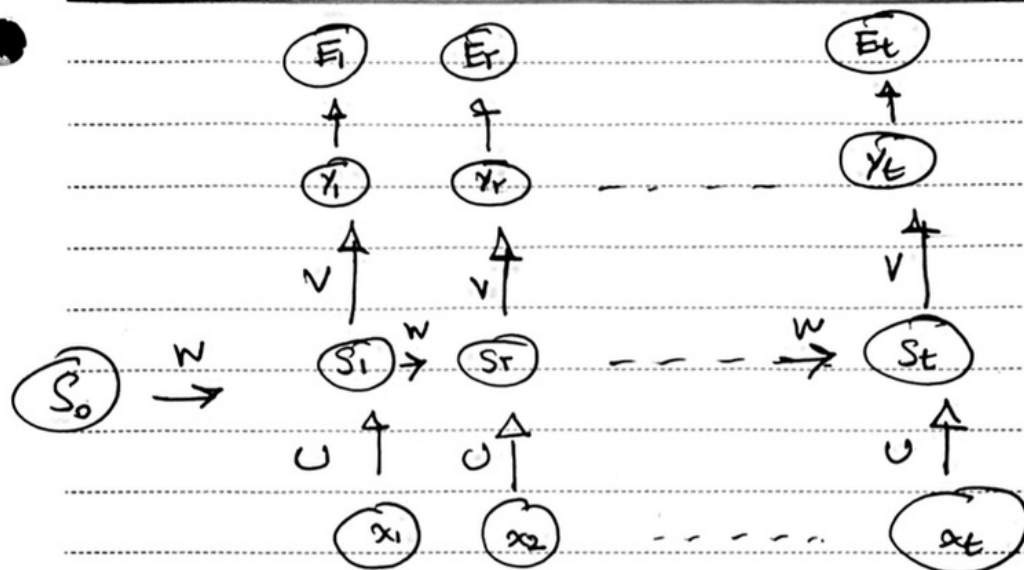


$$E_{total} = - \sum_t \partial_t \hat{y}_t$$

Subject:

Date:



$$S_t = \tanh(W_{St} + U_{xt} + b_s)$$

$$\hat{y}_t = \text{Soft}(V S_t + b_y)$$

$$E_t = - \partial_t \hat{y}_t$$

First  $V$ :

$$\frac{\partial E_t}{\partial V_{ij}} = \frac{\partial E_t}{\partial \hat{y}_t} \times \frac{\partial \hat{y}_t}{\partial q_{te}} \times \frac{\partial q_{te}}{\partial V_{ij}}$$

$$= \frac{\partial E_t}{\partial \hat{y}_t} \times \begin{cases} -\hat{y}_k \hat{y}_l & k \neq l \\ \hat{y}_l (1 - \hat{y}_l) & k = l \end{cases}$$

$$\Rightarrow \frac{\partial E_t}{\partial q_{te}} = \frac{\partial E_t}{\partial \hat{y}_t} \times \hat{y}_l (1 - \hat{y}_l) + \sum_{k \neq l} \left( \frac{\partial E_t}{\partial \hat{y}_k} \right) (-\hat{y}_k \hat{y}_l)$$

$$= \hat{y}_l - \hat{y}_l^2 \rightarrow \sum_k \hat{y}_k = 1 \text{ (one-hot)!!}$$

$$\frac{\partial q_{te}}{\partial V_{ij}} = \frac{\partial}{\partial V_{ij}} (V_{ij} S_{t,j}) = S_{t,j}$$

Subject:

بجواب

Date:

$$\Rightarrow \frac{\partial E_t}{\partial V} = (\hat{y}_t - y_t) \otimes S_t \rightarrow \text{Kxd}$$

! ع.و

$$\Rightarrow \frac{\partial E}{\partial V} = \sum_t (\hat{y}_t - y_t) \otimes S_t$$

for  $w_{ij}$ :

$$\frac{\partial E_t}{\partial w_{ij}} = \underbrace{\frac{\partial E_t}{\partial \hat{y}_t} \times \frac{\partial \hat{y}_t}{\partial q_{te}}}_{\text{سلسلة}} \times \underbrace{\frac{\partial q_{te}}{\partial s_{tm}} \times \frac{\partial s_{tm}}{\partial w_{ij}}}_{\text{سلسلة}}$$

سلسلة

$$\rightarrow = (\hat{y}_t - y_{te})$$

$$\frac{\partial q_{te}}{\partial s_{tm}} = \frac{\partial}{\partial s_{tm}} (V_{lm} \cdot s_{tm}) = V_{lm}$$

$$\frac{\partial s_{tm}}{\partial w_{ij}} = \frac{\partial s_{tm}}{\partial w_{ij}} + \frac{\partial s_{tm}}{\partial s_{t-1}} \frac{\partial s_{t-1}}{\partial w_{ij}} + \frac{\partial s_{tm}}{\partial s_{t-1}} \times 0 \times 0 + \dots$$

السلسلة الأولى هي  $s_{t-1}$  و  $s_{t-1}$  هي  $s_{t-1}$  و  $s_{t-1}$  هي  $s_{t-1}$

$$\Rightarrow \frac{\partial s_{tm}}{\partial w_{ij}} = \sum_{r=0}^t \frac{\partial s_{tm}}{\partial s_{rn}} \times \frac{\partial s_{rn}}{\partial w_{ij}}$$

$$\Rightarrow \frac{\partial E_t}{\partial w_{ij}} = (\hat{y}_t - y_{te}) \times V_{lm} \times \sum_{r=0}^t \frac{\partial s_{tm}}{\partial s_{rn}} \frac{\partial s_{rn}}{\partial w_{ij}}$$

$$\Rightarrow \frac{\partial E}{\partial w_{ij}} = \sum_t \frac{\partial E_t}{\partial w_{ij}}$$

التي هي  $s_{t-1}$  و  $s_{t-1}$  هي  $s_{t-1}$  و  $s_{t-1}$  هي  $s_{t-1}$

هذا هو الشكل النهائي للحل

النتيجة

$$\frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial S_3} \times \frac{\partial S_3}{\partial W} = \frac{\partial E_3}{\partial S_3} \times \frac{\partial S_3}{\partial z_3} \times \frac{\partial z_3}{\partial W}$$

نظير، لم يكون!

$$(\hat{y}_l - y_l) \times V_{lm} \times (1 - S_3^T) \times S_{r,j}$$

$$\Rightarrow \frac{\partial E_3}{\partial W} = V^T \odot (\hat{y}_3 - y_3) \odot (1 - S_3^T) \odot S_r^T$$

~~1x1~~ ((d x d))

نظير؟

$$\frac{\partial E_3}{\partial W} = W^T (V^T \odot (\hat{y}_3 - y_3) \odot (1 - S_3^T)) \odot (1 - S_r^T) \odot S_r^T$$

$$\frac{\partial E_3}{\partial W} = W^T (W^T (V^T \odot (\hat{y}_3 - y_3) \odot (1 - S_3^T)) \odot (1 - S_r^T) \odot (1 - S_1^T) \odot S_0^T)$$



باسمه تعالی

دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

یادگیری عمیق

مقالات انتهایی سوال دوم

## ۱ سوال دوم

### ۱.۱ الف

تابع  $attention$  در واقع یک نگاشت از یک  $query$  و تعدادی پیر  $key - value$  از ورودی به خروجی میبشد، که همگی بصورت بردار هستند.

### ۲.۱ ب

اگر بعد مقداری  $key$  ها زیاد شود، ضرب نقطه ای هزینه محاسباتی و عملیاتی زیادی خواهد داشت و این باعث میشود تابع سافت مکس به گونه ای اشباع شود، که برای حل این مشکل خروجی را تقسیم بر  $\sqrt{d_k}$  میکنیم.

### ۳.۱ ج

۱. در بخش دیکود کردن:  $query$  ها را از لایه دیکودر قبل و پیرهای کی و لیو از ایه انکودر بدست آورده میشوند که باعث دسترسی  $global$  به همه ورودی ها داشته باشد.

۲. در هر انکودر همه ورودی ها از لایه قبلی انکودر آمده اند. و در هر دیکودر نیز همینطور

### ۴.۱ د

برای اینکه ترتیب ورودی ها مشابه نگه داشته شود، بجای استفاده از لایه های کنولوشنی یا  $recurrent$  از مفهومی به نام  $positional encoding$  استفاده میشود، در این کار با توجه به موقعیت هر کلمه در جمله، به هر کلمه یک بردار با بعد  $word embedding$  نسبت داده میشود که فرمول زیر را ارضا میکند: یعنی هر موقعیت کلمه در جمله تناظر یک به یک با یک تابع سینوسی یا فرائنس متفاوت دارد، شبکه ارتباط

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

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

بین کلمه ها را با این بردار ها یاد خواهد گرفت. راهکار دیگر نسبت دادن یک بردار به هر موقعیت است که تفاوت آن با قبلی قابل یادگیری بودن است که نتایج هر دو تقریباً مشابه است همانطور که در شرط  $E$  جدول زیر آمده است:

	$N$	$d_{\text{model}}$	$d_{\text{ff}}$	$h$	$d_k$	$d_v$	$P_{\text{drop}}$	$\epsilon_{\text{ls}}$	train steps	PPL (dev)	BLEU (dev)	params $\times 10^6$
base	6	512	2048	8	64	64	0.1	0.1	100K	4.92	25.8	65
(A)				1	512	512				5.29	24.9	
				4	128	128				5.00	25.5	
				16	32	32				4.91	25.8	
				32	16	16				5.01	25.4	
(B)				16						5.16	25.1	58
				32						5.01	25.4	60
(C)	2									6.11	23.7	36
	4									5.19	25.3	50
	8									4.88	25.5	80
		256			32	32				5.75	24.5	28
		1024			128	128				4.66	26.0	168
			1024							5.12	25.4	53
			4096							4.75	26.2	90
(D)							0.0			5.77	24.6	
							0.2			4.95	25.5	
								0.0		4.67	25.3	
								0.2		5.47	25.7	
(E)		positional embedding instead of sinusoids								4.92	25.7	
big	6	1024	4096	16			0.3		300K	<b>4.33</b>	<b>26.4</b>	213

## ۵.۱

سه معیار برای این مقایسه بررسی میشود، اول مقدار پیچیدگی محاسباتی در هر لایه، دوم مقدار محاسبه با قابلیت پردازش موازی و سوم ارتباط بین ورودی‌های با فاصله زیاد از هم. در رابطه با مورد سوم یکی از مشکلات بزرگ شبکه‌های پردازش زبان این است که با دور شدن استپ‌های زمانی از هم ارتباط آنها از هم کاهش میابد که منجر به کاهش قدرت ماشین در پردازش جمله‌های طولانی میشود که در جدول زیر این سه معیار مقایسه شده اند:

Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types.  $n$  is the sequence length,  $d$  is the representation dimension,  $k$  is the kernel size of convolutions and  $r$  the size of the neighborhood in restricted self-attention.

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

## ۶.۱ و

بخش انکودر شامل ۶ لایه مشابه است که هر یک دو زیر لایه دارا هستند، اولی یک  $self - attention$  چندسر و دومی یک شبکه فولی کانکتد. در هر زیر لایه از یک  $skipconnection$  استفاده شده است که با خروجی هر زیر لایه ورودی را جمع می‌کند، بعد از هر زیر لایه نیز یک  $layernormalization$  قرار داده شده. بعد تمام خروجی ها ثابت و برابر با 512 است در دیکودر نیز ۶ لایه مشابه بالا داریم که زیر لایه های مشابه بالا نیز دارد و علاوه بر این دو یک لایه  $selfattention$  چندسر دیگر نیز قرار داده شده است که به خروجی های انکودر متصل است و نهایتاً  $skip - connection$  ها و  $layer - normalization$  های مشابه نیز داریم.

بهینه ساز مورد استفاده مقاله ادا م است با پارامتر های  $\epsilon = 10^{-9}, \beta_1=0.9, \beta_2=0.98$  و نرخ یادگیری متغیر با زمان با فرمول به شکل زیر داراست:

$$lr_{rate} = d_{model}^{-0.5} \cdot \min(step\_num^{-0.5}, step\_num \cdot warmup\_steps^{-1.5})$$

از دو  $regularizer$  متفاوت استفاده شده، اول دراپ اوت با نرخ 0.1 که یعنی یک دهم وزن ها را در هر ایتريشن دور میندازد، که هم در خروجی زیر لایه ها قرار دارد و هم در هر دو قسمت دیکودر و انکودر پس از حاصل جمع امبدینگ کلمه ها امبدینگ موقعیت ها. همچنین لیبیل اسموسینگ با پارامتر  $\epsilon = 0.1$  نیز استفاده شده است..

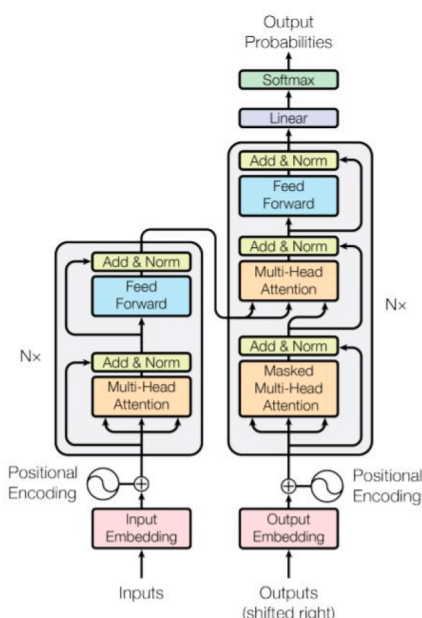
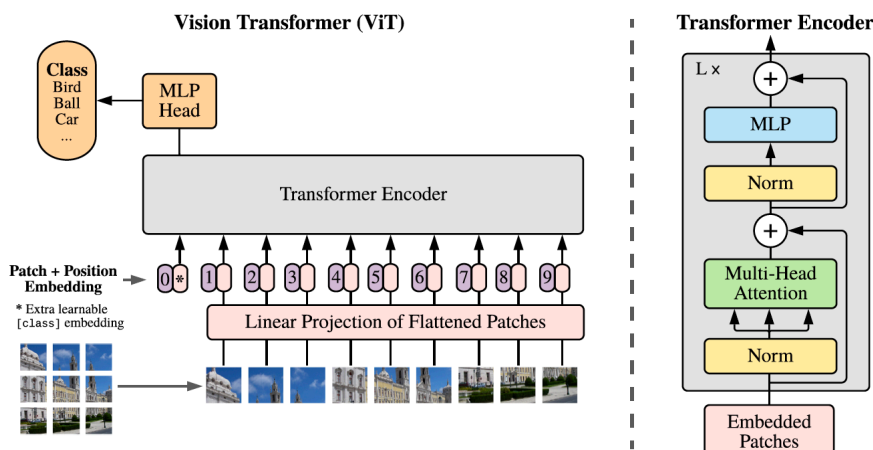


Figure 1: The Transformer - model architecture.

## ۷.۱ ز

در این مقاله با توجه به کاربرد روزافزون مدل های پردازش زبان های طبیعی و اینکه هر حرف را بعنوان یک توکن در نظر میگیرند ، تلاش شده تا با تقسیم یک عکس بهج های مختلف هر کدک از این تیکه ها را بعنوان یک توکن در نظر گرفته ، سپس یک امبدینگ خطی روی آنها زده، ایندکس هر یک را کنرش گذاشته و وارد یک ساختار ترنسفورمر میکنید که به شکل زیر است ، و در نهایت برای امر کسلیفیکیشن آن را ترین میکند که نیازمند لیبل منتظار بالا سمت چپ تصویر میباشد :



حال اگر فرض شود یک عکس با ابعاد  $(H \times W \times C)$  داشته باشیم ، ان را به  $N$  بهج با ابعاد  $(P \times P)$  تقسیم میکنیم که نهایتا به ابعاد زیر در میاید :  $(N \times (P^2 C))$  نهایتا روی این بهج های یک امبدینگ صورت میگیرد تا هر کدام به یک بعد  $D$  تصویر بشوند ، که رابطه آن به شرح زیر است ؛

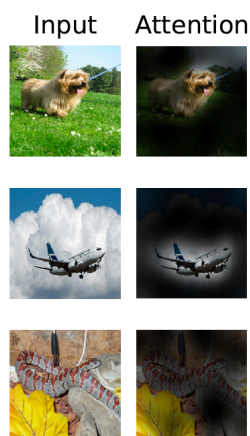
$$\mathbf{z}_0 = [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \dots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{\text{pos}}, \quad \mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \quad \mathbf{E}_{\text{pos}} \in \mathbb{R}^{(N+1) \times D} \quad (1)$$

که  $X_{\text{class}}$  بعنوان  $\text{classificationHead}$  شناخته میشود .

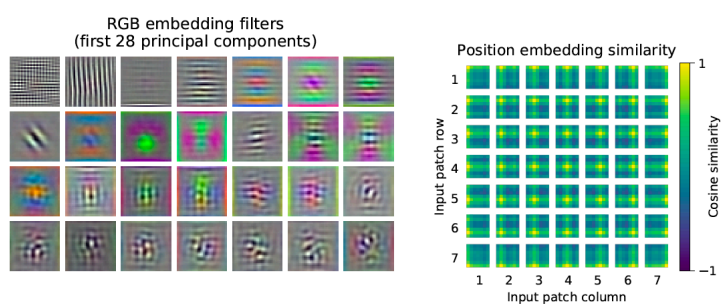
همچنین برای جلوگیری از مشکل  $\text{image} - \text{specific inductive bias}$  که بر خلاف مدل های کانولوشنی برای مدل مارخ میدهد سعی میشود تا ورودی لایه امبدینگ خطی بجای بهج های تصویر اصلی ، بهج های استخراج شده پس از گذرا از یک لایه کانولوشنی باشند .

در ابتدای لایه ترنسفورمر ، بهج ها به یک بعد پایین تر کاهش داده میشوند ، و سپس یک امبدینگ پوزیشن با آنها اضافه میشود ، که قابل ذکر است بهج های به هم نزدیک تر پوزیشن امبدینگ های مشابهی خواهند داشت ، و همچنین  $\text{self} - \text{attention}$  که در لایه ترنسفورمر قرار داده شده است طبق عکس زیر اطلاعات سرتاسر عکس را جمع میکند .





همچنین مفاهیم مثل امبدینگ پوزیشن و امبدینگ فیلترها برای خروجیشان مثال های زیر آمده است :



نهایتاً این ساختار با پری ترین شدن روی دیتاست های بزرگ به نتایج قابل توجه و بعضاً بهتری از مدل های روز دنیا میرسد.