

DDA5001 Machine Learning

Transformer (Part II)

Xiao Li

School of Data Science
The Chinese University of Hong Kong, Shenzhen



Information about Final

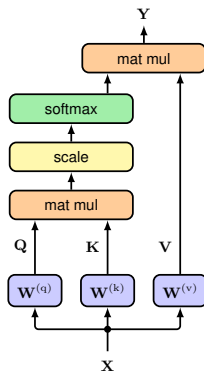
- ▶ Time: December 13 (Saturday), 3:00pm - 5:00pm (2 hours).
- ▶ Venue: Liwen Hall.
- ▶ A closed-book, closed-notes examination. No cheat sheet (this leads to an easier exam). No other things will be allowed, such as calculators.
- ▶ A sample exam will be available later to aid you in your review. Note that do not go overfitting, you need a little generalization. Hence, you are highly recommended to go through our lecture slides and Homework questions.

Recap: Self-Attention

- **Attention** is to consider the corrections between different tokens in a sentence.
- We use self-attention to capture the semantic relationship:

$$Y = \text{attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V.$$

Illustration:



Recap: Multi-Head Attention, Attention Masks

Multi-head attention: To capture different semantic structure by different attentions:

$$\mathbf{Y}(\mathbf{X}) = [\mathbf{H}_1, \dots, \mathbf{H}_H] \mathbf{W}^{(o)}.$$

Each head \mathbf{H}_i is generated by attention.

Attention masks:

- We need attention masks to block attention to **padding tokens** and **future tokens** (for **causality** during training):

$$\mathbf{Y} = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} + \mathbf{M}^{\text{causal}} + \mathbf{M}^{\text{pad}} \right) \mathbf{V},$$

where

$$M_{ij} = \begin{cases} 0, & \text{allowed} \\ -\infty, & \text{blocked} \end{cases}$$

Complexity of Attention

Consider masked attention with n being sequence length.

- ▶ Embedding dimension d is often **fixed**, such as $d = 4096, 8192$.
- ▶ During pre-training, or reasoning modeling training, we often have **long-sequence** training, e.g., n can scale up to 128k. Therefore, we often focus on the scaling with n .

Complexity:

- ▶ **Score matrix**: $S = QK^T \in \mathbb{R}^{n \times n}$ with computation $O(n^2d)$.
- ▶ Masking + softmax on S : Store and softmax $n \times n$ scores with computation $O(n^2)$ (softmax) and storing memory $O(n^2)$ (score matrix and mask).
- ▶ Weighted sum: $Y = \text{softmax}(S)V$ with computation $O(n^2d)$.

Observations:

- ▶ Overall, computation and memory are both **quadratic** in sequence length n .
- ▶ Attention is **not** in linear complexity. One of the major drawback of Transformers.

Flash-Attention

- ▶ **Flash-attention** is a technique to reduce **storage memory** from $O(n^2)$ to $O(nd)$, i.e., nearly linear in n when n is very large.

The main idea:

- ▶ Rewrites attention as a sequence of **smaller block** matrix operations.
- ▶ **Never** materializes **the full** $n \times n$ score matrix in GPU memory.
- ▶ Uses an **online softmax** algorithm to keep numerical stability while processing tiles.

Complexity of flash-attention:

- ▶ **Compute**: Still $O(n^2d)$ (same as standard attention).
- ▶ **Memory**: Reduced from $O(n^2)$ to roughly $O(nd)$.

For short sequences (i.e., small n), the benefit is modest. Flash-attention is very useful to reduce GPU memory when n is very large.

Transformer

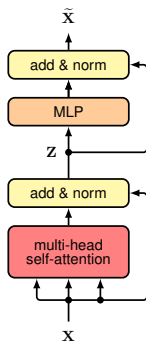
▪

Transformer Layer

Transformer Layer

- ▶ We now need to form a **transformer layer**.
- ▶ We want to add some **residual connection** to improve efficiency.
- ▶ We want to add **normalization** for efficiency and stability as well.
- ▶ We also want to add several MLP layer (feed-forward layer) to have more flexibility.

Overall, a transformer layer has the form:



- ▶ The mathematical modeling is:

$$Z = \text{Norm}(Y(X) + X).$$

$$\tilde{X} = \text{Norm}(\text{MLP}(Z) + Z).$$

- ▶ The dimension of the output \tilde{X} need to be the same as that of the input X .

Norm Method I: Layer Normalization (LayerNorm)

Consider $\mathbf{x} = (x_1, \dots, x_d)$ as one data sample, like one row of \mathbf{X} (i.e., one token). **LayerNorm** first does the following,

$$\hat{x}_i = \frac{x_i - \mu}{\sigma}, \quad \text{where} \quad \mu = \frac{1}{d} \sum_{i=1}^d x_i, \quad \sigma = \sqrt{\frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2}$$

- ▶ The above process is to make the data zero mean and unit variance.

It then adds **trainable** parameters $\gamma = (\gamma_1, \dots, \gamma_d)$ and $\beta = (\beta_1, \dots, \beta_d)$ to scale (γ) and shift (β) the obtained \hat{x}_i :

$$y_i = \gamma_i \hat{x}_i + \beta_i, \quad i = 1, \dots, d.$$

Finally, $\mathbf{y} = (y_1, \dots, y_d)$ will be the layer normalized data sample.

- ▶ LayerNorm can be regarded as giving a new **learned** mean and variance for each data.
- ▶ LayerNorm is important to improving stability of the training process.

Norm Method II: RMSNorm

Another popular norm is **Root Mean Squares Normalization (RMSNorm)**.

- ▶ First compute the root mean square (**RMS**)

$$\text{RMS}(\mathbf{x}) = \sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2 + \varepsilon}.$$

where ε is a small number for numerical stability.

- ▶ Then, apply **normalization** using RMS

$$\hat{x}_i = \frac{x_i}{\text{RMS}(\mathbf{x})}, \quad y_i = \gamma_i \hat{x}_i.$$

γ_i is the trainable scaling factor.

Comparison: LayerNorm both scale and shift, while **RMSNorm only scales**.

- ▶ Older / classic LLMs (GPT-2, GPT-3) uses LayerNorm.
- ▶ More recent very large LLMs uses RMSNorm.

Post-Norm and Pre-Norm

The way we outlined of applying normalization is called **post-norm**:

$$\begin{aligned}Z &= \text{Norm}(\mathbf{Y}(\mathbf{X}) + \mathbf{X}). \\ \widetilde{\mathbf{X}} &= \text{Norm}(\text{MLP}(\mathbf{Z}) + \mathbf{Z}).\end{aligned}$$

- ▶ It post-normalizes output after residual connection.
- ▶ This is the Norm used in the “Attention is All You Need” paper.

However, modern LLMs mainly use **pre-norm** which **pre-normalize the input**:

$$\begin{aligned}Z &= \mathbf{Y}(\text{Norm}(\mathbf{X})) + \mathbf{X}. \\ \widetilde{\mathbf{X}} &= \text{MLP}(\text{Norm}(\mathbf{Z})) + \mathbf{Z}.\end{aligned}$$

As it is found that pre-norm provides a more stable training process.

Complexity of One Transformer Layer

Consider post-norm transformer layer:

$$\mathbf{Z} = \text{Norm}(\mathbf{Y}(\mathbf{X}) + \mathbf{X}), \quad \widetilde{\mathbf{X}} = \text{Norm}(\text{MLP}(\mathbf{Z}) + \mathbf{Z}).$$

Notations. n : sequence length, d : model dimension, d_{MLP} : hidden size of MLP (often $d_{\text{MLP}} \approx 4d$).

- ▶ As we studied, the **attention module** has computation complexity $O(n^2d)$ and memory complexity $O(n^2)$.

In **MLP module**:

- ▶ Computation complexity:

$$O(ndd_{\text{MLP}}) \quad (\text{two linear layers})$$

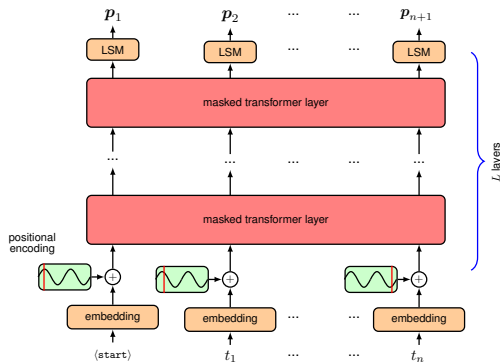
- ▶ Memory for activation values:

$$O(nd_{\text{MLP}}) \quad (\text{hidden activations})$$

So the overall order of complexity of the full Transformer layer is **quadratic in n** .

Transformer Architecture

Decoder Transformer Architecture



- ▶ This is decoder transformer, which is the architecture of GPT, Llama, and many other popular language models.
- ▶ Transformer is to stack transformer layers multiple times (just like one hidden layer to multiple), and add some other modules.

Tokenization

Goal: Map raw text to a sequence of discrete **tokens**

“It is raining.” $\longrightarrow (t_1, t_2, \dots, t_n)$.

- ▶ Is token in the unit of a character or a whole word? None of them.
 - ▶ Characters: very long sequences \Rightarrow inefficient.
 - ▶ Words: huge **vocabulary**, many rare/unknown words.
- ▶ Modern LLMs use **subword tokenization**:
 - ▶ Methods: BPE (GPT2), SentencePiece (Qwen3), WordPiece, etc.
 - ▶ Common words are single tokens; rare words are split:

“unbelievable” \rightarrow ‘‘un’’, ‘‘believ’’, ‘‘able’’

“raining” \rightarrow “rain”, “ing”

- ▶ Each token is represented internally by an **integer id**

$$t_i \in \{1, \dots, V\},$$

where V is the **vocabulary size** (e.g., $V \approx 30\text{k} \sim 128\text{k}$).

Embedding Layer

- ▶ The model cannot work directly with token ids t_i ; it needs vectors.
- ▶ **Embedding matrix shared** by all embedding positions:

$$\mathbf{E} \in \mathbb{R}^{V \times d}, \quad \text{row } \mathbf{E}_{t_i:} = \mathbf{x}_i^\top \in \mathbb{R}^{1 \times d}.$$

Each token id t_i is mapped to an embedding vector \mathbf{x}_i :

$$t_i \longrightarrow \mathbf{x}_i = \mathbf{E}^\top \mathbf{e}_{t_i},$$

where \mathbf{e}_{t_i} is the one-hot vector of t_i , i.e., taking out the corresponding row of \mathbf{E} .

- ▶ Embedding matrix \mathbf{E} is **learned**, rather than designed.
- ▶ After adding **positional encodings** (represent order of words in your sentence), the sequence

$$\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

is fed into the stack of masked transformer layers shown earlier.

The Last Layer: Logistic Regression

- ▶ In the last LSM layer, it means linear transformation + softmax, i.e., multi-class logistic regression on the final learned feature $\mathbf{Z} \in \mathbb{R}^{(n+1) \times d}$, the added token is <start>.
- ▶ The number of class is V , i.e., the vocabulary size. Using classification to choose the next token from the vocabulary.
- ▶ The classification is done using logistic regression with lm_head $\mathbf{W}_{\text{lm}} \in \mathbb{R}^{d \times V}$, i.e., linear classification weight matrix.

Mathematically,

$$\mathbf{P} = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_{n+1}]^\top = \text{softmax}(\mathbf{Z}\mathbf{W}_{\text{lm}}) = \text{softmax}(\mathbf{L}).$$

- ▶ $\mathbf{L} = \mathbf{Z}\mathbf{W}_{\text{lm}} = [\mathbf{l}_1, \dots, \mathbf{l}_{n+1}]^\top \in \mathbb{R}^{(n+1) \times V}$ are called logits corresponding to each input token \mathbf{x}_i .
- ▶ $\mathbf{p}_i \in \mathbb{R}^V$ is the probability distribution over the vocabulary, for sampling the next token i , as we have shifted the position by adding the starting token <start>.

Interpretation: Transformer uses the new attention mechanism. Its overall idea is similar to other neural networks. Trying to extract useful features for linear classification.

Tied Embedding and LM Head

Tied parameters in LLMs:

- ▶ In decoder-only LMs, the **same** matrix (or its transpose) is often used for the input embedding and for the LM head matrix \mathbf{W}_{lm} in the final logistic regression layer.

$$\mathbf{W}_{lm} = \mathbf{E} \quad (\text{or } \mathbf{W}_{lm} = \mathbf{E}^\top, \text{ depending on convention}).$$

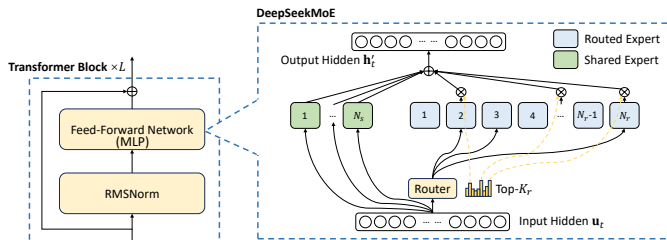
- ▶ This is called Tied LM head.
- ▶ it reduces parameters by $V \times d$. This is often a very large matrix.

This is reasonable as we use the same embedding to transfer tokens to embedding vectors, and then use the same matrix to transfer embedding back to tokens in the vocabulary.

Mixture of Experts (MoE)

Mixture-of-Experts (MoE) Feed-Forward Layer

- ▶ MoE is to expand the **MLP / feed-forward network (FFN)**, not attention.
- ▶ In standard Transformers, each token goes through the **same** MLP layer.
- ▶ In a MoE MLP layer, we have multiple MLPs (**experts**) and a **router** that chooses a small subset of experts for each token.
- ▶ MoE is a simple architecture scaling approach that can expand the representation power of transformer (said by OpenAI people).



Mathematical Definition of MoE

- ▶ Let $\mathbf{u}_t \in \mathbb{R}^d$ be the hidden state / activation value of token \mathbf{x}_t entering a MoE MLP. The router produces a sparse weighting over N_r experts:

$$\mathbf{g}_t = \text{Top-}K_r(\text{softmax}(\mathbf{R}\mathbf{u}_t)) \in \mathbb{R}^{N_r},$$

where $\mathbf{R} \in \mathbb{R}^{N_r \times d}$ is the router matrix and K_r is the number of activated experts per token (e.g., $K_r = 2$).

- ▶ Only the top- K_r entries of \mathbf{g}_t are kept, others are 0, i.e., only the **top- K_r** routed experts are chosen.

Overall, the MoE layer per token is mathematically defined as

$$\mathbf{h}'_t = \underbrace{\mathbf{u}_t}_{\text{skip connection}} + \underbrace{\sum_{i=1}^{N_s} \text{MLP}_i^{(s)}(\mathbf{u}_t)}_{\text{shared MLP}} + \underbrace{\sum_{i=1}^{N_r} \mathbf{g}_t[i] \cdot \text{MLP}_i^{(r)}(\mathbf{u}_t)}_{\text{routed expert}}.$$

↪ Next lecture: Large language models (LLMs).