



Q2

The Transformer applies self-attention mechanism. Here is the corresponding code in the Fairseq-toolkit. Please write the line number(s) for the code(s) which do the scaling with the coefficient.

Solution:

```
In [ ]: """
self.scaling = self.head_dim ** 0.5 # Line 114
q *= self.scaling # Line 602
"""
```

Q3

RMSNorm is another regularizer which can replace the layer normalization. Read the paper and briefly illustrate the calculation process of RMSNorm and its advantages over Layer Normalization

Solution: The RNSNorm is calculated as below

$$\bar{a}_i = \frac{a_i}{\text{RMS}(\vec{a})} g_i, \text{ where } \text{RMS}(\vec{a}) = \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2}, i = 1, 2, \dots, n.$$

where  $g_i$  is the gain parameter used to re-scale the standardized summed inputs, and is set to 1 at the beginning.

One of the advantages of RMSNorm to LayerNorm is that RMSNorm simplifies computation, reducing running time by 6.9% to 40% on different models, while it performs comparable or better performance than LayerNorm on various tasks and neural architectures. Also, RMSNorm is invariant to the scaling of inputs and weights, which stabilizes the layer activations and the model gradients. It also has an implicit learning rate adaptation ability that avoids large-norm weight matrix and improves model convergence.

Appendix

LLM used:

ChatGPT 3.5

请根据以下文字编写相应的代码：

Bind Network.

In Figure 3, we present the details of the bind network, which aims to align the embedding space between ImageBind and LLaMA. Specifically, we denote the  $C_I$ -dimensional global image feature encoded by ImageBind as  $F_I \in \mathbb{R}^{(1 \times C_I)}$ . In the bind network, we first adopt a linear projection layer with a weight matrix  $w_0 \in \mathbb{R}^{(C_I \times C)}$ , formulated as  $F_I^0 = F_I w_0 \in \mathbb{R}^{(1 \times C)}$ , where  $C$  denotes the feature dimension of LLaMA. Inspired by the Feed-Forward Network (FFN) in LLaMA, we then cascade three projection blocks with RMSNorm [43], SiLU activation functions [44], and residual connections [45]. For the  $(i + 1)$ -th block with  $F_I^i$  as input, we formulate the calculation of  $F_I^{(i+1)}$  as (the normalization is omitted for simplicity)

$$F_I^{(i+1)} = F_I^i + (F_I^i w_2 \cdot \text{SiLU}(F_I^i w_1))w_3, 0 \leq i < 3 \quad (1)$$

where  $w_1, w_2 \in \mathbb{R}^{(C \times Ch)}$  and  $w_3 \in \mathbb{R}^{(Ch \times C)}$ , with  $Ch$  denoting the hidden dimension. After the bind network, we obtain the transformed image feature,  $T_I \in \mathbb{R}^{(1 \times C)}$ , which learns to align the embedding space from ImageBind to LLaMA

New Bing

Read the paper and summarize the advantages of RMSNorm when compared with LayerNorm