Enhancing Pre-trained Chinese Character Representation with Word-aligned Attention

Yanzeng Li, Bowen Yu, Mengge Xue, Tingwen Liu*

Institute of Information Engineering, Chinese Academy of Sciences School of Cyber Security, University of Chinese Academy of Sciences {liyanzeng, yubowen, xuemengge, liutingwen}@iie.ac.cn

Abstract

Most Chinese pre-trained models take character as the basic unit and learn representation according to character's external contexts, ignoring the semantics expressed in the word, which is the smallest meaningful utterance in Chinese. Hence, we propose a novel wordaligned attention to exploit explicit word information, which is complementary to various character-based Chinese pre-trained language models. Specifically, we devise a pooling mechanism to align the character-level attention to the word level and propose to alleviate the potential issue of segmentation error propagation by multi-source information fusion. As a result, word and character information are explicitly integrated at the fine-tuning procedure. Experimental results on five Chinese NLP benchmark tasks demonstrate that our method achieves significant improvements against BERT, ERNIE and BERT-wwm.

1 Introduction

Pre-trained language Models (PLM) such as ELMo (Peters et al., 2018), BERT (Devlin et al., 2019), ERNIE (Sun et al., 2019), BERT-wwm (Cui et al., 2019) and XLNet (Yang et al., 2019) have been proven to capture rich language information from text and then benefit many NLP applications by simple fine-tuning, including sentiment classification (Pang et al., 2002), natural language inference (Bowman et al., 2015), named entity recognition (Sang and De Meulder, 2003) and so on.

Generally, most popular PLMs prefer to use attention mechanism (Vaswani et al., 2017) to represent the natural language, such as word-to-word self-attention for English. Unlike English, in Chinese, words are not separated by explicit delimiters. Since without word boundaries information, it is

intuitive to model characters in Chinese tasks directly. However, in most cases, the semantic of a single Chinese character is ambiguous. For example, the character "拍" in word "球拍 (bat)" and "拍卖 (auction)" has entirely different meanings. Moreover, several recent works have demonstrated that considering the word segmentation information can lead to better language understanding, and accordingly benefits various Chinese tasks (Wang et al., 2017; Li et al., 2018; Zhang and Yang, 2018; Gui et al., 2019; Mengge et al., 2019).

All these factors motivate us to expand the character-level attention mechanism in Chinese PLMs to represent the semantics of words ¹. To this end, there are two main challenges. (1) How to seamlessly integrate the segmentation information into character-based attention module of PLM is an important problem. (2) Gold-standard segmentation is rarely available in the downstream tasks, and how to effectively reduce the cascading noise caused by Chinese word segmentation (CWS) tools (Li et al., 2019) is another challenge.

In this paper, we propose a new architecture, named Multi-source Word Aligned Attention (MWA), to solve the above issues. (1) Psycholinguistic experiments (Bai et al., 2008; Meng et al., 2014) have shown that readers are likely to pay approximate attention to each character in one Chinese word. Drawing inspiration from such findings, we introduce a novel word-aligned attention, which could aggregate attention weight of characters in one word into a unified value with the mixed pooling strategy (Yu et al., 2014). (2) For reducing segmentation error, we further extend our word-aligned attention with multi-source segmentation produced by various segmenters and deploy

^{*}Corresponding author

¹Considering the enormous cost of re-training a language model, we hope to incorporate word segmentation information to the fine-tuning process to enhance performance, and leave how to improve the pre-training procedure for a future work.

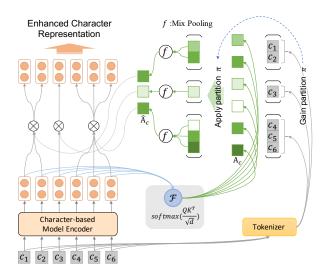


Figure 1: Architecture of Word-aligned Attention

a fusion function to pull together their disparate outputs. As shown in Table 1, different CWS tools may have different annotation granularity. Through comprehensive consideration of multi-granularity segmentation results, we can implicitly reduce the error caused by automatic annotation.

Extensive experiments are conducted on various Chinese NLP tasks including sentiment classification, named entity recognition, sentence pair matching, natural language inference and machine reading comprehension. The results and analysis show that the proposed method boosts BERT, ERNIE and BERT-wwm significantly on all the datasets ².

2 Methodology

2.1 Character-level Pre-trained Encoder

The primary goal of this work is to inject the word segmentation knowledge into character-level Chinese PLMs and enhance original models. Given the strong performance of deep Transformers trained on language modeling, we adopt BERT and its updated variants (ERNIE, BERT-wwm) as the basic encoder in this work, and the outputs from the last layer of encoder are treated as the character-level enriched contextual representations **H**.

2.2 Word-aligned Attention

Although character-level Chinese PLM has remarkable ability to capture language knowledge from text, it neglects the semantic information expressed in the word level. Therefore we apply a word-aligned layer on top of the encoder to integrate the

(Chinese	北京西山森林公园					
enter	thulac	ulac 北京 西山		森林	公园		
8	ictclas	北京	西	山	森林	公园	
Segi	hanlp	北京	西山		森林	公园	

Table 1: Results of different popular CWS tools over "北京西山森林公园(Beijing west mount forest park)".

word boundary information into the representation of characters with an attention aggregation module.

For an input sequence with n characters $S = [c_1, c_2, ..., c_n]$, where c_j denotes the j-th character, CWS tool π is used to partition S into non-overlapping word blocks:

$$\pi(S) = [w_1, w_2, ..., w_m], (m \le n)$$
 (1)

where $w_i = \{c_s, c_{s+1}, ..., c_{s+l-1}\}$ is the *i*-th segmented word with a length of l and s is the index of w_i 's first character in S. We apply self-attention operation with the representations of all input characters to get the character-level attention score matrix $\mathbf{A}_c \in \mathbb{R}^{n \times n}$. It can be formulated as:

$$\mathbf{A}_c = \mathcal{F}(\mathbf{H}) = \operatorname{softmax}(\frac{(\mathbf{K}\mathbf{W}_k)(\mathbf{Q}\mathbf{W}_q)^T}{\sqrt{d}})$$
 (2)

where \mathbf{Q} and \mathbf{K} are both equal to the collective representation \mathbf{H} at the last layer of the Chinese PLM, $\mathbf{W}_k \in \mathbb{R}^{d \times d}$ and $\mathbf{W}_q \in \mathbb{R}^{d \times d}$ are trainable parameters for projection. While \mathbf{A}_c models the relationship between two arbitrarily characters without regard to the word boundary, we argue that incorporating word as atom in the attention can better represent the semantics, as the literal meaning of each individual character can be quite different from the implied meaning of the whole word, and the simple weighted sum in the character level may lose word and word sequence information.

To address this issue, we propose to align \mathbf{A}_c in the word level and integrate the inner-word attention. For ease of exposition, we rewrite \mathbf{A}_c as $[\mathbf{a}_c^1, \mathbf{a}_c^2, ..., \mathbf{a}_c^n]$, where $\mathbf{a}_c^i \in \mathbb{R}^n$ denotes the *i*-th row vector of \mathbf{A}_c , that is, \mathbf{a}_c^i represents the attention score vector of the *i*-th character. Then we deploy π to segment \mathbf{A}_c according to $\pi(S)$. For example, if $\pi(S) = [\{c_1, c_2\}, \{c_3\}, ..., \{c_{n-1}, c_n\}]$, then

$$\pi(\mathbf{A}_c) = [\{\mathbf{a}_c^1, \mathbf{a}_c^2\}, \{\mathbf{a}_c^3\}, ..., \{\mathbf{a}_c^{n-1}, \mathbf{a}_c^n\}]$$
 (3)

In this way, an attention vector sequence is divided into several subsequences and each subsequence represents the attention of one word.

²The source code of this paper can be obtained from https://github.com/lsvih/MWA.

Then, motivated by the psycholinguistic finding that readers are likely to pay similar attention to each character in one Chinese word, we devise an appropriate aggregation module to fuse the innerword attention. Concretely, we first transform $\{\mathbf{a}_c^s,...,\mathbf{a}_c^{s+l-1}\}$ into one attention vector \mathbf{a}_w^i for w_i with the mixed pooling strategy (Yu et al., 2014) ³. Then we execute the piecewise upsampling operation over each \mathbf{a}_w^i to keep input and output dimensions unchanged for the sake of plug and play. The detailed process can be summarized as:

$$\begin{aligned} \mathbf{a}_w^i &= \lambda \text{ Maxpooling}(\{\mathbf{a}_c^s,...,\mathbf{a}_c^{s+l-1}\}) \\ &+ (1-\lambda) \text{ Meanpooling}(\{\mathbf{a}_c^s,...,\mathbf{a}_c^{s+l-1}\}) \end{aligned} \tag{4}$$

$$\hat{\mathbf{A}}_c[s:s+l-1] = \mathbf{e}_l \otimes \mathbf{a}_w^i \tag{5}$$

where $\lambda \in R^1$ is a weighting trainable variable to balance the mean and max pooling, $\mathbf{e}_l = [1,...,1]^T$ represents a l-dimensional all-ones vector, l is the length of w_i , $\mathbf{e}_l \otimes \mathbf{a}_w^i = [\mathbf{a}_w^i,...,\mathbf{a}_w^i]$ denotes the kronecker product operation between \mathbf{e}_l and \mathbf{a}_w^i , $\hat{\mathbf{A}}_c \in \mathbb{R}^{n \times n}$ is the aligned attention matrix. Eqs. 4 and 5 can help incorporate word segmentation information into character-level attention calculation process, and determine the attention vector of one character from the perspective of the whole word, which is beneficial for eliminating the attention bias caused by character ambiguity. Finally, we can obtain the enhanced character representation produced by word-aligned attention as follows:

$$\hat{\mathbf{H}} = \hat{\mathbf{A}}_c \mathbf{V} \mathbf{W}_v \tag{6}$$

where $\mathbf{V} = \mathbf{H}$, $\mathbf{W}_v \in \mathbb{R}^{d \times d}$ is a trainable projection matrix. Besides, we also use multi-head attention (Vaswani et al., 2017) to capture information from different representation subspaces jointly, thus we have K different aligned attention matrices $\hat{\mathbf{A}}_c^k (1 \le k \le K)$ and corresponding representation $\hat{\mathbf{H}}^k$. With multi-head attention architecture, the output can be expressed as follows:

$$\overline{\mathbf{H}} = \operatorname{Concat}(\hat{\mathbf{H}}^1, \hat{\mathbf{H}}^2, ..., \hat{\mathbf{H}}^K) \mathbf{W}_o \tag{7}$$

2.3 Multi-source Word-aligned Attention

As mentioned in Section 1, our proposed wordaligned attention relies on the segmentation results of CWS tool π . Unfortunately, a segmenter is usually unreliable due to the risk of ambiguous and non-formal input, especially on out-of-domain data, which may lead to error propagation and an unsatisfactory model performance. In practice, the ambiguous distinction between morphemes and compound words leads to the cognitive divergence of words concepts, thus different π may provide diverse $\pi(S)$ with various granularities. To reduce the impact of segmentation error and effectively mine the common knowledge of different segmenters, it's natural to enhance the word-aligned attention layer with multi-source segmentation inputs. Formally, assume that there are M popular CWS tools employed, we can obtain M different representations $\overline{\mathbf{H}}^1, ..., \overline{\mathbf{H}}^M$ by Eq. 7. Then we propose to fuse these semantically different representations as follows:

$$\tilde{\mathbf{H}} = \sum_{m=1}^{M} \tanh(\overline{\mathbf{H}}^{m} \mathbf{W}_{g})$$
 (8)

where W_g is a parameter matrix and H denotes the final output of the MWA attention layer.

3 Experiments

3.1 Experiments Setup

To test the applicability of the proposed MWA attention, we choose three publicly available Chinese pre-trained models as the basic encoder: BERT, ERNIE, and BERT-wwm. In order to make a fair comparison, we keep **the same hyper-parameters** (such maximum length, warm-up steps, initial learning rate, etc.) as suggested in BERT-wwm (Cui et al., 2019) for both baselines and our method on each dataset. **We run the same experiment for five times and report the average score** to ensure the reliability of results. Besides, three popular CWS tools: thulac (Sun et al., 2016), ictclas (Zhang et al., 2003) and hanlp (He, 2014) are employed to segment sequence.

The experiments are carried out on five Chinese NLP tasks and six public benchmark datasets:

Sentiment Classification (SC): We adopt ChnSentiCorp⁴ and weibo-100k sentiment dataset⁵ in this task. ChnSentiCorp dataset has about 10k sentences, which express positive or negative emotion. weibo-100k dataset contains 1.2M microblog

³Other pooling methods such as max pooling or mean pooling also works. Here we choose mixed pooling because it has the advantages of distilling the global and the most prominent features in one word at the same time.

https://github.com/pengming617/bert_ classification

⁵https://github.com/SophonPlus/ ChineseNlpCorpus/

Dataset	Task	Max length	Batch size	Epoch	lr* -	Dataset Size		
Dataset	Task	Max length				Train	Dev	Test
ChnSentiCorp	- SC	256	16	3	3×10^{-5}	9.2K	1.2K	1.2K
weibo-100k		128	64	2	2×10^{-5}	100K	~10K	10K
ontonotes	NER	256	16	5	3×10^{-5}	15.7K	4.3K	4.3K
LCQMC	SPM	128	64	3	3×10^{-5}	~239K	8.8K	12.5K
XNLI	NLI	128	64	2	3×10^{-5}	~392K	2.5K	2.5K
DRCD	MRC	512	16	2	3×10^{-5}	27K	3.5K	3.5K

Table 2: Summary of datasets and the corresponding hyper-parameters setting. Reported learning rates* are the initial values of BertAdam.

texts and each microblog is tagged as positive or negative emotion.

Named Entity Recognition (NER): this task is to test model's capacity of sequence tagging. We use a common public dataset Ontonotes 4.0 (Weischedel et al., 2011) in this task.

Sentence Pair Matching (SPM): We use the most widely used dataset LCQMC (Liu et al., 2018) in this task, which aims to identify whether two questions are in a same intention.

Natural Language Inference (NLI): this task is to exploit the contexts of text and concern inference relationships between sentences. XNLI (Conneau et al., 2018) is a cross-language language understanding dataset; we only use the Chinese language part of XNLI to evaluate the language understanding ability. And we processed this dataset in the same way as ERNIE (Sun et al., 2019) did.

Machine Reading Comprehension (MRC): MRC is a representative document-level modeling task which requires to answer the questions based on the given passages. DRCD (Shao et al., 2018) is a public span-extraction Chinese MRC dataset, whose answers are spans in the document.

We implement our model with PyTorch (Paszke et al., 2019), and all baselines are converted weights into PyTorch version. All experiments employ modified Adam (Devlin et al., 2019) as optimizer with 0.01 weight decay and 0.1 warm-up ratio. All pre-trained models are configured to 12 layers and 768 hidden dimension. The detail settings are shown in Table 2.

3.2 Experiment Results

Table 3 shows the performances on five classical Chinese NLP tasks with six public datasets. Generally, our method consistently outperforms all baselines on all five tasks, which demonstrates the effectiveness and universality of the proposed approach. Moreover, the Wilcoxon's test shows that a significant difference (p < 0.05) exits between our model and baseline models.

In detail, on the two datasets of SC task, we observe an average of 0.53% and 0.83% absolute improvement in F1 score, respectively. SPM and NLI tasks can also gain benefits from our enhanced representation. For the NER task, our method obtains 0.92% improvement averagely over all baselines. Besides, introducing word segmentation information into the encoding of character sequences improves the MRC performance on average by 1.22 points and 1.65 points in F1 and Exact Match (EM) score respectively. We attribute such significant gain in NER and MRC to the particularity of these two tasks. Intuitively, Chinese NER is correlated with word segmentation, and named entity boundaries are also word boundaries. Thus the potential boundary information presented by the additional segmentation input can provide better guidance to label each character, which is consistent with the conclusion in (Zhang and Yang, 2018). Similarly, the span-extraction MRC task is to extract answer spans from document (Shao et al., 2018), which also faces the same word boundary problem as NER, and the long sequence in MRC exacerbates the problem. Therefore, our method gets a relatively greater improvement on the DRCD dataset.

3.3 Ablation Study

To demonstrate the effectiveness of our multisource fusion method, we carry out experiments on the DRCD dev set with different segmentation inputs. Besides, we also design two strong baselines by introducing a Transformer layer (1T) and a random tokenizer model (WA_{random}) to exclude the benefits from additional parameters. As shown in Table 4, adding additional parameters by introducing an extra transformer layer can benefit the PLMs. Compared with 1T and WA_{random} , our proposed word-aligned attention gives quite stable improvements no matter what CWS tool we use, which again confirms the effectiveness and rationality of incorporating word segmentation information into character-level PLMs. Another observation is that

Task	S	C	NER	SPM	NLI		RC
Dataset	ChnSenti ^{2,3}	weibo-100k ²	Ontonotes ⁴	LCQMC ^{2,3,4}	$XNLI^{1,2,3,4}$	$DRCD^{2,3}$	[EM F1]
Prev. SOTA [†]	93.1(2019a)	-	74.89(2019b)	85.68(2019c)	67.5(2017d)	75.12(2019e)	87.26(2019e)
BERT	94.72	97.31	79.18	86.50	78.19	85.57	91.16
+MWA	95.34(+0.62)	98.14(+0.83)	79.86(+0.68)	86.92(+0.42)	78.42(+0.23)	86.86(+1.29)	92.22(+1.06)
BERT-wwm	94.38	97.36	79.28	86.11	77.92	84.11	90.46
+MWA	95.01(+0.63)	98.13(+0.77)	80.32 (+1.04)	86.28(+0.17)	78.68(+0.76)	87.00(+2.89)	92.21(+1.75)
ERNIE	95.17	97.30	77.74	87.27	78.04	87.85	92.85
+MWA	95.52 (+0.35)	98.18 (+0.88)	78.78(+1.04)	88.73 (+1.46)	78.71 (+0.67)	88.61 (+0.76)	93.72 (+0.87)

Table 3: Evaluation results regarding each model on different datasets. Bold marks highest number among all models. Numbers in brackets indicate the absolute increase over baseline models. Superscript number 1,2,3,4 respectively represents that the corresponding dataset is also used by BERT (Devlin et al., 2019), BERT-wwm (Wu et al., 2016; Cui et al., 2019), ERNIE (Sun et al., 2019) and Glyce (Meng et al., 2019a), respectively. The results of all baselines are produced by our implementation or retrieved from original papers, and we report the higher one among them. The improvements over baselines are statistically significant (p < 0.05). † denotes the results of previous state-of-the-art models on these datasets without using BERT.

Model	BERT	BERT-wwm	ERNIE
Original	92.06	91.68	92.61
+1T	92.37	92.22	93.42
$+WA_{random}$	91.83	90.33	92.12
$+WA_{thulac}$	92.84	92.73	93.89
$+WA_{ictclas}$	93.05	92.90	93.75
$+WA_{hanlp}$	92.91	93.21	93.91
+MWA	93.59	93.72	94.21

Table 4: F1 results of ablation experiments on the DRCD dev set.

employing multiple segmenters and fusing them together could introduce richer segmentation information and further improve the performance.

3.4 Parameter Scale Analysis

For fair comparison and demonstrating the improvement of our model is not only rely on more trainable parameters, we also conduct experiments on the DRCD dev set to explore whether the performance keeps going-up with more parameters by introducing additional transformer blocks on top of the representations of PLMs.

Model	F1	Param. Number
BERT-wwm	91.68	110M
BERT-wwm+ $1T$	92.23	110M+7.1M
BERT-wwm+ $2T$	91.99	110M+14.2M
BERT-wwm+ $3T$	91.68	110M+21.3M
BERT-wwm+MWA	93.72	110M+7.6M
Robust-BERT-wwm-ext-large	94.40	340M

Table 5: Comparison on the DRCD dev set. The nT denotes the number of additional transformer layers.

In Table 5, +1T denotes that we introduce another one Transformer layer on top of BERT-wwm and +2T means additional 2 layers, M denotes million. As the experimental results showed, when the number of additional layers exceeds 1, the performance starts to decline, which demonstrates

that using an extensive model on top of the PLM representations may not bring additional benefits. We can conclude that MWA doesn't introduce too many parameters, and MWA achieves better performance than +1T under the similar parameter numbers. Besides, we also make comparison with the current best Chinese PLM: Robust-BERT-wwm-ext-large (Cui et al., 2019), a 24-layers Chinese PLM with 13.5 times more pre-training data and 3.1 times more parameters than BERT-wwm, experimental results show that our model can achieve comparable performance, which again confirms the effectiveness of incorporating word segmentation information into character-level PLMs.

4 Conclusion

In this paper, we develop a novel Multi-source Word Aligned Attention model (referred as MWA), which integrates word segmentation information into character-level self-attention mechanism to enhance the fine-tuning performance of Chinese PLMs. We conduct extensive experiments on five NLP tasks with six public datasets. The proposed approach yields substantial improvements compared to BERT, BERT-wwm and ERNIE, demonstrating its effectiveness and universality. Furthermore, the word-aligned attention can also be applied to English PLMs to bridge the semantic gap between the whole word and the segmented Word-Piece tokens, which we leave for future work.

Acknowledgement

We would like to thank reviewers for their insightful comments. This work is supported by the Strategic Priority Research Program of Chinese Academy of Sciences, Grant No. XDC02040400.

References

- Xuejun Bai, Guoli Yan, Simon P Liversedge, Chuanli Zang, and Keith Rayner. 2008. Reading spaced and unspaced chinese text: Evidence from eye movements. *Journal of Experimental Psychology: Human Perception and Performance*, 34(5):1277.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. 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, Wanxiang Che, Ting Liu, Bing Qin, Ziqing Yang, Shijin Wang, and Guoping Hu. 2019. Pre-training with whole word masking for chinese bert. *arXiv preprint arXiv:1906.08101*.
- 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.
- Tao Gui, Ruotian Ma, Qi Zhang, Lujun Zhao, Yu-Gang Jiang, and Xuanjing Huang. 2019. Cnn-based chinese ner with lexicon rethinking. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pages 4982–4988. AAAI Press.
- Tao Gui, Yicheng Zou, Qi Zhang, Minlong Peng, Jinlan Fu, Zhongyu Wei, and Xuanjing Huang. 2019b. A lexicon-based graph neural network for chinese ner. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 1039–1049.
- Han He. 2014. HanLP: Han Language Processing.
- Qiang Huang, Jianhui Bu, Weijian Xie, Shengwen Yang, Weijia Wu, and Liping Liu. 2019c. Multi-task sentence encoding model for semantic retrieval in question answering systems. In 2019 International Joint Conference on Neural Networks (IJCNN), pages 1–8.
- Chia-Hsuan Lee and Hung-Yi Lee. 2019e. Crosslingual transfer learning for question answering.

- Xiaoya Li, Yuxian Meng, Xiaofei Sun, Qinghong Han, Arianna Yuan, and Jiwei Li. 2019. Is word segmentation necessary for deep learning of chinese representations? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3242–3252.
- Yanzeng Li, Tingwen Liu, Diying Li, Quangang Li, Jinqiao Shi, and Yanqiu Wang. 2018. Character-based bilstm-crf incorporating pos and dictionaries for chinese opinion target extraction. In *Asian Conference on Machine Learning*, pages 518–533.
- 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.
- Hongxia Meng, Xuejun Bai, Chuanli Zang, and Guoli Yan. 2014. Landing position effects of coordinate and attributive structure compound words. *Acta Psychologica Sinica*, 46(1):36–49.
- Yuxian Meng, Wei Wu, Fei Wang, Xiaoya Li, Ping Nie, Fan Yin, Muyu Li, Qinghong Han, Xiaofei Sun, and Jiwei Li. 2019a. Glyce: Glyph-vectors for chinese character representations. In *NeurIPS 2019: Thirty-third Conference on Neural Information Processing Systems*, pages 2746–2757.
- Xue Mengge, Yu Bowen, Liu Tingwen, Wang Bin, Meng Erli, and Li Quangang. 2019. Porous lattice-based transformer encoder for chinese ner.
- Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up?: sentiment classification using machine learning techniques. In *Proceedings of the ACL-02 conference on Empirical methods in natural language processing-Volume 10*, pages 79–86. Association for Computational Linguistics.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proc. of NAACL*.
- Erik F Sang and Fien De Meulder. 2003. Introduction to the conll-2003 shared task: language-independent named entity recognition. In *Proceedings of the seventh conference on Natural language learning at HLT-NAACL 2003-Volume 4*, pages 142–147. Association for Computational Linguistics.

- 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*.
- Maosong Sun, Xinxiong Chen, Kaixu Zhang, Zhipeng Guo, and Zhiyuan Liu. 2016. Thulac: An efficient lexical analyzer for chinese. Technical report.
- 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.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Shaonan Wang, Jiajun Zhang, and Chengqing Zong. 2017. Exploiting word internal structures for generic Chinese sentence representation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 298–303, Copenhagen, Denmark. Association for Computational Linguistics.
- Zhiguo Wang, Wael Hamza, and Radu Florian. 2017d. Bilateral multi-perspective matching for natural language sentences. In *Twenty-Sixth International Joint Conference on Artificial Intelligence*, pages 4144–4150.
- Ralph Weischedel, Martha Palmer, Mitchell Marcus, Eduard Hovy, Sameer Pradhan, Lance Ramshaw, Nianwen Xue, Ann Taylor, Jeff Kaufman, Michelle Franchini, et al. 2011. Ontonotes 4.0. *Linguistic Data Consortium LDC2011T03*.
- 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*.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *arXiv preprint arXiv:1906.08237*.
- Dingjun Yu, Hanli Wang, Peiqiu Chen, and Zhihua Wei. 2014. Mixed pooling for convolutional neural networks. In *International Conference on Rough Sets and Knowledge Technology*, pages 364–375. Springer.
- Hua-Ping Zhang, Hong-Kui Yu, De-Yi Xiong, and Qun Liu. 2003. Hhmm-based chinese lexical analyzer ictclas. In *Proceedings of the second SIGHAN workshop on Chinese language processing-Volume 17*, pages 184–187. Association for Computational Linguistics.

Yue Zhang and Jie Yang. 2018. Chinese ner using lattice lstm. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 1554–1564.