

Article

Not peer-reviewed version

---

# Key Technologies of Intelligent Question Answering System for Power System Rules and Regulations Based on Improved BERTserini Algorithm

---

Ming Gao , [Mengshi Li](#) , [Tianyao Ji](#) <sup>\*</sup> , Nanfang Wang , [Guowu Lin](#) , [Qinghua Wu](#)

Posted Date: 23 November 2023

doi: 10.20944/preprints202311.1513.v1

Keywords: intelligent question answering system; improved BERTserini algorithm; rules and regulations; information retrieval



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: 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.

## Article

# Key Technologies of Intelligent Question Answering System for Power System Rules and Regulations Based on Improved BERTserini Algorithm

Ming Gao <sup>1</sup>, Mengshi Li <sup>1</sup>, Tianyao Ji <sup>1,\*</sup>, Nanfang Wang <sup>2</sup>, Guowu Lin <sup>2</sup> and Qinghua Wu <sup>1</sup>

<sup>1</sup> College of Electric Power Engineering, South China University of Technology, Guangzhou 510640, China; 202011002583@mail.scut.edu.cn (M.G.); mengshili@scut.edu.cn (M.L.); wuqh@scut.edu.cn (Q.W.).

<sup>2</sup> School of Electrical Engineering, Guangzhou City University of Technology, Guangzhou, China; wnf2184649296@163.com (N.W.); 16676676716@163.com (G.L.).

\* Correspondence: tyji@scut.edu.cn (T.J.);

**Abstract:** With the continuous breakthrough of natural language processing, the application of intelligent question answering technology in electric power system has attracted wide attention. However, at present, the traditional question answering system has poor performance and is difficult to be applied in engineering practice. This paper proposes an improved BERTserini algorithm for intelligent answering of electric power regulations based on a BERT model. The proposed algorithm is implemented in two stages. The first stage is the text segmentation stage, where a multi-document long text preprocessing technique is utilized that accommodates the rules and regulations text, and then Anserini is used to extract paragraphs with high relevance to the given question. The second stage is answer generation and source retrieval stage, where a two-step fine-tuning based on the Chinese BERT model is applied to generate precise answers based on given questions, while the information regarding documents, chapters, and page numbers of these answers are also output simultaneously. The algorithm proposed in this paper eliminates the necessity for manual organization of professional question-answer pairs, thereby effectively reducing the manual labor cost compared to traditional question answering systems. Additionally, this algorithm exhibits a higher degree of exact match rate and a faster response time for providing answers.

**Keywords:** intelligent question answering system; improved BERTserini algorithm; rules and regulations; information retrieval

## 1. Introduction

Intelligent question answering system is an innovative information service system which integrates natural language processing, information retrieval, semantic analysis and artificial intelligence. The system mainly consists of three core parts, which are question analysis, information retrieval and answer extraction. Through these three parts, the system can provide users with accurate, fast and convenient answering services.

The representative systems of the intelligent question answering system include:

(1) Rule-based algorithms (1960s-1980s). The question-answering system based on this pattern mainly relies on writing a lot of rules and logic to implement the dialogue. ELIZA [1], developed by Joseph Weizenbaum in the 1960s, was the first chatbot designed to simulate a conversation between a psychotherapist and a patient. PARRY [2] is a question-and-answer system developed in the 1970s that simulates psychopaths. The emergence of ELIZA and PARRY provided diverse design ideas and application scenarios for subsequent intelligent question answering systems, thereby promoting the diversification and complexity of dialogue systems. However, the main problem of this model is its lack of flexibility and extensibility. It relies too much on rules or templates set by humans, and consumes a lot of time and manpower. When the questions become complicated, it is difficult to get satisfactory answers through simple rules set by the model.

(2) Statistics-based algorithms (1990s-2000s). The question-answering system based on this model adopts the method of statistical learning to learn patterns and rules from a large number of dialogue data. Common algorithms include Vector Space Model [3] and Conditional Random Fields [4]. ALICE (Artificial Linguistic Internet Computer Entity) [5] is an open-source natural language processing project. The system in question is an open-domain question answering platform capable of addressing queries across a multitude of subjects and domains. Jabberwacky [6] is an early intelligent chatbot employing machine learning and conversational models to enhance its responses continually. These systems are designed to train models that can learn the relationships between questions and answers present in the corpus. Therefore, these models can carry out more natural and smooth dialogue. However, the ability of context understanding and generalization ability is weak, so it is difficult to adapt to model sharing and transfer learning in various professional fields. Moreover, considering statistical models are trained on a large corpus, this kind of model may suffer from data bias when dealing with domain-specific problems and fail to provide accurate answers.

(3) Algorithms based on hybrid technology (2010s-early 2020s). The question-answering system, grounded on this model, can amalgamate diverse techniques encompassing rules, statistics, and machine learning. It leverages multiple input modalities, including speech, image, and text, to interoperate seamlessly. The overarching objective is to facilitate users in accomplishing specific tasks or goals within designated domains, such as booking, traveling, shopping, or ordering food. This synergistic integration of multifarious technologies and input modes fosters a more sophisticated and intelligent dialogue system. Typical question answering systems based on hybrid technology model include Apple's Siri [7], Microsoft's Cortana [8], Amazon's Alexa [9], Facebook's M [10] and Google's Google Assistant [11]. These systems are centered around artificial intelligence and natural language processing technology, aiming to furnish users with personalized and convenient information and services to cater to diverse needs.

The system built based on this pattern has stronger context understanding and personalized customization, but there are two shortcomings: first, the quality of dialogue in such a system is not stable; Secondly, the generalization ability of the model is limited. It is difficult to realize model sharing, transfer learning and answer generation in professional fields. The training of this model requires excessive investment in computing and data resources, and its training and deployment speed is slow.

(4) Algorithms based on pre-trained language (2020s). The model is based on pre-trained language models such as BERT [12], GPT (Generative Pre-trained Transformer) [13], etc. These models are pre-trained on large-scale data and they learn rich language representation and context understanding skills to generate more natural, fluid, and accurate responses. In addition, through the supervised training on domain-specific question answering datasets, the question answering system can answer questions in specialized professional fields. [14] proposed a BERTserini algorithm which improves the exact match rate of the question answering system. In comparison to the original BERT algorithm, the proposed method surpasses its processing byte limit and can provide accurate answers for multi-document long texts.

Although systems built on the BERTserini algorithm perform well on public datasets, there are some problems in the application in professional fields such as electrical power engineering. Considering the low exact match rate and poor answer quality, engineering applications of these models are challenging. The problems are mainly caused by the following aspects.

(1) Lack of model expertise: Language models such as BERT or GPT are usually pre-trained from large amounts of generic corpus collected on the Internet. However, the digital realm offers limited professional resources pertaining to industries like electrical power engineering. As a result, the model has insufficient knowledge reserve when dealing with professional question, which affects the quality of the answers; (2) Differences in document format: There are significant differences between the format of documentation in the electrical power engineering field and that of public datasets. The documents in the electrical power engineering field often exhibit unique formatting, characterized by an abundance of hierarchical headings. It is easy to misinterpret the title as the main content and mistakenly use it as the answer to the question, leading to inaccurate results; (3) Different scenario

requirements: Traditional answering systems do not need to pay attention to the source of answers in the original document. However, a system designed for professional use must provide specific source information for its answers. If such information is not provided, there may arise doubts regarding the accuracy of the response. This further diminishes the utility of the application in particular domains.

This paper proposes an improved BERTserini algorithm to construct an intelligent question answering system in the field of electrical power engineering. The proposed algorithm is divided into two stages:

The first stage is text segmentation. During this phase, the text is segmented and preprocessed. Firstly, a multi-document long text preprocessing method that supports rules and regulations text is proposed. This approach can accurately segment rules and regulations text and generate an index file of answer location information. By doing so, the system can better comprehend the structure of the regulation text, enabling it to locate the answer to the user's question more accurately. Secondly, through the FAQ [15] pre-module, high-frequency questions are intercepted for question pre-processing. This module matches and classifies user-raised questions based on a pre-defined list of common questions, intercepting and addressing high-frequency issues. This reduces the repetition of processing the same or similar problems and enhances the system's response efficiency. Finally, Anserini [16] is employed to extract several paragraphs highly relevant to user problems from multi-document long text. Anserini is an information retrieval tool based on a vector space model that represents a user question as a vector and each paragraph in a multi-document long text as a vector. By calculating the similarity between the user problem vector and each paragraph vector, several paragraphs with high relevance to the user problem can be selected. These paragraphs serve as candidate answers for the system to further analyze and generate the final answer.

The second stage is answer generation and source retrieval stage. During this phase, the Chinese Bert model undergoes Fine-tuning [17], which comprises two steps involving key parameter adjustments. This process enhances the model's comprehension of the relationship between the question and the answer, thereby improving the accuracy and reliability of the generated response. Subsequently, based on the input question, the Bert model extracts several candidate answers from the N paragraphs with the highest similarity to the question, as determined by Anserini. The user can then filter through these multiple relevant paragraphs to identify the answer that best aligns with their query. Finally, the candidate answers are weighted, and the highest-rated answer is outputted along with the chapter and position information of the answer in the original document. This approach facilitates users in quickly locating the most accurate answer while providing pertinent contextual information.

The improved BERTserini algorithm proposed in this paper has three main advantages.

(1) The proposed algorithm implements multi-document long text preprocessing technology tailored for rules and regulations text. Through optimization, the algorithm segments rules and regulations into distinct paragraphs based on its inherent structure and supports answer output with reference to chapters and locations within the document. The effectiveness of this pretreatment technology is reflected in the following three aspects: First, through accurate segmentation, paragraphs that may include questions can be extracted more accurately, thus improving the accuracy of answer generation. Secondly, the original Bert model exhibits a limitation that it outputs the heading of rules and regulations text as the answer frequently. To address this issue, an improved BERTserini algorithm has been proposed. Finally, the algorithm is able to accurately give the location information of answers in the original document chapter. The algorithm enhances the comprehensiveness and accuracy of reading comprehension, generating answers to questions about knowledge and information contained in professional documents related to the field of electric power. Consequently, this leads to a marked improvement in answer quality and user experience for the question answering system.

(2) The proposed algorithm optimizes the training of the corpus in the field of electrical power engineering and fine-tunes the parameters of the large language model. This method eliminates the necessity for manual organization of professional question-answer pairs, knowledge base engineering, and manual template establishment in BERT reading comprehension, thereby effectively reducing labor costs. This enhancement significantly enhances the accuracy and efficiency of the question-answering system.

(3) The proposed algorithm has been developed for the purpose of enhancing question answering systems in engineering applications. This algorithm exhibits a higher degree of exact match rate of questions and a faster response for providing answers.

## 2. Background of the technology

### 2.1. FAQ

Frequently Asked Questions (FAQs) are a collection of frequently asked questions and answers designed to help users quickly find answers to their questions [15]. The key is to build a rich and accurate database of preset questions, which consists of questions and the corresponding answers. They are manually collated from the target documents. The FAQ provides an answer that corresponds to the user's question by matching it with the most similar question.

### 2.2. BM25 algorithm

The Best Match 25 (BM25) algorithm [16]-[17] was initially proposed by Stephen Robertson and his team in 1994 and applied to the field of information retrieval. It is commonly used to calculate the relevance score between documents and queries. The main logic of BM25 is as follows: Firstly, the query statement involves word segmentation to generate morphemes. Then, the relevance score between each morpheme and the search result is calculated.

Finally, by weighting summing the relevance scores of the morpheme with the search results, the relevance score between the retrieval query and the search result documents is obtained. The formula for calculating BM25 algorithm is as follows:

$$Score(D, Q) = \sum_i^n W_i \cdot R(q_i, D) \quad (1)$$

In this context,  $Q$  represents a query statement,  $q_i$  represents a morpheme obtained from  $Q$ . For Chinese, the segmented results obtained from tokenizing query  $Q$  can be considered as morpheme  $q_i$ .  $D$  represents a search result document.  $W_i$  represents the weight of morpheme  $q_i$ , and  $R(q_i, D)$  represents the relevance score between morpheme  $q_i$  and document  $D$ . There are multiple calculation methods for weight parameter  $W_i$ , with Inverse Document Frequency (IDF) being one of the commonly used approaches. The calculation process for IDF is as follows:

$$IDF(q_i) = \log\left(\frac{N - n(q_i) + 0.5}{n(q_i) + 0.5}\right) \quad (2)$$

In the equation,  $N$  represents the total number of documents in the index, and  $n(q_i)$  represents the number of documents that contain  $q_i$ .

Finally, the relevance scoring formula for the BM25 algorithm can be summarized as follows:

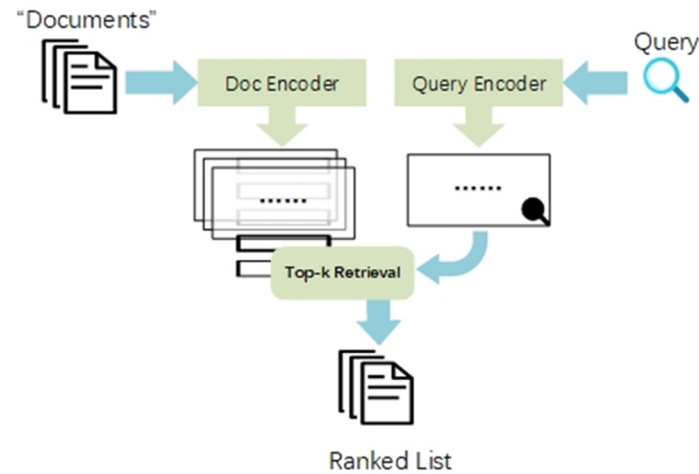
$$Score(D, Q) = \sum_{i=1}^n IDF(q_i) \cdot \frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot \left(1 - b + b \cdot \frac{|D|}{avgdl}\right)} \quad (3)$$

where  $k_1$  and  $b$  are adjustment factors,  $f(q_i, D)$  represents the frequency of morpheme  $q_i$  appearing in document  $D$ ,  $|D|$  denotes the length of document  $D$ , and  $avgdl$  represents the average length of all documents.



### 2.3. Anserini

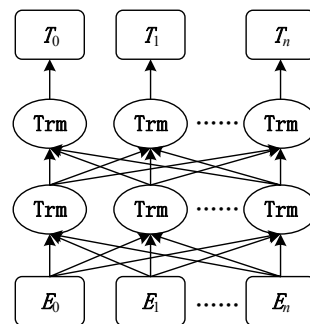
Anserini [18] is an open-source information retrieval toolkit that supports various text-based information retrieval research and applications. The goal of Anserini is to provide an easy-to-use and high-performance toolkit that supports tasks such as full-text search, approximate search, ranking, and evaluation on large-scale text datasets. It enables the conversion of text datasets into searchable index files for efficient retrieval and querying. Anserini incorporates a variety of commonly used text retrieval algorithms, including the BM25 algorithm. With Anserini, it becomes effortless to construct a BM25-based text retrieval system and perform efficient search and ranking on large-scale text collections. The flowchart of the algorithm is illustrated in Figure 1.



**Figure 1.** The flowchart of the Anserini algorithm.

### 2.4. BERT model

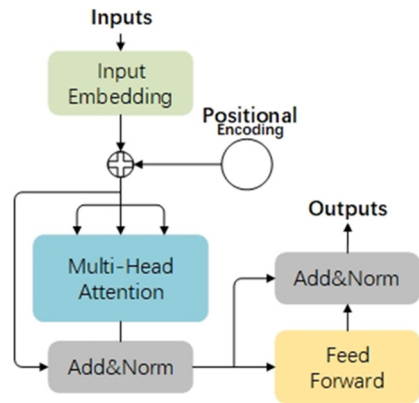
Bidirectional Encoder Representations from Transformers (BERT) [12] is a pre-trained language model proposed by Google in 2018. The model structure is shown in Figure 2. In the model,  $E_i$  represents the encoding of words in the input sentence, which is composed of the sum of three word embedding features. The three word embedding features are Token Embedding, Position Embedding, and Segment Embedding. The integration of these three words embedding features allows the model to have a more comprehensive understanding of the text's semantics, contextual relationships, and sequence information, thus enhancing the BERT model's representational power. The transformer structure in the figure is represented as Trm. The  $T_i$  represents the word vector that corresponds to the trained word  $E_i$ .



**Figure 2.** Architecture of BERT.

BERT exclusively employs the encoder component of the Transformer architecture. The encoder is primarily comprised of three key modules: Positional Encoding, Multi-Head Attention, and Feed-Forward Network. Input embeddings are utilized to represent the input data. Addition and

normalization operations are denoted by “Add&norm”. The fundamental principle of the encoder is illustrated in Figure 3.

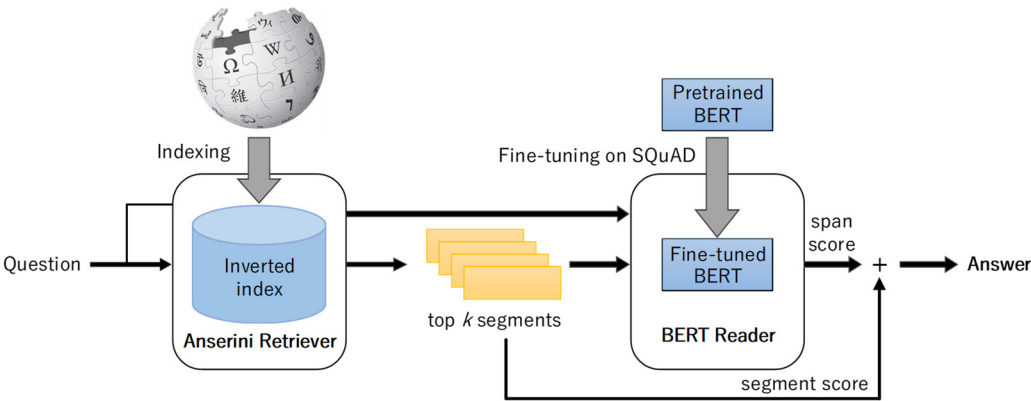


**Figure 3.** Transformer Encoder Principle.

In recent years, several Chinese BERT models have been proposed in the Chinese language domain. Among these, the chinese-BERT-wwm-ext model [19] released by the HIT-iFLYTEK Language Cognitive Computing Lab (HFL) has gained significant attention and serves as a representative example. This model was pre-trained on a corpus of Chinese encyclopedias, news articles, and question-and-answer texts, which contains a total of 5.4 billion words. The model uses the whole-word masking (wwm) strategy.

2.5. BERTserini algorithm

The architecture of BERTserini algorithm [14] is depicted in Figure 4. The algorithm employs the Anserini information extraction algorithm in conjunction with a pre-trained BERT model. In this algorithm, the Anserini retriever is responsible for selecting text paragraphs containing the answer, which are then passed to the BERT reader to determine the answer scope. This algorithm exhibits significant advantages over traditional algorithms. It demonstrates fast execution speed similar to traditional algorithms while also possessing the characteristics of end-to-end matching, resulting in more precise answer results. Furthermore, it supports extracting answers to questions from multiple documents.

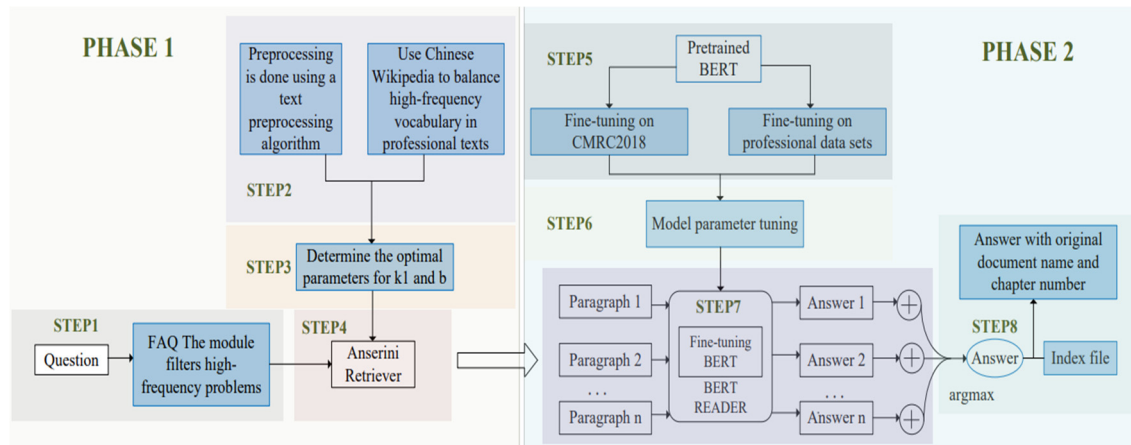


**Figure 4.** Architecture of BERTserini.

3. Improved BERTserini algorithm

3.1. Algorithm description

The improved BERTserini algorithm presented in this paper can be divided into two stages, and the flowchart is illustrated in Figure 5.



**Figure 5.** Flowchart of the proposed algorithm.

### (1) Phase 1: Text Segmentation Stage

The first stage is text segmentation stage, which comprises two key components: (1) Question preprocessing: The FAQ module is utilized to intercept high-frequency questions in advance, thereby achieving question preprocessing. If the FAQ module cannot provide an answer that corresponds to the user's query, then the query is transferred to the subsequent stage of paragraph extraction. Anserini retrieval technology is utilized for paragraph extraction, enabling the rapid extraction of highly relevant paragraphs which are pertinent to user queries within multi-document long text. (2) Document preprocessing: Due to the high degree of keyword overlap in power regulation documents. The paper proposes a multi-document long text preprocessing method supporting regulation texts, which can accurately segment the regulation texts and support the retrieval and tracing of the answer chapters' sources.

#### STEP1 The FAQ module filters out high-frequency problems

The FAQ module is designed to pre-process questions by intercepting and filtering out high-frequency problems. To achieve this, the module requires a default question library that contains a comprehensive collection of manually curated questions and their corresponding answer pairs from the target document. By matching the most similar question to the user's inquiry, the FAQ module can efficiently provide an accurate answer based on the corresponding answer to the question.

The FAQ module employs ElasticSearch, an open-source distributed search and analysis engine, to match user queries in a predefined question library. ElasticSearch is built upon the implementation of Lucene, an open-source full-text search engine library released by the Apache Foundation, and incorporates Lucene's BM25 text similarity algorithm. This algorithm calculates similarity by evaluating word overlap between the user query's text and the default question library, as shown in (3).

The FAQ module will directly return the preset answer to the matched question if the BM25 score returned by ElasticSearch exceeds the predetermined threshold. In cases where the return score falls below this threshold, instead of returning an answer, the question is referred to subsequent steps.

#### STEP2 Text preprocessing and document index generation

The first step involves two tasks. Firstly, Due to the high coincidence rate of keywords in regulations, employing Anserini directly for paragraph retrieval and calculation may lead to the problem of high weight value of high-frequency words. Consequently, when indexing documents, it is advisable to incorporate Chinese Wikipedia text data alongside rules and regulations to counteract the undue influence of excessive word frequency. Secondly, a novel multi-document long text preprocessing algorithm is proposed. This algorithm accurately segments the regulation text while retaining chapter position information of each paragraph and generating the index file. Specific steps include:

Convert documents in pdf or docx format to plain text in txt format.

Remove irrelevant information such as header/footer and page number.



Use regular expressions to extract the title number from the text (for example: 3.3.1), and match the title number to the text.

Use rules to filter out paragraphs in the text such as tables and pictures that are not suitable for machine reading comprehension.

Use Anserini to divide the text title number into words and index the corresponding text.

STEP3 Determine the two parameters  $k_1$  and  $b$

The  $k_1$  and  $b$  parameters utilized in the Anserini module are empirically selected to determine the optimal parameters for this study. A specific methodology is employed, starting from 0.1 within their respective value ranges and incrementing by 0.05 to systematically explore all possible combinations of  $k_1$  and  $b$  values. The selection of the best  $k_1$  and  $b$  values is based on the accuracy assessment of the second stage Bert reading comprehension module.

STEP4 Extract paragraphs and generate paragraph scores

Based on the user's question, Anserini extracts relevant paragraphs from the preprocessed document by filtering out those that are not related to the query. It then matches the question with the paragraphs in the index and selects the top  $N$  paragraphs with the highest relevance to the question. These paragraphs are scored using the BM25 algorithm and recorded them as  $S_{\text{anserini}}$ .

(2) Phase 2: answer generation and source retrieval stage

The second stage is the answer generation and source retrieval stage. After undergoing two steps of fine-tuning and key parameter tuning, the model is capable of extracting accurate answers from  $N$  paragraphs based on the given question. Additionally, the model can output the chapter information of the answer in the original document according to the index file.

STEP5 Select the appropriate Chinese Bert model and fine-tune it

In this research, the Chinese-Bert-WWM-EXT Base model is chosen as the foundational framework. The initial step involves fine-tuning the model using the Chinese Open domain Question answering dataset (CMRC2018). Subsequently, a second round of fine-tuning is conducted by employing the training exam questions related to rules and regulations as specialized datasets.

STEP6 Algorithm parameter tuning

Based on the structural and characteristic features of regulatory documents, the following five crucial parameters of the improved BERTserini algorithm have been optimized:

**paragraph\_threshold.** The paragraph threshold is employed to exclude paragraphs with Anserini scores below this specified limit, thereby conserving computational resources.

**phrase\_threshold.** The answer threshold serves as a filter, excluding responses with a Bert reader score below the specified limit.

**remove\_title.** Removes the paragraph title. If this item is True ( $y=\text{True}$ ,  $n=\text{False}$ ), paragraph headings are not taken into account when the Bert reader performs reading comprehension.

**max\_answer\_length.** The maximum answer length. The maximum length of an answer is allowed to be extracted when the Bert reader performs a reading comprehension task.

**mu.** Score weight is implemented to evaluate both the answer and paragraph using the Bert reader and Anserini extractor, subsequently calculating the final score value of the answer.

STEP7 Extract the answers and give a reading comprehension score

Bert is used to extract the exact answers to the question from the  $N$  paragraphs extracted by Anserini. The sum of the probability of starting and ending positions (logits) for each answer predicted by the model is used as the score of the answer generated by the Bert reading comprehension module. It can be expressed by the following equation:

$$S_{\text{bert}} = \max(\text{start logit}) + \max(\text{end logit}) \quad (4)$$

STEP8 The candidate answers are scored by a comprehensive weighted score, rank the answers by score, output the answer with the highest score, and give the original document name and specific chapter information for the answer.

Use the following equation to calculate the overall weighted score of the answer:

$$S = (1 - \mu) * S_{\text{anserini}} + \mu * S_{\text{bert}} \quad (5)$$

The final score of the answer is calculated by the above formula.  $S_{\text{anserini}}$  represents the BM25 score returned by the Anserini extractor, and  $S_{\text{bert}}$  represents the answer score returned by Bert. The

answers are sorted by the calculated answer score, and the final output is the answer with the highest score. According to the index file, the original document name and chapter information are output together.

### 3.2. Main innovations

(1) Multi-document long text preprocessing method which can process rules and regulations text and support answer provenance retrieval.

In this paper, a multi-document long text preprocessing method is proposed that facilitates answer provenance retrieval and can effectively process the rules and regulations text, which provides a technical path for the construction of intelligent question answering system in specific professional fields. The innovation point of this method is reflected in STEP2. This method divides the rules and regulations into chapters. The original document name of each paragraph and its chapter number information can be preserved. To address the issue of excessive frequency of certain proper nouns, the method incorporates text data from Chinese Wikipedia and performs balance processing. By incorporating a larger corpus, the frequency of a specific proper noun in the text can be effectively diminished, thereby mitigating its influence on the model. This innovative preprocessing method can improve the calculation effect of the subsequent reading comprehension module. The answer can be provided in the original document, including chapter and location information.

(2) Determination of optimal parameters of Anserini and improved BERTserini algorithm

① Determination of the optimal parameters of Anserini. In STEP3, the optimal parameters of Anserini are determined. All possible combinations of  $k_1$  and  $b$  are experimentally tried one by one. And the best value is selected according to the answer performance of the subsequent reading comprehension module questions. The determination of the optimal parameters of Anserini improves the performance of the intelligent question answering system and the exact match of answers (EM).

② Determination of the optimal parameters of the improved BERTserini algorithm: In STEP6, the optimal parameters of the improved BERTserini algorithm are determined. According to the structure and characteristics of regulation documents, five important parameters are optimized. Thus, the algorithm can determine the reasonable threshold of generating candidate answers when the Bert reading comprehension module performs the reading comprehension task. And the answer generation does not take into account the paragraph title and the optimal overall rating weight and other details that constitute high-quality questions and answers.

(3) Fine-tuning of multi-data sets for Bert reading comprehension model

This step is illustrated in STEP5. The Bert model is pre-trained using the CMRC2018 data, and a two-step Fine-tuning was carried out using the existing rules and regulations exam questions. By making full use of data sets in different fields, the accuracy and generalization ability of the model are improved. This method achieves better results in question answering system. At the same time, this method also reduces the time and labor cost required for manual editing of question answering pairs in traditional model training. It also significantly improves Bert's reading comprehension of rules and regulations.

(4) Clever use of FAQ

The clever use of the FAQ is reflected in STEP1. In this paper, the existing rules and regulations are used to train and test questions, which constitutes the questions and answers pairs required by the pre-FAQ module to intercept some high-frequency questions. In this way, a low-cost FAQ module is constructed, which improves the answering efficiency of high-frequency questions, and also improves the exact match rate (EM) of the intelligent question answering system.

## 4. Results analysis of the experiment

### 4.1. Data description.

#### 4.1.1. Document description.

For the present study, a total of 30 documents including regulations, provisions, and operation manuals related to the theme of power safety are selected, such as a company power grid work regulations. The total size of the documents is 30.1 MB, and the intelligent system is required to preprocess all the content within the documents, perform machine reading comprehension, and efficiently answer questions.

#### 4.1.2. Fine-tuning dataset description.

In this study, four datasets are experimented for fine-tuning the Bert model, which include Chinese Machine Reading Comprehension 2018 (CMRC2018) [20], Delta Reading Comprehension Dataset (DRCD) [21], Safety Procedure Test Item data set (SPTI), and a dataset generated through data augmentation based on documentations of a power grid company. The first two datasets are open-source. Based on end-to-end manual evaluation, the results indicate that the model trained using CMRC2018 data performs the best in this study. Therefore, it has been selected as the fine-tuning training dataset. The dataset follows the format of the SQuAD dataset [22]. It consists of a total of 10,142 training samples, 3,219 validation samples, and 1,002 testing samples. The overall size of the dataset is 32.26MB. The SPTI consists of 1020 training and examination questions related to electrical safety regulations.

#### 4.1.3. BERT model description.

In this study, the Chinese-BERT-wwm-ext model released by the HFL is used for training.

#### 4.1.4. Parameter Tuning Explanation for Improved BERTserini Algorithm.

The parameter settings in this study are as follows. `paragraph_threshold=10`, `phrase_threshold=0`, `remove_title=n` (`n=False`, `y=True`), if `remove_title=y`, the paragraph titles will not be considered by the BERT reader algorithm during reading comprehension. `max_answer_length=50`, `mu=0.6`.

The parameter in the BM25 algorithm used in the Anserini module has a value range of (0-1), and the parameter has a value range of (0-3).

### 4.2. Document preprocessing performance.

In accordance with the document pre-processing algorithm proposed, the document format output by Anserini is illustrated in Table 1. Within this context, “text” denotes the output paragraphs obtained from Anserini, “paragraph\_score” represents the specific score assigned to each paragraph, and “docid” indicates the name of the document along with the corresponding section information where the paragraph is situated.

Table 1. Document preprocessing of the Anserini Module.

[
{
“text” : “术语和定义 电力设施 . 应用到电力系统中的发电、变电、输电、 配电和供电有关设备的总称。” ,
“paragraph_score” : 18.094900131225586,
“docid” : “TAG%%××电网有限责任公司电力安全工作规程.pdf%%3.3_1” ,
},
{
“text” : “管理内容与方法 事故事件等级划分标准 . 事故分类。事故分为电 力人身事故、电力设备事故、电力安全事故共 3 类,根据事故后果严重程度从高到 低排序,事故等级依次分为特别重大、重大、较大和一般共 4 级,事故等级划分 标准详见附录一 A.1;” ,
“paragraph_score” : 13.659899711608887,
“docid” : “TAG%%××有限责任公司事故事件管理办法.pdf%%5.1.1_1” ,
},
...
]

4.3. Question-Answering Performance

The Comparison of question-answering performance before and after the improvement of the BERTserini algorithm is presented in Table 2. It can be observed that the original BERTserini algorithm exhibits inaccuracies in extracting the start and end positions of answers when addressing power regulations and standards questions, and even results in incomplete sentences. Compared to the original BERTserini algorithm, the improved BERTserini algorithm proposed in this paper can accurately locate the paragraph containing the correct answer and perform precise answer extraction. Additionally, it removes specific details like paragraph headings during the answering process, adapting to the structural characteristics of professional domain regulatory texts. The answers to certain questions are more accurate and concise than manually generated standard answers.

**Table 2.** Comparison of question-answering performance before and after the improvement of the BERTserini algorithm.

Question	Standard answer	Original BERTserini Algorithm			Improved BERTserini Algorithm		
		Answer	Whether to exact match	Trace the source and results of the answers	Answer	Whether to exact match	Trace the source and results of the answers
被审核单位应于什么时候将相关文件、资料发至各专业审核组？ (When should the audited units send the relevant documents and information to the professional audit teams?)	被审核单位于审核前 5 个工作日。 (The audited unit should send the relevant documents and information 5 working days before the audit.)	被审核单位于审核前 5 个工作日。 (The audited unit should send the relevant documents and information 5 working days before the audit.)	Yes	No	审核前 5 个工作日。 (5 working days before the audit.)	Yes	Yes (《××有限责任公司安全生产风险管理体系审核业务指导书》5.1.5) (×× Power Grid Co., LTD. Safety Production Risk Management System Audit Business Guide 5.1.5)
季度安全生产重点督查内容是什么？ (What is the key inspection content of quarterly safety production?)	每季度末，各级安全监管部应根据年度工作计划、季节特点、重点工作安排等，确定下一季度的督查重点内容。 (At the end of each quarter, the safety supervision department at all levels shall determine the key contents of supervision in the next quarter according to the annual work plan, seasonal characteristics, and key work arrangements.)	季度安全生产重点督查内容。 (Quarterly safety production focus supervision within.)	No	No	每季度末，各级安全监管部应根据年度工作计划、季节特点、重点工作安排等，确定下一季度的督查重点内容。 (At the end of each quarter, the safety supervision department at all levels shall determine the key contents of supervision in the next quarter according to the annual work plan, seasonal characteristics, and key work arrangements.)	Yes	Yes (《××有限责任公司安全生产督查业务指导书》5.2.2) (×× Power Grid Co., LTD. Safety Production Risk Management System Audit Business Guide 5.2.2)



纠正是指什么？ (What does correction mean?)	为消除已发现的不符合所采取的措施。 (Measures taken to eliminate nonconformities that have been identified.)	纠正是指为消除已发现的不符。 (Correction means the elimination of discrepancies that have been found.)	No	No	为消除已发现的不符合所采取的措施。 (Measures taken to eliminate nonconformities that have been identified.)	Yes	Yes (《××有限责任公司安全区代表管理业务指导书》4.5) (×× Power Grid Co., LTD. Safety Zone Representative Management Service Guide 4.5)
公司直属各单位安全监管部负责什么？ (What is the safety supervision department directly under the company responsible for?)	组织编制、发布本单位年度安措计划，监督、考核本单位及所辖县（区）级供电单位安措计划的实施情况。 (To organize the preparation and issuance of the annual safety measure plan of the unit, supervise and evaluate the implementation of the safety measure plan of the unit and the power supply units at the county.) (district) level under its jurisdiction.)	公司直属各单位安全监管部负责组织编制。 (The safety supervision department of each unit directly under the company is responsible for organizing and compiling.)	No	No	组织编制、发布本单位年度安措计划，监督、考核本单位及所辖县（区）级供电单位安措计划的实施情况。 (To organize the preparation and issuance of the annual safety measure plan of the unit, supervise and evaluate the implementation of the safety measure plan of the unit and the power supply units at the county.) (district) level under its jurisdiction.)	Yes	Yes (《××有限责任公司安全技术劳动保护措施管理业务指导书》5.1.2) (×× Power Grid Co., Ltd. Safety zone Representative Management Service Guide 5.1.2)
什么是电力设施？ (What is an electric utility?)	应用到电力系统中的发电、变电、输电、配电和供电有关设备的总称。 (A general term for equipment related to generation, transformation, transmission, distribution and supply used in power systems.)	电力设施应用到电力。 (Utility applied to electricity.)	No	No	应用到电力系统中的发电、变电、输电、配电和供电有关设备的总称。 (A general term for equipment related to generation, transformation, transmission, distribution and supply used in power systems.)	Yes	Yes (《××电网有限责任公司电力安全工作规程》3.3) (×× Power Grid Co., LTD. Power Safety Working Regulations 3.3)

#### 4.4. Comparison of different algorithms.

This paper uses the Exact Match rate (EM), Recall rate (R), and F1 score to measure the question-answering performances of different algorithms. Among them, EM represents the percentage of questions in the question answering system where the answers provided are an exact match with the standard answers.

The specific calculation formula is as follows:

$$EM = \frac{1}{n_{samples}} \sum_{i=0}^{n_{samples}-1} I(\hat{y}_i = y_i) \quad (4)$$

where  $n_{samples}$  represents the total number of samples.  $I(x)$  is an indicator function that takes the value of 1 when  $\hat{y}_i$  is identical to  $y_i$ , and 0 otherwise. As can be seen from the formula, a higher EM value indicates a higher exact match.

Recall rate is to determine the proportion between the number of questions accurately answered by the question-answering system and the total number of questions.

The calculation formula is as follows:

$$R = \frac{TP}{TP + FN} \quad (5)$$

where  $TP$  signifies the accurate count of questions answered correctly by the question-answering system. Conversely,  $FN$  denotes the incorrect count of questions that were responded to inaccurately by the system.

The calculation formula for the F1 score is as follows:

$$F1 = \frac{2(EM * R)}{EM + R} \quad (6)$$

In this study, the initial BERTserini algorithm [17] is used in Algorithm 1. In Algorithm 2, an additional pre-processing algorithm is incorporated based on Algorithm 1. Fine-tuning is conducted using the CMRC2018 dataset, and parameter optimization for the BERTserini algorithm is performed. Algorithm 3 is an extension of Algorithm 2, incorporating the SPTI dataset for fine-tuning. Algorithm 4 is an improved version of the Bertserini algorithm, which is based on Algorithm 3. In addition to incorporating the SPTI dataset for fine-tuning, a pre-processing FAQ module based on short-text similarity calculation is added to filter out frequently asked questions. This module enables more efficient and effective preprocessing of the questions.

As shown in Table 3, the EM value for Algorithm 1 is only 0.261, indicating poor performance. After adopting Algorithm 3, the EM value reaches 0.702, representing an improvement of 62.8%. After adopting the proposed Algorithm 4, the EM value reaches 0.856, demonstrating the best performance. In comparison to Algorithm 1, the proposed algorithm achieved an improvement of 69.5% in terms of EM value, 53.6% in terms of R value, and 63.7% in terms of F1 value. These results demonstrate a practical level of engineering advancement.

**Table 3.** Comparison of different algorithms.

Algorithm m	Content	EM		R		F1	
		Value	Percentage of improvement	Value	Percentage of improvement	Value	Percentage of improvement
Algorithm m 1	Original BERTserini	0.261	—	0.453	—	0.331	—

Algorithm	Document preprocessing + Original BERTserini	0.50		0.78		0.61	
	+Fine-Tuning (CMRC2018) + Parameter tuning	2	48%	3	42.1%	5	46.1%
m 2							
Algorithm	Document preprocessing + Original BERTserini	0.70		0.91		0.79	
	+Fine-Tuning (CMRC2018) + Parameter tuning + Fine-Tuning (SPTI)	2	62.8%	9	50.7%	6	58.4%
m 3							
Algorithm	Document preprocessing + Original BERTserini	0.85		0.97		0.91	
	+Fine-Tuning (CMRC2018) + Parameter tuning + Fine-Tuning (SPTI) +FA	6	69.5%	6	53.6%	2	63.7%
m 4							
Q							

Note: In the table “Percentage of improvement” is calculated based on the values of Algorithm 1 as a reference point. It represents the relative increase in evaluation value (EM, R, and F1) achieved by Algorithm 2, Algorithm 3, and Algorithm 4 compared to Algorithm 1.

4.5. Engineering Application

An intelligent question-answering system for power regulations and standards is constructed based on the proposed improved BERTserini algorithm and experimental data presented in this paper, as shown in Figure 6.

The system provides users with a multi-turn interactive question-answering interface on the topic of power safety, as illustrated in Figure 6a. Users can ask questions by either voice input or manual input. After sending the question, they will receive the system’s response within 400ms. Clicking on the “view details” link below the answer will cause the system to pop up a window displaying the source of the answer, including the name of the original document and the chapter number, as shown in Figure 6b. Clicking on the “full text” link allows users to view the content of the original document where the answer is located, as shown in Figure 6c.



a. Multi-turn interactive question-answering interface b. Knowledge details page c. Full-text source page

**Figure 6.** Intelligent question-answering system built based on the proposed method.

## 5. Conclusions

The improved BERTserini algorithm proposed in this paper is designed for intelligent question-and-answer processing of power regulation documents. In comparison to the conventional BERTserini algorithm, this approach offers the following advantages:

(1) The improved BERTserini algorithm can enhance the accuracy in generating answers for specialized domain problems by fine-tuning in professional data sets and algorithm parameter optimization. Notably, the algorithm eliminates the need for manual sorting of question-answer pairs and manual template construction, thereby significantly reducing labor costs.

(2) The improved BERTserini algorithm supports multi-document long text preprocessing for rules and regulations. This algorithm is capable of answering documents containing 30+ rules and regulations with a length of 30M+ bytes. Additionally, it returns the answer along with the document name and chapter page number information. The average response time of the algorithm is within 400 ms, fully meeting the requirements of engineering applications.

(3) The improved BERTserini algorithm can significantly enhance the exact match rate of intelligent question answering systems when processing rules and regulations text in professional fields. Experimental data indicate that, compared to the original BERTserini algorithm, the exact match rate of the proposed method is improved by 69.5%, R value is increased by 53.6%, F1 value is increased by 63.7%.

The improved BERTserini algorithm has broad applicability in the research and development of intelligent question-answering systems for rules and regulations across various domains. It offers valuable insights and references for constructing similar systems in other areas. Furthermore, by integrating the algorithm with the ongoing technical progress of large language models, the effectiveness of the question-answering system can be continually optimized and refined.

**Author Contributions:** Writing—original draft preparation, M.G.; validation, M.G. and T.J.; formal analysis, M.L.; investigation, Q.W.; writing—review and editing, N.W., G.L.; visualization, T.J.; supervision, M.G.; M.G. and M.L. conceived the idea and provided resources. N.W. and G.L. wrote the manuscript. T.J. and M.G. provided guidelines. M.G., M.L., T.J., N.W., G.L. and Q.W. designed the study and participated in the experiment. All authors have read and agreed to the published version of the manuscript.

**Funding:** Please add: This research received no external funding.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

1. Shum H Y, He X, Li D. From Eliza to XiaoIce: Challenges and Opportunities with Social Chatbots[J]. *Frontiers of Information Technology & Electronic Engineering*, 2018, 19(1): 10-26.DOI:10.1631/FITEE.1700826.
2. Xavier A. Intelligence artificielle et psychiatrie: nocés d’or entre Eliza et Parry[J]. *Linformation Psychiatrique*, 2017, 93(1):51-56. DOI: 10.1684/ipe.2017.1582.
3. Hb B G, Kumar M A, Kp S. Vector Space Model as Cognitive Space for Text Classification[J]. 2017.DOI:10.48550/arXiv.1708.06068.
4. Ma P, Jiang B, Lu Z, et al. Cybersecurity Named Entity Recognition Using Bidirectional Long Short-Term Memory with Conditional Random Fields[J]. *Tsinghua Science and Technology*, 2021, 26(3):259-265. DOI: 10.26599/TST.2019.9010033.
5. Mittal A, Agrawal A, Chouksey A, et al. IJARCCCE A Comparative Study of Chatbots and Humans[J]. 2016.DOI:10.17148/IJARCCCE.2016.53253.
6. Jabberwacky Website [Online]. Available: <http://www.jabberwacky.com/>
7. Bohouta G, Kpuska V Z. Next-Generation of Virtual Personal Assistants (Microsoft Cortana, Apple Siri, Amazon Alexa and Google Home) [C]//IEEE CCWC 2018, The 8th IEEE Annual Computing and Communication Workshop and Conference. IEEE, 2018.
8. Poushneh A. Impact of auditory sense on trust and brand affect through auditory social interaction and control[J]. *Journal of Retailing and Consumer Services*, 2021, 58: 102281. DOI: 10.1016/j.jretconser.2020.102281.
9. Lei X, Tu G H, Liu A X, et al. The Insecurity of Home Digital Voice Assistants -- Amazon Alexa as a Case Study[J]. 2017.DOI:10.48550/arXiv.1712.03327.
10. Straga D. Facebook Messenger as a tool for building relationships with customers: bachelor thesis[J]. 2017.
11. A Berdasco, G López, I Diaz, et al. User Experience Comparison of Intelligent Personal Assistants: Alexa, Google Assistant, Siri and Cortana[J]. 2019.DOI:10.3390/proceedings2019031051.
12. Devlin J, Chang M W, Lee K, et al. Bert: Pre-training of deep bidirectional transformers for language understanding[J]. *arXiv preprint arXiv:1810.04805*, 2018.
13. Dehouche N. Plagiarism in the age of massive Generative Pre-trained Transformers (GPT-3): “The best time to act was yesterday. The next best time is now.”[J]. *Ethics in Science and Environmental Politics*, 2021.DOI:10.3354/esep00195.
14. Yang W, Xie Y, Lin A, et al. End-to-End Open-Domain Question Answering with Bertserini. [J]. *CoRR*,2019, abs/1902.01718.
15. Khamis M A, Ngo H Q, Christopher Ré, et al. FAQ: Questions Asked Frequently[J].*ACM*, 2016.DOI:10.1145/2902251.2902280.
16. A. I. Kadhim, “Term Weighting for Feature Extraction on Twitter: A Comparison Between BM25 and TF-IDF,” 2019 International Conference on Advanced Science and Engineering (ICOASE), Zakho - Duhok, Iraq, 2019, pp. 124-128, doi: 10.1109/ICOASE.2019.8723825.
17. Robertson S E, Walker S, Jones S, et al. Okapi at TREC-3[J]. *Nist Special Publication Sp*, 1995, 109: 109.
18. Yang P, Fang H, Lin J. Anserini: Reproducible Ranking Baselines Using Lucene[J]. *Journal of Data and Information Quality (JDIQ)*, 2018.DOI:10.1145/3239571.
19. Y Cui, W Che, T Liu, B Qin and Z Yang, “Pre-Training With Whole Word Masking for Chinese Bert,” in *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3504-3514, 2021, doi: 10.1109/TASLP.2021.3124365.
20. Cui Y, Liu T, Che W, et al. A span-extraction dataset for Chinese machine reading comprehension[J]. *arXiv preprint arXiv:1810.07366*, 2018.
21. Shao C C, Liu T, Lai Y, et al. DRCD: A Chinese machine reading comprehension dataset[J]. *arXiv preprint arXiv:1806.00920*, 2018.
22. Rajpurkar P, Zhang J, Lopyrev K, et al. Squad: 100,000+ questions for machine comprehension of text[J]. *arXiv preprint arXiv:1606.05250*, 2016.

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