

Topic Sentence Named Entity Recognition: A New Task with Its Dataset and Benchmarks

Anonymous ACL submission

Abstract

In this paper, we focus on a new type of named entity recognition (NER) task called topic sentence NER. A topic sentence means a short and compact sentence that acts as a summary of a long document. For example, a title can be seen as a topic sentence of its article. Topic sentence NER aims to extract named entities in a topic sentence given the corresponding unlabeled document as a reference. This task represents real-world scenarios where full-document NER is too expensive and obtaining the entities only in topic sentences is enough for downstream tasks. To achieve this, we construct a large-scale human-annotated Topic Sentence NER dataset, named TSNER. The dataset contains 12,000 annotated sentences accompanied by their unlabeled document. Based on TSNER, we propose a family of representative and strong baseline models, which can utilize both single-sentence and document-level features. We will make the dataset public in the hope of advancing the research on the topic sentence NER task.

1 Introduction

Named entity recognition is a fundamental Natural Language Processing task, which aims to label each word in sentences with predefined types, such as Person (PER), Organization (ORG), Location (LOC), etc. The results of NER play a crucial role in many downstream NLP tasks, e.g., relation extraction (Bunescu and Mooney, 2005), information retrieval (Chen et al., 2015), and question answering (Yao and Van Durme, 2014).

In this paper, we propose a new type of NER task named Topic Sentence NER, which attempts to recognize entities in topic sentences. A topic sentence is a key sentence for a document or a paragraph, which usually conveys the gist of them in a concise way. An example is shown in Figure 1. The task is defined to extract named entities like ‘悬崖之上(*Impasse*)’ in the topic sentence. The

significance of the topic sentence NER lies in two aspects. First, in many practical scenarios, it is not necessary to obtain all entities in a full-text document. Due to the time and cost of labeling and processing documents, topic sentence NER can be an effective alternative. Second, topic sentence NER is more challenging by nature and it requires new ways to incorporate the heterogeneous inputs. On the one hand, topic sentences are more informative but short in length, making the in-sentence context for NER limited. On the other hand, there are unlabeled documents that can potentially enrich the context of the topic sentence, but it is unclear how to effectively utilize the information for NER.

Given the realistic necessity and challenges of topic sentence NER, in this paper, we focus on addressing such a new kind of NER task. We construct a new dataset named TSNER, representing for Topic Sentence Named Entity Recognition. Specifically, we collect 12,000 online articles in Chinese. The articles are about 9 topics and contain entities of 16 types. For each article, we label the entities in its title and consider the title as the topic sentence of its accompanying document.

Based on the proposed dataset, we establish a family of strong baseline models as benchmarks for topic sentence NER. We consider two categories of models: single-sentence NER model and document enhanced NER model. 1) The former only uses the topic sentence as its input and consists of commonly used models that have achieved SOTA performance on many single-sentence NER datasets. 2) The latter takes both the topic sentence and its corresponding document into consideration. Two challenges have to be tackled for the document-enhanced NER model: capturing dependency-term dependency in a computational efficiency way and distinguishing information helpful for NER from a large unrelated, noisy text. Based on the analysis, we adapt three lines of work for document-enhanced NER: distant supervision,

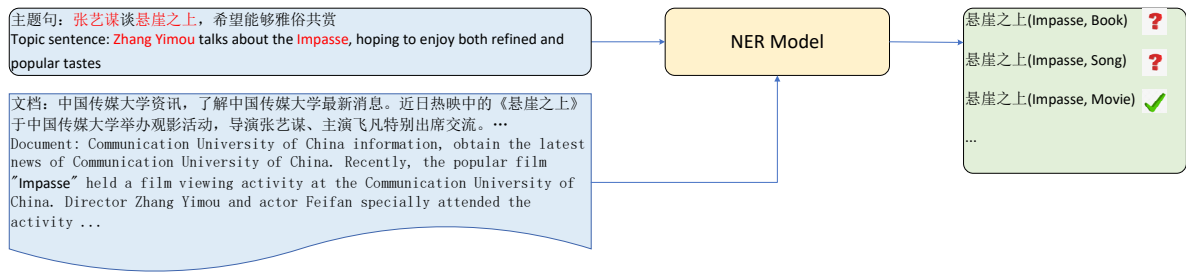


Figure 1: A case of topic sentence NER. The topic sentence is brief and it alone provides limited context. With the help of document information, ‘悬崖之上(*Impasse*)’ can be recognized as an entity of Movie type.

document-level language modeling, and information extraction and fusion.

To the best of our knowledge, this paper is the first to propose and address the topic sentence NER task. Our key contributions are as follows:

- We introduce topic sentence NER, a new NER task focusing on recognizing entities in topic sentences. This task is driven by real-world needs and is of particular research value.
- To better understand the topic sentence NER task, we propose the TSNER dataset, in which each annotated topic sentence is paired with an unlabeled document.
- Based on TSNER, we establish a family of benchmark models and conduct extensive experiments, revealing effective ways to leverage document information for this task.

2 Related work

2.1 Single Sentence NER

Previous works mainly consider the NER task as a single sentence task. Traditional methods try to build the single sentence feature manually and use the CRF model to process the feature (Lafferty et al., 2001). With the advantages of significant performance improvement and eliminating feature engineering, neural network models become prevalent in NER research recently, e.g. FFN (Collobert et al., 2011), LSTM (Lample et al., 2016), CNN (Ma and Hovy, 2016), and pre-trained language model (Devlin et al., 2019). The single sentence NER model can better handle the situation when the entity has abundant context information, which is not satisfied in the topic sentence NER task.

2.2 Document-level NER

Document-level NER extends single-sentence NER to recognize all entities in the whole document.

Gui et al. (2020) introduces a two-stage label refinement approach to improve document-level label consistency. Luoma and Pyysalo (2020) explores the use of cross-sentence information for NER based on BERT. Akbik et al. (2018); Luo et al. (2020) attempts to use a memory network to better address the long-term dependency problem in the document. However, it is hard to apply document NER methods directly to our topic sentence NER task as the document sentences are unlabeled. Besides, the concise writing style of topic sentences makes the task even more challenging.

2.3 Other Document-level NLP models

Our work is also related to other document-level NLP tasks, such as document-level classification, question answering, and coreference resolution. Existing approaches to modeling document information can be summarized into three categories. The first is to chunk a document into smaller pieces of text to be independently processed by single-sentence models, and then to combine the results through a fusion network (Joshi et al., 2019). The second is to shorten the document by selecting only the informative parts of it as the input of the model (Clark and Gardner, 2018; Chen et al., 2017). The third is to develop new model architecture to accommodate the whole document (Beltagy et al., 2020; Gupta and Berant, 2020; Zaheer et al., 2020). Our baseline models for document enhanced NER are derived from these three types of models.

3 Topic Sentence NER

In real-world situations, the results of NER are often used in downstream tasks like relation extraction, information retrieval, and question answering. In these applications, the requirement to recognize all entities in a full-text document is not always necessarily essential, and recognizing entities only

in topic sentences is enough, especially when huge amounts of text have to be processed with a limit of time and cost. For example, the entities in the abstract of a scientific paper are enough for an up-to-date scholar search engine; the entities in a news title are enough for hot event detection and trend analysis. However, such a need for entity recognition on topic sentences has not been put forward and explored in previous NER research.

Compared with regular sentences or documents involved in previous NER tasks, topic sentences exhibit unique linguistic characteristics that makes the NER more challenging. Specifically, topic sentences are often short in length but more informative in that it contains a higher density of entity words. Take the topic sentence shown in Figure 1 as an example. The number of words belonging to entities exceeds 40% of the total number of tokens. Consequently, the word ‘悬崖之上(*Impasse*)’ has a limited context and is difficult to be distinguished as a book, a song, a movie, or a non-entity word. Furthermore, while document can incorporated to enrich the context of topic sentences, there are no ground truth NER labels for the sentences in the document, making previous document-level NER models inapplicable. This calls for a new research direction of context limited and document enhanced NER methods.

Given the realistic necessity and challenges of topic sentence NER, in the remainder of this paper, we will show our initial attempt to address this problem. We will first give the definition of topic sentence NER. Then we will present our constructed dataset and analysis on it. Finally, we will propose a series of benchmark models and compare their experiment results. To the best of our knowledge, this paper is the first to propose and address the topic sentence NER task.

3.1 Task Definition

We formally define topic sentence NER as a sequence labeling task on a topic sentence accompanied by an unlabeled document. The input of topic sentence NER consists of two parts: a topic sentence $x = \{x_1, x_2, \dots, x_t\}$ and an unlabeled document $D = \{s_1, s_2, \dots, s_n\}$. The goal of the task is to assign each token $x_i \in x$ with a label $y_i \in Y$. Y is a set of pre-defined entity tags in BIO or other format.

3.2 Dataset Construction

The data source we used as an initial corpus is a collection of news articles in Chinese, which contains a large variety of entities from different areas. We selected 12,000 articles on nine topics, including tourism, sports, politics, food, culture, economy, movies, entertainment, and games. We designed a NER scheme consisting of 16 commonly used entity types. The names and distribution of the entity types are shown in Table 1. More details of the dataset will be shown in Appendix A and Github ¹.

We employed paid annotators to annotate the dataset. We sent the titles along with the articles to the annotators and instructed them to annotate the entities in the titles with the reference to the documents. All the decisions are made based on the title and article together. We find in many cases the title alone can not be understood by the human at a first glance. After scanning the document, however, one can confidently label the entities in the title. All the annotators are instructed with detailed and formal annotation guidelines to have adequate linguistic knowledge of each entity type. To ensure the quality of TSNER, we randomly selected 10% of the data and examined the results by ourselves. If the sentence-level accuracy of the annotation is lower than 90%, the batch will be re-annotated.

3.3 Dataset Profile

We report some interesting statistics of our dataset compared with several widely-used NER datasets including MSRA (Levow, 2006), OntoNotes (Weischedel et al., 2013), WeiboNER (Peng and Dredze, 2015; He and Sun, 2017)². We calculated the average length and entity rate for each dataset. The results is shown in Table 2. Two unique characteristics of topic sentences can be revealed, as follows.

1) Shorter length: The average length of the topic sentence is 22 and only half of the MSRA dataset. The NER of the short sentence is more complex for less information.

2) More informative: In the topic sentence NER, the rate of entity token accounts for the whole token is 30, which means that the sentence has less context information, which makes more hard for the NER. Besides, Our topic sentences are

¹XXX

²For datasets with multiple languages, we only analyze the part in Chinese. In the future, we will extend our work to other languages.

not post-processed, whereas other datasets often filter the sentences that do not contain entities. It can also demonstrate that the topic sentence contains more information.

The short sentence and high information rate make the topic sentence NER more challenging, and it is important to introduce the document information to help the topic sentence NER. Similarly, the statistics for documents are shown in Table 3. We can further draw the following two challenges to be solved.

1) Long document length: Compared with previous widely used datasets, TSNER provides a long unlabeled document, the length of the document is 1386, which means that the document contains a large noise, and we are required to extract the important information to help the topic NER.

2) Relatedness to topic sentence: The document is highly related to the title. The rate of entity both appear in document and topic sentence accounts for the whole entity is 80%, which means that the correction of the title and document is high, the document can provide useful information for NER.

| Type | Num/Rate | Type | Num/Rate |
|-------------|------------|--------|----------|
| address | 1889 (15%) | name | 630 (5%) |
| ename | 1648 (13%) | book | 622 (5%) |
| food | 1100 (9%) | tvplay | 610 (5%) |
| event | 1087 (8%) | show | 537 (4%) |
| aname | 994 (8%) | scene | 428 (3%) |
| orgnization | 853 (7%) | song | 380 (3%) |
| company | 622 (6%) | gname | 270 (2%) |
| movie | 699 (5%) | game | 259 (2%) |

Table 1: The distribution of different entity types in TSNER train part.

| | TSAvgLen | EntRate | Doc |
|-----------|----------|---------|-----|
| MSRA | 47 | 12.3 | No |
| OntoNotes | 31 | 9.1 | No |
| Weibo NER | 55 | 4.5 | No |
| TSNER | 22 | 30.0 | Yes |

Table 2: A comparison between TSNER and other existing widely-used NER datasets. TSAvgLen means Topics Sentence Average length, and EntRate means the rate of entity token accounts for the whole token.

| | Train | Dev | Test |
|-----------------|--------|-------|-------|
| #sen | 8400 | 1800 | 1800 |
| #char | 185.4k | 38.9k | 39.5k |
| #entity | 12.8k | 2.6k | 2.6k |
| doc avg len | 1386 | 1344 | 1377 |
| Entity Doc Rate | 79.3 | 79.1 | 80.2 |

Table 3: The statistics of TSNER. Entity Doc Rate means the rate of entity both appear in document and topic sentence accounts for the whole entity.

4 Benchmarks

Based on the TSNER, we develop a family of representative and strong baselines. We first present single sentence NER models in Section 4.1. Then we introduce document-enhanced NER models in Section 4.2. The single sentence NER only uses topic sentence as input, while the document-enhanced NER can use both topic sentence and document.

4.1 Single-sentence NER models

BiLSTM-CRF. BiLSTM-CRF (Lample et al., 2016) is a strong baseline that has been widely used in previous works.

Softlexicon. In Chinese NER, explicitly providing word segmentation and word tagging information can be potentially helpful. A series of models have been proposed based on this motivation (Zhang and Yang, 2018; Yang et al., 2019; Li et al., 2020; Ma et al., 2020; Liu et al., 2021). Among them we choose the SoftLexicon (Ma et al., 2020) as our baseline due to its fast speed and competitive performance.

BERT-CRF. The BERT-CRF (Devlin et al., 2019) baseline is chosen as a representative for NER models based on pre-trained language models (PLMs).

WWM-CRF. PLMs share the same problem with other models when processing Chinese text. In order to take into account lexical information, PLMs with enhanced input layers and training techniques have been proposed (Cui et al., 2019, 2020; Diao et al., 2020; Sun et al., 2021). We choose the WWM model (Cui et al., 2019) for its popularity and proved generalization ability.

4.2 Document-enhanced NER Models

Distant supervision. A natural way to leverage the unlabeled document data is to regard it as an

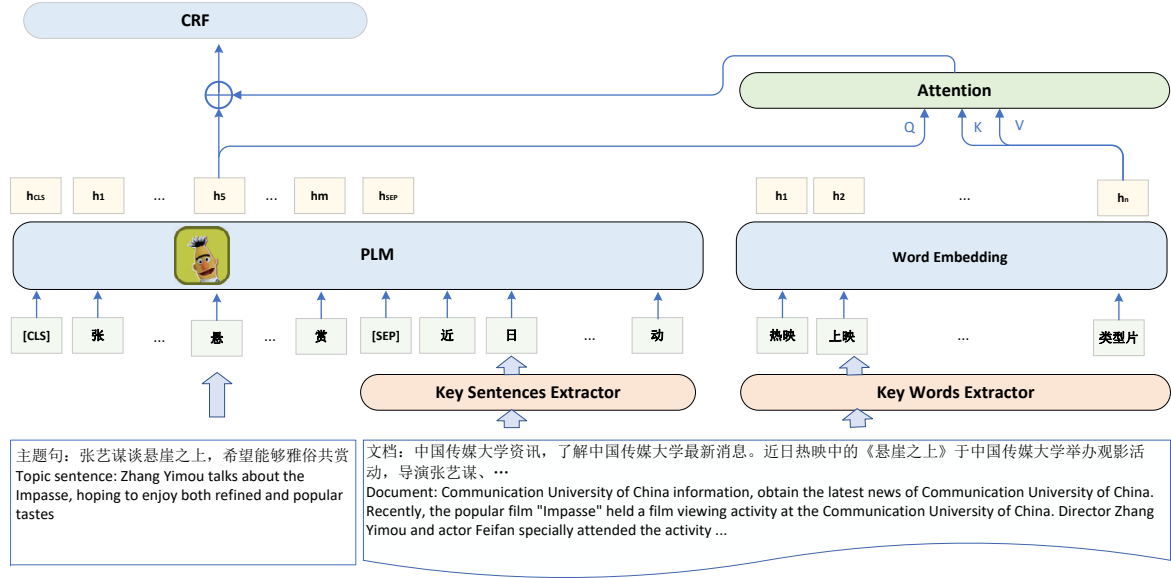


Figure 2: Model architecture of our document gist fusion model. The extracted gist information includes key sentences. The key sentences are encoded together with the topic sentence to provide extra context. The embeddings of the keywords are fused into the hidden states of the topic sentence using an attention mechanism.

in-domain corpus for distantly supervised learning. To do so, we first curated an entity dictionary by extracting all the annotated entities in the train set of TSNER. Then, we use the entity dictionary to match sentences in the documents to obtain distantly supervised data. Finally, the distantly supervised data and human annotated data are mixed together as the training data for BERT-CRF or WWM-CRF. We denote the two models as **BERT-CRF-DS** and **WWM-CRF-DS**. There are dedicated methods to reduce the noise in distantly supervised data that can be explored in the future.

Document-level PLM. Document-level PLMs are supposed to accommodate full document as input and automatically learn to properly utilize its information for downstream tasks. In recent work, several models have been proposed to reduce memory and speed up the training of transformer models (Beltagy et al., 2020; Gupta and Berant, 2020; Zaheer et al., 2020). In this paper, we build NER model for topic sentence based on Longformer (Beltagy et al., 2020), whose attention mechanism is a drop-in replacement for the standard self-attention and combines a local windowed attention with a task motivated global attention. The topic sentence is prepend to the document as the input of Longformer and the global attention is applied on the topic sentence. Finally, we use the output of the Longformer as the input of CRF.

Document gist fusion. While document-level PLMs can encode a full-text document, not all words in the document are helpful for the topic sentence NER task. Incorporating too much unrelated information will bring noise in training. Based on the observation, we propose a document gist fusion model for topic sentence NER. The idea is to first extract gist information from the document using heuristic approaches, and then fuse the gist information into the NER process. We will first describe the model design. The methods to extract gist information will be discussed in the next subsection.

The model architecture is shown in Figure 2. We consider two forms of gist information, i.e., key sentences and keywords. Compared with the document, the key sentences are short enough and can be easily fed into a transformer model. Hence, we append the key sentences to the topic sentence as an additional input to a PLM encoder:

$$H^s = \text{PLM}([x; S])_{[1:m]} \quad (1)$$

where x is the topic sentence with length m , S is the set of selected key sentences from the document. $H^s = \{h_1^s, h_2^s, \dots, h_m^s\}$ is the hidden states of the topic sentence, which corresponds to the first m tokens of the inputs and is augmented with the extra context of the key sentences.

As only a few key sentences are extracted from the document, they may be not enough to cover

all necessary information for recognizing the entities in the topic sentence. We also consider including keywords as a global context that indicates the document’s topic. The keywords are encoded separately by a word embedding layer.

$$H^w = \text{WordEmb}(w) \quad (2)$$

where w is the set of n selected keywords and $H^w = \{h_1^w, h_2^w, \dots, h_n^w\}$ is the embedding for each keyword.

We use an attention network to better modeling the relation between the sentence-level information H^s and the keywords information H^w . The attention mechanism is similar to the attention in Vaswani et al. (2017). We transform $h_i^s \in H^s$ into the attention query q_i , and keywords embedding into both the key k_j and the value v_j , where q_i, k_j , and v_j are in the same dimension. The calculations of the attention layer are as follows:

$$q_i = W^s h_i^s \quad (3)$$

$$k_j = W^w h_j^w \quad (4)$$

$$v_j = W^v h_j^w \quad (5)$$

$$u_{ij} = q_i k_j \quad (6)$$

$$\alpha_{ij} = \frac{\exp(u_{ij})}{\sum_{z=1}^n \exp(u_{iz})} \quad (7)$$

$$r_i = \sum_{j=1}^n \alpha_{ij} v_j \quad (8)$$

Concatenating q_i and r_i we obtain a fused representation of the topic sentences and the gist of the document:

$$f_i = [q_i; r_i] \quad (9)$$

Then f_i will be fed into a CRF layer to output entity labels. Next, we will elaborate on how we designed efficient heuristics to select key sentences and keywords from the document.

4.3 Key sentence and keyword selection

We explore several methods to select the key sentences and key words for our gist fusion model. The key sentence selection is to select a maximum number of N sentences from the document. In order to provide adequate context with a reasonable cost of longer input length, we empirically set $N = 5$ in our study.

First in order. In this strategy, we simply take the sentences from the beginning of the document.

Similarity-based. The idea of this strategy is to select sentences that are semantically similar to the topic sentence based on a similarity metric. Two similarity metrics for sentences are considered. One is Word Mover’s Distance (WWD) (Kusner et al., 2015) based on word embedding. The other is pretrained SBERT (Reimers and Gurevych, 2019), which derives semantic aware sentence embedding from a Siamese BERT network and uses cosine similarity to measure similarity between them.

Noun overlapping. We propose a simple method to select the key sentence based on the co-occurrence of noun words in a topic sentence and its document. Sharing common nouns means that two sentences have a closer relationship, and that they together form a richer context for the common nouns. Specifically, we scan the sentences of the document in the natural order and pick out sentences that share at least one common noun with the topic sentence. In order to increase the diversity, we limit the number of sentences that each noun can associate with to two. When the limit is exceeded, only the two sentences that are more close to the beginning in the document will be kept.

For keyword extraction, we explore the following two methods.

TextRank (Mihalcea and Tarau, 2004) is a graph-based word ranking model inspired by PageRank. It is widely used for selecting informative words from a document.

Yake (Campos et al., 2020) is a more recent and lightweight approach for keyword extraction, which uses statistical features to measure the importance of each word in a document.

By combining the above key sentence and keyword selection methods with the model architecture in Figure 2, we expand our benchmarks with a series of document gist fusion models. The keywords information is only added to PLM-Noun models, which we found in our pilot study to achieve better result.

5 Results and Analysis

In this section, we report the results of various experiments carried on the TSNER dataset. Following the evaluation metrics in previous NER research, we report results in terms of entity-level

| Model | Resource | DEV | | | Test | | |
|-------------------|----------|--------------|--------------|--------------|--------------|--------------|--------------|
| | | P | R | F | P | R | F |
| BiLSTM-CRF | TS | 61.10 | 59.16 | 60.12 | 60.08 | 59.97 | 60.03 |
| SoftLexicon | TS | 69.59 | 59.62 | 64.22 | 70.64 | 61.27 | 65.62 |
| BERT-CRF | TS | 78.06 | 76.69 | 77.37 | 77.49 | 77.62 | 77.56 |
| WWM-CRF | TS | 78.11 | 76.77 | 77.43 | 77.98 | 78.01 | 77.99 |
| BERT-CRF-DS | TS + doc | 78.42 | 77.36 | 77.88 | 78.66 | 78.54 | 78.60 |
| WWM-CRF-DS | TS + doc | 78.51 | 77.47 | 77.99 | 78.71 | 78.60 | 78.66 |
| Longformer | TS + doc | 78.50 | 77.42 | 77.96 | 78.36 | 78.64 | 78.50 |
| WWM-SBERT | TS + doc | 80.35 | 78.66 | 79.49 | 80.55 | 79.70 | 80.12 |
| WWM-First | TS + doc | 81.33 | 79.01 | 80.15 | 81.23 | 79.88 | 80.55 |
| WWM-WWD | TS + doc | 81.38 | 79.21 | 80.28 | 82.38 | 80.14 | 81.24 |
| WWM-Noun | TS + doc | 81.50 | 79.98 | 80.73 | 82.31 | 81.48 | 81.89 |
| WWM-Noun-Yake | TS + doc | 80.46 | 79.81 | 80.13 | 81.79 | 80.72 | 81.25 |
| WWM-Noun-TextRank | TS + doc | 81.47 | 80.38 | 80.92 | 82.47 | 81.69 | 82.08 |

Table 4: The performances of different approaches on TSNER dataset.

(exact entity match) standard micro Precision (P), Recall (R), and F1 score. We will also present our analysis of the results.

5.1 Results

Table 4 shows the results of all benchmark models on TSNER. We summarize the findings into the following conclusions.

1) Introducing the document information can significantly improve the performance of topic sentence NER. For example, compared with the WWM-CRF model, three types of document-enhanced models (WWM-CRF-DS, Longformer, WWM-Noun-TextRank) can improve the F1 score by 0.67%, 0.51%, 4.09% respectively on the test set.

2) For document enhanced models, different models can incorporate different levels of document information and yield different performance. Document gist fusion models achieve better than distant supervision (DS) and Longformer. Even the worst performing document gist fusion model (WWM-SBERT) can outperform WWM-CRF-DS, demonstrating the advantage of understanding the gist of document. Surprisingly, the Longformer model shows the lowest performance. We suppose that Longformer may not be suitable for the NER task. Besides, as the data used for pretraining Longformer is different from BERT or WWM, we may not equally compare Longformer with the other models based BERT or WWM.

3) The performance of different document gist fusion models varies largely. The best model (WWM-Noun-TextRank) surpasses the worst model (WWM-SBERT) by 1.96%. This indicates a research direction on how to better extract useful information from the document. There are also some other interesting findings. First, choosing the most similar sentences may not lead to a better result. In the contrary, sentence selection based on SBERT the performs worst. Second, the ways to select keywords also have an impact on NER. The Yake based method yields a negative effect.

5.2 Error Analysis

Since the document-enhanced model outperforms the single-sentence model in topic sentence NER, in order to better analyze the reasons behind, we counted three types of errors: entity type error, cross-boundary error, and non-overlapping error. The type error means that the boundary of the predicted entity is correct but the predicted type is wrong. The cross-boundary error means that the boundary of golden one overlaps the model prediction. The non-overlapping error means no common words between gold one and model prediction. We show the error analysis of two representative models in Table 6. From the table, we summarize the following two observations.

1) Non-overlapping error type takes up most of the errors, so more attention needs to be paid to it, followed by the entity type error. We find in many

| Topic sentence and document | WWM-CRF | WWM-Noun-TextRank |
|---|---------|-------------------|
| TS: 2019[褚橙] _{Food} 来了 Here comes [Chu orange] _{Food} , 2019 Doc: ...橙子便是来自云南哀牢山的[褚橙]... ...Oranges are [Chu orange] from Ailao Mountain | Name | Food |
| TS: 11月15日, 三分钟[兴化] _{Address} 新鲜事来了! November 15, three minutes of [Xinghua] _{Address} news Doc: ...[兴化]市2019年公开招聘... ...[Xinghua] open recruitment in 2019... | None | Address |

Table 5: Case study. In the topic sentence, the text in brackets is the candidate mention, followed by the golden label. The text in brackets in the document is the sharing common entity between topic sentence and document. Predicted labels in red denote the wrong answer.

| | Type | Cboundary | NOOVER |
|-------------------|------|-----------|--------|
| WWM-CRF | 223 | 224 | 274 |
| WWM-Noun-TextRank | 181 | 221 | 243 |

Table 6: The statistics of different errors that occur in the output of WWM-Noun-TextRank models on the test set. Cboundary means that Cross-Boundary error and NOOVER is non-overlapping error.

cases that delimiters like punctuation marks in the topic sentence can help to recognize the boundary of an entity, but assigning the entity with a correct type is more difficult as the context is limited.

2) Leveraging document information can effectively reduce non-overlapping errors and entity type errors. However, it is unexpected that the document information has little effect on reducing cross-boundary errors.

5.3 Case Study

To clearly show the effectiveness of document-enhanced models for the topic sentence NER task, we analyze two representative cases by comparing the output of WWM-CRF and WWM-Noun-TextRank. The cases and prediction results are shown in Table 5. One type of common error is wrong entity type. The WWM-CRF model tends to predict entity type based on the mentioned words alone. In the first case, WWM-CRF model predicts ‘褚橙(Chu orange)’ as a person name as ‘褚(Chu)’ is a last name in Chinese names. The document-enhanced model can avoid the mistake: the WWM-Noun-TextRank model can refer to the document context to predict it as a food. Another type of common error is missing entities. In the second example, ‘兴化(Xinghua)’ is not recognized by the WWM-CRF model. In contrast, the doc-

ument enhanced model can correctly predict ‘兴化(Xinghua)’ as an address. We suppose that the word ‘市(city)’ in the document acts as a clear clue to guide the model’s prediction.

6 Conclusion and Future Work

In this paper, we propose a new task called topic sentence NER. The task is driven by real-world scenarios where extracting entities in topic sentences instead of the full-text documents is sufficient and economic. While the task is of value and is more challenging than regular NER, it has not been explored in previous research. To address this task, we build a large-scale manually annotated NER dataset, named TSNER. A family of baseline models are also established based on TSNER. We hope our dataset and benchmarks will advancing the research on topic sentence NER.

In the future, the following interesting directions can be explored.

1) When using distant supervision methods, how to leverage the noise in the document and how to model the relation between topic sentence and document are worth exploring.

2) It is promising to build a pre-trained model to learn the relationship between topic sentences and corresponding documents.

3) Strategies to extract explicit information in the document have been proved helpful for topic sentence NER and hence worth being further explored.

We also suggest incorporating more external information into NER other than document information, e.g., knowledge base and visual contents.

579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634

References

Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. [Contextual string embeddings for sequence labeling](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1638–1649. Association for Computational Linguistics.

Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. [Longformer: The long-document transformer](#). *CoRR*, abs/2004.05150.

Razvan Bunescu and Raymond Mooney. 2005. [A shortest path dependency kernel for relation extraction](#). In *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*.

Ricardo Campos, Vítor Mangaravite, Arian Pasquali, Alípio Jorge, Célia Nunes, and Adam Jatowt. 2020. [Yake! keyword extraction from single documents using multiple local features](#). *Inf. Sci.*, 509:257–289.

Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. [Reading wikipedia to answer open-domain questions](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 1870–1879. Association for Computational Linguistics.

Yubo Chen, Liheng Xu, Kang Liu, Daojian Zeng, and Jun Zhao. 2015. Event extraction via dynamic multi-pooling convolutional neural networks. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 167–176.

Christopher Clark and Matt Gardner. 2018. [Simple and effective multi-paragraph reading comprehension](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, pages 845–855. Association for Computational Linguistics.

Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. 2011. Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12(Aug):2493–2537.

Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Shijin Wang, and Guoping Hu. 2020. [Revisiting pre-trained models for chinese natural language processing](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*, pages 657–668. 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](#). *CoRR*, abs/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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.

Shizhe Diao, Jiaxin Bai, Yan Song, Tong Zhang, and Yonggang Wang. 2020. [ZEN: pre-training chinese text encoder enhanced by n-gram representations](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*, pages 4729–4740. Association for Computational Linguistics.

Tao Gui, Jiacheng Ye, Qi Zhang, Yaqian Zhou, Yeyun Gong, and Xuanjing Huang. 2020. [Leveraging document-level label consistency for named entity recognition](#). In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, pages 3976–3982. ijcai.org.

Ankit Gupta and Jonathan Berant. 2020. Gmat: Global memory augmentation for transformers. *arXiv preprint arXiv:2006.03274*.

Hangfeng He and Xu Sun. 2017. [F-score driven max margin neural network for named entity recognition in chinese social media](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 713–718. Association for Computational Linguistics.

Mandar Joshi, Omer Levy, Luke Zettlemoyer, and Daniel S. Weld. 2019. [BERT for coreference resolution: Baselines and analysis](#). 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 2019, Hong Kong, China, November 3-7, 2019*, pages 5802–5807. Association for Computational Linguistics.

Matt J. Kusner, Yu Sun, Nicholas I. Kolkin, and Kilian Q. Weinberger. 2015. [From word embeddings to document distances](#). In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 957–966. JMLR.org.

John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Eighteenth International Conference on Machine Learning*, pages 282–289.

Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016.

| | |
|-----|--|
| 804 | Manzil Zaheer, Guru Guruganesh, Kumar Avinava |
| 805 | Dubey, Joshua Ainslie, Chris Alberti, Santiago On- |
| 806 | tanon, Philip Pham, Anirudh Ravula, Qifan Wang, |
| 807 | Li Yang, et al. 2020. Big bird: Transformers for |
| 808 | longer sequences. In <i>NeurIPS</i> . |
| 809 | Yue Zhang and Jie Yang. 2018. Chinese ner using lattice |
| 810 | lstm . In <i>Proceedings of the 56th Annual Meeting of</i> |
| 811 | <i>the Association for Computational Linguistics (Vol-</i> |
| 812 | <i>ume 1: Long Papers)</i> , pages 1554–1564. Association |
| 813 | for Computational Linguistics. |

| | |
|--|-----|
| A Categories in TSNER | 814 |
| The entity types we used are shown in Table 7. | 815 |
| B Implementation Details | 816 |
| BiLSTM-CRF: The character embedding is pre- | 817 |
| trained on Chinese Giga-Word using word2vec | 818 |
| (Mikolov et al., 2013). The character embedding | 819 |
| dimension is set to 100, the LSTM hidden states | 820 |
| dimension is set to 300 and the initial learning rate | 821 |
| is set to 0.001. The models is trained using 100 | 822 |
| epochs with a batch size of 16. | 823 |
| SoftLexicon: We use the same code ³ from the | 824 |
| paper (Ma et al., 2020). The LSTM-based sequence | 825 |
| modeling layer is used. | 826 |
| Pretrained Language Model: The pre-trained | 827 |
| language model is from huggingface ⁴ . The initial | 828 |
| learning rate of PLM is set to 1×10^{-5} . We fine- | 829 |
| tune models using 20 epochs with a batch size of | 830 |
| 16. | 831 |
| WWM-Noun-TextRank: The word embed- | 832 |
| ding is pre-trained on Chinese Giga-Word using | 833 |
| word2vec (Mikolov et al., 2013). The word embed- | 834 |
| ding dimension is set to 50. The embedding of q , | 835 |
| k , v is 150. | 836 |
| Computing Infrastructure: All experiments | 837 |
| are conducted on an NVIDIA Tesla V100 (32 GB | 838 |
| of memory). | 839 |

³<https://github.com/v-mipeng/LexiconAugmentedNER>

⁴<https://huggingface.co/models>

| Categories | Interpretation | Example |
|--------------------------------------|--|---|
| 地址 Address (address) | 常见的行政区划, 如省, 市, 县, 村, 常见国家名 Common administrative divisions, such as counties, provinces, cities, villages | 北京, 中关村, 中国 Beijing, Zhongguancun, China |
| 景点 Attraction (scene) | 除地址外较小的较具体的地名, 如旅游景点等 In addition to the address, smaller and more specific place names, such as tourist attractions, etc | 长沙公园, 海洋馆, 植物园 Changsha Park, aquarium, botanical garden |
| 娱乐人物 Entertainer (ename) | 与娱乐相关的人物, 包括影视演员, 歌手等 Entertainment related characters, including film and television actors, singers, etc | 胡歌, 彭昱畅, 张学友 Hu Ge, Peng Yuchang, Zhang Xueyou |
| 体育人物 Sports figures (aname) | 主要是运动员等 Mainly athletes, etc | 刘翔, 郭晶晶 Liu Xiang, Guo Jingjing |
| 文创人物 Virtual character (gname) | 游戏, 影视剧, 小说等中的虚拟角色 Virtual characters in games, film and television dramas, novels, etc | 寒冰射手, 李元芳 Ice shooter, Li Yuanfang |
| 其他人物 Other person name (name) | 除娱乐, 体育, 文创的其他人物 Other person name besides Entertainer, Sports figures and Virtual character | 马化腾, 马云 Ma Huateng, Ma Yun |
| 公司 Company (company) | 以盈利为目的的公司 Profit oriented company | 阿里, 腾讯 Ali, Tencent |
| 组织机构 Organizations (organization) | 除公司外的团体, 如兴趣爱好团体, 大学 groups other than companies, such as interest groups, universities | 海淀棋社, 北京大学 Haidian chess club, Peking University |
| 电影 Movies (movie) | 在电影院上线的视频 Videos launched in cinemas | 英雄本色, 纵横四海 A Better Tomorrow, Once A Thief |
| 电视节目 TV programs (tvshow) | 在电视或网络上上线的电视剧, 综艺等 TV dramas and variety shows launched on TV or on the Internet | 琅琊榜, 甄传 Langya list, biography of Zhen Huan |
| 表演 Performance (show) | 需现场观看的节目, 如话剧, 戏曲, 相声, 小品等 programs to be watched on site, such as drama, opera, crosstalk, sketch, etc. | 天仙配, 女驸马 Tianxianpei, daughter-in-law |
| 事件 Events (event) | 大型赛事, 展览, 会议等 major events, exhibitions, conferences, etc. | 东京奥运会 Tokyo Olympic Games |
| 歌曲 Song (song) | 普通歌曲 ordinary song | 我愿意, 吻别 Still Here, Take me to your heart |
| 书名 literature (book) | 小说, 杂志, 文学作品等 novels, magazines, literary works, etc. | 挪威的森林, 飞鸟集 Norwegian Wood, Stray Birds |
| 美食 Food (food) | 各种食物 all kinds of food | 炸鸡腿, 汉堡 fried chicken leg, hamburger |
| 游戏 Games (game) | 各种游戏 all kinds of games | 魔兽, 王者荣耀 Warcraft, Honor of Kings |

Table 7: Categories in TSNER.