

Original Research Article

A Hybrid Question Answering System

ABSTRACT

In this study, we propose a hybrid Question Answering (QA) system for Arabic language. The system consists of four modules. 1) a knowledgebase(KB), 2) an online module, and 3) A Text- to-KB transformer to construct our own knowledge base from web texts. Using these modules, we can query three types of information sources: knowledge bases, constructed knowledge bases from web text. Text-to-KB uses web search results to identify question topic entities, map question words to KB predicates, and to enhance the features of the candidates obtained from the KB. The system scored f-measure of .495 when using KB. The system performed better with f-measure of .573 when using both KB and Text-to-KB module. The system demonstrates higher performance by combining knowledge base and text from external resources.

Keywords: Question Answering System, Information Retrieval, Information Extraction, Knowledge Base

1. INTRODUCTION

Whenever a user needs information about a specific topic, it simply supplies a query to any search engine, e.g. Google. Traditional search engines returns a list of links to documents which may contain the answer. The user has to browse these links and tries to locate the answer. QA systems retrieves specific answers in response to user questions , rather than a lists of links to documents. Two approaches for Question Answering (QA) have evolved: text-centric, and knowledge base-centric. Text-Centric QA systems use collection of text documents to return passages relevant to a user's question and extract candidate answers [1]. The KB-Centric QA systems, which are emerged from the database community, depends on large scale knowledge bases, such as Freebase [2], DBpedia [3], WikiData [4] which store a massive amount of knowledge about various kinds of entities. KBQA systems have been classified into two major approaches: semantic parsing, and information extraction (IE) [5]. The semantic parsing focuses on understanding the question, and tries to parse sentences into their logical forms (semantic representations)[6, 7, 8]. Information Extraction(IE) approaches [9, 10, 11] are based on detecting topic entities in the question, and employing predefined templates for mapping the question to predicates, exploring these entities' neighborhood in a KB. Various QA based on knowledge base (KB) approaches to have been proposed. QA systems are developing from systems based on information retrieval (IR) to ones based on KBs. QA systems based on KBs provides very high precision, but requires curated KBs; However, these KBs cannot include all the information that web text can communicate. To overcome this limitation, other information sources besides curated KBs are required. In this paper, we present a hybrid QA system that utilizes multiple information sources: a curated KB, KB constructed from text, and web text.

2. METHODOLOGY

The following figure shows the architecture of the system.

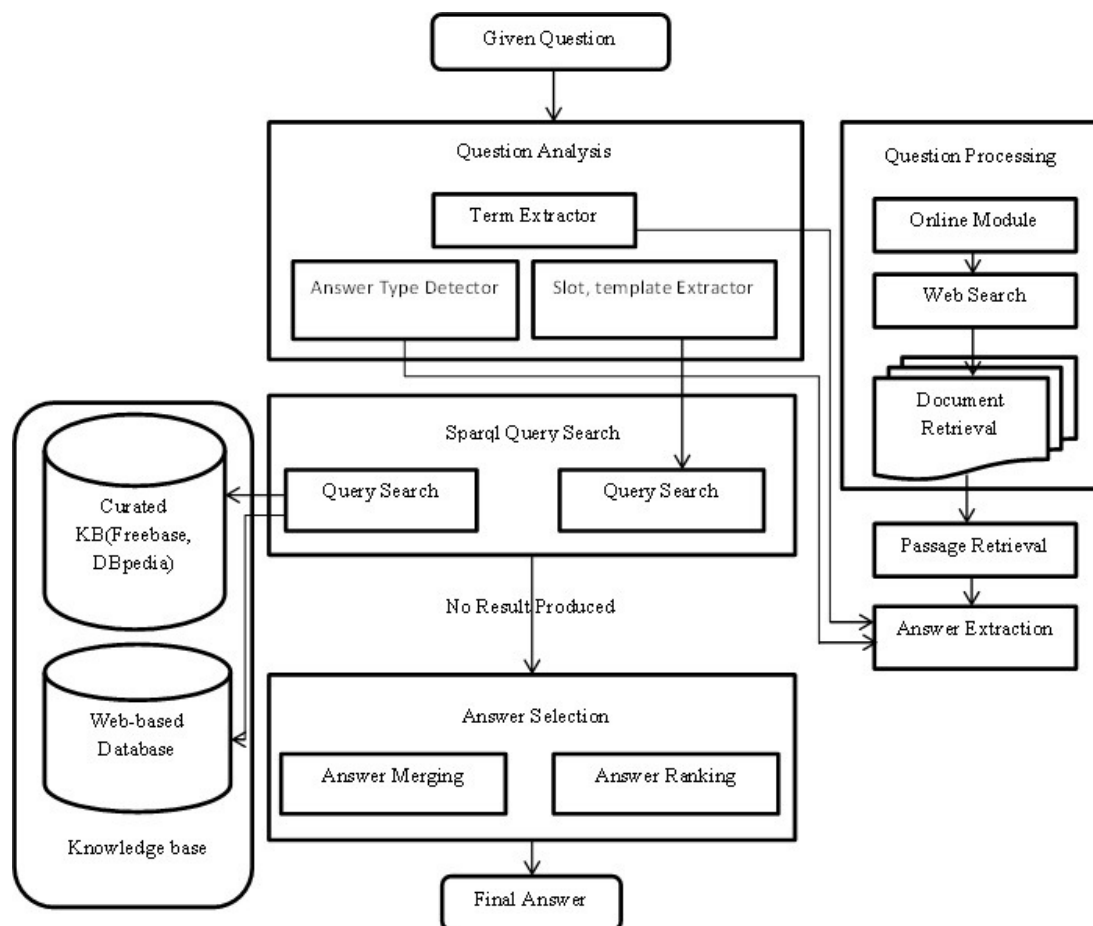


Fig. 1. QA Architecture

2.1 Knowledge base

A Question Answering system based on KB takes a natural language (NL) question as its input and uses structured KBs like DBpedia to retrieve the answer. A KB-based QA system employs structured information sources, so it generates very specific answers. First the NL question is segmented/tokenized into individual words/tokens; then string based methods are employed and NL phrases along with KB node mapping dictionary are automatically generated to match KB vocabulary to the tokens. We generate query candidates by using a limited set of hand-crafted grammar rules to combine tokens into a single unified representation of meaning. In the LSP approach, patterns that consist of regular expression patterns that express the POS ,lexical or chunk type patterns of a NL question and a SPARQL query template are generated. If a match is found, slots in the SPARQL query template are occupied with the word-matched chunks from NL question. However, there is no context information for KB-based QA modules , and therefore it cannot score/rank its answer candidates; instead KB-based module forwards its answer candidate to an answer merging task in the online module and this module rank the answer candidates.

2.2 Online Module

The online module searches text to find answers. The online module performs four tasks (Fig 1): first is question classification. It analyzes the question semantically and identifies the answer type and; the second is the passage retriever. It retrieves relevant passages by segmenting the documents that are related to the user question; the third task is the answer extractor. It extracts answer candidates; the fourth task merges answer candidates from the online module and KB, it then ranks the answer candidates and returns the final list of answers. Context information are used to scores answer

candidates which are the output of the SPARQL[12] not only from online module answer extraction task. Lexical, syntactic and semantic analysis are employed for question processing, which includes extracting terms by a Support Vector Machines (SVM) [13]. Lucene[14] is utilized for indexing web pages dump and for searching and processing relevant documents and passages which contain the answer. After the analysis of passages is performed, sentences in the passages are scored. Named Entities(NEs) which have the same or similar answer types as answer candidates from top-n sentences in passages are extracted. Finally, our system ranks answer candidates from answer extraction task using semantic similarity between question and sentences that include answer candidates and the final answer list is to user is delivered to the user.

2.3 Text-to-KB

The limitation of the KB is that it can only store small amount of information as compared to its original unstructured text. To overcome this problem, we extract triples from unstructured text and store them in a repository. In order to extract triples from unstructured text, we use the semantic role labels of a sentence and the dependency tree. Extraction templates are constructed that specify, for each dependency tree structure pattern, how triples should be extracted. A full document is retrieved to detect sentences that include word tokens that occur in arguments and relation words of each seed triple. Then a dependency tree of the sentence for each seed triple is constructed, sentence pair, and a linear path which contains arguments and relation words is identified. This path with location of arguments and relation words can generate an extraction template. Semantic rule labeling provide similar results that can be converted to triple format. Predicates of the results are considered as relation phrases and each argument and argument modifier are considered as each argument of triples. A small set of rules is also used to convert semantic rule labeling results to triples.

3. RESULTS AND DISCUSSION

The results of existing approaches and our Text-to-KB system are provided in Table 1. The result reported for our QA system is computed using precision, Recall and F1-measure. As we can see, Text2KB significantly improves over the baseline system.

Table 1. System performance using KB only& Using both KB and Text-to-KB module

System	Precision	Recall	F1-measure
Knowledge base	0.635	.406	.495
Knowledge base+ Text-to-KB (Web Search)	0.642	0.519	.573

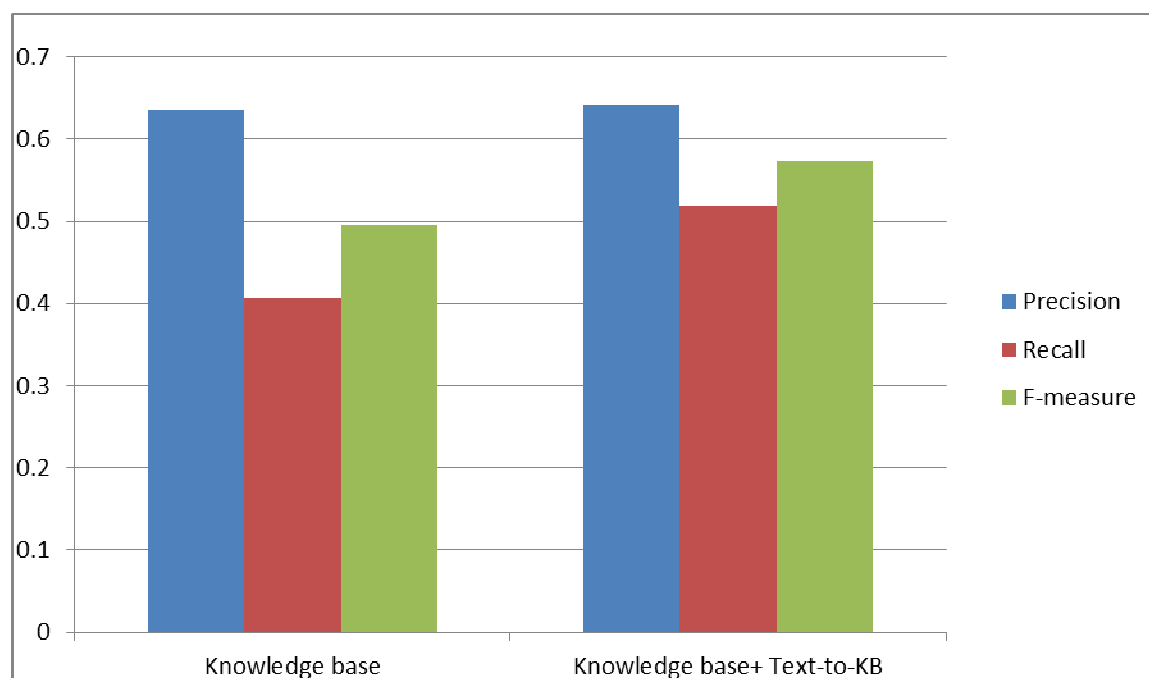


Fig. 2. System performance using KB & KB+ Text-to-KB

We demonstrated that by coupling evidence from knowledge base and text from external resources the system performance can be boosted. The system scored .495 for the f-measure using only KB. The performance of the system the system is proved to be better using both KB and text-to-KB. The computed f-measure using both methods is .573.

4. CONCLUSION

In this paper, we show that unstructured text resources can be used for knowledge base question answering to enhance query understanding, generation of candidate answer and ranking. We focused on two techniques and text information sources: web search results for query understanding and training data for candidate generation and ranking. The features employed are an n-gram of words and POS. The proposed system uses semantic relatedness among question and sentences to rank answer candidates from KB and online module and provide the final answer list to user.

REFERENCES

1. Dang HT, Kelly D, and Lin JJ. Overview of the TREC 2007 question answering track. In Proceedings of TREC, 2007.
2. Bollacker K, Evans C, Paritosh P, Sturge T, and Taylor J. Freebase: A collaboratively created graph database for structuring human knowledge. In Proceedings of the 2008 ACM SIGMOD international conference on Management of data (SIGMOD '08), 2008: 1247-1250, doi:10.1145/1376616.1376746.
3. Auer S, Bizer C, Kobilarov G, Lehmann J, Cyganiak R, and Ives Z. Dbpedia: A nucleus for a web of open data. The Semantic Web. Lecture Notes in Computer Science. Springer, 2007; 4825: 722-735. DOI: 10.1007/978-3-540-76298-0_52.
4. Vrande D and Krotzsch M. Wikidata: A free collaborative knowledgebase. Communications of ACM, Sept. 2014;57(10): 78-85, DOI: 10.1145/2629489.

5. Yao X, Berant J, and Durme BV. Freebase qa: Information extraction or semantic parsing?. In Proceedings of the ACL 2014 Workshop on Semantic Parsing, ACL, 2014, DOI: 10.3115/v1/W14-2416.
6. Berant J, Chou A, Frostig R, and Liang P. Semantic parsing on freebase from question-answer pairs. In Proceedings of EMNLP, 2013.
7. Berant J and Liang P. Semantic parsing via paraphrasing. In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics , ACL, 2014; 1: 1415–1425 ,DOI: 10.3115/v1/P14-1133.
8. Berant J and Liang P. Imitation learning of agenda-based semantic parsers. Transactions of the Association for Computational Linguistics,ACL, 2015; 3: 545–558.
9. Bast H and Haussmann E. More accurate question answering on freebase. In Proceedings of the 24th ACM International on Conference on Information and Knowledge Management(CIKM '15), ACM, 2015; 1431-1440, DOI: 10.1145/2806416.2806472.
10. Yih W-T, Chang M-W, He X, and Gao J. Semantic parsing via staged query graph generation: Question answering with knowledge base. In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing, ACL, 2015; 1321–1331 , DOI: 10.3115/v1/P15-1128.
11. Yao X and Van Durme BV. Information extraction over structured data: Question answering with freebase. In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL, 2014; Vol. 1, DOI: 10.3115/v1/P14-1090.
12. SPARQL. Accessed 27 January 2017. Available : <https://www.w3.org/TR/sparql11-query/>.
13. Schlaefter N, Ko J, Betteridge J, Pathak MA, Nyberg E AND Sautter G . Semantic Extensions of the Ephyra QA System for TREC 2007. In Proceedings of The Sixteenth Text REtrieval Conference(TREC 2007), 2007.
14. Lucene. Accessed 03 February 2017. Available :<http://lucene.apache.org/core>.

APPENDIX