

An Embarrassingly Easy but Strong Baseline for Nested Named Entity Recognition

Hang Yan*, Yu Sun*, Xiaonan Li, Xipeng Qiu†

Shanghai Key Laboratory of Intelligent Information Processing, Fudan University
School of Computer Science, Fudan University
{hyan19, lixn20, xpqiu}@fudan.edu.cn
yusun21@m.fudan.edu.cn

Abstract

Named entity recognition (NER) is the task to detect and classify entity spans in the text. When entity spans overlap between each other, the task is named as nested NER. Span-based methods have been widely used to tackle nested NER. Most of these methods get a score matrix, where each entry corresponds to a span. However, previous work ignores spatial relations in the score matrix. In this paper, we propose using Convolutional Neural Network (CNN) to model these spatial relations. Despite being simple, experiments in three commonly used nested NER datasets show that our model surpasses several recently proposed methods with the same pre-trained encoders. Further analysis shows that using CNN can help the model find more nested entities. Besides, we find that different papers use different sentence tokenizations for the three nested NER datasets, which will influence the comparison. Thus, we release a pre-processing script to facilitate future comparison.¹

1 Introduction

Named Entity Recognition (NER) is the task to extract entities from raw text. It has been a fundamental task in the Natural Language Processing (NLP) field. Previously, this task is mainly solved by the sequence labeling paradigm through assigning a label to each token (Huang et al., 2015; Ma and Hovy, 2016; Yan et al., 2019). However, this method is not directly applicable to the nested NER scenario, since a token may be included in two or more entities. To overcome this issue, the span-based method which assigns labels to each span is introduced (Eberts and Ulges, 2020; Li et al., 2020; Yu et al., 2020).

*Equal contribution.

†Corresponding author.

¹Code is available at https://github.com/yhcc/CNN_Nested_NER

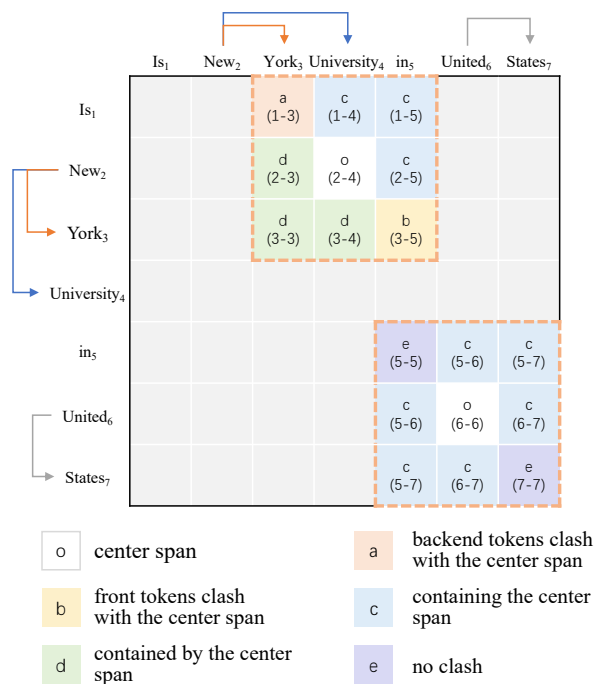


Figure 1: All valid spans of a sentence. We use the start and end tokens to pinpoint a span, for instance, “(2-4)” represents “New York University”. Spans in the two orange dotted squares indicates that the center span can have the special relationship (different relations are depicted in different colors) with its surrounding spans. For example, the span “New York” (2-3) is contained by the span “New York University” (2-4). Therefore, the “(2-3)” span is annotated as “d”.

Eberts and Ulges (2020) use a pooling method over token representations to get the span representation, and then conduct classification on this span representation. Li et al. (2020) transform the NER task into a Machine Reading Comprehension (MRC) form, they use the entity type as the query, and ask the model to select spans that belong to this entity type. Yu et al. (2020) utilize the Biaffine decoder from dependency parsing (Dozat and Manning, 2017) to convert the span classification into classifying the start and end token pairs. However, these work does not take advantage of the spatial

correlations between adjacent spans.

As depicted in Figure 1, spans surrounding a span have special relationships with the center span. It should be beneficial if we can leverage these spatial correlations. In this paper, we use the Biaffine decoder (Dozat and Manning, 2017) to get a 3D feature matrix, where each entry represents one span. After that, we view the span feature matrix as a spatial object with channels (like images) and utilize Convolutional Neural Network (CNN) to model the local interaction between spans.

We compare this simple method with recently proposed methods (Wan et al., 2022; Li et al., 2022; Zhu and Li, 2022; Yuan et al., 2022). To make sure our method is strictly comparable to theirs, we ask the authors for their version of data. Although all of them use the same datasets, we find that the statistics, such as the number of sentences and entities, are not the same. The difference is caused by the usage of distinct sentence tokenization methods, which will influence the performance as shown in our experiments. To facilitate future comparison, we release a pre-processing script for ACE2004, ACE2005 and Genia datasets.

Our contributions can be summarized as follows.

- We find that the adjacent spans have special correlations between each other, and we propose using CNN to model the interaction between them. Despite being very simple, it achieves a considerable performance boost in three widely used nested NER datasets.
- We release a pre-processing script for the three nested NER datasets to facilitate direct and fair comparison.
- The way we view the span feature matrix as a spatial object with channels shall shed some light on future exploration of span-based methods for nested NER task.

2 Proposed Method

In this section, we first introduce the nested NER task, then describe how to get the feature matrix. After that, we present the CNN module to model the spatial correlation on the feature matrix. A general framework can be viewed in Figure 2.

2.1 Nested NER Task

Given an input sentence $X = [x_1, x_2, \dots, x_n]$ with n tokens, the nested NER task aims to extract all

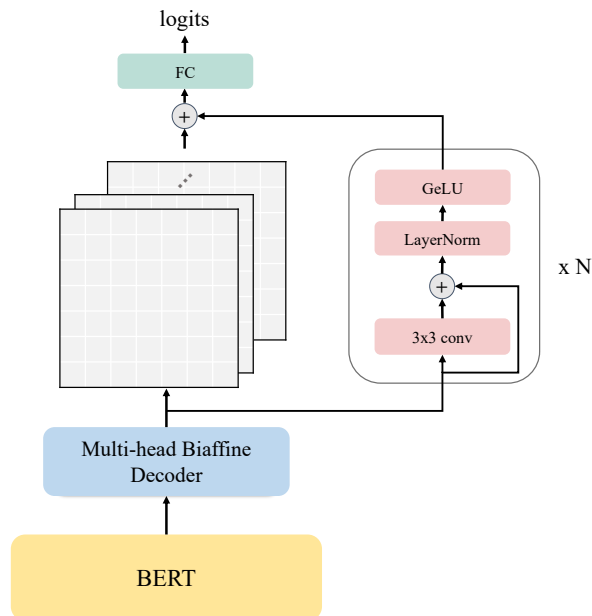


Figure 2: The proposed method in this paper. Use several blocks of CNN to model the spatial correlations between neighbor spans.

entities in X . Each entity can be expressed as a tuple (s_i, e_i, t_i) . s_i, e_i are the start, end index of the entity. $t_i \in \{1, \dots, |T|\}$ is its entity type and $T = \{t_1, \dots, t_n\}$ is entity types. As the task name suggests, entities may overlap with each other, but different entities are not allowed to have crossing boundaries. For a sentence with n tokens, there are $n(n+1)/2$ valid spans.

2.2 Span-based Representation

We follow Yu et al. (2020) to formulate this task into a span classification task. Namely, for each valid span, the model assigns an entity label to it. The method first uses an encoder to encode the input sentence as follows:

$$\mathbf{H} = \text{Encoder}(X),$$

where $\mathbf{H} \in \mathbb{R}^{n \times d}$, and d is the hidden size. Various pre-trained models, such as BERT (Devlin et al., 2019), are usually used as the encoder. For the word tokenized into several pieces, we use max-pooling to aggregate from its pieces' hidden states.

Next, we use a multi-head Biaffine decoder (Dozat and Manning, 2017; Vaswani et al., 2017) to get the score matrix \mathbf{R} as follows:

$$\begin{aligned} \mathbf{H}_s &= \text{LeakyReLU}(\mathbf{H}\mathbf{W}_s), \\ \mathbf{H}_e &= \text{LeakyReLU}(\mathbf{H}\mathbf{W}_e), \\ \mathbf{R} &= \text{MHBiaffine}(\mathbf{H}_s, \mathbf{H}_e) \end{aligned}$$

	# Param. (Million)	ACE2004			ACE2005		
		P	R	F1	P	R	F1
<i>Data from Li et al. (2022)</i>							
W2NER (2022)[BERT-large]	355.4	87.33	87.71	87.52	85.03	88.62	86.79
Ours[BERT-large]	345.1	87.82 ₃₈	87.40 ₂₀	87.61 ₁₈	86.39 ₆₁	87.24 ₃₄	86.82 ₄₅
w.o. CNN[BERT-large]	343.6	86.54 ₄₈	87.09 ₄₁	86.81 ₂₁	84.88 ₂₆	86.99 ₃₃	85.92 ₂₇
<i>Data from Wan et al. (2022)</i>							
SG (2022)[BERT-base]	112.3	86.70	85.93	86.31	84.37	85.87	85.11
Ours[BERT-base]	110.5	86.85 ₆₁	86.45 ₃₆	86.65 ₂₂	84.94 ₄₉	85.40 ₂₇	85.16 ₁₆
w.o. CNN[BERT-base]	109.1	85.79 ₄₆	85.78 ₁₂	85.78 ₂₂	82.91 ₂₁	84.89 ₂₃	83.89 ₁₆
<i>Data from Zhu and Li (2022)</i>							
BS (2022)[RoBERTa-base]	125.6	88.43	87.53	87.98	86.25	88.07	87.15
Ours[RoBERTa-base]	125.6	87.77 ₂₇	88.28 ₃₆	88.03 ₁₄	86.58 ₇₈	87.94 ₄₆	87.25 ₄₈
w.o. CNN[RoBERTa-base]	125.2	86.71 ₂₇	87.40 ₄₂	87.05 ₁₈	85.48 ₃₉	87.54 ₅₉	86.50 ₂₆
<i>Data from this work</i>							
W2NER (2022)[BERT-large]†	355.4	87.17 ₁₁	87.70 ₁₉	87.43 ₁₁	85.78 ₃₀	87.81 ₂₄	86.77 ₂₁
Ours[BERT-large]	345.1	87.98 ₃₀	87.50 ₂₂	87.74 ₁₆	86.26 ₆₅	87.56 ₃₁	86.91 ₂₃
w.o. CNN[BERT-large]	343.6	86.60 ₆₈	86.48 ₃₆	86.54 ₁₉	84.91 ₃₄	87.39 ₂₆	86.13 ₃₀
BS (2022)[RoBERTa-base]†	125.6	87.32 ₄₀	86.84 ₁₆	87.08 ₂₄	86.58 ₃₈	87.84 ₅₉	87.20 ₃₂
Ours[RoBERTa-base]	125.6	87.33 ₄₁	87.29 ₂₅	87.31 ₁₆	86.70 ₂₉	88.16 ₅₄	87.42 ₂₆
w.o. CNN[RoBERTa-base]	125.2	86.09 ₃₆	86.88 ₂₃	86.48 ₁₇	85.17 ₆₇	88.03 ₃₅	86.56 ₃₈

Table 1: Experiment results and the number of parameters for different models in the ACE2004 and ACE2005 datasets. Models in the same block use the same data. The subscript means the standard deviation (e.g 87.73₁₈ means 87.73±0.18). † means our reproduction with their publicly available code.

where $W_s, W_e \in \mathbb{R}^{d \times h}$, h is the hidden size, $\text{MHBiAffine}(\cdot, \cdot)$ is the multi-head Biaffine decoder², and $\mathbf{R} \in \mathbb{R}^{n \times n \times r}$, r is the feature size. Each cell (i, j) in the \mathbf{R} can be seen as the feature vector $\mathbf{v} \in \mathbb{R}^r$ for the span. And for the lower triangle of \mathbf{R} (where $i > j$), the span contains words from the j -th to the i -th (Therefore, one span will have two entries if its length is larger than 1).

2.3 CNN on Feature Matrix

As shown in Figure 1, the cell has relations with cells around. Therefore, we propose using CNN to model these interactions. We repeat the following CNN block several times in our model:

$$\begin{aligned} \mathbf{R}' &= \text{Conv2d}(\mathbf{R}), \\ \mathbf{R}'' &= \text{GeLU}(\text{LayerNorm}(\mathbf{R}' + \mathbf{R})), \end{aligned}$$

where Conv2d, LayerNorm and GeLU are the 2D CNN, layer normalization (Ba et al., 2016) and GeLU activation function (Hendrycks and Gimpel, 2016). The layer normalization is conducted in the feature dimension. A noticeable fact here is that since the number of tokens n in sentences varies, their \mathbf{R} s are of different shape. To make sure results are the same when \mathbf{R} is processed in batch, the 2D CNN has no bias term, and all the paddings in \mathbf{R} are filled with 0.

²The detailed description is in the Appendix A.1.

2.4 The Output

We use a perceptron to get the prediction logits P as follows:³

$$P = \text{Sigmoid}(W_o(\mathbf{R} + \mathbf{R}'') + b),$$

where $W_o \in \mathbb{R}^{|T| \times r}$, $b \in \mathbb{R}^{|T|}$, $P \in \mathbb{R}^{n \times n \times |T|}$. And then, we use golden labels y_{ij} and the binary cross entropy to calculate the loss as:

$$\mathbb{L}_{BCE} = - \sum_{0 \leq i, j < n} y_{ij} \log(P_{ij}),$$

More special details about our proposed method during training and inference procedure are described in Appendix A.

3 Experiment

3.1 Experimental Setup

To verify the effectiveness of our proposed method, we conduct experiments in three widely used nested NER datasets, ACE 2004⁴ (Doddington et al., 2004), ACE 2005⁵ (Walker and Consortium, 2005) and Genia (Kim et al., 2003).

³We did not use the Softmax because in the very rare case (such as in the ACE2005 and Genia dataset), one span can have more than one entity tag.

⁴<https://catalog.ldc.upenn.edu/LDC2005T09>

⁵<https://catalog.ldc.upenn.edu/LDC2006T06>

Besides, we choose recently published papers as our baselines. To make sure our experiments are strictly comparable to theirs, we ask the authors for their versions of data. The data statistics for each paper are listed in the Appendix B. For ACE2004 and ACE2005, although all of them use the same document split as suggested (Lu and Roth, 2015), they use different sentence tokenizations, resulting in different numbers of sentences and entities. To facilitate future research on nested NER, we release the pre-processing code and fix some tokenization issues to avoid including unannotated text and dropping entities. While for the Genia data, there are some annotation conflicts. For examples, one document with the bibliomisc MEDLINE:97218353 is duplicated in the original data, and different work has different annotations on it. We fix these conflicts. We replicate each experiment five times and report its average performance with standard derivation.

	# Param. (Million)	Genia		
		P	R	F1
<i>Data from Li et al. (2022)</i>				
W2NER (2022)	113.6	83.10	79.76	81.39
Ours	112.6	83.18 ₂₄	79.70 ₈	81.40 ₁₁
w.o. CNN	111.1	80.66 ₄	79.76 ₇	80.21 ₅
<i>Data from Wan et al. (2022)</i>				
SG (2022)	112.7	77.92	80.74	79.30
Ours	112.2	81.05 ₄₈	77.87 ₆₅	79.42 ₂₀
w.o. CNN	111.1	78.60 ₄₁	78.35 ₅₂	78.47 ₁₆
<i>Data from Yuan et al. (2022)</i>				
TriAffine (2022)	526.5	80.42	82.06	81.23
Ours	128.4	83.37 ₉	79.43 ₁₅	81.35 ₈
w.o. CNN	111.1	80.87 ₂₃	79.47 ₂₃	80.16 ₁₆
<i>Data from this work</i>				
W2NER [†]	113.6	81.58 ₆₁	79.11 ₄₉	80.32 ₂₃
Ours	112.6	81.52 ₂₁	79.17 ₁₈	80.33 ₁₃
w.o. CNN	111.1	78.59 ₂₈	79.85 ₁₄	79.22 ₁₂

Table 2: Experiment results and the number of parameters for different models in the Genia Dataset. All models use the BioBERT-base (Lee et al., 2020) as encoder. Models in the same block use the same data. The subscript means the standard deviation (e.g 81.40₁₁ means 81.40±0.11). † means our reproduction with their publicly available code.

3.2 Main Results

Results for ACE2004 and ACE2005 are listed in Table 1, and results for Genia is listed in Table 2. When using the same data from previous work, our simple CNN model surpasses the baselines with less or similar number of parameters, which proves

that using CNN to model the interaction between neighbor spans can be beneficial to the nested NER task. Besides, in the bottom block, we reproduce some baselines in our newly processed data to facilitate future comparison. Comparing the last block (processed by us) and the upper blocks (data from previous work), different tokenizations can indeed influence the performance. Therefore, we appeal for the same tokenization for future comparison.

	FEPR	FERE	NEPR	NERE
<i>ACE2004</i>				
Ours	86.9 _{0.2}	87.3 _{0.5}	88.4 _{0.6}	88.8 _{0.9}
w.o. CNN	86.3 _{0.8}	86.8 _{0.3}	89.4 _{0.8}	86.6 _{1.3}
<i>ACE2005</i>				
Ours	86.2 _{0.6}	88.3 _{0.1}	91.4 _{0.5}	89.0 _{0.8}
w.o. CNN	85.2 _{0.7}	87.9 _{0.3}	91.3 _{0.5}	86.2 _{0.8}
<i>Genia</i>				
Ours	81.7 _{0.2}	79.4 _{0.2}	71.7 _{1.6}	75.5 _{1.3}
w.o. CNN	79.0 _{0.3}	80.0 _{0.1}	72.7 _{1.2}	64.8 _{1.0}

Table 3: The precision and recall for flat and nested entities in the test set of three datasets. Compared with models without CNN (“w.o. CNN”), the most improved metric is bold. By using CNN, the recall for nested entities improve significantly. The subscript means the standard deviation (e.g 88.8_{0.9} means 88.8±0.9).

3.3 Why CNN Helps

To study why CNN can boost the performance of the nested NER datasets, we split entities into two kinds. One kind is entities that overlap with other entities, and the other kind is entities that do not. We design 4 metrics NEPR, NERE, FEPR and FERE, which are flat entity precision, flat entity recall, nested entity precision and nested entity recall, respectively.⁶, and list the results in Table 3. Compared with models without CNN, the NERE with CNN improve for 2.2, 2.8 and 10.7 on ACE2004, ACE2005 and Genia respectively. Namely, much of the performance improvement can be ascribed to finding more nested entities. This is expected as the CNN can be more effective for exploiting the neighbor entities when they are nested.

4 Related Work

Previously, four kinds of paradigms have been proposed to solve the nested NER task.

The first one is the sequence labeling framework (Straková et al., 2019), since one token can be

⁶The detailed calculation description of the 4 metrics locates in the Appendix D.

contained in more than one entities, the Cartesian product of the entity labels are used. However, the Cartesian labels will suffer from the long-tail issue.

The second one is to use the hypergraph to efficiently represent spans (Lu and Roth, 2015; Muis and Lu, 2016; Katiyar and Cardie, 2018; Wang and Lu, 2018). The shortcoming of this method is the complex decoding.

The third one is the sequence-to-sequence (Seq2Seq) framework (Sutskever et al., 2014; Lewis et al., 2020; Raffel et al., 2020) to generate the entity sequence. The entity sequence can be the entity pointer sequence (Yan et al., 2021; Fei et al., 2021) or the entity text sequence (Lu et al., 2022). Nevertheless, the Seq2Seq method suffers from the time-demanding decoding.

The fourth one is to conduct span classification. Eberts and Ulges (2020) proposed to enumerate all possible spans within a sentence, and use a pooling method to get the span representation. While Yu et al. (2020) proposed to use the start and end tokens of a span to pinpoint the span, and use the Biaffine decoder to get the scores for each span. The span-based methods are friendly to parallelism and the decoding is easy. Therefore, this formulation has been widely adopted (Wan et al., 2022; Zhu and Li, 2022; Li et al., 2022; Yuan et al., 2022). However, the relation between neighbor spans was ignored in previous work.

5 Conclusion

In this paper, we propose using CNN on the score matrix of span-based NER model. Although this method is very simple, it achieves comparable or better performance than recently proposed methods. Analysis shows exploiting the spatial correlation between neighbor spans through CNN can help model find more nested entities. And experiments show that different tokenizations indeed influence the performance. Therefore, it is necessary to make sure all comparative baselines use the same tokenization. To facilitate future comparison, we release a new pre-processing script for three nested NER datasets.

Limitations

While we discover that simply applying CNN on top of the score matrix of span-based NER model performs well on the nested NER scenario, there are still some limitations that are worth discussing. Firstly, we mainly choose three commonly used

nested NER datasets, which may lack generalization. Secondly, we only focus on nested NER tasks for the spatial relations between spans are more intuitive and common in nested scenario than those in flat NER. However, the principle of using CNN to model the relations is also applicable to spans in the flat NER task. Future work can take flat NER into consideration based on our exploration, and experiments on more datasets.

Acknowledgements

We would like to thank the anonymous reviewers for their insightful comments. We also thank the developers of fastNLP⁷ and fitlog⁸. This work was supported by the National Natural Science Foundation of China (No. 62236004 and No. 62022027) and CCF-Baidu Open Fund.

References

- Lei Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. [Layer normalization](#). *CoRR*, abs/1607.06450.
- 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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.
- George R. Doddington, Alexis Mitchell, Mark A. Przybicki, Lance A. Ramshaw, Stephanie M. Strassel, and Ralph M. Weischedel. 2004. [The automatic content extraction \(ACE\) program - tasks, data, and evaluation](#). In *Proceedings of the Fourth International Conference on Language Resources and Evaluation, LREC 2004, May 26-28, 2004, Lisbon, Portugal*. European Language Resources Association.
- Timothy Dozat and Christopher D. Manning. 2017. [Deep biaffine attention for neural dependency parsing](#). In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- Markus Eberts and Adrian Ulges. 2020. [Span-based joint entity and relation extraction with transformer](#)
-
- ⁷<https://github.com/fastnlp/fastNLP>. FastNLP is a natural language processing python package.
- ⁸<https://github.com/fastnlp/fitlog>. Fitlog is an experiment tracking package.

- pre-training. In *ECAI 2020 - 24th European Conference on Artificial Intelligence, 29 August-8 September 2020, Santiago de Compostela, Spain, August 29 - September 8, 2020 - Including 10th Conference on Prestigious Applications of Artificial Intelligence (PAIS 2020)*, volume 325 of *Frontiers in Artificial Intelligence and Applications*, pages 2006–2013. IOS Press.
- Hao Fei, Donghong Ji, Bobo Li, Yijiang Liu, Yafeng Ren, and Fei Li. 2021. [Rethinking boundaries: End-to-end recognition of discontinuous mentions with pointer networks](#). In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 12785–12793. AAAI Press.
- Dan Hendrycks and Kevin Gimpel. 2016. [Bridging non-linearities and stochastic regularizers with gaussian error linear units](#). *CoRR*, abs/1606.08415.
- Zhiheng Huang, Wei Xu, and Kai Yu. 2015. [Bidirectional LSTM-CRF models for sequence tagging](#). *CoRR*, abs/1508.01991.
- Arzoo Katiyar and Claire Cardie. 2018. [Nested named entity recognition revisited](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 861–871. Association for Computational Linguistics.
- Jin-Dong Kim, Tomoko Ohta, Yuka Tateisi, and Jun’ichi Tsujii. 2003. [GENIA corpus - a semantically annotated corpus for bio-textmining](#). In *Proceedings of the Eleventh International Conference on Intelligent Systems for Molecular Biology, June 29 - July 3, 2003, Brisbane, Australia*, pages 180–182.
- Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. [Biobert: a pre-trained biomedical language representation model for biomedical text mining](#). *Bioinform.*, 36(4):1234–1240.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 7871–7880. Association for Computational Linguistics.
- Jingye Li, Hao Fei, Jiang Liu, Shengqiong Wu, Meishan Zhang, Chong Teng, Donghong Ji, and Fei Li. 2022. [Unified named entity recognition as word-word relation classification](#). In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelfth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, pages 10965–10973. AAAI Press.
- Xiaoya Li, Jingrong Feng, Yuxian Meng, Qinghong Han, Fei Wu, and Jiwei Li. 2020. [A unified MRC framework for named entity recognition](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 5849–5859. Association for Computational Linguistics.
- Wei Lu and Dan Roth. 2015. [Joint mention extraction and classification with mention hypergraphs](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 857–867. The Association for Computational Linguistics.
- Yaojie Lu, Qing Liu, Dai Dai, Xinyan Xiao, Hongyu Lin, Xianpei Han, Le Sun, and Hua Wu. 2022. [Unified structure generation for universal information extraction](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 5755–5772. Association for Computational Linguistics.
- Xuezhe Ma and Eduard H. Hovy. 2016. [End-to-end sequence labeling via bi-directional lstm-cnns-crf](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics.
- Aldrian Obaja Muis and Wei Lu. 2016. [Learning to recognize discontinuous entities](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pages 75–84. The Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *J. Mach. Learn. Res.*, 21:140:1–140:67.
- Jana Straková, Milan Straka, and Jan Hajic. 2019. [Neural architectures for nested NER through linearization](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 5326–5331. Association for Computational Linguistics.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. [Sequence to sequence learning with neural networks](#). In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 3104–3112.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. **Attention is all you need**. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.

C. Walker and Linguistic Data Consortium. 2005. *ACE 2005 Multilingual Training Corpus*. LDC corpora. Linguistic Data Consortium.

Juncheng Wan, Dongyu Ru, Weinan Zhang, and Yong Yu. 2022. **Nested named entity recognition with span-level graphs**. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 892–903. Association for Computational Linguistics.

Bailin Wang and Wei Lu. 2018. **Neural segmental hypergraphs for overlapping mention recognition**. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 204–214. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. **Transformers: State-of-the-art natural language processing**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Hang Yan, Bocao Deng, Xiaonan Li, and Xipeng Qiu. 2019. **TENER: adapting transformer encoder for named entity recognition**. *CoRR*, abs/1911.04474.

Hang Yan, Tao Gui, Junqi Dai, Qipeng Guo, Zheng Zhang, and Xipeng Qiu. 2021. **A unified generative framework for various NER subtasks**. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 5808–5822. Association for Computational Linguistics.

Juntao Yu, Bernd Bohnet, and Massimo Poesio. 2020. **Named entity recognition as dependency parsing**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 6470–6476. Association for Computational Linguistics.

Zheng Yuan, Chuanqi Tan, Songfang Huang, and Fei Huang. 2022. **Fusing heterogeneous factors with triaffine mechanism for nested named entity recognition**. In *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May*

22-27, 2022

, pages 3174–3186. Association for Computational Linguistics.

Enwei Zhu and Jinpeng Li. 2022. **Boundary smoothing for named entity recognition**. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 7096–7108. Association for Computational Linguistics.

A Detailed Proposed Method

A.1 Multi-head Biaffine Decoder

The input of Multi-head Biaffine decoder is two matrices $\mathbf{H}_s, \mathbf{H}_e \in \mathbb{R}^{n \times h}$, and the output is $\mathbf{R} \in \mathbb{R}^{n \times n \times r}$. The formulation of Multi-head Biaffine decoder is as follows

$$\begin{aligned} \mathbf{S}_1[i, j] &= (\mathbf{H}_s[i] \oplus \mathbf{H}_e[j] \oplus \mathbf{w}_{|i-j|})W, \\ \{\mathbf{H}_s^{(k)}\}, \{\mathbf{H}_e^{(k)}\} &= \text{Split}(\mathbf{H}_s), \text{Split}(\mathbf{H}_e), \\ \mathbf{S}_2^{(k)}[i, j] &= \mathbf{H}_s^{(k)}[i]U\mathbf{H}_e^{(k)}[j]^T, \\ \mathbf{S}_2 &= \text{Concat}(\mathbf{S}_2^{(1)}, \dots, \mathbf{S}_2^{(K)}), \\ \mathbf{R} &= \mathbf{S}_1 + \mathbf{S}_2, \end{aligned}$$

where $\mathbf{H}_s, \mathbf{H}_e \in \mathbb{R}^{n \times h}$, h is the hidden size, $\mathbf{w}_{|i-j|} \in \mathbb{R}^c$ is the span length embedding for length $|i - j|$, $W \in \mathbb{R}^{(2h+c) \times r}$, $\mathbf{S}_1 \in \mathbb{R}^{n \times n \times r}$, r is the biaffine feature size, $\text{Split}(\cdot)$ equally splits a matrix in the last dimension, thus, $\mathbf{H}_s^{(k)}, \mathbf{H}_e^{(k)} \in \mathbb{R}^{n \times h_k}$; h_k is the hidden size for each head, and $U \in \mathbb{R}^{h_k \times r_k \times h_k}$, $\mathbf{S}_2 \in \mathbb{R}^{n \times n \times r}$, and $\mathbf{R} \in \mathbb{R}^{n \times n \times r}$.

We do not use multi-head for W , because it does not occupy too many parameters and using multi-head for W harms the performance slightly.

A.2 Training Loss

Unlike previous works that only use the upper triangle part to get the loss (Yu et al., 2020; Zhu and Li, 2022), we use both upper and lower triangles to calculate the loss, as depicted in section 2.4. The reason is that in order to conduct batch computation, we cannot solely compute features from the upper triangle part. Since features from the lower triangle part have been computed, we also use them for the output. The tag for the score matrix is symmetric, namely, the tag in the (i, j) -th entry is the same as that in the (j, i) -th.

		Sentence				Mention				
		#Train	#Dev	#Test	Avg. Len	#Ovlp.	#Train	#Dev	#Test	Avg. Len
ACE2004	W2NER	6,802	813	897	20.12	12,571	22,056	2,492	3,020	2.5
	SG	6,198	742	809	21.55	12,666	22,195	2,514	3,034	2.51
	BS	6,799	829	879	20.43	12,679	22,207	2,511	3,031	2.51
	Ours	6,297	742	824	23.52	12,690	22,231	2,514	3,036	2.64
ACE2005	W2NER	7,606	1,002	1,089	17.77	12,179	24,366	3,188	2,989	2.26
	SG	7,285	968	1,058	18.60	12,316	24,700	3,218	3,029	2.26
	BS	7,336	958	1,047	18.90	12,313	24,687	3,217	3,027	2.26
	Ours	7,178	960	1,051	20.59	12,405	25,300	3,321	3,099	2.40
Genia	W2NER	15,023	1,669	1,854	25.41	10,263	45,144	5,365	5,506	1.97
	SG	15,022	1,669	1,855	26.47	10,412	47,006	4,461	5,596	2.07
	Triaffine	16,692	-	1,854	25.41	10,263	50,509	-	5,506	1.97
	Ours	15,038	1,765	1,732	26.47	10,315	46,203	4,714	5,119	2.0

Table 4: The statistics used in each paper. “W2NER”⁹, “SG”, “BS” and “Triaffine” are from Li et al. (2022), Wan et al. (2022), Zhu and Li (2022) and Yuan et al. (2022), respectively. #Ovlp. means the number of overlapping mentions. Different papers use different sentence tokenization for ACE2004 and ACE2005, resulting in different numbers of sentences in each split. To facilitate future comparison, we release a pre-processing script to prepare ACE2004 and ACE2005. Previously, some entities will be dropped because of sentence tokenization, we avoid sentence tokenization within an entity and resulting in more entities. And for Genia, different papers use different train/dev/test splits. Besides, the Genia data has conflicting annotations, we remove these sentences. The data annotated with “Ours” is obtained by our pre-processing code.

A.3 Inference

When inference, we calculate scores in the upper triangle part as:

$$\hat{P}_{ij} = (P_{ij} + P_{ji})/2,$$

where $i \leq j$. Then we only use this upper triangle score to get the final prediction. The decoding process generally follows Yu et al. (2020)’s method. We first prune out the non-entity spans (none of its scores is above 0.5), then we sort the remained spans based on their maximum entity score. We pick the spans based on this order, if a span’s boundary clashes with selected spans’, it is ignored.

B Data

We list the statistics for each dataset in Table 4.¹⁰ As shown in the table, the number of sentences and even the number of entities are different for each paper on the same dataset. Therefore, it is not fair to directly compare results. For the ACE2004 and ACE2005, we release the pre-processing code to get data from the LDC files. We make sure no entities are dropped because of the sentence tokenization. Thus, the pre-processed ACE2004 and ACE2005 data from this work in Table 4 have the most entities.

¹⁰The number of entities is different from that reported in their paper, because we found some duplicated sentences in their data.

And for Genia, we appeal for the usage of train/dev/test, and we release the data split within the code repository. Moreover, in order to facilitate the document-level NER study, we split the Genia dataset based on documents. Therefore, sentences from train/dev/test splits are from different documents, the document ratio for train/dev/test is 8:1:1. Besides, we find one conflicting document annotation in Genia, we fix this conflict. After comparing different versions of Genia, we find the W2NER (Li et al., 2022) and Triaffine (Yuan et al., 2022) drop the spans with more than one entity tags (there are 31 such entities). Thus, they have less number of nested entities than us. While SG (Wan et al., 2022) includes the discontinuous entities, so they have more number of nested entities than us.

C Implementation Details

We use the AdamW optimizer to optimize the model and the transformers package for the pre-trained model (Wolf et al., 2020). The hyperparameter range in this paper is listed in Table 5.

D FEPR FERE NEPR NERE

We split entities into two kinds based on whether they overlap with other entities, and the statistics for each dataset are listed in Table 6. When calculating the flat entity precision (FEPR), we first get all flat entities in the prediction and calculate their

	ACE2004	ACE2005	Genia
# Epoch	50	50	5
Learning Rate	2e-5	2e-5	7e-6
Batch size	48	48	8
# CNN Blocks	[2, 3]	[2, 3]	3
CNN kernel size	3	3	3
CNN Channel dim.	[120, 200]	[120, 200]	200
# Head	[1, 5]	[1, 5]	4
Hidden size h	200	200	400
Warmup factor	0.1	0.1	0.1

Table 5: The hyper-parameters in this paper.

	# Ent.	# Flat Ent.	# Nested Ent.
ACE2004	3,036	1,614	1,422
ACE2005	3,099	1,913	1,186
Genia	5,119	3,963	1,156

Table 6: The flat and nested entity statistics in the test set of each dataset.

ratio in the gold. For the flat entity recall (FERE), we get all flat entities in the gold and calculate their ratio in the prediction. And we get the nested entity precision (NEPR) and nested entity recall (NERE) similarly.

ACL 2023 Responsible NLP Checklist

A For every submission:

- A1. Did you describe the limitations of your work?
Section "Limitations" (5th section)
- A2. Did you discuss any potential risks of your work?
Section "Limitations" (5th section)
- A3. Do the abstract and introduction summarize the paper's main claims?
"Abstract" and section 1
- A4. Have you used AI writing assistants when working on this paper?
Left blank.

B Did you use or create scientific artifacts?

Left blank.

- B1. Did you cite the creators of artifacts you used?
No response.
- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?
No response.
- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?
No response.
- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?
No response.
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?
No response.
- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.
No response.

C Did you run computational experiments?

Left blank.

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?
No response.

The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

Section 3 and Appendix C

- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?

Section 3

- C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?

Appendix C

D Did you use human annotators (e.g., crowdworkers) or research with human participants?

Left blank.

- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?

No response.

- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?

No response.

- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?

No response.

- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

No response.

- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

No response.