

# Decoder-Only Transformers

Notes on various aspects of Decoder-Only Transformers.

## Contents

<b>1</b>	<b>Decoder-Only Fundamentals</b>	<b>1</b>
1.1	Component Details	2
1.1.1	Embedding Layer and Positional Encodings	3
1.1.2	Layer Norm	3
1.1.3	Causal Attention	3
1.1.4	MLP	5
1.1.5	Language Model Head	6
1.1.6	All Together	6
<b>2</b>	<b>Memory and Activations</b>	<b>8</b>
2.1	No Sharding/Parallelism	8
2.1.1	MLP Activations	8
2.2	Optimizer States and Mixed Precision	9
<b>A</b>	<b>Conventions and Notation</b>	<b>9</b>
<b>B</b>	<b>TODO</b>	<b>10</b>

## 1 Decoder-Only Fundamentals

The Transformers architecture [1], which dominates Natural Language Processing (NLP) as of July 2023, is a relatively simple architecture. There are various flavors and variants of Transformers, but focus here on the decoder-only versions which underlie the GPT models [2–4].

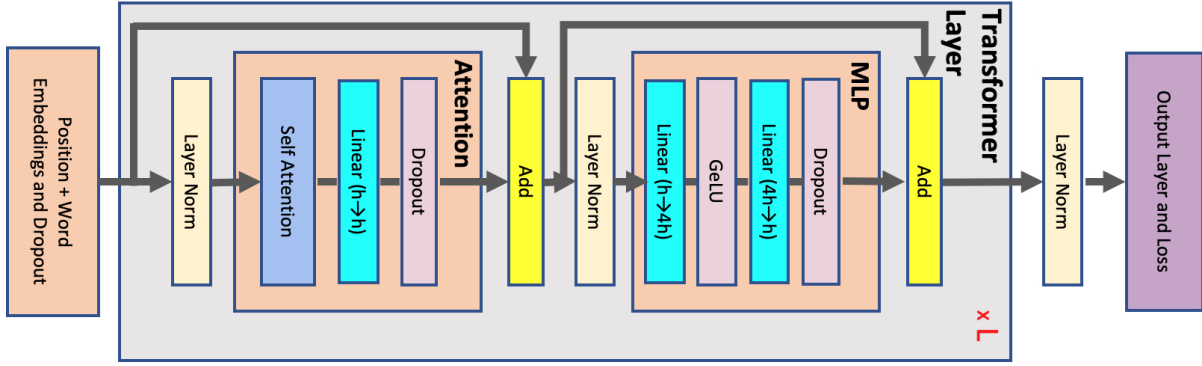
The full decoder-only architecture can be seen in Fig. 1. The parameters which define the network can be found in App. A.

An outline of the mechanics:

1. Raw text is **tokenized** and turned into a series of integers<sup>1</sup> whose values lie in  $\text{range}(V)$ , with  $V$  the vocabulary size.
2. The tokenized text is chunked and turned into  $(B, S)$ -shaped (batch size and sequence length, respectively) integer tensors,  $x_{bs}$ .
3. The **embedding layer** converts the integer tensors into continuous representations of shape  $(B, S, D)$ ,  $z_{bsd}$ , with  $D$  the size of the hidden dimension. **Positional encodings** have

---

<sup>1</sup>There are about 1.3 tokens per word, on average.



**Figure 1.** The full transformers architecture. Diagram taken from [5]

also been added to the tensor at this stage to help the architecture understand the relative ordering of the text.

4. The  $z_{bsd}$  tensors pass through a series of transformer blocks, each of which has two primary components:
  - (a) In the **attention** sub-block, components of  $z_{bsd}$  at different positions ( $s$ -values) interact with each other, resulting in another  $(B, S, D)$ -shaped tensor,  $z'_{bsd}$ .
  - (b) In the **MLP** block, each position in  $z'_{bsd}$  is processed independently and in parallel by a two-layer feed-forward network, resulting once more in a  $(B, S, D)$ -shaped tensor.

Importantly, there are **residual connections** around each of these<sup>2</sup> (the arrows in Fig. 1).

5. Finally, we convert the  $(B, S, D)$ -shaped tensors to  $(B, S, V)$ -shaped ones,  $y_{bsv}$ . This is the role of the **language model head** (which is often just the embedding layer used in an inverse manner.)
6. The  $y_{bsv}$  predict what the next token will be, having seen the **context** of the first  $s$  tokens in the sequence.

Each batch (the  $b$ -index) is processed independently. We omitted **LayerNorm** and **Dropout** layers above, as well as the causal mask; these will be covered below as we step through the architecture in more detail.

### 1.1 Component Details

We break down the various components below.

<sup>2</sup>This gives rise to the concept of the **residual stream** which each transformer block reads from and writes back to repeatedly.

### 1.1.1 Embedding Layer and Positional Encodings

The **embedding** layer is just a simple look up table: each of the `range(V)` indices in the vocabulary is mapped to a  $D$ -dimensional vector via a large  $(V, D)$ -shaped table/matrix. This layer maps  $x_{bs} \rightarrow z_{bsd}$ . In `torch`, this is an `nn.Embedding(V, D)` instance.

To each item in a batch, we add identical **positional encodings** to the vectors above with the goal of adding fixed, position-dependent correlations in the sequence dimension which will hopefully make it easier for the architecture to pick up on the relative positions of the inputs<sup>3</sup>. This layer maps  $z_{bsd} \leftarrow z_{bsd} + p_{sd}$ , with  $p_{sd}$  the positional encoding tensor.

The above components require  $(V + S)D \approx VD$  parameters per layer.

### 1.1.2 Layer Norm

The original transformers paper [1] put **LayerNorm** instances after the **attention** and **MLP** blocks, but now it is common [6] to put them before these blocks<sup>4</sup>.

The **LayerNorm** operations acts over the sequence dimension. Spelling it out, given the input tensor  $z_{bsd}$  whose mean and variance over the  $s$ -index are  $\mu_{bd}$  and  $\sigma_{bd}$ , respectively, the **LayerNorm** output is

$$z_{bsd} \leftarrow \left( \frac{z_{bsd} - \mu_{bd}}{\sigma_{bd}} \right) \times \gamma_d + \beta_d \equiv \text{LayerNorm}_s z_{bsd} \quad (1.1)$$

where  $\gamma_d, \beta_d$  are the trainable scale and bias parameters. In `torch`, this is a `nn.LayerNorm(D)` instance.

Since there are two **LayerNorm** instances in each transformer block, these components require  $2D$  parameters per layer.

### 1.1.3 Causal Attention

The **causal attention** layer is the most complex layer. It features  $H$  triplets<sup>5</sup> of weight matrices<sup>6</sup>  $Q_{df}^h, K_{df}^h, V_{df}^h$  where  $a \in \{0, \dots, H - 1\}$  and  $f \in \{0, \dots, D/H\}$ . From these, we form three different vectors:

$$q_{bsf}^h = z_{bsd} Q_{df}^h, \quad k_{bsf}^h = z_{bsd} K_{df}^h, \quad v_{bsf}^h = z_{bsd} V_{df}^h \quad (1.2)$$

These are the **query**, **key**, and **value** tensors, respectively.

Using the above tensors, we will then build up an **attention map**  $w_{bss'}^h$ , which corresponds to how much attention the token at position  $s$  pays to the token at position  $s'$ . Because we have the goal of predicting the next token in the sequence, we need these weights to be causal: the final

<sup>3</sup>Positional encodings and the causal mask are the only components in the transformers architecture which carry weights with a dimension of size  $S$ ; i.e. they are the only parts that have explicit sequence-length dependence. A related though experiment: you can convince yourself that if the inputs  $z_{bsd}$  were just random noise, the transformers architecture would not be able to predict the  $s$ -index of each such input in the absence of positional encodings.

<sup>4</sup>Which makes intuitive sense for the purposes of stabilizing the matrix multiplications in the blocks

<sup>5</sup> $H$  must divide the hidden dimension  $D$  evenly.

<sup>6</sup>There are also bias terms, but we will often neglect to write them explicitly or account for their (negligible) parameter count.

prediction  $y_{bsv}$  should only have access to information propagated from positions  $x_{bs'v}$  with  $s' \leq s$ . This corresponds to the condition that  $w_{bs's'}^h = 0$  if  $s' > s$ .

These weights come from **Softmax**-ed attention scores, which are just a normalized dot-product over the hidden dimension:

$$w_{bs's'd}^h = \text{Softmax}_{s'} \left( m_{ss'} + q_{bs'f}^h k_{bs'f}^h / \sqrt{D/H} \right), \quad \text{s.t.} \quad \sum_{s'} w_{bs's's'}^h = 1 \quad (1.3)$$

The tensor  $m_{ss'}$  is the causal mask which zeroes out the relevant attention map components above

$$m_{ss'} = \begin{cases} 0 & s \leq s' \\ -\infty & s > s' \end{cases}$$

The  $\sqrt{D/H}$  normalization is motivated by demanding that the variance of the **Softmax** argument be 1 at initialization, assuming that other components have been configured so that the query and key components are i.i.d. from a Gaussian normal distribution<sup>7</sup>.

The weights above are then passed through a dropout layer and used to re-weight the **value** vectors and form the tensors

$$t_{bs'f}^h = \text{Drop} \left( w_{bs's's'}^h v_{bs'f}^h \right) \quad (1.4)$$

and these  $H$  different (B, S, D/H)-shaped tensors are then concatenated along the  $f$ -direction to re-form a (B, S, D)-shaped tensor<sup>8</sup>

$$u_{bsd} = \text{Concat}_{fd}^a (t_{bs'f}^a) . \quad (1.5)$$

Finally, another weight matrix  $O_{d'd}$  and dropout layer transform the output once again to get the final output

$$z_{bsd} = \text{Drop} (u_{bsd'} O_{d'd}) . \quad (1.6)$$

For completeness, the entire operation in condensed notation with indices left implicit is:

$$z \leftarrow \text{Drop} \left( \text{Concat} \left( \text{Drop} \left( \text{Softmax} \left( \frac{(z \cdot Q^h) \cdot (z \cdot K^h)}{\sqrt{D/H}} \right) \right) \cdot z \cdot V^h \right) \cdot O \right) \quad (1.7)$$

where all of the dot-products are over feature dimensions (those of size  $D$  or  $D/H$ ). H nice<sup>9</sup>, but hopefully pedagogically useful, implementation of the attention layer is below:

```

10 class CausalAttention(nn.Module):
11     def __init__(
12         self,
```

<sup>7</sup>However, in [7] it is instead argued that no square root should be taken in order to maximize the speed of learning via SGD.

<sup>8</sup>It is hard to come up with good index-notation for concatenation.

<sup>9</sup>An example optimization: there is no need to form separate  $Q^h, K^h, V^h$  **Linear** layers, one large layer which is later chunked is more efficient

```

13         attn_heads=H,
14         hidden_dim=D,
15         block_size=K,
16         dropout=0.1,
17     ):
18         super().__init__()
19         self.head_dim, remainder = divmod(hidden_dim, attn_heads)
20         assert not remainder, "attn_heads must divide hidden_dim evenly"
21
22         self.Q = nn.ModuleList([nn.Linear(hidden_dim, self.head_dim) for _ in range(attn_heads)])
23         self.K = nn.ModuleList([nn.Linear(hidden_dim, self.head_dim) for _ in range(attn_heads)])
24         self.V = nn.ModuleList([nn.Linear(hidden_dim, self.head_dim) for _ in range(attn_heads)])
25         self.O = nn.Linear(hidden_dim, hidden_dim)
26
27         self.attn_dropout = nn.Dropout(dropout)
28         self.final_dropout = nn.Dropout(dropout)
29         self.register_buffer(
30             "causal_mask",
31             torch.tril(torch.ones(block_size, block_size)[None]),
32         )
33
34     def get_qkv(self, inputs):
35         queries = [q(inputs) for q in self.Q]
36         keys = [k(inputs) for k in self.K]
37         values = [v(inputs) for v in self.V]
38         return queries, keys, values
39
40     def get_attn_maps(self, queries, keys, values, seq_len):
41         norm = math.sqrt(self.head_dim)
42         non_causal_attn_scores = [(q @ k.transpose(-2, -1)) / norm for q, k in zip(queries, keys)]
43         causal_attn_scores = [
44             a.masked_fill(self.causal_mask[:, :S, :S] == 0, float("-inf"))
45             for a in non_causal_attn_scores
46         ]
47         attn_maps = [a.softmax(dim=-1) for a in causal_attn_scores]
48         return attn_maps
49
50     def forward(self, inputs):
51         S = inputs.shape[1]
52         queries, keys, values = self.get_qkv(inputs)
53         attn_maps = self.get_attn_maps(queries, keys, values, S)
54         weighted_values = torch.concatenate(
55             [self.attn_dropout(a) @ v for a, v in zip(attn_maps, values)], dim=-1
56         )
57         z = self.final_dropout(self.O(weighted_values))

```

The parameter count is dominated by the weight matrices which carry  $4D^2$  total parameters per layer.

#### 1.1.4 MLP

The feed-forward network is straightforward and corresponds to

$$z_{bsd} \leftarrow \phi(z_{bsd'} W_{d'e}^0) W_{ed}^1 \quad (1.8)$$

where  $W^0$  and  $W^1$  are  $(D, FD)$ - and  $(FD, D)$ -shaped matrices, respectively (see App. A for notation) and  $\phi$  is a non-linearity<sup>10</sup>. The implementation is straightforward:

```

8  class MLP(nn.Module):
9      def __init__(
10         self,
11         hidden_dim=D,
12         expansion_factor=F,
13         dropout=0.1,
14     ):
15         super().__init__()
16         linear_1 = nn.Linear(hidden_dim, expansion_factor * hidden_dim)
17         linear_2 = nn.Linear(expansion_factor * hidden_dim, hidden_dim)
18         gelu = nn.GELU()
19         drop = nn.Dropout(dropout)
20         self.layers = nn.Sequential(linear_1, gelu, linear_2, drop)
21
22     def forward(self, inputs):
23         z = self.layers(inputs)
24         return z

```

This block requires  $2FD^2$  parameters per layer, only counting the contribution from weights.

### 1.1.5 Language Model Head

The layer which converts the  $(B, S, D)$ -shaped outputs,  $z_{bsd}$ , to  $(B, S, V)$ -shaped predictions over the vocabulary,  $y_{bsv}$ , is the **Language Model Head**. It is a linear layer, whose weights are usually tied to be exactly those of the initial embedding layer of Sec. 1.1.1.

### 1.1.6 All Together

It is then relatively straightforward to tie everything together. In code, we can first create a transformer block like

```

10  class TransformerBlock(nn.Module):
11      def __init__(
12         self,
13         attn_heads=H,
14         block_size=K,
15         dropout=0.1,
16         expansion_factor=F,
17         hidden_dim=D,
18         layers=L,
19         vocab_size=V,
20     ):
21         super().__init__()
22         self.attn_ln = nn.LayerNorm(hidden_dim)
23         self.mlp_ln = nn.LayerNorm(hidden_dim)
24         self.attn = CausalAttention(attn_heads, hidden_dim, block_size, dropout)
25         self.mlp = MLP(hidden_dim, expansion_factor, dropout)
26

```

---

<sup>10</sup>The GeLU non-linearity is common.

```

27     def forward(self, inputs):
28         z = self.attn(self.attn_ln(inputs)) + inputs
29         z = self.mlp(self.mlp_ln(z)) + z
30         return z

```

And then the entire architecture:

```

9  class DecoderOnly(nn.Module):
10     def __init__(
11         self,
12         attn_heads=H,
13         block_size=K,
14         dropout=0.1,
15         expansion_factor=F,
16         hidden_dim=D,
17         layers=L,
18         vocab_size=V,
19     ):
20         super().__init__()
21         self.embedding = nn.Embedding(vocab_size, hidden_dim)
22         self.pos_encoding = nn.Parameter(torch.randn(1, block_size, hidden_dim))
23         self.drop = nn.Dropout(dropout)
24         self.trans_blocks = nn.ModuleList(
25             [
26                 TransformerBlock(
27                     attn_heads,
28                     block_size,
29                     dropout,
30                     expansion_factor,
31                     hidden_dim,
32                     layers,
33                     vocab_size,
34                 )
35                 for _ in range(layers)
36             ]
37         )
38         self.final_ln = nn.LayerNorm(hidden_dim)
39         self.lm_head = nn.Linear(hidden_dim, vocab_size, bias=False)
40         self.lm_head.weight = self.embedding.weight # Weight tying.
41
42     def forward(self, inputs):
43         S = inputs.shape[1]
44         z = self.embedding(inputs) + self.pos_encoding[:, :S]
45         z = self.drop(z)
46         for block in self.trans_blocks:
47             z = block(z)
48         z = self.final_ln(z)
49         z = self.lm_head(z)
50         return z
51

```

## 2 Memory and Activations

In this section we summarize the train-time memory costs of Transformers under various training strategies. At a high level, the memory cost is much more than simply the cost of the model parameters. Significant factors include:

- Optimizer states, like those of Adam
- Mixed precision training costs, due to keeping multiple model copies.
- Activation memory<sup>11</sup>, needed for backpropagation.
- Gradients

Activations require the most detailed analysis, so we start with their analysis.

### 2.1 No Sharding/Parallelism

Start with the simplest case where no parallelism are used. The costs of training then come from the model parameters, the optimizer state, the gradients, and the activations. Every parameter is assumed to be stored in `fp32`, i.e. four bytes per parameter<sup>12</sup>.

We will assume the use of Adam throughout, for simplicity, which stores `two different running averages` per-model parameter. Even in this vanilla scenario, the cost of the optimizer states is significant.

Next, we eventually will get gradients out of back propagation, one-per parameter, and so the gradient cost is also equal to the model weight cost.

Finally, activations, which require more analysis [5]. Unlike the above factors, the activation memory will depend on both the batch size and input sequence length,  $B$  and  $S$ , scaling linearly with both.

#### 2.1.1 MLP Activations

We will count the number of elements which need to be cached. First we cache the  $(B, S, D)$ -shaped inputs into the first MLP layer. These turn into the  $(B, S, FD)$  inputs of the non-linearity, which are then passed into the last `Linear` layer. This sums to the following number of cached activations elements:

$$N_{\text{act}}^{\text{MLP}} = (2F + 1)BS. \quad (2.1)$$

The `surprisingly`

---

<sup>11</sup>Activations refers to any intermediate value which needs to be cached in order to compute backpropagation. We will be conservative and assume that the inputs of all operations need to be stored, though in practice gradient checkpointing and recomputation allow one to trade caching for redundant compute. In particular, flash attention [8] makes use of this strategy.

<sup>12</sup>We will usually express memory in terms of GiB, where  $1 \text{ GiB} \equiv 2^{30}$  bytes.



## 2.2 Optimizer States and Mixed Precision

From the previous section, the pure parameter-count of the decoder-only Transformers architecture is

$$N_{\text{params}} \approx (4 + 2F)LD^2 + VD + \mathcal{O}(LD) \quad (2.2)$$

where the first term comes from the `TransformerBlock` weight matrices, the next term is the embedding matrix, and the neglected terms come from biases, `LayerNorm` instances, and other negligible factors.<sup>13</sup>

The memory cost of the bare parameters are usually dwarfed

## A Conventions and Notation

We loosely follow the conventions of [5] and denote the main Transformers parameters by:

- $B$ : microbatch size
- $K$ : the block size (maximum sequence length<sup>14</sup>)
- $S$ : input sequence length
- $V$ : vocabulary size
- $D$ : the hidden dimension size
- $L$ : number of transformer layers
- $H$ : number of attention heads
- $P$ : pipeline parallel size
- $T$ : tensor parallel size
- $F$ : expansion factor for MLP layer (usually  $F = 4$ )

Where it makes sense, we try to use the lower-case versions of these characters to denote the corresponding indices on various tensors. For instance, an input tensor with the above batch size, sequence length, and vocabulary size would be written as  $x_{bsv}$ , with  $b \in \{0, \dots, B-1\}$ ,  $s \in \{0, \dots, S-1\}$ , and  $v \in \{0, \dots, V-1\}$  in math notation, or as `x[b, s, v]` in code. Typical transformers belong to the regime

$$V \gg D, S \gg L, H \gg P, T. \quad (\text{A.1})$$

As indicated above, we use zero-indexing. We also use `python` code throughout<sup>15</sup> and write all ML code using standard `torch` syntax. To avoid needing to come up with new symbols in math

<sup>13</sup>For the usual  $F = 4$  case, the MLP layers are twice as costly as the `CausalAttention` layers.

<sup>14</sup>In the absence of methods such as ALiBi [9] can be used to extend the sequence length at inference time.

<sup>15</sup>Written in a style conducive to latex, e.g. no type-hints and pegagogy prioritized.

expressions we will often use expressions like  $x \leftarrow f(x)$  to refer to performing a computation on some argument ( $x$ ) and assigning the result right back to the variable  $x$  again.

Physicists often joke (half-seriously) that Einstein’s greatest contribution to physics was his summation notation in which index-sums are implied by the presence of repeated indices and summation symbols are entirely omitted. For instance, the dot product between two vectors would be written as

$$\vec{x} \cdot \vec{y} = \sum_i x_i y_i \equiv x_i y_i \quad (\text{A.2})$$

We use similar notation which is further adapted to the common element-wise deep-learning operations. The general rule is that if a repeated index appears on one side of an equation, but not the other, then a sum is implied, but if the same index appears on both sides, then it’s an element-wise operation. The Hadamard-product between two matrices  $A$  and  $B$  is just

$$C_{ij} = A_{ij} B_{ij} . \quad (\text{A.3})$$

We also put explicit indices on operators such as Softmax to help clarify the relevant dimension, e.g. we would write the softmax operation over the  $b$ -index of some batched tensor  $x_{bvd\dots}$  as

$$s_{bvd\dots} = \frac{e^{x_{bvd\dots}}}{\sum_{v=0}^{V-1} e^{x_{bvd\dots}}} \equiv \text{Softmax}_v x_{bvd\dots} , \quad (\text{A.4})$$

indicating that the sum over the singled-out  $v$ -index is gives unity.

## B TODO

- Tokenizers
- Generation
- Activations
- Flash attention
- BERT family
- Residual stream
- Scaling laws

## References

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” [arXiv:1706.03762 \[cs.CL\]](#). [1](#), [3](#)
- [2] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, *et al.*, “Language models are unsupervised multitask learners,” *OpenAI blog* **1** (2019) no. 8, 9. [1](#)

- [3] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, “Language models are few-shot learners,” [arXiv:2005.14165](#) [cs.CL].
- [4] OpenAI, “Gpt-4 technical report,” [arXiv:2303.08774](#) [cs.CL]. 1
- [5] V. Korthikanti, J. Casper, S. Lym, L. McAfee, M. Andersch, M. Shoeybi, and B. Catanzaro, “Reducing activation recomputation in large transformer models,” [arXiv:2205.05198](#) [cs.LG]. 2, 8, 9
- [6] R. Xiong, Y. Yang, D. He, K. Zheng, S. Zheng, C. Xing, H. Zhang, Y. Lan, L. Wang, and T.-Y. Liu, “On layer normalization in the transformer architecture,” [arXiv:2002.04745](#) [cs.LG]. 3
- [7] G. Yang, E. J. Hu, I. Babuschkin, S. Sidor, X. Liu, D. Farhi, N. Ryder, J. Pachocki, W. Chen, and J. Gao, “Tensor programs v: Tuning large neural networks via zero-shot hyperparameter transfer,” [arXiv:2203.03466](#) [cs.LG]. 4
- [8] T. Dao, D. Y. Fu, S. Ermon, A. Rudra, and C. Ré, “Flashattention: Fast and memory-efficient exact attention with io-awareness,” [arXiv:2205.14135](#) [cs.LG]. 8
- [9] O. Press, N. A. Smith, and M. Lewis, “Train short, test long: Attention with linear biases enables input length extrapolation,” *CoRR* **abs/2108.12409** (2021) , 2108.12409. <https://arxiv.org/abs/2108.12409>. 9