

encod	ler-d	lecod	ler st	ructure
				accarc

Positional Encoding

组件训练结果评估

Transformer

At each step the model is autoregressive, consuming the previously generated symbols as additional input when generating the next

Given z, the decoder then generates an output sequence (y1, ..., ym) of symbols one element at a time. Composed of a stack of N = 6 identical layers, each layer has THREE sub-layers

Addition inserts a third sub-layer, which performs multi-head attention over the output of the encoder stack

Also modify the self-attention sublayer in the decoder stack to prevent positions from attending to subsequent position(防止位置关 注后续位置).

by masking out (setting to $-\infty$) all values in the input of the softmax which correspond to illegal connections.

While the two are similar in theoretical complexity, dot-product attention is much faster and more space-efficient in practice, since it can be implemented using highly optimized matrix multiplication code. We suspect that for large values of dk, the dot products grow large in To illustrate why the dot products magnitude, pushing the softmax get large, assume that the function into regions where it has components of q and k are **Scaled Dot-Product Attention** extremely small gradients 4. To independent random variables counteract this effect, we scale the with mean 0 and variance 1. Then dot products by 1√dk their dot product, $q \cdot k = \sum qiki$, has mean 0 and variance dk it beneficial to linearly project the Multi-head attention allows the queries, keys and values h times model to jointly attend to information from different with different, learned linear projections to dk, dk and dv representation subspaces at dimensions, respectively. different positions. $\mathrm{MultiHead}(Q,K,V) = \mathrm{Concat}(\mathrm{head}_1,\ldots,\mathrm{head}_\mathrm{h})W^O$ **Multi-Head Attention** $\mathbf{W} = \mathrm{Attention} \left(QW_i^Q, KW_i^K, VW_i^V ight)$

add "positional encodings" to the input embeddings at the bottoms of the encoder and decoder stacks. The positional encodings have the same dimension d-model as the embeddings, so that the two can be summed. $W_i^Q \in \mathbb{R}^{d_{ ext{model}} imes d_k}, W_i^K \in \mathbb{R}^{d_{ ext{model}} imes d_k}, W_i^V \in \mathbb{R}^{d_{ ext{model}} imes d_v}, W^O \in \mathbb{R}^{hd_v imes d_{ ext{model}}}$

 $PE_{({
m pos}\ ,2i)} = \sin \Bigl({
m pos}\ /10000^{2i/d_{
m model}} \Bigr) \ PE_{({
m pos}\ ,2i+1)} = \cos \Bigl(pos/10000^{2i/d_{
m model}} \Bigr)$

where pos is the position and i is the dimension

We chose this function because we hypothesized it would allow the model to easily learn to attend by relative positions, since for any fixed offset k, PEpos+k can be represented as a linear function of PEpos.

it may allow the model to extrapolate to sequence lengths longer than the ones encountered during training.

While single-head attention is 0.9 BLEU worse than the best setting, quality also drops off with too many heads. replace our sinusoidal positional encoding with learned positional embeddings, observe nearly identical results to the base model.

reducing the attention key size dk hurts model quality. This suggests that determining compatibility is not easy and that a more sophisticated compatibility function than dot product may be beneficial.

bigger models are better, and dropout is very helpful in avoiding over-fittin