

Formal Algorithms for Transformers

Mary Phuong¹ and Marcus Hutter¹

¹DeepMind

This document aims to be a self-contained, mathematically precise overview of transformer architectures and algorithms (*not results*). It covers what transformers are, how they are trained, what they are used for, their key architectural components, and a preview of the most prominent models. The reader is assumed to be familiar with basic ML terminology and simpler neural network architectures such as MLPs.

Keywords: formal algorithms, pseudocode, transformers, attention, encoder, decoder, BERT, GPT, Gopher, tokenization, training, inference.

Contents

1	Introduction	1
2	Motivation	1
3	Transformers and Typical Tasks	3
4	Tokenization: How Text is Represented	4
5	Architectural Components	4
6	Transformer Architectures	7
7	Transformer Training and Inference	8
8	Practical Considerations	9
A	References	9
B	List of Notation	16

A famous colleague once sent an actually very well-written paper he was quite proud of to a famous complexity theorist. His answer: “I can’t find a theorem in the paper. I have no idea what this paper is about.”

1. Introduction

Transformers are deep feed-forward artificial neural networks with a (self)attention mechanism. They have been tremendously successful in natural language processing tasks and other domains. Since their inception 5 years ago [VSP¹⁷], many variants have been suggested [LWLQ21]. Descriptions are usually graphical, verbal, partial, or incremental. Despite their popularity, it seems no pseudocode has ever been published for any variant. Contrast this to other fields of computer science, even to “cousin” discipline reinforcement learning [MKS¹³, SBB18, EMK²¹].

This report intends to rectify the situation for Transformers. It aims to be a self-contained, com-

plete, precise and compact overview of transformer architectures and formal algorithms (but *not results*). It covers what Transformers are (Section 6), how they are trained (Section 7), what they’re used for (Section 3), their key architectural components (Section 5), tokenization (Section 4), and a preview of practical considerations (Section 8) and the most prominent models.

The essentially complete pseudocode is about 50 lines, compared to thousands of lines of actual real source code. We believe these formal algorithms will be useful for theoreticians who require compact, complete, and precise formulations, experimental researchers interested in implementing a Transformer from scratch, and encourage authors to augment their paper or text book with formal Transformer algorithms (Section 2).

The reader is assumed to be familiar with basic ML terminology and simpler neural network architectures such as MLPs.

In short, a (formally inclined) reader, upon understanding the contents of this document, will have a solid grasp of transformers: they will be ready to read and contribute to the literature on the topic as well as implement their own Transformer using the pseudocode as templates.

2. Motivation

The true story above the introduction describes quite well the feeling we have when browsing

many Deep Learning (DL) papers; unable to figure out the algorithmic suggestions exactly. For practitioners, the papers may be sufficiently detailed, but the precision needed by theoreticians is usually higher. For reasons not entirely clear to us, the DL community seems shy of providing pseudocode for their neural network models. Below we describe the SOTA in DL paper writing and argue for the value of formal algorithms. The reader already convinced about their merit can without loss skip this section.

The lack of scientific precision and detail in DL publications. Deep Learning has been tremendously successful in the last 5 to 10 years with thousands of papers published every year. Many describe only informally how they change a previous model. Some 100+ page papers contain only a few lines of prose informally describing the model [RBC⁺21]. At best there are some high-level diagrams. No pseudocode. No equations. No reference to a precise explanation of the model. One may argue that most DL models are minor variations of a few core architectures, such as the Transformer [VSP⁺17], so a reference augmented by a description of the changes should suffice. This would be true if (a) the changes were described precisely, (b) the reference architecture has been described precisely elsewhere, and (c) a reference is given to this description. Some if not all three are lacking in most DL papers. To the best of our knowledge no-one has even provided pseudocode for the famous Transformer and its encoder/decoder-only variations.

Interfacing algorithms. Equally important are proper explanations of how these networks are trained and used, but sometimes it is not even clear what the inputs and outputs and potential side-effects of the described model are. Of course someone experienced would usually be able to correctly guess, but this is not a particularly scientific approach. The experimental section in publications often does not explain what is actually fed into the algorithms and how. If there is some explanation in the methods section, it is often disconnected from what is described in the experimental section, possibly due to different

authors writing the different sections. The core algorithms for the models should be accompanied by the wrapper algorithms that call them, e.g. (pre)training, fine-tuning, prompting, inference, deployment. Sometimes this is simple, but sometimes the magic happens there. In any case, if these things are not formally stated they remain unclear. Again, if the setup is standard and has been formally explained elsewhere, a simple reference will do.

Source code vs pseudocode. Providing open source code is very useful, but not a proper substitute for formal algorithms. There is a massive difference between a (partial) Python dump and well-crafted pseudocode. A lot of abstraction and clean-up is necessary: remove boiler plate code, use mostly single-letter variable names, replace code by math expressions wherever possible, e.g. replace loops by sums, remove (some) optimizations, etc. A well-crafted pseudocode is often less than a page and still essentially complete, compared to often thousands of lines of real source code. This is hard work no-one seems to be willing to do. Of course a forward process of first designing algorithms and write up pseudocode on paper, and then implementing it is fine too, but few DL practitioners seem to work that way.

Examples of good neural network pseudocode and mathematics and explanations. Multi-Layer Perceptrons (MLPs) are usually well-described in many papers, e.g. [MPCB14, BFT17, JGH18], though also without pseudocode. For a rare text-book example of pseudocode for a non-trivial neural network architecture, see Algorithm S2 of [SGBK⁺21], which constitutes a *complete*, i.e. *essentially executable*, pseudocode of just 25 lines based on a 350-line Python Colab toy implementation, which itself is based on a proper 1000+ line implementation.

This work aims to do the same for Transformers: The whole decoder-only Transformer GPT Algorithm 10 based on attention Algorithms 4 and 5 and normalization Algorithm 6 including training Algorithm 13 and prompting and inference Algorithm 14 all-together is less than 50 lines of pseudocode, compared to e.g. the 2000-

line self-contained C-implementation [Bel21].

[Ala19] is a great blog-post explaining Transformers and [EGKZ21] describes the attention mechanism to sufficient mathematical precision to allow proving properties about it, but neither provides pseudocode. [Elh21] is an attempt to understand Transformers by reverse-engineering the computations they perform and interpreting them as circuits.

Motivation. But does anyone actually need pseudocode and what would they be useful for (we sometimes get asked)? We find the absence of pseudocode in DL and this question quite perplexing, but apparently it requires answering. Providing such pseudocode can be useful for many purposes:

- They can be used as templates and adapted to precisely describe future variations, and therewith set a new standard in DL publishing. We explicitly encourage the reader to copy and adapt them to their needs and cite the original as “adapted from [PH22]”.
- Having all that matters on one page in front of you makes it easier to develop new variations compared to reading prose or scrolling through 1000s of lines of actual code.
- They can be used as a basis for new implementations from scratch, e.g. in different programming languages, without having to wade through and reverse-engineer existing idiosyncratic real source code.
- They may establish notational convention, which eases communication and reading future variations.
- The process of converting source code into pseudocode can exhibit implementation errors (as it e.g. did in [SGBK⁺21]).
- Theoreticians need compact, complete, and precise representations for reasoning and ultimately proving properties about algorithms. They are often unwilling or unable to reverse engineer code, or guess the meaning of words or fancy diagrams.

With this motivation in mind, the following five sections formally describe all aspects of transformer architectures, training, and inference.

3. Transformers and Typical Tasks

Transformers are neural network models that excel at natural language processing, or more generally at modelling sequential data. Two common types of tasks they are used for are *sequence modelling* and *sequence-to-sequence prediction*.

Notation. Let V denote a finite set, called a *vocabulary*, often identified with $[N_V] := \{1, \dots, N_V\}$. This could be words or letters, but typically are sub-words, called tokens. Let $x \equiv x[1 : \ell] \equiv x[1]x[2] \dots x[\ell] \in V^*$ be a sequence of tokens, e.g. a sentence or a paragraph or a document. Unlike in Python, we use arrays starting from 1, and $x[1 : \ell]$ includes $x[\ell]$. For a matrix $M \in \mathbb{R}^{d \times d'}$, we write $M[i, :] \in \mathbb{R}^{d'}$ for the i th row and $M[:, j] \in \mathbb{R}^d$ for the j -th column. We use matrix \times column vector convention more common in mathematics, compared to the default row vector \times matrix in the transformer literature, i.e. our matrices are transposed. See Appendix B for a complete list of notation.

Chunking. The predominant paradigm in machine learning is (still) learning from independent and identically distributed (i.i.d.) data. Even for sequence modelling for practical reasons this tradition is upheld. The training data may naturally be a collection of (independent) articles, but even then, some may exceed the maximal context length ℓ_{\max} transformers can handle. In this case, an article is crudely broken into shorter chunks of length $\leq \ell_{\max}$.

Sequence modelling (DTransformer). Given a vocabulary V , let $x_n \in V^*$ for $n \in [N_{\text{data}}]$ be a dataset of sequences (imagined to be) sampled i.i.d. from some distribution P over V^* . The goal is to learn an estimate \hat{P} of the distribution $P(x)$. In practice, the distribution estimate is often decomposed via the chain rule as $\hat{P}(x) = \hat{P}_\theta(x[1]) \cdot \hat{P}_\theta(x[2] | x[1]) \dots \hat{P}_\theta(x[\ell] | x[1 : \ell - 1])$, where θ consists of all neural network parameters to be learned. The goal is to learn a distribution over a single token $x[t]$ given its preceding tokens $x[1 : t - 1]$ as context.

Examples include e.g. language modelling, RL policy distillation, or music generation.

Sequence-to-sequence (seq2seq) prediction (EDTransformer). Given a vocabulary V and an i.i.d. dataset of sequence pairs $(z_n, x_n) \sim P$, where P is a distribution over $V^* \times V^*$, learn an estimate of the conditional distribution $P(x|z)$. In practice, the conditional distribution estimate is often decomposed as $\hat{P}(x|z) = \hat{P}_\theta(x[1]|z) \cdot \hat{P}_\theta(x[2]|x[1], z) \cdots \hat{P}_\theta(x[\ell]|x[1:\ell-1], z)$.

Examples include translation (z = a sentence in English, x = the same sentence in German), question answering (z = question, x = the corresponding answer), text-to-speech (z = a piece of text, x = a voice recording of someone reading the text).

Classification (ETransformer). Given a vocabulary V and a set of classes $[N_C]$, let $(x_n, c_n) \in V^* \times [N_C]$ for $n \in [N_{\text{data}}]$ be an i.i.d. dataset of sequence-class pairs sampled from $P(x, c)$. The goal in classification is to learn an estimate of the conditional distribution $P(c|x)$.

Examples include e.g. sentiment classification, spam filtering, toxicity classification.

4. Tokenization: How Text is Represented

In the context of natural language tasks, *tokenization* refers to how a piece of text such as “My grandma makes the best apple pie.” is represented as a sequence of vocabulary elements (called *tokens*).

Character-level tokenization. One possible choice is to let V be the English alphabet (plus punctuation). In the example above, we’d get a sequence of length 36: [‘M’, ‘y’, ‘ ’, ...]. Character-level tokenization tends to yield very long sequences.

Word-level tokenization. Another choice would be to let V consist of all English words

(plus punctuation). In the example above, we’d get a sequence of length 7: [‘My’, ‘grandma’, ‘makes’, ...]. Word-level tokenization tends to require a very large vocabulary and cannot deal with new words at test time.

Subword tokenization. This is the method used in practice nowadays: V is a set of commonly occurring word segments like ‘cious’, ‘ing’, ‘pre’. Common words like ‘is’ are often a separate token, and single characters are also included in V to ensure all words are expressible.

There are in fact many ways to do subword tokenization. One of the simplest and most successful ones is Byte Pair Encoding [Gag94, SHB16] used in GPT-2 [RWC¹⁹].

Final vocabulary and text representation. Given a choice of tokenization / vocabulary, each vocabulary element is assigned a unique index in $\{1, 2, \dots, N_V - 3\}$. A number of special tokens are then added to the vocabulary. The number of special tokens varies, and here we will consider three: `mask_token` := $N_V - 2$, used in masked language modelling (see Algorithm 12); `bos_token` := $N_V - 1$, used for representing the beginning of sequence; and `eos_token` := N_V , used for representing the end of sequence. The complete vocabulary has $N_V = |V|$ elements.

A piece of text is represented as a sequence of indices (called *token IDs*) corresponding to its (sub)words, preceded by `bos_token` and followed by `eos_token`.

5. Architectural Components

The following are the neural network building blocks (functions with learnable parameters) from which transformers are made. Full architectures featuring these building blocks are presented in the next section. (By a slight abuse of notation, we identify V with the set $\{1, 2, \dots, N_V\}$.)

Token embedding. The token embedding learns to represent each vocabulary element as a vector in \mathbb{R}^{d_e} ; see Algorithm 1.

Algorithm 1: Token embedding.

Input: $v \in V \cong [N_V]$, a token ID.
Output: $e \in \mathbb{R}^{d_e}$, the vector representation of the token.
Parameters: $W_e \in \mathbb{R}^{d_e \times N_V}$, the token embedding matrix.

1 **return** $e = W_e[:, v]$

Positional embedding. The positional embedding learns to represent a token’s position in a sequence as a vector in \mathbb{R}^{d_e} . For example, the position of the first token in a sequence is represented by a (learned) vector $W_p[:, 1]$, the position of the second token is represented by another (learned) vector $W_p[:, 2]$, etc. The purpose of the positional embedding is to allow a Transformer to make sense of word ordering; in its absence the representation would be permutation invariant and the model would perceive sequences as “bags of words” instead.

Learned positional embeddings require that input sequence length is at most some fixed number ℓ_{\max} (the size of the learned positional embedding matrix must be finite and fixed in advance of training). An intuitive explanation of how this works can be found at [Ala18]. For pseudocode, see Algorithm 2.

Not all transformers make use of *learned* positional embeddings, some use a hard-coded mapping $W_p : \mathbb{N} \rightarrow \mathbb{R}^{d_e}$ instead [Ker21]. Such hard-coded positional embeddings can (theoretically) handle arbitrarily long sequences. The original Transformer [VSP⁺17] uses

$$W_p[2i-1, t] = \sin(t/\ell_{\max}^{2i/d_e}), \\ W_p[2i, t] = \cos(t/\ell_{\max}^{2i/d_e}).$$

for $0 < i \leq d_e/2$.

The positional embedding of a token is usually added to the token embedding to form a token’s initial embedding. For the t -th token of a sequence x , the embedding is

$$e = W_e[:, x[t]] + W_p[:, t]. \quad (1)$$

Attention. Attention is the main architectural component of transformers. It enables a neural

Algorithm 2: Positional embedding.

Input: $\ell \in [\ell_{\max}]$, position of a token in the sequence.
Output: $e_p \in \mathbb{R}^{d_e}$, the vector representation of the position.
Parameters: $W_p \in \mathbb{R}^{d_e \times \ell_{\max}}$, the positional embedding matrix.

1 **return** $e_p = W_p[:, \ell]$

network to make use of contextual information (e.g. preceding text or the surrounding text) for predicting the current token.

On a high level, attention works as follows: the token currently being predicted is mapped to a *query* vector $q \in \mathbb{R}^{d_{\text{attn}}}$, and the tokens in the context are mapped to *key* vectors $k_t \in \mathbb{R}^{d_{\text{attn}}}$ and *value* vectors $v_t \in \mathbb{R}^{d_{\text{value}}}$. The inner products $q^\top k_t$ are interpreted as the degree to which token $t \in V$ is important for predicting the current token q – they are used to derive a distribution over the context tokens, which is then used to combine the value vectors. An intuitive explanation how this achieves attention can be found at [Ala18, Ala19]. The precise algorithm is given in Algorithm 3.

Algorithm 3: Basic single-query attention.

Input: $e \in \mathbb{R}^{d_{\text{in}}}$, vector representation of the current token
Input: $e_t \in \mathbb{R}^{d_{\text{in}}}$, vector representations of context tokens $t \in [T]$.
Output: $\tilde{v} \in \mathbb{R}^{d_{\text{out}}}$, vector representation of the token and context combined.
Parameters: $W_q, W_k \in \mathbb{R}^{d_{\text{attn}} \times d_{\text{in}}}$,
 $b_q, b_k \in \mathbb{R}^{d_{\text{attn}}}$, the query and key linear projections.
Parameters: $W_v \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$, $b_v \in \mathbb{R}^{d_{\text{out}}}$, the value linear projection.

1 $q \leftarrow W_q e + b_q$
2 $\forall t : k_t \leftarrow W_k e_t + b_k$
3 $\forall t : v_t \leftarrow W_v e_t + b_v$
4 $\forall t : \alpha_t = \frac{\exp(q^\top k_t / \sqrt{d_{\text{attn}}})}{\sum_u \exp(q^\top k_u / \sqrt{d_{\text{attn}}})}$
5 **return** $\tilde{v} = \sum_{t=1}^T \alpha_t v_t$

There are many ways the basic attention mechanism is used in transformers. We list some of the most common variants below.

It will be useful to define the softmax function for matrix arguments, as well as a Mask matrix:

$$\text{softmax}(A)[t_z, t_x] := \frac{\exp A[t_z, t_x]}{\sum_t \exp A[t, t_x]}, \quad (2)$$

$$\text{Mask}[t_z, t_x] = \begin{cases} 1 & \text{for bidirectional attention} \\ [[t_z \leq t_x]] & \text{for unidirectional att.} \end{cases} \quad (3)$$

Algorithm 4: $\tilde{V} \leftarrow \text{Attention}(X, Z | \mathcal{W}_{qkv}, \text{Mask})$

```
/* Computes a single (masked) self- or
   cross- attention head. */
Input:  $X \in \mathbb{R}^{d_x \times \ell_x}$ ,  $Z \in \mathbb{R}^{d_z \times \ell_z}$ , vector
   representations of primary and
   context sequence.
Output:  $\tilde{V} \in \mathbb{R}^{d_{\text{out}} \times \ell_x}$ , updated
   representations of tokens in  $X$ ,
   folding in information from
   tokens in  $Z$ .
Parameters:  $\mathcal{W}_{qkv}$  consisting of:
    $W_q \in \mathbb{R}^{d_{\text{attn}} \times d_x}$ ,  $b_q \in \mathbb{R}^{d_{\text{attn}}}$ 
    $W_k \in \mathbb{R}^{d_{\text{attn}} \times d_z}$ ,  $b_k \in \mathbb{R}^{d_{\text{attn}}}$ 
    $W_v \in \mathbb{R}^{d_{\text{out}} \times d_z}$ ,  $b_v \in \mathbb{R}^{d_{\text{out}}}$ .
Hyperparameters:  $\text{Mask} \in \{0, 1\}^{\ell_z \times \ell_x}, \uparrow(3)$ 
1  $Q \leftarrow W_q X + b_q \mathbf{1}^\top$  [[Query  $\in \mathbb{R}^{d_{\text{attn}} \times \ell_x}$ ]]
2  $K \leftarrow W_k Z + b_k \mathbf{1}^\top$  [[Key  $\in \mathbb{R}^{d_{\text{attn}} \times \ell_z}$ ]]
3  $V \leftarrow W_v Z + b_v \mathbf{1}^\top$  [[Value  $\in \mathbb{R}^{d_{\text{out}} \times \ell_z}$ ]]
4  $S \leftarrow K^\top Q$  [[Score  $\in \mathbb{R}^{\ell_z \times \ell_x}$ ]]
5  $\forall t_z, t_x, \text{ if } \neg \text{Mask}[t_z, t_x] \text{ then } S[t_z, t_x] \leftarrow -\infty$ 
6 return  $\tilde{V} = V \cdot \text{softmax}(S / \sqrt{d_{\text{attn}}})$ 
```

Bidirectional / unmasked self-attention. Given a sequence of token representations, this variant applies attention to each token, treating all tokens in the sequence as the context. See Algorithm 4, called with token sequence $Z = X$ and no masking (Mask=1).

Unidirectional / masked self-attention. Given a sequence of token representations, this variant applies attention to each token, treating all preceding tokens (including itself) as the context. Future tokens are masked out, so this causal auto-regressive version can be used for online prediction. See Algorithm 4, called with token sequence $Z = X$ and $\text{Mask}[t_z, t_x] := [[t_z \leq t_x]]$. For this Mask, the output $\tilde{V}[:, 1 : t]$ only depends

on $X[:, 1 : t]$, hence can be used to predict $X[:, t + 1]$.

Cross-attention. Given two sequences of token representations (often in the context of a sequence-to-sequence task), this variant applies attention to each token of the primary token sequence X , treating the second token sequence Z as the context. See Algorithm 4, called with $\text{Mask}=1$. While the output \tilde{V} and input sequences X have the same length ℓ_x , the context sequence Z can have different length ℓ_z .

Multi-head attention. The attention algorithm presented so far (Algorithm 4) describes the operation of a single *attention head*. In practice, transformers run multiple attention heads (with separate learnable parameters) in parallel and combine their outputs; this is called *multi-head attention*; see Algorithm 5

Algorithm 5: $\tilde{V} \leftarrow \text{MHAttention}(X, Z | \mathcal{W}, \text{Mask})$

```
/* Computes Multi-Head (masked) self-
   or cross- attention layer. */
Input:  $X \in \mathbb{R}^{d_x \times \ell_x}$ ,  $Z \in \mathbb{R}^{d_z \times \ell_z}$ , vector
   representations of primary and
   context sequence.
Output:  $\tilde{V} \in \mathbb{R}^{d_{\text{out}} \times \ell_x}$ , updated
   representations of tokens in  $X$ ,
   folding in information from
   tokens in  $Z$ .
Hyperparameters:  $H$ , number of
   attention heads
Hyperparameters:  $\text{Mask} \in \{0, 1\}^{\ell_z \times \ell_x}, \uparrow(3)$ 
Parameters:  $\mathcal{W}$  consisting of
   For  $h \in [H]$ ,  $\mathcal{W}_{qkv}^h$  consisting of:
   |  $W_q^h \in \mathbb{R}^{d_{\text{attn}} \times d_x}$ ,  $b_q^h \in \mathbb{R}^{d_{\text{attn}}}$ ,
   |  $W_k^h \in \mathbb{R}^{d_{\text{attn}} \times d_z}$ ,  $b_k^h \in \mathbb{R}^{d_{\text{attn}}}$ ,
   |  $W_v^h \in \mathbb{R}^{d_{\text{mid}} \times d_z}$ ,  $b_v^h \in \mathbb{R}^{d_{\text{mid}}}$ .
    $W_o \in \mathbb{R}^{d_{\text{out}} \times H d_{\text{mid}}}$ ,  $b_o \in \mathbb{R}^{d_{\text{out}}}$ .
1 For  $h \in [H]$ :
2    $Y^h \leftarrow \text{Attention}(X, Z | \mathcal{W}_{qkv}^h, \text{Mask})$ 
3    $Y \leftarrow [Y^1; Y^2; \dots; Y^H]$ 
4 return  $\tilde{V} = W_o Y + b_o \mathbf{1}^\top$ 
```

Algorithm 6: $\hat{e} \leftarrow \text{layer_norm}(e|\gamma, \beta)$

/* Normalizes layer activations e . */

Input: $e \in \mathbb{R}^{d_e}$, neural network activations.

Output: $\hat{e} \in \mathbb{R}^{d_e}$, normalized activations.

Parameters: $\gamma, \beta \in \mathbb{R}^{d_e}$, element-wise scale and offset.

- 1 $m \leftarrow \sum_{i=1}^{d_e} e[i]/d_e$
- 2 $v \leftarrow \sum_{i=1}^{d_e} (e[i] - m)^2/d_e$
- 3 **return** $\hat{e} = \frac{e - m}{\sqrt{v}} \odot \gamma + \beta$, where \odot denotes element-wise multiplication.

Algorithm 7: Unembedding.

Input: $e \in \mathbb{R}^{d_e}$, a token encoding.

Output: $p \in \Delta(V)$, a probability distribution over the vocabulary.

Parameters: $W_u \in \mathbb{R}^{N_v \times d_e}$, the unembedding matrix.

- 1 **return** $p = \text{softmax}(W_u e)$

Layer normalisation. Layer normalisation explicitly controls the mean and variance of individual neural network activations; the pseudocode is given in Algorithm 6. Some transformers use a simpler and more computationally efficient version of layer normalization setting $m = \beta = 0$, called root mean square layer normalisation, or RMSnorm.

Unembedding. The unembedding learns to convert a vector representation of a token and its context into a distribution over the vocabulary elements; see Algorithm 7. The algorithm describes an independently learned unembedding matrix, but note that sometimes the unembedding matrix is instead fixed to be the transpose of the embedding matrix.

6. Transformer Architectures

This section presents a few prominent transformer architectures, based on attention Algorithms 4 and 5 and using normalization Algorithm 6, in roughly historical order:

- EDT [VSP¹⁷] The original sequence-to-sequence / Encoder-Decoder Transformer,

Algorithms 8, 11 and 15.

- BERT [DCLT19], which is an instance of an encoder-only transformer (encoder-only means that it is derived from the encoder-decoder architecture by dropping the decoder part), Algorithms 9 and 12.
- GPT [RWC¹⁹, BMR²⁰], which is an instance of a decoder-only transformer, Algorithms 10, 13 and 14.

While the main architectural difference between BERT and GPT is in attention masking, they also differ in a number of less important ways: e.g. they use different activation functions and the layer-norms are positioned differently. We included these differences in the pseudocode to stay faithful to the original algorithms, but note that different transformer architectures may adopt these selectively.

To simplify notation, we denote by \mathcal{W} the entire set of parameters (query, key, value and output linear projections) required by a multi-head attention layer:

$$\mathcal{W} := \begin{pmatrix} W_q^h \in \mathbb{R}^{d_{\text{attn}} \times d_x}, & b_q^h \in \mathbb{R}^{d_{\text{attn}}}, & h \in [H] \\ W_k^h \in \mathbb{R}^{d_{\text{attn}} \times d_z}, & b_k^h \in \mathbb{R}^{d_{\text{attn}}}, & h \in [H] \\ W_v^h \in \mathbb{R}^{d_{\text{mid}} \times d_z}, & b_v^h \in \mathbb{R}^{d_{\text{mid}}}, & h \in [H] \\ W_o \in \mathbb{R}^{d_{\text{out}} \times H d_{\text{mid}}}, & b_o \in \mathbb{R}^{d_{\text{out}}} \end{pmatrix} \quad (4)$$

Encoder-decoder / sequence-to-sequence transformer [VSP¹⁷]. This is the very first transformer. It was originally used for sequence-to-sequence tasks (machine translation), which is why it is more complicated than its successors.

The idea behind the architecture is as follows: First, the context sequence is encoded using bidirectional multi-head attention. The output of this ‘encoder’ part of the network is a vector representation of each context token, taking into account the entire context sequence. Second, the primary sequence is encoded. Each token in the primary sequence is allowed to use information from the encoded context sequence, as well as primary sequence tokens that precede it. See Algorithm 8 for more details.

Encoder-only transformer: BERT [DCLT19]. BERT is a bidirectional transformer trained on the task of masked language modelling. Given a piece of text with some tokens masked out, the goal is to correctly recover the masked-out tokens. The original use of BERT was to learn generally useful text representations, which could then be adapted for various downstream NLP tasks. The masking is not performed via the Mask parameter but differently: During training each input token is replaced with probability p_{mask} by a dummy token `mask_token`, and evaluation is based on the reconstruction probability of these knocked-out tokens (see Algorithm 12).

The BERT architecture resembles the encoder part of the seq2seq transformer (hence ‘encoder-only’). It is described in detail in Algorithm 9. It uses the GELU nonlinearity instead of ReLU:

$$\text{GELU}(x) = x \cdot \mathbb{P}_{X \sim N(0,1)}[X < x]. \quad (5)$$

(When called with vector or matrix arguments, GELU is applied element-wise.)

Decoder-only transformers: GPT-2 [RWC⁺19], GPT-3 [BMR⁺20], Gopher [RBC⁺21]. GPT-2 and GPT-3 are large language models developed by OpenAI, and Gopher is a large language model developed by DeepMind. They all have similar architectures and are trained by autoregressive language modelling: Given an incomplete sentence or paragraph, the goal is to predict the next token.

The main difference from BERT is that GPT/Gopher use unidirectional attention instead of bidirectional attention; they also apply layer-norms in slightly different order.

See Algorithm 10 for the pseudocode of GPT-2. GPT-3 is identical except larger, and replaces dense attention in Line 6 by sparse attention, i.e. each token only uses a subset of the full context.

Gopher also deviates only slightly from the GPT-2 architecture: it replaces layer norm in lines 5, 7 and 10 by RMSnorm ($m = \beta = 0$), and it uses different positional embeddings.

Multi-domain decoder-only transformer: Gato [RZP⁺22]. Gato is a multi-modal multi-task transformer built by DeepMind. It is a single neural network that can play Atari, navigate 3D environments, control a robotic arm, caption images, have conversations, and more.

Under the hood, each modality is converted into a sequence prediction problem by a separate tokenization and embedding method; for example images are divided into non-overlapping 16×16 patches, ordered in raster order (left-to-right, top-to-bottom) and processed by a ResNet block to obtain a vector representation.

The actual Gato architecture is then a decoder-only transformer like the one in Algorithm 10, but where Line 2 is replaced with modality-specific embedding code.

7. Transformer Training and Inference

This section lists the pseudocode for various algorithms for training and using transformers:

- **EDTraining()** Algorithm 11 shows how to train a sequence-to-sequence transformer (the original Transformer [VSP⁺17]).
- **ETraining()** Algorithm 12 shows how to train a transformer on the task of masked language modelling (like BERT [DCLT19]).
- **DTraining()** Algorithm 13 shows how to train a transformer on the task of next token prediction (like CPT-x [BMR⁺20] and Gopher [RBC⁺21]).
- **DIference()** Algorithm 14 shows how to prompt a transformer trained on next token prediction (like GPT-x [BMR⁺20]). The temperature parameter τ interpolates between most likely continuation ($\tau = 0$), faithful sampling ($\tau = 1$), and uniform sampling ($\tau = \infty$).
- **EDInference()** Algorithm 15 shows how to use a sequence-to-sequence transformer for prediction.

Gradient descent. The described training Algorithms 11 to 13 use Stochastic Gradient Descent

(SGD) $\theta \leftarrow \theta - \eta \cdot \nabla \text{loss}(\theta)$ to minimize the log loss (aka cross entropy) as the update rule. Computation of the gradient is done via automatic differentiation tools; see [BPRS18, Table 5]. In practice, vanilla SGD is usually replaced by some more refined variation such as RMSProp or AdaGrad or others [Rud16]. Adam [KB15] is used most often these days.

8. Practical Considerations

While the vanilla transformers provided here may work in practice, a variety of “tricks” have been developed over the years to improve the performance of deep neural networks in general and transformers in particular [LWLQ21]:

- **Data preprocessing:** cleaning, augmentation [FGW⁺21], adding noise, shuffling [Lem21] (besides tokenization and chunking).
- **Architecture:** sparse layers, weight sharing (besides attention).
- **Training:** improved optimizers, mini-batches, batch normalization, learning rate scheduling, weight initialization, pre-training, ensembling, multi-task, adversarial (besides layer normalization) [Sut15].
- **Regularization:** weight decay, early stopping, cross-validation, dropout, adding noise [MBM20, TZ22].
- **Inference:** scratchpad prompting, few-shot prompting, chain of thought, majority voting [LAD⁺22].
- **Others.**

A. References

- [Ala18] Jay Alammar. The Illustrated Transformer. <http://jalammar.github.io/illustrated-transformer/>, 2018.
- [Ala19] Jay Alammar. The Illustrated GPT-2 (Visualizing Transformer Language Models). <http://jalammar.github.io/illustrated-gpt2/>, 2019.
- [Bel21] Fabrice Bellard. NNCP v2: Lossless Data Compression with Transformer. <https://bellard.org/libnc/gpt2tc.html>, 2021.
- [BFT17] Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for neural networks. *NeurIPS*, 2017.
- [BMR⁺20] Tom Brown, Benjamin Mann, Nick Ryder, et al. Language models are few-shot learners. *NeurIPS*, 2020.
- [BPRS18] Atilim Gunes Baydin, Barak A. Pearl-
mutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind. Automatic Differentiation in Machine Learning: A Survey. *Journal of Machine Learning Research*, 18(153):1–43, 2018.
- [DCLT19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *ACL*, 2019.
- [EGKZ21] Benjamin L. Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive Biases and Variable Creation in Self-Attention Mechanisms. *arXiv:2110.10090 [cs, stat]*, October 2021.
- [Elh21] Nelson Elhage. A Mathematical Framework for Transformer Circuits. <https://transformer-circuits.pub/2021/framework/index.html>, 2021.
- [EMK⁺21] Yonathan Efroni, Dipendra Misra, Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Provable RL with Exogenous Distractors via Multistep Inverse Dynamics, March 2021.
- [FGW⁺21] Steven Y Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. A survey of data augmentation approaches for NLP. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 968–988, 2021.
- [Gag94] Philip Gage. A new algorithm for data compression. *Dr. Dobbs / C Users Journal*, 12(2):23–38, 1994.
- [JGH18] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *NeurIPS*, 2018.

- [KB15] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015.
- [Ker21] Jonathan Kernes. Master Positional Encoding: Part I. <https://towardsdatascience.com/master-positional-encoding-part-i-63c05d90a0c3>, March 2021.
- [LAD⁺22] Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Sloane, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. Solving Quantitative Reasoning Problems with Language Models. *arXiv:2206.14858 [cs]*, June 2022.
- [Lem21] Chris Lemke. Data preprocessing in NLP. <https://towardsdatascience.com/data-preprocessing-in-nlp-c371d53ba3e0>, July 2021.
- [LWLQ21] Tianyang Lin, Yuxin Wang, Xiangyang Liu, and Xipeng Qiu. A Survey of Transformers, June 2021.
- [MBM20] Reza Moradi, Reza Berangi, and Behrouz Minaei. A survey of regularization strategies for deep models. *Artificial Intelligence Review*, 53(6):3947–3986, August 2020.
- [MKS⁺13] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing Atari with Deep Reinforcement Learning, December 2013.
- [MPCB14] Guido F Montufar, Razvan Pascanu, Kyunghyun Cho, and Yoshua Bengio. On the number of linear regions of deep neural networks. *NeurIPS*, 2014.
- [PH22] M. Phuong and M. Hutter. Formal algorithms for transformers. Technical report, DeepMind, London, UK, 2022. LaTeX source available at <http://arXiv.org>
- [RBC⁺21] Jack W. Rae, Sebastian Borgeaud, Trevor Cai, et al. Scaling language models: Methods, analysis & insights from training gopher. *arXiv:2112.11446*, 2021.
- [Rud16] Sebastian Ruder. An overview of gradient descent optimization algorithms. <https://ruder.io/optimizing-gradient-descent/>, January 2016.
- [RWC⁺19] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 2019.
- [RZP⁺22] Scott Reed, Konrad Żołna, Emilio Parisotto, et al. A generalist agent. *arXiv:2205.06175*, 2022.
- [SBB18] Richard S. Sutton, Andrew G. Barto, and Francis Bach. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, Massachusetts, second edition edition edition, November 2018.
- [SGBK⁺21] Eren Sezener, Agnieszka Grabska-Barwińska, Dimitar Kostadinov, Maxime Beau, Sanjukta Krishnagopal, David Budden, Marcus Hutter, Joel Veness, Matthew Botvinick, Claudia Clopath, Michael Häusser, and Peter E. Latham. A rapid and efficient learning rule for biological neural circuits. Technical report, DeepMind, London, UK, 2021.
- [SHB16] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *54th Annual Meeting of the Association for Computational Linguistics*, pages 1715–1725. Association for Computational Linguistics (ACL), 2016.
- [Sut15] Ilya Sutskever. A Brief Overview of Deep Learning. <http://yyue.blogspot.com/2015/01/a-brief-overview-of-deep-learning.html>, January 2015.
- [TZ22] Yingjie Tian and Yuqi Zhang. A comprehensive survey on regularization strategies in machine learning. *Information Fusion*, 80:146–166, April 2022.
- [VSP⁺17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 2017.

Algorithm 8: $P \leftarrow \text{EDTransformer}(z, x|\theta)$

/* Encoder-decoder transformer forward pass */

Input: $z, x \in V^*$, two sequences of token IDs.

Output: $P \in (0, 1)^{N_V \times \text{length}(x)}$, where the t -th column of P represents $\hat{P}_\theta(x[t+1] | x[1:t], z)$.

Hyperparameters: $\ell_{\max}, L_{\text{enc}}, L_{\text{dec}}, H, d_e, d_{\text{mlp}} \in \mathbb{N}$

Parameters: θ includes all of the following parameters:

- $W_e \in \mathbb{R}^{d_e \times N_V}, W_p \in \mathbb{R}^{d_e \times \ell_{\max}}$, the token and positional embedding matrices.
- For $l \in [L_{\text{enc}}]$:
 - | $\mathcal{W}_l^{\text{enc}}$, multi-head encoder attention parameters for layer l , see (4),
 - | $\gamma_l^1, \beta_l^1, \gamma_l^2, \beta_l^2 \in \mathbb{R}^{d_e}$, two sets of layer-norm parameters,
 - | $W_{\text{mlp}1}^l \in \mathbb{R}^{d_{\text{mlp}} \times d_e}, b_{\text{mlp}1}^l \in \mathbb{R}^{d_{\text{mlp}}}, W_{\text{mlp}2}^l \in \mathbb{R}^{d_e \times d_{\text{mlp}}}, b_{\text{mlp}2}^l \in \mathbb{R}^{d_e}$, MLP parameters.
- For $l \in [L_{\text{dec}}]$:
 - | $\mathcal{W}_l^{\text{dec}}$, multi-head decoder attention parameters for layer l , see (4),
 - | $\mathcal{W}_l^{\text{e/d}}$, multi-head cross-attention parameters for layer l , see (4),
 - | $\gamma_l^3, \beta_l^3, \gamma_l^4, \beta_l^4, \gamma_l^5, \beta_l^5 \in \mathbb{R}^{d_e}$, three sets of layer-norm parameters,
 - | $W_{\text{mlp}3}^l \in \mathbb{R}^{d_{\text{mlp}} \times d_e}, b_{\text{mlp}3}^l \in \mathbb{R}^{d_{\text{mlp}}}, W_{\text{mlp}4}^l \in \mathbb{R}^{d_e \times d_{\text{mlp}}}, b_{\text{mlp}4}^l \in \mathbb{R}^{d_e}$, MLP parameters.

$W_u \in \mathbb{R}^{N_V \times d_e}$, the unembedding matrix.

/* Encode the context sequence */

- 1 $\ell_z \leftarrow \text{length}(z)$
- 2 for $t \in [\ell_z]$: $e_t \leftarrow W_e[:, z[t]] + W_p[:, t]$
- 3 $Z \leftarrow [e_1, e_2, \dots, e_{\ell_z}]$
- 4 **for** $l = 1, 2, \dots, L_{\text{enc}}$ **do**
- 5 | $Z \leftarrow Z + \text{MHAttention}(Z | \mathcal{W}_l^{\text{enc}}, \text{Mask} \equiv 1)$
- 6 | for $t \in [\ell_z]$: $Z[:, t] \leftarrow \text{layer_norm}(Z[:, t] | \gamma_l^1, \beta_l^1)$
- 7 | $Z \leftarrow Z + W_{\text{mlp}2}^l \text{ReLU}(W_{\text{mlp}1}^l Z + b_{\text{mlp}1}^l \mathbf{1}^\top) + b_{\text{mlp}2}^l \mathbf{1}^\top$
- 8 | for $t \in [\ell_z]$: $Z[:, t] \leftarrow \text{layer_norm}(Z[:, t] | \gamma_l^2, \beta_l^2)$
- 9 **end**
- 10 | /* Decode the primary sequence, conditioning on the context */
- 11 for $t \in [\ell_x]$: $e_t \leftarrow W_e[:, x[t]] + W_p[:, t]$
- 12 $X \leftarrow [e_1, e_2, \dots, e_{\ell_x}]$
- 13 **for** $i = 1, 2, \dots, L_{\text{dec}}$ **do**
- 14 | $X \leftarrow X + \text{MHAttention}(X | \mathcal{W}_i^{\text{dec}}, \text{Mask}[t, t'] \equiv [[t \leq t']])$
- 15 | for $t \in [\ell_x]$: $X[:, t] \leftarrow \text{layer_norm}(X[:, t] | \gamma_i^3, \beta_i^3)$
- 16 | $X \leftarrow X + \text{MHAttention}(X, Z | \mathcal{W}_i^{\text{e/d}}, \text{Mask} \equiv 1)$
- 17 | for $t \in [\ell_x]$: $X[:, t] \leftarrow \text{layer_norm}(X[:, t] | \gamma_i^4, \beta_i^4)$
- 18 | $X \leftarrow X + W_{\text{mlp}4}^i \text{ReLU}(W_{\text{mlp}3}^i X + b_{\text{mlp}3}^i \mathbf{1}^\top) + b_{\text{mlp}4}^i \mathbf{1}^\top$
- 19 | for $t \in [\ell_x]$: $X[:, t] \leftarrow \text{layer_norm}(X[:, t] | \gamma_i^5, \beta_i^5)$
- 20 **end**
- 21 | /* Derive conditional probabilities and return */

21 **return** $P = \text{softmax}(W_u X)$

Algorithm 9: $P \leftarrow \text{ETransformer}(x|\theta)$

```
/* BERT, an encoder-only transformer, forward pass */
```

Input: $x \in V^*$, a sequence of token IDs.

Output: $P \in (0, 1)^{N_v \times \text{length}(x)}$, where each column of P is a distribution over the vocabulary.

Hyperparameters: $\ell_{\max}, L, H, d_e, d_{\text{mlp}}, d_f \in \mathbb{N}$

Parameters: θ includes all of the following parameters:

$W_e \in \mathbb{R}^{d_e \times N_v}$, $W_p \in \mathbb{R}^{d_e \times \ell_{\max}}$, the token and positional embedding matrices.

For $l \in [L]$:

| \mathcal{W}_l , multi-head attention parameters for layer l , see (4),

| $\gamma_l^1, \beta_l^1, \gamma_l^2, \beta_l^2 \in \mathbb{R}^{d_e}$, two sets of layer-norm parameters,

| $W_{\text{mlp1}}^l \in \mathbb{R}^{d_{\text{mlp}} \times d_e}$, $b_{\text{mlp1}}^l \in \mathbb{R}^{d_{\text{mlp}}}$, $W_{\text{mlp2}}^l \in \mathbb{R}^{d_e \times d_{\text{mlp}}}$, $b_{\text{mlp2}}^l \in \mathbb{R}^{d_e}$, MLP parameters.

$W_f \in \mathbb{R}^{d_f \times d_e}$, $b_f \in \mathbb{R}^{d_f}$, $\gamma, \beta \in \mathbb{R}^{d_f}$, the final linear projection and layer-norm parameters.

$W_u \in \mathbb{R}^{N_v \times d_e}$, the unembedding matrix.

```

1  $\ell \leftarrow \text{length}(x)$ 
2 for  $t \in [\ell]$  :  $e_t \leftarrow W_e[:, x[t]] + W_p[:, t]$ 
3  $X \leftarrow [e_1, e_2, \dots, e_\ell]$ 
4 for  $l = 1, 2, \dots, L$  do
5    $X \leftarrow X + \text{MHAttention}(X | \mathcal{W}_l, \text{Mask} \equiv 1)$ 
6   for  $t \in [\ell]$  :  $X[:, t] \leftarrow \text{layer_norm}(X[:, t] | \gamma_l^1, \beta_l^1)$ 
7    $X \leftarrow X + W_{\text{mlp2}}^l \text{GELU}(W_{\text{mlp1}}^l X + b_{\text{mlp1}}^l \mathbf{1}^\top) + b_{\text{mlp2}}^l \mathbf{1}^\top$ 
8   for  $t \in [\ell]$  :  $X[:, t] \leftarrow \text{layer_norm}(X[:, t] | \gamma_l^2, \beta_l^2)$ 
9 end
10  $X \leftarrow \text{GELU}(W_f X + b_f \mathbf{1}^\top)$ 
11 for  $t \in [\ell]$  :  $X[:, t] \leftarrow \text{layer_norm}(X[:, t] | \gamma, \beta)$ 
12 return  $P = \text{softmax}(W_u X)$ 

```

Algorithm 10: $P \leftarrow \text{DTransformer}(x|\theta)$

/* GPT, a decoder-only transformer, forward pass */

Input: $x \in V^*$, a sequence of token IDs.

Output: $P \in (0, 1)^{N_V \times \text{length}(x)}$, where the t -th column of P represents $\hat{P}_\theta(x[t+1] | x[1:t])$.

Hyperparameters: $\ell_{\max}, L, H, d_e, d_{\text{mlp}} \in \mathbb{N}$

Parameters: θ includes all of the following parameters:

$W_e \in \mathbb{R}^{d_e \times N_V}$, $W_p \in \mathbb{R}^{d_e \times \ell_{\max}}$, the token and positional embedding matrices.

For $l \in [L]$:

| \mathcal{W}_l , multi-head attention parameters for layer l , see (4),

| $\gamma_l^1, \beta_l^1, \gamma_l^2, \beta_l^2 \in \mathbb{R}^{d_e}$, two sets of layer-norm parameters,

| $W_{\text{mlp1}}^l \in \mathbb{R}^{d_{\text{mlp}} \times d_e}$, $b_{\text{mlp1}}^l \in \mathbb{R}^{d_{\text{mlp}}}$, $W_{\text{mlp2}}^l \in \mathbb{R}^{d_e \times d_{\text{mlp}}}$, $b_{\text{mlp2}}^l \in \mathbb{R}^{d_e}$, MLP parameters.

$\gamma, \beta \in \mathbb{R}^{d_e}$, final layer-norm parameters.

$W_u \in \mathbb{R}^{N_V \times d_e}$, the unembedding matrix.

```

1  $\ell \leftarrow \text{length}(x)$ 
2 for  $t \in [\ell]$  :  $e_t \leftarrow W_e[:, x[t]] + W_p[:, t]$ 
3  $X \leftarrow [e_1, e_2, \dots, e_\ell]$ 
4 for  $l = 1, 2, \dots, L$  do
5   for  $t \in [\ell]$  :  $\tilde{X}[:, t] \leftarrow \text{layer\_norm}(X[:, t] | \gamma_l^1, \beta_l^1)$ 
6    $X \leftarrow X + \text{MHAttention}(\tilde{X} | \mathcal{W}_l, \text{Mask}[t, t'] = [[t \leq t']])$ 
7   for  $t \in [\ell]$  :  $\tilde{X}[:, t] \leftarrow \text{layer\_norm}(X[:, t] | \gamma_l^2, \beta_l^2)$ 
8    $X \leftarrow X + W_{\text{mlp2}}^l \text{GELU}(W_{\text{mlp1}}^l \tilde{X} + b_{\text{mlp1}}^l \mathbf{1}^\top) + b_{\text{mlp2}}^l \mathbf{1}^\top$ 
9 end
10 for  $t \in [\ell]$  :  $X[:, t] \leftarrow \text{layer\_norm}(X[:, t] | \gamma, \beta)$ 
11 return  $P = \text{softmax}(W_u X)$ 

```

Algorithm 11: $\hat{\theta} \leftarrow \text{EDTraining}(z_{1:N_{\text{data}}}, x_{1:N_{\text{data}}}, \theta)$

/* Training a seq2seq model */

Input: $\{(z_n, x_n)\}_{n=1}^{N_{\text{data}}}$, a dataset of sequence pairs.

Input: θ , initial transformer parameters.

Output: $\hat{\theta}$, the trained parameters.

Hyperparameters: $N_{\text{epochs}} \in \mathbb{N}$, $\eta \in (0, \infty)$

```

1 for  $i = 1, 2, \dots, N_{\text{epochs}}$  do
2   for  $n = 1, 2, \dots, N_{\text{data}}$  do
3      $\ell \leftarrow \text{length}(x_n)$ 
4      $P(\theta) \leftarrow \text{EDTransformer}(z_n, x_n | \theta)$ 
5     loss( $\theta$ ) =  $-\sum_{t=1}^{\ell-1} \log P(\theta)[x_n[t+1], t]$ 
6      $\theta \leftarrow \theta - \eta \cdot \nabla \text{loss}(\theta)$ 
7   end
8 end
9 return  $\hat{\theta} = \theta$ 
```

Algorithm 13: $\hat{\theta} \leftarrow \text{DTraining}(x_{1:N_{\text{data}}}, \theta)$

/* Training next token prediction */

Input: $\{x_n\}_{n=1}^{N_{\text{data}}}$, a dataset of sequences.

Input: θ , initial decoder-only transformer parameters.

Output: $\hat{\theta}$, the trained parameters.

Hyperparameters: $N_{\text{epochs}} \in \mathbb{N}$, $\eta \in (0, \infty)$

```

1 for  $i = 1, 2, \dots, N_{\text{epochs}}$  do
2   for  $n = 1, 2, \dots, N_{\text{data}}$  do
3      $\ell \leftarrow \text{length}(x_n)$ 
4      $P(\theta) \leftarrow \text{DTransformer}(x_n | \theta)$ 
5     loss( $\theta$ ) =  $-\sum_{t=1}^{\ell-1} \log P(\theta)[x_n[t+1], t]$ 
6      $\theta \leftarrow \theta - \eta \cdot \nabla \text{loss}(\theta)$ 
7   end
8 end
9 return  $\hat{\theta} = \theta$ 
```

Algorithm 12: $\hat{\theta} \leftarrow \text{ETraining}(x_{1:N_{\text{data}}}, \theta)$

/* Training by masked language modelling */

Input: $\{x_n\}_{n=1}^{N_{\text{data}}}$, a dataset of sequences.

Input: θ , initial encoder-only transformer parameters.

Output: $\hat{\theta}$, the trained parameters.

Hyperparameters: $N_{\text{epochs}} \in \mathbb{N}$, $\eta \in (0, \infty)$, $p_{\text{mask}} \in (0, 1)$

```

1 for  $i = 1, 2, \dots, N_{\text{epochs}}$  do
2   for  $n = 1, 2, \dots, N_{\text{data}}$  do
3      $\ell \leftarrow \text{length}(x_n)$ 
4     for  $t = 1, 2, \dots, \ell$  do
5        $\tilde{x}_n[t] \leftarrow \text{mask\_token or } x_n[t]$ 
       randomly with probability  $p_{\text{mask}}$  or  $1 - p_{\text{mask}}$ 
6     end
7      $\tilde{T} \leftarrow \{t \in [\ell] : \tilde{x}_n[t] = \text{mask\_token}\}$ 
8      $P(\theta) \leftarrow \text{ETraining}(\tilde{x}_n | \theta)$ 
9     loss( $\theta$ ) =  $-\sum_{t \in \tilde{T}} \log P(\theta)[x_n[t], t]$ 
10     $\theta \leftarrow \theta - \eta \cdot \nabla \text{loss}(\theta)$ 
11  end
12 end
13 return  $\hat{\theta} = \theta$ 
```

Algorithm 14: $y \leftarrow \text{DIference}(x, \hat{\theta})$

/* Prompting a trained model and using it for prediction. */

Input: Trained transformer parameters $\hat{\theta}$.

Input: $x \in V^*$, a prompt.

Output: $y \in V^*$, the transformer's continuation of the prompt.

Hyperparameters: $\ell_{\text{gen}} \in \mathbb{N}$, $\tau \in (0, \infty)$

```

1  $\ell \leftarrow \text{length}(x)$ 
2 for  $i = 1, 2, \dots, \ell_{\text{gen}}$  do
3    $P \leftarrow \text{DTransformer}(x | \hat{\theta})$ 
4    $p \leftarrow P[:, \ell + i - 1]$ 
5   sample a token  $y$  from  $q \propto p^{1/\tau}$ 
6    $x \leftarrow [x, y]$ 
7 end
8 return  $y = x[\ell + 1 : \ell + \ell_{\text{gen}}]$ 
```

Algorithm 15: $\hat{x} \leftarrow \text{EDIference}(z, \hat{\theta})$

```
/* Using a trained seq2seq model for
   prediction. */  
Input: A seq2seq transformer and trained
   parameters  $\hat{\theta}$  of the transformer.  
Input:  $z \in V^*$ , input sequence, e.g. a
   sentence in English.  
Output:  $\hat{x} \in V^*$ , output sequence, e.g. the
   sentence in German.  
Hyperparameters:  $\tau \in (0, \infty)$   
1  $\hat{x} \leftarrow [\text{bos\_token}]$   
2  $y \leftarrow 0$   
3 while  $y \neq \text{eos\_token}$  do  
4    $P \leftarrow \text{EDTransformer}(z, \hat{x} | \hat{\theta})$   
5    $p \leftarrow P[:, \text{length}(\hat{x})]$   
6   sample a token  $y$  from  $q \propto p^{1/\tau}$   
7    $\hat{x} \leftarrow [\hat{x}, y]$   
8 end  
9 return  $\hat{x}$ 

---


```

B. List of Notation

Symbol	Type	Explanation
$[N]$	$:= \{1, \dots, N\}$	set of integers $1, 2, \dots, N - 1, N$
i, j	$\in \mathbb{N}$	generic integer indices
V	$\cong [N_V]$	vocabulary
N_V	$\in \mathbb{N}$	vocabulary size
V^*	$= \bigcup_{\ell=0}^{\infty} V^\ell$	set of token sequences; elements include e.g. sentences or documents
ℓ_{\max}	$\in \mathbb{N}$	maximum sequence length
ℓ	$\in [\ell_{\max}]$	length of token sequence
t	$\in [\ell]$	index of token in a sequence
d_{\dots}	$\in \mathbb{N}$	dimension of various vectors
x	$\equiv x[1 : \ell]$	$\equiv x[1]x[2]...x[\ell] \in V^\ell$ primary token sequence
z	$\equiv z[1 : \ell]$	$\equiv z[1]z[2]...z[\ell] \in V^\ell$ context token sequence
$M[i, j]$	$\in \mathbb{R}$	entry M_{ij} of matrix $M \in \mathbb{R}^{d \times d'}$
$M[i, :] \equiv M[i]$	$\in \mathbb{R}^{d'}$	i -th row of matrix $M \in \mathbb{R}^{d \times d'}$
$M[:, j]$	$\in \mathbb{R}^d$	j -th column of matrix $M \in \mathbb{R}^{d \times d'}$
e	$\in \mathbb{R}^{d_e}$	vector representation / embedding of a token
X	$\in \mathbb{R}^{d_e \times \ell_x}$	encoded primary token sequence
Z	$\in \mathbb{R}^{d_e \times \ell_z}$	encoded context token sequence
Mask	$\in \mathbb{R}^{\ell_z \times \ell_x}$	masking matrix, it determines the attention context for each token
$L, L_{\text{enc}}, L_{\text{dec}}$	$\in \mathbb{N}$	number of network (encoder, decoder) layers
l	$\in [L]$	index of network layer
H	$\in \mathbb{N}$	number of attention heads
h	$\in [H]$	index of attention head
N_{data}	$\in \mathbb{N}$	(i.i.d.) sample size
n	$\in [N_{\text{data}}]$	index of sample sequence
η	$\in (0, \infty)$	learning rate
τ	$\in (0, \infty)$	temperature; it controls the diversity-plausibility trade-off at inference
W_e	$\in \mathbb{R}^{d_e \times N_V}$	token embedding matrix
W_p	$\in \mathbb{R}^{d_e \times \ell_{\max}}$	positional embedding matrix
W_u	$\in \mathbb{R}^{N_V \times d_e}$	unembedding matrix
W_q	$\in \mathbb{R}^{d_{\text{attn}} \times d_x}$	query weight matrix
b_q	$\in \mathbb{R}^{d_{\text{attn}}}$	query bias
W_k	$\in \mathbb{R}^{d_{\text{attn}} \times d_z}$	key weight matrix
b_k	$\in \mathbb{R}^{d_{\text{attn}}}$	key bias
W_v	$\in \mathbb{R}^{d_{\text{out}} \times d_z}$	value weight matrix
b_v	$\in \mathbb{R}^{d_{\text{out}}}$	value bias
\mathcal{W}_{qkv}		collection of above parameters of a single-head attention layer
W_o	$\in \mathbb{R}^{d_{\text{out}} \times H d_{\text{mid}}}$	output weight matrix
b_o	$\in \mathbb{R}^{d_{\text{out}}}$	output bias
\mathcal{W}		collection of above parameters of a multi-head attention layer, see eq. (4)
W_{mlp}	$\in \mathbb{R}^{d_1 \times d_2}$	weight matrix corresponding to an MLP layer in a Transformer
b_{mlp}	$\in \mathbb{R}^{d_1}$	bias corresponding to an MLP layer in a Transformer
γ	$\in \mathbb{R}^{d_e}$	layer-norm learnable scale parameter
β	$\in \mathbb{R}^{d_e}$	layer-norm learnable offset parameter
$\theta, \hat{\theta}$	$\in \mathbb{R}^d$	collection of all learnable / learned Transformer parameters