

Chinese Technical Terminology Extraction by Using DC-value and Information Entropy

Zhang Liwei (✉ zhangliwei19810@126.com)

Capital University of Economics and Business

Research Article

Keywords: technical terminology extraction, domain C-value, information entropy

Posted Date: April 12th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1530516/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Chinese Technical Terminology Extraction by Using DC-value and Information Entropy**

2 Zhang Liwei

3 *School of Management and Engineering, Capital University of Economics and Business, 100070,*

4 *Beijing, China*

5 *E-mail: zhangliwei19810@126.com*

6
7 **Abstract**—China's technology is developing rapidly, and the number of patent applications has surged. Therefore, there is an urgent
8 need for technical managers and researchers that how to apply computer technology to conduct in-depth mining and analysis of lots
9 of Chinese patent documents to efficiently use patent information, perform technological innovation and avoid R&D risks. Automatic
10 term extraction is the basis of patent mining and analysis, but many existing approaches focus on extracting domain terms in English,
11 which are difficult to extend to Chinese due to the distinctions between Chinese and English languages. At the same time, some
12 common Chinese technical terminology extraction methods focus on the high-frequency characteristics, while technical domain
13 correlation characteristic and the unithood feature of terminology are given less attention. Aimed at these problems, this paper
14 proposes a Chinese technical terminology method based on DC-value and information entropy to achieve automatic extraction of
15 technical terminology in Chinese patents. The empirical results show that the presented algorithm can effectively extract the technical
16 terminology in Chinese patent literatures and has a better performance than the C-value method, the log-likelihood ratio method and
17 the mutual information method, which has theoretical significance and practical application value.

18 **Keywords**—technical terminology extraction; domain C-value; information entropy

19

20

21

22 **Introduction**

23 Terminology refers to a vocabulary unit describing the knowledge system of the professional domain, which
24 contains abundant professional domain knowledge¹. Terminology epitomizes and loads the core knowledge of a certain
25 technology domain, whose change reflects the development trend of the technology domain, to some extent².
26 Terminology plays an important role in aspects of machine translation, scientific writing, question answering systems,
27 automatic abstracting, knowledge communication, etc. Thus, many countries attach great importance to the construction
28 of terminology corpuses, such as the EURODICAUTOM of European Union, LEXIS of the Language Office of the
29 Federal Republic of Germany, TEAM of Siemens, TERMDOK of Sweden, DANETERM of Copenhagen Business
30 School, the TER MINUM terminology group of Canada, the ROSTERM terminology base of Russian, etc.³ Currently,
31 the terminology in many technology domains mainly rely on artificial construction⁴, which is not only time-consuming,
32 but also has a large cost⁵. Therefore, how to automatically extract terminology has been a concern for a long time.

33 Patent literature is the carrier of science and technology(S&T) information, recording the process of human S&T
34 development. As the world's largest technology information source, patents cover 90%-95% of the world's S&T
35 information⁶. Most of the new inventions, new technologies, new crafts, and new equipment of various countries in
36 various periods are reflected in patent literatures⁷. Currently, China's enormous economic market has attracted the
37 attention of relevant people, both domestically and abroad. At the same time, China's technology is developing rapidly,
38 and the number of patent applications has surged, with the number of patent applications in 2019 and 2020 continuously
39 ranked first in the world. Therefore, how to apply computer technology to conduct in-depth mining and analysis of
40 massive Chinese patent literature to make full use of patent information, perform technological innovation and avoid
41 R&D risks has attracted widespread attention. Automatic term extraction is the basis of patent mining and analysis, but
42 many existing approaches focus on extracting domain terms in English and are difficult to extend to Chinese due to the
43 distinctions between Chinese and English languages. At the same time, some common Chinese technical terminology
44 extraction methods focus on high-frequency characteristics, while technical domain correlation characteristics and the
45 unithood features of terminology receive less attention.

46 In response to the above problems, this paper takes Chinese patent literature as the research object and proposes a
47 method of extracting technical terms that combines grammatical rules and statistical methods to effectively identify
48 technical terms and improve the accuracy of term extraction. The remainder of this paper is organized as follows. In
49 Section 2 we describe existing work on automatic term extraction and focus on the challenges posed by domain-specific
50 and unithood characteristics. In Section 3, the difference between Chinese and English in the process of extracting
51 technical terms is analysed. In Section 4 we present some basic notions associated with terms and the features of patent
52 terms. We develop our proposed methodology for term extraction from Chinese patent literature in Section 5.
53 Experimental evaluations and performance comparisons are given in Section 6. Finally, Section 7 concludes the method
54 proposed in the paper and discusses the areas of future work.

55 **Theoretical Background**

56 Identifying and extracting domain terms from patent literature is a challenging task, which is mainly reflected in two
57 aspects: on one hand, the domain terms in the literature are very professional and rarely appear in the general thesaurus;
58 on the other hand, the phenomenon of term abbreviations, entity inclusion, and mutual reference in the literature are
59 very common, which puts forwards higher requirements for the correctness and completeness of term recognition.
60 Automatic term extraction methods can be summarized into several categories: rule-based methods, statistics-based
61 methods, machine learning-based methods, deep learning-based methods, semantic correlation-based methods,
62 graph-based methods, etc..

63 Rule-based term extraction methods mainly consider the context of the terms, the internal components of the terms

64 and other factors to identify terms, use grammatical rules, semantic rules, etc. to match in the corpus and output
65 multicharacter units that meet the established rules as terms. The common term extraction models mainly focus on
66 language features⁸⁻⁹, syntactic patterns¹⁰⁻¹², and retrieval strategies¹³. The advantages of the method include being
67 concise, intuitive, and having a strong expressive ability. The method can apply expert knowledge, and the accuracy is
68 high when the prior knowledge can match the text. However, this method usually requires an expert knowledge base as
69 a foundation, and whether building a knowledge base manually or automatically, it requires the intervention and
70 supervision of domain experts. At the same time, terms in different fields have different characteristics in terms of word
71 composition. To obtain a better extraction effect, the knowledge base must be continuously updated and adjusted. In
72 view of the shortcomings of this method, such as poor adaptability, excessive manual intervention, inability to identify
73 unknown words, etc., the application of this method has great limitations in terminology extraction.

74 The term extraction methods based on statistics apply various statistical models to measure whether a word string
75 is a term in the sense of probability. The term evaluation measures can be categorized as termhood features and unithood
76 features¹⁴. The main parameters used to compute the termhood and unithood of the candidate terms are frequency¹⁵,
77 TF*IDF¹⁶, C-value/NC-value¹⁷⁻¹⁸, DCFS(Domain Component Feature Set)¹⁹, hypothesis testing (z-test,t-test, chi-square
78 test, etc.)²⁰⁻²¹, likelihood ratio (LR)²²⁻²³[Error! Reference source not found.](#), information entropy²⁴⁻²⁵, mutual information (MI)²⁶⁻²⁷,
79 etc. The advantages of the methods are mainly manifested in the following aspects: they are easy to implement and
80 require less manual intervention; they are adaptable and can be used in different territories; and the unknown words can
81 be identified. The disadvantages are as follows: they are not sufficiently concise and intuitive; they are very dependent
82 on the corpus, and there must be a sufficient corpus to obtain a more ideal result; the accuracy rate is not high, because
83 many related words in the probabilistic sense are not terms; the low frequency terms cannot be identified; and due to the
84 need to perform many calculations, it is easy to cause operational efficiency problems.

85 The methods based on machine learning refer to the extraction of terms through machine training text features and
86 constructing models. This method can compensate for the shortcomings of other methods that cannot identify
87 low-frequency terms and use the data learning model to determine the possibility of whether the word string is a term.
88 Common machine learning methods include the maximum entropy model²⁸ and the conditional random field model²⁹⁻³¹.
89 However, the methods based on machine learning have high requirements on the scale and quality of the training corpus,
90 and a large-scale manual annotation corpus is required as the training data. Moreover, the methods are not yet mature,
91 and more attempts and verifications are needed. There is currently no targeted, complete, and large-scale annotated
92 corpus in patent literature.

93 The term extraction methods based on deep learning primarily combine the latest deep learning technologies to
94 automatically extract terminology. It is a special machine learning method for data representation that can solve the
95 problem of manually selecting the best feature engineering in the extracted terms. Related studies have applied the deep
96 learning method based on neural networks to term extraction; for example, combining Bi-LSTM³²⁻³⁴, CNN³⁵⁻³⁶, etc., to
97 conduct research in order to avoid manual feature extraction and other issues. However, the method highly relies on a
98 large-scale annotated corpus, and manual annotation of the corpus is time-consuming and labour-intensive.

99 Currently, some new methods have appeared in the field of automatic term extraction, such as the term extraction
100 methods based on semantic correlation, the extraction methods based on graphs, and so on. The extraction methods
101 based on semantic correlation mainly use the semantic relationship between phrases to improve the ranking of terms,
102 and thereby increase the accuracy of term extraction. Lahbib et al.³⁷ applied the idea of semantic correlation to the field
103 of bilingual term extraction, and extracted the source-end terms specific to the field. Astrakhantsev et al.³⁸ proposed the
104 KeyConceptRelatedness (KCR) method, which applied key concepts in the field to measure the quality of candidate
105 terms. Yu et al.³⁹ presented CBDLP, a data leakage prevention model based on confidential terms and their context terms.
106 The graph-based term extraction methods are inspired by the ranking method of web page importance in PageRank.
107 Mihalcea et al.⁴⁰ first applied PageRank to the field of natural language processing(NLP), and proposed a TextRank

108 method to extract key words. Semantic Graph-Based Concept Extraction (SGCCE), a novel concept extraction method
109 was proposed by Qiu et al.⁴¹. Khan et al.⁴² presented the Term Ranker method, constructed an undirected weighted
110 graph and improved the score of low-frequency terms.

111 In summary, related methods based on rules, statistics, machine learning, deep learning, etc. have all been used for
112 technical term extraction, and these methods have their own advantages and disadvantages. Based on the existing
113 research, this paper extracts the part-of-speech rules and grammatical rules of the terms in accordance with the strong
114 domain characteristics of patent terms and constructs a Chinese patent term extraction model based on DC-value and
115 information entropy theory.

116 **The Difference between Chinese and English in the Process of Extracting Technical Terms**

117 The biggest difference between Chinese and English is that in English, a "word" is used as the unit, where a single
118 word can express a precise meaning, while in Chinese, the unit is generally a "character", and current Chinese
119 emphasizes that "two-syllable words dominate". That is, it is difficult for each individual character to express a complete
120 meaning. At least two characters are combined to form a word that has an accurate meaning.

121 At the same time, each word in English is divided by "spaces". Therefore, when extracting English terms, it is easy
122 to extract individual words, but when extracting Chinese terms, it is difficult to express a complete meaning for each
123 individual character, so usually words composed of multiple characters are extracted. In addition, English belongs to
124 inflectional language, while Chinese is an isolated language. Thus, there are the following differences between English
125 and Chinese: ① There are relatively rich inflections in English, and the relationship between words is expressed
126 through inflections. ② An inflectional morpheme can express several different grammatical meanings in English. ③
127 The word order is strict in Chinese. Due to the lack of morphological changes in isolated words, there is no
128 morphological sign of what component a word belongs to in a sentence; it is completely determined according to the
129 word order. ④ Function words are very important in Chinese. The relationship between words in isolated languages is
130 often reflected by function words, an important grammatical means.

131 **Terminology and Patent Terminology**

132 *A. Basic Principles of Terminology Structure*

133 Terms are a type of language representation of concepts in a certain technology domain, and may be words, phrases,
134 letters or digital symbols. According to the structure of terms, they can be divided into simple terms and complex
135 terms⁴³. Among them, the simple terms are composed of only one word, for example, "communication" and
136 "information", while complex terms can be broken down into smaller units with an independent meaning; for example,
137 "communication apparatus" is made up of "communication" and "apparatus".

138 *B. Features of Patent Terminology*

139 Because patent literature belongs to S&T literature, the terms extracted from patent literature have general
140 characteristics of S&T terminology. The characteristics are roughly summarized as follows⁴⁴:

- 141 • Existing headwords. There are a few basic terms frequently appearing in a certain technology domain, which are
142 very important and may be headwords. Then you can find that in the domain, many complex terms consist of the
143 headwords in nominal structure or predicate structure. For example, in the password domain, a term that often
144 appears is the word "key", which could be seen as the headword, to constitute the nominal structure, such as
145 "session key", "master key", etc.; or the predicate structure, such as "key management", "key update", etc. Thus, a
146 large number of compound terms are formed. In this technology domain, the word "key" is a headword.
- 147 • Existing nested relationship among terms. Some complex terms are iteratively combined by simple terms, so there
148 is a nested relationship among terms. For example, the nested relationship among "symmetric cryptography
149 algorithms", "cryptography algorithms", and "algorithms" can be seen.

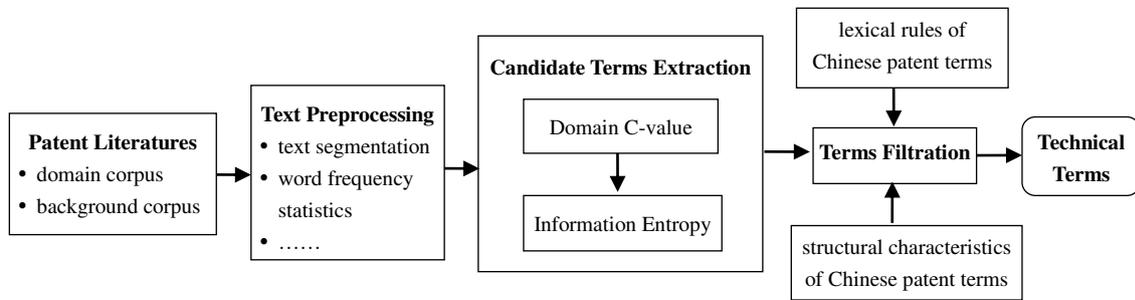
- 150 • Constituting connecting structure by symbols. Terms are composed of symbols ("/", "-", ".", "_", etc.), such as
- 151 "MH/Ni battery", "D-H key exchange protocol", etc.
- 152 • Combining English words with Chinese words to construct terms. Many terms are composed of both Chinese
- 153 words and English words together to form technical terms.
- 154 • Greater difference in length. There are not only existing terms with 2 or 3 characters, such as "电池" and "电动机",
- 155 but also existing terms with lengths greater than 6 or 10, such as "反应式步进电机" and "管式固体氧化物燃料电
- 156 池".
- 157 • Uneven distribution in different domain. Because of the great difference in technical content in different
- 158 technology domains, terms are closely related to technology domains, namely, the terms frequently appear in a
- 159 technology domain but rarely emerge in other technology domains.
- 160 Patents can be products, production methods, or technical schemes⁴⁵. In addition to the general characteristics of
- 161 S&T terminology, patent terminology also has its own uniqueness, which is roughly as follows:
- 162 • The vast majority of patent terminology expresses the specific entity of objects, components, and other objective
- 163 existences. This type of term must include nouns that act as headwords.
- 164 • There exist a few terms representing abstract concepts of crafts and methods. These terms are mainly composed of
- 165 verbs, and a few nouns, for example, "weld", "forge", etc.
- 166 • A term with more characters is, generally speaking, the object mainly described by the patent literature. The type
- 167 of terms represent the latest technology frontier and need to be given significant attention, such as "electronic
- 168 control gasoline injection engine", "plug-in series hybrid electric vehicle", etc.

169 Terminology Extraction Method Based on Domain C-value and Information Entropy

170 A. Framework of Terminology Extraction

171 According to the characteristics of patent literature, the framework of technical terminology extraction is

172 constructed, which is shown in Figure 1.



173 **Figure 1** Technical terminology extraction framework on Chinese patents

174 The terminology extraction system is mainly composed of three parts: the text preprocessing module, the candidate

175 terms extraction module and the terms filtration module.

176 B. Domain C-value (abbreviated as DC-value)

177 The C-value method is a type of hybrid terminology extraction method combining linguistic rules and

178 statistical theory¹⁷. The calculation formula of the C-value is shown in Eq. (1)⁴⁶:

$$179 \quad C\text{-value}(s) = \begin{cases} \log_2 |s| \times f(s) & s \text{ is not nested} \\ \log_2 |s| \times \left(f(s) - \frac{1}{c(b_i)} \sum_{i=1}^{c(s)} f(b_i) \right) & s \text{ is nested} \end{cases} \quad (1)$$

181 where s represents a candidate term, $|s|$ refers to the length of candidate term s , whose value is the number
 182 included by s ; $f(s)$ represents the appearance frequency of s ; b_i represents the candidate terms including s ; and
 183 $c(b_i)$ is the number of b_i .

184 However, the technical terms have the characteristics of domain correlation. The domain terms frequently
 185 appear or only appear in the texts belonging to a certain domain, while they rarely or never appear in other
 186 domains⁴⁷. Therefore, the C-value method is optimized in this paper with the introduction of a background
 187 corpus. Then the corpus is composed of two parts, the domain corpus and background corpus, based on which
 188 the concept of the word frequency distribution rate is proposed and the domain C-value is constructed for the
 189 preliminary extraction of the candidate terms.

190 1) *Word Frequency Distribution Rate*

191
$$df(s) = \frac{sf(s)}{bf(s)} \times 100\% \quad (2)$$

192 where s represents the candidate term; $sf(s)$ represents the frequency of s , appearing in the domain corpus;
 193 $bf(s)$ represents the frequency rate of s , appearing in the background corpus; and $df(s)$ represents the word
 194 frequency distribution rate of s .

195 2) *Domain C-value (DC-value)*

196 DC-value is set as Eq. (3)

197
$$DC - value = \begin{cases} \log_2 |s| \times \lg sf(s) \times \frac{sf(s)}{bf(s) + sf(s)} \times 100\% & s \text{ is not nested} \\ \log_2 |s| \times \frac{sf(s) - \frac{1}{sc(s)} \sum_{i=1}^{sc(s)} sf(b_i)}{bf(s) + sf(s)} \times \lg sf(s) & s \text{ is nested} \end{cases} \quad (3)$$

198 where s represents a candidate term; $|s|$ refers to the length of s ; $sf(s)$ represents the frequency of s
 199 appearing in the domain corpus; b_i represents the extracted candidate terms including s ; $sc(s)$ is the number of
 200 candidate terms including s in the domain corpus; $bf(s)$ represents the frequency of s appearing in the
 201 background corpus; and $\lg sf(s)$ is a high-frequency weighted coefficient, meaning that the more times the
 202 candidate term appears in the domain corpus, the larger the weight in the same ratio⁴⁸.

203 The extraction accuracy and performance of low-frequency words are effectively improved through the
 204 DC-value algorithm. However, the unithood feature is not considered. Aimed at this problem, the method of
 205 information entropy is introduced in subsequent research to ensure the integrity of the obtained terms.

206 C. *Information Entropy Method*

207 Information entropy in information theory represents the uncertainty of random variables. The more
 208 uncertain a random variable is, the larger its entropy value is. In the terminology extraction, the information
 209 entropy is mainly used to calculate the uncertainty of the boundaries of strings. The more uncertain the border of
 210 a string is, the larger the information entropy is. Then the string is more likely to be a complete term⁴⁹⁻⁵⁰.

211 The border uncertainty of strings is measured by computing the left and right information entropy of strings
 212 in this paper. For example, in the following paragraph "本发明提供一种转矩传感器以及动力转向装置。在具

213 有一对解算器的转矩传感器中，能够将上述两解算器的特性用作转矩传感器。”，the string “转矩传感器”
 214 has appeared a total of 3 times. Its left adjacent words successively are "种", "的" and "作", and its right
 215 adjacent words successively are "以", "中" and "。". In the entire corpus, the string "转矩传感器" appears a
 216 total of 27 times. The number of different left adjacent words amounts to 15, and the number of different right
 217 adjacent words is 19. It can be seen that the left and right adjacent words are not fixed. Therefore, it can be
 218 inferred that "转矩传感器" is likely to be a complete phrase, or even a term.

219 In the study of whether the phrase of "转矩传感" is complete or not, the phrase "转矩传感" appears 29
 220 times. The different left adjacent words are 19, while the right ones are only 2. Thus, "转矩传感" is not suitable
 221 to be a complete phrase.

222 Then, the formulas of the left and right information entropy are defined as follows:

$$LE(s) = -\sum_{l \in L} p(ls|s) \log_2 p(ls|s)$$

$$RE(s) = -\sum_{r \in R} p(sr|s) \log_2 p(sr|s)$$

223 where s is the candidate term, l is the left adjacent word of s , ls is the phrase composed of l and s , $p(ls|s)$ means
 224 the conditional probability that l is the left adjacent word of s in the case of the appearance of s , r is the right
 225 adjacent word of s , sr is the phrase consisting of s and r , and $p(sr|s)$ means the conditional probability that r is
 226 the right adjacent word of s in the case of the appearance of s . $LE(s)$ and $RE(s)$ represent the left and right
 227 information entropy of s , respectively. The smaller $LE(s)$ and $RE(s)$ and the more fixed the left and right
 228 adjacent words are, then the less likely it is that s is an independent phrase. To comprehensively evaluate the
 229 possibility of s standing alone as a phrase, the threshold values of the left and right information entropy are set
 230 to filter candidate strings that cannot stand alone as phrases. The setting of the threshold is shown in the
 231 formula:

$$232 \quad RE(s) \geq E_{\min} \quad \text{and} \quad LE(s) \geq E_{\min}$$

233 Where E_{\min} is the threshold value set manually⁵⁰.

234 D. Terminology filtration

235 In order to extract terms more fully and effectively, terminology filtering rules are set through a large
 236 amount of corpus analysis. The lexical rules are as follows:

- 237 • Location words, state words, interjections, and pronouns are not included in the terms;
- 238 • Terms should not begin with conjunctions, auxiliary words, or suffixes;
- 239 • Terms should not end with orientation words, auxiliary words, conjunctions, or prefixes;
- 240 • Nouns or verbs must be contained in terms;
- 241 • Adjectives and adverbs cannot stand alone as terms;
- 242 • Focus on filtering symbols (such as "-", ".", "_", "/", etc.);
- 243 • Focus on filtering the candidate terms containing English marks;
- 244 • The length of every term is less than 15;
- 245 • When a word does not appear in the stop word list and its part of speech is shown in Table 1, it needs to be
 246 filtered as a stop word.

247

Table 1 Part-of-speech tag table of special words

Tag	Description	Tag	Description
t	time word	q	quantifier
m	numeral	s	location word
r	pronoun	o	onomatopoeia
p	preposition	y	modal
d	adverb	z	state word
f	position word	c	conjunction

248

Note: Part-of-speech tagging uses the LTP (Language Technology Platform) part-of-speech tag set.

249 Experiment and Results

250 A. Datasets construction

251 In this paper, the public service platform of Shanghai intellectual property (<http://www.shanghaiip.cn/Search/login.do>) is applied as a patent retrieval database. The attributes of title, abstract, claims and
 252 international patent classification (IPC) are applied to retrieve the relevant patents, in which patents in the
 253 domain of information and communication are used as the domain dataset and patents in the domain of electric
 254 vehicles are used as the background dataset. We respectively selected 30,000 Chinese invention patents from the
 255 field of information and communication and the field of electric vehicles, where the retrieval time range was
 256 from 2010 to 2020. Then, a total of 60,000 items are used to construct a Chinese patent dataset. Among them,
 257 20,000 items are respectively taken from the domain dataset and background dataset, and a total of 40,000 items
 258 are used as the training set. 10,000 items are separately taken from the domain dataset and the background
 259 dataset, and a total of 20,000 items are used as the test set.

261 To generate the initial candidate terms, we used the corresponding analysis tools and relational
 262 corpus-Chinese lexical analysers LTP(Language Technology Platform) and ACE RDC 2005 (Automatic Content
 263 Extraction Relation Detection and Characterization) Chinese corpus to perform data preprocessing for the
 264 domain dataset and background dataset. LTP provides a series of Chinese natural language processing tools that
 265 can be used to perform word segmentation, part-of-speech tagging, and syntactic analysis of Chinese text. The
 266 ACE RDC 2005 Chinese corpus contains three fields information——newswire, broadcast conversations and
 267 newspaper, and includes 85,575 relation instances, in which there are 8,469 positive instances. In this paper, LTP
 268 is applied to segment sentences into words and assign each word a POS tagging; ACE RDC 2005 is then used to
 269 merge synonyms or similar words. After the data preprocessing work has been completed, 50,129 initial
 270 candidate terms are obtained.

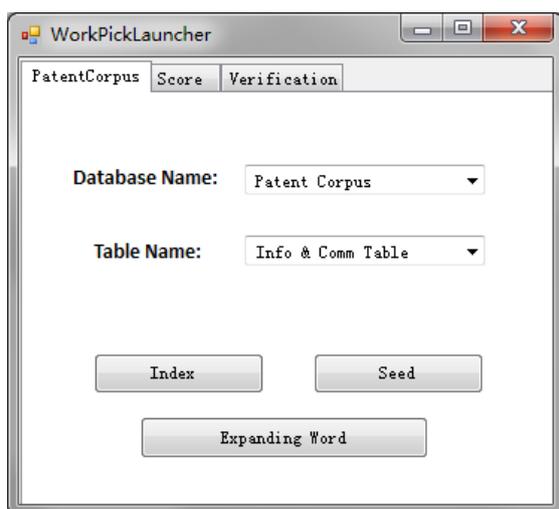
271 Table 2 shows the number of candidate terms after different selection and filtration steps. The results of
 272 each step are based on the results of the previous step.

273

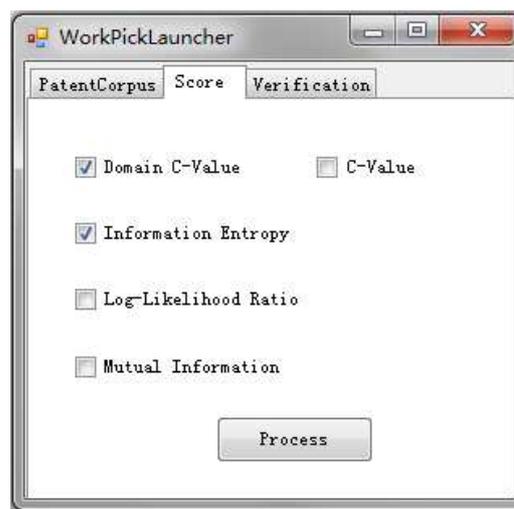
Table 2 The selection and filtration of candidate terms

Step	Term Candidates After Selection and Filtration	Term Candidates Numbers
1	Initial candidate terms	50129
2	Candidate terms after selection (DC-value+information entropy)	10782

274
 275 Finally, we successfully reduced the size of the candidate set from 50,129 to 3,921.
 276 *B. Experimental Results*
 277 For ease of application, a technical terminology extraction tool has been developed. The tool interfaces are
 278 shown in Figure 2 and Figure 3.



279
 280 **Figure 2 Patent corpus selection interface**



281 **Figure 3 Extraction algorithm selection interface**

282 The terms are then extracted by applying the extraction tools based on the methods of DC-value and
 283 information entropy algorithms. The results are shown in Table 3.

Table 3 Terminology extraction results

Candidate Terms	English Translation	Frequency	Word Segmentation	Part Of Speech	Terms?
多媒体子系统	multimedia subsystem	7	多媒体 + 子系统	n+n	Yes
光突发交换	optical Burst Switching	11	光 + 突发 + 交换	d+vi+v	Yes
光路交换	optical circuit switching	13	光 + 路 + 交换	d+n+v	Yes
光分组交换	optical packet switching	13	光 + 分组 + 交换	d+vd+v	Yes
偏振模色散补偿	polarization mode dispersion compensation	16	偏 + 振 + 模 + 色散 + 补偿	d+vg+ng+n+vn	Yes
链路	link	28	链 + 路	ng+n	Yes
媒体接入控制	media access control	17	媒体 + 接入 + 控制	n+vn+vn	Yes
突发光发射	burst mode transmitter	21	突 + 发光 + 发射	d+vi+v	Yes
突发光接收	burst mode receiver	21	突 + 发光 + 接收	d+vi+v	Yes
无线资源调度	radio resources scheduling	23	无线 + 资源 + 调度	b+n+vn	Yes
无线资源管理	radio resources management	23	无线 + 资源 + 管理	b+n+vn	Yes
正交频分复用	orthogonal frequency division multiplexing	21	正 + 交 + 频 + 分 + 复用	d+v+ag+v+vn	Yes
自动交换光网络	automatically switched optical network	25	自动 + 交换 + 光 + 网络	d+v+d+n	Yes
多粒度光交换	multi-granularity optical switching	23	多 + 粒度 + 光 + 交换	m+n+d+v	Yes
多用户	multiuser	25	多 + 用户	m+n	No
多粒度	multi-granularity	27	多 + 粒度	m+n	No
多粒度光	multi-granularity optical	25	多 + 粒度 + 光	m+n+n	No

284 C. Result Analysis

285 Generally, two indicators, P (precision) and R (recall rate), are used to evaluate the effect of the term
 286 extraction. However, in a corpus that has not all been manually tagged, it is difficult to determine the total
 287 number of terms it contains. Therefore, an alternative method is adopted, that is, P is expressed as a percentage
 288 of the number of terms correctly marked by the system to the number of terms extracted by the system; and R is
 289 expressed as the percentage of the number of terms correctly marked by the system to the number of manually
 290 tagged terms⁴⁵.

291 Among them, the number of manually tagged terms were obtained by extracting 175 documents according
 292 to each IPC subcategory in the domain corpus. Finally, a total of 2,625 documents were extracted, and a total of
 293 559 manually tagged terms were obtained.

294
$$P = \frac{\text{correct terms}}{\text{extracted terms}} \times 100\%$$

295
$$R = \frac{\text{correct terms}}{\text{tagged terms}} \times 100\%$$

296 To comprehensively evaluate the effect of the term extraction algorithm, the F-score evaluation index can
 297 be used, which is the harmonic mean of P and R, and the calculation formula is as follows⁸:

298
$$F\text{-Score} = \frac{2 \times P \times R}{P + R}$$

299 In this paper, 60,000 patent documents in the domain of information and communication and in the domain
 300 of electric vehicles were processed through the extraction algorithms based on DC-value and information
 301 entropy. According to the extraction results of technical terminology, the P, R and F-Score indicators are
 302 calculated.

303 To truly reflect the performance of the term extraction method based on the DC-value and information
 304 entropy proposed in this paper, several current mainstream term extraction methods are used for a comparison.
 305 These methods include the C-value, likelihood ratio, and mutual information methods.

306 Part of the contrastive result and the performance comparison between the method proposed in this paper
 307 and the other three methods for the extraction of technical terminology are shown in Table 4 and Table 5:

308 **Table 4 Technical terminology extraction results of four methods**

Candidate Terms	English Translation	Terms Or Not?			
		DC-value + Information Entropy	C-value	Log-Likelihood Ratio	Mutual Information
多媒体子系统	multimedia subsystem	Yes	Yes	Yes	Yes
光突发交换	optical Burst Switching	Yes	Yes	Yes	Yes
光路交换	optical circuit switching	Yes	Yes	Yes	Yes
光分组交换	optical packet switching	Yes	Yes	No	No
偏振模色散补偿	polarization mode dispersion compensation	No	No	No	No
链路	link	Yes	No	No	No
媒体接入控制	media access control	Yes	Yes	Yes	Yes
突发光发射	burst mode transmitter	Yes	Yes	No	Yes
突发光接收	burst mode receiver	Yes	Yes	No	Yes
无线资源调度	radio resources scheduling	Yes	Yes	Yes	Yes

无线资源管理	radio resources management	Yes	Yes	Yes	Yes
正交频分复用	orthogonal frequency division multiplexing	No	Yes	No	No
自动交换光网络	automatically switched optical network	Yes	Yes	Yes	Yes

309 Table 4 shows the extraction results of technical terminology by three different algorithms. Aimed at the
310 same candidate terms, the judging result may be different.

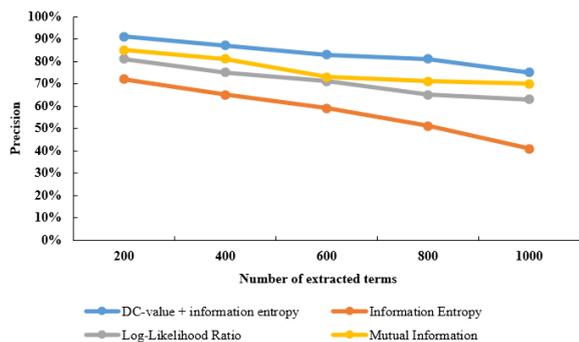
311

Table 5 Performance comparison among the methods

Term extraction method	Precision	Recall rate	F-Score
DC-value + information entropy	82.79%	85.51%	84.13%
Information Entropy	50.16%	31.97%	39.05%
Log-Likelihood Ratio	78.16%	81.32%	79.71%
Mutual Information	80.27%	79.30%	79.78%

312 Table 5 shows that the P, R and F-Score values of the terminology extraction algorithm based on the
313 DC-value and Information Entropy are 82.79%, 85.51% and 84.13%, respectively, which is significantly better
314 than the ones based on C-value, Log-likelihood estimation and mutual information methods. Therefore, the
315 validity of the algorithm proposed in the paper is verified.

316 At the same time, the experiment compared the results of the four methods when 200, 400, 600, 800, and
317 1000 terms were extracted. The experimental results show that as the number of extracted terms increases, the
318 precision is decreasing, the recall rate is increasing, and F-Score is also increasing. The precision and recall rate
319 of the first 1,000 extracted terms among the four methods are compared, as shown in Figure 4 and Figure 5. In
320 precision, the DC-value and information entropy method is 37%, 10% and 6% higher than the information
321 entropy, log-likelihood ratio and mutual information methods, respectively. In recall rate, the DC-value and
322 information entropy method is 49%, 7% and 11% higher than the information entropy, log-likelihood ratio and
323 mutual information methods, respectively.



324

Figure 4. Precision Comparison of Extraction Results

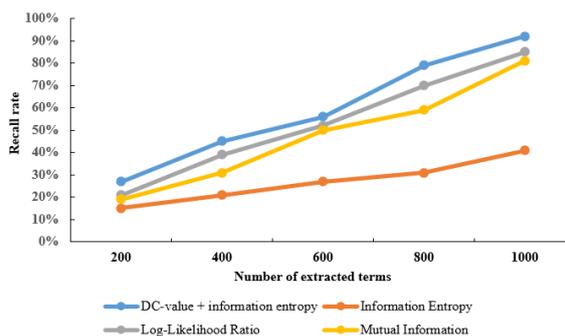


Figure 5. Recall Rate Comparison of Extraction Results

326 Through the analysis of experimental results, the method in this paper has been significantly improved
327 compared with other methods.

328 (1) Compared with the methods based on machine learning, the method in this paper does not require
329 high-quality training corpus, and need not to spend a lot of time for corpus training. At the same time, through
330 empirical studies in the fields of "biodegradable plastics", "carbon capture" and other fields, the effects are

331 similar to the above, verifying that the method is applicable to various professional fields.

332 (2) The extraction effect of combined terms and long terms is better. Due to the introduction of the
333 background corpus, the setting of nested terms, the discrimination of term boundaries, etc., term recognition is
334 more accurate. Through the analysis of the first 1,000 candidate terms extracted, the extraction ratio of terms
335 with 6 characters and above is higher than the extraction ratio of terms with less than 6 characters. For example,
336 the term “生物降解专用树脂(biodegradable special resin)” and the term "高强度导电聚乙烯醇(high-strength
337 conductive polyvinyl alcohol)" are both accurately extracted.

338 **Conclusion**

339 Automatic term extraction is an important issue in natural language processing, and is the basis of patent
340 mining and analysis. China currently attaches great importance to technological development, and Chinese
341 patent applications have surged. Many (S&T) managers and researchers in different organizations urgently need
342 conduct in-depth mining and analysis of massive Chinese patents in order to formulate accurate and effective
343 technology research and development strategies. However, many existing approaches focus on extracting the
344 domain terms in English and are difficult to extend to Chinese due to the distinctions between Chinese and
345 English languages. Therefore, this paper proposed a Chinese patent term extraction method based on DC-value
346 and information entropy to achieve automatic extraction of technical terms in Chinese patents.

347 Based on the traditional C-value method, this paper proposes the concept of the word frequency
348 distribution rate and constructs the DC-value method to measure the termhood of terms. According to the
349 characteristics of the terms, the relationship between terms and the context of terms is considered, and the left
350 and right information entropy are used to calculate the boundary uncertainty of the strings. Through the above
351 work, the selection of technical terms was completed according to the features of termhood and unithood. In
352 addition, through the analysis of the structural features and lexical rules of Chinese patent terms, the filtering of
353 technical terms was completed. The experiments showed that the method in this paper achieved better extraction
354 results.

355 To improve the speed and accuracy of the algorithm, in future work, we will introduce association rules
356 into the term extraction research to calculate the relevance of words, construct the relational structure of words
357 or phrases and obtain domain terms. By deeply exploring the technology of automatic machine learning
358 semantic relations between terms, the effectiveness and intelligence of term extraction can be improved.

359 **Acknowledgement**

360 This study was supported by the National Key R&D Program of China (2021YFE0101300), the National
361 Social Science Foundation (15CTQ031), Research on Patent Quality Evaluation Method (ZBBZK-2021-012),
362 Beijing Social Sciences Foundation (14JGC113), and Beijing Natural Science Foundation (9132004).

363 **Author contributions**

364 ZLW: conceptualization, methodology, software, model development, writing-original draft preparation, and
365 manuscript revision.

366 **Competing interests**

367 The author declare no competing interests.

368

369

REFERENCES

- 370 [1] Gu J, Wang H (2011) Study on term extraction on the basis of Chinese domain texts. *N Technol Lib Inf Serv* 4: 29-34(In Chinese).
- 371 [2] Wang Q, Li Y, Zhang P (2003) Automatic term extraction in the field of information technology. *Terml Std Inf Technol* 1:32-33,37 (In
- 372 Chinese).
- 373 [3] Liang A(2007) On the Development of terminological knowledge engineering. *Terml Std Inf Technol* 2:4-10,15. (In Chinese)
- 374 [4] Lin Y, Chen Z, Sun Q (2011) Computer domain term automatic extraction and hierarchical structure building. *Comput Eng*
- 375 37(2):172-174. (In Chinese)
- 376 [5] Han H, An X (2012) Chinese scientific and technical term extraction by using C-value and unithood measure. *Lib Inf Serv*
- 377 56:85-89(In Chinese) .
- 378 [6] Liu CY, Yang JC (2008) Decoding patent information using patent maps. *Data Sci J* 7: 14-22.
- 379 [7] Kisik S, Kyuwoong K, Sungjoo L (2018) Identifying promising technologies using patents: A retrospective feature analysis and a
- 380 prospective needs analysis on outlier patents. *Technol Forecast Soc Change* 128: 118-132.
- 381 [8] Fu J, Fan X, Mao J, Yu Z (2010) An algorithm of Chinese domain term extraction based on language feature. *Trans B Inst Technol* 30,:
- 382 307-310(In Chinese).
- 383 [9] Tatar S, Cicekli I (2011) Automatic rule learning exploiting morphological features for named entity recognition in Turkish. *J Inf Sci*
- 384 37(2): 137–151.
- 385 [10] Zhang X, Dai Y, Gao Z(2008) Applying syntactic patterns to semantic relation extraction from a terminology dictionary. *Eng Technol* 8:
- 386 43-45(In Chinese).
- 387 [11] Lee J, Yi JS, Son J (2019) Development of automatic-extraction model of poisonous clauses in international construction contracts
- 388 using rule-based NLP. *J Comput Civ Eng* 33: 04019003.
- 389 [12] Shao W, Hua B, Song L (2021) A pattern and POS auto-learning method for terminology extraction from scientific text. *Data Inf*
- 390 *Manag* 5(3): 329–335.
- 391 [13] Déjean H, Gaussier R, Sadat F (2020) Bilingual terminology extraction: An approach based on a multilingual thesaurus spplicable to
- 392 comparable corpora[EB/OL]. <http://www.xrce.xerox.com/content/download/23595/171307/file/dejean.pdf>.
- 393 [14] Kageura K, Umino B (1996) Methods of automatic term recognition: a review.*Terminology* 3: 259-289.
- 394 [15] Salton G, Buckley C (1988) Term-weighting approaches in automatic text retrieval. *Inform Process Manag* 24: 513-523.
- 395 [16] Lossio-Ventura JA, Jonquet C, Roche M, Teisseire M (2014) Biomedical terminology extraction: A new combination of statistical and
- 396 web mining approaches. In: *Proc. of the JADT*. pp:421–432.
- 397 [17] Frantzi K, Ananiadou S, Mima H (2000) Automatic recognition of multi-word terms: The c-value/nc-value method. *Int J Digit Lib* 3:
- 398 115-130.
- 399 [18] Astrakhantsev N (2015) Methods and software for terminology extraction from domain-specific text collection [Ph.D. Thesis]. Institute
- 400 for System Programming of Russian Academy of Sciences.
- 401 [19] Zhang QL, Sui ZF (2007) Measuring termhood in automatic terminology extraction. In *proceedings of the international conference on*
- 402 *natural language processing and knowledge engineering*, IEEE press, Piscataway, NJ. 328-335.
- 403 [20] Hua W, Zhang HY (2007) Extraction of Chinese term based on chi-square test. *Comput. Appl.* 27: 3019-3025(In Chinese).
- 404 [21] Montgomery DC, Runger GC (2018) *Applied Statistics and Probability for Engineers*. 7th, NJ: Wiley. pp:208–211.
- 405 [22] Dunning T. (1993) Accurate methods for the statistics of surprise and coincidence. *Comput Linguist*19: 61-74 .
- 406 [23] Verberne S, Sappeli M, Hiemstra D, Kraaij W (2016) Evaluation and analysis of term scoring methods for term extraction. *Inform*
- 407 *Retrieval J* 19:510-545.
- 408 [24] Dong YY, Li WH, Hu H (2017) Domain term extraction method based on hierarchical combination strategy for Chinese Web
- 409 documents. *J Northwest Polytechnical Univ* 35(4): 729–735 (in Chinese).
- 410 [25] Li L(2013) *The Research of Term and Relation Acquisition Methods for Domain Ontology Learning*, Ph.D. Dissertation. Dalian Univ
- 411 *Technol* pp:63-69 (In Chinese).
- 412 [26] Church KW, Hanks P (1990) Word association norms, mutual information, and lexicography. *Comput Linguist* 16: 22-29.
- 413 [27] Zeng W, Xu S, Zhang Y, Zhai J (2014)*The Research and Analysis on Automatic Extraction of Science and Technology Literature*
- 414 *Terms Lib Inform Technol* 1:51-55(In Chinese).
- 415 [28] Muheyat-N, Kunsaul-T (2016) Research on Automatic Identification of IT Terms in Kazakh. *J China Inform Process* 30: 68-73(In
- 416 Chinese).
- 417 [29] Mozharova VA, Loukachevitch NV (2016) Combining knowledge and CRF-based approach to named entity recognition in
- 418 Russian[C]/*International conference on analysis of images, Social Networks and Texts*. Springer.pp:185-195.
- 419 [30] Wang H,Wang M, Su X(2016) A study on chinese patent terms extraction for ontology learning. *J China Soc Sci Tech Inform.* 35:

- 420 573-585 (In Chinese).
- 421 [31] Zeng D, Sun C, Lin L, Liu B (2017) LSTM- CRF for drug-named entity recognition. *Entropy* 19: 283-295.
- 422 [32] Liu Y, Yin L, Zhang K (2019) Deep transfer learning for technical term extraction-A case study in computer numerical control system.
423 *J Intell* 38: 168-175.
- 424 [33] Miwa M, Bansal M (2016) End-to-End Relation Extraction using LSTMs on Sequences and Tree Structures. In: Proceedings of the 54th
425 annual meeting of the association for computational linguistics. Association for computational linguistics.
- 426 [34] Geng Z, Chen G, Han Y, Lu G, Li F (2020) Semantic relation extraction using sequential and tree-structured LSTM with attention.
427 *Inform Sci* 509: 183–192.
- 428 [35] Khosla K, Jones R, Bowman N (2019) Featureless deep learning methods for automated key-term extraction. Stanford: Stanford
429 University.
- 430 [36] Xu D, Tian Z, Lai R, Kong X, Tan Z, Shi W (2020) Deep learning based emotion analysis of microblog texts. *Inform Fusion* 64: 1-11.
- 431 [37] Lahbib W, Bounhas I, Slimani Y (2018) A possibilistic approach for Arabic domain terminology extraction and translation. In: Proc. of
432 the Int'l Symp. on Computer and Information Sciences. Cham: Springer-Verlag, pp:231–238.
- 433 [38] Astrakhantsev N (2014) Automatic term acquisition from domain-specific text collection by using Wikipedia. Proceedings of the
434 Institute for System Programming of RAS. 26(4): 7-20.
- 435 [39] Yu X, Tian Z, Qiu J, Jiang F (2018) A data leakage prevention method based on the reduction of confidential and context terms for
436 smart mobile devices. *Wirel Commun Mob Com* pp: 1-11.
- 437 [40] Mihalcea R, Tarau P (2004) TextRank: Bringing order into text. In: Proc. of the EMNLP. Stroudsburg: ACL. pp:404–411.
- 438 [41] Qiu J, Chai Y, Tian Z, Du X, Guizani M (2019) Automatic concept extraction based on semantic graphs from big data in smart city.
439 *IEEE T Comput Soc Sys* 7(1): 225-233.
- 440 [42] Khan M T, Ma Y, Kim J (2016) Term ranker: A graph-based re-ranking approach. In: Proc. of the 29th Int'l Florida Artificial
441 Intelligence Research Society Conf. Florida: AAAI Press, pp: 310–315.
- 442 [43] Wang H, Li G (2014) Research of automatic term wxtraction based on association rules. *Lib Inform* 5: 20-25 (In Chinese).
- 443 [44] Chen S, Yu B (2013) Model of automatic term extraction for technology domain. *Systems. Eng Theory Prac* 33: 230-235.
- 444 [45] State Intellectual Property Office of The P.R.C. (2015) What is a invention patent [EB/OL].[2015-05-15].[http://www.sipo.
445 gov.cn/zsjz/cjw/200804/t20080402_367771.html](http://www.sipo.gov.cn/zsjz/cjw/200804/t20080402_367771.html).
- 446 [46] Ji P, Yan XY, Cen Y (2010) A survey of term recognition and extraction for domain-specific Chinese text information processing. *Lib
447 Inf Serv* 54: 124-129 (In Chinese).
- 448 [47] Zhu Q, Leng F (2012) Existing Problems and Developing Trends of Automatic Term Recognition. *Lib. Inf. Serv.* 56:104-109(In
449 Chinese)
- 450 [48] Xiong L, Tan L, Zhong M (2013) An automatic term extraction system of improved C-value based on effective word frequency. *N.
451 Technol. Lib. Inf. Serv.* 9: 54-59(In Chinese).
- 452 [49] Patry A, Langlais P (2005) Corpus-based terminology extraction. In 7th International Conference on Terminology and Knowledge
453 Engineering Copenhagen, Denmark, pp:313-321.
- 454 [50] Li L (2013) The Research of term and relation acquisition methods for domain ontology learning, Ph.D. Dissertation, Dalian Univ.
455 Technol, pp:63-69 (In Chinese).