

Research Article

Attribute Extraction Study in the Field of Military Equipment Based on Distant Supervision

Xindong You,¹ Meijing Yang,¹ Junmei Han,² Jiangwei Ma,¹ Gang Xiao,² and Xueqiang Lv¹ 

¹Beijing Key Laboratory of Internet Culture Digital Dissemination, Beijing Information Science and Technology University, Beijing, China

²National key Laboratory for Complex Systems Simulation, Institute of Systems Engineering, China

Correspondence should be addressed to Xueqiang Lv; lxq@bistu.edu.cn

Received 6 August 2021; Revised 24 September 2021; Accepted 23 October 2021; Published 23 November 2021

Academic Editor: Honghao Gao

Copyright © 2021 Xindong You et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

The effective organization and utilization of military equipment data is an important cornerstone for constructing knowledge system. Building a knowledge graph in the field of military equipment can effectively describe the relationship between entity and entity attribute information. Therefore, relevant personnel can obtain information quickly and accurately. Attribute extraction is an important part of building the knowledge graph. Given the lack of annotated data in the field of military equipment, we propose a new data annotation method, which adopts the idea of distant supervision to automatically build the attribute extraction dataset. We convert the attribute extraction task into a sequence annotation task. At the same time, we propose a RoBERTa-BiLSTM-CRF-SEL-based attribute extraction method. Firstly, a list of attribute name synonyms is constructed, then a corpus of military equipment attributes is obtained through automatic annotation of semistructured data in Baidu Encyclopedia. RoBERTa is used to obtain the vector encoding of the text. Then, input it into the entity boundary prediction layer to label the entity head and tail, and input the BiLSTM-CRF layer to predict the attribute label. The experimental results show that the proposed method can effectively perform attribute extraction in the military equipment domain. The *F1* value of the model reaches 77% on the constructed attribute extraction dataset, which outperforms the current state-of-art model.

1. Introduction

With the continuous development of Internet technology, data from all walks of life is growing rapidly. Organizing these data through knowledge graph technology can effectively improve data utilization efficiency. In the military field, the construction of knowledge graph is not only conducive for the military commanders to quickly and deeply understand certain military equipment but also can be combined with knowledge map and intelligent system for rapid intelligent decision-making assistance [1].

Attribute extraction is an important step in knowledge graph construction, which refers to extracting the attribute name and attribute value of entities from text data. Facing a large amount of text data in the military field, extracting attribute data automatically is one of the keys to study the construction of a military knowledge graph. The traditional

attribute extraction methods are divided into rule-based methods and machine learning-based methods. Zhai and Qiu [2] proposed a rule-based knowledge meta-attribute extraction method based on phrase structure trees. The rule-based method needs to set rules manually according to the data characteristics, so the migration of the method is poor. Jakob and Gurevych [3] fused multiple features and used conditional random fields [4] to extract attributes. However, machine learning methods require a large amount of labelled data and manual features. In recent years, deep learning methods have also been gradually applied to attribute extraction. Toh and Su [5] used a bidirectional recurrent neural network BRNN combined with a conditional random field for attribute value extraction. Cheng et al. [6] used a bidirectional long short-term memory network BiLSTM combined with a gated dynamic attention mechanism for attribute extraction. However, attribute extraction

based on deep learning methods also requires a large amount of annotated data. In the field of weaponry, there is a lack of corresponding annotated datasets. Manual annotation is not only time-consuming but also the level of the annotator will largely affect the quality of the annotated corpus [7]. Through investigation, we found that Baidu Encyclopedia currently contains a large number of weapon and equipment entries. There are a large number of semi-structured and unstructured data in the encyclopedia pages, which contain rich information of entity attributes. We propose a new way of attribute data annotation based on the characteristics of the encyclopedia pages. We annotate the unstructured text data by distant supervision based on the InfoBox data of the encyclopedia pages. At the same time, we convert the attribute extraction task into a sequence annotation task and use the RoBERTa-BiLSTM-CRF-SEL method for attribute data extraction.

In summary, the contribution points of this paper can be divided into the following three points.

- (1) A new way of data annotation is proposed for the characteristics of encyclopedia data. In the annotation process, the subjective is fixed according to the name of the encyclopedia page, and then, its attributes and attribute values are annotated
- (2) Based on Baidu Encyclopedia data, the military domain attribute extraction dataset is automatically constructed by using the idea of distant supervision
- (3) RoBERTa-BiLSTM-CRF-SEL is designed for automatic attribute extraction in the field of weapons. The method obtains entity boundary features through the entity boundary prediction layer. The loss of boundary prediction layer and the loss of attribute prediction layer are weighted and summed as the loss value of the model. In this way, the model entity recognition effect is improved. On the military equipment attribute extraction dataset, the *F1* of the proposed method reaches 0.77, which is better than other existing methods

2. Related Work

Attribute extraction methods can be mainly classified into rule-based methods, machine learning-based methods, and deep learning-based methods. The rule-based approach needs to formulate rules manually for specific situations. This method is simple and usually oriented to specific domains. Although the method has a high accuracy rate, it has a small scope of application and is difficult to migrate to other domains. The method based on machine learning is more flexible, but it needs the support of artificial features and large-scale datasets. The method based on deep learning can automatically mine hidden features between texts through a neural network model, but it also requires large amounts of labelled data for model training and optimization.

In the early studies of attribute extraction, scholars mainly formulated a series of rules to extract attributes. Hu

and Liu [8] extracted commodity attributes from customer reviews by frequent itemset feature extraction. Li et al. [9] presented an automatic method to obtain encyclopedia character attributes, and the speech tagging of each attribute value was used to locate the encyclopedia free text. The rules were discovered by statistical method, and the character attribute information was obtained from encyclopedia text according to rules matching. Yu et al. [10] proposed an approach of extracting maritime information and converting unstructured text into structural data. Ding et al. [11] formed nine types of description rules for attribute extraction by manually constructing rules. They analyzed the quantitative relationship and emotional information of attribute description and finally designed and implemented the academic concept attribute extraction system. Qiao et al. [12] suggested a rule-based character information extraction algorithm. Based on the rules, they researched and developed a character information extraction system and finally realized the automatic extraction of semistructured character attribute information. Kang et al. [1] offered an unsupervised attribute triplet extraction method for the military equipment domain. According to the distribution law of attribute triples in sentences, this method adopts an attribute indicator extraction algorithm based on frequent pattern mining and completes the extraction of attribute triples by setting extraction rules and filtering rules.

In a machine learning-based attribute extraction method, Zhang et al. [13] introduced word-level features in the CRF model and used domain dictionary knowledge as an aid for product attribute extraction. Xu et al. [14] introduced shallow syntactic information and heuristic location information and input them to CRF as features, which effectively improved the attribute extraction performance of the model. Gurumdimma et al. [15] presented the approach to extracting these events based on the dependency parse tree relations of the text and its part of speech (POS). The proposed method uses a machine-learning algorithm to predict events from a text. Cheng et al. [16] broke through the current method of a statistical operation mainly in the scope of sentences in the attribute attribution judgment. They proposed a method of character attribute extraction that is classified from text to sentence with the guidance of text knowledge. Kambhatla [17] employed maximum entropy models to combine diverse lexical, syntactic, and semantic features derived from the text. References [18–20] suggested a weakly supervised automatic extraction method that uses very little human participation to solve the problem of lack of training corpus. Zhang et al. [21] offered a novel composite kernel for relation extraction. The composite kernel consists of two individual kernels: an entity kernel that allows for entity-related features and a convolution parse tree kernel that models syntactic information of relation examples. Liu et al. [22] put a perceptron learning algorithm that fuses global and local features for attribute value extraction of unstructured text. The combination of features makes the model obtain better feature representation ability. Li et al. [23] constructed three kinds of semantic information through word attributes, word dependencies, and word embeddings of words. The three semantic information are

combined with the conditional random field model to realize the extraction of commodity attributes.

In recent years, attribute extraction methods based on deep learning have gradually become mainstream. Wang et al. [24] regarded attribute extraction as a text sequence labelling task. Input the word sequences and lexical sequences into a GRU network, and then, use CRF for sequence label prediction. Xu et al. [25] considered that there is a gap between the meaning of a word expression in general and specialized domains. Therefore, they input both word embeddings from the generic domain and word embeddings from the specialized domain into a convolutional neural network model. The model is used to decide which expression is more preferred to achieve the attribute extraction. For the low performance of slot filling method applied in Chinese entity-attribute extraction at present, He et al. [26] presented a distant supervision relation extraction method based on bidirectional long short-term memory neural network. Wei et al. [27] proposed an attribute extraction-oriented class-convolutional interactive attention mechanism. The target sentence was first input into a bidirectional recurrent neural network to obtain the implicit expression of each word and then underwent class-convolution interactive attention. The force mechanism performed representation learning. To solve the problem that traditional information extraction methods have poor extraction results due to the existence of long and difficult sentences and the diversity of natural language expressions, Wu et al. [28] introduced text simplification as the pre-processing process of extraction. Among them, text reduction is modeled as a sequence-to-sequence (seq2seq) translation process and is implemented with the seq2seq-RNN model in the field of machine translation. Huang et al. [29] proposed a different method, which uses an independent graph based on a neural network as the input and is accompanied by two attention mechanisms to better capture indicative information. Cheng et al. [30] used the advantages of the CRF model to deal with the sequence labelling problem and realized the automatic extraction of journal keywords by integrating the part-of-speech information and the CRF model into the BiLSTM network. Luo et al. [31] proposed a new bidirectional dependency grammar tree to extract the dependency structure features of a given sentence and then combined the extracted grammar features with the semantic features extracted using BiLSTM and finally used CRF for attribute word annotation. Feng et al. [32] introduced an entity attribute value extraction method based on machine reading comprehension model and crowdsourcing verification due to the high noise characteristics of Internet corpus. The attribute extraction task is transformed into a reading comprehension task. Luo et al. [33] introduced a MLBiNet (multilayer bidirectional network) that integrates cross-sentence semantics and associated event information, thereby enhancing the discrimination of events mentioned within. Xi et al. [34] presented bidirectional entity level decoder (BERD) to gradually generate argument role sequences for each entity.

To address the problem of lack of annotation data in the military equipment domain, the attribute extraction dataset

in the military equipment domain is automatically constructed based on distant supervision. The attribute annotation sequence is decoded by RoBERTa model combined with BiLSTM-CRF model, and the entity boundary prediction layer is also added to improve the effect of entity recognition in this paper.

3. Attribute Extraction Methods Based on RoBERTa and Entity Boundary Prediction

The model proposed in this paper is mainly composed of text coding layer, entity boundary prediction layer, and BiLSTM-CRF attribute prediction layer. We first encode the input text through RoBERTa [35] to obtain its hidden layer state vector. Then, input them into the entity boundary prediction layer and the BiLSTM-CRF attribute prediction layer, respectively. At the entity boundary prediction layer, the 0/1 coding method is used to label the entity head and tail, respectively, and then, the `start_loss` and `end_loss` of the two sequence labels are calculated. In the BiLSTM-CRF attribute prediction layer, we take the output result of the entity boundary prediction layer as a feature and splice it with the text vector. Input the splicing results into BiLSTM-CRF to predict the text attribute tag. Next, calculate its loss value `att_loss`. Finally, in the model optimization, we consider the three-loss values together, weigh the summation, and achieve the overall optimization of the model by backpropagation. The model structure diagram is shown in Figure 1.

3.1. Text Encoding Layer. BERT is a pretrained language model proposed by Google in 2018. BERT uses the bidirectional transformer structure as the main framework of the algorithm, which can capture the bidirectional relations in utterances more thoroughly. BERT uses a self-supervised approach to train the model based on a massive corpus, which can learn a good feature representation for words. Therefore, BERT has achieved good results in several downstream tasks such as text classification and sequence annotation. RoBERTa model is an improved version based on the BERT model. Compared with BERT, RoBERTa has improved both the training data and training methods and pretrained the model more adequately.

In terms of training data, RoBERTa uses 160G training text, while BERT only uses 16G training text. RoBERTa also uses a new dataset CCNEWS and confirms that using more data for pretraining can further improve the performance of downstream tasks. At the same time, RoBERTa has increased the batch size. BERT uses 256 batch size. RoBERTa uses a larger batch size in the training process. Researchers have tried batch sizes ranging from 256 to 8000. Liu et al. found through experiments that the performance of certain downstream tasks can be slightly improved after removing the NSP (next sentence prediction, NSP) loss. Therefore, in the training method, RoBERTa deleted the NSP task. In addition, unlike the static masking mechanism of BERT, RoBERTa uses a dynamic masking mechanism to randomly generate a new mask pattern every time. BERT relies on random masks and predicted tokens. The original BERT implementation performs a mask during data

a temporal recurrent neural network, which can better capture the longer distance dependencies in the text. The LSTM model structure is shown in Figure 2.

There are three inputs to the LSTM, which are the hidden layer state vector h_{t-1} at the previous moment, the cell state C_{t-1} at the previous moment, and the input x_t at the current moment. Inside the LSTM, the retention and forgetting of information are decided by three gating mechanisms. The first is the forgetting gate, which is used to decide what information to forget from the cell state. The forgetting gate is used to read h_{t-1} and x_t and outputs data between 0 and 1 to decide which information in C_{t-1} to keep and which to discard, where 1 means fully retained, and 0 means all discarded. The input gate is used to decide which new information is added to the cell state, and the output gate decides which data in the cell state will be output. The calculation formulas of the LSTM model are shown in

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

$$i_t = \sigma(W_i[h_{t-1}, x_t] + b_i), \quad (2)$$

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o), \quad (3)$$

$$\tilde{C}_t = \tanh(W_C[h_{t-1}, x_t] + b_C), \quad (4)$$

$$C_t = f_t \cdot C_{t-1} + i_t \cdot \tilde{C}_t, \quad (5)$$

$$h_t = o_t * \tanh(C_t). \quad (6)$$

LSTM can only encode information in one direction. To effectively use the context information, we use a bidirectional LSTM structure for encoding.

By calculating the hidden layer vector output of the LSTM in both positive and negative directions and splicing them together, the hidden layer state vector of BiLSTM is finally obtained. The formulas are shown in

$$\vec{h}_t = \overrightarrow{\text{LSTM}}(\vec{h}_{t-1}, w_t), \quad (7)$$

$$\overleftarrow{h}_t = \overleftarrow{\text{LSTM}}(\overleftarrow{h}_{t-1}, w_t), \quad (8)$$

$$h_t = \text{concat}(\vec{h}_t, \overleftarrow{h}_t). \quad (9)$$

The conditional random field is a conditional probability distribution model of output $Y = (Y_1, Y_2, \dots, Y_n)$ given a set of input variables $X = (X_1, X_2, \dots, X_n)$. CRF is a serialization annotation algorithm, which can consider the dependencies between tags to obtain the globally optimal tag sequence.

For a set of label prediction sequence Y , its scoring formula is shown in

$$\text{score}(x, y) = \sum_{i=0}^n A_{y_i, y_{i+1}} + \sum_{i=1}^n P_{i, y_i}. \quad (10)$$

Among them, P is an $n \times m$ dimensional matrix, m represents the number of labels to be predicted, and P_{ij}

represents the possibility that input i is the label j . A is the transition matrix, and $A_{i,j}$ represents the probability of transition from label i to label j .

Therefore, for all possible prediction sequence sets Y_x of the input sequence X , the conditional probability is as shown

$$P(y | x) = \frac{e^{\text{score}(x, y)}}{\sum_{\tilde{y} \in Y_x} e^{\text{score}(x, \tilde{y})}}. \quad (11)$$

In training, we optimize the model by maximizing the log-likelihood probability of the correct output label in Equation (12). For prediction, we select the sequence with the highest score as the best prediction sequence, which is calculated as shown in Equation (12).

$$y^* = \arg \max_{\tilde{y} \in Y_x} \text{score}(x, \tilde{y}). \quad (12)$$

Take sentences in the dataset as an example, such as ‘‘On November 11-1989, the USS Abraham Lincoln was officially commissioned at Naval Station Norfolk and integrated into the American Atlantic Fleet,’’ ‘‘November 11-1989’’ would be marked as ‘‘B-FY,’’ ‘‘USS’’ would be marked as ‘‘B-ST,’’ ‘‘Abraham’’ would be marked as ‘‘I-ST,’’ ‘‘Lincoln’’ would be marked as ‘‘I-ST,’’ and ‘‘American’’ would be marked as ‘‘B-GJ’’ (please refer to Chapter 4 for label meaning).

3.4. Loss Value Calculation. In terms of loss value calculation, we take the weighted sum of entity boundary loss value and attribute identification loss value as the final loss value. The loss value is used to optimize the overall parameters of the model (as shown in Figure 3).

The loss value calculation formula is shown in Equation (13), where $\text{loss}_{\text{start}}$ and loss_{end} represent the loss values of entity head recognition and entity tail recognition, respectively, and $\text{loss}_{\text{attribute}}$ represents the loss value generated by the attribute sequence labelling. $\alpha, \beta, \gamma \in [0, 1]$ are hyperparameters that control the weighted summation of the three-loss values.

$$\text{loss} = \alpha \text{loss}_{\text{start}} + \beta \text{loss}_{\text{end}} + \gamma \text{loss}_{\text{attribute}}. \quad (13)$$

4. Experimental Results and Analysis

4.1. Acquisition of Military Equipment Attribute Data

4.1.1. Data Acquisition. The experimental data came from the Baidu Encyclopedia website (<https://baike.baidu.com/>), and the data acquisition process is shown in Figure 4. We cannot directly obtain military-related terms from Baidu Encyclopedia, because the website does not classify and index terms. The military channel of <http://globe.com/> has a summary display of various types of weapons and equipment. We get the names of various military equipment from the military channel of the World Wide Web. Then, we expand the rules, splice them with the links of encyclopedia entries, and finally, get the URL links of the required military equipment-related entries. After obtaining the links of military equipment entries in Baidu Encyclopedia, we analyzed

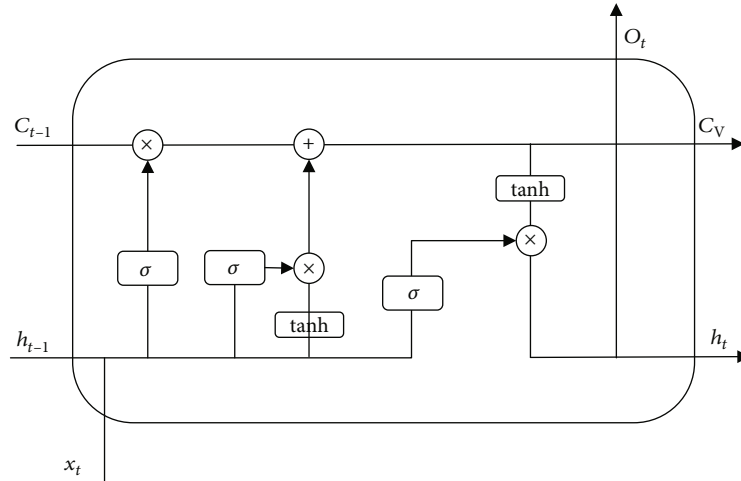


FIGURE 2: LSTM structure diagram.

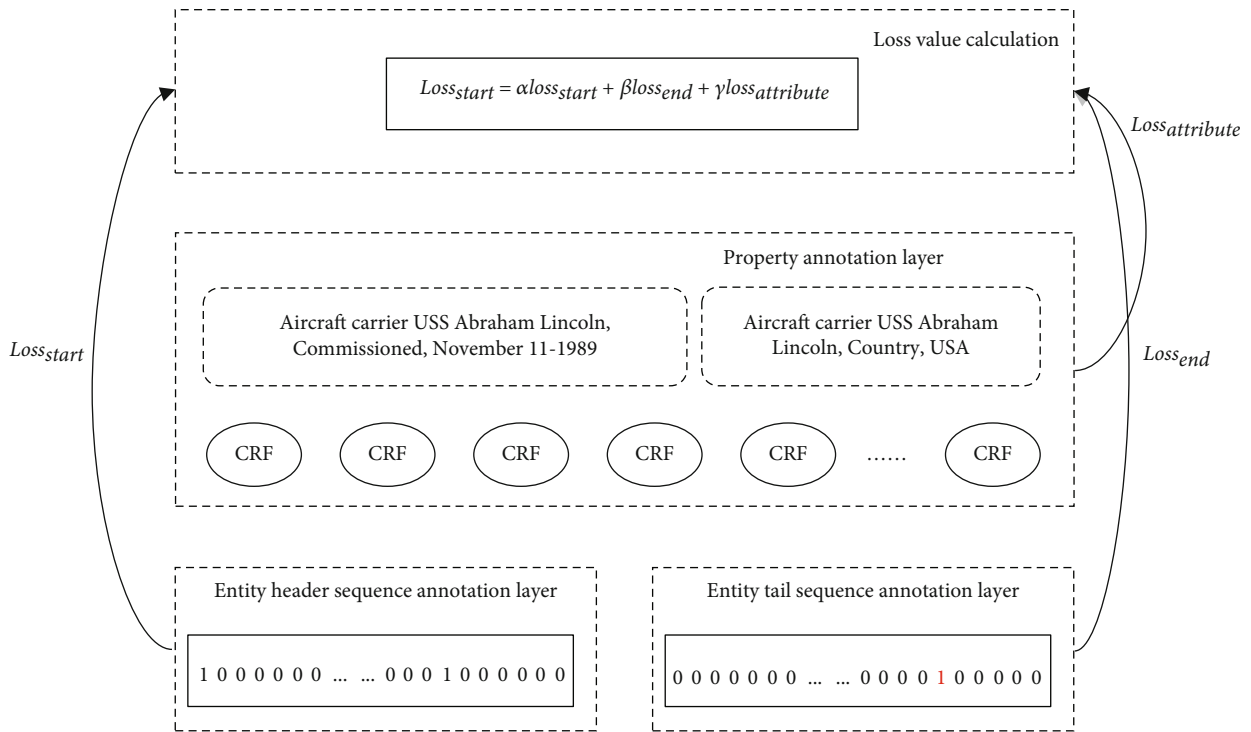


FIGURE 3: Calculation of loss value.

the encyclopedia entry pages and found that the entries mainly consist of entry names, information boxes containing attribute data, and a large amount of unstructured text. We used a crawler to collect the InfoBox data and text data in the Baidu Encyclopedia entry of weapons and equipment and finally collected 1757 encyclopedia data of military equipment.

4.1.2. Data Annotation. Data annotation by manual is not only time-consuming and laborious but also different annotators may have different annotation rules for the same piece of data. Therefore, automatic annotation of data has become the focus of current research. Encyclopedia word data

consists of two main parts, which are attribute data in the information frame and unstructured text description data. Taking the "Nimitz aircraft carrier" as an example, the entry information box of the aircraft carrier contains basic attributes such as "English Name," "Nation," "pretype/level," and "subtype/level." The text data is an introduction to the basic information of the "Nimitz aircraft carrier." Observing its text data, it can be seen that it contains textual expressions of the "English Name," "Nation," and other attribute values of the "Nimitz aircraft carrier."

For this data feature, the data annotation in this paper is based on the distant supervision hypothesis [33]. The distant supervision hypothesis means that when there is a

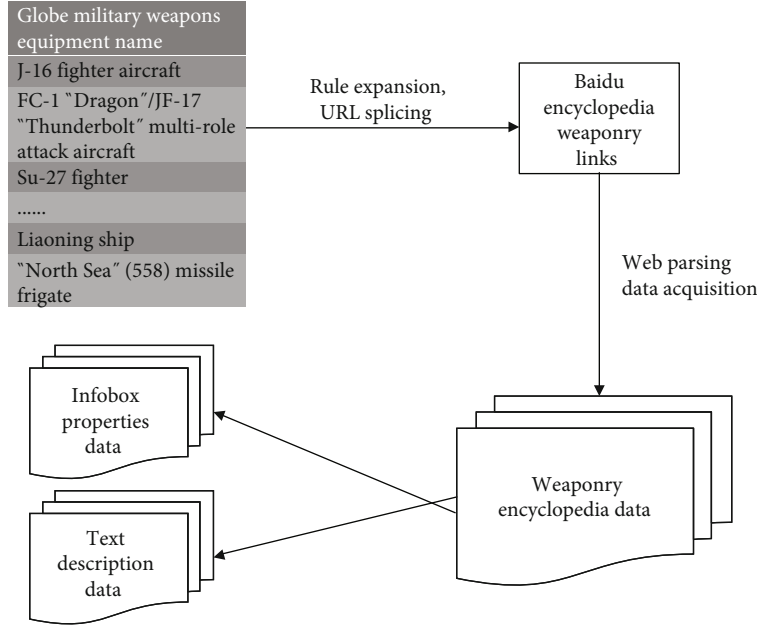


FIGURE 4: Flow chart of attribute data collection.

relationship between two entities, then all sentences containing the pair of entities are considered to express this relationship to some extent. Distant supervision is to provide labels for data with the help of external knowledge bases, to save the trouble of manual labelling [37]. Attributes can also be considered as a type of relationship, so the distant supervision assumption is applied to the annotation of attribute data. Taking the Nimitz aircraft carrier as an example, the information box in Figure 5 shows that the relationship between "Nimitz aircraft carrier" and the attribute "United States" is a "Nation" attribute. Then, based on the distant supervision assumption, all sentences containing "Nimitz aircraft carrier" and "United States" can be labelled with the "Nation" attribute, for example, the sentences "Nimitz Aircraft Carrier (CVN-68) is the first ship of the Nimitz-class aircraft carriers of the United States Navy" and "The Nimitz aircraft carrier started construction in June 1968. It was launched in May 1972 and delivered to the United States Navy in May 1975". Both of these sentences contain the words "Nimitz aircraft carrier" and "United States," and the triad (Nimitz aircraft carrier, nation, United States) can be considered to exist in these two sentences when labelling the data. Suppose a dataset $D = \{s_1, s_2, \dots, s_n\}$, where s_i represents sentence and is unstructured text. Train a model F such that $F(s_i; \theta) = [(e_i^t, e_j^t, e_k^t)]$, where θ represents model parameters, and e_i^t, e_j^t, e_k^t represent the T_{th} entity and its corresponding relationship. The idea of the distant supervision algorithm is to use knowledge base to align plain text for annotation and then perform supervised training.

However, Baidu Encyclopedia website is an open knowledge platform, and the editors of entries are not fixed. Therefore, there is a lack of standardization and unity in the naming of attributes, which leads to a variety of expressions of the same attribute. Since the data in the military field has a certain degree of confidentiality, the field itself has data sparsity. Different attribute expressions can lead to a variety of

data labels. If the labels are too scattered, the annotation data of each type of attribute will be small, which is difficult to obtain a good attribute extraction effect. To merge multiple attribute labels' expressions, we count the distribution of attribute names to select high-frequency words as attribute names. The attribute expressions present in the military equipment data of the encyclopedia website were merged by manual means, and a synonym table of military equipment attribute names was constructed. We merged and normalized the attribute expressions in the encyclopedia data through the synonym table. The synonym table of military equipment attribute names is shown in Table 1.

Normalize the attributes in the information box through the attribute name synonym table in Table 1 to obtain the attribute triplet set. Combined with the introduction text of triplet set and military equipment entries, data annotation is carried out through distant supervision. As shown in Figure 5, the text content in the encyclopedia web pages is expanded and described with the title of the entry as the center. The attribute triad is also composed with the title of the entry as the head entity. Based on the characteristics of this encyclopedia attribute data, a text sentence usually contains a primary entity and multiple attribute values corresponding to the entity. This chapter proposes a new data labelling method. Unlike the previous annotation form, we first annotate subjective according to the entry title and then annotate the attribute values corresponding to the primary entity separately (as shown in Table 2).

The relationship extraction dataset constructed by the distant supervision method often has noisy data (as shown in Table 3). For the relational triad [Obama, born in, United States], "Obama" and "United States" is a relationship of birth, and the distant supervision method is used for "Obama is the 44th president of the United States." The error occurs when the annotation is performed. For the

USS Nimitz (The first ship of the Nimitz-class aircraft carrier CVN-68)



The USS Nimitz (English: Nimitz Aircraft Carrier, port number: CVN-68 ^[1]), is the first ship of the US Navy's Nimitz-class aircraft carrier and a large nuclear-powered aircraft carrier in active service .

The ship is propelled by nuclear power , equipped with 4 elevators , 4 steam catapults and 4 arresting cables , and can eject a combat aircraft every 20 seconds. The model equipment in the carrier combat wing differs according to the nature of the combat mission. It can carry different purpose carrier aircraft to launch attacks on enemy aircraft, ships, submarines and land targets, and protect the maritime fleet. Kampfgruppe with its core usually consists of 4-6 cruisers , destroyers , submarines and supply ship constituted only.

Construction of the USS Nimitz started in June 1968, launched in May 1972, and delivered to the U.S. Navy in May 1975. The ship was first incorporated into the Atlantic Fleet , and its home port was Norfolk on the east coast of the United States. After being transferred to the Pacific Fleet, the station was changed to Everett Naval Base ^[2] .

Chinese name	USS Nimitz	Service time	May 3, 1975
Foreign name	USS Nimitz CVN-68 ^[1]	nation	U.S
Pre-type/level	Kitty Hawk class aircraft carrier and USS Enterprise ^[1]	Launch time	May 13, 1972
		Homeport	Everett Naval Base
Subtype/level	Ford class aircraft carrier ^[1]	Origin of ship...	Admiral Chester W. Nimitz
Development...	1961-1968	Full load drai...	101196 tons

FIGURE 5: Example of encyclopedia data.

triple [Obama, President, United States], the relationship between “Obama” and “United States” is “President,” and the annotation of “Obama was born in the United States” based on the distant supervision method will result in an annotation error. Since the relationship triad is composed of [entity, relationship, entity] and there may be multiple relationships between entities and entities, the distant supervision method often causes mislabelling problems when constructing relationship extraction datasets.

The attribute triad is composed of [entity, attribute, attribute value]. For military equipment data, the attribute values are usually some numerical information with unit agency names. Therefore, there is less possibility of distant supervision mislabelling problem when labelling the data for the entity and attribute values in military equipment data. For example, for the triad (Nimitz aircraft carrier, total load displacement, 101196 tons), the attribute value “101196 tons” is not a common entity, and it is challenging to generate other attribute relationships with the “Nimitz aircraft carrier.” The automatic annotation of the attribute data in the field of military equipment by the distant supervision method does not generate many mislabelling problems, and the correctness of the dataset can be guaranteed to a certain extent.

4.2. Dataset and Evaluation Index Description

(1) Description of dataset

We use the attribute extraction dataset constructed by the distant supervision method to verify method’s effective-

ness. 4291 attribute extraction corpus was constructed through the distant supervision method and the filtering of rules, including 3432 items in the training set, 429 items in the validation set, and 430 items in the test set. The details are shown in Table 4.

The Military Weaponry Dataset is a text corpus of weapons and equipment extracted from the military channel of the World Wide Web. Combined with the relevant knowledge of the encyclopedia website, the attribute extraction dataset is constructed using distant supervision and annotation. In the labeling process, we first determine the type of attribute contained in the sentence and then label the corresponding head entity and tail entity for each attribute [38]. The specific annotation example of the dataset is shown in Table 5. Among them, O stands for irrelevant words, B stands for the beginning of the entity, I stands for the middle part of the entity, ST stands for the subjective, GJ stands for the nation, QX stands for the pretype, ZL stands for weight, WW stands for foreign name, YZ stands for development time, FY stands for service time, TY stands for retirement time, CX stands for subtype, SF stands for first flight time, and DW stands for construction unit.

(2) Evaluation criteria

The experiment uses accuracy rate, recall rate, and F_1 value as evaluation indicators to evaluate the effectiveness of the method. The calculation method is shown in formula (14) to formula (16). Among them, $TP_{\text{attribute}}$ represents the number of attribute labels correctly identified in the forecast

TABLE 1: List of synonyms for attribute names of military equipment.

Name	Synonym			
Nation	Nation	Nation of origin	Nationality	Nation of manufacture
	Equipment nation	Place of birth	Nation of origin	Nation of construction
	Producing nation	Development nation	Manufacturing nation	Country
	R & D nation	Design nation	Origin	Build nation
	Affiliation nation			
Foreign name	Foreign name	Spanish name	German name	Japanese name
	Latin name	Russian name	English alias	Other translated names
	Japanese alias	French name	English scientific name	Korean name
	English name	Alias	Nickname	
Development/construction time	Development time	Manufacturing time	Development date	Development year
	Start to develop	Start development time	Design time	Construction time
	Build date	Construction year	Start time	
Service time	Year of service	Service	During service	Service date
Decommissioning time	Retirement time	End-time	Retirement date	Time to retire
	Retirement years	Retired		
Pretype/level	Pretype/level	Pretype	Predecessor	Former model
Subtype/level	Prestage			
Launch/first flight time	Subtype/level	Subtype	Secondary	
	Launch time	Launch date	Launch	First flight time
	First test flight	First flight	First flight date	Maiden flight
Development/construction unit	Development unit	R & D unit	Development company	Development organization
	Developer	Manufacturer	Production unit	Construction unit
Weight	Weight	Full load drainage	Full load displacement	Standard displacement
	Standard drainage	Displacement		

TABLE 2: Example of data annotation.

Example one	The [458 Sebari/subjective] commenced construction in [1990/construction time], was launched on [August 27, 1991/launch time] and commissioned on [August 4, 1992/commissioning time].
Example two	On [11 November 1989/service time], the [USS Abraham Lincoln/Subjective] was officially commissioned at Naval Station Norfolk as part of the [American/National] Atlantic Fleet.

TABLE 3: Example of distant supervision error labelling.

[Obama, born in, United States]	[Obama] was the 44th president of the [United States].	False
	[Obama] was born in the [United States].	True
[Obama, president, United States]	[Obama] was the 44th president of the [United States].	True
	[Obama] was born in the [United States].	False

TABLE 4: Description of attribute extraction dataset.

Dataset	Military equipment attribute extraction dataset
Training set	3432
Validation set	429
Test set	430

$$P = \frac{TP_{\text{attribute}}}{TP_{\text{attribute}} + FP_{\text{attribute}}}, \quad (14)$$

$$R = \frac{TP_{\text{attribute}}}{TP_{\text{attribute}} + FN_{\text{attribute}}}, \quad (15)$$

$$F_1 = \frac{2 \times PR}{P + R}. \quad (16)$$

output, $FP_{\text{attribute}}$ represents the number of attribute labels incorrectly identified in the forecast output, and $FN_{\text{attribute}}$ represents the number of unidentified attribute labels.

4.3. *Experimental Parameter Settings.* The experimental parameter settings are shown in Table 6. The batch size is

TABLE 5: Data annotation example.

T-44 tank (English: T-44 medium tank) is a medium tank developed by the Soviet Union on the basis of the T-34/85 tank in the mid-1940s.							
T-44	B-ST	Is	O	On	O	Mid-1940s	O
Tank	I-ST	a	O	The	O		
(O	Medium	O	Basis	O		
English	O	Tank	O	Of	O		
:	O	Developed	O	The	O		
T-44	B-WW	By	O	T-34/85	B-QX		
Medium	I-WW	The	O	Tank	I-QX		
Tank	I-WW	Soviet	B-GJ	In	O		
)	O	Union	I-GJ	The	O		

TABLE 6: Attribute extraction parameter settings.

Parameter	Parameter value
Batch size	8
Learning rate	$5e-5$
RoBERTa hidden layer size	768
LSTM hidden layer size	128
Sentence length	256
Training rounds	20
Dropout	0.5

set to 8, the learning rate is set to $5e-5$, the hidden layer size of RoBERTa is set to 768 according to the pretraining model, and the hidden layer size of LSTM is set to 128.

4.4. Description of Comparison Experiments and Analysis of Experimental Results. In this paper, the attribute extraction task is converted into a sequence annotation task. The current mainstream sequence annotation method is BiLSTM-CRF method. To obtain richer text vector information, we adopt RoBERTa for text encoding. At the same time, to be able to increase the entity recognition accuracy and improve the model extraction effect, we also add an entity boundary prediction layer. To verify the effectiveness of the methods, we design a total of five methods as the baseline models for attribute extraction from the perspective of the ablation experiment [39], as shown below. At the same time, we replace RoBERTa with BERT(base) for comparison experiments based on the following 5 methods. The details are shown in Tables 7 and 8.

The first experiment only uses the public pretraining model RoBERTa to label the attribute sequence. RoBERT uses longer time, larger batch size, and more data for training. This model has achieved good results. The $F1$ value of this experiment reached 0.719. The second experiment is RoBERTa+CRF model, which adds a conditional random field model to label the attribute sequence based on RoBERTa. Adding the CRF layer can add some constraints to the final predicted label to ensure that they are valid. These constraints can be automatically learned by the CRF layer from the training dataset during the training process. The third experiment uses RoBERTa+CRF+SEL model.

TABLE 7: Comparison experiment of attribute extraction (RoBERTa).

Model	Precision	Recall	$F1$
RoBERTa	0.662	0.786	0.719
RoBERTa+CRF	0.691	0.775	0.731
RoBERTa+CRF+SEL	0.745	0.780	0.762
RoBERTa+BiLSTM+CRF	0.721	0.751	0.735
RoBERTa+BiLSTM+CRF+SEL	0.740	0.803	0.770

TABLE 8: Attribute extraction comparison experiment (BERT).

Model	Precision	Recall	$F1$
BERT(base)	0.640	0.731	0.682
BERT(base)+CRF	0.670	0.752	0.709
BERT(base)+CRF+SEL	0.746	0.771	0.758
BERT(base)+BiLSTM+CRF	0.687	0.751	0.718
BERT(base)+BiLSTM+CRF+SEL	0.729	0.776	0.752

The entity boundary prediction layer is added on top of RoBERTa+CRF, splice it as features with hidden layer vector, and consider the loss value of entity boundary prediction layer and attribute labelling loss value in the process of model optimization. The accuracy of this experiment reached 0.745, which is the highest compared to the accuracy of the other four experiments. The fourth is RoBERTa+BiLSTM+CRF model, which uses RoBERTa to vectorize the input text. The traditional BiLSTM+CRF model is used to predict the text attribute sequence. The last model adds features of the entity boundary prediction layer based on RoBERTa+BiLSTM+CRF and comprehensively considers the entity boundary prediction layer loss and the BiLSTM-CRF attribute prediction layer loss when optimizing the model. The recall rate of this experiment reached 0.803, and the $F1$ value reached 0.77, which is better than the existing model.

It can be seen from Table 7 that the model effect has been improved to a certain extent after decoding with CRF. This may be that CRF can restrict the predicted label results and ensure that label "I" appears after label "B" with a high probability, which improves the effect of entity

TABLE 9: Comparison of extraction results of different attribute categories.

Attributes	Precision	Recall	<i>F1</i>	Number of training set samples
Subjective	0.80	0.87	0.83	3004
Nation	0.90	0.95	0.92	3032
Foreign name	0.63	0.64	0.64	716
Development time/construction time	0.50	0.48	0.49	311
Service time	0.39	0.51	0.44	423
Decommissioning time	0.50	0.14	0.22	31
Pretype/level	0.38	0.36	0.37	317
Subtype/level	0.26	0.29	0.27	105
Launch time/first flight time	0.49	0.60	0.54	293
Development unit/construction unit	0.54	0.68	0.60	556
Weight	0.70	0.84	0.76	60

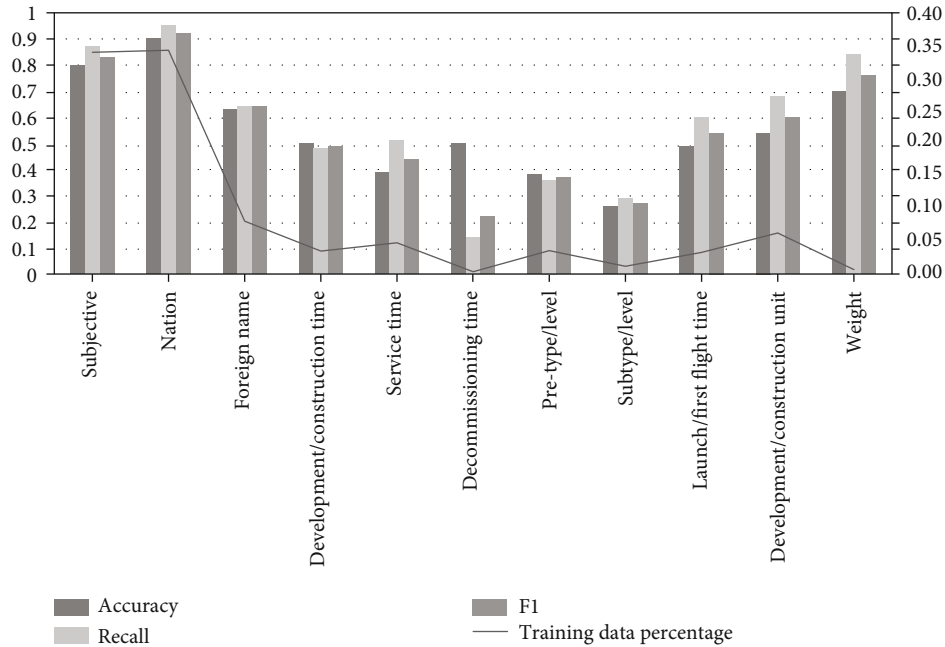


FIGURE 6: Comparison of the combination of the extraction results of different attribute categories for attribute extraction.

recognition to a certain extent. The *F1* value of the model is increased by 0.28 to 0.31 after adding the entity boundary prediction layer SEL. It may be that the increase of the entity boundary prediction layer helps to improve the effectiveness of the entity and attribute value boundary recognition. The overall effect has been improved. After adding BiLSTM for encoding, the model product has a slight improvement compared with the previous one, which may be due to the powerful encoding ability of RoBERTa has more fully obtained the contextual semantic information in the vector. So, the change is smaller after adding BiLSTM.

As shown in Table 8, the results of all five comparison experiments decrease after replacing RoBERTa with BERT(base), which indicates that RoBERTa is more effective in text vector representation and is more suitable for the mil-

itary equipment domain. Comparing BERT(base)+CRF+SEL with BERT(base)+BiLSTM+CRF+SEL, it can be seen that the effect of adding BiLSTM layer may not be greatly improved when the semantic information obtained by contextual encoding through BERT is richer. With the addition of BiLSTM, the recall of the model increase slightly, but the accuracy and *F1* both decrease to some extent.

To further study the recognition effect of RoBERTa+BiLSTM+CRF+SEL method on each different attribute, we test the accuracy, recall, and *F1* of the model on different attributes in the military equipment attribute dataset (as shown in Table 9). At the same time, to show the experimental results more clearly, a combined graph is drawn to show the model extraction effect and the number of samples in the training set. The histogram shows the accuracy, recall, and

F1. The number of samples in the training set is shown by the line graph.

As can be seen from Table 9, the F1 of the model varies considerably for different categories. For example, for the categories of “Subjective” and “Nation,” the F1 reach 0.83 and 0.92. For the categories of “Decommissioning Time,” “Pretype/Level,” and “Subtype/Level,” the F1 are around 0.2 to 0.3. It can be shown in Figure 6 that the recognition effect is better for the categories with more samples in the training set and worse for the categories with fewer samples in the training set. Therefore, it can be speculated that the large difference in the recognition effect of different categories may be caused by the uneven distribution of samples in the training set. For “Nation,” the description of “Nation” often appears in the sentence and every weapon information box in the encyclopedia entry. Therefore, the attribute can obtain more training corpus and a better extraction effect. As for “Decommissioning Time,” many pieces of equipment may be in active service, and there is less description of retirement. Therefore, the available corpus is far smaller than that of other attributes, and the model recognition effect is less satisfactory. The number of training samples for the attribute “Weight” is smaller, but the extraction effect is better. The reason may be that the attribute value of this attribute is usually numerical type, with obvious characteristics and easy to identify.

5. Conclusions

To address the problem of sparse data in the field of weaponry, a distant supervision approach is used to automatically annotate the weaponry data on Baidu Encyclopedia. The method reduces the working time of manual annotation and constructs a weaponry attribute extraction dataset. Based on the characteristics of the encyclopedia data, a data annotation approach is proposed. The annotation of the subjective is performed first, followed by the annotation of the attribute values corresponding to the subjective. For the constructed weapon and equipment dataset, we use the method of RoBERTa+BiLSTM+CRF+SEL to extract attributes and input text sentences into the pretrained attribute extraction model for attribute recognition. The method first uses RoBERTa to vectorize the text and then input it to the entity boundary prediction layer to obtain entity boundary features. This feature is spliced with the hidden layer state vector output by RoBERTa and input to the BiLSTM-CRF attribute prediction layer. Entity attribute triples are predicted through a sequence labeling method. The F1 value of this method is 0.763, which is better than other baseline models.

The attribute extraction of weapons and equipment is one of the important steps to construct the knowledge graph of weapons and equipment. We only consider Baidu Encyclopedia data in terms of data source, which has the problems of insufficient data, less partial attribute data, and unbalanced sample distribution. In future work, we should consider using more sources of data for distantly supervised annotation to expand the data scale. We should also try to solve the problem of uneven data distribution to improve the effect of model extraction.

Data Availability

The data used to support the findings of this study have not been made available because military data is classified.

Conflicts of Interest

The authors declare that they have no conflicts of interest.

Acknowledgments

The authors also acknowledge the Natural Science Foundation of Beijing under grant no. 4212020, National Natural Science Foundation of China under grant no. 62171043, Defense-Related Science and Technology Key Lab Fund Project under grant no. 6412006200404, Qin Xin Talents Cultivation Program of Beijing Information Science and Technology University under grant no. QXTCP B201908, and Research Planning of Beijing Municipal Commission of Education under grant no. KM202111232001.

References

- [1] K. Ruizhi, H. Wenning, K. Cheng, and Z. Donghui, “Attribute extraction for military equipment entity,” *Application Research of Computers*, vol. 33, no. 12, pp. 3721–3724, 2016.
- [2] Z. H. A. I. Jie and Q. I. U. Jiang-nan, “Research on the rule-based knowledge unit attributes extraction method,” *Information Science*, vol. 34, no. 4, pp. 43–47, 2016.
- [3] N. Jakob and I. Gurevych, “Extracting opinion targets in a single and cross-domain setting with conditional random fields,” in *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, pp. 1035–1045, MIT, Massachusetts, USA, 2010.
- [4] J. Lafferty, A. McCallum, and F. C. N. Pereira, “Conditional random fields: probabilistic models for segmenting and labeling sequence data,” in *Proceedings of the 18th International Conference on Machine Learning 2001 (ICML 2001)*, pp. 282–289, 2001.
- [5] Z. Toh and J. Su, “Nlangp at semeval-2016 task 5: improving aspect based sentiment analysis using neural network features,” in *Proceedings of the 10th international workshop on semantic evaluation (SemEval-2016)*, pp. 282–288, San Diego, California, 2016.
- [6] C. Meng, H. Yu, T. Jian, Z. Jiashuo, Z. Bowei, and Y. Jianmin, “Gated dynamic attention mechanism towards aspect extraction,” *Pattern Recognition and Artificial Intelligence*, vol. 32, no. 2, pp. 184–192, 2019.
- [7] M. Jin, Y. Yifan, and C. Wenliang, “Distant supervision for person attribute recognition,” *Journal of Chinese Information Processing*, vol. 34, no. 6, pp. 64–72, 2020.
- [8] M. Hu and B. Liu, “Mining and summarizing customer reviews,” in *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 168–177, Seattle, Washington, U.S.A., 2004.
- [9] H. Li, Y. Yang, H. Yin, and Z. Jia, “Rules-based character attributes extraction from Baidu Encyclopedia,” *Journal of Integration Technology*, vol. 3, 2013.
- [10] C. Yu, Z. Mao, and S. Gao, “An approach of extracting information for maritime unstructured text based on rules,” *Traffic Information and Safety*, vol. 35, no. 2, pp. 40–47, 2017.

- [11] D. Junjun, Y. Zheng, and H. Bolin, "Rule-based attribute extraction of academic concepts," *Information Theory and Practice*, vol. 34, no. 12, pp. 10–14, 2011.
- [12] L. Qiao, C. Li, Z. Zhong, J. Wang, and D. Liu, "Research on people's information extraction based on rules," *Journal of Nanjing Normal University(Natural Science Edition)*, vol. 35, no. 4, pp. 134–139, 2012.
- [13] S. Zhang, W. Jia, Y. Xia, Y. Meng, and H. Yu, "Opinion analysis of product reviews," in *2009 Sixth International Conference on Fuzzy Systems and Knowledge Discovery*, pp. 591–595, Tianjin, China, 2009.
- [14] X. Bing, Z. Tie-Jun, W. Shan-Yu, and Z. De-Quan, "Extraction of opinion targets based on shallow parsing features," *Acta Automatica Sinica*, vol. 37, no. 10, pp. 1241–1247, 2011.
- [15] N. Y. Gurumdimma, D. B. Bisandu, and E. Ojedayo, "Event extraction from textual data," *Journal of Computer Science and Its Application*, vol. 26, no. 1, 2020.
- [16] N. Cheng, Y. Zou, Y. Teng, and M. Hou, "On the method of personal attributes extraction based on textual knowledge and hierarchical classification," *Applied Linguistics*, vol. 1, pp. 125–134, 2019.
- [17] N. Kambhatla, "Combining lexical, syntactic, and semantic features with maximum entropy models for information extraction," in *Proceedings of the ACL Interactive Poster and Demonstration Sessions*, pp. 178–181, Barcelona Spain, 2004.
- [18] Y. Yufei, D. Qi, J. Zhen, and Y. Hongfeng, "Weakly supervised method for attribute relation extraction," *Journal of Computer Applications*, vol. 34, no. 1, pp. 64–68, 2014.
- [19] C. Liwei, F. Yansong, and Z. Dongyan, "Extracting relations from the web via weakly supervised learning," *Journal of Computer Research and Development*, vol. 50, no. 9, pp. 1825–1835, 2013.
- [20] J. Zhen, Y. Yang, and D.-k. He, "Attribute extraction of Chinese online encyclopedia based on weakly supervised learning," *Journal of University of Electronic Science and Technology of China*, vol. 43, no. 5, pp. 758–763, 2014.
- [21] M. Zhang, J. Zhang, J. Su, and G. Zhou, "A composite kernel to extract relations between entities with both flat and structured features," in *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pp. 825–832, Sydney, Australia, 2006.
- [22] L. Qian, D. Wu, L. Yue, X. Cheng, and P. Lin, "Extracting attribute values for named entities based on global feature," *Journal of Computer Research and Development*, vol. 53, no. 4, pp. 941–948, 2016.
- [23] L. Chengliang, Z. Zhongying, L. Chao, Q. Liang, and W. Yan, "Extracting product properties with dependency relationship embedding and conditional random field," *Data Analysis and Knowledge Discovery*, vol. 4, no. 5, pp. 54–65, 2020.
- [24] R. Wang, M. Xianru, and Q. Kong, "Entity-attribute extraction with GRU+CRF method," *Journal of Modern Information*, vol. 38, no. 10, pp. 57–64, 2018.
- [25] H. Xu, B. Liu, L. Shu et al., "Double embeddings and CNN-based sequence labelling for aspect extraction," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 592–598, Melbourne, Australia, 2018.
- [26] Z. He, Z. Zhou, L. Gan, J. Huang, and Y. Zeng, "Chinese entity attributes extraction based on bidirectional LSTM networks," *Journal of Computational Science and Engineering*, vol. 18, no. 1, p. 65, 2019.
- [27] W. Zhenkai, C. Meng, Z. Xiabing et al., "Convolutional interactive attention mechanism for aspect extraction," *Journal of Computer Research and Development*, vol. 57, no. 11, pp. 2456–2466, 2020.
- [28] W. Cheng, W. Chaokun, and W. Muxian, "Entity attributes extraction based on text simplification," *Computer Engineering and Applications*, vol. 56, 2020.
- [29] C. Bin, S. Shuicai, Y. Du, and X. Shibin, "Keyword extraction for journals based on part-of-speech and BiLSTM-CRF combined model," *Data Analysis and Knowledge Discovery*, vol. 5, no. 3, pp. 101–108, 2020.
- [30] H. Luo, T. Li, B. Liu, B. Wang, and H. Unger, "Improving aspect term extraction with bidirectional dependency tree representation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 7, pp. 1201–1212, 2019.
- [31] S. Feng, J. Liu, H. Jiang, and Y. Xiao, "Attribute value extraction method based on machine reading comprehension model and crowdsourcing verification," *Computer Engineering*, vol. 47, no. 5, pp. 97–103, 2021.
- [32] Y. Liu, M. Ott, N. Goyal et al., "Roberta: a robustly optimized bert pretraining approach," <http://arxiv.org/abs/1907.11692>.
- [33] D. Lou, Z. Liao, S. Deng, N. Zhang, and H. Chen, "MLBiNet: a cross-sentence collective event detection network," <http://arxiv.org/abs/2105.09458>.
- [34] X. Xi, W. Ye, S. Zhang, Q. Wang, H. Jiang, and W. Wu, "Capturing event argument interaction via a bi-directional entity-level recurrent decoder," <http://arxiv.org/abs/2107.00189>.
- [35] M. Mintz, S. Bills, R. Snow, and D. Jurafsky, "Distant supervision for relation extraction without labelled data," in *Proceedings of the 47th Annual Meeting of the ACL and the 4th IJCNLP of the AFNLP*, pp. 1003–1011, Suntec, Singapore, 2009.
- [36] Y. Zhu, W. Zhang, Y. Chen, and H. Gao, "A novel approach to workload prediction using attention-based LSTM encoder-decoder network in cloud environment," *EURASIP Journal on Wireless Communications and Networking*, vol. 2019, no. 1, p. 18, 2019.
- [37] J. Xiao, H. Xu, H. Gao, Y. Li, and Y. Li, "A weakly supervised semantic segmentation network by aggregating seed cues: the multi-object proposal generation perspective," *ACM Transactions on Multimedia Computing Communications and Applications*, vol. 17, no. 1s, pp. 1–19, 2021.
- [38] X. Ma, H. Gao, H. Xu, and M. Bian, "An IoT-based task scheduling optimization scheme considering the deadline and cost-aware scientific workflow for cloud computing," *EURASIP Journal on Wireless Communications and Networking*, vol. 2019, no. 1, p. 19, 2019.
- [39] Y. Huang, H. Xu, H. Gao, X. Ma, and W. Hussain, "SSUR: an approach to optimizing virtual machine allocation strategy based on user requirements for cloud data center," *IEEE Transactions on Green Communications and Networking*, vol. 5, no. 2, pp. 670–681, 2021.