

(۱) (الف)

$windowsize=2 \rightarrow$  (بی بزرگنمای اول و آخر) ۴ جمله می باشد

$$= 4 \times 6 + 3 \times 2 + 2 \times 2 = 34$$

۲ جمله بزرگ نمای اول و آخر  
۳ جمله بزرگ نمای اول و آخر

$$34 \times 5 = 170$$

هر جمله می باشد سه صفحه توسعه ممکن است (اریح) :

$$34 + 170 = 204$$

4\*

$$d_{head} = d_{model} / h \rightarrow w_Q, w_K, w_V \in \mathbb{R}^{d_{model} \times d_{head}}$$

\*

لักษن ابعاد از  $d_{model}$  و  $d_{model} \times d_{head}$  باشند که سریار محاسبات هستند حذف ممکن است  
ها می سود. از طرفی سریار سازی باست لذا می تصریح می شود سریار (از این) عادی

\*) صحیح  
\*) خطا، با این تعداد head ها ممکن است  $w_{head}$  متناسب تولید نکند.

(۱) (ج)

$$-\frac{\delta \log P(w_t | w_0)}{\delta V_{w_0}} = -\frac{\delta}{\delta V_{w_0}} \left( u_{w_t}^T v_{w_0} - \log \sum_{k=1}^K p(u_{w_t}^T v_k) \right)$$

$$= -\left( u_{w_t} - \frac{\sum_k u_{w_t} p(u_{w_t}^T v_k)}{\sum_k p(u_{w_t}^T v_k)} \right) = \sum_k u_{w_t} p(w_k | w_t) - u_{w_t}$$

\*)  $c=1$  ، مدل پیش بینی برداری روابط و اتفاقات محلی با بلند نگاه مسکن عالی است. این باست منسوب است  
مدل روابط دستوری بین کلمات را بایار لبیر.

\*)  $c=100$  مبد پیش بینی از کلمات در نظر گرفته شد که باست دستور را مدل روابط معنایی بین  
کلمات را بستر باز نماید

\*)  $c=5$  باست مدلی دو معنایه ای بین دو کلمات قبل است

Output Softmax  $\left( \frac{QK^T}{\sqrt{d}} \right) V$

(الف)

$$\frac{QK^T}{\sqrt{d}} = \frac{1}{\sqrt{d}} X W_Q W_K^T X^T$$

لہا اسی صورت میں FCN کے لیے  $\frac{QK^T}{\sqrt{d}} = WXX^T$  میں داری ہے

$$WXX^T = \frac{1}{\sqrt{d}} X W_Q W_K^T X^T \xrightarrow{M = \frac{1}{\sqrt{d}} W_Q W_K^T} WXX^T = XMX^T$$

$$W = f(M) \rightarrow W_{ij} = \sum_{k,l} M_{kl} \cdot g(X_i, X_j) \rightarrow \text{output} = \text{softmax}(WXX^T)v$$

$$W_Q W_K^T \rightarrow D \times D \rightarrow D^2 \rightarrow D^2 N^2$$

$$WXX^T \rightarrow N \times N \rightarrow N^2$$

$$A = \text{softmax}(QK^T) \quad \text{سیارہ از منحصرہ ماتریس } A \text{ پر دلیل احتمال درجہ نگوچے مقتدار میں دستیاب صورت درد کے باعث میں سودھ ماتریس}$$

$$A_{ij} = \text{softmax}(Q_i K_j^T) \quad \text{سینودھ صحیح ارتکل ٹافرن سپاٹس}$$

کو نیکان بارے توکن کا دیگر دلیل اسی پاسندر درایتھر مربوط ہے اسی نیز کے بندوق اصریح

$$H_i = \text{Attention}(Q_i, K_i, V_i)$$

$$Q_i = X W_i^Q, K_i = X W_i^K, V = X W_i^V$$

$$Q_i = K_i = V_i = X \quad \leftarrow W_i^Q = W_i^K = W_i^V, \text{ creating } \hookrightarrow$$

$$\text{softmax}\left(\frac{X X^T}{\sqrt{d}}\right) \rightarrow \text{redistribution of input}$$

$$\text{softmax}\left(\frac{X X^T}{\sqrt{d}}\right) X \rightarrow \text{rearranged input}$$

Input = (32, 2048, 768)

: Encoder - 1 (S)

$V, K, Q = (32, 2048, 768)$

Multihead Attention output: (32, 2048, 768)

Feed Forward Layer 1: (32, 2048, 384)

Feed Forward Layer 2: (32, 2048, 768)

Input: (32, 2048, 768)

Self Attention output: (32, 2048, 768)

Cross Attention output: (32, 2048, 768)

Feed Forward layer 1: (32, 2048, 384)

Feed forward layer 2: (32, 2048, 768)

: decoder

Token Embedding:  $30000 \times 768 = 23040000$  : Encoder - Y

Self Attention:  $4 \times 768 \times 768 = 2359296$

$\xrightarrow{\times 8} 18874368$

Feed Forward:  $\left\{ \begin{array}{l} \text{Layer 1: } 768 \times 384 = 294912 \\ \text{Layer 2: } 768 \times 384 = 294912 \end{array} \right. + = 587824$

$\xrightarrow{\times 8} 4718592$

Layer Norm.:  $2 \times 768 \times 8 = 12288$

Total = 46645248

$$12 \times 2359296 = 28311552 : \text{self Attention} \quad \text{Decoder}$$

$$12 \times 2359296 = 28311552 : \text{Cross Attention}$$

$$\text{Feed Forward: } 12 \times 589824 = 7071880$$

$$\text{Layer Norm: } 2 \times 768 \times 12 = 78482$$

$$\text{Total} = 63719424$$

"He is sleeping" —  $\mathcal{N}_{\text{Input}} = [\text{Token}_1, \text{Token}_2, \text{Token}_3]$  ↗

embedding  $(32, 2048, 768) \xrightarrow{\text{encoder}}$  Context representation

Decoder input:  $y_{\text{input}}(\text{in context}) = [\text{SOS}] \xrightarrow{\text{self}} \text{generate tokens}$   
Cross, Best match for generating

$\xrightarrow{} [\text{SOS}] \xrightarrow{\text{He} \rightarrow [\text{SOS}]} \xrightarrow{\text{is}} \xrightarrow{\text{[SOS] هو}} \xrightarrow{\text{sleeping}} \xrightarrow{\text{[SOS] هو نائم}} \xrightarrow{\text{Best match}}$

$\xrightarrow{\Sigma \text{Eos}} \xrightarrow{\text{[SOS] هو نائم [Eos]}}$

(۲) که دلیل برین کا ہسپ ہزینہ حساب ادا کے۔ دلیل نیگری Generalization بہر است ذرا  
روابط بین جملے معرفہ مختص یہ کہ task یا معرفہ خاص میں ہے۔ دلیل بعدی شروعاتی تبریز است  
کہ سُنان (LH مدل) ۷۰٪ میں NSP کو RoBERTa استنادہ منکر کرتے اور BERT کو تراوید  
ب) Pooling برداری زیر ایجاد Embedding را بکہ بردار باطنی بتے ہیں ہے، این بردار تھیں ممکن  
کہ لاتھی کلمات داخل جملہ را نشان دهد۔ حد اتفاق این کہ بردار مصنوعی بندہ را بکہ صورت غیر وہ ذخیرہ کرنے  
کہ در ۷۰٪ میں NSP task میں تواند ہے۔