

# Automatic Keyphrase Extraction from Scientific Chinese Medical Abstracts Based on Character-Level Sequence Labeling

Liangping Ding<sup>1,2</sup>, Zhixiong Zhang<sup>1,2,3†</sup>, Huan Liu<sup>1,2</sup>,  
Jie Li<sup>1,2</sup>, Gaihong Yu<sup>1,2</sup>

<sup>1</sup>National Science Library, Chinese Academy of Sciences, Beijing 100190, China

<sup>2</sup>University of Chinese academy of sciences, Beijing 100049, China

<sup>3</sup>Wuhan Library, Chinese Academy of Sciences, Wuhan 430071, China

Citation: Ding, L.P., Zhang, Z.X., Liu, H., Li, J., & Yu, G.H. (2021). Automatic keyphrase extraction from Scientific Chinese medical abstracts based on character-level sequence labeling. *Journal of Data and Information Science*, 6(3), 35–57. <https://doi.org/10.2478/jdis-2021-0013>

Received: Oct. 31, 2020

Revised: Dec. 27, 2020

Accepted: Jan. 15, 2021

## Abstract

**Purpose:** Automatic keyphrase extraction (AKE) is an important task for grasping the main points of the text. In this paper, we aim to combine the benefits of sequence labeling formulation and pretrained language model to propose an automatic keyphrase extraction model for Chinese scientific research.

**Design/methodology/approach:** We regard AKE from Chinese text as a character-level sequence labeling task to avoid segmentation errors of Chinese tokenizer and initialize our model with pretrained language model BERT, which was released by Google in 2018. We collect data from Chinese Science Citation Database and construct a large-scale dataset from medical domain, which contains 100,000 abstracts as training set, 6,000 abstracts as development set and 3,094 abstracts as test set. We use unsupervised keyphrase extraction methods including term frequency (TF), TF-IDF, TextRank and supervised machine learning methods including Conditional Random Field (CRF), Bidirectional Long Short Term Memory Network (BiLSTM), and BiLSTM-CRF as baselines. Experiments are designed to compare word-level and character-level sequence labeling approaches on supervised machine learning models and BERT-based models.

**Findings:** Compared with character-level BiLSTM-CRF, the best baseline model with F1 score of 50.16%, our character-level sequence labeling model based on BERT obtains F1 score of 59.80%, getting 9.64% absolute improvement.

**Research limitations:** We just consider automatic keyphrase extraction task rather than keyphrase generation task, so only keyphrases that are occurred in the given text can be extracted. In addition, our proposed dataset is not suitable for dealing with nested keyphrases.

**Practical implications:** We make our character-level IOB format dataset of Chinese Automatic Keyphrase Extraction from scientific Chinese medical abstracts (CAKE) publicly



† Corresponding author: Zhixiong Zhang (E-mail: zhangzhx@mail.las.ac.cn).

**Research Paper**

available for the benefits of research community, which is available at: <https://github.com/possible1402/Dataset-For-Chinese-Medical-Keyphrase-Extraction>.

**Originality/value:** By designing comparative experiments, our study demonstrates that character-level formulation is more suitable for Chinese automatic keyphrase extraction task under the general trend of pretrained language models. And our proposed dataset provides a unified method for model evaluation and can promote the development of Chinese automatic keyphrase extraction to some extent.

**Keywords** Automatic keyphrase extraction; Character-level sequence labeling; Pretrained language model; Scientific chinese medical abstracts

## 1 Introduction

Automatic keyphrase extraction (AKE) is a task to extract important and topical phrases from the body of a document (Turney, 2000), which is the basis of information retrieval (Jones & Staveley, 1999), text summarization (Zhang, Zincir-Heywood, & Milios, 2004), text categorization (Hulth & Megyesi, 2006), opinion mining (Berend, 2011), and document indexing (Frank et al., 1999). It can help us quickly go through large amounts of textual information to find out the main stating point of the text. Appropriate keyphrases can serve as a highly concise summarization of the text and are beneficial to retrieve text.

Classic keyphrase extraction algorithms usually contain two steps (Hasan & Ng, 2014). The first step is to generate candidate keyphrases, in which plenty of manually designed heuristics are combined to select potential candidate keyphrases. And the second step is to determine which of these candidate keyphrases are correct. One of the shared disadvantages in above-mentioned two-step approaches is that the model performance in second step is based on the quality of candidate keyphrases generated in the first step. So some researchers reformulate keyphrase extraction as a sequence labeling task and validate the effectiveness of this formulation.

Zhang et al. (2008) firstly reformulated keyphrase extraction as a sequence labeling task and constructed a CRF model to extract keyphrases from Chinese text, which skips the step of candidate keyphrase generation. They used 600 documents to train the model and designed lots of features manually. Moreover, they used word-level sequence labeling instead of character-level, tagging the words rather than characters. In Chinese, word is the minimal unit to express semantics. The advantage of word-level formulation is that we can model the relationship among words directly while the disadvantage is that it still depends on the word segmentation results of Chinese tokenizer.

By virtue of automatic extracting features, deep learning methods exceed machine learning methods and gradually become the mainstream in many natural language



processing (NLP) tasks. Transformer (Vaswani et al., 2017), an emerging model architecture for handling long-term dependencies, is a substitute to classic neural networks such as Long Short-Term Memory network. In 2018, Google released BERT (Devlin et al., 2019), which is a language model pretrained on large-scale unannotated text and used Transformer to capture deep semantic and syntactic features in text. In 2019, Sahrawat et al. (2019) regarded AKE as a sequence labeling task and applied lots of pretrained language models including BERT to English automatic keyphrase extraction task, showing the effectiveness of pretrained language model.

Compared to English keyphrase extraction, Chinese keyphrase extraction is facing two challenges: lacking of publicly available annotated dataset and relying on Chinese word segmentation tool. Firstly, supervised methods need ground-truth keyphrases of the text to train the model, while there are few Chinese publicly annotated keyphrase extraction datasets, which makes it difficult to do objective evaluation among different researches. Secondly, English tokens are split by white space while there is no delimiter among Chinese words.

To address the above-mentioned challenges, in this paper, we constructed a high quality dataset for Chinese automatic keyphrase extraction. We formulated keyphrase extraction from scientific Chinese medical abstracts as a character-level sequence labeling task which doesn't rely on Chinese tokenizer. And also we designed experiments to compare the model performance under word-level and character-level sequence labeling formulations, which has not been explored. In addition, for scientific Chinese medical abstracts, English words are interspersed with Chinese words, which increases the difficulty of data preprocessing. So we used Unicode Coding to distinguish English and Chinese, which regards each English word as the elementary unit and each Chinese character as the elementary unit.

Our key contributions are summarized as follows:

1. We regarded AKE from scientific Chinese medical abstracts as a character-level sequence labeling task and fine-tuned the parameters of BERT (Devlin et al., 2019) to make it adapt to our large-scale keyphrase extraction dataset. Our approach skips the step of candidate keyphrase extraction and is independent of Chinese tokenizer. And also we transferred the pretrained language model BERT to downstream Chinese AKE task without complicated manually-designed features.
2. We designed comparative experiments against word-level and character-level sequence labeling formulation for Chinese keyphrase extraction to verify the effectiveness of character-level formulation, especially under the general trend of pretrained language model. The comparative experiments were conducted



on machine learning baseline models and BERT-based model. We found that the performance of character-level formulation is comparable to word-level formulation or even higher for traditional machine learning algorithms while has overwhelming advantages for pretrained language model.

3. We processed data from Chinese Science Citation Database and constructed a large-scale character-level dataset for AKE from scientific Chinese medical abstracts. The dataset was labeled using Inside–Outside–Beginning tagging scheme (IOB format), which is a common tagging format in chunking tasks such as named entity recognition task. Our proposed dataset contains 100,000 abstracts in training set, 6,000 abstracts in development set and 3,094 abstracts in test set. We made our processed large-scale dataset (CAKE) publicly available for the benefits of the research community.

## 2 Related work

### 2.1 Automatic keyphrase extraction

Automatic keyphrase extraction has received lots of attention for more than 20 years. Over this time, existing classic methods usually contain two steps: generating candidate keyphrases and determining which of these candidate keyphrases match ground-truth keyphrases. In the first step, candidate keyphrases generation relies on some heuristics such as extracting n-grams that appear in external knowledge base (Grineva et al., 2009; Medelyan et al., 2009), extracting phrases that satisfy pre-defined lexical patterns (Barker & Cornacchia, 2000; Hulth, 2003; Le et al., 2016; Wang et al., 2016). The classic approaches in the second step can be divided into two categories: unsupervised approaches and supervised approaches.

Unsupervised approaches can be divided into four types: statistics-based approaches (Campos et al., 2018), graph-based approaches (Grineva et al., 2009; Mihalcea & Tarau, 2004), embedding-based approaches (Liu et al., 2009, 2010), and language model-based approaches (Tomokiyo & Hurst, 2003). Graph-based methods are the most popular ones while statistics-based methods still hold the attention of the research community (Papagiannopoulou & Tsoumakas, 2019).

As for Statistics-based approaches, these approaches don't need any training corpus and they are based on statistical features of the given text such as word frequency (Luhn, 1957), TF\*IDF (Salton et al., 1975), PAT-tree (Chien, 1997), and word co-occurrences (Matsuo & Ishizuka, 2004). And it's suitable for one single document because no prior information is needed. In 1995, Cohen (1995) used N-gram statistical information to automatically index the document. It didn't use any stop list, stemmer or domain-specific external information, allowing for easy



application in any language or domain with slight modification. In 1997, Chien (1997) used PAT-tree and mutual information between words to extract Chinese keyphrases. In 2009, Carpena et al. (2009) considered word frequency and spatial distribution features that keywords are clustered whereas irrelevant words distribute randomly in text. These statistical approaches are usually easy to transfer to a new domain because no prior information is applied.

As for graph-based approaches, keyphrase extraction is a ranking problem substantially. The model scores each candidate for its likelihood of being a ground-truth keyphrase and returns top-ranked keyphrases by setting a threshold. There are lots of popular unsupervised learning algorithms for keyphrases extraction, such as TextRank (Mihalcea & Tarau, 2004), LexRank (Erkan & Radev, 2004), TopicRank (Bougouin et al., 2013), SGRank (Danesh et al., 2015), and SingleRank (Wan & Xiao, 2008).

As for supervised approaches, classic keyphrase extraction is formulated as a binary classification problem (Frank et al., 1999; Turney, 2002) to determine whether the potential candidate keyphrases match ground-truth keyphrases for the text or not. Traditional machine learning algorithms such as Naïve Bayes (Witten et al., 2005), maximum entropy (Li et al., 2004), decision trees (Turney, 2000), SVM (Zhao et al., 2011), bagging (Hulth, 2003), and boosting (Hulth et al., 2001) rely heavily on complicated manually-designed features which can be broadly divided into two categories: within collection features and external resource-bases features (Hasan & Ng, 2014). Within collection features use textual features within training data and can be further divided into statistical features such as term frequency (Hulth, 2003), TF\*IDF (Salton & Buckley, 1988), syntactic features such as some linguistic patterns (Kim & Kan, 2009), and structural features such as location that keyphrases occur in (Wang et al., 2016). External resource-based features consist of lexical knowledge bases such as Wikipedia (Grineva et al., 2009; Medelyan et al., 2009), document citations (Caragea et al., 2014), hyperlinks (Kelleher & Luz, 2005). These methods have some weaknesses. The prediction for each candidate keyphrase is independent to that of others, which means that the model can't capture the connection among keyphrases.

These two-step keyphrase extraction approaches have some drawbacks. Firstly, error propagation. The candidate keyphrases generation errors occurring in the first step will be passed to the second step and influence the performance of the downstream methods. Secondly, the model performance relies heavily on some heuristic settings such as threshold, external resources (Wikipedia, domain ontology, lexicon dictionary etc.), and filtration patterns of POS tags, which make it difficult to transfer to a new domain. Thirdly, it's not able to find an optimal N value (number



**Research Paper**

of keyphrases to extract for the text) based on article contents so it is usually set to a fixed parameter which results in keyphrase extraction performance varying with the value for  $N$ . Fourthly, the number of keyphrases is same among text, ignoring the physical truth and bringing lots of redundant keyphrases or losing lots of important keyphrases. Finally, in the second step, the model just analyzes the semantic and syntactic properties of candidate keyphrases separately while losing the meaning of the whole text.

Zhang et al. (2008) first reformulated keyphrase extraction to a sequence labeling task, and utilized user-defined tagging scheme to annotate each word in Chinese text and indicated its chunk belonging. And they used Conditional Random Field model, which shows great performance in sequence labeling task. They designed lots of manually-designed features such as POS tagging, TF\*IDF, and other location features. Li et al. (2013) also used word-level sequence labeling model to extract keyphrases in automotive field for Chinese text. Casting keyphrase extraction as a sequence labeling task bypasses the step of candidate keyphrases generation and provides a unified method for automatic keyphrase extraction. Moreover, in sequence labeling, keyphrases are correlated to each other instead of being independent units.

Supervised machine learning methods require precise feature engineering and they rely heavily on manually-designed features, which are time-consuming. Using deep learning method to automatically extract features has become the mainstream of many natural language processing tasks. There are some practices for English AKE. In 2016, Zhang et al. (2016) cast keyphrase extraction as a sequence labeling task and proposed a joint-layer recurrent neural network model to extract keyphrases from tweets, which doesn't need complicated feature engineering. In 2019, Sahrawat et al. (2019) constructed a BiLSTM-CRF model and used contextualized word embedding from pretrained language models to initialize the embedding layer. They evaluated model performance on three English benchmark datasets: Inspec (Hulth, 2003), SemEval-2010 (Kim et al., 2010), SemEval-2017 (Augenstein et al., 2017) and their model achieved state-of-the-art results on these three benchmark datasets.

Compared with English AKE, Chinese AKE is more complicated owing to the characteristic that there is no delimiter among Chinese words. So there is an additional step in most Chinese AKE models: using Chinese tokenizer to segment words. For traditional two-step keyphrase extraction models, generating Chinese candidate keyphrases needs to use Chinese tokenizer to segment words first. For Chinese AKE models based on sequence labeling, existing methods still use word-level tagging, restricted by the segmentation results of Chinese tokenizer.





## 2.2 Sequence labeling based on BERT

With the improvement of computer hardware and the increase of available data, deep learning based methods gradually occupy the dominant position in the field of natural language processing. Although deep neural networks can learn highly nonlinear features, they are prone to over-fitting without large amount of annotated data. And the objective functions of almost all deep learning architectures are highly non-convex function of the parameters, with the potential for many distinct local minima in the model parameter space (Erhan et al., 2010). Thus, how to initialize parameters has been a problem that puzzles researchers. The breakthrough came in 2006 with the algorithms for deep belief networks (Hinton, Osindero, & Teh, 2006) and stacked auto-encoders (Bengio et al., 2007), which are all based on a similar approach: greedy layer-wise unsupervised pre-training followed by supervised fine-tuning.

Compared with traditional supervised learning tasks that randomly initialize parameters then learn language representations directly from annotated text, pretraining-finetuning mode not only capture the syntactic and semantic features of tokens from large-scale unannotated text but also provide a good initial point for the downstream task, improving the generalization ability of the downstream supervised learning task.

Recently, BERT (Devlin et al., 2019), short for Bidirectional Encoder Representations from Transformers, which is a pretrained language model receiving widespread concern and is believed to be a milestone in NLP. BERT was pretrained on large-scale unlabeled data from BooksCorpus and English Wikipedia, containing more than 3.3 billion tokens in total. Using BERT to fine-tune the downstream supervised tasks breaks the record for 11 NLP tasks including sentence classification, named entity recognition, natural language inference etc., which proves the feasibility of pretraining-finetuning mode. Using pretrained language models (Howard & Ruder, 2018; Peters et al., 2018; Radford et al., 2018) has become a standard component of SOTA (state-of-the-art) model architecture in many natural language processing tasks.

Most previous works for sequence labeling are built upon different combinations of LSTM and CRF (Giorgi & Bader, 2018; Habibi et al., 2017; Wang et al., 2019). Since the release of BERT, some researchers have shown the effectiveness of applying BERT or BERT-based models to sequence labeling task such as named entity recognition task. The underlying architecture of BERT is Transformer, which performs strongly on various tasks depending on its capability to capture long term frequency. Lee et al. (2019) introduced BioBERT, which was pretrained on largescale biomedical corpora using the model architecture same with BERT. They tested



BioBERT on several publicly datasets for named entity recognition such as NCBI disease, BC5CDR. The results showed that BioBERT outperforms the state-of-the-art models on six of nine datasets.

In this paper, we combined the benefits of formulating keyphrase extraction from Chinese medical abstracts as a character-level sequence labeling task and the advantage of pretraining-finetuning mode, which can not only avoid errors occurring in Chinese tokenizer, but also extract features automatically rather than using complicated manually-designed features.

### 3 Methodology

### 3.1 Task definition

We cast Chinese keyphrase extraction as a character-level sequence labeling task and used IOB format as the input format to the model. This task can be formally stated as:

Let  $d = \{w_1, w_2, \dots, w_n\}$  be an input text, where  $w_i$  represents the  $i^{th}$  element. If the input text is mixed up with Chinese and English, the element is a character for Chinese and a word for English. Assign each  $w_i$  in the text one of the three class labels  $Y = \{K_B, K_I, K_O\}$ , where  $K_B$  denotes that  $w_i$  locates at the beginning of a keyphrase,  $K_I$  denotes that  $w_i$  locates in the inside or end of a keyphrase, and  $K_O$  denotes that  $w_i$  is not a part of all keyphrases. For example, there is a sentence “X 连锁先天性肾上腺发育不良患儿的临床及 NR0B1 基因突变分析 (Clinical and NR0B1 gene mutation analysis in children with X-linked congenital adrenal dysplasia)” and the keyphrases in this sentence are “X 连锁先天性肾上腺发育不良 (X-linked congenital adrenal dysplasia)” and “NR0B1 基因 (NR0B1 gene).”

After IOB format transformation, the character-level tagging result of this sentence is shown in Figure 1. As we can see, we split the sentence according to the language which regarded each English word as the elementary unit and each Chinese character as the elementary unit. This character-level formulation avoids errors of Chinese tokenizer, which has been a troublesome problem in Chinese keyphrase extraction.

X	linked	congenital	adrenal	dysplasia	children	clinical	and	nr0b1	gene	mutation	analysis														
X	连	锁	先	天	性	肾	上	腺	发	育	不	良	患	儿	的	临	床	及	nr0b1	基	因	突	变	分	析
B	I	I	I	I	I	I	I	I	I	I	I	I	O	O	O	O	O	O	B	I	I	O	O	O	O

Figure 1. An example of character-level sequence labeling.

### 3.2 Keyphrase extraction evaluation measures

Although there is a suit of evaluation measures for sequence labeling task, in automatic keyphrase extraction, what we really care about is whether we can extract



correct keyphrases of the provided text. So we used precision, recall, and F1-score based on actual matching keyphrases against the ground-truth keyphrases for evaluation as used by previous studies (Kim et al., 2010).

Traditionally, automatic keyphrase extraction system have been assessed by using the proportion of Top-N candidates that exactly match the ground-truth keyphrases. For keyphrase extraction based on sequence labeling, there is no need for N value and so we just used the keyphrases predicted by the model to evaluate the AKE performance. But we needed to firstly recognize the keyphrases from IOB format before evaluation. We concatenated characters between label “B” and the last adjacent label “I” behind label “B” as predicted keyphrase.

We denoted the total number of predicted keyphrases as  $r$ , the number of predicted keyphrases matching with ground-truth keyphrases as  $c$ , the number of ground-truth keyphrases as  $s$ . The evaluation measures were defined as follows:

$$\text{Precision: } P = \frac{c}{r}$$

$$\text{Recall: } R = \frac{c}{s}$$

$$\text{F1-score: } F = \frac{2 \times P \times R}{P + R}$$

### 3.3 Dataset construction

We collected data from Chinese Science Citation Database, which is a database contains more than 1,000 kinds of excellent journals published in mathematics, physics, chemistry, biology, medicine, and health etc. We set some constraints to restrict data to Chinese medical data as well as no incomplete and duplicated records included to ensure the quality of data. The constraints were listed as follows:

- (1) According to Chinese Library Classification (CLC), the CLC code of medical data starts with the capital letter “R”. So we restricted data to records that the metadata field of CLC code starts with the capital letter “R”.
- (2) The metadata field of language was set to Chinese.
- (3) The metadata fields of title, abstract, and keyphrases were not null. Here, keyphrases refer to author-assigned keyphrases.

Statistics showed that there were 757,277 records meeting the above-mentioned constraints in total. The title and the abstract of each article were concatenated as the source input text. Furthermore, there are two types of keyphrases: extractive keyphrases and abstractive keyphrases. Extractive keyphrases refer to keyphrases



Research Paper

that are present in the source input text while abstractive keyphrases refer to keyphrases that are absent in the source input text. Because we formulated keyphrase extraction as a character-level sequence labeling task and can only extract keyphrases that are present in the source input text, we just considered the extractive keyphrases.

For a given text, we expected that all author-assigned keyphrases are extractive keyphrases, so we can annotate as many extractive keyphrases as possible. To achieve that, we firstly matched each author-assigned keyphrase with the given text to see if all author-assigned keyphrases can be found in the text. Then we limited our dataset to records that all author-assigned keyphrases are extractive keyphrases. After filtration, there were 169,094 records in total. We aim to construct a large-scale dataset for our deep neural network model because although deep neural networks can learn highly non-linear features, they are prone to over-fitting compared with traditional machine learning methods.

We chose 100,000 records as our training set, 6,000 records as our development set and 3,094 records as our test set. Training set was used for training the keyphrase extraction model. Development set was used in the training process to monitor the generalization error of the model and to tune hyper-parameters. Test set was used to test the performance of the model. Note that there was no overlap among data sets. Next, we processed these three data sets using IOB format to make them suitable for modeling sequence labeling task.

In this paper, we are going to compare word-level and character-level formulation for Chinese keyphrase extraction. So we constructed datasets for character-level and word-level sequence labeling separately. For the generation of word-level dataset, we used Chinese tokenizer Jieba<sup>①</sup> to segment words. And the tagging process was almost the same with that of character-level dataset construction except that we tagged the words rather than characters. An example of word-level sequence labeling is shown in Figure 2.

X	linked	congenital	adrenal	dysplasia	children	clinical	and	nr0b1	gene mutation	analysis	
X	连锁	先天性	肾上腺	发育不良	患儿	的	临床	及	nr0b1	基因突变	分析
B	I	I	I	I	O	O	O	O	B	I	O

Figure 2. An example of word-level sequence labeling.

- For character-level IOB format generation, we did some preprocessing steps:
- (1) Using Unicode Coding to distinguish Chinese and English. To address the problem that English words and Chinese words are mixed together in Chinese medical abstracts, we used Unicode Coding to distinguish English and



<sup>①</sup> <https://github.com/fxsjy/jieba>

Chinese. Our proposed data sets can greatly deal with the split of English words and Chinese characters, in which English word and Chinese character is the minimal unit respectively.

- (2) Converting from all half width to full half width. Punctuations in Chinese medical text include two format: full width and half width. Authors may neglect the format of punctuations, which causes the problem that keyphrases can't match with the abstract. For example, the authors might provide the keyphrase “er:yag 激光” (er:yag laser), but they used “er:yag 激光” (er:yag laser) in the abstract in which the colon was in full width format. So we transformed all half width punctuations to full width punctuations except full stop.
- (3) Dealing with special characters. There are lots of special characters in scientific Chinese medical abstracts and sometimes there are space characters next to these special characters while sometimes not. To unify the format, we dropped all space characters next to special characters.
- (4) Lowercase. We transformed all English words to their lowercase format.

After preprocessing, we did the tagging process, in which we matched keyphrases with the source input text to find the locations of keyphrases present in the text and tagged the characters within the locations with either label “B” or label “I” and characters not within the locations with label “O”. For the first character in the keyphrase, tag it with label “B” and for the characters other than the first character in the keyphrase, tag them with label “I”.

	X	linked			congenital			adrenal			dysplasia			children	
Text:	X	连	锁	先	天	性	肾	上	腺	发	育	不	良	患	儿
Location:	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
Tag:	B	I	I	I	I	I	I	I	I	I	I	I	I	O	O

Figure 3. An example of character-level iob format generation.

Figure 3 is an example of character-level IOB format generation. In this example, the keyphrase is “X 连锁先天性肾上腺发育不良” (X-linked congenital adrenal dysplasia). We matched the keyphrase and returned the location between 0 and 12. So we tagged the character in location 0 with label “B” and the characters located between 1 and 12 with label “I”. Other characters not within the location were tagged with label “O”.

Note that there were two special occasions in our tagging process and we applied some tricks:



Research Paper

(1) Given two author-assigned keyphrases of the input text, if there is a containment relationship between the location span of two keyphrases, we use Maximum Matching Rule to tag the longest keyphrase. For example:

**Text:** “穴位注射罗哌卡因分娩镇痛对产妇产程的影响” (Effect of acupoint injection of ropivacaine labor analgesia on maternal labor)

This text has two author-assigned keyphrases: “分娩” (Childbirth) and “分娩镇痛” (Labor analgesia). The location span of “分娩” (Childbirth) is between 8 and 9 while the location span of “分娩镇痛” (Labor analgesia) is between 8 and 11. So we tagged the characters within the longest keyphrase “分娩镇痛” (Labor analgesia) with label “B” or “I”.

(2) If the first few characters of a keyphrase is equal to the last few characters of the other keyphrase and this keyphrase appears after the other keyphrase in a given text, we will concatenate these two keyphrases by their common characters. For example:

**Text:** “术中经食管超声心动图对心脏瓣膜置换术后即刻人工瓣膜功能异常的诊断价值” (Diagnostic value of intraoperative transesophageal echocardiography for abnormal prosthetic valve function in the immediate postoperative period after heart valve replacement)

This text has two author-assigned keyphrases: “人工瓣膜” (Prosthetic Valves) and “瓣膜功能异常” (Abnormal valve function). These two keyphrases share common characters “瓣膜” (Valves) and appear next to each other in the text. Then we will tag the keyphrase “人工瓣膜功能异常” (Abnormal prosthetic valve function) instead of “人工瓣膜” (Prosthetic Valves) or “瓣膜功能异常” (Abnormal valve function). This step determines that our dataset is suitable for flat keyphrase extraction rather than nested keyphrase extraction, which means that each character will be assigned only one label.

To examine the quality of our data sets, we counted the number of recognized keyphrases, the number of correct recognized keyphrases, and the number of ground-truth keyphrases in our generated data sets. And we used evaluation measures mentioned in section 3.2 to see the IOB generation performance. The IOB generation results for character-level and word-level are summarized in Table 1 and Table 2 separately.

Table 1. Character-level IOB generation results on data sets.

Data Set	P	R	F	Number of Recognized keyphrases	Number of Correct Recognized Keyphrases	Number of Ground-truth Keyphrases
Training Set	99.18%	99.42%	99.3%	416,013	409,371	408,373
Development Set	99.13%	99.54%	99.34%	25,942	26,169	26,061
Test Set	99.15%	99.56%	99.36%	13,344	13,458	13,403



Table 2. Word-level IOB generation results on data sets.

Data Set	P	R	F	Number of Recognized keyphrases	Number of Correct Recognized Keyphrases	Number of Ground-truth Keyphrases
Training Set	91.15%	96.93%	93.96%	395,852	434,266	408,373
Development Set	91.35%	97.03%	94.11%	25,287	27,680	26,061
Test Set	90.99%	97.11%	93.95%	13,016	14,305	13,403

As we can see, the F1-score of each character-level generated data set was higher than the corresponding word-level generated data set for more than 5 percent. For character-level data sets, owing to the above-mentioned tricks that we applied to IOB generation, the evaluation measures don't reach to 100%. But the character-level IOB generation results on all three data sets still show that our data sets are of good quality. For word-level sequence labeling data sets, the segmentation error of the Chinese tokenizer is a critical reason that the evaluation measures are lower than that of character-level. Take the example mentioned in section 3.1 as an example, the word-level tagging result is shown in Table 2. There was one incorrect keyphrase “nr0b1 基因突变” (nr0b1 gene mutation) which was supposed to be “nr0b1 基因” (nr0b1 gene). Except for tagged incorrect keyphrases, there might be missing keyphrases because of segmentation error for word-level sequence labeling.

### 3.4 Model architecture

We initialized our sequence labeling keyphrase extraction model with pretrained BERT model. The architecture of BERT is based on a multi-layer bidirectional Transformers (Vaswani et al., 2017). Instead of the traditional left-to-right language modeling objective, BERT was pretrained on two tasks: predicting randomly masked tokens and predicting whether two sentences follow each other. Our sequence labeling keyphrase extraction model followed the same architecture as BERT and was optimized on scientific Chinese medical abstracts. We used a feed-forward neural network which acted as a linear classifier layer on top of the representations from the last layer of BERT to compute character level IOB probabilities. Our model architecture is shown in Figure 4.

For a given token, its input representation was constructed by summing the Wordpiece embedding (Wu et al., 2016), segment embedding, and position embedding. The first token of each sequence was always the special token [CLS]. The segment embedding is useful in sentence pairs task such as question answering to differentiate sentence. Sentence pairs were separated by a special token [SEP] and sentence **A** embedding was added to each token in the first sentence while sentence **B** embedding was added to each token in the second sentence. Our task is a single sentence task, so we only used sentence **A** embedding. The position



## Research Paper

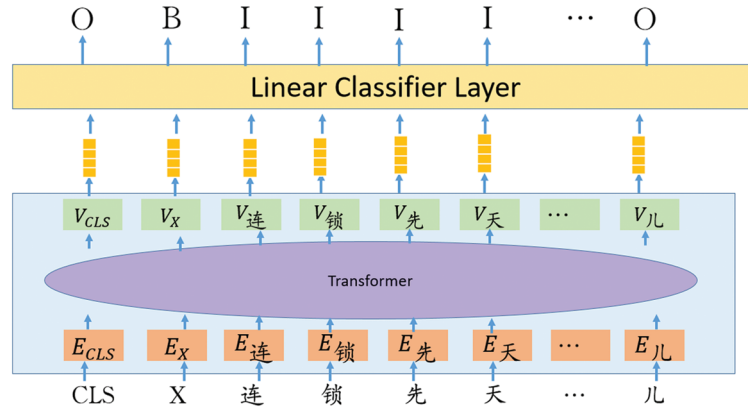


Figure 4. Character-level sequence labeling keyphrase extraction model architecture.

embedding was used to indicate the location of the token in the text and can only take the length lower than 512. A visual representation of our character-level input representations is given in Figure 5.

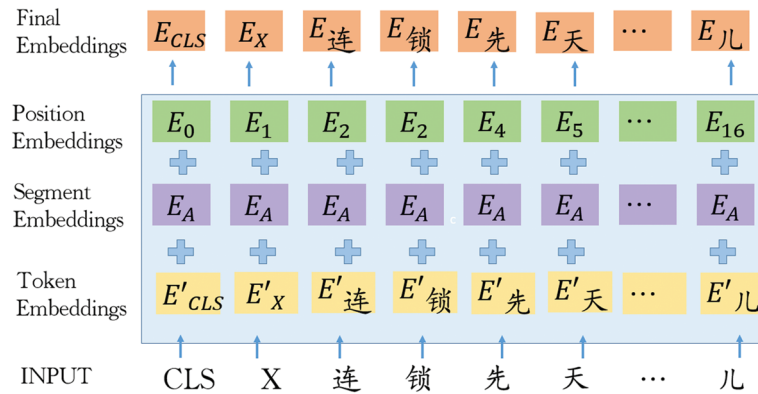


Figure 5. Input representations of character-level sequence labeling keyphrase extraction model.



In addition, BERT can only take the input with the maximum length of 512. Owing to this limitation, some source input text will be truncated, causing the problem that the model might predict some single character as keyphrases. In most cases, single Chinese character makes no sense. We found that some single Chinese characters are meaningful including chemical elements in The Periodic Table such as “氢” (hydrogen), “氦” (helium), organs such as “胃” (stomach), “脾” (spleen) and animals such as “鼠” (mouse), “鸡” (chicken). So we designed a user-defined lexicon to store meaningful Chinese characters for further filtration.

## 4 Experiments & results

### 4.1 Experimental design

In this paper, we firstly conducted unsupervised baseline experiments to demonstrate that traditional unsupervised two-step keyphrase extraction methods are sensitive to N value and the lexicon scale, which depends on precise manual settings. Then before we used sequence labeling formulation to Chinese keyphrase extraction task, we designed comparative experiments using word-level and character-level formulation on supervised machine learning baseline methods and BERT-based methods to verify the effectiveness of character-level. Finally, we compared the best unsupervised baseline model, the best character-level machine learning baseline model, and our character-level BERT-based sequence labeling keyphrase extraction model to prove the strength of sequence labeling formulation and pre-trained language model.

Regarding to unsupervised baselines, we used some traditional approaches including term frequency, TF\*IDF, and TextRank. As we know, the performance of traditional unsupervised approaches varies with the value for N (number of top ranked keyphrases), which is a parameter set manually. And traditional unsupervised Chinese keyphrase extraction relies on Chinese tokenizer to generate candidate keyphrases. Usually, user-defined lexicon will make a great difference to the results of Chinese word segmentation.

So we designed two groups of experiments using control variable method for unsupervised baselines according to N value and lexicon scale. Group 1 kept the same lexicon scale and compared the performance of baseline approaches at different N value of 3 and 5 to ensure the stability of the baseline approaches. Group 2 kept the same N value and compared the performance of baseline approaches when the lexicon scale for the Chinese tokenizer is different to test the transferability of baseline approaches. We set two kinds of lexicon scales, one using all ground-truth keyphrases in training set, development set, and test set as lexicon, the other just using ground-truth keyphrases in training set.

Regarding to supervised machine learning baselines, we cast keyphrase extraction as a sequence labeling task instead of a binary classification task and used CRF, BiLSTM, BiLSTM-CRF algorithms as machine learning baselines.

### 4.2 Experimental settings

As for unsupervised baseline approaches, we used Jieba for Chinese word segmentation. Before generating candidate keyphrases, we did some preprocessing steps, such as removing stop words and some special characters. We restricted candidate keyphrases within our user-defined lexicon and noun phrases.





**Research Paper**

Of the three machine learning baseline approaches, CRF was trained by regularized maximum likelihood estimation and we used Viterbi algorithm to find the optimal sequence of labels. BiLSTM and BiLSTM-CRF were trained with Stochastic Gradient Descent (SGD). The learning rate was set to  $5e-4$  and the model was trained for 15 epochs with early stopping. The hidden layers were set to 512 units and the embedding size was 768 in both models. In addition, the batch size was set to 64.

For our BERT-based keyphrase extraction model, due to system memory constraints, the batch size was set to 7 and we used SGD to optimize Cross Entropy Loss. The initial learning rate was set to  $5e-5$  and gradually decreased to  $5e-8$  as the training progresses and the model was trained for 3 epochs.

In this paper, we used F1-score to evaluate model performance, which is the weighted average of precision and recall, taking both precision and recall into account.

### 4.3 Unsupervised baseline experiments

As for traditional unsupervised baseline experiments, we conducted two groups of comparative experiments according to N value and lexicon scale as what we have mentioned in section 4.1.

For the group of N value experiments, we restricted the lexicon scale to whole lexicon, which contains author-assigned keyphrases in all the training set, development set, and test set as user-defined lexicon for Jieba word segmentation. Table 3 provides the results of N value comparison experiments of baseline approaches. Increasing the N value will improve the recall but lower the precision. We found that the F1-score of baseline approaches varied with the N value, but TF\*IDF achieved best performance among all baseline models no matter the N value. And when the N value was 3, the F1-score of TF\*IDF was 44.59%, which was higher than that when N value was 5.

Table 3. N-value comparative experiments of unsupervised baseline approaches.

Method	Top 3 Candidate Keyphrases			Top 5 Candidate Keyphrases		
	P	R	F	P	R	F
Term Frequency	47.66%	33.36%	39.24%	37.53%	43.78%	40.42%
TF*IDF	54.14%	37.90%	<b>44.59%</b>	40.37%	47.11%	43.48%
TextRank	43.13%	30.19%	35.52%	33.29%	38.84%	35.85%

For the group of lexicon scale experiments, we restricted N value to 3 to compare baseline approaches at different lexicon scales. Table 4 presents the results of lexicon scale comparative experiments of baseline approaches. As we can see, for all unsupervised baseline approaches, the performance of using lexicon that only



contains keyphrase in training set for Jieba word segmentation dropped at least 7% compared to that of using whole lexicon. The results showed that traditional unsupervised keyphrase extraction models for Chinese medical abstracts had poor transferability so when transferring to a new domain and no lexicon can be used, the keyphrase extraction performance might be poor.

Table 4. Lexicon scale comparative experiments of unsupervised approaches.

Method	P	R	F
Term Frequency (whole lexicon)	47.66%	33.36%	39.24%
Term Frequency (training set lexicon)	37.31%	26.11%	30.72%
TF*IDF (whole lexicon)	54.14%	37.90%	<b>44.59%</b>
TF*IDF (training set lexicon)	42.18%	29.53%	34.74%
TextRank (whole lexicon)	43.13%	30.19%	35.52%
TextRank (training set lexicon)	34.37%	24.06%	28.30%

#### 4.4 Word-level and character-level sequence labeling comparative experiments

We used word-level and character-level sequence labeling dataset separately to train and evaluate supervised machine learning baseline models and BERT-based models.

##### 4.4.1 Supervised machine learning baseline models

The F1-score evaluation metrics of word-level and character-level comparative experiments on machine learning baseline models are listed in Table 5. As we can see, word-level sequence labeling formulation was better than character-level sequence labeling formulation for CRF and BiLSTM algorithms while a little bit lower than character-level sequence labeling formulation for BiLSTM-CRF algorithms. The reason might be that BiLSTM-CRF is a more powerful model to capture the contextual relationship among characters to make up for the disadvantage that character-level formulation doesn't model the relationship among words directly.

Table 5. Word-level and character-level comparative experiments of supervised machine learning baselines.

Method	Word-Level	Character-Level
CRF	47.90%	46.37%
BiLSTM	44.35%	38.38%
BiLSTM-CRF	49.86%	<b>50.16%</b>

##### 4.4.2 BERT-based models

The precision, recall, and F1-score evaluation metrics of word-level and character-level sequence labeling comparative experiments on BERT-based models are listed



**Research Paper**

in Table 6. For word-level sequence labeling formulation, we just used the hidden state corresponding to the first character of the word as input to the linear classifier, which is the same approach used in (Devlin et al., 2019) for named entity recognition task. We found that the precision for word-level was extremely lower than character-level and the F1-score of word-level sequence labeling formulation was more than 20% lower than character-level formulation. Detailed analysis is conducted for this result. We assumed that Chinese BERT uses Wordpiece tokenizer which will tokenize each Chinese word into characters in the pretraining process. So Chinese BERT is character-level and has learned good semantic representation of Chinese characters through pretraining, which can maximize the advantages of the character-level sequence labeling formulation and avoid its shortcomings.

Table 6. Word-level and character-level comparative experiments of BERT-based models.

Metrics	Word-Level	Character-Level
P	26.88%	60.33%
R	54.93%	59.28%
F	36.10%	<b>59.80%</b>

#### 4.5 BERT-based character-level experiments

From the results of the above word-level and character-level comparative experiments, we decided to apply character-level formulation to our BERT-based Chinese keyphrase extraction model and the best character-level machine learning baseline model is BiLSTM-CRF. We compared the best unsupervised method TF\*IDF with our character-level sequence labeling BiLSTM-CRF model and found that sequence labeling formulation was beneficial for Chinese keyphrase extraction task. And we used character-level BiLSTM-CRF to compare with our character-level BERT-based model. The performance results are summarized in Table 7. Compared with BiLSTM-CRF, our BERT-based model achieved F1-score of 59.80%, exceeding that of baseline approach by 9.64%, which showed that the pretrained language model captured rich features that are useful for downstream keyphrase extraction task. And we removed single Chinese characters that were not in the user-defined lexicon. After removal, the keyphrase extraction performance of our adjusted model reached to 60.56%.

Table 7. Performance evaluation of keyphrase extraction.

Method	P	R	F
TF*IDF (Baseline)	54.14%	37.90%	44.59%
BiLSTM-CRF (Baseline)	42.55%	61.09%	50.16%
BERT-based Model (our model)	60.33%	59.28%	59.80%
Adjusted Model (our model)	61.95%	59.22%	<b>60.56%</b>



And we compared the predicted keyphrases with author-assigned ground-truth keyphrases and found that some predicted phrases were concatenation of author-assigned keyphrases. For example, there are two author-assigned keyphrases “卒中” (Stroke) and “抑郁” (Depression), while our model extracted keyphrases “卒中后抑郁” (Post-stroke depression). Another example, there are two author-assigned keyphrases “急性肠胃炎” (acute gastroenteritis) and “食源性疾病” (foodborne disease), while our model extracted keyphrases “食源性胃肠炎” (Foodborne gastroenteritis). These examples indicated that as though our model got the F1-score of 59.80%, our model can achieve good practical application performance. In addition, it also indicated that the calculation of evaluation measure is an issue we need to consider further. Using the proportion of predicted phrases that exactly match the ground-truth keyphrases to assess the model is actually not appropriate because there are some biases for author-assigned keyphrases and sometimes the phrases predicted by our model are also concise descriptions for the text.

## 5 Discussion & conclusion

In this paper, we formulated automatic keyphrase extraction as a character-level rather than word-level sequence labeling task and used pretrained language model BERT to fine-tune our keyphrase extraction model on scientific Chinese medical abstracts. Through our experimental work, we proved the benefits of this formulation with this architecture, which bypasses the step of Chinese tokenizer and leverages the power of pretrained language model. In addition, we also designed comparative experiments to verify that character-level formulation is more suitable for Chinese keyphrase extraction task under the trend of pretrained language model.

Our approach only dealt with keyphrase extraction rather than keyphrase generation, so it can just handle extractive keyphrases. In the future, we plan to build keyphrase generation model to extract keyphrases. And also we will explore the solutions to solve the limitation of BERT’s maximum sentence length to avoid being truncated. We expect some of the findings in this paper will provide valuable experiences for automatic keyphrase extraction and other NLP problems like document summarization, term extraction etc.

## Acknowledgments

This work is supported by the project “Research on Methods and Technologies of Scientific Researcher Entity Linking and Subject Indexing” (Grant No. G190091) from the National Science Library, Chinese Academy of Sciences and the project “Design and Research on a Next Generation of Open Knowledge Services System and Key Technologies” (2019XM55).



## Author Contributions

Liangping Ding (dingliangping@mail.las.ac.cn): performing the research; writing and revising the manuscript. Zhixiong Zhang (zhangzhx@mail.las.ac.cn): Instructing the research and instructing the structure of the manuscript. Huan Liu (liuhuan@mail.las.ac.cn): revising the manuscript and revising the structure of the manuscript. Jie Li (lijie201909@mail.las.ac.cn): revising the manuscript. Gaihong Yu (yugh@mail.las.ac.cn): revising the manuscript.

## References

- Augenstein, I., Das, M., Riedel, S., Vikraman, L., & McCallum, A. (2017). Semeval 2017 task 10: Scienceie-extracting keyphrases and relations from scientific publications. ArXiv Preprint ArXiv:1704.02853.
- Barker, K., & Cornacchia, N. (2000). Using noun phrase heads to extract document keyphrases. Conference of the Canadian Society for Computational Studies of Intelligence, 40–52.
- Bengio, Y., Lamblin, P., Popovici, D., & Larochelle, H. (2007). Greedy layer-wise training of deep networks. Advances in Neural Information Processing Systems, 153–160.
- Berend, G. (2011). Opinion expression mining by exploiting keyphrase extraction. In Proceedings of the 5th International Joint Conference on Natural Language Processing, 1162–1170, Chiang Mai, Thailand.
- Bougouin, A., Boudin, F., & Daille, B. (2013). Topicrank: Graph-based topic ranking for keyphrase extraction. In Proceedings of International Joint Conference on Natural Language, 543–551, Nagoya, Japan.
- Campos, R., Mangaravite, V., Pasquali, A., Jorge, A.M., Nunes, C., & Jatowt, A. (2018). A text feature based automatic keyword extraction method for single documents. European Conference on Information Retrieval, 684–691.
- Caragea, C., Bulgarov, F.A., Godea, A., & Gollapalli, S.D. (2014). Citation-enhanced keyphrase extraction from research papers: A supervised approach. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, 1435–1446.
- Carpene, P., Bernaola-Galván, P., Hackenberg, M., Coronado, A., & Oliver, J. (2009). Level statistics of words: Finding keywords in literary texts and symbolic sequences. Physical Review E, 79(3), 035102.
- Chien, L.F. (1997). PAT-tree-based keyword extraction for Chinese information retrieval. In Proceedings of the 20th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, 50–58.
- Cohen, J.D. (1995). Highlights: Language-and domain-independent automatic indexing terms for abstracting. Journal of the American Society for Information Science, 46(3), 162–174.
- Dai, A.M., & Le, Q.V. (2015). Semi-supervised sequence learning. Advances in Neural Information Processing Systems, 3079–3087.
- Danesh, S., Sumner, T., & Martin, J.H. (2015). Sgrank: Combining statistical and graphical methods to improve the state of the art in unsupervised keyphrase extraction. Proceedings of the Fourth Joint Conference on Lexical and Computational Semantics, 117–126.
- Devlin, J., Chang, M.W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. ArXiv:1810.04805 [Cs]. <http://arxiv.org/abs/1810.04805>



- Erhan, D., Bengio, Y., Courville, A., Manzagol, P.A., Vincent, P., & Bengio, S. (2010). Why does unsupervised pre-training help deep learning? *Journal of Machine Learning Research*, 11(Feb), 625–660.
- Erkan, G., & Radev, D.R. (2004). Lexrank: Graph-based lexical centrality as salience in text summarization. *Journal of Artificial Intelligence Research*, 22, 457–479.
- Frank, E., Paynter, G., Witten, I., Gutwin, C., & Nevill-Manning, C. (1999). Domain-Specific Keyphrase Extraction. In *Proceeding of 16th International Joint Conference on Artificial Intelligence*, Stockholm, Sweden, 668–673.
- Giorgi, J.M., & Bader, G.D. (2018). Transfer learning for biomedical named entity recognition with neural networks. *Bioinformatics*, 34(23), 4087–4094.
- Grineva, M., Grinev, M., & Lizorkin, D. (2009). Extracting key terms from noisy and multitheme documents. In *Proceedings of the 18th International Conference on World Wide Web*, 661–670.
- Habibi, M., Weber, L., Neves, M., Wiegandt, D.L., & Leser, U. (2017). Deep learning with word embeddings improves biomedical named entity recognition. *Bioinformatics*, 33(14), i37–i48.
- Hasan, K.S., & Ng, V. (2014). Automatic keyphrase extraction: A survey of the state of the art. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1262–1273. <https://doi.org/10.3115/v1/P14-1119>
- Hinton, G.E., Osindero, S., & Teh, Y.W. (2006). A fast learning algorithm for deep belief nets. *Neural Computation*, 18(7), 1527–1554. <https://doi.org/10.1162/neco.2006.18.7.1527>
- Howard, J., & Ruder, S. (2018). Universal language model fine-tuning for text classification. *ArXiv:1801.06146 [Cs, Stat]*. <http://arxiv.org/abs/1801.06146>
- Hulth, A. (2003). Improved automatic keyword extraction given more linguistic knowledge. In *Proceedings of the 2003 Conference on Empirical Methods in Natural Language*, 216–223.
- Hulth, A., Karlgren, J., Jonsson, A., Boström, H., & Asker, L. (2001). Automatic keyword extraction using domain knowledge. *International Conference on Intelligent Text Processing and Computational Linguistics*, 472–482.
- Hulth, A., & Megyesi, B.B. (2006). A study on automatically extracted keywords in text categorization. In *Proceedings of the 21st International Conference on Computational Linguistics and the 44th Annual Meeting of the Association for Computational Linguistics*, 537–544.
- Jones, S., & Staveley, M.S. (1999). Phrasier: A system for interactive document retrieval using keyphrases. In *Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, 160–167.
- Kelleher, D., & Luz, S. (2005). Automatic hypertext keyphrase detection. *IJCAI*, 5, *Proceedings of the Nineteenth International Joint Conference on Artificial Intelligence*, Edinburgh, Scotland, UK, 1608–1609.
- Kim, S.N., & Kan, M.Y. (2009). Re-examining automatic keyphrase extraction approaches in scientific articles. In *Proceedings of the Workshop on Multiword Expressions: Identification, Interpretation, Disambiguation and Applications*, 9–16.
- Kim, S.N., Medelyan, O., Kan, M.Y., & Baldwin, T. (2010). Semeval-2010 task 5: Automatic keyphrase extraction from scientific articles. In *Proceedings of the 5th International Workshop on Semantic Evaluation*, 21–26.



**Research Paper**

- Le, T.T.N., Le Nguyen, M., & Shimazu, A. (2016). Unsupervised keyphrase extraction: Introducing new kinds of words to keyphrases. *Australasian Joint Conference on Artificial Intelligence*, 665–671.
- Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C.H., & Kang, J. (2019). BioBERT: A pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, btz682. <https://doi.org/10.1093/bioinformatics/btz682>
- Liu, Z.Y., Huang, W.Y., Zheng, Y.B., & Sun, M.S. (2010). Automatic keyphrase extraction via topic decomposition. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, 366–376.
- Liu, Z.Y., Li, P., Zheng, Y.B., & Sun, M.S. (2009). Clustering to find exemplar terms for keyphrase extraction. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1-Volume 1*, 257–266.
- Luhn, H.P. (1957). A statistical approach to mechanized encoding and searching of literary information. *IBM Journal of Research and Development*, 1(4), 309–317.
- Matsuo, Y., & Ishizuka, M. (2004). Keyword extraction from a single document using word co-occurrence statistical information. *International Journal on Artificial Intelligence Tools*, 13(01), 157–169.
- Medelyan, O., Frank, E., & Witten, I.H. (2009). Human-competitive tagging using automatic keyphrase extraction. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 3-Volume 3*, 1318–1327.
- Mihalcea, R., & Tarau, P. (2004). TextRank: Bringing order into text. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, 404–411.
- Papagiannopoulou, E., & Tsoumakas, G. (2019). A review of keyphrase extraction. *ArXiv:1905.05044 [Cs]*. <http://arxiv.org/abs/1905.05044>
- Peters, M., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., & Zettlemoyer, L. (2018). Deep Contextualized Word Representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2227–2237. <https://doi.org/10.18653/v1/N18-1202>
- Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding with unsupervised learning. Technical Report, OpenAI.
- Sahrawat, D., Mahata, D., Kulkarni, M., Zhang, H., Gosangi, R., Stent, A., Sharma, A., Kumar, Y., Shah, R.R., & Zimmermann, R. (2019). Keyphrase Extraction from Scholarly Articles as Sequence Labeling using Contextualized Embeddings. *ArXiv Preprint ArXiv:1910.08840*.
- Salton, G., & Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. *Information Processing & Management*, 24(5), 513–523.
- Salton, G., Yang, C.S., & Yu, C.T. (1975). A theory of term importance in automatic text analysis. *Journal of the American Society for Information Science*, 26(1), 33–44.
- Tomokiyo, T., & Hurst, M. (2003). A language model approach to keyphrase extraction. In *Proceedings of the ACL 2003 Workshop on Multiword Expressions: Analysis, Acquisition and Treatment*, 33–40.
- Turney, P.D. (2000). Learning algorithms for keyphrase extraction. *Information Retrieval*, 2(4), 303–336.
- Turney, P.D. (2002). Learning to extract keyphrases from text. *ArXiv Preprint Cs/0212013*.





- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 5998–6008.
- Wan, X., & Xiao, J. (2008). Single document keyphrase extraction using neighborhood knowledge. *AAAI*, 8, 855–860.
- Wang, M., Zhao, B., & Huang, Y. (2016). Ptr: Phrase-based topical ranking for automatic keyphrase extraction in scientific publications. *International Conference on Neural Information Processing*, 120–128.
- Wang, X., Zhang, Y., Ren, X., Zhang, Y.H., Zitnik, M., Shang, J.B., Langlotz, C., & Han, J.W. (2019). Cross-type biomedical named entity recognition with deep multi-task learning. *Bioinformatics*, 35(10), 1745–1752.
- Witten, I.H., Paynter, G.W., Frank, E., Gutwin, C., & Nevill-Manning, C.G. (2005). Kea: Practical automated keyphrase extraction. In *Design and Usability of Digital Libraries: Case Studies in the Asia Pacific* (pp. 129–152). IGI global.
- Wu, Y., Schuster, M., Chen, Z., Le, Q.V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., & Macherey, K. (2016). Google's neural machine translation system: Bridging the gap between human and machine translation. *ArXiv Preprint ArXiv:1609.08144*.
- Zhang C.Z., Wang H.L., Liu Y., Wu D., Liao Y., & Wang B. (2008). Automatic keyword extraction from documents using conditional random fields. *Journal of Computational Information Systems*, 4(3), 1169–1180.
- Zhang, Q., Wang, Y., Gong, Y., & Huang, X.J. (2016). Keyphrase extraction using deep recurrent neural networks on twitter. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 836–845.
- Zhang, Y., Zincir-Heywood, N., & Milios, E. (2004). World wide web site summarization. *Web Intelligence and Agent Systems: An International Journal*, 2(1), 39–53.
- Zhao, W.X., Jiang, J., He, J., Song, Y., Achananuparp, P., Lim, E.P., & Li, X. (2011). Topical keyphrase extraction from twitter. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, 379–388.
- Li, L.S., Dang, Y.Z., Zhang, J., & Li, D. (2013). Term extraction in the automotive field based on conditional random fields. *Journal of Dalian University of Technology*, 53(2), 267–272.
- Li, S.J., Wang, H.F., Yu, S.W., & Xin, C.S. (2004). Application research of maximum entropy model for keyword automatic indexing. *Chinese Journal of Computers*, 27(9), 1192–1197.



This is an open access article licensed under the Creative Commons Attribution-NonCommercial-NoDerivs License (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

