# Research on Open Domain Question Answering System

Zhonglin Ye[1], Zheng Jia[1(✉)], Yan Yang[1], Junfu Huang[1], and Hongfeng Yin[2]

[1] School of Information Science and Technology,
Southwest Jiaotong University, Chengdu, China
`zhonglin_ye@foxmail.com`, `{zjia,yyang}@swjtu.edu.cn`,
`990504422@qq.com`
[2] DOCOMO Innovations Incorporation, Palo Alto, USA
`hongfeng_yin@yahoo.com`

**Abstract.** Aiming at open domain question answering system evaluation task in the fourth CCF Natural Language Processing and Chinese Computing Conference (NLPCC2015), a solution of automatic question answering which can answer natural language questions is proposed. Firstly, SPE (Subject Predicate Extraction) algorithm is presented to find answers from the knowledge base, and then WKE (Web Knowledge Extraction) algorithm is used to extract answers from search engine query result. Experimental data provided in the evaluation task includes the knowledge base and questions in natural language. The evaluation result shows that MRR is 0.5670, accuracy is 0.5700, and average F1 is 0.5240, and indicates the proposed method is feasible in open domain question answering system.

**Keywords:** Automatic question answering · Open domain · Natural language understanding · Information extraction

## 1    Introduction

Google, Baidu, Bing and other search engines return hyperlinks containing the keywords of user queries, which do not give users a simple and direct answer, and users need to browse web pages to find answers they need. Although search engines can help users find answers to a certain extent, but users may need to click many hyperlinks of pages. Driven by the mobile Internet, the automatic question answering system based on domain knowledge base can directly get an intuitive and accurate answer, so it becomes an important focus of research.

This paper presents a solution of open domain question answering (QA). The method mainly includes three parts: (1) SPE (Subject Predicate Extraction) algorithm; (2) WKE (Web Knowledge Extraction) algorithm; (3) answer format standardization. SPE algorithm extracts subject and predicate in questions, searching answers from the knowledge base according to the subject and predicate of questions. SPE solves the problem of QA in limited domain, such as People, Time and Geography. For other

domains, WKE algorithm is used to extract answers from the unstructured texts of search engine results, and can effectively solves the problem of QA in open domain. In order to improve the accuracy of question answering system, answers are turned into a standardized format. The knowledge base is from Baidubaike, and the quantity of it is about 4 million triples which are in the format of <Subject, Predicate, Object>. Testing data are 1000 natural language questions which are from Microsoft's Bing query log and generated from the knowledge base.

The rest chapters of this paper are arranged as follows: the second chapter is about the related work. The third chapter introduces the proposed method. The fourth chapter is the experiment. The last chapter concludes the whole paper and looks forward to the future research.

## 2      Related Work

QA is to provide a quick, direct and accurate answer for the natural language questions [1]. The existing QA technologies can be divided into three categories according to the source of answers: (1) QA based on search engine; (2) QA based on community, (3) QA based on knowledge base [1]. In QA based on search engine, the answers are extracted from the web pages of search result [2]. QA based on community calculates the similarity between the questions raised by users and questions stored in database asked in the past to get results [3]. The main works of QA based on knowledge base is semantic analysis [4,5,6,7,8,9] and knowledge base building [10], such as Poon [11], Yahya [12] and Berant [13,14,15] put forward a method to build the QA system based on semantic analysis. First, they extracted the subject entities and predicate words of questions. Then the questions are converted into SPARQL structured query language. Finally, answers are acquired through retrieving knowledge base by SPARQL. Bordes [16] and Yao [17] proposed a way of information retrieval to build QA system. Fader [18] and Kushman [19] presented a method of building QA system based on open information extraction. Bao [20] and Comas [21] raised an approach of constructing QA system based on translating question into answer on the basis of knowledge base.

## 3      Methods

The procedure of open domain QA is shown in Figure 1 as follows:

As shown in Figure 1, the system is mainly composed of four parts: question classification, SPE algorithm, and WKE algorithm and answer format standardization. Firstly, natural language questions are classified. Then, SPE algorithm is applied to extract subject and predicate of the questions, and retrieve the answers from the knowledge base. If there are no answers, WKE algorithm is applied to search the questions in search engine to get the web pages containing the questions, to analyze the texts of search results, and extract answers from the texts. Finally, the answers are turned into a standard format.
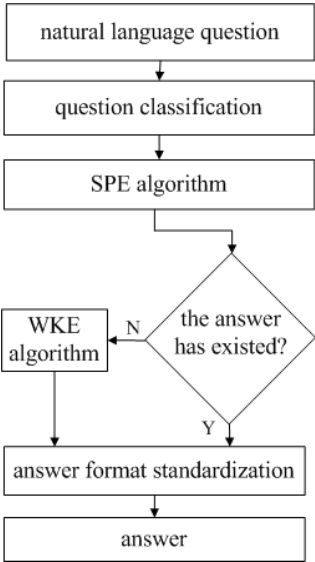
**Fig. 1.** Procedure of open domain question answering

## 3.1    Question Classification

The questions are divided into four categories: 人物|People[1], 地点|Geography, 时间|Time and Others. In this paper, the category of questions is got through feature words matching. Some of feature words for each category are described in Table 1.

**Table 1.** Feature word examples

| Category | Feature words |
|---|---|
| 人物|People | 谁|who，名字|name，姓名|name，中文名|Chinese name，英文名|English name，别名|alias，叫什么|what's called |
| 地点|Geography | 在哪|where，哪里|where，地方|address，位置|position |
| 时间|Time | 哪一年|which year，什么时候|when，时间|time，诞生|born，时期|date，哪年|which year，何时|when |
| Other | 无|null |

If any question contains the feature words, the category of the question can be acquired. The category has two functions: one is used to identify the predicate of the question in SPE algorithm, and the other is to extract answers in WKE algorithm.

---

[1] The content after "|" is the English translation of Chinese words.

### 3.2    SPE Algorithm

The knowledge triples consist of Subject, Predicate, and Object. The questions usually contain the words describing Subject and Predicate but the Object is absent. The goal of SPE algorithm is to extract the words describing Subject and Predicate and to search the knowledge base to find the Object. The procedure of SPE algorithm is shown in Figure 2:
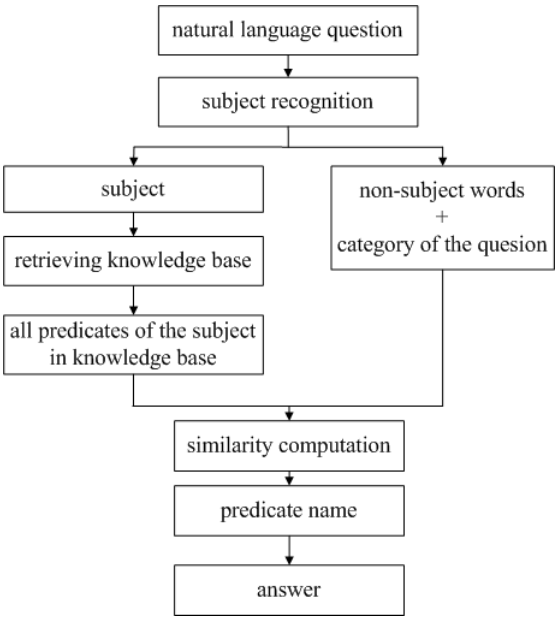


**Fig. 2.** SPE algorithm description

Illustrated as Figure 2, the steps of SPE algorithm are shown as below.

(1) Identify the subject entity name based on CRF (Conditional Random Field) algorithm and SWJTU Chinese Word Segmentation System.

(2) If the question contains a subject, then retrieve knowledge base to find all triples of the subject.

(3) The words except the subject which are called non-subject words are combined with category of the question.

(4) Compute the similarity [22] between each predicate of the subject and the combination of non-subject words and categories. Select predicate name which similarity value is the highest.

(5) According to the subject and predicate, the object can be got from triples in the knowledge base and the object is the answer of the question.

For example, the question is 侏罗纪世界什么时候上映|When will Jurassic World release.

The feature word is 什么时候|when, so the category of the question is 时间|Time. The subject is recognized as 侏罗纪世界|Jurassic World based on CRF. The non-subject word is上映|release. The combination of category and non-subject word is 上映时间|release time. There is a predicate named "上映时间|release time" in the triple <侏罗纪世界|Jurassic World, 上映时间|release time, 2015-06-12> of the subject 侏罗纪世界|Jurassic World, so the object "2015-06-12" can be found as the answer of the question. The detailed descriptions of main steps are as follows.

### 3.2.1    Subject Recognition

We firstly use SWJTU Chinese Word Segmentation System [23] to recognize the entities in questions. This system can achieve good performance on recognizing several categories of entities such as Person, Place, and Organization etc. However it cannot recognize many other categories such as Film and Music etc. To realize good performance of entity recognition, then we use the CRF algorithm [24] to recognize the entities of Film, Music, Book, Game and Application categories.

CRF named entity recognition algorithm requires the establishment of a training model which is used to predict the new entities. Because the better the quality of training data is, the higher performance the algorithm becomes, it needs to collect a lot of high-quality questions manually labeled as training data. We collected a large number of questions through crawling on the Internet or writing by hand, and labeled these questions manually. The training data are got mainly from the following resources.

(1) Web. We crawled the data from Baidu Knows, Douban, 360 search engine and so on.
(2) Question collection online system. We developed a question collection online system and about more than 300 students participated in writing and labeling questions.
(3) Mobile assistant voice data..We bought the mobile assistant voice data from Datatang Corporation. The number of the questions is 145371 and we labeled the data manually.

The total number of the questions is about more than 30000. The average F-measure value of recognizing subject named entities for the five categories based on CRF is 92.44%.

### 3.2.2    Retrieve Knowledge Base

Once the subject entity in questions is recognized, the triples in knowledge base containing this subject can be acquired by retrieving knowledge base. A predicate dictionary of the subject can be created according to the triples.

For example, the question is 梁启超的生日是什么时候|When is Qichao Liang's birthday. The subject is recognized as a person梁启超|Qichao Liang. By retrieving the knowledge base, we can get the knowledge triples as shown in Table 2.

The predicate dictionary of Qichao Liang is created from triples as follows:

[中文名|Chinese name, 外文名|English name, 别名|alias, 国籍|nationality, 民族|volk, 出生地|birthplace, 出生日期|birthdate, 逝世日期|death time, 职业|career, 信仰|faith, 主要成就|achievement, 代表作品| representative works].

Table 2. The triple examples of subject "梁启超|Qichao Liang"

| No. | <Subject, Predicate, Object> | No. | <Subject, Predicate, Object> |
|---|---|---|---|
| 1 | <梁启超|Qichao Liang，中文名|Chinese name，梁启超|Qichao Liang > | 7 | <梁启超|Qichao Liang，出生日期|birthdate，1873年2月23日|February 23, 1873 > |
| 2 | <梁启超|Qichao Liang，外文名|English name，Qichao Liang > | 8 | <梁启超|Qichao Liang，逝世日期|death time，1929年1月19日|January 19, 1929 > |
| 3 | <梁启超|Qichao Liang，别名|alias，卓如|Zhuo Ru > | 9 | <梁启超|Qichao Liang，职业|career，思想家|thinker> |
| 4 | <梁启超|Qichao Liang，国籍|nationality，中国|China> | 10 | <梁启超|Qichao Liang，信仰|faith，儒学|Confucianism > |
| 5 | <梁启超|Qichao Liang，民族|volk，汉族|Han> | 11 | <梁启超|Qichao Liang，主要成就|achievement，公车上书| Gong Che Shang Shu 、戊戌变法|Wu Xu Reformation> |
| 6 | <梁启超|Qichao Liang，出生地|birthplace，广东省新会市|Xinhui city, Guangdong province> | 12 | <梁启超|Qichao Liang，代表作品| representative works，《中国近三百年学术史》|《scholastic history of the past 300 years in China》、《中国历史研究法》|《historiography of Chinese history》> |

### 3.2.3    Similarity Computation

There are many different expressions for a predicate in natural language questions. After acquiring the subject in questions and the predicate dictionary in knowledge base, the similarity between the words in questions and predicates in dictionary is computed to determine the predicate in questions. We combine the non-subject words with the category of the questions into new words.

For example, 生日|birthday is the non-subject word of the question "梁启超的生日是什么时候|When is Qichao Liang's birthday" and the question category is 时间|Time. These two words are combined into a new word "生日时间|birthday time". Then we use the method proposed in [22] to compute the semantic similarity between "生日时间|birthday time" and the words in the predicate dictionary of "梁启超|Qichao Liang", and we choose the word with the highest similarity and more than 0.5 as the predicate in the question. When the predicate of a question has been determined, the object which is the answer can be found,

### 3.3    WKE Algorithm

WKE algorithm is a real-time knowledge extraction algorithm based on search engine. It mainly consists of two parts: 1) structured knowledge extraction from Baidu Zhixin; 2) unstructured text knowledge extraction from web pages. Knowledge extraction procedure is shown in Figure 3.

Firstly, we use extract structured knowledge from Baidu Zhixin. If we cannot get the answer from Baidu Zhixin or Baidu Zhixin doesnot return any result in search engine, we combine the title and abstract of web pages in search engine result into an unstructured text and extract answers from it.
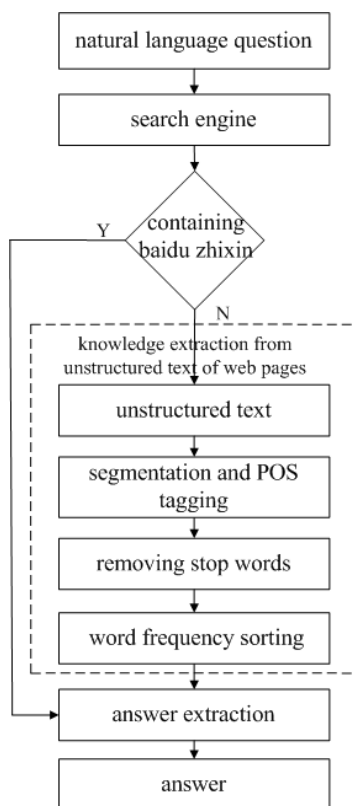
**Fig. 3.** WKE knowledge extraction procedure

### 3.3.1    Knowledge Extraction from Baidu Zhixin

Baidu Zhixin is a QA system integrated into Baidu search engine, which can answer simple semantic question. For example, if input the natural language question "成都有多少人口|What is the population of Chengdu" in Baidu search engine, the search result of Baidu Zhixin is shown in Figure 4:



**Fig. 4.** Baidu Zhixin search result

Baidu Zhixin is a semi-structured data, so we can extract the answer by analyzing the structure of Baidu Zhixin. The structure is different according to different types of queries. Some flags marking the start of answers in Baidu Zhixin source code are shown in Table 3.

**Table 3.** Flags in source code of Baidu Zhixin

| Flags | Question examples |
| --- | --- |
| op_exactqa_s_answer | 成都有多少人口| what is the population of Chengdu |
| c-text c-text-important | 春秋五霸是谁|who are the Five Overloads in Spring and Autumn Period |
| op_exactqa_body | 京城四美是谁|who are the Four Beauty in Capital City |
| op_definitive_answer_po_item_img | 海贼王四皇是谁| who are the Four Emperor in One Piece |

Baidu Zhixin can only answer some of questions. If the questions cannot be solved by Baidu Zhixin, we need to analyze the unstructured text and extract answers from it.

### 3.3.2    Knowledge Extraction from Unstructured Texts

The title and abstract information of the first page of the search engine are formed an unstructured text which may contain the answer with high probability. For example, the search result of the question "谷歌创始人是谁|who is the founder of Google" is shown in Figure 5.



**Fig. 5.** Search result example

The process of extracting knowledge from unstructured text is as follows:

(1)   Get the title and abstract content of each link in the first page.
(2)   Segment and tag POS for the unstructured text.
(3)   Remove stop words.
(4)   Count word frequency and sort the words
(5)   Extract the word with the highest frequency according to the category of questions.

We use the term frequency-inverse document frequency (TF-IDF) to count word frequency at step 4. To gain the useful feature of answer, WKE algorithm defines the categories of questions as Time, Geography, People and Others, which are also the types of answers. When the category of the question is Time, the answer of the question is a time. When the category is People, the answer of the question is a person name. We extract the answer according to the POS tag of words. The POS tag of Time is "t", the tag of People is "nr" and the tag of Geography is "ns". Other category is not limited to any POS tag. Therefore, the answer extraction and features selection rules are summarized as following:

(1) In Geography category questions, if multiple words have "ns" POS tag and appear in the text continuously, they can be combined into a complete geographic name. Otherwise the word tagged as "ns" with the highest frequency is extracted as the answer.
(2) In People category questions, if the question contains "五霸|The Five Wverloads" , "四杰|Four Outstanding Poets", "四大|Four Celebrities" , "四少|Four Talented People", "六君子|Six Nobles",, and the words in the text have "nr" POS tag and appear in successive, these person names are extracted as the answer. Otherwise, the word tagged as "nr" with the highest frequency is extracted as the answer.
(3) In Time category questions, if the unstructured text contains the words with "t" POS tag, the word with the highest frequency is extracted. If the unstructured text doesn't contain the words with "t", the word with the POS tag of "m" which refers to quantifiers and the highest frequency is extracted as the answer.
(4) For the above three categories, if there is not any part of speech needed, any word with the highest frequency is as the answer.
(5) Other question category, the word with the highest frequency is the answer.

### 3.4    Answer Format Standardization

Answer format standardization refers to change the format of answer which is tagged as Time or Geography according to the words contained in questions.

For the Geographic category questions, if the questions contain the characters such as "省|province", "国家|country", "州|state", or "市|city", the part in answers after these above characters is removed.

For the Time category questions, if the questions contain about the characters such as "年|year", the part in answers after the above character is removed.

# 4    Experiment

## 4.1    Data Set

The data set used in this paper is published by NLPCC2015 evaluation task and includes knowledge base and testing questions. The knowledge base is huge which has about 48 million triples given in the form as shown in Table 4:

**Table 4.** The triple knowledge information in Knowledge base

| Subject | Predicate | Object |
|---|---|---|
| 成都\|Chengdu | 中文名称\|Chinese name | 成都\|Chengdu |
| 成都\|Chengdu | 外文名\|English name | chengdu |
| 成都\|Chengdu | 别名\|alias | 蓉城\| Rong city、芙蓉城\| lotus city、锦官城\|the city of brocade、天府之国\| the land of abundance |
| 成都\|Chengdu | 所属地区\|subordinate regions | 西南地区\| southwest |
| 成都\|Chengdu | 面积\|area | 12390平方公里\|12390 square kilometers |
| 成都\|Chengdu | 人口\|population | 1417万（2013年）\|1.417 millions (2013) |
| 成都\|Chengdu | 常住人口\| resident population | 1435万（2013年）\|1.435 millions (2013) |
| 成都\|Chengdu | 气候条件\| climate condition | 亚热带季风性湿润气候\|subtropical monsoon humid climate |
| ----- | ----- | ----- |

The number of the testing questions is 1000, which contains 421 questions of People category, 179 questions of Geography category, 168 questions of Time category, and the remains belong to Others category. Some of question examples of each category are shown in Table 5:

**Table 5.** Question examples of each category

| Category | Question example | Category | Question example |
|---|---|---|---|
| People | 初唐四杰是谁?\|who is the four distinguished poet in Primary Tang? | Time | 什么时候开始的南水北调工程？\| when did the South-North Water Transfer Project start？ |
| Geography | 武当山在哪里?\|where is the Wudang Mountain? | Other | 新加坡的国花是什么？\|what is the national flower of Singapore？ |

## 4.2    Evaluation Metric

This evaluation uses the MRR (Reciprocal Rank Mean), accuracy and averaged F1 values to measure the performance of open domain QA, and the related calculation formula is defined as following:

$$MRR = \frac{1}{|Q|}\sum_{i=1}^{|Q|}\frac{1}{rank_i} \qquad (1)$$

$|Q|$ denotes the total number of testing questions set, $rank_i$ denotes position of the first correct answer, if the all answer is not correct, $\dfrac{1}{rank_i}$ is set to 0.

$$Accuracy = \frac{1}{|Q|}\sum_{i=1}^{|Q|}\delta(C_i, A_i) \qquad (2)$$

$C_i$ denotes test answer, $A_i$ denotes standard answer, if $A_i$ contains one answer of $C_i$ at least, $\delta(C_i, A_i)$ equals to 1, otherwise it equals to 0.

$$AveragedF1 = \frac{1}{|Q|}\sum_{i=1}^{|Q|}F_i \qquad (3)$$

If $A_i$ are not fully existed in $C_i$, $F_i$ equals to 0, otherwise, $F_i$ can be calculated by using the following formula:

$$F_i = \frac{2*\dfrac{\#(C_i, A_i)}{|C_i|}*\dfrac{\#(C_i, A_i)}{|A_i|}}{\dfrac{\#(C_i, A_i)}{|C_i|}+\dfrac{\#(C_i, A_i)}{|A_i|}} \qquad (4)$$

$\#(C_i, A_i)$ denotes the total number of same answer between $C_i$ and $A_i$, $|C_i|$ and $|A_i|$ is the total number of the answers in $|C_i|$ and $|A_i|$.

## 4.3    Experimental Results

The methods presented in this paper are described below: (A) SPE algorithm based on knowledge base; (B) SKEBZ (Structured Knowledge Extraction Based on Baidu Zhixin); (C) UKEWP (Unstructured Knowledge Extraction Based on Web Page). After the QA system receives a natural language question, firstly, we firstly use method (A) to answer the question, and then uses method (B), finally, uses method (C). In order to evaluate the performance of each method, the averaged F1, accuracy and MRR values of each algorithm are given for the 1000 testing data, the results is shown in Table 6:

**Table 6.** Performance comparison of the three methods proposed in this paper

| Question answering methods | MRR | Accuracy | Average F1 |
| --- | --- | --- | --- |

| | | | |
|---|---|---|---|
| SPE | 0.1601 | 0.1730 | 0.1312 |
| SKEBZ | 0.3600 | 0.3690 | 0.3314 |
| UKEWP | 0.4841 | 0.4970 | 0.4479 |
| SPE+SKEBZ | 0.3914 | 0.4022 | 0.3692 |
| SPE+UKEWP | 0.5317 | 0.5540 | 0.5066 |
| SKEBZ+UKEWP | 0.5390 | 0.5481 | 0.5211 |

The results of the methods proposed in the evaluation task in NLPCC2015 are shown in Table 7:

**Table 7.** Evaluation results in the task

| Registered team | MRR | Accuracy | Average F1 |
|---|---|---|---|
| Team-1 | 0.1430 | 0.1660 | 0.1196 |
| **Team-2** | **0.5675** | **0.5700** | **0.5240** |
| Team-3 | 0.3360 | 0.4130 | 0.2990 |

Our methods achieve the best performance that MRR is 0.5675, accuracy is 0.57, and average F1 is 0.5240 in all teams.

## 5    Conclusion

Aiming at NLPCC2015 open domain QA evaluation task of complex sentence structure, the paper presents a solution of open domain QA. Firstly, question is divided into "People", "Time", "Geography" and "Other" category, and then SPE algorithm is proposed to search answers in knowledge base. If the answer is not retrieved in knowledge base, WKE algorithm is adopted to answer the questions. In addition, the paper studies how to use question context information to standardize the answer format in order to enhance the accuracy of QA. The proposed methods achieve good performance that MRR is 0.5672, accuracy is 0.57 and average F1 value is 0.5240 in NLPCC2015 open domain evaluation, and the evaluation results fully proves the feasibility of method proposed in this paper.

## References

1. Mooney, R.J.: Learning for semantic parsing. In: Gelbukh, A. (ed.) CICLing 2007. LNCS, vol. 4394, pp. 311–324. Springer, Heidelberg (2007)
2. Filman, R.E., Sangam, P.: Searching the Internet. IEEE Internet Computing **2**(4), 21–23 (1998)
3. Jeon, J., Croft, W.B., Lee, J.H.: Finding simliar question in large question and answer archives. In: Proceedings of the ACM Fourteen Conference on Inference and Knowledge Management, pp. 84−90 (2005)

4. Zettlemonyer, L.S., Collins, M.: Learning to map sentence to logical form: Structured classification with probabilistic catagorical grammars. In: Proceedings of the 21th Conference on Uncertainty in Artificial Intelligence, pp. 658−666 (2005)

5. Wong, Y.W., Mooney, R.J.: Learning for semantic parsing with statistical machine translation. In: Proceedings of the Main Conference on Human Language Technology Conference of the North American Chapter of the Association of Computational Linguistics, pp. 439−446 (2006)

6. Wong, Y.W., Mooney, R.J.: Generation by inverting a semantic parser that uses statistical machine translation. In: Proceeding of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics, Rochester, pp. 172−179 (2007)

7. Zettlemoyer, L., Collions, M.: Online learning of relaxed CCG grammars for parsin to logical form. In: Proceeddings of ENNLP, pp. 678−687 (2007)

8. Kwiatkowski, T., Luke, S., ZettleMoyer, S.G.: Lexical generalization in CCG grammer induction for semantic parsing. In: Proceeddings of ENNLP, pp. 1512−1523 (2011)

9. Kwiatkowski, T., Luke, S., ZettleMoyer, S.G.: Inducing probabilistic CCG grammar form with higherorder unification. In: Proceeddings of EMNLP, pp. 1223−1233(2010)

10. Freebase (2015). http://www.freebase.com

11. Poon, H.F., Domingos, P.: Unsupervised semantic parsing. In: Proceeddings of EMNLP, pp. 1−10 (2009)

12. Yahya, M., Berberich, K., Elbassuoni, S.: Natural language question for the web of data. In: Proceedings of EMNLP, pp. 379−390 (2012)

13. Berant, J., Chou, A., Roy, F., et al.: Freebase QA: Information extraction or semantic parsing. In: The 2014 Conference on Empirical Methods on Natural Language Processing, pp. 1511−1527 (2014)

14. Berant, J., Chou, A., Roy, F., et al.: Semantic parsing on freebase from question-answer pairs. In: The 2013 Conference on Empirical Methods on Natural Language Processing, pp. 153−1544 (2013)

15. Berant, J., Liang, P.: Semantic parsing via paraphrasing. In: The 52nd Annual Meeting of the Association for Computational Linguistics, pp. 479-485 (2014)

16. Bordes, A., Chopra, S., Weston, J.: Question answering with subgraph embeddings. In: The 2014 Conference on Empirical Methods on Natural Language Processing, pp. 1535−1545 (2014)

17. Yao, X., Durme, B.: Information extraction over structured data: question answering with freebase. In: The 52nd Annual Meeting of the Association for Computational Linguistics, pp. 753−770 (2014)

18. Fader, A., Luke, Z., Oren, E.: Open question answering over curated and extracted knowledge bases. In: Proceedings of the Conference on Knowledge Discovery and Data Mining (KDD), pp. 1256−1265 (2014)

19. Kushman, N., Artzi, Y., Luke, Z., et al.: Learning to automatically solve algebra word problems. In: The 52nd Annual Meeting of the Association for Computational Linguistics, pp. 1042−1061 (2014)

20. Bao, J., Duan, N., Zhou, M., et al.: Knowledge-based question answering as machine translation. In: The 52nd Annual Meeting of the Association for Computational Linguistics, pp. 1272−1294 (2014)

21. Cristina, E.B., Pere, R., Comas.: Full machine translation for factoid question answering. In: Proceedings of EACL, pp. 20−29. ACL, Stroudsburg (2012)
22. Wei, C.Y., Zhan, Q., Fan, X.Z., et al.: Event Information Enhanced Question Semantic Representation for Chinese Question Answering System. Journal of Chinese Information Processing **1**(29), 147–154 (2015)
23. SWJTU Chinese Word Segmentation System. http://ics.swjtu.edu.cn
24. Cho, H.C., Okazaki, N., Miwa, M., et al.: Named Entity Recognition with Multiple Segment Representation. Information Processing Management **49**(4), 954–965 (2013)