

Article

A New Chinese Named Entity Recognition Method for Pig Disease Domain Based on Lexicon-Enhanced BERT and Contrastive Learning

Cheng Peng^{1,2,3} , Xiajun Wang^{1,4}, Qifeng Li^{1,2,3,*}, Qinyang Yu^{1,2,3}, Ruixiang Jiang^{1,2,3}, Weihong Ma^{1,2,3} , Wenbiao Wu^{1,2,3}, Rui Meng^{1,2,3}, Haiyan Li^{1,2,3}, Heju Huai^{1,2,3}, Shuyan Wang^{1,2,3} and Longjuan He⁵

¹ Information Technology Research Center, Beijing Academy of Agriculture and Forestry Sciences, Beijing 100097, China; pengc@nercita.org.cn (C.P.); 202221108012240@stu.hubu.edu.cn (X.W.); yuqy@nercita.org.cn (Q.Y.); jiangrx@nercita.org.cn (R.J.); mawh@nercita.org.cn (W.M.); wuwb@nercita.org.cn (W.W.); mengr@nercita.org.cn (R.M.); botanylihaiyan@163.com (H.L.); huaihj@nercita.org.cn (H.H.); wangsy@nercita.org.cn (S.W.)

² National Innovation Center of Digital Technology in Animal Husbandry, Beijing 100097, China

³ National Engineering Research Center for Information Technology in Agriculture, Beijing 100097, China

⁴ Faculty of Resources and Environmental Science, Hubei University, Wuhan 430061, China

⁵ Institute of Agricultural Economics and Development, Chinese Academy of Agricultural Sciences, Beijing 100081, China; helongjuan@caas.cn

* Correspondence: liqf@nercita.org.cn

Featured Application: Our work provides reliable technical support for the information extraction of pig diseases in Chinese. It can be applied to other domain-specific fields, thereby facilitating seamless adaptation for named entity identification across diverse contexts.

Abstract: Named Entity Recognition (NER) is a fundamental and pivotal stage in the development of various knowledge-based support systems, including knowledge retrieval and question-answering systems. In the domain of pig diseases, Chinese NER models encounter several challenges, such as the scarcity of annotated data, domain-specific vocabulary, diverse entity categories, and ambiguous entity boundaries. To address these challenges, we propose PDCNER, a Pig Disease Chinese Named Entity Recognition method leveraging lexicon-enhanced BERT and contrastive learning. Firstly, we construct a domain-specific lexicon and pre-train word embeddings in the pig disease domain. Secondly, we integrate lexicon information of pig diseases into the lower layers of BERT using a Lexicon Adapter layer, which employs char–word pair sequences. Thirdly, to enhance feature representation, we propose a lexicon-enhanced contrastive loss layer on top of BERT. Finally, a Conditional Random Field (CRF) layer is employed as the model’s decoder. Experimental results show that our proposed model demonstrates superior performance over several mainstream models, achieving a precision of 87.76%, a recall of 86.97%, and an F1-score of 87.36%. The proposed model outperforms BERT-BiLSTM-CRF and LEBERT by 14.05% and 6.8%, respectively, with only 10% of the samples available, showcasing its robustness in data scarcity scenarios. Furthermore, the model exhibits generalizability across publicly available datasets. Our work provides reliable technical support for the information extraction of pig diseases in Chinese and can be easily extended to other domains, thereby facilitating seamless adaptation for named entity identification across diverse contexts.

Keywords: pig disease; Chinese named entity recognition; lexicon-enhanced BERT; contrastive learning; small sample



Citation: Peng, C.; Wang, X.; Li, Q.; Yu, Q.; Jiang, R.; Ma, W.; Wu, W.; Meng, R.; Li, H.; Huai, H.; et al. A New Chinese Named Entity Recognition Method for Pig Disease Domain Based on Lexicon-Enhanced BERT and Contrastive Learning. *Appl. Sci.* **2024**, *14*, 6944. <https://doi.org/10.3390/app14166944>

Academic Editor: Tobias Meisen

Received: 22 July 2024

Revised: 4 August 2024

Accepted: 6 August 2024

Published: 8 August 2024



Copyright: © 2024 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Pig diseases, such as African swine fever or highly pathogenic porcine reproductive and respiratory syndrome (PRRS), have become major issues threatening the healthy development of the pig industry in China. These diseases have not only caused significant

economic losses to society but have also impacted food safety. In the context of large-scale and intensive pig breeding practices, it is of great significance to establish intelligent diagnostic and preventive measures for pig diseases. Early prevention and timely diagnosis are pivotal for maintaining swine health and mitigating potential losses.

Named Entity Recognition (NER) assumes a critical role in this endeavor by identifying specific entities within textual corpora, serving as the cornerstone for numerous downstream tasks in natural language processing. These tasks include but are not limited to information retrieval, intelligent question answering, and knowledge graph construction. The extraction of entity information from the pig disease domain, including diseases, disease sites, symptoms, drugs, etc., can rapidly extract key information from vast quantities of unstructured data. This process provides reliable core data for intelligent diagnosis, dialogue consultation, drug recommendation among veterinarians and farmers, and other application scenarios related to pig diseases. Additionally, it offers veterinarians more comprehensive and accurate diagnostic information, thereby improving the efficiency and accuracy of diagnoses.

However, the existing entity recognition methods mostly focus on the recognition of a person, location, organization, etc. Given the pressing need to bolster disease surveillance and management in swine, there arises an urgent imperative to develop specialized NER methodologies tailored to the specific lexicon of pig disease terminology in Chinese.

The early NER methods include rule-based recognition methods and statistics-based machine learning recognition methods. In recent years, with the rapid development of neural networks, methods of deep learning are more suitable for the task of NER and have become the mainstream method [1–5].

The deep learning-based NER method can learn more complex features and achieve good results. In contrast to the rule-based and statistics-based approaches, deep learning-based NER methods do not necessitate an abundance of artificial features. Therefore, deep learning-based methods have attracted wide interest from researchers. Common deep learning models include convolutional neural network (CNN), recurrent neural network (RNN), graph neural network (GNN), deep neural network (DNN), generative adversarial network (GAN), long short-term memory network (LSTM), Transformer and BERT (bi-directional encode representation from transformers), and so on [1,6]. Deep learning models have become the mainstream method and achieve state-of-the-art results in NER. However, the scalability of deep learning models applied in a specific domain remains a significant challenge. The primary challenges are the presence of specialized vocabulary and domain-specific knowledge, coupled with the scarcity of labeled datasets. Additionally, Chinese named entity recognition faces difficulties such as ambiguous boundaries and entity ambiguity [1,2].

The lexicon-based NER method can effectively avoid segmentation errors and improve the accuracy of entity boundary recognition by integrating potential word information into feature vectors. Currently, a large number of lexicon-enhanced Chinese entity extraction methods have been proposed, with better performance than methods based on character embedding or word embedding [7,8]. Lattice-LSTM [9] has achieved new benchmark results on several public Chinese NER datasets. However, the Lattice-LSTM model architecture is complex, which limits its application in many industrial areas requiring real-time NER responses. A convolutional neural network-based method that incorporates lexicons using a rethinking mechanism was proposed, which can model all the characters and potential words that match the sentence in parallel [10]. A lexicon-based graph neural network with global semantics was proposed to tackle word ambiguities. In this model, the lexicon knowledge is used to connect characters to capture the local composition, while a global relay node can capture global sentence semantics and long-range dependency [11]. A Lexicon-Enhanced BERT (LEBERT) for Chinese sequence labeling was put forward [12]. The model integrates external lexicon knowledge into BERT layers directly with a Lexicon Adapter layer and achieves better performance than both lexicon-enhanced models

and the BERT baseline in Chinese datasets. More character–word association models have been proposed, such as SoftLexicon [13], FLAT [14], and PLTE [15].

The pre-trained NER method effectively leverages deep bidirectional contextual information. It demonstrates superior performance with shorter training times, reduced labeling data requirements, and improved results compared to traditional models. Currently, BERT [16] is widely used, followed by ELMo [17], RoBERTA [18], ERNIE [19], ALBERT [20], and others. At present, the pre-trained models and lexicon are integrated by utilizing their respective strengths. Li proposed Flat-Lattice Transformer for Chinese NER, which converts the lattice structure into a flat structure consisting of spans [14]. Li proposed the LEBERT-BiLSTM-CRF model for elementary mathematics text NER, which integrates external lexicon knowledge into BERT layers directly with a lexicon adapter layer and performs better than other NER models [21].

Contrastive learning acquires feature representations of samples by comparing positive and negative samples in feature space. This approach has garnered significant attention in the fields of computer vision (CV) and natural language processing (NLP). The ConSERT (Contrastive Framework for Self-supervised Sentiment Representation Transfer) and SimCSE (Simple Contrastive Learning of Sentiment Embedding) models, which use different data enhancement methods and comparative learning loss function to learn the representation of sentences, obtain SOTA results on the task of text semantic similarity [22,23]. Contrastive learning with prompt guiding for few-shot NER (COPNER) was proposed and outperforms state-of-the-art models with a significant margin in most cases. This method introduces category-specific words that COPNER composed out of prompts as supervised signals for contrastive learning to optimize entity token representation [24]. Moreover, Named Entity Recognition in low-resource scenarios based on contrastive learning has also received considerable attention [25–27]. He proposed a novel prompt-based contrastive learning method for few-shot NER without template construction and label word mappings [25]. Li proposed a multi-task learning framework CLINER for few-shot NER [26].

In the field of livestock husbandry, text mining, Named Entity Recognition (NER), intelligent question-and-answer systems, and artificial intelligence (AI) technologies have been gradually applied. However, this field faces numerous challenges, including the prevalence of technical terms, complex knowledge structures, fine knowledge granularity, and a lack of labeled datasets [28]. Seok created a BERT-DIS-NER model that adds a CRF layer to BERT for disease named entity recognition and used syllable unit-based named entity recognition that can reflect the characteristics of disease names. The F1-score is 0.81, trained with human data and fine-tuned with animal data [29]. Kung designed and implemented an intelligent knowledge question-and-answer system for pig farming based on bi-GRU and SNN methods, combined with the LTSM deep learning method [30].

NER methods have found extensive applications in the agricultural domain and other specific fields [31–36]. Nonetheless, there remains an apparent gap in current research concerning the accurate recognition of named entities within the domain of pig diseases in Chinese. Pig disease data are characterized by complex entities, fuzzy boundaries, and domain-specific vocabulary, which encompasses specialized terminologies drawn from the domains of animal husbandry and veterinary science.

Furthermore, the resources in the field of pig diseases are confined and dispersed, exacerbating the scarcity of publicly available benchmark corpora and labeled datasets specific to this domain in Chinese. While considerable research has been devoted to NER systems in human medicine [37–40], it remains impractical to directly transfer such models to the domain of pig diseases due to the domain-specific rules and vocabulary governing this domain. Hence, named entity recognition in the field of pig diseases needs to be further explored. A model of Pig Disease Chinese Named Entity Recognition (PDCNER) is proposed in this paper. The main contributions of the paper are as follows:

- (1) We propose a simple yet effective NER model that integrates enhanced lexicon and contrastive learning for the complex pig disease domain, making the model more sensitive to texts in this domain and improving predictions for entities. The lexicon-enhanced BERT facilitates the direct integration of external lexicon knowledge of pig diseases into BERT layers via a Lexicon Adapter layer.
- (2) To enrich the semantic feature representation and improve performance under data scarcity conditions, we propose a lexicon-enhanced contrastive loss layer on top of the BERT encoder. Experimental results on small sample scenarios and common public datasets demonstrate that our model outperforms other models.
- (3) Given the lack of an annotated corpus for the pig disease domain, we collected and annotated a new Chinese corpus and annotated datasets consisting of 7518 entities. To address the insensitivity of word segmentation caused by the specialization of the pig disease domain, we constructed a lexicon for identifying specific terms in pig diseases using frequency statistics methods under the guidance of veterinarians.

The remainder of the paper is organized as follows. Section 2 introduces the dataset and method proposed in this paper. Section 3 provides a detailed description of our experiments and analyzes the results. Finally, the conclusions are presented in Section 4.

2. Materials and Methods

2.1. Materials

2.1.1. Corpus Collection and Pre-Processing

Due to the lack of NER public benchmark datasets in the pig disease domain, a new Chinese pig disease corpus was constructed and annotated under the guidance of animal disease experts and veterinarians. To ensure the quality of the data, we collected information on pig diseases from professional books, published standards, Baidu Encyclopedia, and official websites. The data source details are mentioned in Appendix A.

After the data acquisition, we performed basic data pre-processing steps. Firstly, optical character recognition (OCR) technology was used to convert the books and standards into text format. Secondly, useless data, such as garbled characters or special symbols, were manually deleted and wrong words were modified in raw text. Thirdly, the duplicate data and invalid data were removed. Ultimately, a comprehensive and effective text corpus containing 1.45 million characters was obtained (Corpus I of pig disease).

2.1.2. Analysis of Corpus

First, we analyze the characteristics of the corpus. The corpus of pig diseases contains many specialized terms, such as disease names, medicine names, and body parts. Examples include “腹膜炎 (peritonitis)”, “盐酸头孢塞呋 (ceftazidime hydrochloride)”, and “膀胱黏膜 (bladder mucosa)”. Additionally, boundary ambiguity and nested entities are common problems in named entity recognition (NER). For instance, the phrase “非洲猪瘟 (African swine fever)” encompasses two distinct entities: the location entity “非洲 (Africa)” and the disease entity “猪瘟 (swine fever)”. This ambiguity results in inaccuracies during the word segmentation process, thereby increasing the challenges associated with recognizing these entities.

2.1.3. Corpus Annotation

According to the characteristics of the corpus and referring to existing labeling standards in fields such as human medicine, we developed labeling criteria in collaboration with veterinary experts. These criteria include non-overlapping labeling, non-nested labeling, and covering as many types as possible. Finally, a multi-round iteration approach was used to label and refine the data, ensuring the accuracy and consistency of the labeling.

On the basis of the above work, 152,596 characters were selected from Corpus I to form Corpus II for entity labeling. The Label Studio tool and BIEO labeling method were used to label entities. B represents the start position of the entity, I represents the inside of the entity, E represents the end position of the entity, and O represents the other. Five pig disease

attributes, including pig type, disease name, body parts, symptom, and medicine, were labeled in the corpus text. Finally, the annotated pig disease corpus containing 7518 entities was obtained under the guidance of pig disease experts. The statistical information of labeled entities is presented in Table 1.

Table 1. Statistics of annotated entities.

Category	Category Definition	Examples	Numbers	Proportion of the Total
Type	Name of different types of pig	妊娠母猪, 仔猪 (Pregnant sows, piglets)	735	9.78%
Disease	Name of pig disease	猪丹毒, 胸膜炎 (Porcine erysipelas, pleurisy)	958	12.74%
Body parts	Body position, organs, and system of pigs	心脏, 巨噬细胞 (Heart, Macrophages)	2063	27.44%
Symptom	External performance caused by diseases	气喘, 咳嗽, 水肿 (Asthma, cough, swollen)	2973	39.55%
Medicine	Medications for treating diseases	替米考星, 克林霉素 (Timicosin, clindamycin)	789	10.49%
Total			7518	100%

2.1.4. Construction of Lexicon and Pre-Training Word Embedding

We constructed a lexicon in the pig disease domain based on Corpus II and professional books. Firstly, the most commonly used professional terms were extracted from Corpus II by word segmentation and frequency statistics. Then, some professional words in the glossary of professional books were manually added to the dictionary under the guidance of veterinarians, such as “猪副嗜血杆菌病 (*Haemophilus parasuis*)”, “噻苯达唑 (*Thiabendazole*)”, and “内阿米巴原虫 (*Entamoeba* spp.)”. Finally, the pig disease lexicon comprising 2391 professional terms was obtained for understanding the specific words and technical terms. Subsequently, this lexicon was incorporated into the built-in dictionary of Jieba to avoid the incorrect segmentation of words.

To obtain a high-quality embedded representation of pig diseases, we trained Corpus I and the lexicon. The Gensim tool (v.3.8.3) was used to train the Word2Vec model with a word vector dimension of 200. The construction process of the lexicon and the pre-training of word embeddings are illustrated in Figure 1.

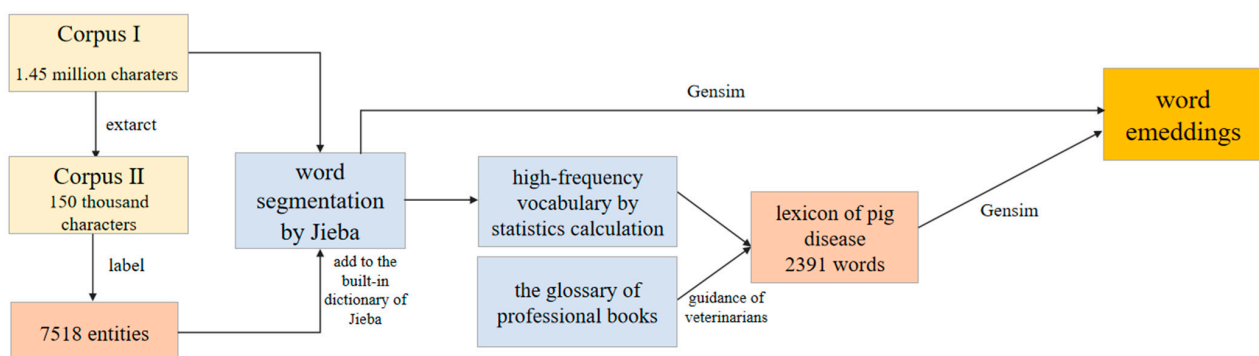


Figure 1. Construction process of lexicon and pre-training word embedding.

2.2. Methods

The structure of the PDCNER model is shown in Figure 2. Firstly, the Chinese sentences in the pig disease corpus are converted into a character–word pair sequence, and both Chinese character features and lexicon features are used as inputs. Secondly, a lexicon adapter is added between Transformer layers, which is used to dynamically extract the most relevant matching items. The word of each character uses the bi-linear attention mechanism from character to word, and the lexicon adapter is applied between adjacent

Transformers in BERT. The lexicon features and BERT representations are fully interacted through multi-layer encoders in BERT, so that lexicon knowledge can be effectively integrated into BERT. The contrastive loss layer is above the Lexicon-Enhanced BERT encoder, ensuring that similar samples are as close as possible, while dissimilar samples are as far apart as possible. Embeddings of the same type of entity are treated as positive samples, whereas embeddings of different types of entities are treated as negative samples. Considering the correlation between consecutive labels, a Conditional Random Field (CRF) layer is employed to label the sequence.

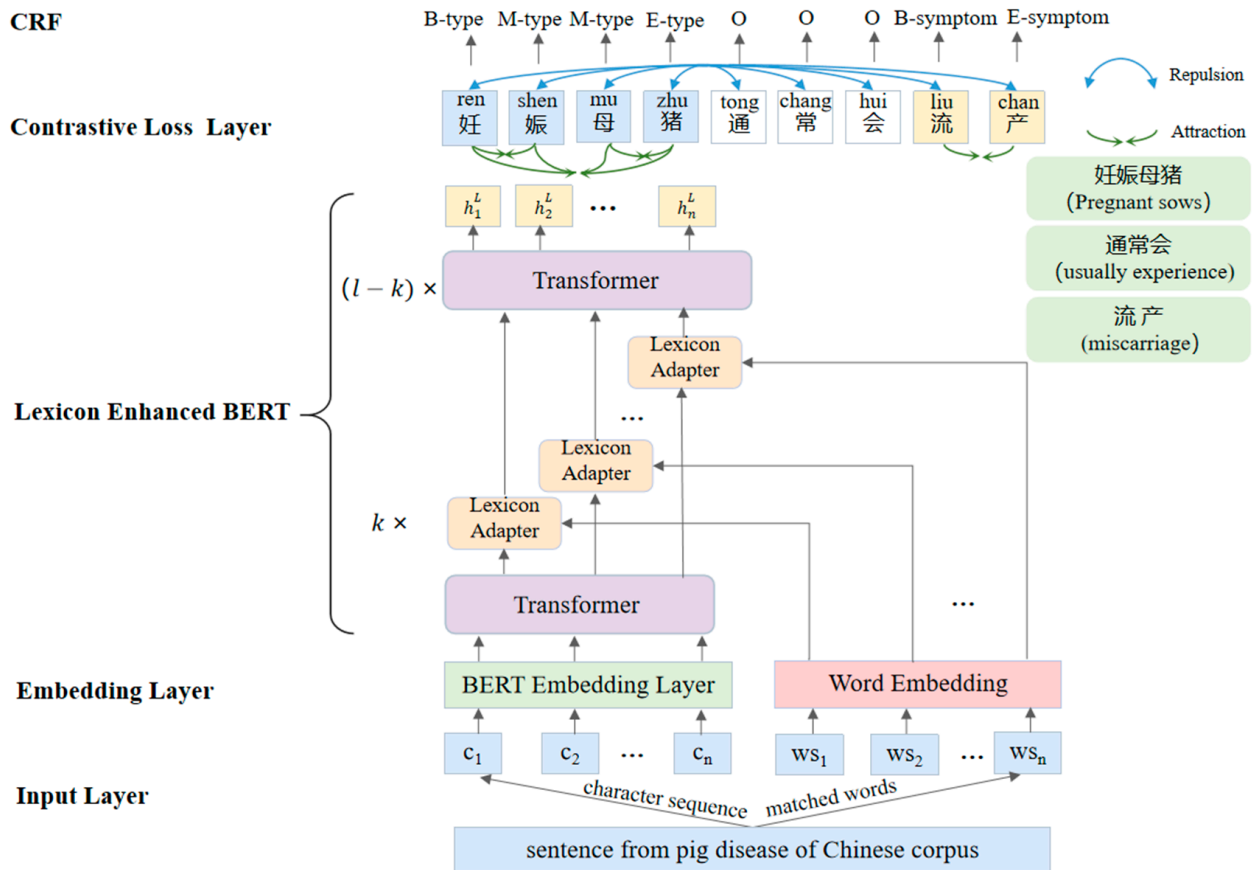


Figure 2. Structure of PDCNER.

2.2.1. Char–Word Pair Sequence

According to the Lexicon-Enhanced BERT in [12], we firstly expand the character sequence into a sequence of character–word pairs for applying lexical information of pig disease.

A Chinese sentence with n characters, $S_C = \{c_1, c_2, \dots, c_n\}$. We identify all potential words within the sentence by comparing the character sequence with the lexicon of pig disease. To achieve this, we first create a Trie data structure based on the lexicon. Then, we examine all character subsequences within the sentence, matching them with the Trie to identify all potential words. For instance, in the truncated sentence “非洲猪瘟 (African swine fever)”, we can identify five distinct words: “非洲 (Africa)”, “非洲猪 (African pig)”, “猪瘟 (swine fever)”, “瘟 (epidemic disease)”, and “非洲猪瘟 (African swine fever)”. In the field of pig disease, African swine fever is a complete disease name and should not be separated. Following this, for each matched word, we associate it with the characters that compose it. In conclusion, we pair each character with its associated words and transform the Chinese sentence into a sequence of character–word pairs, represented as:

$$S_{Cw} = \{(c_1, ws_1), (c_2, ws_2), \dots, (c_n, ws_n)\} \tag{1}$$

where c_i denotes the i -th character in the sentence, and ws_i signifies the words matched and assigned to c_i .

2.2.2. Lexicon Adapter

Using the lexicon adapter proposed in LEBERT [12], the pig disease lexicon information is directly injected into BERT for integrating lexical features.

For the i -th character in a character–word sequence, the input is (h_i^c, x_i^{ws}) , and h_i^c represents the character vector, the output of a transformation layer in BERT. $x_i^{ws} = \{x_{i1}^w, x_{i2}^w, \dots, x_{im}^w\}$ represents a group of word embeddings. The j -th word in x_i^{ws} is represented as follows:

$$x_{ij}^w = e^w(w_{ij}) \quad (2)$$

where e^w is the pre-trained word embedding list and w_{ij} represents the j -th word in ws_i .

To align these two different representations, a nonlinear transformation is used for each word vector:

$$v_{ij}^w = W_2 \left(\tanh(W_1 x_{ij}^w + b_1) \right) + b_2 \quad (3)$$

where $W_1 \in \mathbb{R}^{d_c \times d_w}$, $W_2 \in \mathbb{R}^{d_c \times d_c}$, and d_w and d_c represent the dimension of word embedding and BERT's hidden size, respectively. b_1 and b_2 are scaler bias.

Each character is associated with a variety of words, but the degree of contribution from each word differs. For instance, in the field of pig diseases, the words “非洲 (Africa)” and “猪瘟 (swine fever)” are more important than “非洲猪 (African pigs)” and “瘟 (epidemic disease)”. In order to find the most relevant words, the character–word attention mechanism is used. The correlation of each word is calculated as follows:

$$a_i = \text{softmax} \left(h_i^c W_{\text{attn}} V_i^T \right) \quad (4)$$

where $W_{\text{attn}} \in \mathbb{R}^{d_c \times d_c}$ is the bi-linear attention mechanism. All v_{ij}^w are assigned to the i -th character $V_i = (v_{i1}^w, \dots, v_{im}^w)$, where m is the total number of assigned words.

Lastly, the lexicon information is integrated into the vector representation of the character.

$$\tilde{h}_i = h_i^c + \sum_{j=1}^m a_{ij} v_{ij}^w \quad (5)$$

2.2.3. Lexicon-Enhanced BERT

The lexicon adapter is attached between transformer layers in BERT so that the knowledge of the pig disease lexicon can be injected into BERT. A sequence of characters $\{c_1, c_2, \dots, c_n\}$ is input into the input embedding of BERT and then $E = \{e_1, e_2, \dots, e_n\}$ is obtained by adding token, segmentation, and position embedding. After that, E is input into the Transformer encoders. Each layer is as follows.

$$G = \text{LayerNormalization} \left(H^{l-1} + \text{Multiheadattention} \left(H^{l-1} \right) \right) \quad (6)$$

$$H^l = \text{LayerNormalization} (G + \text{FFN}(G))$$

where $H^l = \{h_1^l, h_2^l, \dots, h_n^l\}$ represents the output of the l -th layer and $H^0 = E$. FFN represents the two-layer feed-forward network with RELU as the hidden activation function.

The lexicon information was input between the k -th and $(k+1)$ -th layer Transformer. $H^k = \{h_1^k, h_2^k, \dots, h_n^k\}$ were first obtained after k consecutive Transformer layers. Subsequently, each character–word pair (h_i^k, x_i^{ws}) was processed through the lexicon adapter to obtain a new hidden layer representation and the i -th pair was converted into \tilde{h}_i^k accordingly.

$$\tilde{h}_i^k = \text{Lexicon Adapter} \left(h_i^k, x_i^{ws} \right) \quad (7)$$

$H^k = \{h_1^k, h_2^k, \dots, h_n^k\}$ are input into the remaining $(L-K)$ Transformer, as there are 12 layers of Transformers in BERT. Finally, the output of the L -th Transformer H^L used for the name entity recognition task is obtained.

2.2.4. Lexicon-Enhanced Contrastive Learning

The normalized temperature-scaled cross-entropy loss, denoted as NT-Xent, was used as our contrastive loss function [41]. For each training iteration, we randomly select N texts from the datasets to form a mini-batch, which yields $2N$ feature representations. The model is then trained to identify each data point's corresponding pair among the $2(N - 1)$ negative samples present within the batch.

$$\mathcal{L}_{i,j} = -\log \frac{\exp(s(r_i, r_j)/T)}{\sum_k^{2N} I_{[k \neq i]} \exp(s(r_i, r_k)/T)} \quad (8)$$

where r_i, r_j represent the embedding for entities of the same type, whereas r_i, r_k denote the embedding for entities of different types. $s()$ refers to the cosine similarity function, where T acts as the temperature parameter, and I serves as an indicator function. Ultimately, we calculate the final contrastive loss by averaging the classification losses of all $2N$ instances within the batch.

2.2.5. Construction of Positive and Negative Pairs

Construction of positive and negative samples is very important for contrastive learning, which aims to learn better feature representations of data. We construct them according to the following steps.

1. Positive sample pairs based on the same label token

We take tokens of the same entity type as positive sample pairs.

Direct Matching: Within a sequence or document, directly identify all tokens labeled with the same entity type and use them as positive sample pairs. For example, if both “fattening pigs (育肥猪)” and “breastfeed sows (哺乳母猪)” are labeled as “type” entities in a sentence, they form a positive sample pair.

Cross-Document Matching: Identify all tokens labeled with the same entity type across different paragraphs or sentences. Although these tokens may appear in different contexts, they can still be considered positive sample pairs due to their shared entity type.

2. Negative sample pairs based on entities of different types

We take tokens of different entity types as negative sample pairs. Specifically, we select tags labeled with different entity types as negative sample pairs. For example, in different sentences, “Porcine blue ear disease (猪蓝耳病)” is labeled as a “disease” entity, and “盐酸头孢塞夫 (ceftazidime hydrochloride)” is labeled as a “medicine” entity, thus constituting a negative sample pair. Additionally, entities labeled as ‘O’ can also form negative sample pairs with entities labeled with the five specified types.

3. Results

3.1. Evaluation

To identify a named entity for pig diseases, it is necessary to correctly identify both the boundaries of the entity and its corresponding categories. The proposed PDCNER model is evaluated using standard measures, Precision (P), Recall (R), and F1, which are computed using Equations (9)–(11).

$$P = \frac{T_P}{T_P + F_P} \times 100\% \quad (9)$$

$$R = \frac{T_P}{T_P + F_n} \times 100\% \quad (10)$$

$$F1 = \frac{2PR}{P + R} \times 100\% \quad (11)$$

where T_p represents the number of positive samples that are accurately predicted, while F_p denotes the count of positive samples that are inaccurately predicted. Additionally, F_n stands for the number of negative samples that are incorrectly predicted.

3.2. Experimental Settings

Experiments. The hardware environment that the experimental research relied on was Intel(R) Xeon(R) Silver4116 CPU@2.10 GHz (Intel, Santa Clara, CA, USA), GPU@NVIDIA Tesla P100 (NVIDIA, Santa Clara, CA, USA). The software environment was Python3.8 and tensorflow 2.0.0. The model parameters were set as follows: based on BERT_{BASE} [16] version, with 12 transformer layers, 768 hidden layers, and 12 multi-head attention mechanisms. The lexicon corresponds to the vocabulary of pre-trained word embeddings in field of pig disease. During the training process, we incorporate the Lexicon Adapter between the first and second Transformer layers of BERT [12]. Meanwhile, the parameters of BERT and the pre-trained word embedding were fine-tuned. The hyperparameters are shown in Table 2. This study randomly divided the datasets into training, validation, and test sets according to a ratio of 7:2:1.

Table 2. Hyperparameters of the proposed model.

Hyper Param	Value	Hyper Param	Value
batch_size	16	dropout	0.5
learning rate	0.00001	optimizer	Adam
epoch	20	maximum sentence length	256

3.3. Results

3.3.1. Comparison with Baseline Models

We evaluate the effectiveness of PDCNER against widely used NER models on the pig disease corpus. The overall findings of our experiments are presented in Table 3.

Table 3. Comparison of experimental results for different NER models.

Model Category	Model	P (%)	R (%)	F1 (%)
baseline model without pre-training	BILSTM_CRF	71.58	67.51	69.49
pre-trained model	BERT-BiLSTM-CRF	80.29	84.73	82.45
	BERT-CRF	81.62	84.39	82.98
	BERT-CNN-CRF	82.44	80.55	81.48
	BERT-WWM-ext-BiLSTM-CRF	81.73	85.47	83.56
	RoBERTa-BiLSTM-CRF	81.64	85.31	83.43
pre-trained model with lexicon	BERT-BiLSTM-CRF-SoftLexicon	82.99	84.73	83.85
	LEBERT	87.18	86.45	86.81
	PDCNER (ours)	87.76	86.97	87.36

The bold in the last row denotes the results we obtained in our model.

According to Table 3, several inferences can be drawn. Firstly, the F1-scores of pre-trained models are higher than those of models without pre-training by more than 10%, such as the BILSTM_CRF model. This indicates that pre-trained models, with their deeper network structures, have learned more language features and enhanced their ability to recognize entities. Secondly, the pre-trained model incorporating lexicon information demonstrates better performance than models without lexicon integration. This improvement is primarily due to the integration of dictionary information specific to the pig disease domain, which effectively captures entity boundaries and word information. Thirdly, the

recognition performance of our model significantly outperforms other models, achieving a micro-average F1-score of 87.36% in recognizing five major entities. This experimental result indicates that our model can effectively improve entity recognition performance in the domain of pig disease.

Comparative analysis with the results of BERT-BiLSTM-CRF reveals notable improvement in the precision, recall, and F1-score of PDCNER, with improvements of 7.47 percentage points, 2.24 percentage points, and 4.91 percentage points, respectively.

Through comparative evaluation utilizing the same datasets and downstream model, PDCNER demonstrates superior accuracy in identifying pig disease entities compared to LEBERT, exhibiting improvements in precision, recall, and F1-score by 0.58 percentage points, 0.52 percentage points, and 0.55 percentage points, respectively.

3.3.2. The Recognition Effect on Different Entities

For better understanding of the proposed approach, we evaluate the PDCNER model separately on the five entities type, disease, body parts, symptom, and medicine, which are presented in Figure 3. We found that the F1-scores for type, disease, and medicine all exceeded 90%, with the F1-score for type being the highest at 95.41%. Conversely, the lowest F1-score was for the entity of symptom, at 80.92%.

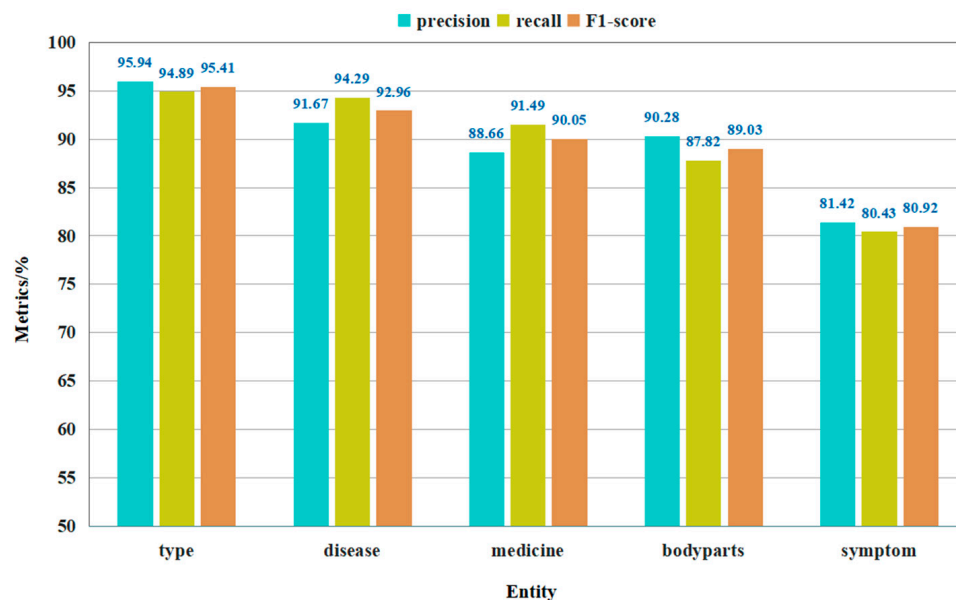


Figure 3. Precision, recall, and F1-score of PDCNER in recognizing five major entities.

On the other hand, the F1-scores of disease entities and medicine entities were 92.96% and 90.05%, respectively. Both disease and medicine entities include a large number of technical terms, yet the method proposed in this paper achieves a good recognition effect on these two entities. The results demonstrate that PDCNER fully utilizes both Chinese character features and lexicon knowledge in the pig disease domain at the input level, and the lexicon adapter can effectively leverage pig disease knowledge.

3.3.3. The Recognition Effect on Small Sample

In order to verify the reliability and robustness of PDCNER in the condition of scarce data for entity recognition, we used 1%, 10%, 30%, and 50% of labeled samples for experimentation. The results can be found in Figure 4 and Table 4. The result shows that the PDCNER model has obvious improvement compared to BERT-BiLSTM-CRF and LEBERT.

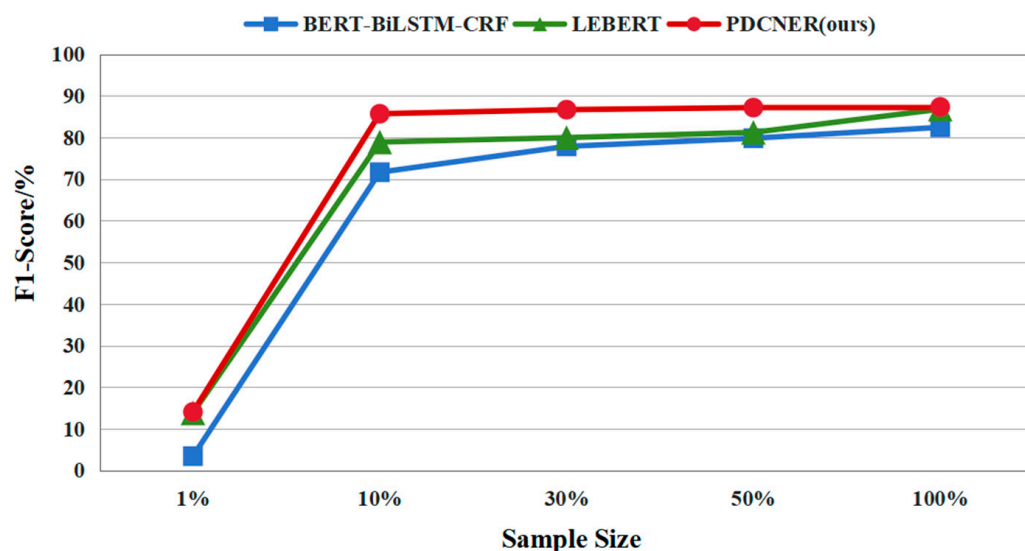


Figure 4. F1-score of 3 models with different sample size.

Table 4. Results of the small-sample experiments.

Model	1%			10%			30%		
	P	R	F1	P	R	F1	P	R	F1
BERT-BiLSTM-CRF	31.79	1.80	3.41	67.92	75.84	71.66	74.59	81.44	77.87
LEBERT	17.39	11.43	13.79	74.81	83.47	78.91	74.05	80.14	76.97
PDCNER (ours)	18.18	11.43	14.04	85.00	86.44	85.71	87.22	86.17	86.69
Model	50%			100%					
	P	R	F1	P	R	F1			
BERT-BiLSTM-CRF	79.18	83.45	81.26	80.29	84.73	82.45			
LEBERT	81.09	83.4	82.23	87.18	86.45	86.81			
PDCNER (ours)	87.68	86.71	87.19	87.76	86.97	87.36			

The bold in the last row denotes the results we obtained in our model.

The F1-score of PDCNER reaches 85.71% when the sample size is 10%, which is only 1.65% lower than that of the full sample. As the sample size increases to 30%, the F1-score of the PDCNER model further improves to 86.69%, showing a marginal decrease of only 0.67% compared to the full sample. Moreover, it outperforms the BERT-BiLSTM-CRF and LEBERT models by 14.05% and 6.8%, respectively, with only 10% samples available. These results demonstrate the PDCNER model's capability to achieve higher recognition accuracy even under data scarcity scenarios.

3.3.4. The Recognition Effect on Public Datasets

To assess the generalization capability of PDCNER, we conducted evaluations across three public benchmark datasets: Weibo, Ontonotes, and Resume. As illustrated in Table 5, the PDCNER model achieved the highest F1-scores across all three datasets. Specifically, the F1 values of PDCNER were 9.78%, 4.33%, and 0.82% higher than those of the BERT-BiLSTM-CRF model for the Weibo, Ontonotes, and Resume datasets, respectively. Additionally, compared to the LEBERT model, the PDCNER method showed improvements of 2.53%, 0.37%, and 0.06% for the same datasets. These results indicate that PDCNER not only exhibits superior performance on the pig disease corpus but also demonstrates remarkable generalization capability across different domains.

Table 5. F1-score for each model on public datasets.

Model	Weibo			Ontonotes			Resume		
	P	R	F1	P	R	F1	P	R	F1
BERT-BiLSTM-CRF	71.29	67.1	69.13	84.11	80.2	82.11	96.56	97.23	95.89
LEBERT	78.2	74.64	76.38	89.76	82.68	86.07	95.89	97.42	96.65
PDCNER (ours)	81.42	76.56	78.91	89.68	83.43	86.44	96.07	97.36	96.71

The bold in the last row denotes the results we obtained in our model.

4. Discussion

Our experimental results demonstrate that the proposed PDCNER method effectively utilizes pig disease feature information by seamlessly integrating it into the BERT architecture. Specifically, the Lexicon Adapter is inserted between the first and second transformer layers within BERT, facilitating the infusion of pig disease lexicon knowledge into the model's representations.

In addition to the lexicon-enhanced BERT's contribution to improvement, contrastive learning has also played a significant role. The lexicon-enhanced contrastive learning method improves the model's capacity for semantic representation of text based on the normalized temperature-scaled cross-entropy loss function. This loss function enhances the model's ability to identify similar entities by minimizing the distance between positive samples (i.e., the same type of entities) while simultaneously maximizing the distance from negative samples (i.e., different types of entities). Consequently, this approach reduces the risk of the model incorrectly classifying different entities. Furthermore, it enables the model to adjust the weighting of distance measurements, allowing it to focus more on distinguishing subtle differences between different entities during the training process.

Judging from the recognition results of different entities, despite the sample sizes of type, disease, and medicine entities being only about 10%, their F1 values are all above 90%. In contrast, although symptoms have the largest sample size, accounting for 39.55%, their F1 value is the lowest. This disparity arises primarily because the boundaries of pig type and disease entities are very clear, whereas the boundaries of symptom entities are more ambiguous. For instance, type entities typically end with terms such as "pigs (猪)" (e.g., sick pigs (患病猪), conservation pigs (保育猪), fattening pigs (育肥猪)), and disease entities usually end with terms like "disease (病)," "inflammation (炎)," and "plague (瘟)" (e.g., porcine blue ear disease (猪蓝耳病), necrotic enteritis (坏死性肠炎), African swine fever (非洲猪瘟)). In contrast, symptom entities are often composed of complex verbs, conjunctions, and modifiers, such as "feed intake continuing to decrease (采食量持续下降)" and "continuous spasmodic cough (连续痉挛性咳嗽)". Moreover, the average length of complex symptom entities is approximately eight Chinese characters, which contributes to a lower overall recognition rate.

5. Conclusions

High-quality extraction of entities related to pig diseases is critical for intelligent consultation, question answering, technical recommendations, and other application scenarios.

In this study, we constructed a corpus, labeled datasets, and lexicon for Chinese named entity recognition specific to pig diseases, encompassing 152,596 characters, 7518 entities, and 2391 professional terms. To tackle the challenges of entity identification in the pig disease domain, such as the scarcity of annotated data, numerous technical terms, and fuzzy boundaries, we propose the PDCNER model. This model integrates lexicon information from the pig disease domain into BERT's Transformer layers at the lower level and employs contrastive learning to enhance representation quality and generalization capability. The results indicate that the PDCNER model surpasses the performance of BERT-BiLSTM-CRF and other mainstream models, achieving precision, recall, and F1-score of 87.76%, 86.97%, and 87.36%, respectively. This demonstrates high-quality entity recognition in

the field of pig diseases. Moreover, small-sample experiments confirm that our model is more suitable than other models for completing the named entity recognition task in data-scarce scenarios. Experiments on public datasets also verify its generalization ability. Our approach provides a reference for improving NER performance in domain-specific applications. It can be easily applied to other fields, facilitating seamless adaptation for named entity identification across diverse contexts.

In future work, we plan to focus on the following aspects:

- (1) We aim to enhance the identification of more fine-grained entity types in the pig disease domain, including appearance symptoms and anatomical symptoms.
- (2) We will further optimize the pig disease corpus and construct a more extensive professional vocabulary to enhance NER performance.
- (3) We will explore the application of our model in other animal diseases, such as chicken or cow diseases, and investigate methods to handle nested entities and discontinuous entities based on our approach.

Author Contributions: Conceptualization, C.P. and Q.L.; Data curation, R.M., H.L., S.W. and L.H.; Formal analysis, W.W. and H.H.; Investigation, S.W.; Methodology, C.P. and X.W.; Software, X.W.; Validation, Q.Y., R.J. and W.M.; Writing—original draft, C.P. and X.W.; Writing—review and editing, C.P. and Q.L. All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded by National Science and Technology Major Project (2021ZD0113802).

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

Data Availability Statement: As the datasets used in this manuscript will be used for other technical research, they are available from the corresponding author upon reasonable request. The codes have been open-sourced at <https://github.com/tufeifei923/pdner> (accessed on 5 August 2024).

Acknowledgments: We thank the editors and the anonymous reviewers for their valuable suggestions.

Conflicts of Interest: The authors declare no conflicts of interest.

Appendix A. Data Source Details

No	Type	Example
1	Professional books	<p>Jeffrey J. Zimmerman, Locke A. Kerri, Alejandro Raminez. etc., Editor-in-Chief. Hanchun Yang, main translation. Disease of Swine. North United publishing Media Co., Ltd., Liaoning science and technology publishing house: Beijing, China, 2022.</p> <p>Yousheng Xu. Primary color atlas of scientific pig raising and pig disease prevention and control. China Agricultural Publishing House: Beijing, China, 2017.</p> <p>Changyou Li, Xiaocheng Li. Prevention and control technology of swine epidemic disease. China Agricultural Publishing House: Beijing, China, 2015.</p> <p>Jianxin Zhang. Diagnosis and control of herd pig epidemic disease. He'nan Science and Technology Press: Zhengzhou, China, 2014.</p> <p>Chaoying Luo, Guibo Wang. Prevention and treatment of pig diseases and safe medication. Chemical industry press: Beijing, China, 2016, etc.</p>
2	Standard specification	<p>Technical Specification for Quarantine of Porcine Reproductive and Respiratory Syndrome (SN/T 1247-2022), Diagnostic Techniques for Mycoplasma Pneumonia in Swine (NY/T 1186-2017), Diagnostic Techniques for Infectious Pleuropneumonia in Swine (NY/T 537-2023), Diagnostic Techniques for Swine Dysentery (NY/T 545-2023), Technical Specification for Quarantine of Porcine Rotavirus Infection (SN/T 5196-2020), etc.</p>
3	Technological specification	<p>Technical specification for prevention and control of highly pathogenic blue ear disease in pigs, technical specification for prevention and control of foot-and-mouth disease, technical specification for prevention and control of classical swine fever, etc.</p>
4	Policy paper	<p>Ministry of Agriculture and Rural Affairs "List of Class I, II and III Animal Diseases", The Ministry of Agriculture issued the "Guiding Opinions on Prevention and Control of Highly Pathogenic Porcine Blue Ear Disease (2017–2020)", Notice of National Guiding Opinions on Prevention and Control of Classical Swine Fever (2017–2020), etc.</p>
5	Relevant industry website	<p>China Veterinary Website (https://www.cadc.net.cn/sites/MainSite/, 10 December 2023), Big Animal Husbandry Website (https://www.dxumu.com/, 10 December 2023), Huinong Website (https://www.cnbnb.com/, 30 January 2024), etc.</p>

References

1. Li, J.; Sun, A.X.; Han, J.L.; Li, C.L. A Survey on Deep Learning for Named Entity Recognition. *IEEE Trans. Knowl. Data Eng.* **2022**, *34*, 50–70. [[CrossRef](#)]
2. Cheng, J.R.; Liu, J.X.; Xu, X.B.; Xia, D.W.; Liu, L.; Sheng, V. A review of Chinese named entity recognition. *KSII Trans. Internet Inf. Syst.* **2021**, *15*, 2012–2030.
3. Mi, B.G.; Fan, Y. A review: Development of named entity recognition (NER) technology for aeronautical information intelligence. *Artif. Intell. Rev.* **2022**, *56*, 1515–1542.
4. Liu, P.; Guo, Y.; Wang, F.; Li, G. Chinese named entity recognition: The state of the art. *Neuro Comput.* **2022**, *473*, 37–53. [[CrossRef](#)]
5. Qiu, X.; Sun, T.; Xu, Y.; Shao, Y.; Dai, N.; Huang, X. Pre-trained models for natural language processing: A survey. *Sci. China Technol. Sci.* **2020**, *63*, 1872–1897. [[CrossRef](#)]
6. Kang, Y.L.; Sun, L.B.; Zhu, R.B.; Li, M.Y. Survey on Chinese named entity recognition with deep learning. *J. Huazhong Univ. Sci. Technol. (Nat. Sci. Ed.)* **2022**, *50*, 44–53.
7. Huang, S.B.; Sha, Y.P.; Li, R.S. A Chinese named entity recognition method for small-scale dataset based on lexicon and unlabeled data. *Multimed. Tools Appl.* **2023**, *82*, 2185–2206. [[CrossRef](#)]
8. Dang, X.; Wang, L.; Dong, X.; Li, F.; Deng, H. Improving Low-Resource Chinese Named Entity Recognition Using Bidirectional Encoder Representation from Transformers and Lexicon Adapter. *Appl. Sci.* **2023**, *13*, 10759. [[CrossRef](#)]
9. Zhang, Y.; Yang, J. Chinese NER using lattice LSTM. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, Melbourne, Australia, 15–20 July 2018; Volume 1, pp. 1554–1564.
10. Gui, T.; Ma, R.; Zhang, Q.; Zhao, L.; Jiang, Y.G.; Huang, X. CNN-Based Chinese NER with lexicon rethinking. In Proceedings of the 28th International Joint Conference on Artificial Intelligence, Macao, China, 10–16 August 2019; Volume 8, pp. 4982–4988.
11. Gui, T.; Zou, Y.; Zhang, Q.; Peng, M.; Fu, J.; Wei, Z.; Huang, X.J. A Lexicon-Based Graph Neural Network for Chinese NER. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), Hong Kong, China, 3–7 November 2019; pp. 1040–1050.
12. Liu, W.; Fu, X.; Zhang, Y.; Xiao, W. Lexicon enhanced Chinese sequence labeling using BERT adapter. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, Bangkok, Thailand, 1–6 August 2021.
13. Ma, R.; Peng, M.; Zhang, Q.; Huang, X. Simplify the usage of lexicon in Chinese NER. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Online, 5–10 July 2020; pp. 5951–5960.
14. Li, X.; Yan, H.; Qiu, X.; Huang, X. FLAT: Chinese NER using flat-lattice transformer. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Online, 5–10 July 2020; pp. 6836–6842.
15. Mengge, X.; Yu, B.; Liu, T.; Zhang, Y.; Meng, E.; Wang, B. Porous lattice transformer encoder for Chinese NER. In Proceedings of the 28th International Conference on Computational Linguistics, Barcelona, Spain, 8–13 December 2020; pp. 3831–3841.
16. Devlin, J.; Chang, M.W.; Lee, K.; Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv* **2018**, arXiv:1810.04805.
17. Peters, M.E.; Neumann, M.; Iyyer, M.; Gardner, M.; Clark, C.; Lee, K.; Zettlemoyer, L. Deep contextualized word representations. *arXiv* **2018**, arXiv:1802.05365.
18. Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; Stoyanov, V. RoBERTa: A robustly optimized BERT pretraining approach. *arXiv* **2019**, arXiv:1907.11692.
19. Sun, Y.; Wang, S.; Li, Y.; Feng, S.; Wu, H. ERNIE: Enhanced representation through knowledge integration. *arXiv* **2019**, arXiv:1904.09223v1.
20. Lan, Z.; Chen, M.; Goodman, S.; Gimpel, K.; Sharma, P.; Soricut, R. ALBERT: A lite BERT for self-supervised learning of language representations. In Proceedings of the 8th International Conference on Learning Representations (ICLR), Addis Ababa, Ethiopia, 26–30 April 2020.
21. Li, S.; Bai, Z.Q.; Zhao, S.; Jiang, G.S.; Shan, L.L.; Zhang, L. A LEBERT-based model for named entity recognition. In Proceedings of the 2021 3rd International Conference on Artificial Intelligence and Advanced Manufacture (AIAM), ACM International Conference Proceeding Series, Manchester, UK, 23–25 October 2021; pp. 980–983.
22. Yan, Y.M.; Li, R.M.; Wang, S.R.; Zhang, F.Z.; Wu, W.; Xu, W. ConSERT: A contrastive framework for self-supervised sentence representation transfer. *arXiv* **2021**, arXiv:2105.11741.
23. Gao, T.; Yao, X.; Chen, D. SimCSE: Simple Contrastive Learning of Sentence Embeddings. EMNLP. In Proceedings of the Conference on Empirical Methods in Natural Language Processing, Punta Cana, Dominican Republic, 7–11 November 2021; pp. 6894–6910.
24. Huang, Y.; He, K.; Wang, Y.; Zhang, X.; Gong, T.; Mao, R.; Li, C. COPNER: Contrastive learning with prompt guiding for few-shot named entity recognition. In Proceedings of the 29th International Conference on Computational Linguistics, Gyeongju, Republic of Korea, 12–17 October 2022; pp. 2515–2527.
25. He, K.; Mao, R.; Huang, Y.; Gong, T.; Li, C.; Cambria, E. Template-Free Prompting for Few-Shot Named Entity Recognition via Semantic-Enhanced Contrastive Learning. In *IEEE Transactions on Neural Networks and Learning Systems*; IEEE: Piscataway, NJ, USA, 2023. [[CrossRef](#)]
26. Li, X.W.; Li, X.L.; Zhao, M.K.; Yang, M.; Yu, R.G.; Yu, M.; Yu, J. CLINER: Exploring task-relevant features and label semantic for few-shot named entity recognition. *Neural Comput. Appl.* **2023**, *36*, 4679–4691. [[CrossRef](#)]

27. Chen, P.; Wang, J.; Lin, H.F.; Zhao, D.; Yang, Z.H.; Wren, J. Few-shot biomedical named entity recognition via knowledge-guided instance generation and prompt contrastive learning. *Bioinformatics* **2023**, *39*, btad496. [[CrossRef](#)]
28. Sahadevan, S.; Hofmann-Apitius, M.; Schellander, K.; Tesfaye, D.; Fluck, J.; Friedrich, C.M. Text mining in livestock animal science: Introducing the potential of text mining to animal sciences. *J. Anim. Sci.* **2012**, *90*, 3666–3676. [[CrossRef](#)] [[PubMed](#)]
29. Oh, H.S.; Lee, H. Named Entity Recognition for Pet Disease Q&A System. *J. Digit. Contents Soc.* **2022**, *23*, 765–771.
30. Kung, H.Y.; Yu, R.W.; Chen, C.H.; Tsai, C.W.; Lin, C.Y. Intelligent pig-raising knowledge question-answering system based on neural network schemes. *Agron. J.* **2021**, *113*, 906–922. [[CrossRef](#)]
31. Zhang, D.; Zheng, G.; Liu, H.; Ma, X.; Xi, L. AWdpCNER: Automated Wdp Chinese Named Entity Recognition from Wheat Diseases and Pests Text. *Agriculture* **2023**, *13*, 1220. [[CrossRef](#)]
32. Veena, G.; Kanjirangat, V.; Gupta, D. AGRONER: An unsupervised agriculture named entity recognition using weighted distributional semantic model. *Expert Syst. Appl.* **2023**, *229*, 120440. [[CrossRef](#)]
33. Zhang, L.; Nie, X.; Zhang, M.; Gu, M.; Geissen, V.; Ritsema, C.J.; Niu, D.; Zhang, H. Lexicon and attention-based named entity recognition for kiwifruit diseases and pests: A Deep learning approach. *Front. Plant Sci.* **2022**, *13*, 1053449. [[CrossRef](#)] [[PubMed](#)]
34. Guo, X.C.; Lu, S.H.; Tang, Z.; Bai, Z.; Diao, L.; Zhou, H.; Li, L. CG-ANER: Enhanced contextual embeddings and glyph features-based agricultural named entity recognition. *Comput. Electron. Agric.* **2022**, *194*, 106776. [[CrossRef](#)]
35. Huang, B.; Lin, Y.; Pang, S.; Fu, L. Named Entity Recognition in Government Audit Texts Based on ChineseBERT and Character-Word Fusion. *Appl. Sci.* **2024**, *14*, 1425. [[CrossRef](#)]
36. Guo, Y.; Feng, S.; Liu, F.; Lin, W.; Liu, H.; Wang, X.; Su, J.; Gao, Q. Enhanced Chinese Domain Named Entity Recognition: An Approach with Lexicon Boundary and Frequency Weight Features. *Appl. Sci.* **2024**, *14*, 354. [[CrossRef](#)]
37. Jia, Y.C.; Zhu, D.J. Medical Named Entity Recognition Based on Deep Learning. *Comput. Syst. Appl.* **2022**, *31*, 70–81. (In Chinese)
38. Du, J.; Hao, Y.; Song, F. Research and Development of Named Entity Recognition in Chinese Electronic Medical Record. *Acta Electron. Sin.* **2022**, *50*, 3030–3053.
39. Cao, L.L.; Wu, C.C.; Luo, G.; Guo, C.; Zheng, A.N. Online biomedical named entities recognition by data and knowledge-driven model. *Artif. Intell. Med.* **2024**, *150*, 102813. [[CrossRef](#)] [[PubMed](#)]
40. Zhai, Z.; Fan, R.; Huang, J.; Xiong, N.; Zhang, L.; Wan, J.; Zhang, L. A Named Entity Recognition Method Based on Knowledge Distillation and Efficient Global Pointer for Chinese Medical Texts. *IEEE Access* **2024**, *12*, 83563–83574. [[CrossRef](#)]
41. Chen, T.; Kornblith, S.; Norouzi, M.; Hinton, G. A simple framework for contrastive learning of visual representations. *arXiv* **2020**, arXiv:2002.05709.

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.