





دانشگاه شهید باهنر کرمان

دانشکده فنی و مهندسی

گروه مهندسی کامپیوتر

پروژه کارشناسی

آشنایی با ترنسفورمرها

استاد راهنما: دکتر مهدی افتخاری

توسط: امیرحسین امیرماهانی و غلامرضا دار

تاریخ: ۸/۶/۱۴۰۰

چکیده

چندین سال است که LSTM ها و GRU ها روش های اصلی برای انجام تقریباً تمام کارهای مربوط به پردازش زبان طبیعی و درک زبان^۱ هستند. مانند هر روشی، این روش ها هم مشکلات و ضعف های خود را دارند.

در این پروژه در مورد پیدایش ترنسفورمر^۲ ها، یک روش بر مبنای توجه^۳ به عنوان جایگزین بسیار عالی برای LSTM ها و GRU ها بحث کردیم. ترنسفورمر ها نشان داده اند که یک معماری بسیار گسترش پذیر و قدرتمند هستند. در زمان نوشتن این پروژه هر چند ماه یکبار بهبودی اساسی در ترنسفورمرها رخ می داد. بررسی کردن این تغییر ها نیز بخشی از بحث ما بود.

در این پروژه در مورد نسخه های مختلف ترنسفورمر ها از جمله ... Linformer, Performer, ... صحبت کردیم و اینکه چگونه کارایی ترنسفورمرها را چندین برابر می کنند. در مورد ایده های این نسخه های مختلف صحبت کردیم و همچنین نقاط ضعف و قوت هر کدام را نیز بیان کردیم.

در نهایت به طور عملی با کار کردن با ترنسفورمرها آشنا شدیم و طی چندین مثال قابلیت های آنها برای حل مشکلات روزمره با قدرت بسیار زیاد را نشان دادیم. همچنین راجع به آینده شاخه ی پردازش زبان به طور کلی بحث شد و اینکه چگونه ترنسفورمرها می توانند همه چیز را تغییر دهند!

^۱ Language Understanding

^۲ Transformer

^۳ Attention

۶	۱- فصل اول: مقدمه و پیشگفتار.....
۷	۱-۱- پیشگفتار.....
۷	۲-۱- شبکه های RNN.....
۸	۳-۱- شبکه های LSTM.....
۹	۴-۱- شبکه های GRU.....
۹	۵-۱- نگرشی نو.....
۱۰	۲- فصل دوم: نگاهی عمیق تر به معماری ترنسفورمرها.....
۱۱	۲-۱- معماری ترنسفورمرها.....
۱۳	۲-۲- سلف اتنشن.....
۱۴	۳-۲- نحوه عملکرد سلف اتنشن.....
۱۵	۲-۳-۱- بدست آوردن بردار های Q,K,V.....
۱۷	۲-۴- اتنشن چند سره.....
۱۹	۳- فصل سوم: چالش های معماری ترنسفورمر.....
۲۱	۳-۱- ترنسفورمرهای خطی.....
۲۲	۳-۱-۱- نتایج بنچمارک ها.....
۲۲	۳-۱-۲- ترنسفورمر خطی در مرحله استنتاج.....
۲۳	۳-۱-۳- کاربردهای ترنسفورمر خطی.....
۲۳	۳-۲- پرفورمر.....
۲۶	۳-۲-۱- توجه بر اساس پیشندها.....
۲۷	۳-۲-۲- مقایسه ها.....
۲۷	۳-۲-۳- حوزه های کاربرد پرفورمر.....
۲۹	۴- فصل چهارم: آشنایی عملی با ترنسفورمرها (و محصولات Hugging face).....
۳۰	۴-۱- کمپانی هاگینگ فیس.....

۳۰.....	۲-۴- کتابخانه Transformers
۳۳.....	۳-۴- مثال عملی FineTune کردن ترنسفورمر
۳۴.....	۵- فصل پنجم: جمع بندی و نتیجه گیری
۳۵.....	۵-۱- جمع بندی نهایی
۳۶.....	مراجع

شکل ۱-۱: ساختار شبکه عصبی RNN	۷
شکل ۲-۱: ساختار یک واحد LSTM	۸
شکل ۱-۲: ساختار معماری ترنسفورمر	۱۱
شکل ۲-۲: شمای کلی ترنسفورمر	۱۲
شکل ۳-۲: نگاهی جزئی تر به ساختمان ترنسفورمر	۱۲
شکل ۴-۲: نحوه ی کار مکانیزم سلف اتنشن	۱۴
شکل ۵-۲: بدست آوردن بردار های Q,K,V	۱۵
شکل ۶-۲: مراحل بعدی محاسبه سلف اتنشن	۱۶
شکل ۷-۲: استفاده تابع سافت مکس در محاسبه ی سلف اتنشن	۱۶
شکل ۸-۲: نحوه محاسبه ی ماتریس ها در اتنشن چند سره	۱۷
شکل ۹-۲: نحوه ی محاسبه ماتریس Z	۱۸
شکل ۱-۳: روش های توجه اندک از اساس کار شبکه های عصبی گرافی نیز ایده گرفتند که با روابط خاص بصورت گراف توجه ظاهر می شوند.	۲۰
شکل ۳-۲: نمودار های بالا نشان میدهند که می توان نتیجه محاسبات ترنسفورمر ها را با ماتریس های درجه کمتر تخمین زد. باید توجه داشت که اگر داده ها در ابعاد مختلف ماتریس بطور یکنواخت توزیع شده باشند یک نمودار خطی قابل انتظار بود. توی این شکل مشخص است که اکثر تاثیر داده ها در ۱۲۸ بعد نخست وجود دارد و در ۳۸۴ بعد، بعدی تاثیر کمی وجود دارد. این نمودار نتیجه آزمایش بر روی ترنسفورمر های ۱۲ و ۲۴ لایه و روی مجموعه دادگان IMDB و WIKI103 بوده است.	۲۱
شکل ۳-۳: بررسی اساس کار ترانسفورمر	۲۳
شکل ۴-۳: خطای مشاهده شده در کرنل مثلثاتی	۲۵
شکل ۵-۳: مرتبه پیچیدگی پرفورمر	۲۶
شکل ۶-۳: توجه براساس پیشوندها	۲۶
شکل ۷-۳: نمودار های دقت و زمان اجرای الگوریتم های مختلف	۲۷
شکل ۸-۳: نتیجه استفاده پرفورمر در مبحث ژنتیک	۲۷
شکل ۱-۴: بخشی از مدل های موجود در هاگینگک فیس	۳۱
شکل ۲-۴: استفاده از کتابخانه ترنسفورمر	۳۱

شکل ۴-۳: استفاده از کتابخانه ترنسفورمر در مبحث آنالیز احساس ۳۲

فهرست جدول ها

- جدول ۱-۳: بنچمارک‌های ترنسفورمر خطی ۲۲
- جدول ۲-۳: این جدول نشان دهنده یک عمل استنباط با استفاده از یک GPU تسلا وی ۱۰۰ با ۱۶ گیگ وی رم است. جدول سمت چپ نشان دهنده زمان صرفه جویی شده نسبت به استفاده از ترنسفورمر معمولی و جدول سمت راست نشون دهنده حافظه صرفه جویی شده نسبت به استفاده از ترنسفورمر معمولی است. ۲۳

۱- فصل اول: مقدمه و پیشگفتار

۱-۱- پیشگفتار

از زمان اختراع اولین کامپیوتر همیشه یک چالش موجود بود که کامپیوترها بتوانند آسان تر با انسان ها ارتباط برقرار کنند. در این زمینه تلاش های زیادی شد و سعی شد که با تعریف نمادهای بصری ارتباط بین کامپیوتر و انسان را تسهیل کنند؛ اما هنوز ارتباط زبانی انسان با ماشین به عنوان یک مسئله و چالش باقی مانده بود. ازین رو برنامه نویسان و طراحان الگوریتم تلاش های زیادی کردند. بعد از ۱۹۵۸ که شبکه های عصبی پا به عرصه وجود گذاشتند. تلاش های بسیاری برای ارتباط زبانی با کامپیوتر از طریق شبکه های عصبی و یادگیری ماشین، صورت گرفت. هر کدام ازین روش ها مزایا و معایبی داشتند که در ادامه به بررسی آن ها می پردازیم.

منابع این گزارش به صورت مجموعه ای از فیلم ها جمع آوری شده است تا علاقمندان این حوزه براحتی و با صرف زمان کمتر بتوانند با مسائل روز دنیا آشنا شوند و در حل آن ها تلاش ورزند.^۴

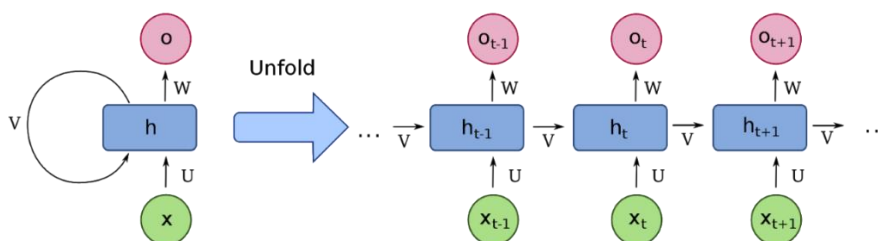
۱-۲- شبکه های RNN

در سال ۱۹۸۶ میلادی یک شبکه عصبی به نام RNN^۵ و برای کاربردهای زیر معرفی شد [1]:

۱. ترجمه زبانی

۲. خلاصه سازی متن

۳. پیش بینی بخش بعدی متن



شکل ۱-۱: ساختار شبکه عصبی RNN

در این تصویر نحوه کارکرد این شبکه، به تصویر کشیده شده است.

^۴ [Transformer and language models - YouTube](#)

^۵ Recurrent neural network

این شبکه، از چندین سلول تشکیل میشود؛ که هر سلول دانش اکتسابی اش را به سلول بعدی انتقال میدهد. بدین صورت میتوان تاثیر توالی کلمات و اثر وابستگی کلمات به کلمات ماقبل خودشان را در نتیجه نهایی مشاهده کرد. اما این روش ایرادتی نیز داشت [1].

۱. نیاز هر سلول به اطلاعات سلول قبل

۲. زمان اجرای زیاد

۳. عدم امکان استفاده از قدرت موازی سازی

۴. مشکل محوشدگی گرادیان

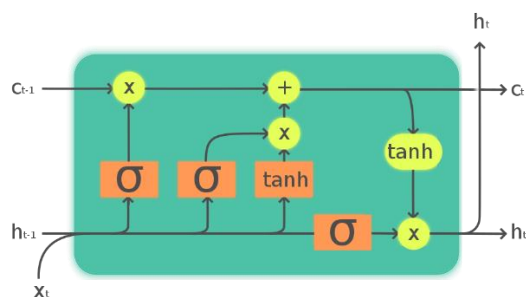
۱-۳- شبکه های LSTM

بعد از یازده سال و در سال ۱۹۹۷ میلادی شبکه های عصبی Lstm معرفی شدند [2].

این شبکه ها در سال ۲۰۰۷ توانستند انقلابی در زمینه تشخیص گفتار آغاز کنند.

LSTM ها رکورد های بهبود ترجمه ماشینی مدل سازی زبان و پردازش چند زبانه را شکستند همچنین ترکیبشان با شبکه های CNN^۷ زیرنویس خودکار تصاویر را بهبود بخشید [2]. این شبکه ها از RNN ها سریعتر بودند اما همچنان با مشکل زمان اجرای زیاد دست و پنجه نرم میکردند، ازین رو برای تسریع آن ها از شتاب دهنده های سخت افزاری استفاده میشد.

LSTM ها تا حدی توانستند مشکل محوشدگی گرادیان (به صفر رسیدن تاثیر خطا در لایه های اولی) را نیز حل کنند؛ زیرا واحد های LSTM اجازه می دهند که شیب ها جریان یابند. با این حال این واحد ها هنوز با مشکل انفجار گرادیان و مشکل زمان زیاد آموزش روبرو بودند [2].



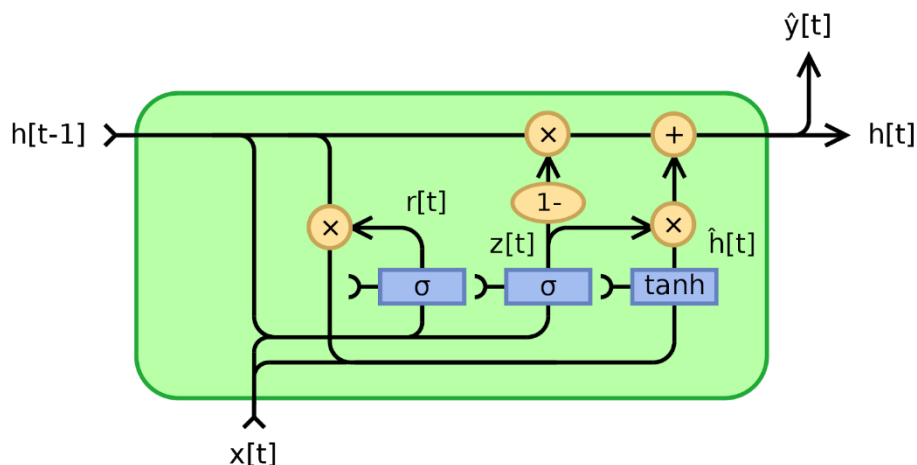
شکل ۱-۲: ساختار یک واحد LSTM

^۶ Long Short-Term Memory

^۷ Convolutional Neural Network

۱-۴- شبکه های GRU

در سال ۲۰۰۴ شبکه های GRU^۱ معرفی شدند [3]. ای شبکه ها بسیار به LSTM شبیه اند با این تفاوت که علاوه بر کاهش تعداد پارامترها، برای جلوگیری از انفجار گرادین، یک گیت فراموشی به سلول های LSTM افزودند. کارایی این شبکه ها بسیار شبیه شبکه های LSTM اند ولی تحقیقات نشان داده اند که GRU در برخی مجموعه های کوچکتر و کم تکرار عملکرد بهتری دارد [3].



شکل ۱-۳: ساختار یک واحد GRU

۱-۵- نگرشی نو

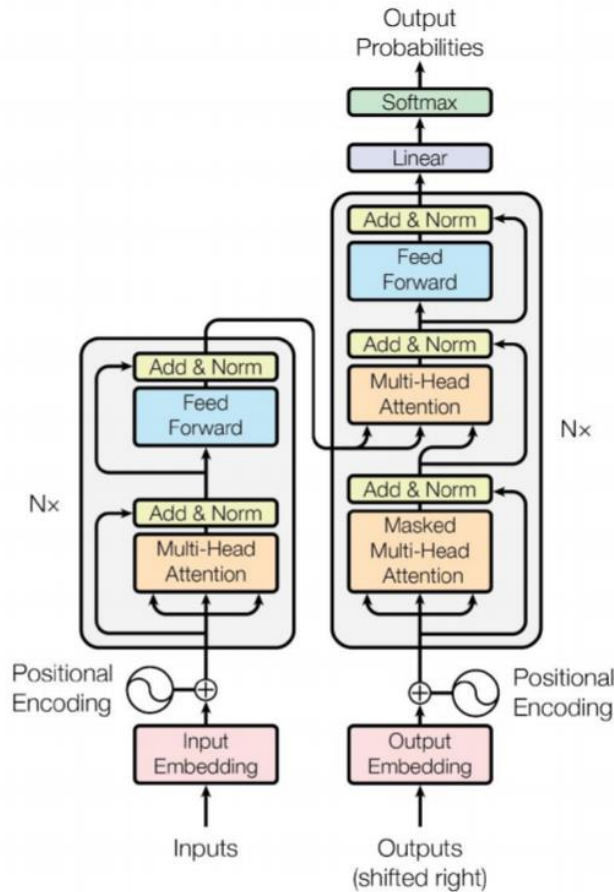
در سال ۲۰۱۷ یک معماری جدید برای شبکه های عصبی معرفی شد که مسئله ارتباط کلمات یک جمله را با دید دیگری بررسی میکرد، دیدی که باعث یک انقلاب جدید در مدل های زبانی شد و توانست انسان را به آرزوی ارتباط زبانی با کامپیوتر و حتی تسهیل ارتباط با انسان های دیگر نزدیک تر کند [4]. در ادامه این گزارش به بررسی ساختار این معماری، کاربردهای آن و مشکلات آن خواهیم پرداخت.

^۱ Gated recurrent unit

۲- فصل دوم: نگاهی عمیق تر به معماری ترنسفورمرها

۲-۱- معماری ترنسفورمرها

در فصل قبل با ترنسفورمرها و کاربردهای آنها تا حدی آشنا شدیم. هدف این فصل بررسی کردن دقیق تر معماری و ساختار داخلی ترنسفورمرها است.



شکل ۲-۱: ساختار معماری ترنسفورمر

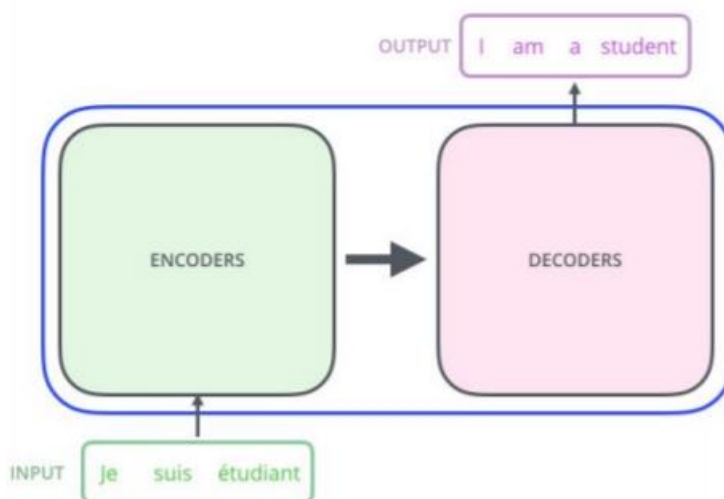
همانطور که در شکل ۲-۱ مشاهده می‌شود، ترنسفورمر از تعدادی بلاک و عملیات مختلفی بین این بلاک‌ها تشکیل شده است [4]. در ادامه این فصل به شرح جزئیات مربوط به این بلاک‌ها می‌پردازیم.

از دید بیرون و به عنوان مثال در کار ترجمه یا Machine Translation ترنسفورمر یک بلاک است که بر اساس ورودی که به آن می‌دهیم، کلمات متناظر خروجی را یکی پس از دیگری حدس می‌زند [5].



شکل ۲-۲: شمای کلی ترنسفورمر

اگر کمی بیشتر وارد ترنسفورمر شویم متوجه میشویم که ترنسفورمر از تعدادی Encoder و تعدادی Decoder که به هم متصل اند ساخته شده است. واحد Decoders در واقع استکی از انکدرها است که به صورت سری به هم متصل شده اند، این انکدرها از نظر ساختار داخلی دقیقا مشابه یکدیگر هستند اما وزن های مستقلی دارند. به طور مشابه دیکدرز هم استکی از دیکدرهای مشابه به هم متصل است. اتصال بین Encoders و Decoders به این صورت است که خروجی آخرین یا بالاترین انکدر به تک تک دیکدرها ارسال شده.



شکل ۲-۳: نگاهی جزئی تر به ساختمان ترنسفورمر

تعداد واحد های انکدر و دیکدر در مقاله Attention is all you need برابر ۶ در نظر گرفته شده است [4]. اما همانطور که میدانید این عدد یک هایپر پارامتر^۹ است و قابل تغییر می باشد.

اگر درون یک واحد انکدر را نگاه کنیم، خواهیم دید که شامل یک لایه Self Attention و یک لایه Feed Forward است. با مفهوم Attention در بخش های قبل آشنا شدیم و در این بخش بیشتر با Self Attention آشنا خواهیم شد. سلف

^۹ Hyperparameter

اتنشن این اجازه را می‌دهد که در حالی که یک کلمه ورودی در حال انکد شدن است، بتوانیم به سایر کلمات ورودی نگاه کنیم. خروجی لایه سلف اتنشن به لایه Feed Forward پاس داده خواهد شد.

یک واحد Decoder نیز بسیار شبیه به یک واحد انکدر است با این تفاوت که یک لایه میانی Encoder Decoder Attention دارد که به دیکدر کمک میکند توجه خود را روی قسمت مناسبی از خروجی متمرکز کند.

خوب است کمی بیشتر با نحوه کار انکدر ها آشنا شویم. اولین انکدر کلمات ورودی را به صورت لیستی از بردار ها دریافت میکند، اما ورودی انکدر های بعدی خروجی انکدر قبل از خودشان است که از لحاظ ابعاد با ورودی برابر هستند اما هر کدام معنی درونی خود را دارند.

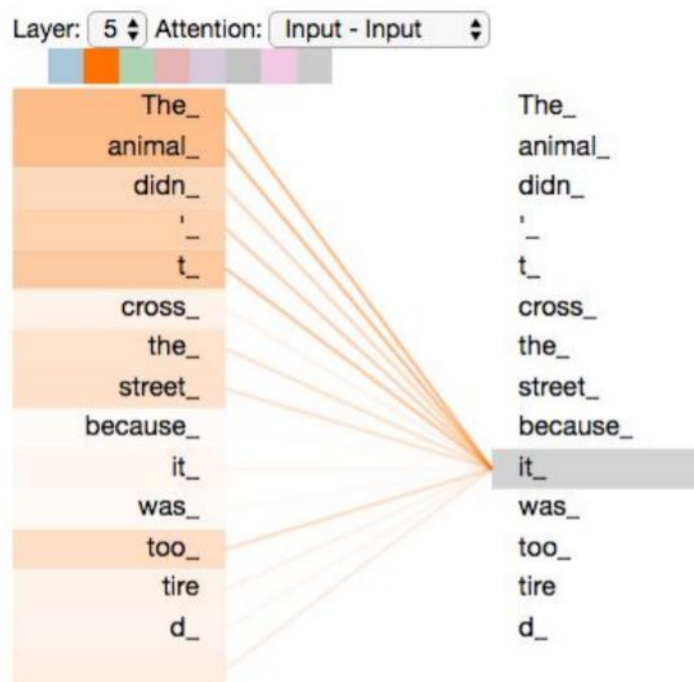
این کلمات ورودی به یک لایه Self Attention می‌روند و تعدادی ارتباطات بین کلمات بوجود می‌آید ولی بعد از خروج از این لایه مسیر مستقل خود را به واحد های کوچکتر Feed Forward پیش میگیرند. در این قسمت کلمات مختلف به هم وابسته نیستند که یکی از ویژگی های خیلی مهم ترنسفورمر ها است که قابلیت Parallelize شدن را به آنها میدهد [4].

۲-۲- سلف اتنشن^{۱۰}

نوبت آن است به صورت عمیق تر مفهوم سلف اتنشن را بررسی کنیم. فرض کنید جمله The animal didn't cross the street because it was too tired را داریم و میخواهیم بدانیم کلمه "it" به چه اشاره میکند. ما به عنوان انسان خیلی راحت میتوانیم متوجه شویم که منظور animal است و این قسمت از جمله در واقع Animal was too tired بوده است. اما این کار برای ماشین مخصوصا اگر فاصله بین animal و it زیاد باشد اصلا ساده نیست.

حال کاربرد سلف اتنشن این است که زمانی که مدل در حال پردازش یک کلمه است، میتواند به سایر کلمات نگاه بی‌اندازد و با توجه به آنها انکدینگ بهتری از کلمه مورد نظر تولید کند. بنابراین مشابه Hidden State ها در RNN ها، سلف اتنشن به ترنسفورمر کمک میکند درک بهتری از بقیه کلمات مرتبط درون کلمه فعلی Bake کند [5].

^{۱۰} Self Attention



شکل ۲-۴: نحوه ی کار مکانیزم سلف اتنشن

در شکل ۲-۴ میبینیم که ستون سمت چپ و راست هر دو جمله ورودی هستند و در ستون سمت راست کلمه *it* انتخاب شده است و در ستون سمت چپ میزان توجه به سایر کلمات ورودی دیده می شود. همانطور که انتظار میرفت در این مثال بیشترین توجه به کلمه *Animal* بوده است.

تصویری که در صفحه قبل دیدید توسط یک ابزار به اسم *Tensor 2 Tensor* تولید شده است که توسط تیم *Google Brain* توسعه داده شده و یک لایبرری مفید برای انجام تحقیقات در زمینه *Machine learning* می باشد.

۲-۳- نحوه عملکرد سلف اتنشن

نحوه عملکرد سلف اتنشن به این شکل است که علاوه بر فضای توصیفی^{۱۱} که از کلمات ورودی تولید شده است، سه بردار دیگر به ازای هر کلمه به اسامی *Query, Key, Value* در نظر میگیرد [4].

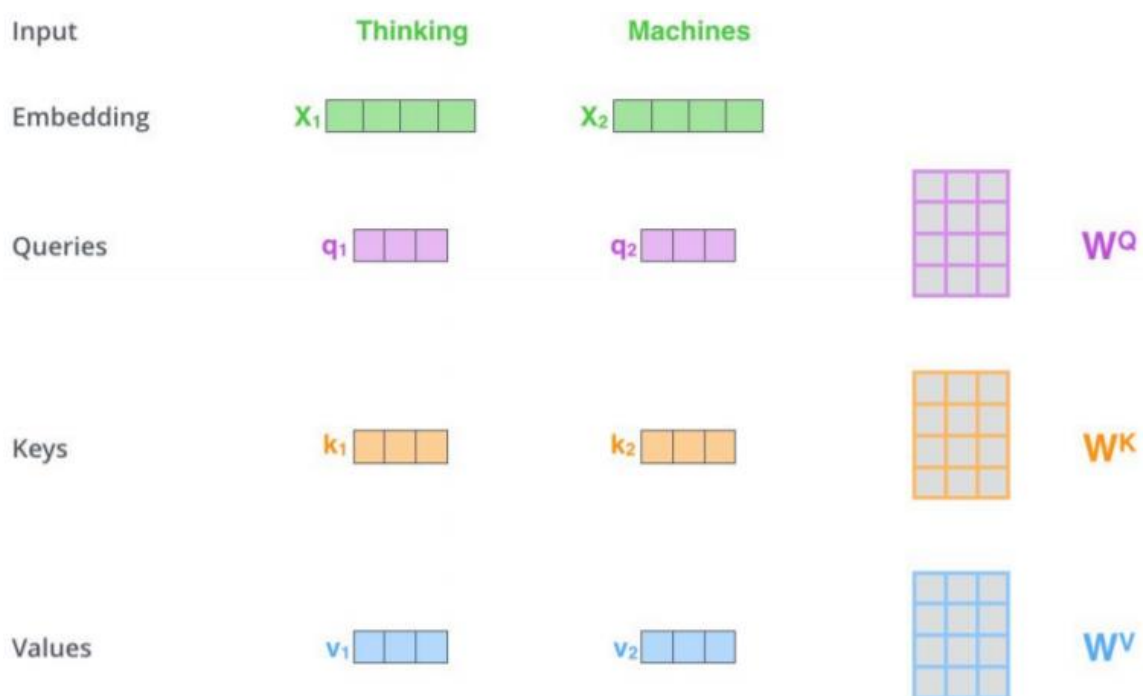
در مقاله *Attention is all you need* اندازه ورودی ۵۱۲ در نظر گرفته شده است [4] و اندازه این سه بردار ۶۴ است. نکته حائز اهمیتی که در مقاله ذکر شده به آن اشاره شده این است که این سه بردار لزوما نیاز نیست از اندازه ورودی کمتر باشند اما

^{۱۱} Embedding

کمتر بودن اندازه آنها کمک میکند در مراحل بعد Multi head self attention به صورت Constant time محاسبه شود. راجع به Multi head attention در ادامه توضیح خواهیم داد.

۲-۳-۱- بدست آوردن بردار های Q,K,V

برای محاسبه این سه بردار، سه ماتریس به نام های W_v , W_k , W_q در نظر میگیریم که در ابتدا با مقادیر تصادفی initialize شده اند. با ضرب کردن این ماتریس ها در بردار ورودی به ترتیب Q,K,V بدست می آیند. در مرحله Train شدن شبکه مقادیر این ماتریس ها آپدیت میشود و شبکه سعی میکند رفته رفته مقادیر مناسب تری برای این ماتریس ها پیدا کند. در واقع میتوان گفت این ماتریس ها ماتریس های وزن هستند.



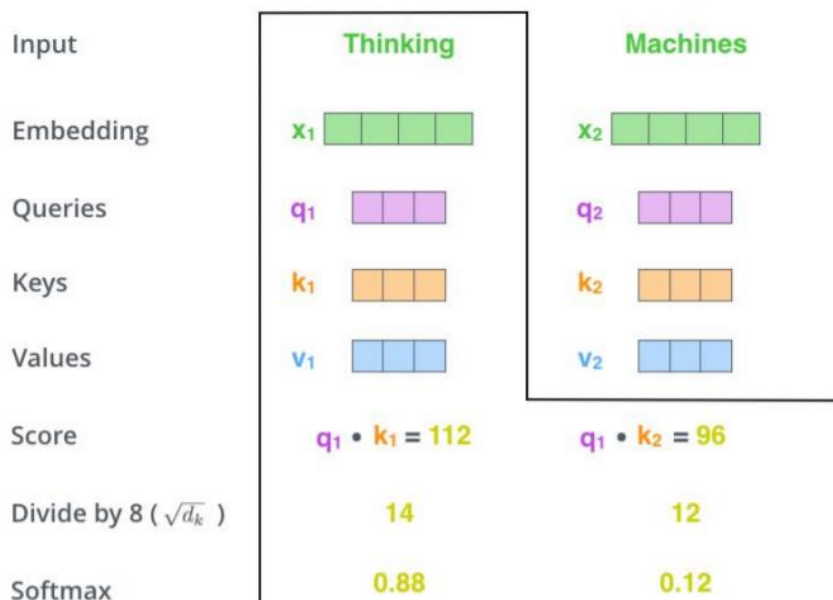
شکل ۲-۵: بدست آوردن بردار های Q,K,V

مرحله بعد برای محاسبه سلف اتنشن بدست آوردن یک Score است. اسکور برای هر کلمه ورودی با توجه به یک کلمه خاص، از ضرب کردن Query آن کلمه خاص در Key هر کلمه بدست می آید (در اینجا منظور از ضرب ضرب داخلی است).

اگر به شکل ۲-۵ دقت کنید بهتر متوجه این موضوع می شویم. در این مثال Query کلمه it در Key همه کلمات ضرب می شود و Score آن کلمه با توجه به it محاسبه می شود.

مرحله بعد تقسیم کردن اسکور بر یک عدد است، در مقاله پیشنهاد شده این عدد جزر dimensionality سه بر دار کمکی باشد [4] قاندا تا این عدد هم میتواند تغییر کند اما نتیجه استفاده از این عدد پایدار شدن گرادیان ها بوده است [5].

مرحله بعد عبور دادن این مقدار از تابع سافت مکس است که باعث می شود اعدادمان در بازه ۰ و ۱ قرار بگیرند و جمع همه آنها برابر ۱ شود که برای مقصود ما بسیار مناسب است.



شکل ۲-۶: مراحل بعدی محاسبه سلف اتنشن

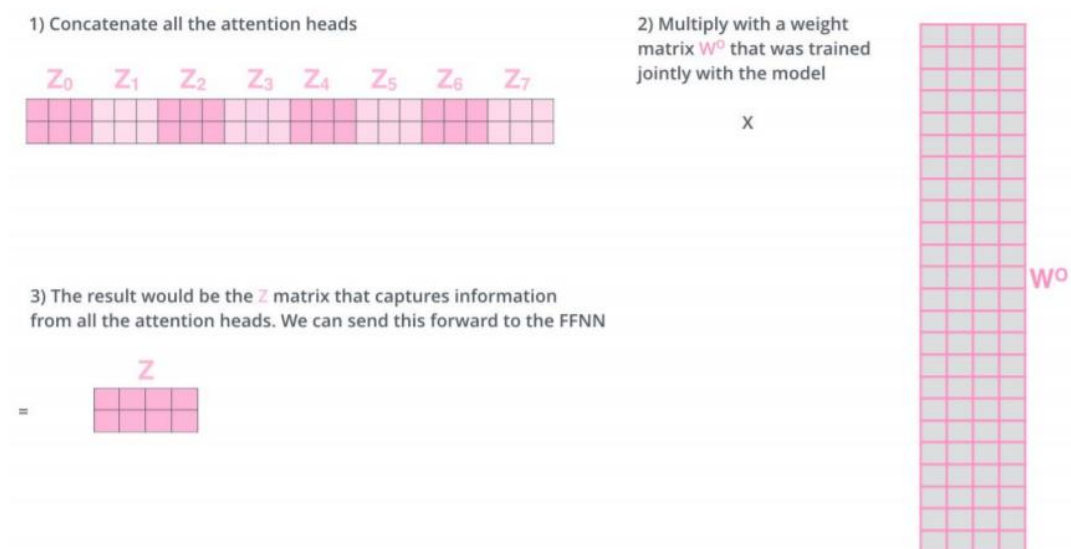
در نهایت برای ساده تر کردن محاسبات ، فرمول محاسبه سلف اتنشن را به صورت ماتریسی در می آوریم.

فرمول نهایی در نهایت یک خروجی به نام Z می دهد و ما این خروجی را به لایه Feed forward می دهیم [5].

$$\text{softmax} \left(\frac{Q \times K^T}{\sqrt{d_k}} \right) V = Z$$

شکل ۲-۷: استفاده تابع سافت مکس در محاسبه ی سلف اتنشن

کردن این خروجی ها بدین صورت است که ابتدا مانند شکل ۹-۲ ۸ خروجی را به هم Concatenate میکنیم سپس در ماتریس W^o ضرب میکنیم، حاصل Z نهایی ما است که آماده ارسال به لایه Feed Forward می باشد [5].



شکل ۹-۲: نحوه ی محاسبه ماتریس Z

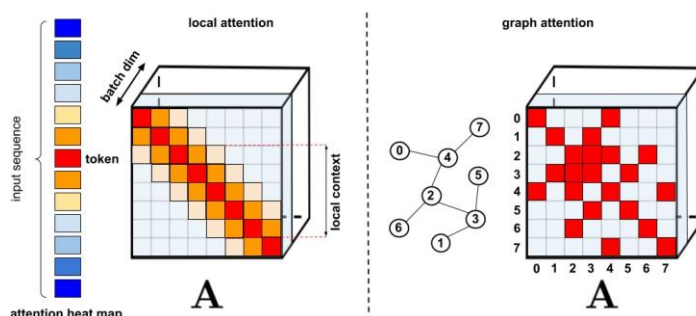
۳- فصل سوم: چالش های معماری ترنسفورمر

یکی از مشکلات اصلی معماری ترنسفورمر مکانیزم self-attention آن هاست [6]. که بدین صورت عمل میکند: مقدار بازنمایی هر توکن با بازدید از همه ی توکن ها بروز رسانی میشود تا بتواند داده ها را برای مدت زیادی ذخیره و نگه داری کند البته این امر باعث میشود تا پیچیدگی زمانی این الگوریتم از مرتبه نمایی و وابسته به اندازه ورودی باشد [6]. نمایی بودن پیچیدگی ترنسفورمر باعث افزایش زمان اجرا، افزایش مصرف برق و نیاز به سخت افزار های بیشتر میگردد. که از نقاط ضعف این الگوریتم است.

تلاش هایی برای بهبود پیچیدگی ترنسفورمرها شد که بطور خلاصه بررسی میکنیم:

۱. مکانیزم توجه اندک^{۱۳}

روش کار این مکانیزم بدین صورت است که هر توکن بجای توجه و محاسبات بازنمایی برای همه توکن ها به تعدادی توکن توجه کند. اینکار پیچیدگی نمایی ترنسفورمر ها را به مرتبه $O(n\sqrt{n})$ رساند [6]. البته این روش به صرفه نبود زیرا در ازای کاهش ۲ درصدی کارایی تنها ۲۰ درصد افزایش سرعت داشت و نمیتوانست بخوبی مسئله را حل کند.



شکل ۳-۱: روش های توجه اندک از اساس کار شبکه های عصبی گرافی نیز ایده گرفتند که با روابط خاص بصورت گراف توجه ظاهر می شوند.

متأسفانه ، روشهای توجه اندک هنوز هم محدودیتهای زیادی دارند.

۱. به عملیات ضرب ماتریس پراکنده کارآمد نیاز دارند که در همه شتاب دهنده ها موجود نیست.
۲. معمولاً تضمین های نظری دقیق برای قدرت تمثال خود ارائه نمیدهند.
۳. فقط برای مدل های ترنسفورمر و pre-training بهینه شدند.
۴. معمولاً لایه های توجه بیشتری را برای جبران نمایش های پراکنده روی هم میگذارند ، که استفاده از آن ها را با سایر مدل های pre-trained دشوار می کند، زیرا که نیاز به آموزش مجدد و مصرف قابل توجه انرژی دارند.

^{۱۳} sparsity attention mechanism

اگر از همه محدودیت‌هایی گفتیم بگذریم باز دو مشکل اساسی این مکانیزم این است که گاهی اوقات به اندازه‌ای که نیاز است از ورودی انتخاب نمی‌کند و مشکل بوجود می‌آورد همچنین نمی‌تواند از سربار یکسری عملگرهایی که خیلی سربار دارند کم کند نمونش softmax که برای نرمال سازی یا سامانه‌های پیشنهادگر زیاد استفاده می‌شود

۲. Reformer

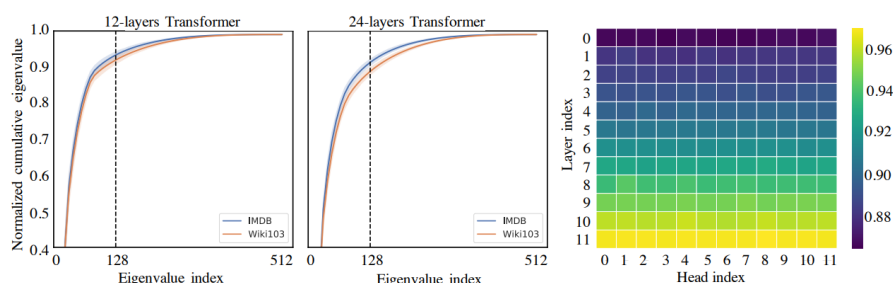
از دیگر تلاش‌های کاهش پیچیدگی می‌توان مقاله reformer را نام برد [7]. این مقاله از روش Locally sensitive hashing برای کاهش پیچیدگی استفاده می‌شود و پیچیدگی را به مرتبه $O(n \log n)$ کاهش می‌دهد [7]. این روش نیز خالی از ایراد نبود؛ از ایرادات این روش این بود که تأثیر چندانی برای رشته‌های کوتاه نداشت و تأثیر آن فقط در رشته‌هایی با طول بلند تر از ۲۰۴۸ چشم گیر و قابل توجه بود. دیگر مشکل این روش این بود که با استفاده از روش‌های رمزنگاری پیچیده کارایی را کاهش می‌داد.

این مشکلات باعث شد تا محققان به دنبال راه حل دیگری باشند.

۳-۱- ترنسفورمرهای خطی^{۱۴}

ترنسفورمر خطی توسط تیم فیسبوک در سال ۲۰۲۰ در قالب یک مقاله معرفی شد [6].

ایده اصلی این مقاله این بود که ماتریس توجه در ترنسفورمرها می‌تواند به کمک یک ماتریس درجه کمتر تخمین زده شود. در این مقاله از روش Singular value decomposition [8] برای ماتریس‌های Key, Value کمک گرفته شد تا بتوانند درجه ماتریس‌ها را کاهش و در نهایت پیچیدگی محاسبات را کاهش دهند. این مهم باعث شد تا ترنسفورمرهای خطی هم در مرحله آموزش و هم در مرحله استنتاج پیچیدگی خطی داشته باشند. نام گذاری اینگونه ترنسفورمرها نیز به دلیل مرتبه پیچیدگی آن‌ها بود (Linformer (Linear + Transformer).



شکل ۳-۲: نمودارهای بالا نشان می‌دهند که می‌توان نتیجه محاسبات ترنسفورمرها را با ماتریس‌های درجه کمتر تخمین زد. باید توجه داشت که اگر داده‌ها در ابعاد مختلف ماتریس بطور یکنواخت توزیع شده باشند یک نمودار خطی قابل انتظار بود. توی این شکل مشخص است که اکثر تأثیر داده‌ها در ۱۲۸ بعد نخست وجود دارد و در ۳۸۴ بعد، بعدی تأثیر کمی وجود دارد.

این نمودار نتیجه آزمایش بر روی ترنسفورمرهای ۱۲ و ۲۴ لایه و روی مجموعه داده‌گان IMDB و WIKI103 بوده است.

^{۱۴} Linformer

درجه ماتریس:

در جبر خطی درجه یک ماتریس تعداد بعد فضای برداری است که توسط ستون های آن ماتریس پوشش داده میشود به بیان ساده تر درجه ماتریس تعداد سطرها یا ستون های مستقل خطی یک ماتریس است [9].

$$\begin{bmatrix} 1 & 0 & 1 \\ -2 & -3 & 1 \\ 3 & 3 & 0 \end{bmatrix} \quad \begin{bmatrix} 1 & -1 \\ 1 & -1 \\ 0 & 0 \\ 2 & -2 \end{bmatrix}$$

به عنوان مثال در ماتریس سمت راست ستون وسط را میتوان از تفریق ستون های چپ و راست بدست آورد. پس تنها ۲ ستون مستقل خطی دارد و درجه آن ۲ است.

یا در ماتریس سمت چپ هر دو ستون از روی یکدیگر با قرینه کردن درایه ها میتوان بدست آورد. پس درجه آن ۱ است.

۳-۱-۱-۱- نتایج بنچمارک ها

جدول ۳-۱: بنچمارک های ترنسفورمر خطی

n	Model	SST-2	IMDB	QNLI	QQP	Average
512	Liu et al. (2019), RoBERTa-base	93.1	94.1	90.9	90.9	92.25
	Linformer, 128	92.4	94.0	90.4	90.2	91.75
	Linformer, 128, shared kv	93.4	93.4	90.3	90.3	91.85
	Linformer, 128, shared kv, layer	93.2	93.8	90.1	90.2	91.83
	Linformer, 256	93.2	94.0	90.6	90.5	92.08
	Linformer, 256, shared kv	93.3	93.6	90.6	90.6	92.03
	Linformer, 256, shared kv, layer	93.1	94.1	91.2	90.8	92.30
512	Devlin et al. (2019), BERT-base	92.7	93.5	91.8	89.6	91.90
	Sanh et al. (2019), Distilled BERT	91.3	92.8	89.2	88.5	90.45
1024	Linformer, 256	93.0	93.8	90.4	90.4	91.90
	Linformer, 256, shared kv	93.0	93.6	90.3	90.4	91.83
	Linformer, 256, shared kv, layer	93.2	94.2	90.8	90.5	92.18

نتیجه بنچمارک هایی که ترنسفورمر خطی بر روی آن ها اجرا شده است این مهم را نشان میدهد که ترنسفورمر خطی در همه مواقع سریعتر از ترنسفورمر معمولی عمل کرده است [6]؛ ولی معمولاً نتایج بدتری به همراه داشته است. هرچند که در مواردی خاص توانسته از ترنسفورمر معمولی هم نتیجه بهتری داشته باشد. در این بنچمارک ها نشان داده شد که کارایی ترنسفورمر خطی به ابعاد وابسته نیست و به یک مقدار ثابت وابسته است که اثباتی بر خطی بودن آن است [6].

۳-۱-۲- ترنسفورمر خطی در مرحله استنتاج

ترنسفورمر خطی علاوه بر خطی عمل کردن در مرحله آموزش، در مرحله استنباط هم خطی عمل میکند و بهبود های زیادی را در زمینه استنتاج یا استفاده کردن از مدل آموزش داده شده میتوانیم شاهد باشیم.

جدول ۳-۲: این جدول نشان دهنده یک عمل استنباط با استفاده از یک GPU تسلا وی ۱۰۰ با ۱۶ گیگ وی رم است. جدول سمت چپ نشان دهنده زمان صرفه جویی شده نسبت به استفاده از ترنسفورمر معمولی و جدول سمت راست نشون دهنده حافظه صرفه جویی شده نسبت به استفاده از ترنسفورمر معمولی است.

length n	projected dimensions k					length n	projected dimensions k				
	128	256	512	1024	2048		128	256	512	1024	2048
512	1.5x	1.3x	-	-	-	512	1.7x	1.5x	-	-	-
1024	1.7x	1.6x	1.3x	-	-	1024	3.0x	2.9x	1.8x	-	-
2048	2.6x	2.4x	2.1x	1.3x	-	2048	6.1x	5.6x	3.6x	2.0x	-
4096	3.4x	3.2x	2.8x	2.2x	1.3x	4096	14x	13x	8.3x	4.3x	2.3x
8192	5.5x	5.0x	4.4x	3.5x	2.1x	8192	28x	26x	17x	8.5x	4.5x
16384	8.6x	7.8x	7.0x	5.6x	3.3x	16384	56x	48x	32x	16x	8x
32768	13x	12x	11x	8.8x	5.0x	32768	56x	48x	36x	18x	16x
65536	20x	18x	16x	14x	7.9x	65536	60x	52x	40x	20x	18x

۳-۱-۳- کاربردهای ترنسفورمر خطی

از کاربردها و تاثیر های ترنسفورمر خطی در بلند مدت میتوانیم موارد زیر را نام ببریم [6]:

۱. کار کردن روی کارهای مربوط به تصویر
۲. انجام کارهایی که صرفا بخاطر هزینه بر بودن قابل انجام نبودن
۳. ذخیره انرژی و کمک به محیط زیست (به دلیل بهینگی و مصرف کمتر انرژی برق)

از معایب احتمالی این ساختار میتوان به مشکلات اخلاقی و استفاده نادرست از این قدرت اشاره کرد. مثل استفاده از این شبکه در شکستن رمزها و ... [6]

۳-۲- پرفورمر^{۱۵}

پرفورمر در تاریخ ۳۰ سپتامبر ۲۰۲۰ و توسط تیم تحقیقاتی گوگل به عنوان یک مقاله کنفرانسی ارائه شد [10].

شیوه کار پرفورمر [11]:

time complexity: $\mathcal{O}(L^2d)$

dimensions: $(L \times L)(L \times d) = (L \times d)$

$$\begin{pmatrix} L \times d & d \times L & L \times d \\ Q & K^T & V \end{pmatrix}$$

dimensions: $(L \times d)(d \times L) = (L \times L)$

time complexity: $\mathcal{O}(L^2d)$

time complexity: $\mathcal{O}(d^2L)$

dimensions: $(L \times d)(d \times d) = (L \times d)$

$$\begin{pmatrix} L \times d & d \times L & L \times d \\ Q & (K^T V) \end{pmatrix}$$

dimensions: $(d \times L)(L \times d) = (d \times d)$

time complexity: $\mathcal{O}(d^2L)$

شکل ۳-۳: بررسی اساس کار ترانسفورمر

اساس کار ترنسفورمر سه ماتریس key, query, value است که طبق شکل سمت چپ ضرب میشوند و با توجه به ابعادشون و بقیه مراحل سربار سنگینی برای شبکه ایجاد میکنند محققین گوگل با بررسی ابعاد این ماتریس ها به این نتیجه رسیدند که چون

^{۱۵} performer

d از L کوچکتر است؛ اگر مثل تصویر سمت راست بتوانیم ضرب را انجام دهیم خیلی سربار کاهش پیدا میکند ولی در مرحله محاسبه softmax بخاطر ذات نمایی این الگوریتم تاثیر سربار چندین برابر میشود. پس رهایی از الگوریتم softmax و یک تخمین خوب برای مقدارش میتواند مشکل را حل کند [10].

بعد از محاسبات ریاضی بسیار زیاد محققین توانستن با ضرب دوتا مقدار K' و Q' نتیجه تابع softmax را تقریب بزنند. اینکار سربار نمایی الگوریتم به ازای هر درایه ماتریس را حذف میکرد.

$$\text{softmax} \left(\frac{\mathbf{QK}^T}{\sqrt{d}} \right) = \mathbf{D}^{-1} \mathbf{A} \quad (1-3)$$

$$\text{with: } \mathbf{A} = \exp(\mathbf{QK}^T / \sqrt{d}), \quad \mathbf{D} = \text{diag}(\mathbf{A} \mathbf{1}_L)$$

$$\mathbf{A}(i, j) = K(\mathbf{q}_i, \mathbf{k}_j) = \exp(\mathbf{q}_i \mathbf{k}_j^T) = \phi(\mathbf{q}_i)^T \phi(\mathbf{k}_j) \quad (2-3)$$

اگر ماتریس نتیجه تابع softmax باشد، برای تخمین زدن آن میتوان از کرنل استفاده کرد که این کرنل مقداری متغیر است و میتواند از نوع گوسی یا مثلثاتی یا... باشد.

$$\phi(\mathbf{x}) = \frac{h(\mathbf{x})}{\sqrt{m}} (f_1(w_1^T \mathbf{x}), \dots, f_1(w_m^T \mathbf{x}), \dots, f_l(w_1^T \mathbf{x}), \dots, f_l(w_m^T \mathbf{x})) \quad (3-3)$$

$$h(\mathbf{x}) = 1, l = 2, f_1 = \sin, f_2 = \cos, \mathcal{D} = \mathcal{N}(0, \mathbf{I}_d) \quad (4-3)$$

برای محاسبه ϕ ها این رابطه موجود است که اثبات آن در مقاله بطور کامل آورده شده است. و اگر در مورد ویژگی های فوریه تصادفی اطلاعاتی داشته باشید درک آن برایتان ساده تر خواهد بود. کاربرد ϕ ها این است که مسئله را از بُعد خطی به یک بُعد دیگر ببرد. همچنین می دانیم که اگر ϕ_q را در ϕ_k ضرب کنیم یک تقریب خوب از تابع softmax بدست خواهد آمد [11].

در رابطه محاسبه ϕ چندین امگا داریم که بردارهای تصادفی از یک توزیع احتمال مشخص مثل d اند که مثلا در این رابطه از صفر تا i است که میتوان گفت یک توزیع نرمال است.

h یک تابع قطعی از ورودی است که با تقسیم بر \sqrt{m} نرمال میشود. اینجا h را فعلا یک در نظر گرفته شده است و D نرمال سازی خطی بوده است. انتخاب h و f یا توابع نوع کرنل را مشخص میکنند که مثلا اینجا گوسی است [11].

$$h(\mathbf{x}) = \exp\left(\frac{\|\mathbf{x}\|^2}{2}\right) \quad (5-3)$$

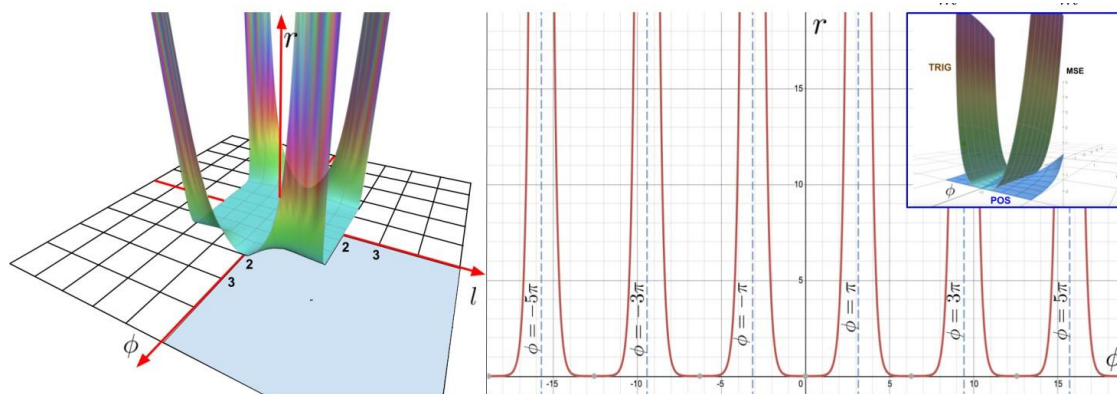
با کمی تغییر تابع h میتوانیم یک کرنل برای softmax بسازیم با مشکلاتی روبرو خواهیم شد، مهمترین مشکل آن این است که تابع softmax خروجی همیشه مثبت دارد در حالی که این تخمین ممکن است خروجی منفی نیز داشته باشد.

برای رفع این مشکل یک رابطه جدید برای کرنل softmax پیشنهاد شد اگر بخواهیم این کرنل را با همان رابطه ϕ که داشتیم بدست بیاوریم ثابت ها و توابع مجهول به صورت زیر در خواهند آمد [11].

$$\begin{aligned} K_{SM}(\mathbf{x}, \mathbf{y}) &= \exp(\mathbf{x}^\top \mathbf{y}) = \mathbb{E}_{w \sim \mathcal{N}(0, \mathbf{I}_d)} [\exp(w^\top \mathbf{x} - \frac{\|\mathbf{x}\|^2}{2}) \exp(w^\top \mathbf{y} - \frac{\|\mathbf{y}\|^2}{2})] \\ &= \mathbb{E}_{w \sim \mathcal{N}(0, \mathbf{I}_d)} [\exp(-\frac{\|\mathbf{x}\|^2}{2}) \exp(-\frac{\|\mathbf{y}\|^2}{2}) \exp(w^\top \mathbf{x}) \exp(w^\top \mathbf{y})] \end{aligned} \quad (6-3)$$

$$h(\mathbf{x}) = \exp(-\frac{\|\mathbf{x}\|^2}{2}), l = 1, f_1 = \exp, \mathcal{D} = \mathcal{N}(0, \mathbf{I}_d) \quad (7-3)$$

بعد از انجام محاسبات یک نمودار از نتیجه ترسیم شد که یک مشکل اساسی را نشان میداد که اگر با همون کرنل قبلی یا اصطلاحاً کرنل مثلثاتی کار کنیم خطا در اطراف صفر به سمت بی نهایت میرود، در صورتی که برای کرنل جدید خطا به سمت صفر میرفت.



شکل ۳-۴: خطای مشاهده شده در کرنل مثلثاتی

در مقاله کرنل های دیگری نیز از تابع softmax اثبات شده اند [10].

بعد از انجام این محاسبات پیچیده سربار محاسباتی از شبکه حذف میشود و بجای تابع softmax خواهیم داشت:

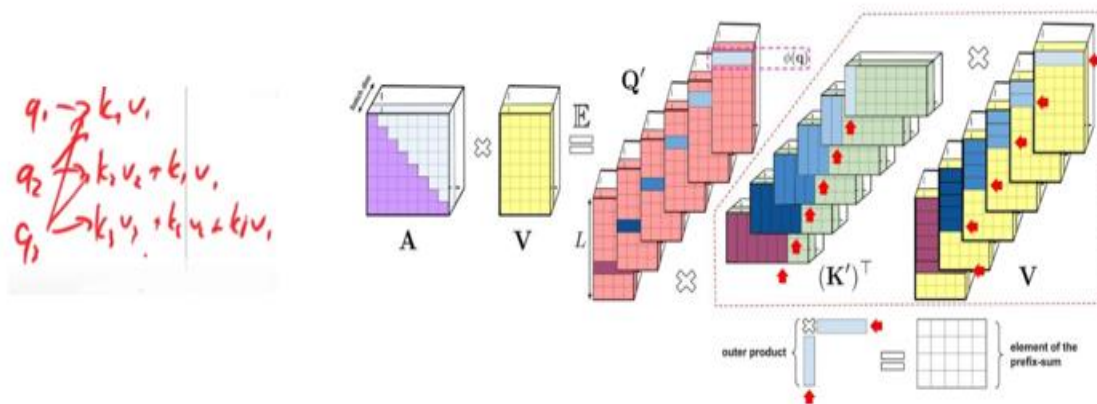
$$\begin{array}{c}
 \text{time complexity: } \mathcal{O}(Lmd) \\
 \uparrow \\
 \text{dimensions: } (L \times m)(m \times d) = (L \times d) \\
 \begin{array}{ccc}
 L \times m & m \times L & L \times d \\
 Q' & (K'^T & V)
 \end{array} \\
 \downarrow \\
 \text{dimensions: } (m \times L)(L \times d) = (m \times d) \\
 \downarrow \\
 \text{time complexity: } \mathcal{O}(Lmd)
 \end{array}$$

شکل ۳-۵: مرتبه پیچیدگی پرفورمر

که این تابع از مرتبه پیچیدگی $\mathcal{O}(Lmd)$ است و براحتی میتوان گفت که پیچیدگی آن خطی است [11].

۳-۲-۱- توجه بر اساس پیشوندها

خب ممکن است که بخواهیم این روش را در مدل های عمومی تری مانند GPT پیاده کنیم؛ در این مدل ها چون ماتریس A خلوت^{۱۶} است ما از روش جمع پیشوندی استفاده میکنیم که رابطه آن در شکل زیر آورده شده است.

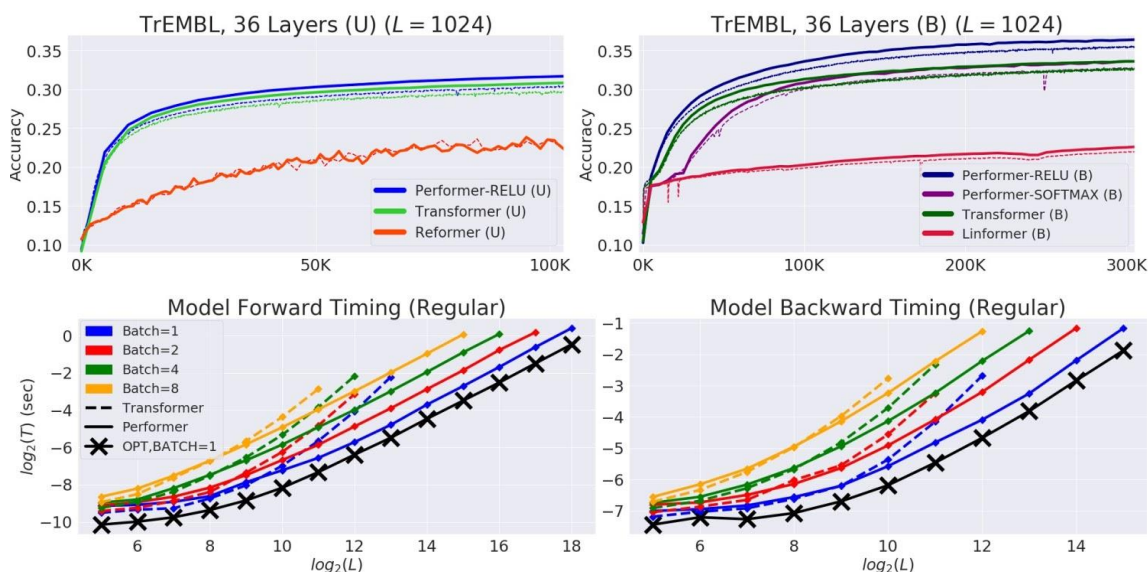


شکل ۳-۶: توجه بر اساس پیشوندها

برتری پرفورمر به نسبت ترنسفورمر خطی در این مورد این است که پرفورمر جداسازی ماتریس ها از یکدیگر را بهتر انجام داده است [10].

^{۱۶} Sparse

۲-۲-۳- مقایسه ها

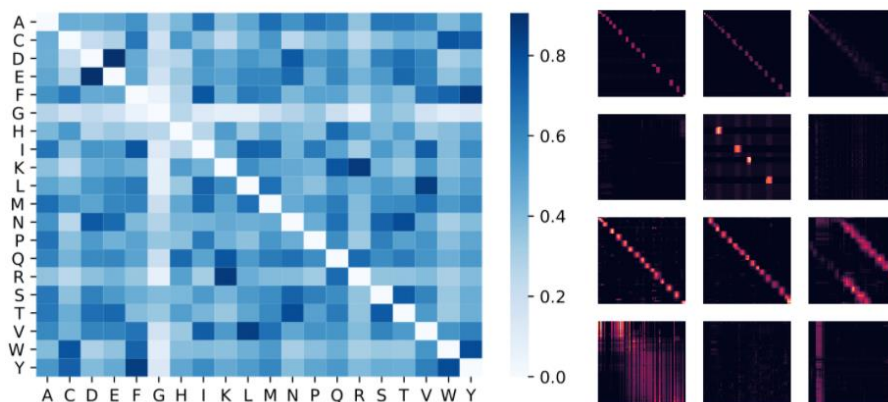


شکل ۷-۳: نمودارهای دقت و زمان اجرای الگوریتم های مختلف

نمودارهای بالا دقت و زمان اجرای الگوریتم های مختلف را باهم مقایسه کرده اند که در مجموع میتوان گفت که پرفورمر عملکرد بسیار خوبی از خودش نشان داده و زمان اجرای کمتری دارد .

۳-۲-۳- حوزه های کاربرد پرفورمر

۱. بیولوژی و داروسازی
۲. حفظ محیط زیست (زمان اجرای کمتر و مصرف کمتر انرژی)
۳. تحقیقات بر روی شبکه های ترنسفورمر
۴. سازگاری با روش ها و مدل های موجود



شکل ۸-۳: نتیجه استفاده پرفورمر در مبحث ژنتیک

یکی از زمینه های کاربرد پرفورمر که در مقاله بررسی شده است، مبحث بیولوژی و ژنتیک می باشد [10] که همیشه چالش طول بلند ژنوم ها در این حوزه خیلی از شبکه ها را رفوزه میکرد. پرفورمر با استفاده از تخمین و کرنل هایش توانسته به این چالش غلبه کند و به نتایج خوبی در این زمینه برسد .

۴- فصل چهارم: آشنایی عملی با ترنسفورمرها (و محصولات Hugging face)

۴-۱- کمپانی هاگینگ فیس^{۱۷}

در این بخش می‌خواهیم با پیاده سازی و استفاده عملی از ترنسفورمر ها بیشتر آشنا شویم. ابتدا خوب است کمی با کمپانی هاگینگ فیس و خدماتشان آشنا شویم. در بخش بعدی به صورت عملی یک مدل ترنسفورمر ترین شده را برای یک تسک خاص (کامنت های IMDB) فاین تون^{۱۸} (تنظیم) میکنیم.

شرکت هاگینگ فیس یکی از شرکت های اصلی فعال در زمینه NLP و بخصوص ترنسفورمر هاست. یکی از خدمات اصلی و مهم این شرکت در دسترس قرار دادن انواع مدل های ترنسفورمر از قبل ترین شده که توسط شرکت های بزرگی مانند google, facebook و ... منتشر شده اند است.

همچنین این شرکت سرویس های Hosted inference را هم برای مدل ها ارائه میدهد تا نیاز نداشته باشید خودتان به فکر دیپلوی کردن و پیکربندی محیط برای استفاده از نتیجه مدلتان باشید. این سرویس یک نسخه رایگان دارد که تا ۳۰۰۰۰ کاراکتر ورودی در ماه را پوشش میدهد و نسخه های غیر رایگان با امکانات بیشتر هم موجود است.

یکی دیگر از سرویس هایی که این شرکت ارائه میدهد سرویس نوپای Auto NLP است. هدف این سرویس این است که در آینده ای نه چندان دور به راحتی وارد کردن چند دستور بتوانیم مدل های NLP خود را به آسانی ترین و دیپلوی کنیم.

۴-۲- کتابخانه Transformers

اصلی ترین محصول این شرکت کتابخانه Transformer برای پایتون است. این لایبرری که از PyTorch و TensorFlow که دو فریم ورک اصلی ماشین لرنینگ در پایتون هستند پشتیبانی میکند، امکانات بسیار زیادی برای کار با ترنسفورمر ها در اختیار علاقه مندان قرار میدهد.

^{۱۷} HuggingFace

^{۱۸} Fine-Tune

بخشی از مدل های موجود در این لایبرری را در شکل ۴-۱ مشاهده میکنید.

Model	Tokenizer slow	Tokenizer fast	PyTorch support	TensorFlow support	Flax Support
ALBERT	✓	✓	✓	✓	✗
BART	✓	✓	✓	✓	✓
BERT	✓	✓	✓	✓	✓
Bert Generation	✓	✗	✓	✗	✗
BigBird	✓	✓	✓	✗	✓
BigBirdPegasus	✗	✗	✓	✗	✗
Blenderbot	✓	✗	✓	✓	✗
BlenderbotSmall	✓	✗	✓	✓	✗
CLIP	✓	✓	✓	✗	✓
CTRL	✓	✗	✓	✓	✗

شکل ۴-۱: بخشی از مدل های موجود در هاگینگ فیس

برای شروع به کار با کتابخانه Transformer نیاز است از طریق pip آن را نصب کنید. پس از نصب یکی از ساده ترین روش ها برای شروع به کار با ترنسفورمر ها استفاده از پایپلاین^{۱۹} ها هستند. پایپلاین ها در واقع مسیر های از قبل تعیین شده ای هستند که برای تسک های پر کاربرد NLP در این کتابخانه قرار داده شده اند.

```
>>> from transformers import pipeline

>>> question_answerer = pipeline("question-answering")

>>> context = r"""
... Extractive Question Answering is the task of extracting an answer from a text given a question. An example of a
... question answering dataset is the SQuAD dataset, which is entirely based on that task. If you would like to fine-tune
... a model on a SQuAD task, you may leverage the examples/pytorch/question-answering/run_squad.py script.
... """
```

شکل ۴-۲: استفاده از کتابخانه ترنسفورمر

توصیه می شود فایل notebook [12] پیوست شده به گزارش را برای دیدن جزئیات پیاده سازی و مشاهده خروجی کار بررسی کنید. چندین پایپلاین پر کاربرد در این notebook با پارامتر های مختلف بررسی شده اند. در ادامه چند مثال از پایپلاین را مشاهده میکنید.

^{۱۹} Pipeline

```

In [ ]:
classifier = pipeline('sentiment-analysis')

In [ ]:
classifier('i love you')

Out[ ]:
[{'label': 'POSITIVE', 'score': 0.9998656511306763}]

In [ ]:
classifier('i hate you')

Out[ ]:
[{'label': 'NEGATIVE', 'score': 0.9991129040718079}]

```

شکل ۴-۳: استفاده از کتابخانه ترنسفورمر در مبحث آنالیز احساس

```

In [ ]:
mlm = pipeline("fill-mask")

In [ ]:
ml = mlm(f"HuggingFace is creating a {mlm.tokenizer.mask_token} that the
for i in ml:
    print(i["score"], "\t", i["token_str"])

0.17927460372447968    tool
0.1134939044713974    framework
0.05243545398116112    library
0.03493543714284897    database
0.02860247902572155    prototype

```

```

In [ ]:
mlpers = mlmpersian(f"[MASK] رتبه من در کنکور ارشد 1400 برابر با")

print("رتبه من در کنکور ارشد 1400 برابر با *** خواهد شد")
for i in mlpers:
    print(i["score"], "\t", i["token_str"])

رتبه من در کنکور ارشد 1400 برابر با *** میشود
0.0288712065666914    ۶۳
0.01915029250085354    ۲
0.01827997900545597    ۹۰
0.017009330913424492    \
0.01648857071995735    ۷

```

شکل ۴-۴: استفاده از کتابخانه ترنسفورمر در مبحث ماسکینگ

۴-۳- مثال FineTune کردن ترنسفورمر

همانطور که گفته شد ترنسفورمر ها از زبان و گرامر درک دارند (Language Understanding) یکی از قدرتمند ترین کار هایی که میتوان با ترنسفورمر ها انجام داد این است که از دانش آنها نسبت به زبان به عنوان زیر بنا برای سایر تسک های NLP استفاده کرد. این کار را میتوان با فاین تون کردن لایه آخر (Classifier) انجام داد. به این ترتیب میتوان با Freeze کردن وزن های درونی ترنسفورمر و ترین کردن لایه آخر آن روی دیتاستی معین، مدل را برای آن تسک بخصوص آماده کرد.

در مثالی که پیوست شده ما از ترنسفورمر ترین شده ی DistilBERT استفاده میکنیم و بدون نیاز به ترین کردن مجدد ترنسفورمر و تنها با ترین کردن مدل بر روی دیتاست نسبتاً کوچکی از کامنت های IMDB ، توانایی دسته بندی این گونه کامنت ها را بدست آوریم [13]. توضیحات تکمیلی در مورد این کار در ویدیو چهارم وجود دارد.

جهت دسترسی به نسخه اصلی این notebook که توسط کمپانی HuggingFace آماده شده به آدرس زیر مراجعه کنید.

https://colab.research.google.com/github/huggingface/notebooks/blob/master/transformers_doc/tensorflow/custom_datasets.ipynb

۵- فصل پنجم: جمع بندی و نتیجه گیری

۵-۱- جمع بندی نهایی

در دنیای امروزی که محاسبات پیچیده و سنگین و پردازش بر روی داده های حجیم جزو جدایی ناپذیر رشد و توسعه میزنس ها است، یافتن راه هایی برای کاهش هزینه ها و افزایش کارایی این اعمال بسیار مورد توجه واقع شده. چند سالی است که با قدرتمند تر شدن GPU ها و قابلیت های پردازش موازی آنها شبکه های عصبی عمیق بسیار اهمیت پیدا کرده اند و بسیاری از شبکه های عصبی به عنوان جهشی بزرگ در فناوری یاد میکنند [4].

ترنسفورمر ها جهشی بزرگ و رو به جلو در تمام این زمینه ها هستند. قابلیت پردازش میلیون ها داده و نگهداری میلیون ها پارامتر قابل ترین، قابلیت موازی سازی بسیار زیاد و عوامل دیگر باعث شده اند توجه زیادی به ترنسفورمر ها شود و بسیاری از سرویس های مهم و بزرگ امروزی مانند موتور جستجوی google به سرعت هر چه تمام تر از این تکنولوژی استقبال کنند و از ترنسفورمر ها به عنوان بخش اصلی عملیات خود استفاده کنند.

همچنین تحقیقات نشان داده بر خلاف چیزی که به نظر می رسد ترنسفورمر ها نه تنها در زمینه های NLP بلکه در بسیاری زمینه های دیگر مانند Computer Vision و ... نیز بسیار موثر واقع می شوند.

علاوه بر اینها تحقیقات اخیر روی ترنسفورمر ها حاکی از آن است که قابلیت گسترش و افزایش کارایی ترنسفورمر ها تا چندین و چند برابر هم مهیاست و روز به روز نتایج جدید و شگفت انگیز تری از آنها مشاهده میکنیم.

اگر به مقالات چاپ شده در مورد ترنسفورمر ها دقت کنید، خواهید دید که ترنسفورمر ها روز به روز در حال رشد و قوی تر شدن هستند. معماری های جایگزینی که در این پروژه بررسی شدند مانند Linformer [6] و Performer [10] هر کدام به تنهایی چندین برابر ترنسفورمر ها را قوی تر کرده اند و علاوه بر آن خود این معماری ها هم عاری از نقص نیستند و امکان بهبود دادن به این معماری های جدیدتر هم به میزان قابل توجهی امکان پذیر است.

انجام کارهای پژوهشی در زمینه ترنسفورمر ها نیز بسیار داغ است و مملو از سوال پاسخ داده نشده می باشد. در نهایت میتوان گفت که آینده از آن Transformer هاست.

لازم است یادآور شویم که علاوه بر گزارشی که در حال مطالعه آن هستید، تعدادی اسلاید (Google slide) و چهار ویدیو یوتوب و تعدادی فایل Google Colab برای کد های عملی فصل چهار نیز موجود میباشد که در قسمت منابع قابل دسترس هستند.

- [1] Wikipedia. Recurrent neural network. https://en.wikipedia.org/wiki/Long_short-term_memory.
- [2] Wikipedia. Long short-term memory.
https://en.wikipedia.org/wiki/Recurrent_neural_network .
- [3] Wikipedia. Gated recurrent unit. https://en.wikipedia.org/wiki/Gated_recurrent_unit .
- [4] Vaswani, Ashish & Shazeer, Noam & Parmar, Niki & Uszkoreit, Jakob & Jones, Llion & Gomez, Aidan & Kaiser, Lukasz & Polosukhin, Illia. (2017). Attention Is All You Need. arXiv:1706.03762 .
- [5] Alammam, Jay. Illustrated-transformer. <http://jalammar.github.io/illustrated-transformer/> .
- [6] Wang, Sinong & Z. Li, Belinda & Khabisa, Madian & Fang, Han & Ma, Hao. (2020). Linformer: Self-Attention with Linear Complexity. arXiv:2006.04768v3 .
- [7] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In International Conference on Learning Representations, 2020
- [8] Wikipedia. Singular value decomposition.
https://en.wikipedia.org/wiki/Singular_value_decomposition .
- [9] Kriventsov, Stan. Linformer. Linformer: Self-Attention with Linear Complexity (paper review). <https://medium.com/deep-learning-reviews/linformer-self-attention-with-linear-complexity-paper-review-70aac1d197d9> .
- [10] Choromanski, Krzysztof & Likhoshesterov, Valerii & Dohan, David & Song, Xingyou & Gane, Andreea & Sarlos, Tamas & Hawkins, Peter & Davis, Jared & Mohiuddin, Afroz & Kaiser, Lukasz & Belanger, David & Colwell, Lucy & Weller, Adrian. (2020). Rethinking Attention with Performers. arXiv:2009.14794v1.
- [11] Yannic Kilcher. Rethinking Attention with Performers (Paper Explained).
https://www.youtube.com/watch?v=xJrKIPwVwGM&list=PL3naYV_Z4HeEjbvIRa59x8khw aZ7gcytN&index=22 .
- [12] Amirmahani, Amir Hassan & Dar, Gholamreza. Introduction to Hugging Face Transformer library.ipynb.
- [13] Amirmahani, Amir Hassan & Dar, Gholamreza. Fine Tuning Transformer using IMDB dataset.ipynb.
- [14] Vu, Kevin & Corp, Exxact. A Deep Dive Into the Transformer Architecture – The Development of Transformer Models. <https://www.kdnuggets.com/2020/08/transformer-architecture-development-transformer-models.html>.
- [15] Manzil Zaheer, Guru Guruganesh, Avinava Dubey, Joshua Ainslie, Alberti, Chris & Ontanon, Santiago & Pham, Philip & Ravula, Anirudh & Wang, Qifan & Yang, Li & Ahmed, Amr. (2020). Big Bird: Transformers for Longer Sequences. arXiv:2007.14062 .

[16] Campagnola, Chiara. From Transformers to Performers: Approximating Attention.
<https://towardsdatascience.com/from-transformers-to-performers-approximating-attention-69c88af0b11f>.

[17] The A.I. Hacker - Michael Phi. Illustrated Guide to Transformers Neural Network: A step by step explanation.
https://www.youtube.com/watch?v=4Bdc55j80l8&list=PL3naYV_Z4HeEjbvIRa59x8khwaZ7gcytN&index=2 .

[18] Yannic Kilcher. Linformer: Self-Attention with Linear Complexity (Paper Explained).
https://www.youtube.com/watch?v=-_2AF9Lhweo&list=PL3naYV_Z4HeEjbvIRa59x8khwaZ7gcytN&index=20&ab_channel=YannicKilcher .

[19] Transformers and Language Models - YouTube Playlist.
https://youtube.com/playlist?list=PL3naYV_Z4HeEjbvIRa59x8khwaZ7gcytN .

Abstract

For many years, LSTMs and GRUs have been the go-to method for almost every NLP²⁰ and language understanding task. Just like any method, LSTMs and GRUs have their own flaws and weaknesses.

In this work we discussed the creation of Transformers, an Attention based alternative to LSTMs and GRUs , that has shown to be extremely scalable and performant.

We also researched about state of the art variations of Transformers like Performers and Linformers to name a few, and compared them against regular Transformers, discussed the ideas behind them and found out about their strengths and weaknesses.

Finally, we learned about implementing Transformers to help us with various real life problems and talked about the future of the field in genereal and how transformers can change everything!



Shahid Bahonar University of Kerman

Faculty of Engineering

Department of Computer Engineering

B.Sc. Thesis

Introduction to Transformers

Supervisor : Dr M.Eftekhari

By : Gholamreza Dar and Amirhasan Amirmahani

Date : 8/29/2021