10
11     src_pad_m = (
12         src == self.config.tokenizer["special_symbols"]["[PAD]"]
13     ).transpose(0, 1).to(self.config.device)
14     tgt_pad_m = (
15         tgt == self.config.tokenizer["special_symbols"]["[PAD]"]
16     ).transpose(0, 1).to(self.config.device)
17
18     return src_m, tgt_m, src_pad_m, tgt_pad_m
```

These masks are then used by the model to ignore padding tokens in the input data, and to ensure that the prediction for each token is only dependent on the previous tokens in the sequence.

ENCODER-DECODER ATTENTION

Seen in [Figure 2](#) where the output of the encoder (left) is connected to the decoder (right). The attention mechanism here allows the decoder to focus on different parts of the input sequence when generating each token in the output sequence. It does this by combining the output of the encoder and the output of the masked multi-head attention mechanism as input, thus being able to consider both the context of the input sequence and the previously generated tokens in the output sequence.

Each of these sub-layers also includes a residual connection (or skip connection) around it, followed by layer normalization. [1]

The decoder's output is a sequence of vectors, each of which represents a token in the output sequence. The final linear layer and softmax function are then applied to each vector to produce a probability distribution over the target vocabulary, which is used to generate the output sequence.

ATTENTION

See [section regarding attention](#) above for the theoretical explanation of attention. The following is the implementation from Karpathy [9]:

```
1 class Attention(torch.nn.Module):
2     def __init__(self, config):
3         super().__init__()
4
5         self.c_attn = torch.nn.Linear(config.n_embd, 3 * config.n_embd,
6                                         bias=config.bias)
7         self.c_proj = torch.nn.Linear(config.n_embd, config.n_embd,
8                                         bias=config.bias)
9
10        self.attn_dropout = torch.nn.Dropout(config.dropout)
11        self.resid_dropout = torch.nn.Dropout(config.dropout)
12        self.n_head = config.n_head
13        self.n_embd = config.n_embd
14        self.dropout = config.dropout
15
16        self.register_buffer(
17            "bias",
18            torch.tril(
19                torch.ones(config.block_size, config.block_size)
20            ).view(1, 1, config.block_size, config.block_size)
21        )
22
23    def forward(self, x):
24        B, T, C = x.size()
25
26        q, k, v = self.c_attn(x).split(self.n_embd, dim=2)
27        k = k.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
28        q = q.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
29        v = v.view(B, T, self.n_head, C // self.n_head).transpose(1, 2)
30
31        att = (q @ k.transpose(-2, -1)) * (1.0 / math.sqrt(k.size(-1)))
32        att = att.masked_fill(self.bias[:, :, :T, :T] == 0, float('-inf'))
33        att = torch.nn.functional.softmax(att, dim=-1)
34        att = self.attn_dropout(att)
35        y = att @ v
36        y = y.transpose(1, 2).contiguous().view(B, T, C)
37
```

```

38         y = self.resid_dropout(self.c_proj(y))
39         return y

```

MULTI-LAYER PERCEPTRON

In the transformer architecture seen in [Figure 2](#), each attention block also consists of a feed-forward layer. This layer, often referred to as "multi-layer perceptron" (MLP) can be represented through the class:

```

1 class MLP(nn.Module):
2     def __init__(self, config):
3         super().__init__()
4         self.c_fc = nn.Linear(config.n_embd, 4 * config.n_embd, bias=
5             config.bias)
6         self.gelu = nn.GELU()
7         self.c_proj = nn.Linear(4 * config.n_embd, config.n_embd, bias=
8             config.bias)
9         self.dropout = nn.Dropout(config.dropout)
10
11     def forward(self, x):
12         x = self.c_fc(x)
13         x = self.gelu(x)
14         x = self.c_proj(x)
15         x = self.dropout(x)
16         return x

```

Which each element is passed through after the attention mechanism has been performed, such that each block can be concisely contained in the class:

```

1 class Block(nn.Module):
2     def __init__(self, config):
3         super().__init__()
4         self.ln_1 = LayerNorm(config.n_embd, bias=config.bias)
5         self.attn = Attention(config)
6         self.ln_2 = LayerNorm(config.n_embd, bias=config.bias)
7         self.mlp = MLP(config)
8
9     def forward(self, x):
10         x = x + self.attn(self.ln_1(x))
11         x = x + self.mlp(self.ln_2(x))
12         return x

```

The two classes being fetched from Karpathy [9].

See [22, 23, 24, 9, 16] among others for further implementations of the transformer.

Translation model implementation

| Henlo.

Generation model implementation

| Henlo.



GPT-2 finetuning implementation

| Henlo.

References

- [1] Ashish Vaswani et al. *Attention Is All You Need*. 2023. arXiv: [1706.03762 \[cs.CL\]](#).
- [2] George R. Mangun. *The Neuroscience of Attention: Attentional Control and Selection*. Oxford University Press, Jan. 2012. ISBN: 9780195334364. DOI: [10.1093/acprof:oso/9780195334364.001.0001](#). URL: <https://doi.org/10.1093/acprof:oso/9780195334364.001.0001>.
- [3] Rick Merritt. *What is a Transformer Model?* 2022. URL: <https://blogs.nvidia.com/blog/what-is-a-transformer-model/>.
- [4] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. *Neural Machine Translation by Jointly Learning to Align and Translate*. 2016. arXiv: [1409.0473 \[cs.CL\]](#).
- [5] Alec Radford et al. “Language Models are Unsupervised Multitask Learners”. In: (2019).
- [6] Jong Chul Ye. *Geometry of Deep Learning: A Signal Processing Perspective*. Springer Nature Singapore, 2022. ISBN: 9789811660467. DOI: <https://doi.org/10.1007/978-981-16-6046-7>.
- [7] Jay Alammar. *The Illustrated Transformer*. 2018. URL: jalammar.github.io/illustrated-transformer/.
- [8] Mary Phuong and Marcus Hutter. *Formal Algorithms for Transformers*. 2022. arXiv: [2207.09238 \[cs.LG\]](#).
- [9] Andrej Karpathy. *nanoGPT*. 2023. URL: <https://github.com/karpathy/nanogpt>.
- [10] Andrej Karpathy. *minbpe*. 2024. URL: <https://github.com/karpathy/minbpe>.
- [11] Thomas Wolf et al. “Transformers: State-of-the-Art Natural Language Processing”. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Online: Association for Computational Linguistics, Oct. 2020, pp. 38–45. URL: <https://www.aclweb.org/anthology/2020.emnlp-demos.6>.
- [12] Wikipedia contributors. *Byte pair encoding*. 2024. URL: https://en.wikipedia.org/wiki/Byte_pair_encoding.
- [13] OpenAI. *tiktoken*. 2024. URL: <https://github.com/openai/tiktoken>.
- [14] Wilson L. Taylor. ““Cloze Procedure”: A New Tool for Measuring Readability”. In: *Journalism Quarterly* 30.4 (1953), pp. 415–433. DOI: [10.1177/107769905303000401](#). eprint: <https://doi.org/10.1177/107769905303000401>. URL: <https://doi.org/10.1177/107769905303000401>.
- [15] TensorFlow. *word2vec*. 2024. URL: <https://www.tensorflow.org/text/tutorials/word2vec>.
- [16] Huggingface. *How do Transformers work?* URL: <https://huggingface.co/learn/nlp-course/chapter1/4>.
- [17] Jacob Devlin et al. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. 2019. arXiv: [1810.04805 \[cs.CL\]](#).
- [18] Albert Gu and Tri Dao. *Mamba: Linear-Time Sequence Modeling with Selective State Spaces*. 2023. arXiv: [2312.00752 \[cs.LG\]](#).
- [19] Tom B. Brown et al. *Language Models are Few-Shot Learners*. 2020. arXiv: [2005.14165 \[cs.CL\]](#).
- [20] Zihang Dai et al. *Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context*. 2019. arXiv: [1901.02860 \[cs.LG\]](#).
- [21] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. *Reformer: The Efficient Transformer*. 2020. arXiv: [2001.04451 \[cs.LG\]](#).
- [22] PyTorch. *Language Modeling with nn.Transformer and torchtext*. 2024. URL: https://pytorch.org/tutorials/beginner/transformer_tutorial.html.
- [23] Sasha Rush et al. *The Annotated Transformer*. 2022. URL: <http://nlp.seas.harvard.edu/annotated-transformer/#embeddings-and-softmax>.

- [24] Kevin Ko. *Transformer: PyTorch Implementation of "Attention Is All You Need"*. 2019. URL: <https://github.com/hyunwoongko/transformer/tree/master>.

Source code

The source code behind this report can be found [here](#).

Deep learning libraries

During implementation, [PyTorch](#) was mostly used.

Methods

During creation of the [codebase](#), [GitHub Copilot](#) along with [Mistral](#) was used –mostly helping with tensor dimension issues.

Prompts equivalent to

```
How do I reshape the tensor [...] to correspond with this [...]
Modify this code such that shapes match [...]
Smooth and plot the contents of a csv-file using pandas.
```

were used. *I.e.*, sparring when I was stuck on something.

During writing of this report, **no** artificial intelligence tools were used to generate text. [Mistral](#) was however used when either shortening or correcting/rephrasing certain paragraphs/sentences.

Prompts equivalent to

```
Shorten the paragraph [...]
Highlight errors in this paragraph [...]
```

were used.