



# YNNER: Yi Language Named Entity Recognition Dataset

Chengxian Wang

Minzu University of China, Beijing, China

wcx15801417563@163.com

**Abstract.** Named entity recognition is an important task in the field of natural language processing, used to identify entities in text and classify them into pre-defined types. The research on Yi language named entity recognition is still in its early stages both domestically and internationally. Currently, there is no publicly available comprehensive dataset for Yi language named entity recognition, which has hindered the progress in this field. This paper constructed a named entity recognition dataset (Yi language news named entity recognition, YNNER), and manually annotates the names of person, places, and institutions. Then, the named entity recognition model is used to carry out experimental and comparative analysis on the dataset. The experimental results show that the  $F_1$  values of all models are above 70%, which proved the validity and availability of the dataset constructed. This paper aims to promote the research and development of Yi language named entity recognition, provide dataset and baseline models for this field, and expand related research.

**Keywords:** Natural language processing; Yi language; Deep learning; Named entity recognition.

## 1 Introduction

Named Entity Recognition (NER) is an important and basic work in natural language processing, which aims to identify entity objects such as person, place and organization, etc. Recognition accuracy directly affects the performance of upstream tasks such as knowledge graph [1], machine translation [2] and question answering system [3]. Its development has gone through the early methods based on dictionaries and rules, to the Hidden Markov model(HMM), the Maximum Entropy Markov model(MEMM) and the Conditional Random Field model(CRF) [4-5], and then to the current models such as Recurrent Neural Network(RNN) and Convolutional Neural Network model(CNN) [4]. Although many named entity recognition methods have been proposed, there are still many problems to be solved due to the randomness, complexity and variability of named entities. This is mainly due to the lack of high-quality datasets. The lack of datasets will directly limit the training and deployment effect of the model. Data annotation is still a time-consuming and expensive task, especially in some specific domains, which requires domain experts to perform data annotation, which is a great challenge.

Therefore, how to construct high-quality datasets quickly, accurately and economically is an important problem to be solved in current named entity recognition technology.

This paper constructed a Yi language named entity recognition dataset YNNER, which is derived from news corpus. This paper introduced the preparation of the dataset, the labeling system, and the construction method in detail. Three models, namely the traditional CRF model, the Conditional Random Field model of Bidirectional Long Short-Term Memory Network(BiLSTM-CRF) based on RNN variant and the Conditional Random Field model of Iterative Dilated Convolutional Neural Network(IDCNN-CRF) based on CNN variant, were used to experiment and analyze the dataset.

## 2 Related Work

The research on named entity recognition mainly focuses on how to construct high-quality datasets. Many domestic researchers have carried out in-depth research. Among them, the Chinese named entity recognition dataset is relatively mature. MSRA dataset [5] has been constructed by Microsoft Research Asia in China. The data adopts a variety of annotation methods and the annotation quality of the dataset is high, which is usually used as the standard for Chinese named entity recognition research and evaluation. Weibo dataset [6] is constructed by Sina Corporation of China, which contains large-scale Chinese social media data. With a lot of noise and linguistic variability, the annotation quality is poor. Peking University and Microsoft Research Asia collaborated to create People's Daily dataset [7], which has been widely used in various named entity recognition model research. In addition to general datasets, many specialized datasets have been developed. Du et al. [8] constructed a named entity recognition dataset in the field of military command and control to provide a basis for the construction of knowledge graph. Shah et al. [9] constructed a named entity recognition dataset in the financial field, and designed a weakly supervised entity recognition strategy based on this dataset, which achieved good recognition results. Zhang et al. [10] constructed a named entity recognition dataset of construction documents, introduced the design process, labeling specification, data format, etc., and tested it through the CRF model, and the F1 value reached 87.9%.

Compared with Chinese named entity recognition, the difficulty of constructing minority language named entity recognition datasets is that native speakers or linguists are required to participate in the annotation process, and entity objects cannot be randomly divided. It mainly focuses on Tibetan, Uyghur and Mongolian. Researchers from Xizang University [11] built a Tibetan general named entity recognition dataset, and many scholars have studied entity recognition technology on this dataset and achieved good results. Scholars from Mongolia Normal University [12] constructed named entity recognition datasets in the field of education and literature, and evaluated a variety of deep learning methods to achieve a practical level. Researchers from Xinjiang University [13] constructed a Uyghur named entity recognition dataset to provide a basis for machine translation tasks. Yi language is relatively backward compared with Tibetan,

Uyghur, Mongolian and other languages, and there is no public dataset related to named entity recognition.

### 3 Dataset Construction

#### 3.1 Data Collection

In this paper, news articles are collected from <https://lsrbywb.ls666.com> Liangshan Daily news media platform and pre-processed. Preprocessing includes deleting redundant data, removing incorrect data (such as redundant Spaces, Chinese characters, and incorrect symbols), and converting punctuation marks.

#### 3.2 Annotation Process, Data Format and Annotation Specification

In order to ensure the quality and accuracy of the data, a computer professional, a Yi language linguistic researcher and a linguistics professional were used as annotators to annotate the data according to the Chinese MSRA [5] annotation specification. The dataset was divided into three parts, and the annotation results of the other two persons were checked each other after annotation. The inconsistency is recorded, and finally it is determined by computer professionals to retain or delete. At the same time, repeated self-examination and audit were also carried out.

**Data format** Since each entity sentence may consist of two or more characters, when generating experimental data: In this paper, the BIO [14] (Begin-In-Out) marking pattern is used to determine whether the character is part of a certain type of entity by marking each character. Each type of entity is divided into starting position (B-) and non-starting position (I-), and the non-entity characters are uniformly marked as O. The final complete tag set  $\text{TagSet} = \{O, B\text{-PER}, I\text{-PER}, B\text{-LOC}, I\text{-LOC}, B\text{-ORG}, I\text{-ORG}\}$  contains seven kinds of tags. These labels are used to determine the entity category to which each word belongs for named entity recognition. The defined set of annotations is shown in Table 1.

**Table 1.** Annotation set of Yi language named entities

Label	Meaning
B-PER	Name start character
I-PER	Name is not a starting character
B-LOC	Place starting character
I-LOC	Place is not a starting character
B-ORG	Organization starting character
I-ORG	Organization is not a starting character
O	non-entity

**Annotation specification** (1) Entity types: In the process of manually labeling the corpus of Yi language, three types of entity names such as person names, place names and organization names are used for all the corpus in this paper, and non-named entities do not need to be labeled. (2) Labeling unit: refer to the MSRA named recognition



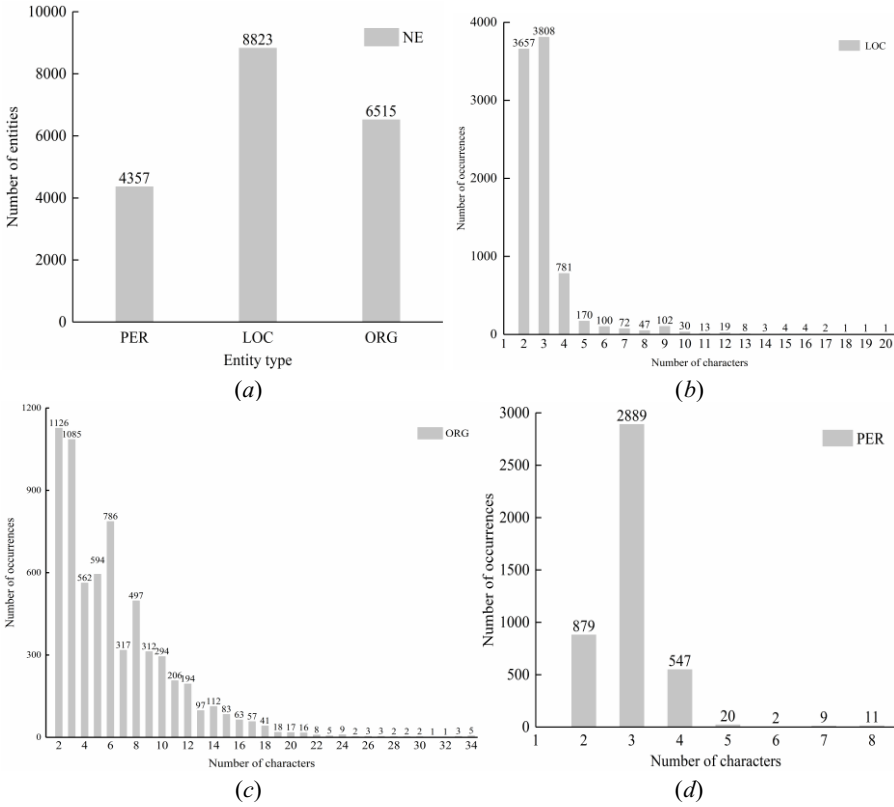


Fig. 1. Distribution of the number of characters and occurrences of entities

### 4 NER Model

Yi language is a monosyllabic language. Compared with other languages, the number of place names and organization names in Yi language is huge, while the transliteration entity is more, and the length is not limited. Yi language is a morphologically underdeveloped type language, similar to Chinese. In this paper, character-level named entity recognition model is selected according to the entity characteristics of Yi language.

In order to further explore and analyze the performance of Yi language named entity recognition on the dataset constructed in this paper, this paper refers to the methods of Tibetan, Uygur and Mongolian recognition research, and finally selects three groups of representative named entity recognition models, namely CRF, BiLSTM-CRF, IDCNN-CRF model. The model architecture diagram is shown in Fig 2.

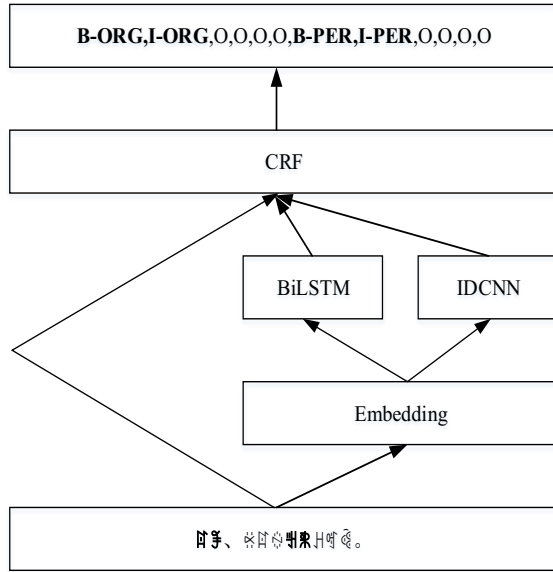


Fig. 2. Model architecture

### 4.1 CRF Model

Conditional random fields are a statistics-based sequence labeling model proposed by Lafferty et al. [15], CRF model is briefly introduced in this paper, and details are given in reference [15].

We regard the Yi language named entity recognition problem as a sequence labeling problem, and generate a first-order linear chain CRF based on an undirected graph  $G=(V, E)$ .  $V$  is the set of random variable  $Y=\{Y_i|1\leq i\leq n\}$ ,  $n$  need to enter a sentence marking unit,  $E=\{(Y_{i-1}, Y_i) | 1\leq i\leq n\}$  or less or less when linear chain consisting of  $n-1$  side. For each sentence  $x$ , define two non-negative factors:

For each edge:

$$\exp\left(\sum_{k=1}^K \lambda_k f_k(y_{i-1}, y_i, x)\right)$$

For each node:

$$\exp\left(\sum_{k=1}^{K'} \lambda'_k f'_k(y_i, x)\right)$$

Where  $f_k$  is a binary feature function and  $K$  and  $K'$  are the number of features defined at each edge and corresponding node. Given a sequence  $x$  that needs to be labeled, the conditional probability of its corresponding labeled sequence  $y$  is given by Eq. 1.

$$P(y|x) = \frac{1}{Z(x)} \exp \left( \sum_{i,k} \lambda_k f_k(y_{i-1}, y_i, x) + \sum_{i,k} \lambda'_k f'_k(y_i, x) \right) \tag{1}$$

$Z(x)$  is the normalization function. Given the training set  $D$ , the parameters of the trained model are used to maximize the conditional likelihood value. When given to tag sequence  $x$ , its corresponding sequence by the parameter  $\text{argmax}_y P(y|x)$ .

### 4.2 BiLSTM-CRF Model

BiLSTM -CRF is a combination of BiLSTM and CRF. The CRF model is introduced above, and the BiLSTM model is mainly introduced below. First, the LSTM model is introduced.

LSTM model [16] is a variant of RNN. The proposed model solves the problem of gradient disappearance and gradient explosion of RNN and the problem of insufficient long-distance information learning ability. Based on the original RNN, LSTM adds input gate, output gate, forget gate and cell state. The structure of the model is shown in Fig. 3.

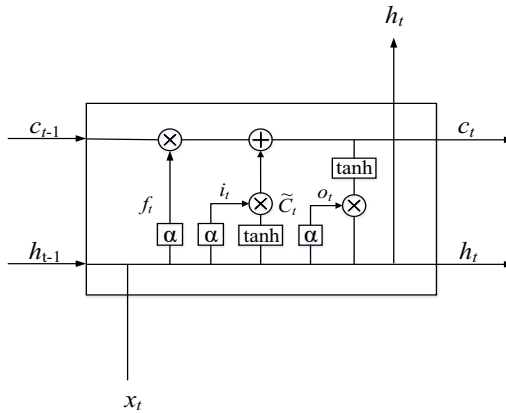


Fig. 3. LSTM structure

The main function of the forgetting gate  $f_t$  is to screen the cell state of the previous layer and determine whether to discard information or retain information, and its calculation is shown in Eq. 2.

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \tag{2}$$

Where,  $h_{t-1}$  represents the hidden layer state at the previous moment,  $x_t$  represents the input at the current moment,  $\sigma$  represents the sigmoid activation function, and  $W_f$  and  $b_f$  represent the weight and bias of the forgetting gate, respectively.

The input gate  $i_t$  pair controls what part of the new information is saved into the cell state. Combined with the temporary state of the cell, the current cell state  $\tilde{c}_t$  is calculated. As shown in Eq 3-Eq 5.

$$c_t = f_t \otimes c_{t-1} + i_t \otimes \tilde{c}_t \quad (3)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \quad (4)$$

$$\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c) \quad (5)$$

Where,  $c_{t-1}$  and  $c_t$  represent the previous layer and the current cell state, respectively.  $\tilde{c}_t$  is a temporary cell state;  $f_t \otimes c_{t-1}$  and  $i_t \otimes \tilde{c}_t$  denote deleted information and added information respectively;  $i_t$  is the value of the input gate;  $W_i$  and  $W_c$  are the weights of the input gate and cell state, respectively,  $b_i$  and  $b_c$  are the biases of the input gate and cell state;

The output gate determines the current hidden layer state by combining the hidden state at the previous moment, the current input value and the current cell state, as shown in Eq 6-Eq 7.

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \quad (6)$$

$$h_t = o_t * \tanh(c_t) \quad (7)$$

The disadvantage of LSTM is that the information can only be transmitted in the positive direction, and it cannot well consider all the information of the context, which affects the accuracy of entity recognition.

BiLSTM is an improvement for the LSTM model by concatenating the forward LSTM with the reverse LSTM. In the forward learning, it learns the historical information from left to right, and in the reverse learning, it learns the future information from from to do. Better learn above and below.

The BiLSTM model usually uses Softmax as the output, but lacks consideration of the dependency between labels, and the recognition results may have problems such as BB and OI. CRF can consider the context and the relationship between adjacent labels, which makes the prediction of label sequences more accurate and coherent. Finally, the label sequence output by the CRF layer is used as the final output of the model.

### 4.3 IDCNN-CRF Model

IDCNN-CRF combines IDCNN and CRF model, and its main purpose is to increase the receptive field of the model without increasing the model parameters and maintaining the speed of the model. In the following, IDCNN is mainly introduced.

IDCNN was proposed by Yu et al. [17] in 2015, with the main purpose of increasing the perception field. In the classical convolutional neural network, the convolution kernel slides over continuous regions, while the dilated convolution adds an expansion



width above the classical convolution, and the data in the middle of the expansion width will be skipped during the convolution operation, and the size of the convolution kernel remains the same, so that a convolution kernel with the same size can obtain a wider input matrix data and increase the perception field of the convolution kernel. See Fig 4 for a schematic of the dilated convolution. The three graphs represent the convolution operation of three layers, respectively. The Fig 4(a) is a normal convolution operation with a convolution kernel size of  $3 \times 3$ . In Fig 4(b), the dilation width of the convolution is 2, and above the convolution in (a), the receptive field of view is increased to  $7 \times 7$ ; The dilation width is 4 in Fig 4(c), above the convolution operation in (b), at which point the convolution receptive field of view is equivalent to expanding to  $15 \times 15$ .

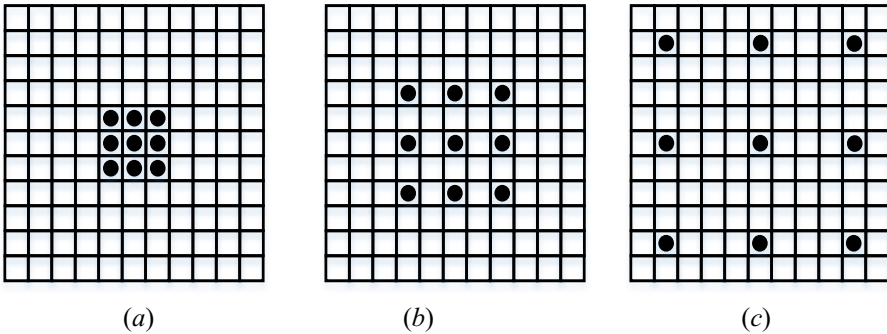


Fig. 4. Dilated convolution

Dilated convolution was originally applied in image processing. Strubell et al. [18] introduced dilated convolution into the field of natural language processing and proposed IDCNN model with remarkable effect. In IDCNN, the receptive field increases exponentially with the increase of the number of layers, but the parameters only increase linearly, so that the receptive field can quickly cover all the input sequences. The sentence is input into IDCNN, and the features are extracted through the convolutional layer, and then connected to the CRF layer through the mapping layer. Finally, the label sequence output by the CRF layer is used as the final output of the model, that is, each entity in the input sequence is labeled.

## 5 Experimental Results and Analysis

### 5.1 Dataset and Evaluation Metrics

At present, there is no public dataset for Yi language named entity recognition. The Yi language named entity recognition dataset (YNNER) constructed in this paper is used for experiments. The resume dataset of this paper contains 11247 news texts and 19695 entities. It contains 4,357 names of person, 8,823 names of places and 6,515 names of organizations. The dataset is divided into training, validation, and test sets according to 8:1:1.

In this experiment, Precision ( $P$ ), recall ( $R$ ) and F1-measure ( $F_1$ ) values are used as evaluation indicators to Measure the named entity recognition performance of the

model [17]. Among them,  $F_1$  value is used as the main reference index and is calculated as shown in Eq.8- Eq.10.

$$P = \frac{TP}{TP + FP} \quad (8)$$

$$R = \frac{TP}{TP + FN} \quad (9)$$

$$F_1 = \frac{2 \times P \times R}{P + R} \quad (10)$$

## 5.2 Parameter Setting

The experimental part of this paper aims to evaluate the performance of three different models on the named entity recognition task. It contains three mainstream models: CRF, BiLSTM-CRF, and IDCNN-CRF.

To ensure the reliability of the results, the maximum sequence length is set to 150, the training Epoch is set to 50, the Hidden\_dim is set to 200, the Batch\_size is set to 32, the Dropout is set to 0.25, the learning rate is set to 0.001, and the optimizer uses Adam. Nums is set to 2 and Filter\_nums is set to 64 in IDCNN model.

The hardware environment of this experiment is set as Windows 10(64bit) operating system, 32GB memory, Inter(R) Core(TM) 5-10300H CPU @2.5Hz processor. The GPU is NVIDIA GeForce GTX 3060 Ti 12GB, and python3.8 and Pytorch1.8.0+cu111 are used for programming and training.

## 5.3 Experiment and Result Analysis

The experimental part of this paper aims to evaluate the performance of three different models on the named entity recognition task. It includes CRF model, BiLSTM-CRF model, IDCNN-CRF model. The experimental results are shown in Table 3 and Table 4.

**Table 3.** Experimental results of the model on different data sets

Model	$P$	$R$	$F_1$
CRF	72.11	68.44	70.22
BiLSTM-CRF	79.21	75.45	77.28
IDCNN-CRF	80.42	77.36	79.33

As can be seen from the results in Table 3, the performance of BiLSTM-CRF and IDCNN-CRF models on YNNER data set is better than that of the traditional CRF model, and the  $F_1$  value can reach more than 75 respectively. Both models have good global optimization ability and perform well.

**Table 4.** Comparison of  $F_1$  values for different categories of entity experiments of each model

Model	PER	LOC	ORG
CRF	74.20	71.3	66.15
BiLSTM-CRF	81.34	75.22	68.50
IDCNN-CRF	82.23	76.54	69.35

According to Table 4, among the three types of entity labels, name recognition is the best, which is because the name context information is clear. The recognition accuracy of place names is slightly lower than that of personal names, and it is found that some of them are misrecognized as personal names, which is caused by polysemy. Institution names have the lowest recognition accuracy, which is due to the fact that institution names usually consist of multiple entities, such as containing places, and the boundaries are difficult to determine.

Error analysis is a useful tool to better understand the strengths and weaknesses of a model or dataset. In this paper, the best performing model BiLSTM-CRF extracts 100 wrong instances in the test set and manually checks them. The entity type accounted for the largest proportion of incorrect identification (60%). This is followed by entity boundary misidentification (40%). There are other reasons for errors, such as lack of training samples, etc. Some typical examples of the majority misclassifications are also listed in order to better understand these errors.

- 1. Entity type error:** This type of error occurs because the model encounters an out-of-vocabulary word. Without training, unknown words will be treated as non-entities or their relationships will be incorrectly predicted. With 成都的春天是木棉花的季节，春天是木棉花的季节，春天是木棉花的季节。(already work in the mail deliverer bent wood mingle six years, says this time before the Spring Festival, every year is the busiest time of the year.) as an example, the model predicts the 木棉花 was forecast to the entity. The rare place names or unknown words in the training set are incorrectly recognized because the model has not been sufficiently trained.
- 2. Entity boundary identification error:** When a place name or organization name composed of multiple words appears in an entity, the combination of multiple words will make it difficult to judge the boundary. With 四川省凉山州西昌市西大街107号，凉山州50.2亩，1284名。四川省凉山州西昌市西大街107号，凉山州50.2亩，1284名。(Sichuan Yi Language School and Liangshan Civil Cadre School implement the system of one school, co-construction and co-management of province and prefecture. The two schools currently have 107 employees, a campus area of 50.2 mu, and 1284 secondary vocational students. Rabbi Ashi VISITED the campus and extended Lantern Festival greetings and best wishes to the faculty and staff.) as an example, the model of the 四川省 marked as place name, but in 四川省凉山州西昌市 such by multiple places and institutions of mixed character combinations institutions unable to properly identify the border in the name.

Considering the characteristics of the annotated data set and the preliminary experimental results, it can be seen that the characteristics of Yi language named entity recognition, such as imbalance, entity nesting and the influence of unknown place names, need to adopt a variety of strategies and methods to improve the accuracy of the algorithm. This is worthy of further research in the future. In contrast, the reasons for the better performance of Chinese and English on named entity recognition tasks mainly include the wealth of data resources, clear rules of language structure and the maturity of natural language processing technology. However, these advantages may be weakened when facing Yi language, so different strategies and methods are needed to improve the accuracy of named entity recognition.

## 6 Conclusion

Aiming at the problem of lacking high-quality annotated corpus in the field of Yi language named entity recognition, this paper constructs an entity naming dataset for Yi language, verifies the validity and availability of the dataset by three mainstream named entity recognition methods, and analyzes the causes of recognition errors in depth. The dataset can provide data support for Yi language named entity recognition work. The dataset is collected from the real data of online news media, and contains three types of entity characteristics in Yi language text. This dataset provides important data support for the research of Yi language named entity recognition algorithms, which can be used for model training, testing and evaluation, and provides a good data support for further research in this field. Through experiments, the validity and practicability of the dataset are proved, and it is expected to provide important reference value for the development of Yi language natural language processing.

There are several shortcomings in this paper. Firstly, due to the shortcomings of using BIO annotation method, such as fuzzy entity boundary, inability to express the end of the entity, coupling with the entity type, and the characteristics of imbalance in Yi language named entity recognition, entity nesting, more entity phrases, and the influence of unknown place names, the next stage of research will use more accurate annotation to improve the annotation quality of Yi language named entity recognition dataset. Secondly, because the language structure and lexical features of Yi language are different from other languages, the next stage of algorithm design will focus on the characteristics of Yi language.

## References

1. Ke J, Wang W, Chen X, et al. Medical entity recognition and knowledge map relationship analysis of Chinese EMRs based on improved BiLSTM-CRF. *Computers and Electrical Engineering*, 2023, No.108: 108709.
2. He C Y, Zhang J J. Neural Machine Translation Model Incorporating Bilingual Named Entity. *Journal of Chinese Information Processing*. 2023, No.37(12): 44-53.
3. Bao J Y, Yu J H, Xu N, et al. Research on the improved method of named entity recognition in Q & A system. *Journal of Data Acquisition & Processing*, 2020, (5): 930-941.

4. Zhao J G, Qian Y R, Wang K, et al. Survey of Chinese named entity recognition research. *Computer Engineering and Applications*, 2024, No.60(1): 15-24.
5. Sun Z, Li X. Named entity recognition model based on feature fusion. *Information*, 2023, 14(2): 133-146.
6. Yan Y, Zhu P, Cheng D, et al. Adversarial multi-task learning for efficient Chinese named entity recognition. *ACM Transactions on Asian and Low-Resource Language Information Processing*, 2023, 22(7): 1-19.
7. Yin Z G, Lu J F. Research on Chemical Named Entity Recognition Combining Radical Features. *COMPUTER & DIGITAL ENGINEERING*, 2023, No.51(4): 809-816.
8. Du X M, Yuan Q B, Yang F, et al. Construction of Named Entity Recognition Corpus in Field of Military Command and Control Support. *Computer Science*, 2022, No.49(201): 133-139.
9. Shah A, Vithani R, Gullapalli A, et al. Finer: Financial named entity recognition dataset and weak-supervision model. *arXiv preprint arXiv: 2302. 11157*, 2023.
10. Zhang Q, Xue C, Su X, et al. Named entity recognition for Chinese construction documents based on conditional random field. *Frontiers of Engineering Management*, 2023, No.10(2): 237-249.
11. Ge L N M, Qun N, Xiang X C R, et al. Tibetan named entity recognition method combined with word segmentation features. *Plateau Scientific Research*, 2023, No.7(4): 106-114.
12. Wang Y R, Lin M, Li Y L. BERT Mongolian Word Embedding Learning. *Computer Engineering and Applications*, 2023, No.59(2): 129-134.
13. Mai H M T·M M T. Study on Uyghur Named Entity Recognition and Related Problems. Xinjiang University, 2018.
14. Mnih V, Heess N, Graves A. Recurrent models of visual attention. *Advances in neural information processing systems*, 2014, 27-36.
15. Lafferty J, McCallum A, Pereira F. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. *Icml. 2001*, No.1(2): 3.
16. Chen X C, Qiu X P, Zhu C X. Long Short-Term Memory Neural Networks for Chinese Word Segmentation. *Proceeding soft the 2015 Conference on Empirical Methods in Natural Language Processing*. 2015. 1197-1206.
17. Yu F, Koltun V. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.
18. Strubell E, Verga P, Belanger D, et al. Fast and accurate entity recognition with iterated dilated convolutions. *arXiv preprint arXiv:1702.02098*, 2017.

**Open Access** This chapter is licensed under the terms of the Creative Commons Attribution-NonCommercial 4.0 International License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits any noncommercial use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

