

ADAPTING PRETRAINED LANGUAGE MODELS FOR LONG DOCUMENT CLASSIFICATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Pretrained Language Models (LMs) have shown excellent results in achieving human-like performance on many language tasks. However, the most powerful LMs have one significant drawback: a fixed-sized input. Due to this constraint, these LMs are unable to utilize the full input of long documents. In this paper, we introduce a new framework to handle documents of arbitrary lengths. We investigate the addition of different mechanisms to handle an extended input size, and we demonstrate how attention can be used for parameter updates based on the most discriminating segments of input. We perform extensive validating experiments on patent and Arxiv datasets, both of which have long text. We show our method consistently outperforms state-of-the-art results reported in recent literature.

1 INTRODUCTION

Neural Network based Language Models (LMs) have seen a flurry of work, where new design and implementation improvements have advanced state-of-the-art performance in a variety of natural language tasks over the past few years (Devlin et al., 2018; Dai et al., 2019; Radford et al., 2019; Yang et al., 2019; Liu et al., 2019). LMs are powerful tools because they process a collection of unlabeled text and learn a rich embedding of natural language without supervision. This representation can be re-purposed on subsequent tasks such as classification and sentiment analysis (Korde & Mahender, 2012). This technique is essential for reaching state-of-the-art performance, as LM-based systems are able to achieve much better results than techniques that only use a small, labeled dataset. Modern LMs achieve their success by utilizing a powerful mechanism called “The Transformer” (Vaswani et al., 2017). The transformer learns strong dependencies between its inputs and has the ability to be stacked as many times as hardware can handle. This mechanism allows LMs to take in relatively large, yet still *fixed*, sized input.

For the largest LMs, the input size can reach up to four thousand tokens; however, this is still a limitation as they cannot process arbitrarily long documents. On many natural language tasks, this fixed input size is sufficient. For example, reading comprehension tasks are often used to analyze the quality of LMs, contain relatively few words on average, and do not have sufficiently long dependencies (Wang et al., 2019). Furthermore, there are multiple tasks where the input data is too long and must be truncated before being processed (Lee & Hsiang, 2019). Truncation is unsuitable because long complex text often contains inter-referential pieces of information. For instance, reading the final chapter of a book after all the previous ones takes on a different meaning compared to reading the same text by itself.

Solving the problem of arbitrarily long input requires more than a cursory glance. A first intuition may be to take a pretrained LM, separate the text into segments, place a Recursive Neural Network (RNN) after the embedding, and simply pass in the segments sequentially. This seems reasonable as RNNs have been used in the past for sequential text-based tasks, e.g. sentiment analysis (Socher et al., 2011). Unfortunately, using RNNs in this way causes two problems. First, RNNs are typically trained via backpropagation through time, making them prone to the problem of vanishing or exploding gradients (Pascanu et al., 2012). While many techniques exist to deal with this issue (Williams & Peng, 1990; Mujika et al., 2018), they do not solve the secondary problem of significant memory requirements for LM parameter updates. Each forward pass through the LM produces a set of parameter gradients. With multiple forward passes, the number of stored gradients quickly

grows as a segmented document increases in length. To further complicate the issue, the number of parameters in transformer-based networks is quadratic in relation to the fixed input size.

In order to solve these problems, we look at the brain’s ability to guide attention and provide behavioral updates as a biological inspiration. The brain processes an extraordinary amount of data, yet from moment to moment, much of that information is filtered out. What becomes filtered is not arbitrary, but is directly influenced by one’s objective. This filtering can be observed when an individual is tasked with counting the number of basketball passes in a video, yet fails to register a large gorilla that appears in center frame (Chabris & Simons, 2011). Even outside the moment to moment, changes in an individual’s behavior is guided by the structure of their values, self-selected or otherwise. An individual’s value structure imposes a framework for determining significance of events (Peterson, 1999). An event that may have been insignificant in the past can go on to take a new meaning once a new value has been gained or once an old one has changed (Laudet et al., 2006). This increase in valence may cause a behavioral change, a reorientation of goals, or a shift in the interpretation of experiences.

Unfortunately, when it comes to modeling this selective learning in Artificial Intelligence systems, common gradient-based methods fall short. A typical neural network model will update every parameter based on all of the inputs to minimize an objective function. While this is desired for many applications, full input based learning can cause issues for others (Pascanu et al., 2012). Ke et al. (2018) attempt to solve this issue via selective attentive backtracking, though they focus on remembering long-term dependencies rather than utilizing the full context of the input. In this work, we introduce a method that performs objective-based filtering during learning, but still utilizes the entire input.

We focus on two domains, specifically for classification: patents and scientific papers. Patent reading is a typical activity for lawyers trying to find relevant documents. Since 2003, the number of patents filed has increased nearly every year (WIPO, 2018); and using an automated system to perform classification is a continuously growing area of interest (Trappey et al., 2006). Scientific papers are another significant area of investigation. Many papers are uploaded to the internet every day, and their automatic categorization is becoming a necessity. A statistics paper, for example, may not be categorized as machine learning by its authors, but could be of interest to the machine learning community. Additionally, Tshitoyan et al. (2019) show material science concepts can be learned from scientific papers by a language model. They demonstrated materials for functional applications could be recommended several years before their initial discovery.

In this work, we use an attention mechanism to discover the significant portions of text to be used in updating a pretrained LM. We find attention-based updates to be essential because the most significant portion of an input sequence may occur anywhere throughout the document. While the datasets we study often start with highly discriminative features (titles and abstracts), we conduct experiments to show our attention mechanism can find the important parts of text even when it does not occur in the first segment. In either case, updating a LM’s parameters from only the first input segment performs well— often better than just using a baseline of the original LM with the input truncated to fit the max size. However, we find using an attention mechanism consistently improves performance and achieves the best results in our experimented language models. Therefore, our contributions are as follows:

1. We introduce a new framework for performing inference over arbitrary length documents.
2. We perform extensive validating experiments of our methods, showing how our attention-based framework consistently outperforms alternative methods.
3. We demonstrate the attention mechanism can be utilized to perform selective language model parameter updates.

2 RELATED WORKS

Language model pretraining is a popular method for tackling many natural language understanding tasks. Bidirectional Encoder Representations from Transformers (BERT) (Devlin et al., 2018) and Generative Pre-Training (GPT) (Radford et al., 2018) are two well-known language model pretraining methods that we utilize in this work. BERT is trained by selectively conditioning on most of

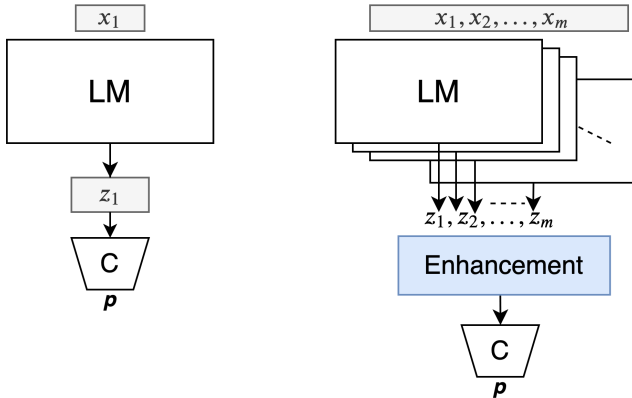


Figure 1: **Left:** The base language model (Base BERT/GPT) for classification. **Right:** The framework of language model enhancements for classification.

the input sequence, masking the rest, and attempting to predict the masked tokens. BERT is concurrently trained with a next sentence prediction task. The token used for the next sentence prediction task is reused for classification after BERT has been trained. GPT trains via next word prediction by conditioning on a sequence of text and trying to predict the next token. By training on Wikipedia and books, GPT is able to generate novel sequences of text. GPT is able to perform classification by appending a special classifying token at the end of an input, encoding the token into a latent representation, and using the representation as input to a new linear classifier set specifically for the task. These language model pretraining methods became state-of-art on natural language processing benchmarks such as GLUE (Wang et al., 2019) and SQUAD (Rajpurkar et al., 2016), achieving close to human-level performance.

As both BERT and GPT are based on transformers, their computational and storage costs scale quadratically with the input sequence length. This limits their application to mostly relatively short pieces of texts. To the best of our knowledge, we are not aware of any works applying these pretrained language models for long document classification. Next, we discuss a few other deep learning-based text classification approaches, with special attention to the classification of scientific papers and patents.

Dai & Le (2015) consider pretraining recurrent neural networks with large corpus of texts, and show improved performance on several text classification tasks. Kim (2014) introduced convolutional neural networks (CNN) for text/sentence classification. Yang et al. (2016) introduced a hierarchical attention mechanism for document classification that attends to interesting sentences and words in a document. The length of documents considered in latter two works are relatively short, with the corpus consisting of mostly individual sentences or online reviews.

For patent classification, Li et al. (2018) present a deep learning algorithm called Deep Patent, which is based on convolutional neural networks and continuous skip-gram embedding. They were the first to apply deep learning to large scale real-world patent classification. Lee & Hsiang (2019) used the pretrained BERT model to classify patents at the section and subclass level, only taking the title and abstract, or the first claim, as input.

For scientific paper classification, He et al. (2019) introduce a relatively large scientific paper dataset and perform classification through a multi-network approach. They introduce a RNN Attention-based reinforcement learning scheme. They select short text sequences to be parsed by a CNN. The representation learned by the CNN is sequentially fed back into the RNN for subsequent text selections and eventual classification.

Very recently, Cohan et al. (2019) construct a joint sentence representation that allows BERT transformer layers to directly utilize contextual information from all words in every sentence. However, their task is classification at the sentence level for a single text input.

3 METHODS

Our approach for classifying long documents is to divide the long document into a sequence of segments, each of which is short enough to be processed by a pretrained language model. The in-

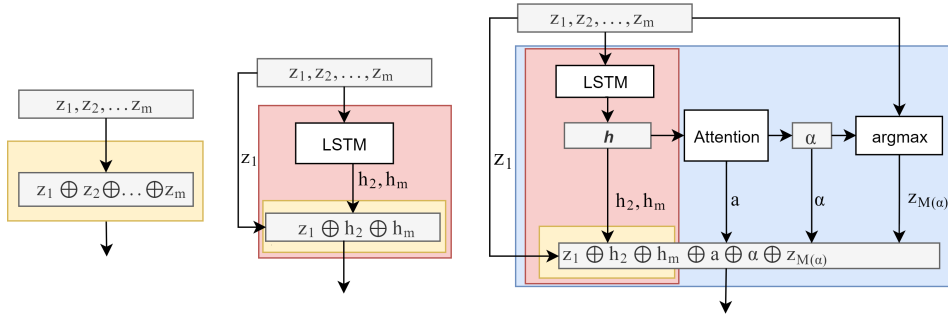


Figure 2: Three enhancements for base LM. Colors are used to represent the progressive enhancement of the combination strategies. **Left:** The concatenation model (Cat-BERT/GPT). **Center:** The RNN based model (RNN-BERT/GPT). **Right:** The full attention model (ATT-BERT/GPT).

formation from the language model’s representation of each segment is utilized in different manners to produce a classifier (see Figure 1(right)). We develop multiple strategies to combine these representations, progressively enhancing the combinations: starting from simple concatenation to an attention-based strategy. In this work, one of our primary contributions is to investigate the effectiveness of these different enhancements.

Let $\mathbf{x} = (x_1, x_2, \dots, x_m)$ be a document, where x_i is a fixed-length sequence of tokens (segment), and m the number of segments in the document. Let $y \in \mathbb{Y}$ be the respective labels in a k -class classification problem. We use $\mathbf{z}_i = \text{LM}(x_i) \in \mathbb{R}^d$ to denote the d -dimensional latent representation, of the segment x_i , for classification (e.g. the representation of the “CLS” token in BERT) from a selected language model LM. Let $C_{\mathbf{W}, \mathbf{b}}(v) = \sigma(\mathbf{W}v + \mathbf{b})$ be a linear classifier followed by the softmax function, and let \mathbf{p} the vector of probabilities of \mathbf{x} being assigned to each class.

Base Language Model (Base LM) In the usual application of deep language models such as BERT and GPT for text classification, the input text is truncated at a fixed length (256, 512, etc) due to limits in the size of the model. This corresponds to our basic model:

$$\mathbf{p} = C_{\mathbf{W}, \mathbf{b}}(\text{LM}(x_1)), \quad (1)$$

where $\mathbf{W} \in \mathbb{R}^{k \times d}$, $\mathbf{b} \in \mathbb{R}^k$, and we assume the segment length of x_1 equals the input size limit of the language model. This model is depicted in Figure 1(left).

Next we describe the three progressive enhancements to the base LM, shown in Figure 2.

Concatenated Language Model (Cat-LM) The first enhancement is a natural extension to improve the basic model by including information from more segments x_2, \dots, x_m . A very simple way to do this is to concatenate the representation $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_m$ before the classification layer. This leads to the model:

$$\mathbf{p} = C_{\mathbf{W}, \mathbf{b}}(\text{LM}(x_1) \oplus \text{LM}(x_2) \cdots \oplus \text{LM}(x_m)), \quad (2)$$

where $\mathbf{W} \in \mathbb{R}^{k \times md}$, $\mathbf{b} \in \mathbb{R}^k$. This model is difficult to perform backpropagation directly because we cannot hold m copies of the LM parameters in memory at the same time. We solve this problem by stopping the backpropagation paths of some of the segments, namely $\mathbf{z}_2, \dots, \mathbf{z}_m$. Section 3.1 discusses this approximation.

RNN-augmented Language Model (RNN-LM) The third model we want to consider is one that summarizes the information from $\mathbf{z}_1, \dots, \mathbf{z}_m$ using a bidirectional LSTM (Hochreiter & Schmidhuber, 1997). Let $(\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_m) = \text{biLSTM}(\mathbf{z}_1, \dots, \mathbf{z}_m)$ be the q dimensional hidden state representations from a bidirectional LSTM, where $\mathbf{h}_i \in \mathbb{R}^q$. The biLSTM-based model can be written as:

$$\mathbf{p} = C_{\mathbf{W}, \mathbf{b}}(\text{LM}(x_1) \oplus \mathbf{h}_2 \oplus \mathbf{h}_m), \quad (3)$$

where $\mathbf{W} \in \mathbb{R}^{k \times (d+2q)}$, $\mathbf{b} \in \mathbb{R}^k$. For this model we also stop the gradient computation at $\mathbf{z}_2, \dots, \mathbf{z}_m$, and do not backpropagate beyond the LSTM parameters.

We do not include $\mathbf{z}_2, \dots, \mathbf{z}_m$ as input to the classifier as the size of \mathbf{W} and \mathbf{b} scale directly with m , and with a large m the number of parameters would grow counter-productively. However, we still

include z_1 as the text from that segment typically contains the most discriminative features. By including z_1 , there is a direct path for gradient updates to the language model. Additionally, we found that including h_1 , instead of h_2 , in the classifier consistently produced worse results. We believe this is from the forward direction of h_1 not providing any contextual information. Comparatively, the backwards direction of h_2 provides enough context that, when concatenated with z_1 , achieves the best performance.

Attention-based Language Model (ATT-LM) For our attention-based model, we utilize the structure of RNN-LM and add the attention mechanism as described in Yang et al. (2016). We define our attention variables as follows:

$$\begin{aligned} \mathbf{u}_i &= \tanh(\mathbf{W}_s \mathbf{h}_i + \mathbf{b}_s), i = 1, \dots, m \\ \alpha_i &= \sigma(\mathbf{u}_i^T \mathbf{u}_s) \\ \mathbf{a} &= \sum_{i=1}^m \alpha_i \mathbf{h}_i, \end{aligned} \tag{4}$$

where $\mathbf{W}_s \in \mathbb{R}^{2q \times q}$, $\mathbf{b}_s \in \mathbb{R}^q$, $\mathbf{u}_s \in \mathbb{R}^q$ are learned attention parameters for attention over segments. Let $M(\boldsymbol{\alpha}) = \arg \max_{1 \leq i \leq m} \alpha_i$ be the index of the segment that gives the highest attention weight. These definitions give us our equation for the attention-based model by concatenating a set of relevant features:

$$\mathbf{p} = C_{\mathbf{W}, \mathbf{b}}(\text{LM}(x_1) \oplus \mathbf{h}_2 \oplus \mathbf{h}_m \oplus \mathbf{a} \oplus \boldsymbol{\alpha} \oplus \text{LM}(x_{M(\boldsymbol{\alpha})})), \tag{5}$$

where $\mathbf{W} \in \mathbb{R}^{k \times (2d+3q+m)}$, $\mathbf{b} \in \mathbb{R}^k$ are the parameters of the linear classifier. For this model, we stop the gradient computations paths for all z_i 's apart from z_1 and the selected $z_{M(\boldsymbol{\alpha})}$.

For all four methods, we let n be the number of training documents and use the negative log likelihood as the loss function:

$$L = - \sum_{i=1}^n \log p_{y_i}, \tag{6}$$

3.1 PARAMETER UPDATES

We run mini-batch stochastic gradient descent for parameter updates. We treat the parameters of the language models and the other parameters (weight matrix \mathbf{W} for classification, LSTM parameters, etc) differently. We perform full gradient computation on non-LM parameters, and only approximate gradient computation for the LM's parameters by stopping backpropagation on selected segments (see the model description above and the Appendix Section A).

4 EXPERIMENTS

Implementation Details. We use PyTorch (Paszke et al., 2017) to conduct all our experiments. For our pretrained language models, we use the HuggingFace (2019) implementation of the Base-BERT model (110M parameters) and use the GPT model of similar size (117M parameters). Due to computational constraints, we do not use BERT-Large or larger GPT models. We use BERT's lower-case tokenizer and GPT's tokenizer with an added classification token. Both models take a fixed input size of 256 tokens, contain 12 transformer blocks, and have a hidden size of 768 neurons. We apply dropout (Srivastava et al., 2014) ($p = .1$) before the final linear layer.

Unless otherwise noted we use a learning rate of $2e^{-5}$, use a scheduled ADAM optimizer (Kingma & Ba, 2015), train for 3 epochs over each training dataset, use a training mini-batch size of 32 documents, and set all other hyper parameters to their default values. While the number of segments m does not need to be fixed, for ease of processing we set $m = 8$ for the patent datasets as patents have a relatively constrained length and 8 allows for minimal padding to be used. We also use $m = 16$ segments for the Arxiv dataset as one document typically contains 6k words, many of which are removed as non-meaningful or are intentionally truncated as part of the bibliography.

Next, we describe the different datasets we use in our experiments.

Method \ Dataset	Arxiv-4	Arxiv-11	section	subclass	wireless	inverted wireless
Li et al. (2018) Deep Patent	-	-	-	<43	-	-
Lee & Hsiang (2019) PatentBert	-	-	80.98	66.80	-	-
He et al. (2019) Local Word Glimpses	94.18	80.47	-	-	-	-
Base-GPT	96.59	84.62	83.32	67.29	89.82	87.69
Base-BERT	97.06	87.42	83.85	68.31	90.21	87.72
Cat-GPT	96.82	80.03	83.43	66.17	89.34	88.80
Cat-BERT	97.06	87.34	83.99	68.34	90.64	89.39
RNN-GPT	96.98	85.31	83.52	67.72	90.16	89.19
RNN-BERT	97.62	87.72	83.99	68.72	90.51	89.41
ATT-GPT	97.62	85.94	83.66	68.13	90.31	90.08
ATT-BERT	97.70	87.96	84.13	69.01	90.69	90.25

Table 1: Micro F1 results on our datasets.

Patents. Patents can be broken down into multiple levels of resolution according to the International Patent Classification System (IPC): *Section, Class, Subclass, ... etc.* The most broad category, *Section*, has eight labels (A-H). For instance, *Section A* is concerned with *Human Necessities*, while *Section H* is concerned with *Electricity*. We also perform classification experiments on a more detailed level of categorization: *Subclass*, which contains 638 labels. Patents were gathered from the Google Patents Public Dataset via SQL queries.

We gathered all documents from the United States Patent Office (USPTO) from 2006-2014 as our training set and use all patents from 2015 as our test set. We have 1,917,334 training and 296,724 testing documents, where 15,172 and 1,835 documents were respectively skipped for missing abstracts. The text of a patent is composed of different parts: title, abstract, and a list of claims. For our purposes, we consider one patent to be first the title, then the abstract, followed by each claim in order—claim 1, claim 2, ... until the last claim.

(Inverted) Wireless Patents. We selected a subset of patent data to perform additional experiments. We chose the wireless (H04) *Class* due to its large number of training and test examples (the second most of all the *Class* data). The *Class* with the most examples (computing) was not chosen due to its large imbalance in *Subclasses* (dominated by one label with over 75% examples). We use the wireless *Class* to construct an inverted patent dataset, where a single patent starts with its last claim, up to the first claim in the reverse order, then the abstract, and lastly the title. It is commonly believed that the abstract and the first claim is the most useful in classifying a patent. We create this dataset to present information in reverse order of relevance to test models that bias towards the beginning of documents (e.g., models that truncate beyond a fixed number of tokens). After processing, the wireless *Class* contains 250,982 training and 42,892 testing documents, where 15,172 and 97 documents were skipped, respectively, for missing abstracts.

Arxiv papers. We study the long document dataset provided by He et al. (2019). It consists of 33,388 papers, from 11 different categories, downloaded from the scientific article hosting website Arxiv. The least occurring category is “math.AC” with 2885 documents, and the most occurring is “cs.DS” with 4136 documents. We call this dataset Arxiv-11. They also provide a subset of the data using four categories and 12,195 documents, which we refer to as Arxiv-4. All downloaded pdf documents were converted to txt files, with no document less than 1,000 words. We randomly sample 90% for training and use the remaining 10% for test.

5 RESULTS

Patents and Arxiv Datasets. We report our main results in Table 1, where the numbers are micro-F1 scores. We make three observations. First, we compare our methods to the previous work using DeepPatent (Li et al., 2018), PatentBert (Lee & Hsiang, 2019) and Local Word Glimpses (He

Method \ Param	3 epochs 2e-5 lr	3 epochs 3e-5 lr	3 epochs 4e-5 lr	3 epochs 5e-5 lr	8 epochs 2e-5 lr	8 epochs 5e-5 lr
Base GPT	96.19	96.59	96.43	96.19	96.27	96.35
Base BERT	96.90	97.06	96.43	96.27	96.58	97.06
Cat-GPT	93.88	94.28	95.63	95.31	96.74	96.82
Cat-BERT	96.82	96.98	96.27	96.66	96.90	97.06
RNN-GPT	96.03	96.82	96.59	96.98	96.90	96.74
RNN-BERT	97.62	96.74	97.06	96.74	96.51	96.90
ATT-GPT	95.79	97.38	97.62	97.62	96.43	97.14
ATT-BERT	97.06	97.70	96.27	95.87	96.98	97.14

Table 2: The effect of changing the learning rate, and epoch, hyper-parameters for all models on the Arxiv-4 dataset.

et al., 2019). Table 1 shows the Base LMs perform well. We note our Base BERT implementation outperforms PatentBert as they do not combine the title, abstract, and claims– but only the title and abstract or just the first claim. Second, we compare Base LMs against the three enhancements across 4 variants of the patent dataset and 2 variants of the Arxiv dataset. The RNN-based and attention-based LMs show consistent improvement over base LMs, while the simple concatenation-based LM is not consistent. Furthermore, the usage of attention is superior in all cases (highlighted in the last line). Lastly, it is also interesting to note that BERT-based models consistently outperform GPT-based models. This is likely due to the use of bidirectional contexts in BERT.

Inverted Patents. As shown in recent work (Lee & Hsiang, 2019; Li et al., 2018), only the abstract or first claim on a patent is needed for good classification performance. To analyze the effect of attention discovering the location of discriminative content, we invert the structure of the patents in our wireless dataset and train new models. Comparing the results from the second last column (wireless) and last column (inverted wireless) in Table 1, we see that Base LMs, which only take into account of first 256 tokens, suffer from a drop of F1 scores of more than 2.0. Cat-LMs and RNN-LMs reduce the gap in F1 scores to about 1.0, while the attention-based models perform the best, with the gap between wireless and inverted wireless reducing to smaller than 0.5.

Exploration of Hyper-parameters. Next, we investigate the effects of different hyper-parameters on the various models using a small set of training data. For this analysis, we use the Arxiv-4 dataset. While the attention-based models seem to do well, table 2 demonstrates the unpredictable nature of using different learning rates along with different training epochs. This behavior aligns with the claims of Devlin et al. (2018), who also found fine-tuning on small datasets sometimes lead to unstable results.

Training and Evaluation Time. Lastly, we compare the training time and evaluation speed of our models on the Arxiv-11 dataset. As shown in Table 3, all of the enhanced variants of the LMs require nearly 3-4x training time and over 2x to evaluate. Considering there are substantially more operations required to process the full input text via multiple forward passes, this slow down is better than expected. The difference in training time between the attention-based models and non-attention-based is surprising, given the fact that the LM’s parameters must be updated with more than one set of gradients. This points towards further gradient computations being feasible for architectures and hardware that can handle the additional required memory.

	Base LM	Cat-LM	RNN-LM	ATT-LM
BERT Training time	1.000	2.711	3.079	3.610
BERT Evaluation time	1.000	2.712	2.940	3.147
GPT Training time	1.000	3.088	3.551	3.558
GPT Evaluation time	1.000	2.520	2.689	2.839

Table 3: An analysis of the different model run times as a factor of the baseline method.

	Base BERT	Cat-BERT	RNN-BERT	ATT-BERT	Ablated ATT-BERT
x_1 gradients	11.92	12.96	11.92	81.05	81.11
No x_1 gradients	11.92	46.72	70.59	81.08	81.11

Table 4: BERT model results when shuffling every x_1 in the Arxiv-11 dataset.

5.1 ABLATION EXPERIMENT

In order to understand the effects of the attention, we introduce an ablated model that allows us to investigate how attention influences and guides backpropagation. The model we use is similar to ATT-BERT, but we remove $LM(x_1)$, h_2 , and h_m from the classifier. Therefore it can be written as:

$$p = C_{W,b}(a \oplus \alpha \oplus LM(x_{M(\alpha)})), \quad (7)$$

with $W \in \mathbb{R}^{k \times (q+m+d)}$, $b \in \mathbb{R}^k$ being the parameters of the linear classifier. This model performs well, but not better than ATT-BERT. The full results of this model can be seen in table 5 in the appendix. But here, we use this model to carry out two additional experiments.

Shuffling Experiment. First, we perform a shuffled input experiment to examine the effect of using attention to guide language model parameter updates. We use all the same setup as the experiments for the Arxiv-11 column in Table 1, except for each training iteration, we randomly permute each segment x_1 between all mini-batch examples. This means methods that use x_1 will no longer be able to rely upon those gradients to give an informative update to the model’s parameters. Table 4 shows how the non-attention-based methods guess the maximum occurring class when gradients are updated using x_1 . The table also shows how these methods perform poorly without gradients to update the language model parameters. However, the attention-based methods are able to perform well despite the loss of information from shuffling, and ATT-BERT is relatively unaffected by the loss of gradients from x_1 .

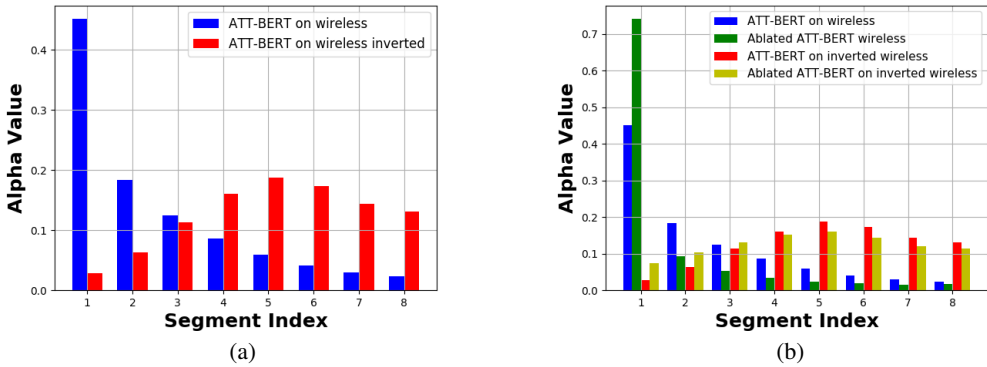


Figure 3: (a) An analysis of the effect on attention α values between the original wireless dataset, and the inverted wireless dataset. (b) A comparison of α values with the ablated model.

Attention α Comparisons. Second, we measure and compare the α values on the wireless and inverted wireless patent dataset. We average the α value for each of the eight input segments over the entire test set. Figure 3 (a) comparatively shows the α values of our attention-based model on the wireless and inverted wireless dataset; this clearly demonstrates how the attention mechanism is able to accurately pick up on the important sections of text. Figure 3(b) shows a comparison between ATT-BERT and the ablated model on both the wireless and inverted wireless datasets. The α values for both models on the inverted dataset seem relatively similar, with the ablated model placing heavier emphasis on the first few segments. But for the wireless dataset, the effect of including the

gradients of x_1 can be seen on, where the α_1 value is over 0.7 on the ablated model. This means enabling backpropagation to occur for input x_1 has a positive effect for ATT-BERT and, moreover, that the title and abstract in a patent are of high discriminative importance.

6 CONCLUSIONS

In this work, we achieved state-of-the-art results on multiple long document classification tasks by utilizing pretrained language models with attention-based enhancements. With language modeling continuing to see improvements every month, we showed how different models can be utilized with our method. We performed numerous experiments to clearly demonstrate the value added by using attention to learn from important segments. We showed that the additional gradient computation as a result of attention is marginal when compared to the consistent improvement of results. We analyzed the effects of the attention mechanism through the loss of input information via both shuffling experiment and a carefully constructed dataset augmentation.

REFERENCES

- Christopher F Chabris and Daniel Simons. *The invisible gorilla: And other ways our intuitions deceive us*. Harmony, 2011.
- Arman Cohan, Iz Beltagy, Daniel King, Bhavana Dalvi, and Daniel S. Weld. Pretrained language models for sequential sentence classification. *CoRR*, abs/1909.04054, 2019. URL <https://arxiv.org/pdf/1909.04054.pdf>.
- Andrew M Dai and Quoc V Le. Semi-supervised sequence learning. In *Advances in Neural Information Processing Systems*, pp. 3079–3087, 2015.
- Zihang Dai, Zhilin Yang, Yiming Yang, William W Cohen, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Jun He, Liqun Wang, Liu Liu, Jiao Feng, and Hao Wu. Long document classification from local word glimpses via recurrent attention learning. *IEEE Access*, 7:40707–40718, 2019. URL <https://doi.org/10.1109/ACCESS.2019.2907992>.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.
- HuggingFace. Pytorch-transformers. <https://github.com/huggingface/pytorch-transformers>, 2019.
- Nan Rosemary Ke, Anirudh Goyal GOYAL, Olexa Bilaniuk, Jonathan Binas, Michael C Mozer, Chris Pal, and Yoshua Bengio. Sparse attentive backtracking: Temporal credit assignment through reminding. In *Advances in Neural Information Processing Systems*, pp. 7640–7651, 2018.
- Yoon Kim. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*, 2014.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR*, 2015.
- Vandana Korde and C Namrata Mahender. Text classification and classifiers: A survey. *International Journal of Artificial Intelligence & Applications*, 3(2):85, 2012.
- Alexandre B Laudet, Keith Morgen, and William L White. The role of social supports, spirituality, religiousness, life meaning and affiliation with 12-step fellowships in quality of life satisfaction among individuals in recovery from alcohol and drug problems. *Alcoholism treatment quarterly*, 24(1-2):33–73, 2006.
- Jieh-Sheng Lee and Jieh Hsiang. PatentBERT: Patent classification with fine-tuning a pre-trained BERT model. *CoRR*, abs/1906.02124, 2019. URL <http://arxiv.org/abs/1906.02124>.
- Shaobo Li, Jie Hu, Yuxin Cui, and Jianjun Hu. DeepPatent: patent classification with convolutional neural networks and word embedding. *Scientometrics*, 117(2):721–744, 2018. doi: 10.1007/s11192-018-2905-5.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Asier Mujika, Florian Meier, and Angelika Steger. Approximating real-time recurrent learning with random kronecker factors. In *Advances in Neural Information Processing Systems*, pp. 6594–6603, 2018.
- Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. Understanding the exploding gradient problem. *CoRR*, abs/1211.5063, 2, 2012.

- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017.
- J.B. Peterson. *Maps of Meaning: The Architecture of Belief*. Routledge, 1999. ISBN 9780415922227.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. URL https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language_understanding_paper.pdf, 2018.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019. URL <https://d4mucfpxsywv.cloudfront.net/better-language-models/language-models.pdf>.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2383–2392, 2016.
- Richard Socher, Jeffrey Pennington, Eric H Huang, Andrew Y Ng, and Christopher D Manning. Semi-supervised recursive autoencoders for predicting sentiment distributions. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 151–161. Association for Computational Linguistics, 2011.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- Amy JC Trappey, Fu-Chiang Hsu, Charles V Trappey, and Chia-I Lin. Development of a patent document classification and search platform using a back-propagation network. *Expert Systems with Applications*, 31(4):755–765, 2006.
- Vahe Tshitoyan, John Dagdelen, Leigh Weston, Alexander Dunn, Ziqin Rong, Olga Kononova, Kristin A Persson, Gerbrand Ceder, and Anubhav Jain. Unsupervised word embeddings capture latent knowledge from materials science literature. *Nature*, 571(7763):95, 2019.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pp. 5998–6008, 2017.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *7th International Conference on Learning Representations, ICLR, 2019*.
- Ronald J Williams and Jing Peng. An efficient gradient-based algorithm for on-line training of recurrent network trajectories. *Neural computation*, 2(4):490–501, 1990.
- WIPO. World intellectual property indicators. *Geneva: World Intellectual Property Organization*, 2018.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. *arXiv preprint arXiv:1906.08237*, 2019.
- Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alexander J. Smola, and Eduard H. Hovy. Hierarchical attention networks for document classification. In *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2016.

Method \ Dataset	Arxiv-4	Arxiv-11	section	subclass	wireless	inverted wireless
Ablated ATT-BERT	97.38	87.22	84.01	68.87	90.31	89.98
ATT-BERT	97.70	87.96	84.13	69.01	90.69	90.25

Table 5: Micro F1 results of our ablated BERT model compared to the ATT-BERT.

A GRADIENT COMPUTATION FOR LANGUAGE MODEL PARAMETERS

Let us denote our language model representation be $z_i = \text{LM}(x_i) = f_{\theta}(x_i)$, where θ are the language model parameters. Let h_{ϕ} be the function we compute on top of the LM representations, e.g., the classifier, LSTM, etc. The parameter ϕ can contain the classification weights W and the LSTM weights. The models considered in this paper can be written as:

$$\mathbf{p} = h_{\phi}(f_{\theta}(x_1), \dots, f_{\theta}(x_m))$$

Coupled with the loss function l (log loss) and the target label y , we have

$$l(\mathbf{p}, y) = l(h_{\phi}(f_{\theta}(x_1), \dots, f_{\theta}(x_m)), y)$$

Computing the gradient over the LM parameters θ , by chain rule we have

$$\begin{aligned} \frac{\partial}{\partial \theta} l(\mathbf{p}, y) &= \frac{\partial}{\partial \mathbf{z}} l(\mathbf{z}, y) \Big|_{\mathbf{z}=h_{\phi}(\dots)} \left[\frac{\partial}{\partial \mathbf{u}} h_{\phi}(\mathbf{u}, f_{\theta}(x_2), \dots, f_{\theta}(x_m)) \frac{\partial}{\partial \theta} f_{\theta}(x_1) \right. \\ &\quad \left. + \frac{\partial}{\partial \mathbf{u}} h_{\phi}(f_{\theta}(x_1), \mathbf{u}, \dots, f_{\theta}(x_m)) \frac{\partial}{\partial \theta} f_{\theta}(x_2) + \dots + \frac{\partial}{\partial \mathbf{u}} h_{\phi}(f_{\theta}(x_1), \dots, \mathbf{u}) \frac{\partial}{\partial \theta} f_{\theta}(x_m) \right] \end{aligned} \quad (8)$$

By stopping the gradient computation over x_2, \dots, x_m , we are dropping the terms related to $\frac{\partial}{\partial \theta} f_{\theta}(x_i)$, $i \geq 2$, from the above formula. In optimization, we say \mathbf{q} is a descent direction if $\langle \mathbf{q}, \partial_{\theta} l \rangle < 0$. The negative gradient $-\partial_{\theta} l$ is clearly a descent direction. The gradient above in Equation (8) is in the form $\partial_{\theta} l = c \sum_{i=1}^m \mathbf{g}_i$, where \mathbf{g}_i are the gradient term of the i th segment in the equation. By dropping terms in the backpropagation, we are assuming either the contribution from \mathbf{g}_1 (or $\mathbf{g}_{M(\alpha)}$ in the attention model of Equation 5) dominates the contributions from other segments, or the gradients \mathbf{g}_i from different segments point towards similar directions, so that the truncated gradient is still a descent direction.

In some cases, if the storage of partial gradients $\frac{\partial}{\partial \mathbf{u}} h_{\phi}(z_1, \dots, \mathbf{u}, \dots, z_m)$ is feasible, it is possible to compute the full gradient in two passes over the mini-batch. For example, in the basic Cat-BERT model for classification, this partial gradient is \mathbf{W}_i scaled by the derivative of softmax, where i is the index of the segment. We can compute these \mathbf{W}_i with the derivative of the softmax as scaling factors in a first pass, and accumulate gradients over $f_{\theta}(x_i)$ scaled by these factors in a second pass over the mini-batch.