# RESEARCH ON SEMANTIZATION OF UNSTRUCTURED TEXT IN MEDICAL FIELD

Yao Liu[1], Zhifang Sui[2] and Yi Huang[1]

[1]Institute of Scientific and Technical Information of China
No. 15, Fuxing Road, Haidian District, Beijing 100038, P. R. China
liuy@istic.ac.cn

[2]Institute of Computational Linguistics
Peking University
No. 5, Yiheyuan Road, Haidian District, Beijing 100871, P. R. China

ABSTRACT. *In recent years, with the rapid development of computer technology and the popularity of the Internet, huge information began to emerge in the form of unstructured text in various fields. Faced with such a huge and rapidly growing amount of data, how to quickly and effectively obtain specific knowledge of the real needs out from the mass is a challenge faced by researchers. Researches in medical field are comparatively few than other fields. Chinese medical text has its own unique characteristics, and we cannot copy the foreign Natural Language Processing methods and technologies. This paper presents a method of processing medical unstructured text. The application results show that the proposed methods can effectively process semantization of medical unstructured text.*
**Keywords:** Unstructured text, NLP, Semantic role labeling, Discourse annotation

1. **Introduction.** In recent years, with the rapid development of computer technology and the popularity of the Internet, lots of information began to emerge in the form of unstructured text in various fields. As one of the top research fields, medical research delivered a huge amount of data online. More than 600,000 papers are published in medical field every year. There are more than 20 million papers in MEDLINE [1]. With serious challenges brought by information explosion, we are in urgent need of some technologies and tools to help us quickly find valuable information out of the mass. Nowadays, semantization is becoming area of research focus, being widely used around the globe. In China, Ou [2] investigated the semantic representation of Chinese Thesaurus. Liu et al. [3] proposed a method of using semantic annotation and ontology to conduct semantization of library resources. Bai and Qiao [4] studied the semantization of bibliography data based on ontology and linked data. However, there are relatively few studies in Chinese medical literature than those in other fields, for the text is often with longer sentences and more complex structures. Also, we cannot copy methods and technologies in medical field used in foreign languages [5]. Therefore, the need to work out a systematic method of processing Chinese medical unstructured text is demanding. Based on the characteristics of unstructured Chinese medical text, this paper proposes a method of processing unstructured Chinese medical text by conducting segmentation, POS tagging, semantic annotation, semantic role labeling, and discourse annotation to achieve semantization of unstructured text in medical field.

The rest of this paper is organized as follows. Section 2 introduces idea and framework. Section 3 describes the key methods and technologies. Section 4 describes the applications of the proposed methods and technologies. Finally, a brief conclusion and future work are given in Section 5.

2. **Framework.** From a linguistic point of view, natural language understanding is a hierarchical process that includes the following five levels: speech analysis, lexical analysis, syntactic analysis, semantic analysis and pragmatic analysis. Among them, semantic information is the core part of the process, and it is one of the main objectives in the field of natural language processing research as well. Based on the characteristics of unstructured Chinese medical text, this paper proposes a systematic method of processing unstructured Chinese medical text by segmentation, POS tagging, semantic annotation, semantic role labeling, and discourse annotation to achieve semantization of unstructured text in medical field.

3. **Key Methods and Technologies.**

3.1. **Word segmentation & POS tagging.** Unlike English, there is no space to separate Chinese words in a sentence. Segmentation is the very first step to process any Chinese text. Since there are many long terms in medical field, we constructed a professional field dictionary containing about 40 thousand terms with the help of Online Semantic Dictionary Construction Platform (OSDCP) to help cut sentences into words. Then, word segmentation and POS tagging are conducted with WSPTT shown as Figure 1.



FIGURE 1. GUI of WSPTT and its output results

3.2. **Semantic annotation.** In order to make computer understand the unstructured medical text at the content level, we need to process the text fragments at a semantic level. Semantic annotation cannot be effectively conducted without domain ontologies, because in a domain ontology there are many concepts as well as relations between them, which can be made good use of in the annotating process. By using Online Ontology Construction Platform (OOCP), we construct a medical ontology. It is assumed that there are text fragments related to a concept in those paragraphs. Semantic score is being used to measure how closely a sentence or paragraph is related to a concept in the domain ontology. Then, based on the occurrence of a concept and its properties in a sentence or paragraph, semantic annotation is performed [6].
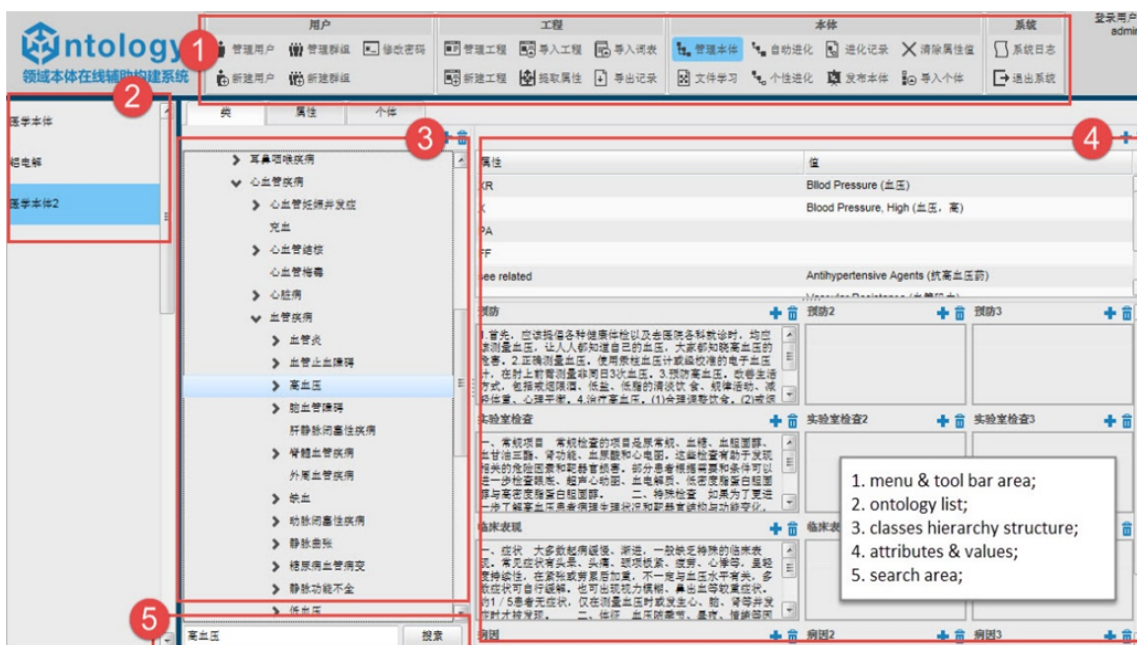
FIGURE 2. GUI of OOCP and its results

3.3. **Semantic role labeling.** Semantic role labeling is another important factor in the semantization process. Most existing Chinese semantic role labeling systems are generally dependent on complete syntactic analysis. Since the performance of Chinese parser is still relatively low, especially in medical field, the result of automated semantic role labeling is unsatisfactory [7]. In view of this, shallow syntactic parsing is introduced instead of complete syntactic analysis as the input for Chinese semantic role labeling. In order to fully reflect the characteristics of the Chinese language with fine-grained description of the words in particular, we refer to the "*Modern Chinese Lexicon Structure* [8]", to extract feature characters from verbs and nouns as well as relevant features of word formation.

Character set in Chinese is relatively closed collection with long-term stability. Word formation in Chinese is much like syntactic structure. With the information, it is much easier to find relations between words and new word discovering. For example, in dealing with a new word, the POS and semantic role of word in a sentence can be inferred from the preceding word and back word. Based on the analysis above, the features having a positive effect in syntactic analysis and semantic role labeling include: center morpheme, morpheme type, word formation, and sequence information. In the parsing stage, in order to facilitate a fairer semantic role labeling, there is no introduction of external resources, and we use rules to generate "pseudo center morpheme" as a feature to provide the basis for parsing [9]. Chunk definition [10] is adopted, and a total of 12 types of chunk are chosen.

The following is a detailed description of the method:

(1) Input: correct word segmentation and POS tagging result;

(2) Shallow syntactic parsing: adopting chunk definition, combined with IOB2 notation, the shallow parsing problem is transformed into a sequence labeling problem;

(3) Semantic role labeling: based on the identification sequence of arguments and arguments classification, SRL is conducted.

3.4. **Discourse annotation.** Discourse annotation aims to reveal the deep relations beyond the sentence level, and present the inner relations of the complicated discourses through normalized annotation process [11]. Taking consideration of the Chinese medical characteristics, we propose an easily operated discourse annotation process based on the traditional discourse theories and their corresponding corpus. We describe the deep

TABLE 1. Types of chunk in SRL

| tag | meaning |
| --- | --- |
| ADJP | Adjective phrase |
| ADVP | Adverb phrase |
| CLP | CLP Classifier |
| DNP | Deg "的" phrase |
| DP | Determiner phrase |
| DNP | Deg "得" phrase |
| LOC | Location phrase |
| LST | List marker |
| NP | Noun phrase |
| PP | Prepositional phrase |
| QP | Quantifier phrase |
| VP | Verb phrase |



FIGURE 3. GUI of discourse annotation platform

structure and meaning of the discourse in three aspects: content, relation and reference. Content tagging is to mark the fragments with clear meaning. General and domain specific tags are introduced to ensure the flexibility of the content tagging: general tags are regarded as the default tags, and domain tags can be expanded as needed within the medical fields. Relation annotation is to describe the logical relation between the adjacent components in a discourse. Reference tagging is to mark pronouns, nouns, etc. By cutting the medical unstructured text into paragraphs, sentences and clauses, the annotating process adopts a bottom-up way, from small to large linguistic unit, to tag the relations, and marks the core and subsidiary ones, until the formation of a discourse annotation tree is comparatively complete [12].

4. **Application.** QA system allows users to use casual language to organize questions, unlike traditional search engines that require users to input the keywords in a rather compelled way. With the information of semantic annotation, semantic role labeling, and

discourse annotation, we are able to match query words to a deep semantic structure and relation, so as to obtain a much more precise answer to present to users.

"*Dermatology and Venereology (Fifth Edition)*" is chosen as raw material of the corpus. With all the tools and platforms above, the corpus is being processed accordingly. From the Good Doctor Dermatology Online and Baidu Knowledge Dermatology Category, we select 100 typical questions in medical field. In this QA system application, the results are shown as Figures 4 and 5.

As can be seen from above, when a user inputs a question in the search box, the system analyzes the query text, decomposes the text into triples, conducts a search in the database, and presents the user an interactive dialog with choices to help further filter the results so as to provide the precise answer to the question. By clicking word links relevant to the question in the result page, a user can also browse related information in
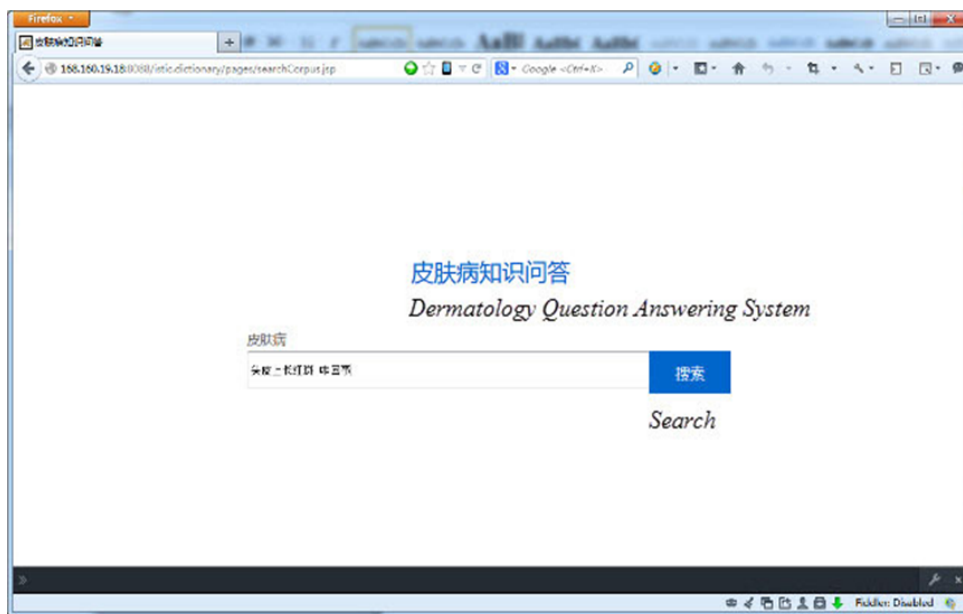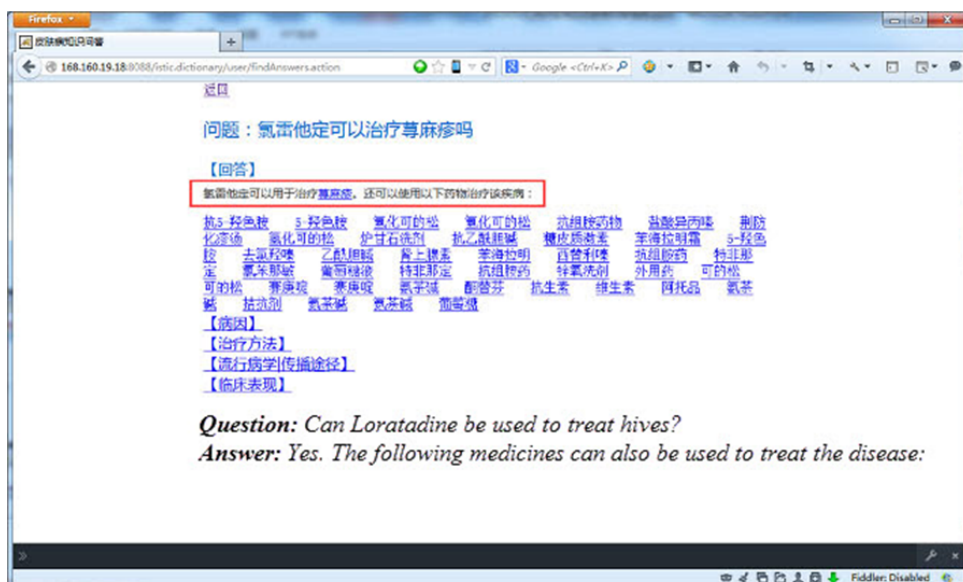


FIGURE 4. Homepage of QA system



FIGURE 5. Query result of QA system

his interests. It can be concluded that our methods can effectively process semantization of medical unstructured text.

5. **Conclusions.** Based on the characteristics of Chinese medicine unstructured text, the paper proposes to process semantization of unstructured text in medical field by segmentation, POS tagging, semantic annotation, semantic role labeling, and discourse annotation to achieve semantization of unstructured text in medical field. The experimental results show that the proposed methods can effectively process semantization of medical unstructured text. In the future we intend to proceed along two lines in parallel: on one hand, to broaden the scope of semantization by adopting some other technologies and tools; on the other hand, to improve some key algorism to meet the needs of semantization in other fields other than medical field.

## REFERENCES

[1] V. Nebot and R. Berlanga, Exploiting semantic annotations for open information extraction: An experience in the biomedical domain, *Knowledge and Information Systems*, vol.38, no.2, pp.365-389, 2014.

[2] S. Ou, Semantic representation of Chinese thesauri, *Library and Information Service*, vol.59, no.16, pp.110-118, 2015.

[3] Y. Liu, X.-F. Chen, Z. Sui, Y. Hu and Q. Zhao, Research on semantic method of library resources' organizing, *ICIC Express Letters*, vol.5, no.4(A), pp.1011-1017, 2011.

[4] H. Bai and X. Qiao, Study of semantic bibliography base on ontology and linked data, *New Technology of Library and Information Service*, vol.26, no.9, pp.18-27, 2010.

[5] J. Zhang and Y. Liu, Medical literature's semantic status and solution, *International Journal of Knowledge and Language Processing*, 2013.

[6] Y. Liu, H. Shi, D. Zheng and Y. Huang, Study on semantic annotation for professional literature, *ICIC Express Letters, Part B: Applications*, vol.5, no.5, pp.1383-1389, 2014.

[7] L. Qiu, L. Wu, K. Zhao and C. Hu, Improving Chinese dependency parsing with auto-extracted dependency triples, *Int. J. of Asian Lang. Proc.*, vol.22, no.2, pp.75-84, 2012.

[8] S. Yu, H. Duan, X. Zhu, B. Swen and B. Chang, Specification for corpus processing at Peking University: Word segmentation, POS tagging and phonetic notation, *Journal of Chinese Language and Computing*, vol.13, no.2, pp.121-158, 2003.

[9] X. Wang, *The Research and Implementation on Chinese Semantic Role Labeling Based on Lightweight Syntactic Information (in Chinese)*, Ph.D. Thesis, Peking University, 2012.

[10] W. Chen, Y. Zhang and H. Isahara, An empirical study of Chinese chunking, *Proc. of the COLING/ACL on Main Conference Poster Sessions*, 2006.

[11] Y. Zhou and N. Xue, PDTB-style discourse annotation of Chinese text, *Proc. of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers*, vol.1, 2012.

[12] Y. Wang, *The Development and Application of Discourse Annotation Platform*, Ph.D. Thesis, Peking University, 2014.