

Attention Is All You Need

Introduction

In practice, we compute the attention function on a set of queries simultaneously, packed together into a matrix Q . The keys and values are also packed together into matrices K and V . We compute the matrix of outputs as:

$$\text{Attention}(Q, K, V) = \text{softmax}(QK^T \sqrt{d_k})V$$

The two most commonly used attention functions are additive attention [2], and dot-product (multiplicative) attention. Dot-product attention is identical to our algorithm, except for the scaling factor of $1/\sqrt{d_k}$.

With a single attention head, averaging inhibits this.

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W_O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

$$i \in \{1, \dots, h\}$$

$$i \in \{1, \dots, h\}$$

$$i \in \{1, \dots, h\}$$

Where the projections are parameter matrices W_i^Q, W_i^K, W_i^V

$$i \in \{1, \dots, h\}$$

$$i \in \{1, \dots, h\}$$

$$i \in \{1, \dots, h\}$$

and $W_O \in \mathbb{R}^{h d_v \times d_{\text{model}}}$.

In this work we employ $h = 8$ parallel attention layers, or heads. For each of these we use $d_k = d_v = d_{\text{model}}/h = 64$. Due to the reduced dimension of each head, the total computational cost is similar to that of single-head attention with full dimensionality.

ReLU activation in between.

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2 \quad (2)$$

While the linear transformations are the same across different positions

In this work, we use sine and cosine functions of different frequencies:

$$\text{PE}(\text{pos}, 2i) = \sin(\text{pos}/10000^{2i/d_{\text{model}}})$$

$$\text{PE}(\text{pos}, 2i+1) = \cos(\text{pos}/10000^{2i/d_{\text{model}}})$$

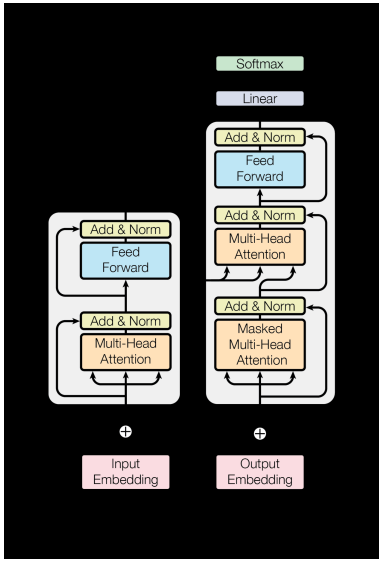
where pos is the position and i is the dimension. That is, each dimension of the positional encoding

We used the Adam optimizer [20] with $\beta_1 = 0.9$, $\beta_2 = 0.98$ and $\epsilon = 10^{-9}$. We varied the learning rate over the course of training, according to the formula:

$$\text{lr}_{\text{rate}} = d^{-0.5}$$

$$\text{model} \cdot \min(\text{step_num}^{-0.5}, \text{step_num} \cdot \text{warmup_steps}^{-1.5}) \quad (3)$$

This corresponds to increasing the learning rate linearly for the first warmup_steps training steps, and decreasing it thereafter proportionally to the inverse square root of the step number. We used $\text{warmup_steps} = 4000$.



$$\bar{V}_S = \bar{A} \bar{V}_R + \bar{B} \bar{I}_R \quad \text{and} \quad \bar{I}_S = \bar{C} \bar{V}_R + \bar{D} \bar{I}_R$$

We have,

$$\bar{A} = \bar{D} = \cosh \gamma l = \cosh(\sqrt{yz} l) = \cosh \sqrt{(ly)(lz)} = \cosh(\sqrt{YZ})$$

$$\bar{B} = Z_C \sinh \gamma l = \sqrt{\frac{Z}{Y}} \sinh(\sqrt{yz} l)$$

$$\bar{B} = \sqrt{\frac{lz}{ly}} \sinh(\sqrt{(ly)(lz)}) = \sqrt{\frac{Z}{Y}} \sinh(\sqrt{YZ})$$

$$\bar{C} = \frac{1}{Z_C} \sinh \gamma l = \sqrt{\frac{Y}{Z}} \sinh(\sqrt{yz} l)$$

∴

$$\bar{C} = \sqrt{\frac{ly}{lz}} \sinh(\sqrt{(ly)(lz)}) = \sqrt{\frac{Y}{Z}} \sinh(\sqrt{YZ})$$

As $\bar{A} = \bar{D}$, we can say that the network is symmetrical.

Consider,

$$\begin{aligned} \bar{A} \bar{D} - \bar{B} \bar{C} &= [(\cosh \sqrt{YZ})(\cosh \sqrt{YZ})] - \left[\sqrt{\frac{Z}{Y}} \sinh(\sqrt{YZ}) \sqrt{\frac{Y}{Z}} \sinh(\sqrt{YZ}) \right] \\ &= \cosh^2(\sqrt{YZ}) - \sinh^2(\sqrt{YZ}) = 1 \end{aligned}$$

As $\bar{A} \bar{D} - \bar{B} \bar{C} = 1$ we can say that the network is reciprocal also.