

BERT: Bidirectional Encoder Representations from Transformers for Understanding Language

Toutanova Kristina Lee Kenton Chang Ming-Wei Devlin Jacob
Language AI Google
{jacobdevlin, mingweichang, kentonl, kristout}@google.com

Abstract

ما یک مدل جدید برای نمایش زبانی معرفی می‌کنیم با نام **BERT**، که مخفف Bidirectional Encoder Representations from Transformers است. برخلاف مدل‌های اخیر نمایش زبان (LSTM)، BERT به گونه‌ای طراحی شده است که نمایش‌های دوجهته عمیقی را از متون بدون برچسب، با شرط بندی هم‌زمان بر بافت‌های سمت چپ و راست در تمام لایه‌ها، پیش‌آموزش دهد. در نتیجه، مدل از پیش‌آموزش داده شده BERT می‌تواند تنها با افزودن یک لایه خروجی اضافی برای طیف گسترده‌ای از وظایف، نظیر پاسخ‌گویی به پرسش و استنتاج زبانی، به مدل‌های پیشرفته‌ی روز تبدیل شود، بی‌آنکه نیاز به تغییرات عمده در معماری خاص هر وظیفه باشد.

BERT از نظر مفهومی ساده و از نظر تجربی قدرتمند است. این مدل در یازده وظیفه پردازش زبان طبیعی (NLP) نتایج جدیدی در سطح بهترین روش‌های موجود به دست آورده است، از جمله ارتقای امتیاز GLUE به ۸۰/۵٪ (افزایش مطلق ۷/۷ واحد درصد)، افزایش دقت MultiNLI به ۸۶/۷٪ (افزایش مطلق ۴/۶ واحد درصد)، افزایش امتیاز SQuAD v1.1 در معیار Test F1 به ۹۳/۲٪ (افزایش مطلق ۱/۵ واحد درصد) و همچنین افزایش امتیاز SQuAD v2.0 Test F1 به ۸۳/۱٪ (افزایش مطلق ۵/۱ واحد درصد).

Introduction ۱

پیش‌آموزش مدل زبان برای بهبود بسیاری از وظایف پردازش زبان طبیعی (LSTM) مؤثر نشان داده شده است. این وظایف شامل وظایف سطح جمله مانند استنتاج زبان طبیعی (LSTM) و بازنویسی (LSTM) است که هدف آنها پیش‌بینی روابط بین جملات با تجزیه و تحلیل جامع آنها است، و همچنین وظایف سطح توکن مانند تشخیص موجودیت‌های نامگذاری شده و پاسخ به سوال، که در آنها مدل‌ها برای تولید خروجی دقیق در سطح توکن مورد نیاز هستند (LSTM).

دو راهکار موجود برای به کارگیری نمایش‌های زبانی از پیش‌آموزش دیده در وظایف پایین‌دستی وجود دارد: مبتنی بر ویژگی و تنظیم دقیق (fine-tuning). رویکرد مبتنی بر ویژگی، مانند مدل ELMo (LSTM)، از معماری‌های خاص هر وظیفه استفاده می‌کند که نمایش‌های از پیش‌آموزش دیده را به عنوان ویژگی‌های اضافی در نظر می‌گیرند. در مقابل،

رویکرد تنظیم دقیق، مانند مدل Generative Pre-trained Transformer (OpenAI GPT) (LSTM)، تنها تعداد کمی پارامتر وابسته به وظیفه را معرفی می‌کند و با تنظیم تمام پارامترهای از پیش‌آموزش دیده، روی وظایف پایین‌دستی آموزش داده می‌شود. هر دو رویکرد در مرحله‌ی پیش‌آموزش از تابع هدف یکسانی استفاده می‌کنند، جایی که مدل‌های زبانی یک‌جهته برای یادگیری نمایش‌های کلی زبان به کار می‌روند. ما استدلال می‌کنیم که روش‌های فعلی توانایی نمایش‌های از پیش‌آموزش دیده را به ویژه در رویکردهای تنظیم دقیق محدود می‌کنند. محدودیت اصلی آن است که مدل‌های زبانی استاندارد یک‌جهته هستند و این امر دامنه‌ی معماری‌هایی را که می‌توان در مرحله‌ی پیش‌آموزش به کار برد، کاهش می‌دهد. برای نمونه، در مدل OpenAI GPT، نویسندگان از معماری چپ‌به‌راست استفاده کرده‌اند که در آن هر توکن تنها می‌تواند به توکن‌های پیشین خود در لایه‌های خودتوجهی Transformer (LSTM) توجه کند. چنین محدودیتی برای وظایف در سطح جمله بهینه نیست و می‌تواند هنگام استفاده از رویکردهای مبتنی بر تنظیم دقیق در وظایفی مانند پاسخ‌گویی به پرسش، که در آن‌ها ترکیب بافت از هر دو جهت ضروری است، اثر منفی شدیدی داشته باشد.

در این مقاله، ما با معرفی BERT، روش‌های مبتنی بر تنظیم دقیق را بهبود می‌دهیم. BERT مخفف Bidirectional Encoder Representations from Transformers است، محدودیت ذکر شده‌ی یک‌جهته بودن را با بهره‌گیری از هدف پیش‌آموزش به نام masked language model (MLM) رفع می‌کند؛ هدفی که از آزمون Cloze (LSTM) الهام گرفته است. در مدل زبانی پوشیده، برخی از توکن‌های ورودی به صورت تصادفی ماسک می‌شوند و هدف، پیش‌بینی شناسه‌ی واژگانی اصلی توکن ماسک شده تنها بر اساس بافت آن است. برخلاف پیش‌آموزش مدل‌های زبانی چپ‌به‌راست، هدف MLM امکان ترکیب بافت چپ و راست را فراهم می‌کند، که در نتیجه می‌توان یک Transformer دوجهته (دوجهته) عمیق را پیش‌آموزش داد. علاوه بر مدل زبانی پوشیده، ما از وظیفه‌ی دیگری به نام next sentence prediction نیز استفاده می‌کنیم که به صورت هم‌زمان نمایش‌های جفت جمله را پیش‌آموزش می‌دهد.

● ما اهمیت پیش‌آموزش دوجهته را برای نمایش‌های زبانی نشان می‌دهیم. برخلاف ؟ که از مدل‌های زبانی

یک‌جهته برای پیش‌آموزش استفاده می‌کند، BERT با بهره‌گیری از مدل‌های زبانی پوشیده (Masked Lan- guage Models) امکان ایجاد نمایش‌های دوجهته و عمیق از پیش‌آموزش‌دیده را فراهم می‌سازد. این رویکرد همچنین در تضاد با روش ؟ است که از ترکیب سطحی دو مدل زبانی چپ‌به‌راست و راست‌به‌چپ آموزش‌داده‌شده به‌صورت مستقل استفاده می‌کند.

- ما نشان می‌دهیم که نمایش‌های ازپیش‌آموزش‌دیده نیاز به بسیاری از معماری‌های پیچیده و خاص هر وظیفه را کاهش می‌دهند. BERT نخستین مدل نمایش مبتنی بر تنظیم دقیق است که در مجموعه‌ی گسترده‌ای از وظایف در سطح جمله و سطح توکن به عملکرد در حد بهترین روش‌های روز دست می‌یابد و از بسیاری از معماری‌های خاص وظیفه پیشی می‌گیرد.

- BERT در یازده وظیفه‌ی پردازش زبان طبیعی (NLP) پیشرفت قابل‌توجهی نسبت به روش‌های پیشین ارائه می‌دهد. کد منبع و مدل‌های ازپیش‌آموزش‌داده‌شده در <https://github.com/google-research/bert> در دسترس هستند.

۲ Work Related

تاریخچه‌ای طولانی در زمینه‌ی پیش‌آموزش نمایش‌های زبانی عمومی وجود دارد. در این بخش، به‌طور خلاصه به بررسی رایج‌ترین رویکردهای مورد استفاده در این حوزه می‌پردازیم.

۱.۲ رویکردهای بدون نظارت مبتنی بر ویژگی

یادگیری نمایش‌هایی از واژه‌ها که بتوانند در دامنه‌های گوناگون به‌کار گرفته شوند، طی دهه‌ها یکی از حوزه‌های فعال پژوهش بوده است. این تلاش‌ها شامل روش‌های غیرعصبی (؟؟؟) و همچنین روش‌های مبتنی بر شبکه‌های عصبی (؟؟) می‌شوند. تعبیه‌های واژه‌ای ازپیش‌آموزش‌دیده، بخش جدایی‌ناپذیر سامانه‌های پردازش زبان طبیعی مدرن محسوب می‌شوند و نسبت به تعبیه‌هایی که از ابتدا آموزش داده می‌شوند، بهبود قابل‌توجهی ارائه می‌دهند (؟؟). برای پیش‌آموزش بردارهای تعبیه‌ی واژه، اهداف مدل‌سازی زبان از چپ به راست (؟؟) و نیز اهدافی که به تمایز واژه‌های درست از نادرست در بافت چپ و راست کمک می‌کنند (؟؟) به‌کار رفته‌اند.

این رویکردها به سطوح درشت‌دانه‌تری نیز تعمیم یافته‌اند، از جمله تعبیه‌های جمله‌ای (؟؟) یا تعبیه‌های بند (؟؟). برای آموزش نمایش‌های جمله، پژوهش‌های پیشین از اهدافی مانند رتبه‌بندی جمله‌های بعدی کاندید (؟؟)، تولید واژه‌های جمله‌ی بعدی از چپ به راست با توجه به نمایش جمله‌ی پیشین (؟؟)، یا اهداف برگرفته از خودرمزگذارهای حذف نویز (؟؟) بهره گرفته‌اند.

مدل ELMo و نسخه‌های پیشین آن (؟؟) پژوهش‌های سنتی تعبیه‌ی واژه را در بُعدی متفاوت گسترش می‌دهند. این مدل‌ها ویژگی‌های وابسته به بافت را از دو مدل زبانی چپ‌به‌راست و راست‌به‌چپ استخراج می‌کنند. نمایش

بافتی هر توکن از به‌هم‌پیوستن نمایش‌های چپ‌به‌راست و راست‌به‌چپ به‌دست می‌آید. هنگامی که تعبیه‌های بافتی واژه در معماری‌های خاص وظیفه ادغام می‌شوند، ELMo عملکرد به‌روز و پیشرفته‌ای را در چندین معیار مهم NLP از جمله پاسخ‌گویی به پرسش‌ها (؟؟)، تحلیل احساسات (؟؟)، و شناسایی نام موجودیت‌ها (؟؟) ارائه می‌دهد (؟؟).

؟ روشی را برای یادگیری نمایش‌های بافتی پیشنهاد کردند که در آن مدل، یک واژه را بر اساس بافت چپ و راست آن با استفاده از LSTM پیش‌بینی می‌کند. مشابه ELMo، مدل آن‌ها نیز مبتنی بر ویژگی بوده و به‌طور عمیق دوسویه نیست. همچنین، ؟ نشان دادند که وظیفه‌ی cloze (پیش‌بینی واژه‌های حذف‌شده) می‌تواند موجب افزایش پایداری مدل‌های تولید متن شود.

۲.۲ رویکردهای تنظیم دقیق بدون ناظر

مشابه رویکردهای مبتنی بر ویژگی، نخستین پژوهش‌ها در این زمینه تنها پارامترهای تعبیه‌ی واژه را از متون بدون برچسب پیش‌آموزش می‌دادند (؟؟).

در سال‌های اخیر، رمزگذارهای جمله یا سند که نمایش‌های بافتی از توکن‌ها تولید می‌کنند، از داده‌های بدون برچسب پیش‌آموزش یافته و سپس برای وظایف پایین‌دستی نظارت‌شده تنظیم دقیق شده‌اند (؟؟؟). مزیت اصلی این رویکردها آن است که تنها تعداد اندکی از پارامترها باید از ابتدا یاد گرفته شوند. تا حدی به دلیل همین مزیت، مدل OpenAI GPT (؟؟) توانست در بسیاری از وظایف سطح جمله از مجموعه‌داده‌ی معیار GLUE (؟؟) به نتایج برتر پیشین دست یابد. برای پیش‌آموزش این مدل‌ها، از اهداف مدل‌سازی زبانی چپ‌به‌راست و خودرمزگذار (auto-encoder) استفاده شده است (؟؟؟).

۳.۲ یادگیری انتقالی از داده‌های نظارت‌شده

پژوهش‌هایی نیز انجام شده که نشان می‌دهند انتقال دانش از وظایف نظارت‌شده با مجموعه‌داده‌های بزرگ — مانند استنتاج زبان طبیعی (؟؟) و ترجمه‌ی ماشینی (؟؟) — می‌تواند مؤثر باشد. در حوزه‌ی بینایی رایانه‌ای نیز پژوهش‌ها اهمیت یادگیری انتقالی از مدل‌های بزرگ ازپیش‌آموزش‌دیده را نشان داده‌اند، به‌طوری‌که یکی از راهکارهای مؤثر، تنظیم دقیق مدل‌هایی است که با مجموعه‌داده‌ی ImageNet (؟؟؟) پیش‌آموزش یافته‌اند.

۳ مدل BERT

در این بخش، مدل BERT و جزئیات پیاده‌سازی آن را معرفی می‌کنیم. چارچوب ما شامل دو مرحله است: پیش‌آموزش و تنظیم دقیق. در مرحله‌ی پیش‌آموزش، مدل بر روی داده‌های بدون برچسب و در قالب چندین وظیفه‌ی پیش‌آموزش مختلف آموزش می‌بیند. در مرحله‌ی تنظیم دقیق، مدل BERT ابتدا با پارامترهای ازپیش‌آموزش‌یافته مقداردهی اولیه می‌شود و سپس تمامی پارامترها با استفاده از داده‌های برچسب‌دار وظیفه‌ی پایین‌دستی تنظیم دقیق می‌شوند. هر وظیفه‌ی پایین‌دستی، مدل تنظیم‌شده‌ی مخصوص به خود را دارد، هرچند همه‌ی

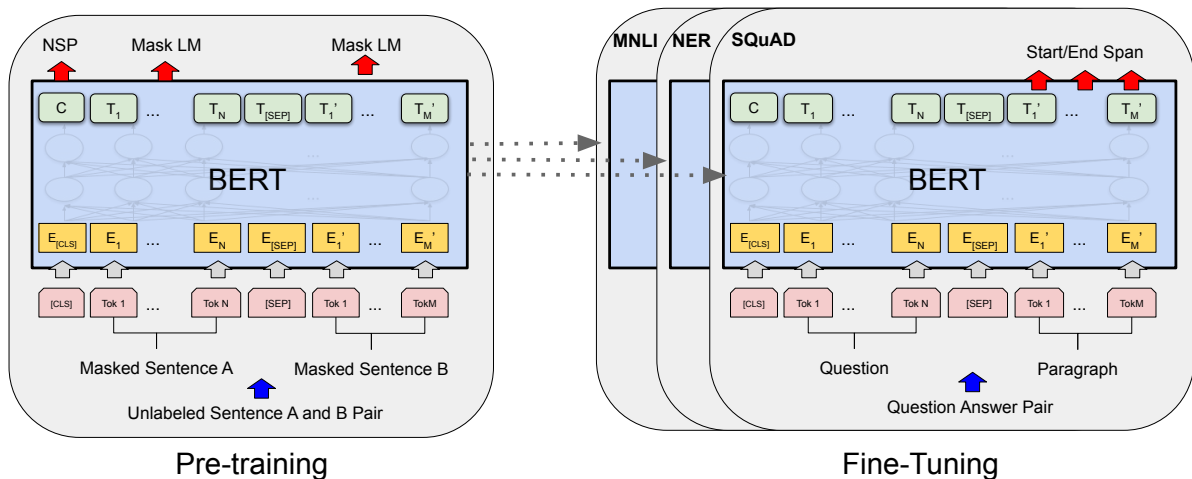


Figure ۱: فرآیند کلی پیش‌آموزش و تنظیم دقیق (fine-tuning) در مدل BERT. به‌جز لایه‌های خروجی، در هر دو مرحله‌ی پیش‌آموزش و تنظیم دقیق از معماری‌های یکسانی استفاده می‌شود. پارامترهای از پیش‌آموزش دیده برای مقداردهی اولیه‌ی مدل‌ها در وظایف پایین‌دستی مختلف به‌کار می‌روند. در طی تنظیم دقیق، تمام پارامترها به‌روزرسانی می‌شوند. نماد [CLS] یک نشانگر ویژه است که در ابتدای هر نمونه ورودی افزوده می‌شود، و [SEP] یک توکن جداساز ویژه است (برای مثال، جهت جداسازی پرسش‌ها و پاسخ‌ها).

آن‌ها با پارامترهای پیش‌آموزش یافته‌ی یکسان مقداردهی اولیه می‌شوند. مثال پرسش و پاسخ در شکل ۱ به عنوان نمونه‌ی جاری در این بخش به کار گرفته می‌شود.

یکی از ویژگی‌های متمایز BERT، معماری یکپارچه‌ی آن در میان وظایف مختلف است. تفاوت میان معماری پیش‌آموزش و معماری نهایی وظایف پایین‌دستی بسیار اندک است.

معماری مدل معماری مدل BERT یک رمزگذار Transformer دوسویه‌ی چندلایه است که بر اساس پیاده‌سازی اصلی معرفی‌شده در [1] ساخته شده و در کتابخانه‌ی `tensorflow/tensor2tensor`^۱ منتشر شده است. از آن‌جا که استفاده از Transformer ها امروزه رایج شده و پیاده‌سازی ما تقریباً مشابه نسخه‌ی اصلی است، از ارائه‌ی توضیحات پیش‌زمینه‌ای مفصل خودداری می‌کنیم و خوانندگان را به منبع اصلی [2] «The Annotated Transformer»^۲ ارجاع می‌دهیم.

در این پژوهش، تعداد لایه‌ها (یعنی بلوک‌های Transformer) را با L ، اندازه‌ی نهان را با H و تعداد سرهای خودتوجهی را با A نمایش می‌دهیم.^۳ نتایج اصلی ما بر روی دو اندازه‌ی مدل گزارش می‌شوند: **BASE BERT** (با $A=12$, $H=768$, $L=12$) و در مجموع ۱۱۰ میلیون پارامتر) و **LARGE BERT** (با $A=16$, $H=1024$, $L=24$) و در مجموع ۳۴۰ میلیون پارامتر).^۴

مدل **BASE BERT** به‌منظور مقایسه، با همان اندازه‌ی مدل

نمایش‌های ورودی/خروجی برای اینکه BERT بتواند طیف گسترده‌ای از وظایف پایین‌دستی را انجام دهد، نمایش ورودی ما قادر است به‌طور غیرمبهم هم یک جمله‌ی منفرد و هم یک جفت جمله (مثلاً پرسش، پاسخ) را در یک توالی توکن نمایش دهد. در سرتاسر این پژوهش، «جمله» می‌تواند هر بازه‌ی دلخواهی از متن متوالی باشد و لزوماً یک جمله‌ی زبان‌شناختی واقعی نیست. «توالی» به توالی ورودی توکن‌ها به BERT گفته می‌شود، که می‌تواند یک جمله یا دو جمله‌ی بسته‌بندی‌شده کنار هم باشد.

ما از WordPiece embeddings (۴) با یک واژه‌نامه ۳۰,۰۰۰ توکنی استفاده می‌کنیم. توکن اول هر توالی همیشه یک توکن ویژه‌ی طبقه‌بندی ([CLS]) است. وضعیت نهان نهایی متناظر با این توکن به عنوان نمایش تجمیعی توالی برای وظایف طبقه‌بندی به کار می‌رود. جفت جملات در یک توالی واحد بسته‌بندی می‌شوند. ما جملات را به دو روش متمایز می‌کنیم. اول، آن‌ها را با یک توکن ویژه ([SEP]) جدا می‌کنیم. دوم، به هر توکن یک نمایش یادگرفته‌شده اضافه می‌کنیم که مشخص می‌کند آیا توکن متعلق به جمله‌ی A است یا جمله‌ی B. همان‌طور که در شکل ۱ نشان داده شده، نمایش ورودی را با E ، بردار نهان نهایی توکن ویژه [CLS] را با $C \in \mathbb{R}^H$ و بردار نهان نهایی توکن ورودی i^{th} را با $T_i \in \mathbb{R}^H$ نمایش می‌دهیم.

^۴ در ادبیات پژوهشی، معمولاً به Transformer دوسویه «رمزگذار Transformer» و به نسخه‌ی محدود به بافت چپ «رمزگشای Transformer» گفته می‌شود، چرا که نسخه‌ی دوم برای تولید متن به کار می‌رود.

^۱ <https://github.com/tensorflow/tensor2tensor>
^۲ <http://nlp.seas.harvard.edu/2018/04/03/attention.html>
^۳ در تمامی موارد، اندازه‌ی لایه‌ی پیش‌خور یا فیلتر را برابر $4H$ در نظر می‌گیریم؛ یعنی ۳۰۷۲ برای $H=768$ و ۴۰۹۶ برای $H=1024$.

برای هر توکن، نمایش ورودی آن با جمع کردن نمایش‌های توکن، بخش و موقعیت متناظر ساخته می‌شود. یک تصویرسازی از این ساختار در شکل ۲ دیده می‌شود.

۱.۳ پیش‌آموزش BERT

بر خلاف ؟ و ؟، ما از مدل‌های زبان سنتی چپ-به-راست یا راست-به-چپ برای پیش‌آموزش BERT استفاده نمی‌کنیم. در عوض، BERT را با استفاده از دو وظیفه بدون نظارت پیش‌آموزش می‌کنیم که در این بخش شرح داده شده‌اند. این مرحله در بخش چپ شکل ۱ نشان داده شده است.

وظیفه #۱: مدل زبان ماسک‌شده (Masked LM) به‌طور شهودی، منطقی است که باور کنیم یک مدل عمیق دوجتهی به‌طور قابل توجهی قدرتمندتر از یک مدل چپ-به-راست یا ترکیب سطحی مدل‌های چپ-به-راست و راست-به-چپ است. متأسفانه، مدل‌های زبان شرطی استاندارد تنها می‌توانند چپ-به-راست یا راست-به-چپ آموزش ببینند، زیرا شرطی‌سازی دوجتهی باعث می‌شود هر واژه به‌طور غیرمستقیم «خودش را ببیند» و مدل به‌سادگی بتواند واژه هدف را در یک زمینه‌ی چندلایه‌ای پیش‌بینی کند.

برای آموزش یک نمایش عمیق دوجتهی، ما به سادگی درصدی از توکن‌های ورودی را به‌طور تصادفی ماسک می‌کنیم و سپس آن توکن‌های ماسک‌شده را پیش‌بینی می‌کنیم. به این فرآیند «مدل زبان ماسک‌شده» (MLM) می‌گوییم، هرچند در منابع ادبیات غالباً به آن وظیفه‌ی Cloze گفته می‌شود (۴). در این حالت، بردارهای نهان‌هایی متناظر با توکن‌های ماسک‌شده به یک خروجی softmax روی واژگان داده می‌شوند، مشابه یک مدل زبان استاندارد. در تمام آزمایش‌های ما، ۱۵٪ از تمام توکن‌های WordPiece هر توالی به‌طور تصادفی ماسک می‌شوند. بر خلاف denoising auto-encoders (۴)، ما تنها واژه‌های ماسک‌شده را پیش‌بینی می‌کنیم و کل ورودی را بازسازی نمی‌کنیم.

اگرچه این کار به ما امکان ایجاد یک مدل پیش‌آموزش دوجتهی را می‌دهد، اما یک نکته منفی دارد: ایجاد ناسازگاری بین پیش‌آموزش و ریزتنظیم، زیرا توکن [MASK] در هنگام ریزتنظیم ظاهر نمی‌شود. برای کاهش این مشکل، ما همیشه واژه‌های «ماسک‌شده» را با توکن واقعی [MASK] جایگزین نمی‌کنیم. تولیدکننده داده‌های آموزشی به‌طور تصادفی ۱۵٪ از موقعیت‌های توکن را برای پیش‌بینی انتخاب می‌کند. اگر توکن i -ام انتخاب شد، ما توکن i -ام را به یکی از موارد زیر جایگزین می‌کنیم: (۱) توکن [MASK] در ۸۰٪ مواقع، (۲) یک توکن تصادفی در ۱۰٪ مواقع، (۳) توکن i -ام بدون تغییر در ۱۰٪ مواقع. سپس، T_i برای پیش‌بینی توکن اصلی با استفاده از cross entropy loss به کار می‌رود. ما انواع مختلف این روش را در ضمیمه ۲.۲ مقایسه کرده‌ایم.

وظیفه #۲: پیش‌بینی جمله بعدی (NSP) بسیاری از وظایف پایین‌دستی مهم مانند پرسش و پاسخ (QA) و استنتاج زبان طبیعی (NLI) بر اساس فهم رابطه بین دو جمله هستند که به‌طور مستقیم توسط مدل‌سازی زبان ثبت نمی‌شود.

برای آموزش مدلی که روابط بین جملات را درک کند، ما پیش‌آموزش را برای وظیفه‌ی دودویی «پیش‌بینی جمله بعدی» انجام می‌دهیم که می‌توان آن را به‌سادگی از هر مجموعه‌ی متنی تک‌زبان ایجاد کرد. به‌طور مشخص، هنگام انتخاب جملات A و B برای هر نمونه‌ی پیش‌آموزش، ۵۰٪ مواقع B جمله واقعی بعد از A است (برچسب‌گذاری شده به‌عنوان IsNext) و ۵۰٪ مواقع جمله‌ای تصادفی از مجموعه انتخاب می‌شود (برچسب‌گذاری شده به‌عنوان NotNext). همان‌طور که در شکل ۱ نشان داده شده، بردار C برای پیش‌بینی جمله بعدی (NSP) به کار می‌رود. با وجود سادگی، نشان می‌دهیم در بخش ۱.۵ که پیش‌آموزش برای این وظیفه برای QA و NLI بسیار مفید است. ^۶ وظیفه NSP ارتباط نزدیکی با اهداف یادگیری نمایش در ؟ و ؟ دارد. با این حال، در کارهای پیشین، تنها نمایش‌های جمله به وظایف پایین‌دستی منتقل می‌شوند، در حالی که BERT تمام پارامترها را برای مقداردهی اولیه‌ی مدل وظیفه انتهایی انتقال می‌دهد.

داده‌های پیش‌آموزش روند پیش‌آموزش عمدتاً از ادبیات موجود درباره پیش‌آموزش مدل زبان پیروی می‌کند. برای مجموعه داده‌های پیش‌آموزش از BooksCorpus (۸۰۰ میلیون واژه) (۴) و ویکی‌پدیای انگلیسی (۲,۵۰۰ میلیون واژه) استفاده می‌کنیم. برای ویکی‌پدیا تنها بخش‌های متنی استخراج شده و لیست‌ها، جداول و سرفصل‌ها نادیده گرفته می‌شوند. استفاده از یک مجموعه داده سطح سند به جای مجموعه داده سطح جمله‌ی مرتب‌شده مانند Billion Word Benchmark (۴) برای استخراج توالی‌های طولانی متوالی حیاتی است.

۲.۳ BERT Fine-tuning

Fine-tuning ساده است، زیرا مکانیزم خود-توجه (self-attention) در معماری Transformer به BERT امکان می‌دهد تا بسیاری از وظایف پایین‌دستی—چه شامل یک متن باشند و چه جفت متن—را با جایگزینی ورودی‌ها و خروجی‌های مناسب مدل کند. برای کاربردهایی که شامل جفت متن هستند، یک الگوی رایج این است که ابتدا جفت متن‌ها را به‌صورت مستقل کدگذاری کرده و سپس از توجه متقابل دوطرفه cross (bidirectional) attention استفاده کنند، مانند ؟؟. در مقابل، BERT از مکانیزم خود-توجه برای یکپارچه‌سازی این دو مرحله استفاده می‌کند: زیرا کدگذاری یک جفت متن الحاق‌شده (concatenated) با خود-توجه، به‌طور مؤثری شامل توجه متقابل دوطرفه بین دو جمله می‌شود.

برای هر وظیفه، ما به‌سادگی ورودی‌ها و خروجی‌های مربوط به آن وظیفه را به BERT متصل کرده و تمام پارامترها را به‌صورت fine-tune end-to-end می‌کنیم. در ورودی، جمله A و جمله B از مرحله پیش‌آموزش (pre-training) معادل موارد زیر هستند: (۱) جفت جملات در بازنویسی معنایی، (paraphrasing) (۲) جفت فرضیه—

^۵ مدل نهایی به دقت ۹۷٪-۹۸٪ در NSP دست می‌یابد.
^۶ بردار C بدون ریزتنظیم، نمایش جمله معنی‌داری نیست، زیرا با NSP آموزش داده شده است.

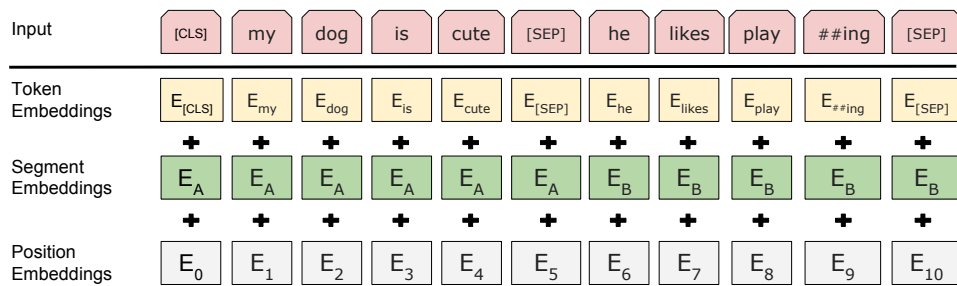


Figure ۲: نمایش ورودی BERT. نمایش های توکن، نمایش های بخش بندی و نمایش های موقعیت هستند.

برای همه وظایف، GLUE از اندازه بسته ۳۲ استفاده کرده و مدل را به مدت ۳ اپوک روی داده ها ریزتنظیم می کنیم. برای هر وظیفه، بهترین نرخ یادگیری ریزتنظیم را (از میان $5e-5$ ، $3e-5$ ، $4e-5$ و $2e-5$) روی مجموعه Dev انتخاب کردیم. علاوه بر این، برای LARGE BERT متوجه شدیم که ریزتنظیم گاهی روی داده های کوچک ناپایدار است، بنابراین چندین راه اندازی تصادفی انجام دادیم و بهترین مدل را روی مجموعه Dev انتخاب کردیم. در راه اندازی های تصادفی، از همان نقطه بررسی پیش آموزش شده استفاده می کنیم، اما جابه جایی داده های ریزتنظیم و مقداردهی اولیه لایه دسته بندی متفاوت انجام می دهیم.^۹

نتایج در جدول ۱ ارائه شده است. هم BASE BERT و هم LARGE BERT تمام سیستم ها را در تمام وظایف با فاصله قابل توجهی پشت سر گذاشته اند و به ترتیب ۵٪ و ۰.۷٪ بهبود دقت متوسط نسبت به وضعیت قبلی ارائه می دهند. توجه داشته باشید که BASE BERT و GPT OpenAI از نظر معماری مدل تقریباً یکسان هستند، به جز ماسک گذاری توجه. برای بزرگترین و پرکارترین وظیفه، GLUE، یعنی BERT MNLI، به بهبود ۶.۴٪ در دقت مطلق دست یافته است. در جدول رده بندی رسمی GLUE^{۱۰}، LARGE BERT امتیاز ۵.۸۰ را کسب کرده است، در حالی که GPT OpenAI در زمان نگارش این مقاله امتیاز ۸.۷۲ را دارد.

۲.۴ SQuAD v۱.۱

مجموعه داده Answering Question Stanford (SQuAD Dataset) (v۱.۱) شامل ۱۰۰ هزار جفت پرسش/پاسخ جمع آوری شده از طریق جمع سپاری است (؟). در این وظیفه، با داشتن یک پرسش و یک متن از ویکی پدیا که شامل پاسخ است، هدف پیش بینی بازه متنی پاسخ در متن داده شده می باشد.

همان طور که در شکل ۱ نشان داده شده است، در وظیفه پاسخ به پرسش، سوال و پاراگراف ورودی را به صورت یک توالی بسته بندی شده نمایش می دهیم، به طوری که سوال از A embedding و پاراگراف از B embedding استفاده

پیش فرض (hypothesis-premise) در استنتاج (entail-ment)، (۳) جفت پرسش-متن (question-passage) در پاسخ دهی به پرسش (question-answering)، (۴) یک جفت متن-تبهگون (degenerate) در طبقه بندی متن یا برچسب گذاری دنباله (sequence tagging). خروجی، نمایش های توکن ها به یک لایه خروجی برای وظایف سطح-توکن، (token-level) مانند برچسب گذاری دنباله یا پاسخ دهی به پرسش، ارسال می شوند؛ و نمایش [CLS] به یک لایه خروجی برای طبقه بندی، مانند استنتاج یا تحلیل احساسات (sentiment analysis)، فرستاده می شود.

در مقایسه با پیش آموزش، fine-tuning هزینه نسبتاً کمی دارد. تمام نتایج ارائه شده در این مقاله را می توان حداکثر در یک ساعت با یک واحد پردازشی ابری (Cloud TPU) یا چند ساعت با یک GPU، با شروع از همان مدل پیش آموزش دیده، بازتولید (replicate) کرد.^۷ جزئیات مربوط به هر وظیفه را در زیربخش های مربوطه در بخش ۴ توصیف می کنیم. برای جزئیات بیشتر می توانید به پیوست ۵.۱ مراجعه کنید.

۴ experiments

در این بخش، نتایج fine-tuning مدل BERT را روی ۱۱ وظیفه پردازش زبان طبیعی ارائه می دهیم.

۱.۴ GLUE

معیار ارزیابی عمومی درک زبان (GLUE) (؟) مجموعه ای از وظایف متنوع درک زبان طبیعی است. توضیحات دقیق تر مجموعه داده های GLUE در پیوست ۱.۱ آورده شده است. برای ریزتنظیم روی GLUE، دنباله ورودی (برای یک جمله یا جفت جمله) را همان طور که در بخش ۳ توضیح داده شد، نمایش می دهیم و از بردار پنهان نهایی $C \in \mathbb{R}^H$ مربوط به اولین توکن ورودی ([CLS]) به عنوان نمایش کلی استفاده می کنیم. تنها پارامترهای جدیدی که در طول ریزتنظیم معرفی می شوند، وزن های لایه دسته بندی $W \in \mathbb{R}^{K \times H}$ هستند، که در آن K تعداد برچسب ها است. ما از بردار C و وزن W برای محاسبه یک خطای دسته بندی استاندارد استفاده می کنیم، یعنی $\log(\text{softmax}(CW^T))$.

^۹ توزیع مجموعه داده های GLUE شامل برچسب های Test نمی شود و ما تنها یک بار ارسال ارزیابی به سرور GLUE برای هر یک از BASE BERT و LARGE BERT انجام دادیم.

^{۱۰} <https://gluebenchmark.com/leaderboard>

^۷ به عنوان مثال، مدل BERT برای مجموعه SQuAD را می توان در حدود ۳۰ دقیقه با یک TPU Cloud آموزش داد تا به نمره F۱ برابر با ۹۱٪ در مجموعه توسعه (Dev) دست یابد.

^۸ See (۱۰) <https://gluebenchmark.com/faq>

Average	RTE	MRPC	STS-B	CoLA	SST-2	QNLI	QQP	MNLI-(m/mm)	System
-	5k.2	5k.3	7k.5	5k.8	6vk	10Ak	363k	392k	
0.74	7.61	0.86	0.81	0.35	2.93	3.82	1.66	1.6/80.80	SOTA Pre-OpenAI
0.71	8.56	9.84	3.73	0.36	4.90	8.79	8.64	1.4/76.76	BiLSTM+ELMo+Attn
1.75	0.56	3.82	0.80	4.45	3.91	4.87	3.70	4.1/81.82	GPT OpenAI
6.79	4.66	9.88	8.85	1.52	5.93	5.90	2.71	4.6/83.84	BASE BERT
1.82	1.70	3.89	5.86	5.60	9.94	7.92	1.72	9.7/85.86	LARGE BERT

The (https://gluebenchmark.com/leaderboard) server evaluation the by scored results. Test GLUE : 1 Table different slightly is column "Average" The examples. training of number the denotes task each below number are GPT OpenAI and BERT ^set. WNLI problematic the exclude we since score. GLUE official the than reported are correlations Spearman MRPC. and QQP for reported are scores F1 task. single single-model. of one as BERT use that entries exclude We tasks. other the for reported are scores accuracy and STS-B. for components. their

۳.۴ SQuAD ۲.۰

وظیفه SQuAD ۲.۰ تعریف مسئله SQuAD ۱.۱ را گسترش می‌دهد و اجازه می‌دهد که هیچ پاسخ کوتاهی در پاراگراف ارائه‌شده وجود نداشته باشد، که مسئله را واقع‌گرایانه‌تر می‌کند.

ما از رویکرد ساده‌ای برای گسترش مدل BERT نسخه SQuAD ۱.۱ برای این وظیفه استفاده می‌کنیم. سوال‌هایی که پاسخی ندارند، به عنوان سوال‌هایی در نظر گرفته می‌شوند که بازه پاسخ آن‌ها شروع و پایان آن روی توکن [CLS] است. فضای احتمالات برای موقعیت‌های شروع و پایان بازه پاسخ گسترش می‌یابد تا شامل موقعیت توکن [CLS] شود. برای پیش‌بینی، امتیاز بازه بدون پاسخ را با

$$s_{null} = S \cdot C + E \cdot C$$

با امتیاز بهترین بازه غیر تهی

$$s_{i,j} = \max_{j \geq i} S \cdot T_i + E \cdot T_j$$

مقایسه می‌کنیم. یک پاسخ غیر تهی پیش‌بینی می‌شود اگر

$$s_{i,j} > s_{null} + \tau$$

که آستانه τ روی مجموعه Dev انتخاب می‌شود تا F1 به حداکثر برسد. ما برای این مدل از داده‌های TriviaQA استفاده نکردیم. Fine-tuning برای ۲ epoch با نرخ یادگیری 5e-5 و اندازه batch برابر با ۴۸ انجام شد. نتایج نسبت به ورودی‌های قبلی leaderboard و کارهای منتشرشده برتر (؟؟) در جدول ۳ نشان داده شده‌اند، با این تفاوت که سیستم‌هایی که BERT را به عنوان یکی از مؤلفه‌های خود استفاده کرده‌اند، حذف شده‌اند. ما بهبود ۱.۵+ F1 نسبت به سیستم برتر قبلی مشاهده می‌کنیم.

۴.۴ SWAG

داده مجموعه Genera- Adversarial With Situations (SWAG) شامل ۱۱۳ هزار مثال تکمیل جمله-جفتی

TriviaQA-Wiki است که از ۴۰۰ توکن اول اسناد تشکیل شده‌اند و حداقل یکی از پاسخ‌های ممکن ارائه‌شده در آن‌ها موجود است.

می‌کند. تنها در هنگام fine-tuning یک بردار شروع $S \in \mathbb{R}^H$ و یک بردار پایان $E \in \mathbb{R}^H$ معرفی می‌کنیم. احتمال اینکه واژه i شروع بازه پاسخ باشد، به صورت حاصلضرب داخلی بین T_i و S محاسبه شده و سپس یک softmax روی تمام واژه‌های پاراگراف اعمال می‌شود:

$$P_i = \frac{e^{S \cdot T_i}}{\sum_j e^{S \cdot T_j}}$$

فرمول مشابه برای انتهای بازه پاسخ استفاده می‌شود. امتیاز یک بازه پیشنهادی از موقعیت i تا موقعیت j به صورت

$$S \cdot T_i + E \cdot T_j$$

تعریف شده است و بازه‌ای با بیشترین امتیاز که $j \geq i$ باشد به عنوان پیش‌بینی انتخاب می‌شود. هدف آموزش، مجموع لگاریتم احتمال موقعیت‌های شروع و پایان صحیح است. ما برای ۳ epoch با نرخ یادگیری 5e-5 و اندازه batch برابر با ۳۲، fine-tuning انجام می‌دهیم.

جدول ۲ نتایج برتر در leaderboard و همچنین نتایج سیستم‌های منتشرشده را نشان می‌دهد (؟؟؟؟). نتایج برتر SQuAD leaderboard توصیف سیستم‌های عمومی به روز را ندارند^{۱۱} و می‌توانند از هر داده عمومی هنگام آموزش استفاده کنند. بنابراین ما در سیستم خود از افزایش داده متواضعانه استفاده می‌کنیم؛ ابتدا روی TriviaQA (؟) fine-tuning انجام می‌دهیم و سپس روی SQuAD ادامه می‌دهیم.

بهترین سیستم ما از سیستم برتر leaderboard با ۵.۱+ F1 در ترکیب مدل‌ها و ۳.۱+ F1 به عنوان یک سیستم واحد پیشی می‌گیرد. در واقع، مدل BERT تنها ما، سیستم ترکیبی برتر را از نظر Score F1 پشت سر می‌گذارد. بدون داده‌های fine-tuning TriviaQA، تنها ۱.۰-۴.۰ F1 کاهش داریم و هنوز هم از تمام سیستم‌های موجود به طور قابل توجهی پیشی می‌گیریم^{۱۲}.

^{۱۱} QANet در ؟ توصیف شده است، اما سیستم پس از انتشار به طور قابل توجهی بهبود یافته است.

^{۱۲} داده‌های TriviaQA که استفاده کردیم شامل پاراگراف‌هایی از

Test		Dev		(System) سامانه
F1	EM	F1	EM	
برترین سامانه‌های جدول امتیازات (۱۰ دسامبر ۲۰۱۸)				
۹۱/۲	۸۲/۳	-	-	انسان (Human)
۹۱/۷	۸۶/۰	-	-	#۱ تجمیعی (Ensemble) - nlnet
۹۰/۵	۸۴/۵	-	-	#۲ تجمیعی (Ensemble) - QANet
منتشر شده‌ها (Published)				
۸۵/۸	-	۸۵/۶	-	BiDAF+ELMo (تکی)
۸۸/۵	۸۲/۳	۸۷/۹	۸۱/۲	Reader R.M. (تجمیعی)
روش ما (Ours)				
-	-	۸۸/۵	۸۰/۸	BASE BERT (تکی)
-	-	۹۰/۹	۸۴/۱	LARGE BERT (تکی)
-	-	۹۱/۸	۸۵/۸	LARGE BERT (تجمیعی)
۹۱/۸	۸۵/۱	۹۱/۱	۸۴/۲	LARGE BERT (تکی) + TriviaQA
۹۳/۲	۸۷/۴	۹۲/۲	۸۶/۲	LARGE BERT (تجمیعی) + TriviaQA

Table ۲: نتایج مجموعه داده SQuAD 1.1. مدل تجمیعی BERT شامل ۷ سامانه است که از نقاط بررسی (Checkpoints) پیش‌تمرین متفاوت و بذره‌های تنظیم متفاوت در ریزتنظیم استفاده می‌کنند.

Test		Dev		(System) سامانه
F1	EM	F1	EM	
برترین سامانه‌های جدول امتیازات (۱۰ دسامبر ۲۰۱۸)				
۸۹/۵	۸۶/۹	۸۹/۰	۸۶/۳	انسان (Human)
۷۸/۰	۷۴/۸	-	-	#۱ تکی - MIR-MRC (F-Net)
۷۷/۱	۷۴/۲	-	-	#۲ تکی - nlnet
منتشر شده‌ها (Published)				
۷۴/۹	۷۱/۴	-	-	unet (تجمیعی)
۷۴/۴	۷۱/۴	-	-	SLQA+ (تکی)
روش ما (Ours)				
۸۳/۱	۸۰/۰	۸۱/۹	۷۸/۷	LARGE BERT (تکی)

Table ۳: نتایج مجموعه داده SQuAD 2.0. در این جدول، سامانه‌هایی که از BERT به عنوان یکی از اجزای خود استفاده کرده‌اند، لحاظ نشده‌اند.

Test	Dev	System
۷.۵۲	۹.۵۱	ESIM+GloVe
۲.۵۹	۱.۵۹	ESIM+ELMo
۰.۷۸	-	GPT OpenAI
-	۶.۸۱	BASE BERT
۳.۸۶	۶.۸۶	LARGE BERT
۰.۸۵	-	Human (expert)†
۰.۸۸	-	Human (۵ annotations)†

Table ۴: Human† accuracies. Test and Dev SWAG re- as samples. ۱۰۰ with measured is performance paper. SWAG the in ported

۱.۵ تأثیر وظایف Pre-training

ما اهمیت دو جهتی بودن عمیق BERT را با ارزیابی دو هدف pre-training نشان می‌دهیم، با استفاده از همان داده‌های pre-training، طرح fine-tuning و ابرپارامترها همانند BASE BERT:

بدون NSP: یک مدل دو جهتی که با استفاده از masked “next LM” (MLM) آموزش داده شده اما بدون وظیفه

است که استدلال مبتنی بر common-sense را ارزیابی می‌کند (؟). برای هر جمله، وظیفه این است که محتمل‌ترین ادامه را از میان چهار گزینه انتخاب کنیم.

هنگام fine-tuning روی داده مجموعه SWAG، چهار توالی ورودی ساخته می‌شود که هر کدام شامل الحاق جمله داده شده (جمله A) و یک ادامه ممکن (جمله B) هستند. تنها پارامتر خاص وظیفه که معرفی می‌شود یک بردار است که حاصل ضرب داخلی آن با نمایش توکن [CLS] یعنی C، امتیاز هر گزینه را نشان می‌دهد و با لایه softmax نرمال می‌شود. مدل برای ۳ epoch با نرخ یادگیری $5 - 2e$ و اندازه batch برابر با ۱۶ fine-tuning می‌شود. نتایج در جدول ۴ ارائه شده‌اند. BASE BERT سیستم baseline نویسندگان ESIM+ELMo را با ۲۷٪ و GPT OpenAI را با ۳٪.۸ پیشی می‌گیرد.

۵ مطالعات Ablation

در این بخش، آزمایش‌های ablation را روی جنبه‌های مختلف BERT انجام می‌دهیم تا اهمیت نسبی آن‌ها را بهتر درک کنیم. مطالعات ablation اضافی را می‌توان در ضمیمه پ یافت.

SQuAD (F1)	SST-2 (Acc)	Set Dev		QNLI (Acc)	MNLI-m (Acc)	Tasks
		MRPC (Acc)	QNLI (Acc)			
۵.۸۸	۷.۹۲	۷.۸۶	۴.۸۸	۴.۸۴		BASE BERT
۹.۸۷	۶.۹۲	۵.۸۶	۹.۸۴	۹.۸۳		NSP No
۸.۷۷	۱.۹۲	۵.۷۷	۳.۸۴	۱.۸۲		NSP No & LTR
۹.۸۴	۶.۹۱	۷.۷۵	۱.۸۴	۱.۸۲		BiLSTM +

Table ۵: Ablation study on pre-training tasks. The model trained is NSP. The architecture is BERT. The task is LTR. The prediction is sentence next the without LM left-to-right as trained is NSP. No & GPT. OpenAI like prediction, sentence next the out on BiLSTM initialized randomly a adds BiLSTM. The fine-tuning of top LTR the of top tuning.

sentence prediction (NSP) است.

LTR و بدون NSP: یک مدل فقط با زمینه چپ که با استفاده از LM استاندارد Left-to-Right (LTR) آموزش داده شده، به جای MLM محدودیت فقط چپ همچنین در fine-tuning اعمال شد، زیرا حذف آن موجب ایجاد ناسازگاری بین pre-train و fine-tune می‌شد که عملکرد در وظایف downstream را کاهش می‌داد. علاوه بر این، این مدل بدون وظیفه NSP پیش‌آموزش داده شد. این مستقیماً قابل مقایسه با GPT OpenAI است، اما با استفاده از مجموعه داده آموزشی بزرگ‌تر ما، نمایش ورودی ما و طرح fine-tuning ما.

ابتدا تأثیر وظیفه NSP را بررسی می‌کنیم. در جدول ۵ نشان می‌دهیم که حذف NSP عملکرد را به‌طور قابل توجهی در MNLI، QNLI و SQuAD ۱.۱ کاهش می‌دهد. سپس تأثیر آموزش نمایش‌های دو جهتی را با مقایسه “بدون NSP” با “LTR” و بدون “NSP” ارزیابی می‌کنیم. مدل LTR در تمام وظایف نسبت به مدل MLM عملکرد ضعیف‌تری دارد و کاهش‌های قابل توجهی در MRPC و SQuAD مشاهده می‌شود.

برای SQuAD، واضح است که مدل LTR در پیش‌بینی توکن‌ها عملکرد ضعیفی خواهد داشت، زیرا حالات پنهان سطح توکن هیچ زمینه‌ای از سمت راست ندارند. برای تلاش جدی در تقویت سیستم، LTR یک BiLSTM با مقداردهی اولیه تصادفی بر روی آن اضافه کردیم. این کار نتایج را در SQuAD به‌طور قابل توجهی بهبود می‌بخشد، اما نتایج همچنان بسیار ضعیف‌تر از مدل‌های دو جهتی پیش‌آموزش داده‌شده هستند. استفاده از BiLSTM عملکرد را در وظایف GLUE کاهش می‌دهد.

ما می‌دانیم که همچنین می‌توان مدل‌های جداگانه LTR و RTL آموزش داد و هر توکن را به‌صورت الحاق دو مدل نمایش داد، همانند کاری که ELMo انجام می‌دهد. با این حال: (a) این دو برابر هزینه یک مدل دو جهتی واحد است؛ (b) این برای وظایفی مانند QA غیرشهودی است، زیرا مدل RTL نمی‌تواند پاسخ را بر اساس سؤال شرطی کند؛ (c) این به‌طور قطعی کمتر از یک مدل دو جهتی عمیق قدرتمند است، زیرا نمی‌تواند از هر دو زمینه چپ و راست در هر لایه استفاده

کند.

۲.۵ تأثیر اندازه مدل

در این بخش، تأثیر اندازه مدل بر دقت وظایف فاین‌تیونینگ را بررسی می‌کنیم. ما چندین مدل BERT با تعداد متفاوتی از لایه‌ها، واحدهای پنهان و سرهای توجه آموزش دادیم، در حالی که سایر ابرپارامترها و روش آموزش همان‌طور که قبلاً توضیح داده شد، حفظ شدند.

نتایج در برخی از وظایف انتخابی GLUE در جدول ۶ نشان داده شده است. در این جدول، میانگین دقت مجموعه Dev از ۵ شروع تصادفی فاین‌تیونینگ گزارش شده است. می‌توان دید که مدل‌های بزرگ‌تر به‌طور قطعی باعث بهبود دقت در تمامی چهار مجموعه داده می‌شوند، حتی برای MRPC که تنها ۳۶۰۰ نمونه آموزشی برچسب‌گذاری شده دارد و تفاوت قابل توجهی با وظایف پیش‌آموزشی دارد. همچنین شاید تعجب‌آور باشد که ما قادر به دستیابی به چنین بهبودهای قابل توجهی بر روی مدل‌هایی هستیم که نسبت به ادبیات موجود، قبلاً نسبتاً بزرگ هستند. به عنوان مثال، بزرگ‌ترین ترنسفورمر مورد بررسی در ؟ دارای $L=۶$ ، $H=۱۰۲۴$ ، $A=۱۶$ و ۱۰۰ میلیون پارامتر برای انکودر است، و بزرگ‌ترین ترنسفورمر یافت شده در ادبیات دارای $L=۶۴$ ، $H=۵۱۲$ ، $A=۲$ و ۲۳۵ میلیون پارامتر است (؟). در مقایسه، BASE BERT دارای ۱۱۰ میلیون پارامتر و LARGE BERT دارای ۳۴۰ میلیون پارامتر است.

مدت‌هاست که مشخص شده افزایش اندازه مدل منجر به بهبود مداوم در وظایف مقیاس بزرگ مانند ترجمه ماشینی و مدل‌سازی زبان می‌شود، همان‌طور که توسط پریپلکسی داده‌های آموزشی نگه‌داشته‌شده در جدول ۶ نشان داده شده است. با این حال، ما بر این باوریم که این نخستین کار است که به‌طور قانع‌کننده نشان می‌دهد مقیاس‌دهی به اندازه‌های مدل بسیار بزرگ نیز منجر به بهبودهای قابل توجه در وظایف با مقیاس بسیار کوچک می‌شود، به شرط آنکه مدل به اندازه کافی پیش‌آموزش دیده باشد. ؟ نتایج مختلطی درباره تأثیر افزایش اندازه bi-LM پیش‌آموزشی از دو لایه به چهار لایه در وظایف پایین‌دستی ارائه دادند و ؟ به‌طور گذرا ذکر کردند که افزایش اندازه بعد پنهان از ۲۰۰ به ۶۰۰ مفید بود، اما افزایش بیشتر تا ۱۰۰۰ به بهبودهای بیشتری منجر نشد. هر دوی این کارهای پیشین از روش مبتنی بر ویژگی استفاده کرده‌اند — ما فرض می‌کنیم زمانی که مدل مستقیماً روی وظایف پایین‌دستی فاین‌تیون می‌شود و تنها از تعداد بسیار کمی پارامتر اضافی با مقداردهی تصادفی استفاده می‌کند، مدل‌های خاص وظیفه می‌توانند از نمایش‌های پیش‌آموزش‌دیده بزرگ‌تر و بیانگرتر بهره‌مند شوند حتی زمانی که داده‌های وظیفه پایین‌دستی بسیار کم باشد.

۳.۵ روش مبتنی بر ویژگی با BERT

تمام نتایج BERT که تاکنون ارائه شده‌اند، از روش فاین‌تیونینگ استفاده کرده‌اند، جایی که یک لایه ساده دسته‌بندی به مدل پیش‌آموزش‌دیده اضافه شده و همه پارامترها به‌طور مشترک روی یک وظیفه پایین‌دستی فاین‌تیون می‌شوند.

دقت مجموعه‌ی توسعه (Dev Set Accuracy)				ابزارامتراها (Hyperparams)		
SST-۲	MRPC	MNLI-m	(ppl) LM	#A	#H	#L
۸۸/۴	۷۹/۸	۷۷/۹	۵/۸۴	۱۲	۷۶۸	۳
۹۰/۷	۸۲/۲	۸۰/۶	۵/۲۴	۳	۷۶۸	۶
۹۱/۳	۸۴/۸	۸۱/۹	۴/۶۸	۱۲	۷۶۸	۶
۹۲/۹	۸۶/۷	۸۴/۴	۳/۹۹	۱۲	۷۶۸	۱۲
۹۳/۳	۸۶/۹	۸۵/۷	۳/۵۴	۱۶	۱۰۲۴	۱۲
۹۳/۷	۸۷/۸	۸۶/۶	۳/۲۳	۱۶	۱۰۲۴	۲۴

Table ۶: مطالعه‌ی حذفی (Ablation) بر روی اندازه‌ی مدل BERT. #L = تعداد لایه‌ها؛ #H = اندازه‌ی بردارهای پنهان؛ #A = تعداد سرهای توجه (Attention Heads). عبارت “LM” (ppl) نشان‌دهنده‌ی پیچیدگی مدل زبانی پوشانده‌شده (Masked Perplexity) LM بر روی داده‌های آموزشی کنارگذاشته‌شده است.

که امکان استفاده موفقیت‌آمیز از همان مدل پیش‌آموزش‌دیده برای مجموعه گسترده‌ای از وظایف NLP را فراهم می‌کند.

مراجع

پیوست مقاله: “BERT: پیش‌آموزش ترنسفورمرهای دوجهته عمیق برای درک زبان”

ما پیوست را به سه بخش سازماندهی کرده‌ایم:

- جزئیات اضافی پیاده‌سازی BERT در پیوست آ ارائه شده است؛
 - جزئیات اضافی مربوط به آزمایش‌های ما در پیوست ب ارائه شده است؛ و
 - مطالعات تحلیل تفکیکی (ablation) اضافی در پیوست پ ارائه شده است.
- ما مطالعات تحلیل تفکیکی اضافی برای BERT ارائه می‌دهیم که شامل موارد زیر است:
- تأثیر تعداد گام‌های آموزش؛ و
 - تحلیل تفکیکی برای روش‌های مختلف ماسک‌گذاری.

آ BERT for Details Additional

۱.۱ Tasks Pre-training the of Illustration

We provide examples of pre-training tasks in the following.

Procedure Masking the and LM Masked

Assuming the sentence is “the dog is hairy and during random masking the procedure we chose the ۴-th token (which corresponds to the word “hairy”) to mask. The procedure is illustrated further below.

- ۸۰٪ of the time: Replace the word with [MASK]
- hairy is dog my e.g., token, [MASK] the [MASK] is dog my

با این حال، روش مبتنی بر ویژگی‌ها، که در آن ویژگی‌های ثابت از مدل پیش‌آموزش‌دیده استخراج می‌شوند، مزایای خاص خود را دارد. اول، همه وظایف نمی‌توانند به راحتی با معماری رمزگذار ترنسفورمر نمایش داده شوند و بنابراین نیاز به اضافه کردن یک معماری مدل خاص وظیفه دارند. دوم، مزایای محاسباتی مهمی وجود دارد، زیرا می‌توان نمایش پرهزینه داده‌های آموزشی را یک بار پیش‌محاسبه کرد و سپس آزمایش‌های متعددی با مدل‌های ارزان‌تر روی این نمایش اجرا نمود.

در این بخش، دو روش را با اعمال BERT روی وظیفه شناسایی موجودیت‌های نام‌دار (NER) مجموعه داده ۲۰۰۳-CoNLL (۴) مقایسه می‌کنیم. در ورودی BERT، از مدل WordPiece با حفظ حروف استفاده می‌کنیم و حداکثر زمینه سند ارائه شده توسط داده‌ها را شامل می‌کنیم. مطابق رویه استاندارد، این وظیفه را به صورت برچسب‌گذاری فرموله می‌کنیم، اما از لایه CRF در خروجی استفاده نمی‌کنیم. نمایش اولین زیر-توکن به عنوان ورودی برای دسته‌بند سطح توکن روی مجموعه برچسب‌های NER استفاده می‌شود.

برای بررسی اثر فاین‌تیونینگ، روش مبتنی بر ویژگی را با استخراج فعال‌سازی‌ها از یک یا چند لایه بدون فاین‌تیون هیچ پارامتری از BERT اعمال می‌کنیم. این تعبیه‌های متنی به عنوان ورودی به یک BiLSTM دو لایه ۷۶۸ بعدی با مقداردهی تصادفی قبل از لایه دسته‌بندی استفاده می‌شوند.

نتایج در جدول ۷ ارائه شده‌اند. LARGE BERT عملکردی رقابتی با روش‌های پیشرو دارد. بهترین روش، نمایش‌های توکن را از چهار لایه مخفی بالایی ترنسفورمر پیش‌آموزش‌دیده متصل می‌کند، که تنها ۳۰.۰ F۱ پایین‌تر از فاین‌تیونینگ کل مدل است. این نشان می‌دهد که BERT برای هر دو روش فاین‌تیونینگ و مبتنی بر ویژگی مؤثر است.

۶ نتیجه‌گیری

بهبودهای تجربی اخیر ناشی از یادگیری انتقالی با مدل‌های زبانی نشان داده‌اند که پیش‌آموزش غنی و بدون نظارت، بخشی اساسی از بسیاری از سیستم‌های درک زبان است. به‌ویژه، این نتایج حتی به وظایف با منابع کم اجازه می‌دهد تا از معماری‌های عمیق یک‌جهته بهره‌مند شوند. سهم اصلی ما تعمیم بیشتر این یافته‌ها به معماری‌های دوجهته عمیق است،

سامانه (System)	F1 مجموعه‌ی توسعه (Dev)	F1 مجموعه‌ی آزمون (Test)
ELMo (?)	۹۵/۷	۹۲/۲
CVT (?)	-	۹۲/۶
CSE (?)	-	۹۳/۱
رویکرد تنظیم دقیق (Fine-tuning approach)		
LARGE BERT	۹۶/۶	۹۲/۸
BASE BERT	۹۶/۴	۹۲/۴
رویکرد مبتنی بر ویژگی (BERT (BASE (BERT		
تعبیه‌ها (Embeddings)	۹۱/۰	-
لایه‌ی ماقبل آخر (Second-to-Last Hidden)	۹۵/۶	-
لایه‌ی آخر (Last Hidden)	۹۴/۹	-
مجموع وزنی چهار لایه‌ی آخر (Weighted Four Last Sum Hidden)	۹۵/۹	-
اتصال چهار لایه‌ی آخر (Concat Four Last Hidden)	۹۶/۱	-
مجموع وزنی هر ۱۲ لایه (Weighted All Sum 12 Layers)	۹۵/۵	-

Table ۷: نتایج تشخیص موجودیت نام‌دار (NER) روی مجموعه‌ی CoNLL-۲۰۰۳. ابرپارامترها با استفاده از مجموعه‌ی توسعه (Dev) انتخاب شده‌اند. مقادیر گزارش شده‌ی Dev و Test میانگین ۵ اجرای تصادفی با همان ابرپارامترها هستند.

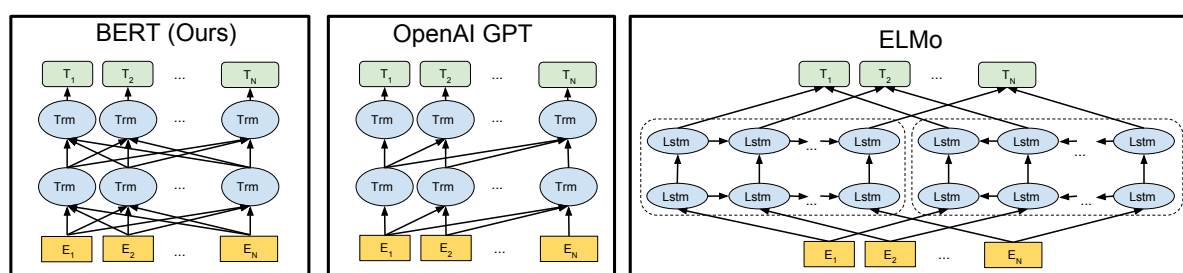


Figure ۳: Differences in pre-training model architectures. BERT uses a bidirectional Transformer. GPT uses a unidirectional Transformer. ELMo uses a unidirectional LSTM. The diagram illustrates the flow of information from input tokens (E1 to EN) through the model layers to output tokens (T1 to TN). BERT (Ours) uses a bidirectional flow, while GPT and ELMo use unidirectional flows.

Compared to standard language model training, the proposed model only masks the input tokens, which is a significant improvement. In the pre-training phase, the model is trained on a batch of 15% masked tokens. This approach allows the model to learn the relationships between tokens in a more efficient manner. The model's performance is evaluated using the F1 score, which is a combination of precision and recall. The results show that the proposed model outperforms the baseline models in terms of F1 score.

The next section illustrates the prediction task. The model is trained to predict the next token in a sequence given the previous tokens. This is a common task in natural language processing and is used to evaluate the model's ability to understand the context of a sentence.

● ۱۰% of the time: Replace the word with a random word.
e.g., my hairy dog is dog
apple is dog

● ۱۰% of the time: Keep the word.
e.g., my hairy dog is dog
The hairy is observed actual the towards representation word.

The advantage of this procedure is that it allows the model to learn the relationships between tokens in a more efficient manner. The model is trained to predict the next token in a sequence given the previous tokens. This is a common task in natural language processing and is used to evaluate the model's ability to understand the context of a sentence. The results show that the proposed model outperforms the baseline models in terms of F1 score.

our in pretraining up speed To length. quence
se- with model the pre-train we experiments.
steps. the of 90% for 128 of length quence
se- of steps the of 10% rest the train we Then.
embed- positional the learn to 512 of quence
dings.

Procedure Fine-tuning 3.1

hyperparameters model most fine-tuning. For
excep- the with pre-training. in as same the are
number and rate. learning size. batch the of tion
was probability dropout The epochs. training of
hyperparam- optimal The .1. at kept always
the found we but task-specific. are values eter
well work to values possible of range following
tasks: all across

22,16 :size Batch •
2e-5 2e-5, 5e-5 : (Adam) rate Learning •
4,3,2 : epochs of Number •

(e.g. sets data large that observed also We
less far were examples) training labeled 100k+
small than choice hyperparameter to sensitive
it so fast. very typically is Fine-tuning sets. data
search exhaustive an run simply to reasonable is
model the choose and parameters above the over
set. development the on best performs that

and ELMo BERT of Comparison 4.1 GPT OpenAI

pop- recent in differences the studies we Here
including models learning representation ular
com- The BERT. and GPT OpenAI ELMo.
are architectures model the between parisons
ad- in that Note .3 Figure in visually shown
BERT differences. architecture the to dition
approaches. fine-tuning are GPT OpenAI and
approach. feature-based a is ELMo while
pre-training existing comparable most The
trains which GPT. OpenAI is BERT to method
text large a on LM Transformer left-to-right a
decisions design the of many fact. In corpus.
as it make to made intentionally were BERT in
meth- two the that so possible as GPT to close
ar- core The compared. minimally be could ods
bi-directionality the that is work this of gument
Sec- in presented tasks pre-training two the and
empiri- the of majority the for account 1.3 tion
are there that note do we but improvements. cal
BERT how between differences other several
trained: were GPT and

examples. lowing

Input = [SEP] store [MASK] to went man the [CLS]
[SEP] milk [MASK] gallon a bought he
Label = IsNext

Input = [SEP] store the to [MASK] man the [CLS]
[SEP] birds ##less flight are [MASK] penguin
Label = NotNext

Procedure Pre-training 2.1

we sequence. input training each generate To
which corpus. the from text of spans two sample
are they though even “sentences” as to refer we
(but sentences single than longer much typically
receives sentence first The also). shorter be can
the receives second the and embedding A the
actual the is B time the of 50% embedding. B
time the of 50% and A follows that sentence next
the for done is which sentence. random a is it
sam- are They task. prediction” sentence “next
to- 512 ≤ is length combined the that such pled
Word- after applied is masking LM The kens.
rate masking uniform a with tokenization Piece
to given consideration special no and .10% of
pieces. word partial
sequences 256 of size batch with train We
to- 128,000 = tokens 512 * sequences 256)
ap- is which steps. 1,000,000 for kens/batch)
word billion 3.3 the over epochs 40 proximately
1e- of rate learning with Adam use We corpus.
of decay weight $L2, \beta_2 = 0.999, \beta_1 = 0.9, 4$
10,000 first the over warmup rate learning .01
We rate. learning the of decay linear and steps.
We layers. all on 1.0 of probability dropout a use
stan- the than rather (9) activation gelu a use
train- The GPT. OpenAI following .relu dard
like- LM masked mean the of sum the is loss ing
prediction sentence next mean the and lihood
likelihood.

4 on performed was BASEBERT of Training
TPU 16) configuration Pod in TPUs Cloud
was LARGE BERT of Training 13 total). chips
chips TPU 64) TPUs Cloud 16 on performed
com- to days 4 took pre-training Each total).
plete.

ex- disproportionately are sequences Longer
se- the to quadratic is attention because pensive

<https://cloudplatform.googleblog.com/2018/06/Cloud-TPU-now-offers-preemptible-pricing-and-global-availability.html>

originally were which of descriptions the
:¶ in summarized

Infer- Language Natural Multi-Genre **MNLI**
entailment crowdsourced large-scale, a is ence
sentences, of pair a Given .(¶) task classification
sen- second the whether predict to is goal the
neutral or ,*contradiction* ,*entailment* an is tence
one. first the to respect with

clas- binary a is Pairs Question Quora **QQP**
if determine to is goal the where task sification
semantically are Quora on asked questions two
.(¶) equivalent

Inference Language Natural Question **QNLI**
Answer- Question Stanford the of version a is
a to converted been has which (¶) Dataset ing
ex- positive The .(¶) task classification binary
do which pairs sentence) (question, are amplex
ex- negative the and answer, correct the contain
same the from sentence) (question, are amplex
answer, the contain not do which paragraph

is Treebank Sentiment Stanford The **SST-¶**
con- task classification single-sentence binary a
re- movie from extracted sentences of sisting
senti- their of annotations human with views
.(¶) ment

Acceptabil- Linguistic of Corpus The **CoLA**
task, classification single-sentence binary a is ity
English an whether predict to is goal the where
.(¶) not or “acceptable” linguistically is sentence

Similar- Textual Semantic The **STS-B**
sentence of collection a is Benchmark ity
other and headlines news from drawn pairs
score a with annotated were They .(¶) sources
two the similar how denoting δ to \backslash from
meaning, semantic of terms in are sentences

Cor- Paraphrase Research Microsoft **MRPC**
ex- automatically pairs sentence of consists pus
human with sources, news online from tracted
pair the in sentences the whether for annotations
.(¶) equivalent semantically are

bi- a is Entailment Textual Recognizing **RTE**
with but MNLI, to similar task entailment nary
^{¶¶}.(¶) data training less much

re- fine-tuning single-task report only we that Note^{¶¶}
could approach fine-tuning multitask A paper, this in sults
ex- For further, even performance the push potentially
RTE on improvements substantial observe did we ample,
MNLI, with training multi-task from

BooksCorpus the on trained is GPT •
on trained is BERT words): (¶••M
and words) (¶••M BooksCorpus the
words). (¶••M Wikipedia

and ([SEP]) separator sentence a uses GPT •
in- only are which ([CLS]) token classifier
learns BERT time: fine-tuning at troduced
embeddings B/A sentence and [CLS] , [SEP]
pre-training, during

batch a with steps $\backslash M$ for trained was GPT •
trained was BERT words: $\backslash \backslash \backslash \backslash$ of size
 $\backslash \backslash \backslash \backslash$ of size batch a with steps $\backslash M$ for
words.

δ e- of rate learning same the used GPT •
BERT experiments: fine-tuning all for δ
learning fine-tuning task-specific a chooses
devel- the on best the performs which rate
set, opment

we differences, these of effect the isolate To
 $\backslash \cdot \delta$ Section in experiments ablation perform
im- the of majority the that demonstrate which
two the from coming fact in are provements
they bidirectionality the and tasks pre-training
enable.

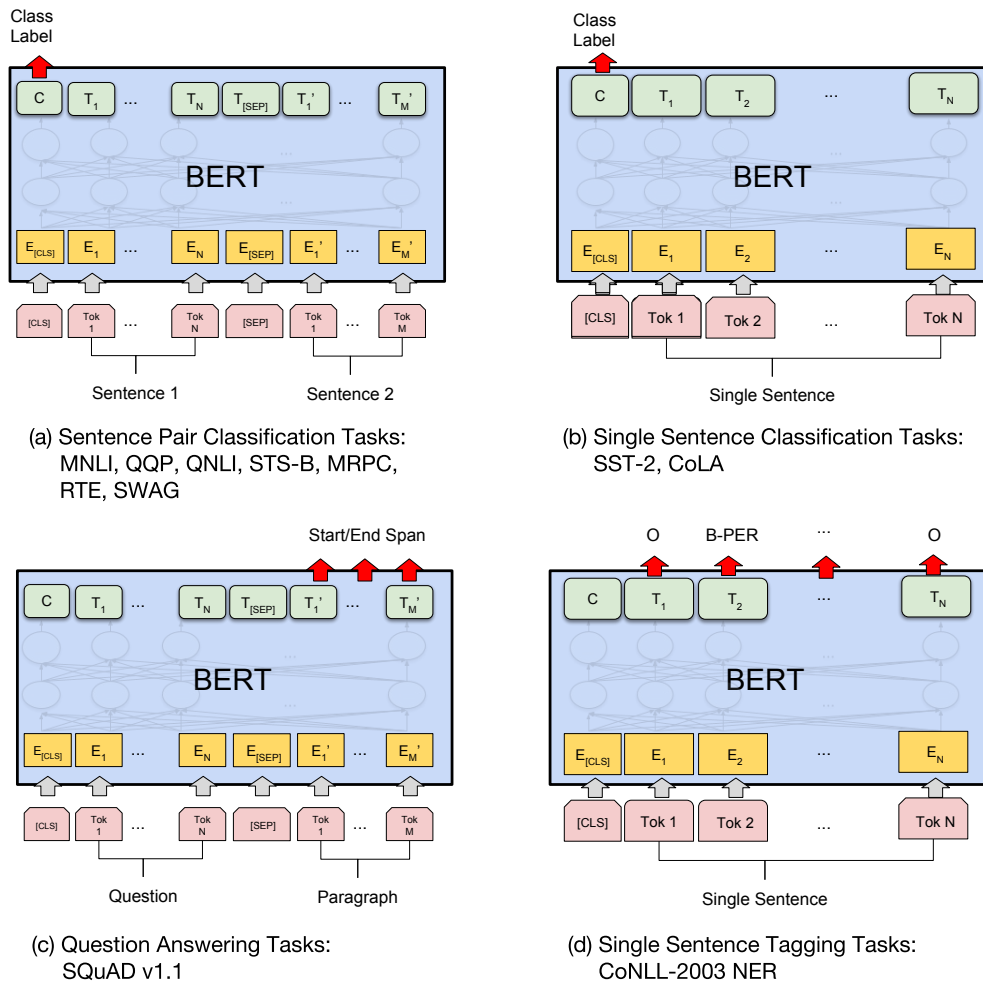
on Fine-tuning of Illustrations $\delta \cdot \backslash$ Tasks Different

dif- on BERT fine-tuning of illustration The
Our $\cdot \backslash$ Figure in seen be can tasks ferent
incorporat- by formed are models task-specific
layer, output additional one with BERT ing
be to need parameters of number minimal a so
and (a) tasks, the Among scratch, from learned
(d) and (c) while tasks sequence-level are (b)
represents E figure, the In tasks, token-level are
contex- the represents T_i embedding, input the
spe- the is [CLS] , i token of representation tual
is [SEP] and output, classification for symbol cial
non-consecutive separate to symbol special the
sequences, token

Setup Experimental Detailed ب

GLUE the for Descriptions Detailed $\backslash \cdot \backslash$ Experiments, Benchmark

from obtained are \backslash Table in results GLUE Our
<https://gluebenchmark.com/leaderboard>
<https://blog.openai.com/> and
GLUE The \cdot language-unsupervised
datasets, following the includes benchmark



Tasks. Different on BERT Fine-tuning of Illustrations : Figure

128,000) pre-training of amount large a achieve to steps) 1,000,000 * words/batch accuracy? fine-tuning high almost achieves BASE BERT Yes. Answer: when MNLI on accuracy additional 0%. 100k to compared steps 1M on trained steps.

con- pre-training MLM Does Question: 2 since pre-training, LTR than slower verge each in predicted are words of 15% only word? every than rather batch con- does model MLM The Answer: model. LTR the than slower slightly verge accuracy absolute of terms in However, the outperform to begins model MLM the immediately. almost model LTR

lan- natural small a is NLI Winograd WNL web- GLUE The .(?) dataset inference guage con- the with issues are there that notes page sys- trained every and 10 dataset. this of struction per- has GLUE to submitted been that's tem accuracy baseline 1.95 the than worse formed therefore We class. majority the predicting of For GPT. OpenAI to fair be to set this exclude the predicted always we submission. GLUE our class. majority

Studies Ablation Additional پ

Steps Training of Number of Effect 1. پ

after accuracy Dev MNLI presents 5 Figure been has that checkpoint a from fine-tuning answer to us allows This steps. k for pre-trained questions: following the

such need really BERT Does Question: 1

token. random other
 rep- table the of part left the in numbers The
 strategies specific the of probabilities the resent
 uses (BERT pre-training MLM during used
 paper the of part right The .(۱۰% ، ۱۰% ، ۸۰%
 feature- the For results. set Dev the represents
 layers ۴ last the concatenate we approach. based
 to shown was which features. the as BERT of
 .۳.۵ Section in approach best the be
 fine-tuning that seen be can it table the From
 strate- masking different to robust surprisingly is
 the only using expected. as However. gies.
 applying when problematic was strategy Mask
 Interest- NER. to approach feature-based the
 performs strategy Rnd the only using ingly.
 well. as strategy our than worse much

۲.۲ Masking Different for Ablation Procedures

uses BERT that mention we ، ۱.۳ Section In
 tokens target the masking for strategy mixed a
 language masked the with pre-training when
 an is following The objective. (MLM) model
 different of effect the evaluate to study ablation
 strategies. masking
 strate- masking the of purpose the that Note
 pre- between mismatch the reduce to is gies
 sym- [MASK] the as fine-tuning. and training
 stage. fine-tuning the during appears never bol
 and MNLI both for results Dev the report We
 fine-tuning both report we NER. For NER.
 the expect we as approaches. feature-based and
 feature-based the for amplified be will mismatch
 chance the have not will model the as approach
 representations. the adjust to

نتایج مجموعه‌ی توسعه (Dev Set)		
NER		MNLI
(Feature-based)	تنظیم دقیق (Fine-tune)	تنظیم دقیق (Fine-tune)
۹۴/۹	۹۵/۴	۸۴/۲
۹۴/۰	۹۴/۹	۸۴/۳
۹۴/۶	۹۵/۲	۸۴/۱
۹۴/۷	۹۵/۲	۸۴/۴
۹۴/۶	۹۴/۸	۸۳/۷
۹۴/۶	۹۴/۹	۸۳/۶

Table ۸: تحلیل تفکیکی (Ablation) بر روی استراتژی‌های مختلف ماسک‌گذاری.

the In .۸ Table in presented are results The
 to- target the replace we that means Mask table.
 Same MLM: for symbol [MASK] the with ken
 Rnd is: as token target the keep we that means
 an- with token target the replace we that means

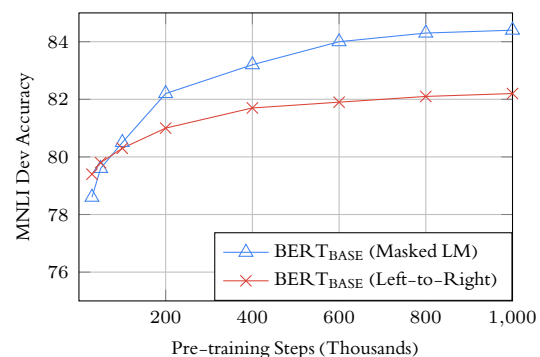


Figure ۵: Ablation of number of training steps. This shows the MNLI accuracy after fine-tuning. The x-axis is the number of steps k for trained model parameters that have been pre-