# ERNIE-Sparse: Robust Efficient Transformer Through Hierarchically Unifying Isolated In-Formation

## Anonymous authors

Paper under double-blind review

### Abstract

Sparse Transformer has recently attracted a lot of attention since the ability for reducing the quadratic dependency on the sequence length. In this paper, we argue that two factors could affect the robustness and causing performance degradation of the Sparse Transformer. The first factor is *information bottleneck sensitivity*, which is caused by the key feature of Sparse Transformer - only a small number of global tokens can attend to all other tokens. The second factor is sparse pattern sensitivity, which is caused by different token connections in different sparse patterns. To address these issues, we propose a well-designed model, named ERNIE-SPARSE. It consists of two distinctive parts: (i) a Hierarchical Sparse Transformer (HST) mechanism, which introduces special tokens to sequentially model local and global information. This method is not affected by bottleneck size and improves model robustness and performance. (ii) Sparse-Attention-Oriented Regularization (SAOR) method, the first robust training method designed for Sparse Transformer, which increases model robustness by forcing the output distributions of transformers with different sparse patterns to be consistent with each other. To evaluate the effectiveness of ERNIE-SPARSE, we perform extensive evaluations. Firstly, we perform experiments on a multi-modal long sequence modeling task benchmark, Long Range Arena (LRA). Experimental results demonstrate that ERNIE-SPARSE significantly outperforms a variety of strong baseline methods including the dense attention and other efficient sparse attention methods and achieves improvements by 2.7% (55.01% vs. 57.78%). Secondly, to further show the effectiveness of our method, we pretrain ERNIE-SPARSE and verified it on 3 text classification and 2 QA downstream tasks, achieve improvements on classification benchmark by 0.83% (91.63% vs. 92.46%), on QA benchmark by 3.27% (74.7% vs. 71.43%). Experimental results continue to demonstrate its superior performance.

## **1** INTRODUCTION

Transformer (Vaswani et al., 2017) architecture is a key component for many pretrained language models such as BERT (Devlin et al., 2018), RoBERTa Liu et al. (2019), XLNet (Yang et al., 2019), ERNIE (Sun et al., 2020b), ALBERT (Lan et al., 2019), ELECTRA(Clark et al., 2020), T5 (Raffel et al., 2020). Self-attention is one of the most important modules in transformer. It eliminates the sequential dependency constraints of recurrent neural networks by introducing interactions between each token pair to capture contextual information. However, the self-attention's computational complexity and memory cost grow quadratically with sequence length, which comes with complexity  $O(N^2)$  for processing contexts of N inputs.

One way to optimize self-attention complexity is introducing sparsity into attention layers (Child et al., 2019; Qiu et al., 2019; Beltagy et al., 2020) by having each token attend to only a subset of tokens in the whole sequence. Recent sparse-attention works (Child et al., 2019; Beltagy et al., 2020; Ainslie et al., 2020; Zaheer et al., 2020) introduce global tokens that can attend to the whole sequence. Those global tokens are used as a form of memory to strengthen global information. While this method reduces the complexity of full self-attention, there are two issues with Sparse Transformer that affect robustness.

The first issue is *information bottleneck sensitivity* as show in Figure 1. Information bottleneck is a phenomenon caused by the low number of global tokens — the model had to encapsulate the entire input sequence into those global tokens. If the size of the information bottleneck becomes smaller, performance can suffer. The second issue is *sparse pattern sensitivity*. Sparse pattern sensitivity means that transformer with different sparse patterns may produce different outputs for the same input. For example, in Figure 2, the left sparse attention divides the input into two blocks, while the right into three by shifting the pattern two tokens to the upper left corner, resulting in the change of attention topology change causes transformer output different values even for the same input. This will leads to over sensitivity to pattern change and less robustness. To increase model robustness, existing works include adversarial training (Goodfellow et al., 2015; Miyato et al., 2019; Liu et al., 2020) and regularization methods (Krizhevsky et al., 2012; Srivastava et al., 2014). However, those methods are designed for dense models. To the best of our knowledge, we are the first to study robust training method for Sparse Transformers and propose a simple but efficient solution to improve the robustness.

To resolve the aforementioned issues, we propose a well-designed efficient model ERNIE-SPARSE. Our method proposes two major techniques including Hierarchical Sparse Transformer (HST) and Self-Attention Oriented Regularization (SAOR). The first technique HST is designed for *information bottleneck sensitivity* problem based on hierarchical attention mechanism. HST introduces special tokens to represent local information and global information in two sequential steps, providing an additional way for modeling global information that is not affected by bottleneck size. The second technique SAOR, is designed for *sparse pattern sensitivity* problem. The main idea is to force models with different sparse patterns to have consistent outputs for two identical inputs. Except for robustness, another benefit of this method is it only needs to be used in the training phase and does not affect the inference of the model.

To evaluate ERNIE-SPARSE's ability to model long documents, we first perform extensive evaluations on a long sequence modeling task benchmark (Tay et al., 2020c): Long Range Arena (LRA). Experiment results demonstrate that our model significantly outperforms existing methods and provides a new robust solution to the long range sequence modeling. Our proposed method achieves improvements by average 2.77 points (55.01% vs. 57.78%) on five long sequence tasks of LRA. We also conduct experiments on 3 long document classification and 2 question answering tasks using pretraining and finetuning paradigm, further shows the effectiveness of ERNIE-SPARSE. Extensive experiments on those 10 tasks demonstrate the effectiveness and robustness of our approach in both cold start and pretrain-then-finetune paradigm.

# 2 RELATED WORK

**Sparse Transformer** Sparse attention is widely adopted to solve the long range sequence modeling problem. A simple version is (Qiu et al., 2019) that splits the sequence into blocks and perform attention only within block. This mechanism is also called local attention because only local tokens within the block can attend to each other. To improve the connection between tokens, global attention is introduced by (Child et al., 2019; Beltagy et al., 2020; Ainslie et al., 2020; Zaheer et al., 2020). The global attention mechanism mainly relies on specifying some tokens as global tokens. Those global tokens are used as a form of memory to strengthen global information. However, we observed that the mechanism of global and local attention in sparse attention would cause an information bottleneck, which would affect the information flow between local tokens and lead to the degradation of effect. This phenomenon was also observed in paper (Alon & Yahav, 2021) for Graph Neural Network (Gori et al., 2005). To solve this problem, we propose to use hierarchical mechanism which will be discussed below.

**Hierarchical Transformer** Hierarchical learning has been suggested by many researchers and it has been empirically shown to be effective in numerous diverse tasks of natural language processing (Zhang et al., 2019; Liu & Lapata, 2019; Rohde et al., 2021; Wu et al., 2021). In this paper, we propose to apply a hierarchical mechanism in Sparse Transformer and provide a new perspective from information flow to describe its beneficiary for Sparse Transformer.

**Robust Training** Although state-of-the-art deep neural networks have achieved high performance on various NLP tasks, the architectures used for these tasks have been shown to be unstable to small modifications of input texts (Ebrahimi et al., 2018; Jin et al., 2020; Sun et al., 2020a). Adversarial

training (Szegedy et al., 2014; Goodfellow et al., 2015; Miyato et al., 2019; Liu et al., 2020) is proposed to improve model robustness by constructing adversarial input. However, current robust training methods are mostly designed for dense models and don't take into account the characteristic of sparse model. Recently, (Liang et al., 2021) proposed a simple regularization strategy which forces the output distributions of different sub models generated by dropout to be consistent with each other. Inspired by this work, we propose to apply different sparse attention patterns for the same input data and encourage the model to generate the same prediction to increase Sparse Transformer's robustness.

#### 3 Method

In this section, we first revisit the mechanism of Sparse Transformer in Section 3.1, and describe the two components of sparse attention (i.e. global attention and local attention). In addition, we provide a new perspective for understanding Sparse Transformers from information flow. At the end of Section 3.1, we discuss the problems existing in the mechanism of the Sparse Transformer. Next, we introduce the motivation, formulation and benefits of the HST in Section 3.2. Finally, a robust and regularization method designed for Sparse Transformer, SAOR, is illustrated in Section 3.3.

## 3.1 REVISITING SPARSE TRANSFORMER

Given a sequence of input  $x = (t_1, ..., t_n)$  of length n, dense attention can be formulated as:

$$\begin{pmatrix} \mathbf{Q} \\ \mathbf{K} \\ \mathbf{V} \end{pmatrix} = \mathbf{H} \begin{pmatrix} \mathbf{W}_{\mathbf{q}} \\ \mathbf{W}_{\mathbf{k}} \\ \mathbf{W}_{\mathbf{v}} \end{pmatrix} + \begin{pmatrix} \mathbf{b}_{\mathbf{q}} \\ \mathbf{b}_{\mathbf{k}} \\ \mathbf{b}_{\mathbf{v}} \end{pmatrix},$$
Attention = Softmax ( $\mathbf{Q}\mathbf{K}^{T}$ )  $\mathbf{V}$ ,

where  $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{n \times d}$  is linearly mapped from  $\mathbf{H} \in \mathbb{R}^{n \times d}$ , d is the size of the hidden vector,  $\mathbf{H} = (h_1^l, ..., h_n^l)$  for layer  $l \in [1, L]$ , L is the number of model layers. For simplicity, the scale factor for  $\mathbf{Q}$  is omitted in the equation. Computation complexity for dense transformer is  $O(n^2)$ . In dense attention, tokens are connected to each other while in sparse attention, tokens are only partially connected.

As shown in Figure 1 (a), sparse attention mainly consists of global attention and local attention. In this Figure, Q, K, V are divided into 4 blocks respectively,  $(t_1, t_2 \sim t_3, t_4 \sim t_5, t_6 \sim t_7)$ . The light blue part denotes the global attention meaning the global tokens from this part can attend to all other tokens. The dark blue part express local attention denoting that the local tokens can only attend to the tokens within the same block. It can be seen from Figure 1 (b) that there is a bottleneck in the information flow of Sparse Transformer. For sparse attention mechanism, tokens in different attention blocks within the same layer are invisible to each other. For example,  $t_3$  can't attend to  $t_4$  and vice versa. For these tokens, the only way to attend to each other is through global tokens as relay nodes. As shown Figure 1 (b), the information of  $t_3$  and  $t_4$  first flows to  $t_1$  at layer l, and then distributes back to  $t_3$  and  $t_4$  at layer l + 1. This mechanism causes the current layer local token information to be carried by only a few global tokens leading to a bottleneck in the flow of information. According to our observation in section 5.1, with the information bottleneck size decreases, the performance degrades.

Another problem caused by sparse attention mechanism is susceptible to the sparse attention pattern change. As shown in Figure 2, for the sparse pattern on the left,  $t_1 \sim t_4$  and  $t_5 \sim t_8$  form two blocks of local attention blocks. For the one on the right,  $t_1 \sim t_2$ ,  $t_3 \sim t_6$ , and  $t_7 \sim t_8$  form three blocks. In other words, the change of pattern will lead to the change of attention topology between tokens. This change will result in different outputs corresponding to different patterns while dense attention does not have this problem. We can see that sparse attention is over sensitive to the pattern change.

For solving those two problems, we propose HST to improve the information flow resuling in solve the bottleneck sensitivity problem, and SAOR to make the Sparse Transformer more robust.

#### 3.2 HST: HIERARCHICAL SPARSE TRANSFORMER

HST mainly contains three key elements, i.e. representative tokens insertion, hierarchical attention and the usage for those representative tokens in the last layer for downstream task.



Figure 1: The comparison between Sparse Transformer (ST) and Hierarchical Sparse Transformer (HST). (a) Sparse Transformer mainly consists of global attention and local attention. (b) The bottleneck that existed in ST is harmful for information flow: all the sequence information is compressed into a fixed size vector. (c) In HST, representative tokens are inserted into local attention for hierarchical attention. (d) Information flow demonstration for HST: interaction between representative nodes can increase the path of global information interaction and solve the bottleneck problem in ST.

As aforementioned, the input sample consists of n tokens  $(t_1, t_2, ..., t_n)$ . We set g as the number of global tokens and w as the block size meaning the number of tokens for each local attention. We propose to insert m representative tokens into input sequence, where  $m = (n - g) \mod w$ . Those tokens are inserted to the start position of each block. Similar to BERT (Devlin et al., 2018), we use [CLS] for those representative tokens. Thus the encoder input is as follows:

$$\mathbf{H}^{0} = \underbrace{\mathbf{E}(t_{1}); \cdots; \mathbf{E}(t_{g})}_{global}; \underbrace{\mathbf{E}(r_{i}); \mathbf{E}(t_{g+w(i-1)+1}); \cdots; \mathbf{E}(t_{g+wi})}_{i-th\ local}, \tag{1}$$

where  $i \in [1, m]$ ,  $\mathbf{E}(t) \in \mathbb{R}^d$  is the embedding lookup for token t.  $\mathbf{E}(r) \in \mathbb{R}^d$  is the representation token embedding. Then we use Sparse Transformer to encode the sequence as follows:

$$\mathbf{H}_{s}^{l} = \text{SparseTransformer}\left(\left[\mathbf{H}^{l-1}\right]\right)$$

where  $\mathbf{H}_{s}^{l} \in \mathbb{R}^{n \times d}$  is the hidden output from one layer Sparse Transformer. To better interact globally, representative tokens are extracted from  $\mathbf{H}_{s}^{l}$  for dense attention as shown in Figure 1 (c). Formally, the hierarchical attention is calculated by:

$$\mathbf{R}_{s}^{l} = \left(r_{1}^{l} \cdots r_{m}^{l}\right),$$
$$\mathbf{R}^{l} = \text{Attention}\left(\mathbf{R}_{s}^{l}\right)$$

where  $\mathbf{R}_{s}^{l} \in \mathbb{R}^{m \times d}$  is the matrix of representative token's hidden states for layer  $l, r_{1}^{l} \cdots r_{m}^{l}$  are extracted from  $\mathbf{H}_{s}^{l}$ . Attention is the dense attention described in 3.1. After dense attention, these representative token hidden vectors in  $\mathbf{R}^{l}$  are distributed back to  $\mathbf{H}_{s}^{l}$  so we get the final hidden vectors  $\mathbf{H}^{l}$  of *l*-th layer. The whole process can be seen more clearly in Figure 1 (c), after hierarchical attention (green matrix), the green dots (denoting representative tokens) are distributed back to the list of tokens. The information flow of those two steps attention is shown in Figure 1 (d) showing that richer global information interaction path are created. For example,  $t_{3}$  and  $t_{5}$  can complete an interaction by representative nodes in addition to global tokens. As this method is only related to the number of local attention blocks m and not affected by bottleneck size, the bottleneck sensitivity problem is solved.

Note that the Attention module will introduce additional weights  $W_q$ ,  $W_k$  and  $W_v$  mentioned in Section 3.1, so how to initialize these three weights needs to be discussed. One option is to randomly initialize these three weights, or we can use the weight of the SparseTransformer to warm start. During training, we should also consider whether to share the weights in Attention with SparseTransformer. Those details will be discussed in Experiment section.

For the downstream task, it is efficient to use the representative tokens as follow:

$$\mathbf{O}^{L} = \text{Pooling} \left( \mathbf{R}^{L} \right),$$
$$\mathcal{P} \left( y \mid x \right) = \text{Softmax}(\mathbf{O}^{L} \mathbf{W}_{o}),$$



Figure 2: The overall framework of our proposed SAOR. We take local sparse pattern for illustration. The picture shows that one input x will go through the model twice and obtain two distributions, while the left sparse pattern is default and right one shows a rolled version sparse pattern.

where Pooling can be MAX, MEAN,  $\mathbf{O}^{L} \in \mathbb{R}^{d}$ ,  $\mathbf{W}_{o} \in \mathbb{R}^{d \times c}$ , c is the downstream task class number. y is the downstream task label,  $\mathcal{P}(y \mid x)$  denotes the predicted probability. Note that for operator Pooling, in addition to pooling representative tokens,  $\text{CLS}_{g}$  can also be used as  $\mathbf{O}^{L}$  where  $\text{CLS}_{g}$  is the first global token.

#### 3.3 SAOR: SPARSE-ATTENTION-ORIENTED REGULARIZATION

Different with dense attention, sparse attention doesn't guarantee interaction between all tokens as the design of sparse attention is to allow tokens can only attend to some of other tokens. This leads to a problem with robustness: different sparse patterns for the same input can affect the output of the attention module. This is a problem that dense attention does not have. In other word, the Sparse Transformer is susceptible to the sparse attention pattern change. To this end, we introduce a regularization method to improve the Sparse Transformer's robustness.

Our method is designed based on the motivation that for the same input, the output of different patterns should be the same. Specifically, we force the model output of the default sparse attention transformer and shifted-window version transformer output to be consistent with each other. As shown in Figure 2, the left shows a L layers Sparse Transformer with the default sparse attention, and the right shows a transformer with shifted attention pattern. The right pattern is shifted two tokens to the upper left compared to the default pattern. We call this method Sparse-Attention-Oriented Regularization (SAOR). Concretely, given the input data x at each training step, we feed x to go through the forward pass of the network twice with different sparse attention patterns. Therefore, we can obtain two distributions of the model predictions, denoted  $\mathcal{P}_1(y \mid x)$  and  $\mathcal{P}_2(y \mid x)$  in a model. Thus the distributions of  $\mathcal{P}_1(y \mid x)$  and  $\mathcal{P}_2(y \mid x)$  are different for the same input data pair (x, y). Then at this training step, in order for those two distributions to be close to each other, our SAOR method tries to minimize the bidirectional Kullback-Leibler (KL) divergence between these two output distributions for the same sample (x, y). Formally speaking,

$$\mathcal{L}_{SAOR} = \frac{1}{2} \left[ \mathcal{D}_{KL} \left( \mathcal{P}_1 \left( y \mid x \right) \| \mathcal{P}_2 \left( y \mid x \right) \right) + \mathcal{D}_{KL} \left( \mathcal{P}_2 \left( y \mid x \right) \| \mathcal{P}_1 \left( y \mid x \right) \right) \right].$$

Assume that the learning objective of target task is negative log-likelihood, the loss of the two forward passes is:

$$\mathcal{L}_{NLL} = -\log \mathcal{P}_1(y \mid x) - \log \mathcal{P}_2(y \mid x),$$

as a result, the total loss should be :

$$\mathcal{L} = \mathcal{L}_{NLL} + \alpha \mathcal{L}_{SAOR},$$

where  $\alpha$  is the loss coefficient of  $\mathcal{L}_{SAOR}$ . Our method naturally makes transformer outputs of different attention patterns close to each other. From this point of view, adding the regularization SAOR to the final loss will make Sparse Transformer more robust.

In addition, since changing the sparse pattern is complicated in practical implementation, we recommend a simple and fast implementation method to achieve the effect of shifting attention pattern approximately. To demonstrate the quick trick, we take the transformers in Figure 2 as an example. For simulating the right sparse pattern in Figure 2, we only need to roll the model input by two tokens and keep the sparse pattern as same as the default left sparse attention pattern, e.g. we roll  $x = (t_1, t_2, \dots, t_6, t_7, t_8)$  to  $x' = (t_7, t_8, t_1, t_2, \dots, t_6)$ . Then we get the corresponding model outputs for x and x' and we can regularize them by making the outputs distance closer.

Models	ListOps	Text	Retrieval	Image	Pathfinder	Avg
Local Attention	15.82	52.98	53.39	41.46	66.63	46.06
Linear Trans. (Katharopoulos et al., 2020)	16.13	65.90	53.09	42.34	75.30	50.55
Reformer (Kitaev et al., 2020)	<u>37.27</u>	56.10	53.40	38.07	68.50	50.67
Sparse Trans. (Child et al., 2019)	17.07	63.58	<u>59.59</u>	44.24	71.71	51.24
Sinkhorn Trans.(Tay et al., 2020b)	33.67	61.20	53.83	41.23	67.45	51.39
Linformer (Wang et al., 2020)	35.70	53.94	52.27	38.56	76.34	51.36
Performer (Choromanski et al., 2021)	18.01	<u>65.40</u>	53.82	42.77	77.05	51.41
Synthesizer (Tay et al., 2020a)	36.99	61.68	54.67	41.61	69.45	52.88
Longformer (Beltagy et al., 2020)	35.63	62.85	56.89	42.22	69.71	53.46
Transformer (Vaswani et al., 2017)	36.37	64.27	57.46	42.44	71.40	54.39
BigBird (Zaheer et al., 2020)	36.05	64.02	59.29	40.83	74.87	<u>55.01</u>
ERNIE-Sparse	37.75	64.47	62.64	45.28	78.77	57.78

Table 1: Experimental results on the long range arena (LRA) benchmark. The highest score for each dataset is highlighted in bold and the second place is underlined.

## 4 EXPERIMENTS

To evaluate our approach and show its performance, we first conduct experiments on a long context sequence modeling benchmark, LRA, which is consisted of 5 multi-modal tasks including logical inference, natural language and image tasks. LRA is a benchmark for cold start models and does not require the model to be pre-trained, so LRA is a suitable benchmark for testing the model structure designs. To further evaluate the ability of ERNIE-SPARSE in the pretrain-then-finetune paradigm, we also follow (Beltagy et al., 2020) and (Zaheer et al., 2020) to pretrain ERNIE-SPARSE and test the pretrained ERNIE-SPARSE on 3 natural language classification and 2 question answering tasks.

#### 4.1 LONG-CONTEXT SEQUENCE MODELING

We first evaluate the effectiveness and efficiency of ERNIE-SPARSE on the LRA benchmark recently introduced by (Tay et al., 2021), which is designed to evaluate efficient transformer models under the long-context scenario. This multi-modal dataset contains two image tasks (Krizhevsky et al., 2009; Linsley et al., 2018), two text-based tasks (Maas et al., 2011; Radev et al., 2013), and one mathematical operation inference task (Nangia & Bowman, 2018). At the same time, LRA also set a robust baseline, in which the hyperparameters and the model were fixed. We run each experiment five times with different random seeds and report the average accuracy.

The result of ERNIE-SPARSE on the LRA tasks are reported in Table 1. First, we note that ERNIE-SPARSE achieves strong results on all tasks consistently compared to the transformer model and significantly outperforms all the other baseline methods and achieve best score in terms of the average accuracy. By taking a closer look at the accuracy for each task, ERNIE-SPARSE wins over baseline models on four out of five tasks. Notably, ERNIE-SPARSE can work well on both image and text and the math inference data sets.

## 4.2 PRETRAINING AND FINETUNING

#### 4.2.1 PRETRAINING

In the training task, we follow (Liu et al., 2019) and pretrain ERNIE-SPARSE using the Mask Language Model (MLM) training object. This task involves predicting tokens that are randomly masked out. For the training samples, we train ERNIE-SPARSE with maximum sequence length of 4096 and train it for 1 million steps. Samples with sequence length less than 4096 will be concatenated as one sample to improve training efficiency. For those samples longer than max sequence length, we truncate them to 4096. Statistics of pretraining data can be found in A.1. Following (Beltagy et al., 2020), ERNIE-SPARSE warm-starts from the public RoBERTa checkpoint and we compare the performance in MLM task in terms of bits per character (BPC). As shown in Table 3, BigBird, Longformer, ERNIE-SPARSE are all better than RoBERTa whose max sequence length is 512. Among those methods, ERNIE-SPARSE performs best.

Models	Arxiv <sup>1</sup>	IMDB	Hyperpartisan
	F1	Acc	F1
RoBERTa	86.86	95.00±0.2	87.80±0.8
BigBird	87.50	95.20±0.2	92.20±1.7
ERNIE-Sparse	89.05	95.53±0.1	92.81±1.4

Table 2: Performance of various models on development set of benchmark natural language understanding tasks.

Models

Longformer BigBird

Reproduce

Longformer

**ERNIE-SPARSE** 

BigBird

Setting	BPC
RoBERTa (Liu et al., 2019)	1.846
Longformer (Beltagy et al., 2020)	1.705
BigBird (Zaheer et al., 2020)	1.678
ERNIE-Sparse	1.674

Table 3: MLM BPC for ERNIE-SPARSE and other models.

Table 4: Model comparison for WikiHop and TriviaQA.

WikiHop

Acc 75.0

75.9

75.2

72.3

75.8

**TriviaQA** 

EM

67.0

68.6

71.8

F1

75.2

79.5

75.0

73.4

76.5

## 4.2.2 TEXT CLASSIFICATION

To test ERNIE-SPARSE on downstream tasks, we first select three text classification tasks: Arxiv paper categories classification (He et al., 2019), IMDB reviews classification (Maas et al., 2011) and Hyperpartisan news detection (Kiesel et al., 2019). The experiment was repeated 5 times for both datasets, and the mean and standard deviation were listed in the table. ERNIE-SPARSE's hyperparameters are recorded in the appendix A.2.2. Note that all linear transformation weights of hierarchical attention are shared with the weights of the previous Sparse Transformer attention. Table 2 summarizes the results of ERNIE-SPARSE. From this table, it shows that ERNIE-SPARSE surpasses all baselines on the text classification datasets. For Arxiv, ERNIE-SPARSE surpasses baseline by a large margin. For IMDB and Hyperpartisan, the performance gain continues demonstrating that ERNIE-SPARSE is capable of utilizing information from long document input.

## 4.2.3 QUESTION ANSWERING

For QA tasks, we choose WikiHop (Welbl et al., 2018) and TriviaQA (Joshi et al., 2017). In ERNIE-SPARSE model, following (Beltagy et al., 2020), we concatenate the answer / question / candidates with special tokens along with the context. As the global tokens is important for QA tasks, the use SOAR requires careful control over which tokens should be shifted. For implementation, we only roll tokens in local attention to create samples to avoid disturbing the global information in the front of sentence.

The results of WikiHop and TriviaQA are shown in Table 4. The first two rows are copied from (Beltagy et al., 2020) and (Zaheer et al., 2020). We have reproduced the results and record them in the third and fourth row. Note that we first tried to reproduce Longformer's score. It can be seen from this table that the score reproduced is basically the same as the score in the (Beltagy et al., 2020). Secondly, we tried to directly replace the warm start model with BigBird's checkpoint and tuned the hyperparameters to get the best reproduce result and the score was recorded in the table. From this table, we see that ERNIE-SPARSE surpasses the baseline and achieves the best results.

## 4.3 ABLATION

**Effect of proposed components: HST and SAOR** Table 5 shows the performance of ERNIE-SPARSE by ablating each proposed component. All models are reported with mean results of 5 runs,

<sup>&</sup>lt;sup>1</sup>The reported numbers from (Zaheer et al., 2020) on Arxiv dataset are not comparable with ours because they did not release the train/test split of the data. So we re-tested RoBERTa and BigBird's open source checkpoint on our Arxiv train/test split. We open source our Arxiv dataset split: https://github.com/anonymous

Models	ListOps	Text	Retrieval	Image	Pathfinder	Avg
#0 ERNIE-SPARSE	37.75	64.47	62.64	45.28	78.77	57.60
#1 w/o HST	35.97	63.25	62.24	42.76	78.79	56.60
#2 w/o SAOR	37.36	63.54	61.87	42.77	75.25	56.16
#3 w/o SAOR + ws in HST	37.41	62.55	60.32	42.62	76.55	55.89
#4 w/o SAOR + R-Drop	37.68	61.73	60.59	43.48	78.11	56.32

Table 5: Performance of ERNIE-SPARSE by ablating each proposed components. **ws** means the linear weights for hierarchical and sparse attention are shared.

the hyperparameters keep the same as the official recommendation. From the last column in Table 5, we see that the HST improves ERNIE-SPARSE with 1.0 percent point on average (#1 vs. #0). As the default setting is weight not sharing for linear mapping weight in hierarchical attention with the sparse attention layer, we add a group of experiment #3 to explore the effect of weight sharing on experimental results. We see that with weight sharing, the results drop by 0.27 points on average by (#3 vs #2). By comparing #2 to #0, we see that SAOR is beneficial in modeling the long sequence bringing 1.44 points improvement on average. Except for those observation on average scores, let us take a closer to the how HST and SAOR perform in different types of data. Since HST is expected to be more useful on tasks that require global information, let's look at ListOps, Text and Image first. It can be seen that after removing HST, the average decrease is 1-2 points (ListOps, Text and Image in #1 vs. #0). Especially, for ListOps, the task need the model to see all characters in order to do valid logical inference, HST have shown its importance (35.97 vs. 37.75). Except for the global information, some tasks require the model has more stability. For example, the image task Pathfinder, dropped 3.5 points after removing SAOR (75.25 vs. 78.77) showing the effect of SAOR. Except for those component discussed above, the relation between SAOR and R-Drop (#4) will be discuss later.

Effect of pooling method for representative tokens In the experiment, we have explored three different pooling methods while keeping other settings unchanged. The results are shown in Table 6. MEAN and MAX represent mean pooling and max pooling, respectively.  $CLS_g$ -Only refers to use the first token CLS only. As presented in Table 6,  $CLS_g$ -Only is the worst on average score, and MEAN performs best. Take a a closer look at the score each task, it shows that MAX performs best at Image and MEAN performs best at Retrieval, indicating that MEAN is

Setting	Image Acc	Retrieval Acc	Avg
ERNIE-SPARSE (MEAN)	45.28	62.64	53.96
ERNIE-SPARSE (MAX)	46.51	58.98	52.75
$ERNIE\text{-}SPARSE (CLS_g\text{-}Only)$	42.88	60.03	50.75

Table 6: Ablation for the pooling method of representative tokens in ERNIE-SPARSE for downstream tasks.

a stable method for predicting strong results and for image tasks, the MAX can be considered as a candidate for hyperparameters. For CLS, the score drops significantly at Image task and maintain a good performance on Retrieval, indicating that the contextual global information is more important for image tasks.

Effect of SAOR vs. R-Drop In this study, we specifically investigate the importance of Dropout (Srivastava et al., 2014) in those experiments. As the Dropout technique is commonly used in deep neural network training and (Liang et al., 2021) use Dropout to constrain the outputs of two subnetworks, we ablate Dropout to compare SAOR and R-Drop (Liang et al., 2021) as shown by (#4 vs. #0) in Table 5. #1 is the SAOR-only version ERNIE-SPARSE result, which means for getting the regularization version model output  $\mathcal{P}_2(y \mid x)$  in 3.3, we roll the input (which is equivalent to shift sparse pattern) and use the Dropout at the same time, #4 means that we only use Dropout to get the regularization version model output  $\mathcal{P}_2(y \mid x)$ . We see that SAOR achieve the best for the average score (#4 vs. #0). Moreover, for specific task in Table 5, our method is more efficient in 3 out of 5, indicating not only SAOR's effectiveness, but also that SAOR and Dropout are compatible and complementary.

#### 5 DISCUSSION

#### 5.1 BOTTLENECK ANALYSIS

In this section, we analyze the impact of bottleneck size on performance of Sparse Transformer and HST respectively. All the hyperparameter configuration follows (Tay et al., 2020c). We ran each experiment twice and average the results. The experimental results can be seen in Figure 3.



Figure 3: Phenomenon of bottleneck and the effectiveness of **HST**: Accuracy across global token size in the LRA benchmark. Performance degradation of ST are observed on all datasets when bottleneck size decrease (blue line). With our proposed HST method, the performance is better (red line). Error bar denotes standard deviation.

Models	ListOps	ListOps Adv	Text	Text Adv	Retrieval	Retrieval Adv	Image	Image Adv	Pathfinder	Pathfinder Adv
ERNIE-SPARSE w/o SAOR	<b>37.40</b> 37.30	<b>18.90</b> 9.65	<b>63.14</b> 62.54	<b>63.17</b> 62.41	<b>61.89</b> 61.84	<b>61.80</b> 61.52	<b>44.82</b> 43.59	<b>42.78</b> 40.12	<b>78.22</b> 76.39	<b>67.29</b> 51.00

Table 7: Performance of w/ and w/o **SAOR**: We construct an adversarial dataset by shift word to test the robustness of ERNIE-SPARSE and the sensitivity for sparse pattern shift.

As discussed in section 3.1, the number of global tokens determines the size of a bottleneck in Sparse Transformer and we recorded the trend of Sparse Transformer and HST in LRA dataset by changing the size of bottleneck. In this figure, blue and red line denotes Sparse Transformer and HST respectively. As shown in Figure 3, as the number of global tokens decreases, the score of Sparse Transformer decreases (blue line in each subfigure). This trend indicates that in Sparse Transformer, the global token causes a bottleneck for information flow, which in turn makes the model sensitive to the bottleneck size. For our method, HST can avoid the bottleneck of Sparse Transformer by introducing hierarchical attention and using representative token for global information interaction. As can be seen in this figure, the task score of HST (red line) in each subfigure can be maintained in the same interval without a downward trend, which proves the effectiveness of HST proposed by us.

#### 5.2 SYNTHETIC DATASET: WORD ORDER SHIFT ATTACK

As mentioned in 3.1, in addition to bottleneck, another problem we found affecting the robustness of the Sparse Transformer is the sensitivity to sparse patterns. To further demonstrate this phenomenon, we create an adversarial dataset to observe the performance of Sparse Transformer and ERNIE-SPARSE. Concretely, we constructed a adversarial version of the LRA by rolling tokens and we consider the roll(x) operation as an attack method named Word Order Shift Attack for the Sparse Transformer. As shown in Table 7, all datasets with Adv are corresponding adversarial version datasets, for example, ListOps Adv is an adversarial dataset for ListOps. Firstly, it was observed that all the models in Table 7 show different degrees of decline on the adversarial dataset, which verifies that Sparse Transformer is sensitive to this kind of attack. Secondly, by comparing the result on adversarial dataset, we can see that ERNIE-SPARSE surpasses the w/o SAOR version by a large margin, which indicates that SAOR method can improve the robustness of the model and reduce the sensitivity of the model to Word Order Shift Attack.

#### 6 CONCLUSION

In this paper, we propose ERNIE-SPARSE, a robust training method for Sparse Transformer. It contains two key features: HST mechanism and SAOR regularization method. Firstly, the advantages and disadvantages of Sparse Transformer are analyzed from the perspective of information flow, and then the over-sensitivity of Sparse Transformer to bottleneck is found. Then representative token is introduced to solve the Sparse Transformer's dependence on bottleneck and improve the performance of Sparse Transformer. On the other hand, we also find that the Sparse Transformer is sensitive to the sparse pattern of attention, and then we convert this sensitivity into a regularization method to enhance model robustness. At last, ERNIE-SPARSE, outperforms a variety of strong baseline methods including the full-rank attention and other efficient sparse and dense attention methods. Moreover, we also conducted pretraining and finetune experiments, which are validated in two downstream benchmarks, text classification and QA. The results also prove the validity of ERNIE-SPARSE.

#### REFERENCES

- Joshua Ainslie, Santiago Ontañón, Chris Alberti, Philip Pham, Anirudh Ravula, and Sumit Sanghai. Etc: encoding long and structured data in transformers. *arXiv e-prints*, pp. arXiv–2004, 2020.
- Uri Alon and Eran Yahav. On the bottleneck of graph neural networks and its practical implications. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=i800Ph0CVH2.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- Krzysztof Marcin Choromanski, Valerii Likhosherstov, David Dohan, Xingyou Song, Andreea Gane, Tamás Sarlós, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser, David Benjamin Belanger, Lucy J. Colwell, and Adrian Weller. Rethinking attention with performers. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021. URL https://openreview.net/forum?id= Ua6zuk0WRH.
- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. ELECTRA: pre-training text encoders as discriminators rather than generators. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net, 2020. URL https://openreview.net/forum?id=rlxMH1BtvB.
- 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.
- Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. Hotflip: White-box adversarial examples for text classification. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 2: Short Papers*, pp. 31–36. Association for Computational Linguistics, 2018. doi: 10.18653/v1/P18-2006. URL https://aclanthology.org/P18-2006/.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In Yoshua Bengio and Yann LeCun (eds.), 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015. URL http://arxiv.org/abs/1412.6572.
- Marco Gori, Gabriele Monfardini, and Franco Scarselli. A new model for learning in graph domains. In *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, volume 2, pp. 729–734. IEEE, 2005.
- 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. doi: 10.1109/ACCESS.2019.2907992. URL https://doi.org/10.1109/ACCESS.2019.2907992.
- Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. Is BERT really robust? A strong baseline for natural language attack on text classification and entailment. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pp. 8018– 8025. AAAI Press, 2020. URL https://aaai.org/ojs/index.php/AAAI/article/ view/6311.
- Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In Regina Barzilay and Min-Yen Kan (eds.), *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 August 4, Volume 1: Long Papers*, pp. 1601–1611. Association for Computational Linguistics, 2017. doi: 10.18653/v1/P17-1147. URL https://doi.org/10.18653/v1/P17-1147.

- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 5156–5165. PMLR, 2020. URL http: //proceedings.mlr.press/v119/katharopoulos20a.html.
- Johannes Kiesel, Maria Mestre, Rishabh Shukla, Emmanuel Vincent, Payam Adineh, David P. A. Corney, Benno Stein, and Martin Potthast. Semeval-2019 task 4: Hyperpartisan news detection. In Jonathan May, Ekaterina Shutova, Aurélie Herbelot, Xiaodan Zhu, Marianna Apidianaki, and Saif M. Mohammad (eds.), *Proceedings of the 13th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2019, Minneapolis, MN, USA, June 6-7, 2019*, pp. 829–839. Association for Computational Linguistics, 2019. doi: 10.18653/v1/s19-2145. URL https://doi.org/10.18653/v1/s19-2145.
- Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020.* OpenReview.net, 2020. URL https://openreview.net/forum?id=rkgNKkHtvB.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Technical report*, 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25:1097–1105, 2012.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. ALBERT: A lite BERT for self-supervised learning of language representations. *CoRR*, abs/1909.11942, 2019. URL http://arxiv.org/abs/1909.11942.
- Xiaobo Liang, Lijun Wu, Juntao Li, Yue Wang, Qi Meng, Tao Qin, Wei Chen, Min Zhang, and Tie-Yan Liu. R-drop: Regularized dropout for neural networks, 2021.
- Drew Linsley, Junkyung Kim, Vijay Veerabadran, Charles Windolf, and Thomas Serre. Learning long-range spatial dependencies with horizontal gated recurrent units. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, pp. 152–164, 2018. URL https://proceedings.neurips.cc/paper/2018/hash/ ec8956637a99787bd197eacd77acce5e-Abstract.html.
- Xiaodong Liu, Hao Cheng, Pengcheng He, Weizhu Chen, Yu Wang, Hoifung Poon, and Jianfeng Gao. Adversarial training for large neural language models. *CoRR*, abs/2004.08994, 2020. URL https://arxiv.org/abs/2004.08994.
- Yang Liu and Mirella Lapata. Hierarchical transformers for multi-document summarization. *arXiv* preprint arXiv:1905.13164, 2019.
- 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.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In Dekang Lin, Yuji Matsumoto, and Rada Mihalcea (eds.), *The 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference, 19-24 June, 2011, Portland, Oregon, USA*, pp. 142–150. The Association for Computer Linguistics, 2011. URL https: //aclanthology.org/P11-1015/.
- Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: A regularization method for supervised and semi-supervised learning. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(8):1979–1993, 2019. doi: 10.1109/TPAMI.2018.2858821. URL https: //doi.org/10.1109/TPAMI.2018.2858821.

- Nikita Nangia and Samuel R. Bowman. Listops: A diagnostic dataset for latent tree learning. In Silvio Ricardo Cordeiro, Shereen Oraby, Umashanthi Pavalanathan, and Kyeongmin Rim (eds.), Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 2-4, 2018, Student Research Workshop, pp. 92–99. Association for Computational Linguistics, 2018. doi: 10.18653/v1/n18-4013. URL https://doi.org/10.18653/v1/n18-4013.
- Jiezhong Qiu, Hao Ma, Omer Levy, Scott Wen-tau Yih, Sinong Wang, and Jie Tang. Blockwise self-attention for long document understanding. *arXiv preprint arXiv:1911.02972*, 2019.
- Dragomir R. Radev, Pradeep Muthukrishnan, Vahed Qazvinian, and Amjad Abu-Jbara. The ACL anthology network corpus. *Lang. Resour. Evaluation*, 47(4):919–944, 2013. doi: 10.1007/s10579-012-9211-2. URL https://doi.org/10.1007/s10579-012-9211-2.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. J. Mach. Learn. Res., 21:140:1–140:67, 2020. URL http://jmlr.org/papers/v21/20-074.html.
- Tobias Rohde, Xiaoxia Wu, and Yinhan Liu. Hierarchical learning for generation with long source sequences. *arXiv preprint arXiv:2104.07545*, 2021.
- 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.
- Lichao Sun, Kazuma Hashimoto, Wenpeng Yin, Akari Asai, Jia Li, Philip S. Yu, and Caiming Xiong. Adv-bert: BERT is not robust on misspellings! generating nature adversarial samples on BERT. *CoRR*, abs/2003.04985, 2020a. URL https://arxiv.org/abs/2003.04985.
- Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Hao Tian, Hua Wu, and Haifeng Wang. Ernie 2.0: A continual pre-training framework for language understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 8968–8975, 2020b.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In Yoshua Bengio and Yann LeCun (eds.), 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings, 2014. URL http://arxiv.org/abs/ 1312.6199.
- Yi Tay, Dara Bahri, Donald Metzler, Da-Cheng Juan, Zhe Zhao, and Che Zheng. Synthesizer: Rethinking self-attention in transformer models. *CoRR*, abs/2005.00743, 2020a. URL https: //arxiv.org/abs/2005.00743.
- Yi Tay, Dara Bahri, Liu Yang, Donald Metzler, and Da-Cheng Juan. Sparse sinkhorn attention. In Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event, volume 119 of Proceedings of Machine Learning Research, pp. 9438–9447. PMLR, 2020b. URL http://proceedings.mlr.press/v119/tay20a.html.
- Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena: A benchmark for efficient transformers. *arXiv preprint arXiv:2011.04006*, 2020c.
- Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena : A benchmark for efficient transformers. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021. URL https://openreview.net/ forum?id=qVyeW-grC2k.
- Trieu H. Trinh and Quoc V. Le. A simple method for commonsense reasoning. *CoRR*, abs/1806.02847, 2018. URL http://arxiv.org/abs/1806.02847.

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pp. 5998–6008, 2017. URL https://proceedings.neurips.cc/paper/2017/hash/ 3f5ee243547dee91fbd053c1c4a845aa-Abstract.html.
- Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. CoRR, abs/2006.04768, 2020. URL https://arxiv.org/abs/2006. 04768.
- Johannes Welbl, Pontus Stenetorp, and Sebastian Riedel. Constructing datasets for multi-hop reading comprehension across documents. *Trans. Assoc. Comput. Linguistics*, 6:287–302, 2018. URL https://transacl.org/ojs/index.php/tacl/article/view/1325.
- Chuhan Wu, Fangzhao Wu, Tao Qi, and Yongfeng Huang. Hi-transformer: Hierarchical interactive transformer for efficient and effective long document modeling. *arXiv preprint arXiv:2106.01040*, 2021.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32, 2019.
- Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for longer sequences. *Advances in Neural Information Processing Systems*, 33, 2020.
- Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. Defending against neural fake news. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pp. 9051–9062, 2019. URL https://proceedings.neurips.cc/paper/2019/hash/ 3e9f0fc9b2f89e043bc6233994dfcf76-Abstract.html.
- Xingxing Zhang, Furu Wei, and Ming Zhou. Hibert: Document level pre-training of hierarchical bidirectional transformers for document summarization. *arXiv preprint arXiv:1905.06566*, 2019.
- Yukun Zhu, Ryan Kiros, Richard S. Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In 2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015, pp. 19–27. IEEE Computer Society, 2015. doi: 10.1109/ICCV.2015.11. URL https://doi.org/10.1109/ICCV.2015.11.

## A APPENDIX

## A.1 PRETRAINING

#### A.1.1 PRETRAINING DATASET

For ERNIE-SPARSE pretraining, we use Wikipedia (English Wikipedia dump<sup>2</sup>; 12GB), BookCorpus (Zhu et al., 2015) (4.6GB), Realnews (Zellers et al., 2019) (7.4GB) and Stories (Trinh & Le, 2018) (11GB). For pretraining, we also sample 5% training data as the validation set to monitor the training process. Table 8 shows statistics of the pretraining data.

Source	Tokens	Avg doc len
Wikipedia	3.0B	515
BookCorpus	1.2B	23K
Realnews	1.8B	3.0K
Stories	2.7B	8.7K

 Table 8: Pretraining data statistics

Parameter	ERNIE-SPARSE
$\alpha$ of $L_{SAOR}$	0
learning rate	3e-5
batch size	256
weight decay	0.1
warmup steps	10k
total steps	1m
max seq length	4096
embedding dim	768
#head	12
#layer	12
activation layer	gelu
dropout	0.1
attn dropout	0.1

Table 9: Hyperparameters for the ERNIE-SPARSE for Pretraining

#### A.1.2 PRETRAINING HYPERPARAMETERS

We split any document longer than 4096 into multiple documents and we joined multiple documents that were much smaller than 4096. During the pre-training phase, we only use mask language model for training tasks. Specifically, we mask 15% of tokens in these four datasets, and train ERNIE-SPARSE to predict the mask. We warm start ERNIE-SPARSE from RoBERTa's checkpoint. The hyperparameters for these ERNIE-SPARSE are given in Table 9. We use a learning rate warmup over the first 10,000 steps, and polynomial decay of the learning rate. Notably, attention weight in HST are shared with sparse attention.

## A.2 TASKS

To evaluate ERNIE-SPARSE, we chose three benchmarks, including LRA, and text classification, as well as question answering. The latter two need to follow the pretraining and finetuning paradigm. Table 10 lists the data distribution, task type, evaluation metric of each dataset.

## A.2.1 HYPERPARAMETERS FOR LRA

Table 11 gives the detail list of hyperparameters used to get results shown in Table 1.

## A.2.2 HYPERPARAMETERS FOR CLASSIFICATION AND QA

Table 12 gives the detail list of hyperparameters used to get results shown in Table 2 and Table 4.

Comus	Toolz	Split	#Somplo	Length in percentile				#Labol	Motrico
Corpus	Task	Spiit	manple "Sample	50%	90%	95%	max	"Laber	withins
			Long Ra	nge Aren	a (LRA)				
		Train	96k	954	1646	1800	1999		Acc
ListOps	Logical Reasoning	Dev	2k	960	1648	1813	1999	10	
		Test	2k	947	1657	1803	1999		
Text	Sentiment Classification	Train	25k	979	2615	3431	13704	2	Acc
тел	Sentiment Classification	Dev	25k	962	2543	3333	12988	2	All
		Train	147k	7648	13467	20495	72885		
Retrieval	Retrieval Retrieval	Dev	18k	7665	13359	19928	72885	2	Acc
		Test	17k	7702	15955	22427	50012		
		Train	45k	-	-	-	1024		
Image	Category Classification	Dev	5k	-	-	-	1024	10	Acc
		Test	10k	-	-	-	1024		
	Image Reasoning	Train	160k	-	-	-	1024		
PathFinder		Dev	20k	-	-	-	1024	2	Acc
		Test	20k	-	-	-	1024		
-			Text	Classifica	ation		•		
Amiri	Cotogory Classification	Train	33k	14733	34209	43951	1121751	11	Miono E1
AIXIV	Category Classification	Test	3.3k	14710	32417	40965	850540	11	MICTO F1
IMDD	Santimant Classification	Train	25k	215	569	748	3084	2	Miono E1
IIVIDD	Sentiment Classification	Test	25k	212	550	724	2778	2	MICTO F1
		Train	516	536	1517	1990	5560		
Hyperpartisan	News Classification	Dev	65	520	1535	1971	2637	2	Micro F1
		Test	65	637	1771	1990	5560		
	ł		Quest	ion Answ	ering	1			
TriviaOA	Question Answering	Train	110k	4576	5027	5166	10091	Span	Macro F1 & EM
111110211		Dev	14k	4577	5026	5169	10210	Span	
WikiHon	Question Answering	Train	43k	1313	3001	3685	19747	Candidates	Acc
,, ikii iop		Dev	5.1k	1413	3184	3871	17004	Canuluates	

Table 10: Downstream tasks statistics. Samples of tasks Image and PathFinder are all  $32 \times 32$  images.

Hyperparameter	ListOps	Text	Retrieval	Image	Pathfinder				
Hyperparameters for H	IST and SAO	R							
HST pooling	{ MIN, MEAN, MAX}								
$\alpha$ of $\mathcal{L}_{SAOR}$		{	$\{0.5, 5, 10\}$						
#roll tokens of $\mathcal{L}_{SAOR}$			$\{2, 8, 16\}$						
Fixed hyperparameters provided by LRA (Tay et al., 2021)									
learning rate	5e-2	5e-2	5e-2	5e-4	1e-3				
batch size	32	32	32	256	512				
weight decay	1e-1	1e-1	1e-1	0	0				
warmup	1000	8000	8000	175	312				
max seq length	2000	1000	4000	1024	1024				
embedding dim	512	256	128	32	64				
#head	8	4	4	1	2				
#layer	4	4	4	1	4				
Q/K/V dim	512	256	128	32	32				
MLP dim	1024	1024	512	64	64				
dropout	0.1	0.1	0.1	0.3	0.2				
attn dropout	0.1	0.1	0.1	0.2	0.1				
lr decay	root square	root square	root square	cosine	cosine				

Table 11: The upper part is the hyperparameter related to ERNIE-SPARSE, while the lower part is the fixed hyperparameter provided by LRA and cannot be changed.

#### A.3 COMPLEXITY ANALYSIS

In this section, we analyze the complexity of ERNIE-SPARSE. Note that the two techniques of ERNIE-SPARSE, HST and SAOR, are general and can be applied with any type of sparse patterns. Considering that only HST of the two techniques affects time complexity, we will only test the performance of HST here. In order to test the complexity, we will follow the following process: firstly we select a sparse attention and record this sparse pattern's performance and secondly we combine

<sup>&</sup>lt;sup>2</sup>https://dumps.wikimedia.org/enwiki/

Hyperparameter	Arxiv	IMDB	Hyperpartisan	WikiHop	TriviaQA
HST pooling	mean	mean	mean	mean	mean
$\alpha \text{ of } \mathcal{L}_{SAOR}$	10	0.1	0	10	3
#roll tokens of $\mathcal{L}_{SAOR}$	8	8	0	8	8
learning rate	6e-5	1e-5	3e-5	3e-5	3e-5
batch size	48	64	16	48	32
epoch	10	20	20	30	10
warmup	10%	10%	10%	200 (steps)	10%
max seq len	4096	2048	1024	4096	4096
#global token			128		
local window size			192		
#random token			192		
Optimizer			Adam		

Table 12: Hyperparameters of classification and question answering tasks.



Figure 4: Runtime and memory of full self-attention, ERNIE-SPARSE and BigBird (Zaheer et al., 2020).

the HST with this sparse pattern to show the adding computation complexity. For the sparse pattern we choose to use the recently proposed pattern from BigBird (Zaheer et al., 2020) for comparison. And then we test ERNIE-SPARSE with same test settings. For those comparison, we carry out the experiment on a V100 32GB GPU. The result is shown in Figure 4. ERNIE-SPARSE 's memory usage scales linearly with the sequence length, unlike the full self-attention mechanism that runs out of memory for long sequences on current GPU. ERNIE-SPARSE and BigBird are almost on par. The speed test data shows that the performance is consistent. The difference between ERNIE-SPARSE and BigBird is negligible in memory test. Given that ERNIE-SPARSE outperformed BigBird on the test set as mentioned in Section 4, the performance of ERNIE-SPARSE was remarkable.