# Pre-Training with Whole Word Masking for Chinese BERT

Yiming Cui<sup>†‡</sup>, Wanxiang Che<sup>†</sup>, Ting Liu<sup>†</sup>, Bing Qin<sup>†</sup>, Ziqing Yang<sup>‡</sup>, Shijin Wang<sup>‡§</sup>, Guoping Hu<sup>‡</sup>

†Research Center for Social Computing and Information Retrieval (SCIR),

Harbin Institute of Technology, Harbin, China

<sup>‡</sup>State Key Laboratory of Cognitive Intelligence, iFLYTEK Research, China

<sup>§</sup>iFLYTEK AI Research (Hebei), Langfang, China

†{ymcui, car, tliu, qinb}@ir.hit.edu.cn

<sup>‡§</sup>{ymcui, zqyang5, sjwang3, qphu}@iflytek.com

## Abstract

Bidirectional Encoder Representations from Transformers (BERT) has shown marvelous improvements across various NLP tasks. Recently, an upgraded version of BERT has been released with Whole Word Masking (WWM), which mitigate the drawbacks of masking partial WordPiece tokens in pre-training BERT. In this technical report, we adapt whole word masking in Chinese text, that masking the whole word instead of masking Chinese characters, which could bring another challenge in Masked Language Model (MLM) pre-training task. The proposed models are verified on various NLP tasks, across sentence-level to document-level, including machine reading comprehension (CMRC 2018, DRCD, CJRC), natural language inference (XNLI), sentiment classification (ChnSentiCorp), sentence pair matching (LCQMC, BQ Corpus), and document classification (THUCNews). Experimental results on these datasets show that the whole word masking could bring another significant gain. Moreover, we also examine the effectiveness of the Chinese pre-trained models: BERT, ERNIE, BERTwwm, BERT-wwm-ext, RoBERTa-wwm-ext, and RoBERTa-wwm-ext-large. 1

## 1 Introduction

Bidirectional Encoder Representations from Transformers (BERT) (Devlin et al., 2019) has become enormously popular and proven to be effective in recent NLP studies which utilizes large-scale unlabeled training data and generates enriched contextual representations, showing its powerful performance on various natural language processing tasks. As we traverse several popular machine reading comprehension benchmarks, such as SQuAD (Rajpurkar et al., 2018), CoQA

(Reddy et al., 2019), QuAC (Choi et al., 2018), NaturalQuestions (Kwiatkowski et al., 2019), RACE (Lai et al., 2017), we can see that most of the top performing models are based on BERT and its variants (Cui et al., 2017; Dai et al., 2019; Zhang et al., 2019a; Ran et al., 2019).

Recently, the authors of BERT have released an updated version of BERT, which is called Whole **Word Masking**. The whole word masking mainly mitigates the drawbacks in original BERT that, if the masked WordPiece token (Wu et al., 2016) belongs to a whole word, then all the WordPiece tokens (which forms a complete word) will be masked altogether. This will explicitly force the model to recover the whole word in Masked Language Model (MLM) pre-training task, instead of just recovering WordPiece tokens, which is much more challenging. Along with the strategy, they also provide pre-trained English models (BERTlarge-wwm) for the community, which is beneficial for the researcher to design more powerful models based on them.

Before Devlin et al. (2019) releasing BERT with whole word masking, Sun et al. (2019) had proposed Enhanced Representation through kNowledge IntEgration (ERNIE) with a similar spirit and trained on not only Wikipedia data but also community QA, Baike (similar to Wikipedia), etc.<sup>2</sup> It was tested on various NLP tasks and showed consistent improvements over BERT.

In this technical report, we adapt the whole word masking strategy in Chinese BERT to verify its effectiveness. The model was pre-trained on the latest Wikipedia dump in Chinese (both Simplified and Traditional Chinese proportions are kept). Note that, we did not exploit additional data in our model, and aim to provide a more general base for

<sup>&</sup>lt;sup>1</sup>We release all the pre-trained models: https://github.com/ymcui/Chinese-BERT-wwm

<sup>&</sup>lt;sup>2</sup>Tsinghua university has also released a model called ERNIE (Zhang et al., 2019b) but was not trained on Chinese. In this paper, ERNIE refers to the model by Sun et al. (2019).

developing NLP systems in Simplified and Traditional Chinese. Extensive experiments are conducted on various Chinese NLP datasets, ranging from sentence-level to document-level, which include machine reading comprehension, sentiment classification, sentence pair matching, natural language inference, document classification, etc. The results show that the proposed model brings another gain over BERT and ERNIE in most of the tasks, and we provide several useful tips for using these pre-trained models, which may be helpful in the future research.

The contributions of this technical report are listed as follows.

- We adapt the whole word masking in Chinese BERT and release the pre-trained models for the community.
- Extensive experiments are carried out to better demonstrate the effectiveness of BERT, ERNIE, and BERT-wwm.
- Several useful tips are provided on using these pre-trained models on Chinese text.

# 2 Chinese BERT with Whole Word Masking

## 2.1 Methodology

We strictly follow the original whole word masking codes and did not change other components, such as the percentage of word masking, etc. An example of the whole word masking is depicted in Figure 1.

#### 2.2 Data Processing

We downloaded the latest Wikipedia dump<sup>3</sup>, and pre-processed with WikiExtractor.py as suggested by Devlin et al. (2019), resulting 1,307 extracted files. Note that, we use both Simplified and Traditional Chinese in this dump. After cleaning the raw text (such as removing html tagger) and separating the document, we obtain 13.6M lines in the final input text. In order to identify the boundary of Chinese words, we use LTP<sup>4</sup> (Che et al., 2010) for Chinese Word Segmentation (CWS). We use official create\_pretraining\_data.py to convert raw input text to the pre-training examples,

which was provided in BERT GitHub repository. We generate two sets of pre-training examples with a maximum length of 128 and 512, as suggested by Devlin et al. (2019), for computation efficiency and learning long-range dependencies.

## 2.3 Pre-Training Details

We assume the whole word masking is a remedy for the BERT to know the word boundary and should be a 'patch' rather than a brand new model. Under this assumption, we did NOT train our model from scratch but from the official BERTbase (Chinese). We train 100K steps on the samples with a maximum length of 128, batch size of 2,560, an initial learning rate of 1e-4 (with warmup ratio 10%). Then, we train another 100K steps on a maximum length of 512 with a batch size of 384 to learn the long-range dependencies and position embeddings. Note that, the input data of the two phases should be changed according to the maximum length. Instead of using original AdamWeightDecayOptimizer in BERT, we use LAMB optimizer (You et al., 2019) for better scalability in large batch. <sup>5</sup> The pre-training was done on Google Cloud TPU v3 with 128G HBM.6

## 2.4 Fine-Tuning on Downstream Tasks

It is straightforward to use this model, as only one step is needed: replace original Chinese BERT<sup>7</sup> with our model, without changing config and vocabulary file.

## 3 Experiments

We carried out extensive experiments on various natural language processing tasks, covering a wide spectrum of text length, i.e. from sentence-level to document-level. Specifically, we choose the following popular Chinese datasets, including the ones that were also used in BERT and ERNIE. We adopt additional datasets for testing their performance in a wider range.

• Machine Reading Comprehension (MRC): CMRC 2018 (Cui et al., 2019), DRCD (Shao et al., 2018), CJRC (Duan et al., 2019)

<sup>3</sup>https://dumps.wikimedia.org/zhwiki/ latest/

<sup>4</sup>http://ltp.ai

<sup>&</sup>lt;sup>5</sup>For further tests and TensorFlow codes on LAMB optimizer, please refer to: https://github.com/ymcui/LAMB\_Optimizer\_TF

<sup>6</sup>https://cloud.google.com/tpu/

<sup>7</sup>https://storage.googleapis.com/bert\_
models/2018\_11\_03/chinese\_L-12\_H-768\_
A-12.zip

```
[Original Sentence]
使用语言模型来预测下一个词的probability。
[Original Sentence with CWS]
使用语言模型来预测下一个词的probability。

[Original BERT Input]
使用语言[MASK]型来[MASK]测下一个词的pro [MASK] ##lity。
[Whold Word Masking Input]
使用语言[MASK] [MASK] 来 [MASK] [MASK]下一个词的[MASK] [MASK] [MASK]。
```

Figure 1: Examples of the whole word masking in BERT.

- Natural Language Inference (NLI): XNLI (Conneau et al., 2018)
- Sentiment Classification (SC): ChnSenti-Corp<sup>8</sup>
- Sentence Pair Matching (SPM): LCQMC (Liu et al., 2018), BQ Corpus (Chen et al., 2018)
- **Document Classification (DC)**: THUC-News (Li and Sun, 2007)

In order to make a fair comparison, for each dataset, we keep the same hyper-parameters (such maximum length, warm-up steps, etc) and only tune the initial learning rate from 1e-4 to 1e-5 for each model. We run the same experiment for ten times to ensure the reliability of results. The best initial learning rate is determined by selecting the best average development set performance. We report the maximum, and average scores to both evaluate the peak and average performance of these models. For detailed hyper-parameter settings, please see Table 1.

In this technical report, we focus on comparing existing Chinese pre-trained models: BERT, ERNIE, and our models including BERT-wwm, BERT-wwm-ext, RoBERTa-wwm-ext, RoBERTa-wwm-ext-large. The model comparisons are depicted in Table 2.

We carried out all experiments under Tensor-Flow framework (Abadi et al., 2016). Note that, ERNIE only provides PaddlePaddle version<sup>9</sup>, so we have to convert the weights into TensorFlow version, where we obtain similar results on XNLI dataset which verifies that the conversion is successful.

## 3.1 Machine Reading Comprehension

Machine Reading Comprehension (MRC) is a representative document-level modeling task which requires to answer the questions based on the given passages. We mainly test these models on three datasets: CMRC 2018, DRCD, and CJRC.

- CMRC 2018: A span-extraction machine reading comprehension dataset, which is similar to SQuAD (Rajpurkar et al., 2016) that extract a passage span for the given question.
- **DRCD**: This is also a span-extraction MRC dataset, but in Traditional Chinese.
- CJRC: Similar to CoQA (Reddy et al., 2019), which has yes/no questions, no-answer questions and span-extraction questions. The data is collected from Chinese law judgment documents. Note that, we only use small-train-data.json for training. The development and test set are collected inhouse (does not publicly available due to the license issue and is not the same as the official competition).

The results are depicted in Table 3, 4, 5. As we can see, BERT-wwm yields significant improvements on CMRC 2018 and DRCD, which demonstrate its effectiveness on modeling long sequences. Also, we find that ERNIE does not show a competitive performance on DRCD, which indicate that it is not suitable for processing Traditional Chinese text. After examining the vocabulary of ERNIE, we discovered that the Traditional Chinese characters are removed<sup>10</sup>, and thus, resulting in an inferior performance. When it comes to CJRC, where the text is written in professional ways regarding Chinese laws, BERT-wwm shows moderate improvement over BERT and ERNIE,

 $<sup>^{8}</sup>$ https://github.com/pengming617/bert\_classification

<sup>9</sup>https://github.com/PaddlePaddle/LARK/
tree/develop/ERNIE

<sup>&</sup>lt;sup>10</sup>Not checked thoroughly, but we could not find some of the common Traditional Chinese characters.

Dataset	Task	MaxLen	Batch	Epoch	Train #	Dev #	Test #	Domain
CMRC 2018	MRC	512	64	2	10K	3.2K	4.9K	Wikipedia
DRCD	MRC	512	64	2	27K	3.5K	3.5K	Wikipedia
CJRC	MRC	512	64	2	10K	3.2K	3.2K	law
XNLI <sup>†‡</sup>	NLI	128	64	2	392K	2.5K	5K	various
ChnSentiCorp <sup>‡</sup>	SC	256	64	3	9.6K	1.2K	1.2K	various
LCQMC <sup>‡</sup>	SPM	128	64	3	240K	8.8K	12.5K	Zhidao
BQ Corpus	SPM	128	64	3	100K	10K	10K	QA
THUCNews	DC	512	64	3	50K	5K	10K	news

Table 1: Hyper-parameter settings and data statistics in different task.  $^{\dagger}$  represents the dataset was also evaluated by BERT (Devlin et al., 2019).  $^{\ddagger}$  represents the dataset was also evaluated by ERNIE (Sun et al., 2019). The dataset without any marks represent new benchmarks on these models.

	BERT	BERT-wwm	ERNIE
Pre-Train Data	Wikipedia	Wikipedia	Wikipedia +Baike+Tieba, etc.
Sentence #	2	24M	173M
Vocabulary #	21,128		18,000 (17,964)
Hidden Activation	G	eLU	ReLU
Hidden Size/Layers	Size/Layers		8 & 12
Attention Head #			12

Table 2: Comparisons of Chinese pre-trained models.

CMRC 2018	D	Dev		Test		Challenge	
CMIRC 2016	EM	<b>F1</b>	EM	<b>F1</b>	EM	F1	
BERT	65.5 (64.4)	84.5 (84.0)	70.0 (68.7)	87.0 (86.3)	18.6 (17.0)	43.3 (41.3)	
ERNIE	65.4 (64.3)	84.7 (84.2)	69.4 (68.2)	86.6 (86.1)	19.6 (17.0)	44.3 (42.8)	
BERT-wwm	66.3 (65.0)	85.6 (84.7)	70.5 (69.1)	87.4 (86.7)	21.0 (19.3)	47.0 (43.9)	
BERT-wwm-ext	67.1 (65.6)	85.7 (85.0)	71.4 (70.0)	87.7 (87.0)	24.0 (20.0)	47.3 (44.6)	
RoBERTa-wwm-ext	67.4 (66.5)	87.2 (86.5)	72.6 (71.4)	89.4 (88.8)	26.2 (24.6)	51.0 (49.1)	
RoBERTa-wwm-ext-large	68.5 (67.6)	88.4 (87.9)	74.2 (72.4)	90.6 (90.0)	31.5 (30.1)	60.1 (57.5)	

Table 3: Results on CMRC 2018 (Simplified Chinese). The average score of 10 independent runs is depicted in brackets. Overall best performance is depicted in boldface.

DDCD	D	ev	Test		
DRCD	EM	<b>F1</b>	EM	<b>F1</b>	
BERT	83.1 (82.7)	89.9 (89.6)	82.2 (81.6)	89.2 (88.8)	
ERNIE	73.2 (73.0)	83.9 (83.8)	71.9 (71.4)	82.5 (82.3)	
BERT-wwm	84.3 (83.4)	90.5 (90.2)	82.8 (81.8)	89.7 (89.0)	
BERT-wwm-ext	85.0 (84.5)	91.2 (90.9)	83.6 (83.0)	90.4 (89.9)	
RoBERTa-wwm-ext	86.6 (85.9)	92.5 (92.2)	85.6 (85.2)	92.0 (91.7)	
RoBERTa-wwm-ext-large	89.6 (89.1)	94.8 (94.4)	89.6 (88.9)	94.5 (94.1)	

Table 4: Results on DRCD (Traditional Chinese).

CIDC	D	ev	Test		
CJRC	$\mathbf{EM}$	<b>F1</b>	EM	<b>F1</b>	
BERT	54.6 (54.0)	75.4 (74.5)	55.1 (54.1)	75.2 (74.3)	
ERNIE	54.3 (53.9)	75.3 (74.6)	55.0 (53.9)	75.0 (73.9)	
BERT-wwm	54.7 (54.0)	75.2 (74.8)	55.1 (54.1)	75.4 (74.4)	
BERT-wwm-ext	55.6 (54.8)	76.0 (75.3)	55.6 (54.9)	75.8 (75.0)	
RoBERTa-wwm-ext	58.7 (57.6)	79.1 (78.3)	59.0 (57.8)	79.0 (78.0)	
RoBERTa-wwm-ext-large	62.1 (61.1)	82.4 (81.6)	62.4 (61.4)	82.2 (81.0)	

Table 5: Results on CJRC.

but not that salient, indicating that further domain adaptation is needed for non-general domains. Also, in professional domains, the performance of Chinese word segmentor may also decrease and will, in turn, affect the performance of ERNIE/BERT-wwm, which rely on Chinese segmentation.

## 3.2 Natural Language Inference

Following BERT and ERNIE, we use Chinese proportion of XNLI to test these models. The results show that ERNIE outperforms BERT/BERT-wwm significantly overall and BERT-wwm shows competitive performance on the test set.

XNLI	Dev	Test
BERT	77.8 (77.4)	77.8 (77.5)
ERNIE	79.7 (79.4)	78.6 (78.2)
BERT-wwm	79.0 (78.4)	78.2 (78.0)
BERT-wwm-ext	79.4 (78.6)	78.7 (78.3)
RoBERTa-wwm-ext	80.0 (79.2)	78.8 (78.3)
RoBERTa-wwm-ext-large	82.1 (81.3)	81.2 (80.6)

Table 6: Results on XNLI.

#### 3.3 Sentiment Classification

We use ChnSentiCorp, where the text should be classified into positive or negative label, for evaluating sentiment classification performance. We can see that ERNIE achieves the best performance on ChnSentiCorp, followed by BERT-wwm and BERT. When it comes to Sina Weibo, BERT-wwm shows better performance in terms of maximum and average scores on the test set. As ERNIE was trained on additional web text, it is beneficial to model non-formal text and capture the sentiment in social communication text, such as Weibo.

ChnSentiCorp	Dev	Test
BERT	94.7 (94.3)	95.0 (94.7)
ERNIE	95.4 (94.8)	95.4 (95.3)
BERT-wwm	95.1 (94.5)	95.4 (95.0)
BERT-wwm-ext	95.4 (94.6)	95.3 (94.7)
RoBERTa-wwm-ext	95.0 (94.6)	95.6 (94.8)
RoBERTa-wwm-ext-large	95.8 (94.9)	95.8 (94.9)

Table 7: Results on ChnSentiCorp.

## 3.4 Sentence Pair Matching

We adopt Large-scale Chinese Question Matching Corpus (LCQMC) and BQ Corpus for testing sentence pair matching task. As we can see that ERNIE outperforms BERT/BERT-wwm on LCQMC data. Though the peak performance of BERT-wwm is similar to BERT, the average score is relatively higher, indicating its potential in achieving higher scores (subject to the randomness). However, on BQ Corpus, we find BERT-wwm generally outperforms ERNIE and BERT, especially averaged scores.

## 3.5 Document Classification

THUCNews is a dataset that contains Sina news in different genres, which is a part of THUCTC.<sup>11</sup> In this paper, specifically, we use a version that contains 50K news in 10 domains (evenly distributed), including sports, finance, technology, etc.<sup>12</sup> As we can see that, BERT-wwm and BERT outperform ERNIE again on long sequence modeling task, demonstrating their effectiveness.

## 4 Useful Tips

As we can see, these pre-trained models behave differently in different natural language processing tasks. Due to the limited computing resources,

Sentence Pair	LCC	QMC	BQ Corpus	
Matching	Dev	Test	Dev	Test
BERT	89.4 (88.4)	86.9 (86.4)	86.0 (85.5)	84.8 (84.6)
ERNIE	89.8 (89.6)	87.2 (87.0)	86.3 (85.5)	85.0 (84.6)
BERT-wwm	89.4 (89.2)	87.0 (86.8)	86.1 (85.6)	85.2 (84.9)
BERT-wwm-ext	89.6 (89.2)	87.1 (86.6)	86.4 (85.5)	85.3 (84.8)
RoBERTa-wwm-ext	89.0 (88.7)	86.4 (86.1)	86.0 (85.4)	85.0 (84.6)
RoBERTa-wwm-ext-large	90.4 (90.0)	87.0 (86.8)	86.3 (85.7)	85.8 (84.9)

Table 8: Results on LCQMC and BQ Corpus.

THUCNews	Dev	Test
BERT ERNIE	97.7 (97.4) 97.6 (97.3)	97.8 (97.6) 97.5 (97.3)
BERT-wwm	98.0 (97.6)	97.8 (97.6)
BERT-wwm-ext	97.7 (97.5)	97.7 (97.5)
RoBERTa-wwm-ext	98.3 (97.9)	97.8 (97.5)
RoBERTa-wwm-ext-large	98.3 (97.7)	97.8 (97.6)

Table 9: Results on THUCNews.

we could not do exhaustive experiments on these datasets. However, we still have some (possibly) useful tips for the readers, where the tips are solely based on the materials above or our experience in using these models.

- Initial learning rate is the most important hyper-parameters (regardless of BERT or other neural networks), and should ALWAYS be tuned for better performance.
- As shown in the experimental results, BERT and BERT-wwm share almost the same best initial learning rate, so it is straightforward to apply your initial learning rate in BERT to BERT-wwm. However, we find that ERNIE does not share the same characteristics, so it is STRONGLY recommended to tune the learning rate.
- As BERT and BERT-wwm were trained on Wikipedia data, they show relatively better performance on the formal text. While, ERNIE was trained on larger data, including web text, which will be useful on casual text, such as Weibo (microblogs).
- In long-sequence tasks, such as machine reading comprehension and document classification, we suggest using BERT or BERTwwm.

- As these pre-trained models are trained in general domains, if the task data is extremely different from the pre-training data (Wikipedia for BERT/BERT-wwm), we suggest taking another pre-training steps on the task data, which was also suggested by (Devlin et al., 2019).
- As there are so many possibilities in pretraining stage (such as initial learning rate, global training steps, warm-up steps, etc.), our implementation may not be optimal using the same pre-training data. Readers are advised to train their own model if seeking for another boost in performance. However, if it is unable to do pre-training, choose one of these pre-trained models which was trained on a similar domain to the down-stream task.
- When dealing with Traditional Chinese text, use BERT or BERT-wwm.

## 5 Disclaimer

The experiments only represent the empirical results in certain conditions and should not be regarded as the nature of the respective models. The results may vary using different random seeds, computing devices, etc. Note that, as we have not been testing ERNIE on PaddlePaddle, the results in this technical report may not reflect its true performance (Though we have reproduced several results on the datasets that they had tested.).

#### 6 Conclusion

In this technical report, we utilize the whole word masking strategy for Chinese BERT and release the pre-trained model for the research community. The experimental results indicate that the proposed pre-trained model yields substantial improvements on various NLP tasks, compared to BERT and ERNIE. We hope the release of the pretrained models could further accelerate the natural language processing in the Chinese research community.

## Acknowledgments

Yiming Cui would like to thank TensorFlow Research Cloud (TFRC) program for supporting this research.

#### References

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. 2016. Tensorflow: A system for large-scale machine learning. In 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16), pages 265–283.
- Wanxiang Che, Zhenghua Li, and Ting Liu. 2010. Ltp: A chinese language technology platform. In *Proceedings of the 23rd International Conference on Computational Linguistics: Demonstrations*, pages 13–16. Association for Computational Linguistics.
- Jing Chen, Qingcai Chen, Xin Liu, Haijun Yang, Daohe Lu, and Buzhou Tang. 2018. The BQ corpus: A large-scale domain-specific Chinese corpus for sentence semantic equivalence identification. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4946–4951, Brussels, Belgium. Association for Computational Linguistics.
- Eunsol Choi, He He, Mohit Iyyer, Mark Yatskar, Wentau Yih, Yejin Choi, Percy Liang, and Luke Zettlemoyer. 2018. QuAC: Question answering in context. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2174–2184, Brussels, Belgium. Association for Computational Linguistics.
- Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel R. Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. Xnli: Evaluating crosslingual sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Yiming Cui, Zhipeng Chen, Si Wei, Shijin Wang, Ting Liu, and Guoping Hu. 2017. Attention-over-attention neural networks for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 593–602. Association for Computational Linguistics.
- Yiming Cui, Ting Liu, Wanxiang Che, Li Xiao, Zhipeng Chen, Wentao Ma, Shijin Wang, and Guoping Hu. 2019. A span-extraction dataset for chinese

- machine reading comprehension. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing*. Association for Computational Linguistics.
- Zihang Dai, Zhilin Yang, Yiming Yang, William W Cohen, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. 2019. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv* preprint arXiv:1901.02860.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Xingyi Duan, Baoxin Wang, Ziyue Wang, Wentao Ma, Yiming Cui, Dayong Wu, Shijin Wang, Ting Liu, Tianxiang Huo, Zhen Hu, et al. 2019. Cjrc: A reliable human-annotated benchmark dataset for chinese judicial reading comprehension. In *China National Conference on Chinese Computational Linguistics*, pages 439–451. Springer.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. Race: Large-scale reading comprehension dataset from examinations. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 796–805. Association for Computational Linguistics.
- Jingyang Li and Maosong Sun. 2007. Scalable term selection for text categorization. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*.
- Xin Liu, Qingcai Chen, Chong Deng, Huajun Zeng, Jing Chen, Dongfang Li, and Buzhou Tang. 2018. Lcqmc: A large-scale chinese question matching corpus. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1952–1962.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for SQuAD. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 784–789, Melbourne, Australia. Association for Computational Linguistics.

- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392. Association for Computational Linguistics.
- Qiu Ran, Peng Li, Weiwei Hu, and Jie Zhou. 2019. Option comparison network for multiple-choice reading comprehension. *arXiv preprint arXiv:1903.03033*.
- Siva Reddy, Danqi Chen, and Christopher D Manning. 2019. Coqa: A conversational question answering challenge. *Transactions of the Association for Computational Linguistics*, 7:249–266.
- Chih Chieh Shao, Trois Liu, Yuting Lai, Yiying Tseng, and Sam Tsai. 2018. Drcd: a chinese machine reading comprehension dataset. *arXiv preprint arXiv:1806.00920*.
- Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao Tian, and Hua Wu. 2019. Ernie: Enhanced representation through knowledge integration. *arXiv* preprint arXiv:1904.09223.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- Yang You, Jing Li, Jonathan Hseu, Xiaodan Song, James Demmel, and Cho-Jui Hsieh. 2019. Reducing bert pre-training time from 3 days to 76 minutes. *arXiv preprint arXiv:1904.00962*.
- Shuailiang Zhang, Hai Zhao, Yuwei Wu, Zhuosheng Zhang, Xi Zhou, and Xiang Zhou. 2019a. Dual comatching network for multi-choice reading comprehension. *arXiv preprint arXiv:1901.09381*.
- Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019b. ERNIE: Enhanced language representation with informative entities. In *Proceedings of ACL 2019*.

## **A Supplemental Material**

## A.1 Learning Rate of Fine-Tuning

Here we list the initial learning rates of fine-tuning different tasks.

	BERT	ERNIE	BERT-wwm*
CMRC 2018	3e-5	3e-5	8e-5
DRCD	3e-5	3e-5	8e-5
CJRC	4e-5	4e-5	8e-5
XNLI	3e-5	3e-5	5e-5
ChnSentiCorp	2e-5	2e-5	5e-5
LCQMC	2e-5	2e-5	3e-5
BQ Corpus	3e-5	3e-5	5e-5
THUCNews	2e-5	2e-5	5e-5

Table 10: Best initial learning rate for different task. \* represents all related models.