

Question Answering Using Structured and Unstructured Data

Doctoral thesis proposal

Denis Savenkov

Dept. of Math & Computer Science

Emory University

denis.savenkov@emory.edu

March, 2016

Abstract

Over more than half a century of research, the area of automatic Question Answering (QA) has progressed from small single domain systems to IBM Watson, who defeated best human competitors in the Jeopardy! TV show [42]. However, many of our questions are still left unanswered, and we still have a lot to do to move beyond 10 blue links in search results [37] as for most of the questions users still have to dig into the retrieved documents or post questions to the community question answering (CQA) websites. Questions come in different flavors, some are asking about a certain fact and can be answered with a short phrase, such as entity name, date or number. Such questions are typically referred to as *factoid*, as opposed to rest of the questions, which are often called *non-factoid*. In my thesis I focus on three topics in QA: 1). combination of structured and unstructured data to improve factoid question answering; 2). improving question summarization, candidate scoring and answer generation using recent advances in neural network research and better utilization of source web document structure; 3). interactions between a question answering system and real people.

Text document collections and knowledge bases (KB) are very effective in answering certain types of factoid questions, but they are also complimentary to each other. I propose to combine these data sources via semantic annotation of KB entity mentions, which effectively extends the knowledge base with additional unstructured information, often missing or complimentary to the KB data.

Non-factoid question answering is somewhat harder as it deals with a more diverse set of question and answer types. In my thesis I propose to improve performance of different stages of QA system pipeline by better utilization of the structure of a web page where a candidate answer is extracted from, and using deep learning techniques, inspired by recent successes in machine translation [8], text summarization [83], automatic caption generation [61] and answer sentence scoring [105].

The focus of the last part of the thesis is on the interaction between human and a search or question answering system. Unfortunately, there always be cases, when a machine cannot provide a good answer to the question. In such cases, a QA system may come back to the user with some suggestions on how a user can solve his search problem, or with a clarification question aiming at resolving certain ambiguities. Alternatively, a machine can consult with the crowd in order to get the answer or help it decide on certain alternatives.

In summary, the goal of the proposed research is to improve the performance of question answering over a variety of different questions a user might have, and to study some reply strategies in case a system fails to deliver a good response. I believe, that results of the proposed work will be useful for the future research in improving automatic question answering.

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Research Questions	3
1.3	Research Plan	4
1.3.1	Step 1 (Chapter 3)	4
1.3.2	Step 2 (Chapter 4)	4
1.3.3	Step 3 (Chapter 5)	4
1.3.4	Research Timeline	4
1.4	Contributions and Implications	4
2	Related Work	5
2.1	Factoid question answering	5
2.1.1	Text-based question answering	5
2.1.2	Knowledge base question answering	7
2.1.3	Hybrid question answering	8
2.2	Non-factoid question answering	10
2.3	User interactions	11
2.4	Summary of Related Work	11
3	Structured and Unstructured Data for Factoid Question Answering	12
3.1	Relation Extraction for Knowledge Base Completion	12
3.1.1	Relation extraction from Question-Answer pairs	12
3.2	Semantic Text Annotations for Hybrid Question Answering	12
3.2.1	Approach	13
3.2.2	Evaluation	14
3.3	Summary	15
4	Non-factoid Question Answering	16
4.1	The architecture of the system	16
4.1.1	Question Analysis	16
4.1.2	Structure of the Web Page for Candidate Generation and Scoring	17
4.1.3	Answer Summarization	17
4.2	Evaluation	18
4.3	Summary	18

5	Human Interaction with Question Answering Systems	20
5.1	Search Hints for Complex Informational Tasks	20
5.2	Clarification Questions	20
5.3	Using the Wisdom of a Crowd for Question Answering	20
6	Summary and Discussion	21
	Bibliography	22

1 Introduction

1.1 Motivation

It has long been a dream to communicate with a computer as with another human being using natural language speech and text. Nowadays, we are coming closer to this dream as natural language interfaces become more and more popular. Our phones are already reasonably good in recognizing speech, and personal assistants, such as Apple Siri, Google Now, Microsoft Cortana, Amazon Alexa, *etc.*, help us with everyday tasks and answer some of our questions. Chat bots are arguably considered “the next big thing”, and a number startups developing this kind of technologies has emerged in Silicon Valley and around the world¹.

Questions are a natural way of communicating ones information needs. Users of modern search engines are already used to receiving an answer panel in response to some of their questions². However, there is still a big number of questions, for which users have to click on 10 blue links and mine the answer from the huge amount of information contained in retrieved documents. Therefore, in my thesis I focus on helping users get answers to their questions by improving question answering methods and interactions between a human and a computer system.

It’s common to divide questions into *factoid* and *non-factoid*. Factoid questions are inquiring about certain facts and can be answered with a short phrase (or list), *i.e.* entity name, date or number. An example of a factoid question is “*What book did John Steinbeck wrote about the people in the dust bowl?*” (answer: “*The Grapes of Wrath*”). Besides this type of questions, people often ask for recommendations, opinion, instructions, definitions, reason *etc.*. These questions are typically combined under an umbrella of non-factoid questions.

There are multiple different types of information sources, that various QA systems use to find the answers to user questions. These data sources can be classified into unstructured (*e.g.* raw natural language text), semi-structured (*e.g.* tables) and structured (*e.g.* knowledge bases). Each of these types has certain advantages and limitations (Table 1.1).

Two major paradigms for factoid question answering are knowledge base question answering (KBQA) and text-based question answer (TextQA). Information contained in a huge volume of text data on the web can be relatively easily queried using terms and phrases from the original question in order to retrieve sentences that might contain the answer. However, each sentence encode very limited amount of information about mentioned entities and aggregating information over unstructured data is quite problematic. On the other hand, modern large scale knowledge bases, such as Freebase [17], dbPedia [6], YAGO [73], WikiData [103], aggregate information about millions of entities into a graph of [subject, predicate, object] RDF triples. The problem with KBs is that they are inherently incomplete and miss a lot of entities, facts and predicates. In addition, triple data representation format complicates retrieval of KB concepts relevant to question phrases. The focus of the proposed research in factoid question answering lies on the idea of combining structured KB and unstructured text data, which can help a QA system to overcome these drawbacks. One way to improve the situation with knowledge base incompleteness is to

¹<http://time.com/4194063/chatbots-facebook-messenger-kik-wechat/>

²<https://www.stonetemple.com/rich-answers-in-search/>

Table 1.1: Pros and cons of structured and unstructured data sources for factoid and non-factoid question answering

	unstructured data	structured data
factoid questions	<p>Text</p> <ul style="list-style-type: none"> + easy to match against question text + cover a variety of different information types - each text phrase encodes a limited amount of information about mentioned entities 	<p>Knowledge Bases</p> <ul style="list-style-type: none"> + aggregate all the information about entities allow complex queries over this data using special languages (e.g. SPARQL) - hard to translate natural language questions into special query languages - KBs are incomplete (missing entities, facts and properties)
non-factoid questions	<p>Text</p> <ul style="list-style-type: none"> + contain relevant information to a big chunk of user needs - hard to extract semantic meaning of a paragraph to match against the question (lexical gap) 	<p>Question-Answer pairs</p> <ul style="list-style-type: none"> + easy to find a relevant answer by matching the corresponding questions - cover a smaller subset of user information needs

extract missing information from other data sources, *e.g.* [24, 25, 35, 38, 49, 64]. In my thesis I focus on one particular data source, that didn’t receive enough attention in the relation extraction literature, namely question-answer pairs. Section 3.1 will describe our experiments and results in utilizing this data to improve knowledge base coverage. Unfortunately, relation extraction isn’t perfect either and there are both precision and recall losses. Alternatively, in my thesis I propose a new hybrid approach to question answering, which leverages a combination of text and knowledge base data to improve every stage of question answering process. More specifically, I propose to use semantic annotation of KB entity mentions in text documents as a bridge between data sources (Section 3.2).

The main challenge in non-factoid question answering lies in the diversity of question and answer types. One of the most effective strategies is to reuse answers to previously asked questions, which could be found, for example, in CQA archives [89]. Unfortunately, it’s not always possible to find a similar question, that has already been answered. Many information needs are unique or contain unique details, which makes it impossible to reuse old answers. Alternative strategies include ranking text passages extracted from retrieved web documents. One of the main challenges of this approach is estimating semantic similarity between the question and an answer candidate [93]. Therefore, one would benefit from knowing what kind of questions could a paragraph of text answer. This information can often be inferred from the structure of a web page, *e.g.* forums, FAQ pages, or estimated using title, subtitle and other page elements. Therefore, one of the questions I’m going to focus in my thesis is how to effectively use the structure of web page to predict whether an extracted passage of text answer the given question.

However, ranking isn’t the only important part of the question answering pipeline. A system can only rank and return a good answer if it was able to retrieve relevant information from a

collection. Non-factoid questions, especially those that people post on CQA websites are often long, which makes it problematic to use directly as search queries. Previous research has studied certain question transformation strategies [1, 22, 70], however the focus was on shorter factoid questions. In my thesis I would like to focus on the problem of query generation for non-factoid questions using some recent advances in deep learning. Another promising direction of research, which I'm going to explore in my thesis, is answer generation, *i.e.* by summarizing the information a system could retrieve. Different answer candidates might be complementary to each other, answer different parts of the question or have different opinions on the subject.

NEED A PARAGRAPH ABOUT USERS AND QA.

1.2 Research Questions

Research questions I proposed addressed in my thesis are the following:

1. How to effectively combine unstructured text and structured knowledge base data to improve factoid question answering?
 - (a) What types of questions can be answered using text, KB or a combination of both?
 - (b) How does semantic annotation of unstructured data compare to information extraction for question answering? (information extraction for KB construction vs open information extraction vs unstructured data annotation)
 - (c) How does a combination of structured and unstructured data sources improve each of the main QA system components: question analysis, candidate generation, evidence extraction and answer selection?
2. How to improve question understanding and answer generation for non-factoid question answering using recent advances in deep learning?
3. What kind of information about the structure of a web page can help to score an extracted passage as a candidate answer to the given question?
4. How can we improve user success with question answering by providing them with search strategy hints, clarifications or by employing the crowd?

1.3 Research Plan

1.3.1 Step 1 (Chapter 3)

1.3.2 Step 2 (Chapter 4)

1.3.3 Step 3 (Chapter 5)

1.3.4 Research Timeline

1.4 Contributions and Implications

The key contributions of the proposed research are: 1. New hybrid KB-text question answering algorithm, that is based on graph search, which includes both KB links as well as text search edges to follow. 2. New labelled dataset for question answering (???) 3. New features for ranking answer candidates ???

2 Related Work

The field of automatic questions answering has a long history of research and dates back to the days when the first computers appear. By the early 60s people have already explored multiple different approaches to question answering and a number of text-based and knowledge base QA systems existed at that time [90, 91]. In 70s and 80s the development of restricted domain knowledge bases and computational linguistics theories facilitated the development of interactive expert and text comprehension systems [5, 88, 112, 111]. The modern era of question answering research was motivated by a series of Text Retrieval Conference (TREC¹) question answering shared tasks, which was organized annually since 1999 [102]. A comprehensive survey of the approaches from TREC QA 2007 can be found in [31]. An interested reader can refer to a number of surveys to track the progress made in automatic question answering over the years [51, 4, 106, 63, 80, 3, 48].

The main focus of research in automatic question answering was on factoid questions. However, recently we can observe an increased interest in non-factoid question answering, and as an indicator in 2015 TREC started a LiveQA shared task track², in which the participant systems had to answer various questions coming from real users of Yahoo! Answers³ in real time.

In the rest of the chapter I will describe related work in factoid (Section 2.1) and non-factoid (Section 2.2) question answering with the focus on data sources used.

2.1 Factoid question answering

Since the early days of automatic question answering researches explored different sources of data, which lead to the development of two major approaches to factoid question answering: text-based (TextQA) and knowledge base question answering (KBQA) [90]. We will first describe related work in TextQA (Section 2.1.1), then introduce KBQA (Section 2.1.2) and in Section 2.1.3 present existing techniques for combining different information sources together.

2.1.1 Text-based question answering

A traditional approach to factoid question answering over text document collections, popularized by TREC QA task, starts by querying a collection with possibly transformed question and retrieving a set of potentially relevant documents, which are then used to identify the answer. Information retrieval for question answering has certain differences from traditional IR methods [62], which are usually based on keyword matches. A natural language question contains certain information, that is not expected to be present in the answer (*e.g.* the keyword *who*, *what*, *when*, *etc.*), and the answer statement might use language that is different from the question (lexical gap problem). On the other side, there is a certain additional information about expected answer statement, that a QA system might infer from the question (*e.g.* we expect to see in a number in response to the “how many” question). One way to deal with this problem is to transform the question in certain

¹<http://trec.nist.gov>

²<http://trec-liveqa.org/>

³<http://answers.yahoo.com/>

ways before querying a collection [1, 22]. Raw text data might be extended with certain semantic annotations by applying part of speech tagger, semantic role labeling, named entity recognizer, *etc.*. By indexing these annotations a question answering system gets an opportunity to query collection with additional attributes, inferred from the question [16, 125].

The next stage in TextQA is to select sentences, that might contain the answer. One of the mostly used benchmark datasets for the task, proposed in [108], is based on TREC QA questions and sentences retrieved by participating systems⁴. The early approaches for the task used simple keyword match strategies [56, 94]. However, in many cases keywords doesn't capture the similarity in meaning of the sentences very well and researches started looking on syntactic information. Syntactic and dependency tree edit distances and kernels allow to measure the similarity between the structures of the sentences [81, 87, 50, 124, 107]. Recent improvements on the answer sentence selection task come are associated with the deep learning techniques, *e.g.* recursive neural networks using sentence dependency tree [57], convolutional neural networks [128, 84], recurrent neural networks [98, 105]. Another dataset, called WikiQA [120], raises a problem of answer triggering, *i.e.* detecting cases when the retrieved set of sentences don't contain the answer.

To provide a user with the concise answer to his factoid question QA systems extract the actual answer phrase from retrieved sentences. This problem is often formulated as a sequence labeling problem, which can be solved using structured prediction models, such as CRF [124], or as a node labeling problem in an answer sentence parse tree [74].

Unfortunately, passages include very limited amount of information about the candidate answer entities, *i.e.* very often it doesn't include the information about their types (person, location, organization, or more fine-grained CEO, president, basketball player, *etc.*), which is very important to answer question correctly, *e.g.* for the question "*what country will host the 2016 summer olympics?*" we need to know that **Rio de Janeiro** is a city and **Brazil** is the country and the correct answer to the question. Therefore, a lot of effort has been put into developing answer type typologies [54, 53] and predicting and matching expected and candidate answer types from the available data [67, 68, 79]. Many approaches exploited external data for this task, I will describe some of this efforts in Section 2.1.3.

Very large text collections, such as the Web, contain many documents expressing the same information, which makes it possible to use a simpler techniques and rely on redundancy of the information. AskMSR QA system was one of the first to exploit this idea, and achieved very impressive results on TREC QA 2001 shared task [21]. The system starts by transforming a question into search queries, extracts snippets of search results from a web search engine, and consider word n-grams as answer candidates, ranking them by frequency. A recent revision of the AskMSR QA system [100] introduced several improvements to the original system, *i.e.* named entity tagger for candidate extraction, and additional semantic similarity features for answer ranking. It was also observed, that modern search engines are much better in returning the relevant documents for question queries and query generation step is no longer needed. Another notable systems, that used the web as the source for question answering are MULDER[65], Aranea [70], and a detailed analysis of what affects the performance of the redundancy-based question answering systems can be found in [69].

⁴A table with all known benchmark results and links to the corresponding papers can be found on [http://aclweb.org/aclwiki/index.php?title=Question_Answering_\(State_of_the_art\)](http://aclweb.org/aclwiki/index.php?title=Question_Answering_(State_of_the_art))

2.1.2 Knowledge base question answering

Earlier in the days knowledge bases were relatively small and contained information specific to a particular domain, *e.g.* baseball statistics [47], lunar geology [112], geography [129]. However, one of the main challenges in KBQA is mapping between natural language phrases to the database concepts, which raises a problem of domain adaption of question answering systems.

Recent development of large scale knowledge bases (*e.g.* dbPedia [6], Freebase [17], YAGO [95], WikiData⁵) shifted the attention towards open domain question answering. Knowledge base question answering approaches can be evaluated on an annual Question Answering over Linked Data (QALD⁶) shared task, and some popular benchmark dataset, such as Free917 [27] and WebQuestions [11]. A survey of some of the proposed approaches can be found in [101].

A series of QALD evaluation campaigns has started in 2011, and since then a number of different subtasks have been offered, *i.e.* since QALD-3 includes a multilingual task, and QALD-4 formulated a problem of hybrid question answering. These tasks usually use dbPedia knowledge base and provide a training set of questions, annotated with the ground truth SPARQL queries. The hybrid track is of particular interest to the topic of this dissertation, as the main goal in this task is to use both structured RDF triples and free form text available in dbPedia abstracts to answer user questions.

The problem of lexical gap and lexicon construction for mapping natural language phrases to knowledge base concepts is one of the major difficulties in KBQA. The earlier systems were mainly trained from question annotated with the correct parse logical form, which is expensive to obtain. Such an approach is hard to scale to large open domain knowledge bases, which contain millions of entities and thousands of different predicates. An idea to extend a trained parser with additional lexicon, built from the Web and other resources, has been proposed by [28]. However, most of the parses of the question produce different results, which means that it is possible to use question-answer pairs directly [11]. PARALEX system ([41]) construct a lexicon from a collection of question paraphrases from WikiAnswers⁷. A somewhat backwards approach was proposed in ParaSempre model of [12], which ranks candidate structured queries by first constructing a canonical utterance for each query and then using a paraphrasing model to score it against the original question. Another approach to learn term-predicate mapping is to use patterns obtained using distant supervision [76] labeling of a large text corpus, such as ClueWeb [122]. Such labelled collections can also be used to train a KBQA system, as demonstrated by [82]. Such an approach is very attractive as it doesn't require any manual labeling and can be easily transferred to a new domain. However, learning from statements instead of question answer pairs has certain disadvantages, *e.g.* question-answer lexical gap and noise in distant supervision labeling. Modern knowledge bases also contain certain surface forms for their predicates and entities, which makes it possible to convert KB RDF triples into questions and use them for training [18]. Finally, many systems work with distributed vector representations for words and RDF triples and use various deep learning techniques for answer selection. A common strategy is to use a joint embedding of text and knowledge base concepts. For example, character n-gram text representation as input to a convolutional neural network can capture the gist of the question and help map phrases to entities and predicates [126]. Joint embeddings can be trained using multi-task learning, *e.g.* a system can learn to embed a question and candidate answer subgraph using question-answer pairs and question paraphrases at the same time ([18]). Memory Networks, developed by the

⁵<http://www.wikidata.org>

⁶www.scit-ec.uni-bielefeld.de/qald/

⁷<https://answers.wikia.com/>

Facebook AI Lab, can also be used to return triples stored in network memory in a response to the user question [19]. This approach uses embeddings of predicates and can answer relatively simple questions, that do not contain any constraints and aggregations. To extend deep learning framework to more complex questions, [34] use multi-column convolutional neural network to capture the embedding of the entity path, context and type.

As for the architecture of KBQA systems, two major approaches have been identified: semantic parsing and information extraction. Semantic parsing starts from question utterances and work to produce the corresponding semantic representation, *e.g.* logical form. The model of [11] uses a CCG parser, which can produce many candidates on each level of parsing tree construction. A common strategy is to use beam search to keep top-k options on each parsing level or agenda-based parsing [13], which maintains current best parses across all levels. An alternative information extraction strategy was proposed by [123], which can be very effective for relatively simple questions. A comparison of this approaches can be found in [122]. The idea of the information extraction approach is that for most of the questions the answer lies in the neighborhood of the question topic entity. Therefore, it is possible to use a relatively small set of query patterns to generate candidate answers, which are then ranked using the information about how well involved predicates and entities match the original question.

Question entity identification and disambiguation is the key component in such systems, they cannot answer the question correctly if the question entity isn't identified. Different systems used NER to tag question entities, which are then linked to a knowledge base using a lexicon of entity names [11, 12, 116]. However, NER can easily miss the right span and the whole system would fail to produce the answer. Recently, most of the state-of-the-art system on WebQuestions dataset used a strategy to consider a reasonable subset of token n-grams, each of which can map to zero or more KB entities, which are disambiguated on the answer ranking stage [121, 9, 99]. Ranking of candidates can be done using a simple linear classification model [121] or a more complex gradient boosted trees ranking model [9, 99].

Some questions contain certain conditions, that require special filters or aggregations to be applied to a set of entities. For example, the question “*who won 2011 heisman trophy?*” contains a date, that needs to be used to filter the set of heisman trophy winners, the question “*what high school did president bill clinton attend?*” requires a filter on the entity type to filter high schools from the list of educational institutions, and “*what is the closest airport to naples florida?*” requires a set of airports to be sorted by distance and the closest one to be selected. Information extraction approaches either needs to extend the set of candidate query templates used, which is usually done manually, or to attach such aggregations later in the process, after the initial set of entities have been extracted [99]. An alternative strategy to answer complex questions is to extend RDF triples as a unit of knowledge with additional arguments and perform question answering over n-tuples [127]. In [110] authors propose to start from single KB facts and build more complex logical formulas by combining existing ones, while scoring candidates using paraphrasing model. Such a template-free model combines the benefits of semantic parsing and information extraction approaches.

2.1.3 Hybrid question answering

A natural idea of combining available information sources to improve question answering has been explored for a long time. WordNet lexical database [75] was among the first resources, that were adapted by QA community [55, 78], and it was used for such tasks as query expansion

and definition extraction. Wikipedia⁸, which can be characterized as an unstructured and semi-structured (infoboxes) knowledge base, quickly became a valuable resource for answer extraction and verification [2, 23]. Developers of the Aranea QA system noticed that structured knowledge bases are very effective in answering a significant portion of relatively simple questions [70]. They designed a set of regular expressions for popular questions that can be efficiently answered from a knowledge base and fall back to regular text-based methods for the rest of the questions.

One of the major drawbacks of knowledge bases is their incompleteness, which means that many entities, predicates and facts are missing from knowledge bases, which limits the number of questions one can answer using them. One approach to increase the coverage of knowledge bases is to extract information from other resources, such as raw text[76, 60, 49], web tables [24], *etc.*. However, the larger the knowledge base gets, the more difficult it's to find a mapping from natural language phrases to KB concepts. Alternatively, open information extraction techniques ([38]) can be used to extract a surface form-based knowledge base, which can be very effective for question answering. Open question answering approach of [40] combines multiple structured (Freebase) and unstructured (OpenIE) knowledge bases together by converting them to string-based triples. User question can be first paraphrased using paraphrasing model learned from WikiAnswers data, then converted to a KB query and certain query rewrite rules can be applied, and all queries are ranked by a machine learning model.

SPOX tuples, proposed in [117], encode subject-predicate-object triples along with certain keywords, that could be extracted from the same place as RDF triple. These keywords encode the context of the triple and can be used to match against keywords in the question. The method attempts to parse the question and uses certain relaxations (removing SPARQL triple statements) along with adding question keyphrases as additional triple arguments. As an extreme case of relaxation authors build a query that return all entities of certain type and use all other question terms to filter and rank the returned list.

However, by applying information extraction to raw text we inevitably lose certain portion of the information due to recall errors, and extracted data is also sometimes erroneous due to precision errors. In [115], authors propose to use textual evidence to do answer filtering in a knowledge base question answering system. On the first stage with produce a list of answers using traditional information extraction techniques, and then each answer is scored using its Wikipedia page on how well it matches the question. Knowledge bases can also be incorporated inside TextQA systems. Modern KBs contain comprehensive entity types hierarchies, which were utilized in QuASE system of [96] for answer typing. In addition, QuASE exploited the textual descriptions of entities stored in Freebase knowledge base as answer supportive evidence for candidate scoring. However, most of the information in a KB is stored as relations between entities, therefore there is a big potential in using all available KB data to improve question answering.

Another great example of a hybrid question answering system is IBM Watson, which is arguably the most important and well-known QA systems ever developed so far. It was designed to play the Jeopardy TV show⁹. The system combined multiple different approaches, including text-based, relation extraction and knowledge base modules, each of which generated candidate answers, which are then pooled together for ranking and answer selection. The full architecture of the system is well described in [42] or in the full special issue of the IBM Journal of Research and Development [43]. YodaQA [10] is an open source implementation of the ideas behind the IBM Watson system.

⁸<http://www.wikipedia.org>

⁹<https://en.wikipedia.org/wiki/Jeopardy!>

2.2 Non-factoid question answering

During earlier days of research non-factoid questions received relatively little attention. TREC QA tasks started to incorporate certain categories of non-factoid questions, such as definition questions, during the last 4 years of the challenge. One of the first non-factoid question answering system was described in [93] and was based on web search using chunks extracted from the original question. The ranking of extracted answer candidates was done using a translation model, which showed better results than n-gram based match score.

The growth of the popularity of community question answering (CQA) websites, such as Yahoo! Answers, Answers.com, *etc.*, contributed to an increased interest of the community to non-factoid questions. Some questions on CQA websites are repeated very often and answers can easily be reused to answer new questions, [72] studies different types of CQA questions and answers and analyzes them with respect to answer re-usability. A number of methods for similar question retrieval have been proposed [14, 89, 36, 58]. WebAP is a dataset for non-factoid answer sentence retrieval, which was developed in [119]. Experiments conducted in this work demonstrated, that classical retrieval methods doesn't work well for this task, and multiple additional semantic (ESA, entity links) and context (adjacent text) features have been proposed to improve the retrieval quality.

Candidate answer passages ranking problem becomes even more difficult in non-factoid questions answering as systems have to deal with larger piece of text and need to "understand" what kind of information is expressed there. One of the first extensive studies of different features for non-factoid answer ranking can be found in [97], who explored information retrieval scores, translation models, tree kernel and other features using tokens and semantic annotations (dependency tree, semantic role labelling, *etc.*) of text paragraphs. Alignment between question and answer terms can serve as a good indicator of their semantic similarity. Such an alignment can be produced using a machine learning model with a set of features, representing the quality of the match [109]. Alignment and translation models are usually based on term-term similarities, which are often computed from a monolingual alignment corpus. This data can be very sparse, and to overcome this issue [45] proposed higher-order lexical semantic models, which estimates similarity between terms by considering paths of length more than 1 on term-term similarity graph. An alternative strategy to overcome the sparseness of monolingual alignment corpora is to use the discourse relations of sentences in a text to learn term association models [86].

Questions often have some metadata, such as categories on a community question answering website. This information can be very useful for certain disambiguations, and can be encoded in the answer ranking model [130]. The structure of the web page, from which the answers are extracted can be very useful as well. Wikipedia articles have a good structure, and the information encoded there can be extracted in a text-based knowledge base, which can be used for question answering [92]. Information extraction methods can also be useful for the more general case of non-factoid question answering. For example, there is a huge number of online forums, FAQ-pages and social media, that contain question-answer pairs, which can be extracted to build a collection to query when a new question arrives [29, 59, 118, 33, 66].

Most of the approaches from TREC LiveQA 2015 combined similar question retrieval and web search techniques [113, 85, 104]. Answers to similar questions are very effective for answering new questions [85]. However, we a CQA archive doesn't have any similar questions, we have to fall back to regular web search. The idea behind the winning system of CMU [104] is to represent each answer with a pair of phrases - clue and answer text. Clue is a phrase that should be similar to the given question, and the passage that follows should be the answer to this question.

2.3 User interactions

IN PROGRESS...

An interesting approach for knowledge base construction through dialog with the user has been proposed by [52].

A very nice crowdsourcing method to obtain answers to tail information needs was proposed by [15]. Question query-url pairs are first mined from query logs, and then the wisdom of a crowd is used to extract and save answers to these questions.

Certain questions cannot be answered by machines only due to reasons such as lack of the appropriate data. Crowdsourcing was explored as one of the options to bridge this knowledge gap and assist the machine in matching, ranking and result aggregation [44]. Wisdom of a crowd can also be exploited to answer more difficult quiz questions [7].

[20] retrieves Wikipedia statements that support user answers for opinion questions.

A number of works have focused on studying and estimating different factors of user satisfaction with question answering systems [77, 71]

Interactive TREC ciQA...

[32] analyzes how clarification questions can be detected and used to improve QA performance. As far as I understand, the clarification question comes from the same user, not the system.

2.4 Summary of Related Work

Most previous work in ...

3 Structured and Unstructured Data for Factoid Question Answering

There are multiple ways to marry unstructured and structured data for joint question answering: convert structured data to unstructured format or vice versa, convert all data to certain intermediate representation or to leave them as is and link the data sources. In my thesis I focus on two approaches: relation extraction for knowledge base completion, and semantic annotation of text for hybrid question answering.

3.1 Relation Extraction for Knowledge Base Completion

The information on the web is stored in multiple different forms, such as natural language statements, tables and infoboxes, images *etc.*. In this work I focus on yet another source of information: question-answer pairs. Community question answering archives contain hundreds of millions of question and corresponding answers. Information expressed in these pairs might be hard to extract or not exist at all in other formats.

3.1.1 Relation extraction from Question-Answer pairs

PUT SUMMARY OF MY NAACL STUDENT RESEARCH WORKSHOP PAPER.

3.2 Semantic Text Annotations for Hybrid Question Answering

Converting unstructured information into structured form by extracting knowledge from text suffers from certain quality losses. Existing relation extraction tools aren't perfect, in particular due to recall losses a lot of information is left behind. Moreover, extractions contain certain level of incorrect information due to precision losses. These errors cap the upper bound on the question answering system performance.

Here I propose to utilize the synergy of structured and unstructured data, and exploit the advantages of each of them to overcome the limitations of the other. More particularly, I propose to annotate and index mentions of knowledge base entities in text documents. Such a representation induces a special kinds of edges to the knowledge base, and allows one to traverse this edges in both directions. These links open up many opportunities for QA reasoning, *e.g.* retrieving all the information about the entity by going from a mention to a KB entity, finding relations between entities by retrieving text passages that mention both of them, extracting candidate evidence by retrieving passages that mention question and answer entities along with some question terms, and so on.

— THE NEXT PIECE IS OLDER, NEED TO REVIEW

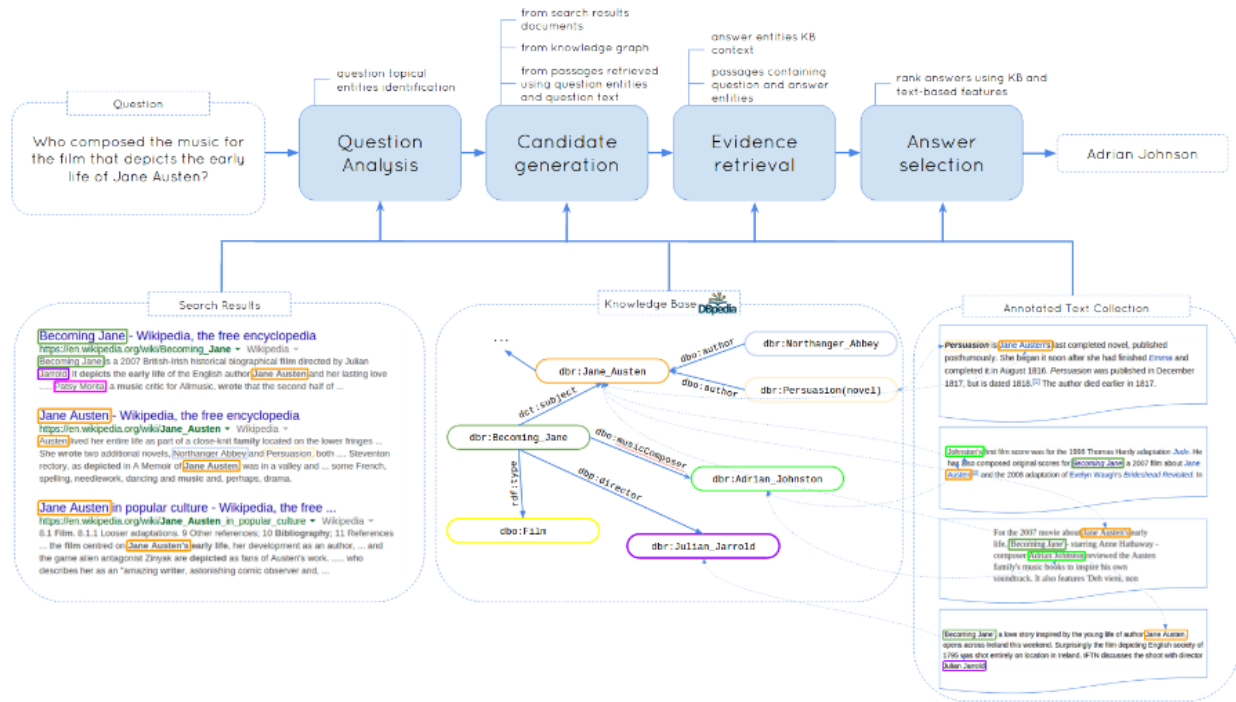


Figure 3.1: Architecture of a hybrid factoid question answering system, that uses a combination of structured knowledge base and unstructured text data

Question answering from text corpora typically starts by retrieving a set of potentially relevant documents using the question (or some transformation of the question [1]) as the query, and then extracting entities, phrases, sentences or paragraphs believed to be the answer to the question. However, the information available in the retrieved pieces of text is very limited and often not enough to decide whether it can be the answer to the given question. For example, below is one of the questions from TREC QA 2007 dataset:

“What republican senators supported the nomination of Harriet Miers to the Supreme Court?”

A candidate answer sentence *“Minority Leader Harry Reid had already offered his open support for Miers.”* mentions a senator “Harry Reid” and clearly says about his support of the nomination. However, “Harry Reid” is not a correct answer to the question because he is a member of the Democratic party. This information is not available in the answer candidate sentence, but it is present as one of the properties in Freebase: [Harry Reid, political_party, Democratic party]¹. Therefore, by looking into the knowledge available about the mentioned entities a QA system can make a better judgment about the candidate answer.

3.2.1 Approach

The architecture of the hybrid QA model I propose is presented on Figure 3.1. Here are the main stages of the question answering process:

¹Actually, in Freebase the entities are connected by a path of length 2 through a mediator node. The predicates on the path are: /government/politician/party and /government/political_party_tenure/party

- **Pre-processing:** identify mentions of KB entities in text document collection and index the documents text and mentions in separate fields
- **Topical entity identification:** search the text collection using question (or reformulated question [1]) as a query and use an approach similar to [30] to detect question topical entities
- **Candidate generation from text:** extract candidate answer (or intermediate answer) entities with evidence from the retrieved text documents using existing techniques, e.g. [100].
- **Candidate generation from KB:** explore the KB neighborhood of question topical entities and entities extracted from text documents on the previous step
- **Candidate generation from KB & Text:** use entity and text index to find entities mentioned near question topical entity and question terms in the document collection
- **KB evidence extraction:** match neighborhood of answer entities (entity type and other entities) against the question to get additional evidence
- **Text evidence extraction:** estimate the similarity between the collection text fragments mentioning question and answer entities and the question text
- **Rank candidate:** rank candidate answers using evidence extracted from the KB as well as from text

Let’s consider an example question “*Who composed the music for the film that depicted the early life of Jane Austen?*” from the QALD dataset² (Figure 3.1). Even though it’s quite easy to identify the “**Jane Austen**” entity in the question, the knowledge base (dbPedia in this example) cannot help us to determine which movie is being referred to. However, there are a lot of documents on the web, that do mention the “**Becoming Jane**” movie and say what is it about. Unfortunately, extracting the name of the composer from these documents is quite challenging, but this task can be easily accomplished by checking the value of the **musicComposer** property in the knowledge base. At the end, for each candidate answer entity, we have all the KB information and passages that mention this entity as evidence to help with the correct answer selection.

3.2.2 Evaluation

Knowledge base QA

I THINK HERE WE CAN INCLUDE OUR RESULTS ON TEXT2KB.

Most of the recent work on knowledge base question answering and semantic parsing have been evaluated on the WebQuestions dataset [11], which contains a collection of question text and correct answer entities. The questions were collected using Google Suggest API and answers crowdsourced using Amazon Mechanical Turk³. The proposed approach will be compared against the previous results⁴ on this dataset. Again, web can be used as a text collection which can be queried using Bing Search API. Relation extraction patterns can be mined using distant supervision from ClueWeb collection using publicly available dataset of Freebase annotations [46].

New factoid question answering dataset. However, WebQuestions dataset has certain limitations, e.g. questions mined using Google Suggest API have very similar structure and lexicon,

²<http://greententacle.techfak.uni-bielefeld.de/cunger/qald/>

³<http://mturk.com/>

⁴<http://goo.gl/sePBja>

and to find the answer to the mined questions users were asked to use the question entity Freebase profile page, which only include entities connected directly with a predicate or through a mediator node. Therefore most of the state-of-the-art results on the dataset use a small number of predefined logical form patterns. On the other hand CQA websites have a fraction of factoid questions with provided text answers. Here I propose to use to construct a new dataset for question answering over Freebase by selecting a subset of QnA pairs with at least one entity in question and answer and some reasonable filtering heuristics and manual validation using crowdsourcing (e.g. through Amazon Mechanical Turk). Existing systems need to be retrained and tested on the new dataset to compare against the proposed model.

Text-based QA

WE WILL ANNOTATE TREC DATASETS WITH ENTITIES!!!!

TREC QA datasets served as a benchmark for various question answering systems. Therefore, to evaluate the proposed approach for question answering over text enriched with the structured data I propose to test it on dataset derived from TREC QA and compare against existing strong baselines, including the most related approaches [40, 96]. The proposed system can use the web as the corpus and query it using Bing Search API⁵. Freebase and Reverb extractions [39] are examples of schema-based and open knowledge bases that can be used for the experiments. The metrics used for evaluation typically include accuracy and mean reciprocal rank (MRR).

For non-factoid question answering this year TREC pioneered a new question answering track - TREC LiveQA⁶, which targets questions asked by real users of Yahoo! Answers website. This year the deadline for system submission was on August 31 and my system trained on CQA QnA pairs participated in the challenge. The results will be available on the TREC Conference in November 2015. Organizers plan to continue with another TREC LiveQA task next year and this is going to be a good estimation of the effectiveness of the proposed techniques on hard real user questions.

3.3 Summary

In this section we considered two different ways of combining unstructured and structured data to improve factoid question answering. Relation extraction from question-answer pairs aims at filling some gaps in KB fact coverage, whereas semantic annotations of text documents provides a way to incorporate information available in unstructured text documents for reasoning along with KB data to improve the performance of factoid question answering.

Factoid questions represent just a part of user information needs. Many problems require more elaborate response, such as a sentence, list of instructions or in general a passage of text. Such questions are usually referred to as non-factoid questions and they will be the focus of the next Chapter.

⁵<https://datamarket.azure.com/dataset/bing/searchweb>

⁶<http://trec-liveqa.org/>

4 Non-factoid Question Answering

In this chapter I summarize the proposed work in developing a non-factoid question answering system. In particular, Section 4.1 described the general architecture of the system, and the following chapters describe the proposed improvements to different stages of the pipeline.

4.1 The architecture of the system

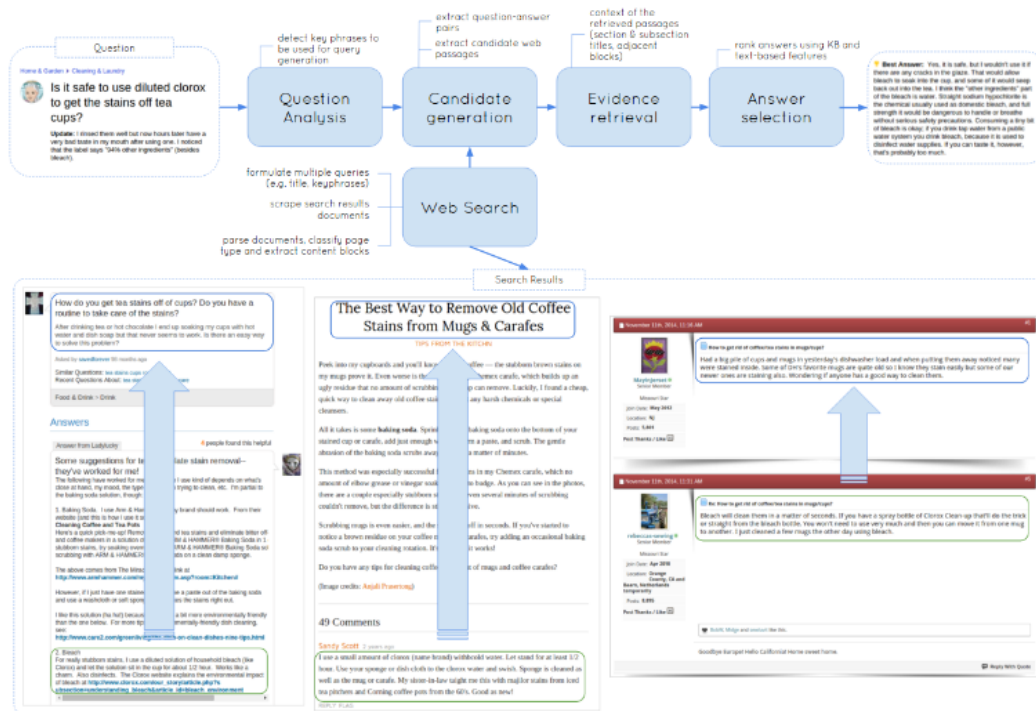


Figure 4.1: Using web page structure information for non-factoid question answering

The architecture of the system I develop is somewhat standard and contains 4 main stages: question analysis, candidate generation, evidence retrieval and answer selection. Figure 4.1 depicts the question-answering pipeline.

4.1.1 Question Analysis

Questions that users post to CQA websites often contain title and body and can be relatively long and contain multiple important and not so important context information and details. The success of the retrieval-based question answering system depends on the information it finds in a collection. Extra long search queries are not very efficient and can return few or zero results. Therefore, there is a problem of summarizing the user question. Some strategies often used include considering first part of the question, *e.g.* title of the question only. However, often the most important part

of the question isn't the title, or the title doesn't contain all the crucial information, *e.g.*

Title: *Diet please please help asap?*

Body: *I want to lose weight, at least 4 stone but I don't know what I should eat :(what should I have for breakfast lunch and dinner? Should I exercise a little? Please help!!*

To summarize the question I propose to use recent advances in the field of deep learning, in particular a model similar to [83]. In more detail, I propose to train neural network based summarization model to generate the summary of the question, which will be used as a query to retrieve similar questions. Such a model can be specifically trained to maximize retrieval performance on a collection of question-answer pairs, retrieved by systems in LiveQA TREC 2015 and labeled by NIST assessors.

4.1.2 Structure of the Web Page for Candidate Generation and Scoring

The diversity and complexity of non-factoid questions pose additional challenges for automatic question answering systems. To answer such questions a system often needs to provide a whole paragraph of text, *e.g.* TREC LiveQA'15 limits the answers to a maximum of 1000 characters. Therefore, a candidate answer becomes much longer, which requires additional attention on candidate generation and ranking stages.

Existing techniques usually extract one or more consecutive sentences not exceeding the maximum answer length, pool them together and rank using certain feature representation, by large ignoring the context information from the page where the answer was extracted from. In the system I develop I propose to utilize the structure of the web page for both candidate generation and scoring.

To generate better candidates I'm going to utilize the structure of retrieved web pages. The previous analysis showed [85], that many search results retrieved for the question are FAQs, forums or other community question answering websites. From such resources it's beneficial to extract question answer pairs, which can be done either by designing wrapper for popular websites, utilizing semantic annotations, such as <https://schema.org> or by applying some of the automatic question-answer extraction methods [29]. For other resources, page segmentation techniques, similar to [26], can split a page into semantically information blocks, which will help to make sure that candidate don't contain disjoint and unreadable information.

After a candidate passage is extracted, we need to score it as a potential answer. Often, a passage taken out of context is hard to understand even for human. Therefore, I propose to include the context information, *i.e.* some features, representing the web page where the passage was taken from and its location there (*e.g.* adjacent passages, which are often related [119]).

4.1.3 Answer Summarization

Unlike factoid questions, where evidence from all retrieved passages is usually aggregates, a traditional non-factoid question-answering system simply returns the top scoring passage as the answer. However, such an answer can contain a lot of redundant or irrelevant information, whereas other good candidate passages may contain supplemental information or different opinion. Therefore, an idea to summarize passages and generate the final answer sounds natural in this scenario. The winning approach from TREC LiveQA'15 included a combination strategy, that simply put to-

gether multiple top scoring passages while the answer doesn't exceed the maximum length [104]. In my thesis I'm planning to explore the answer summarization problem in more detail. More specifically, I propose to explore deep learning techniques, that lie in the core of recent successes in text generation, *e.g.* image caption generation and machine translation models [8, 114]. Answer summarization problem can be posed as answer generation problem using recurrent neural network, that has information about the question and retrieved passages, and it can be trained using existing CQA question-answer pairs and retrieved passages, that can be assumed as the source of the answers. Alternative and simpler approach is to solve this problem as sequential answer selection problem, where a model is trained to predict the next sentence in the answer. Such a model can be trained on a collection of questions and answer sentences, which will provide the model information on both answer discourse coherence and relevance.

4.2 Evaluation

Evaluation of complete non-factoid question answering systems is complicated due to the variability of answer language, the quality of which is impossible to estimate using manually created answer patterns (as is the case for factoid TREC QA dataset). A manual judgment of answers is needed, and luckily TREC LiveQA 2016 is offering such an opportunity. The model I'm developing will participate in the shared task, which will allow us to evaluate it against other competing approaches.

The analysis of individual components can be performed using labeled data from TREC LiveQA 2015, which includes passages extracted from Yahoo! Answers as well as regular web documents. In more detail, the performance of the answer summarization module will be estimated by similar questions retrieval performance, *e.g.* Precision @ N since we are interested in retrieving more relevant answers in TopN rather than good ranking within the retrieved group. To make results reproducible a collection of Yahoo! Answers QnA pairs from the WebScope¹ and Lucene IR library will be used for question retrieval.

The dataset to evaluate the effectiveness of using web page structure for answer scoring will be derived from TREC LiveQA 2015 labels, which include a big number of passages, that were generated from regular web pages. Therefore, the problem of evaluation of answer scoring methods can be posed as passage ranking problem and metrics, such as Precision@1, can be used as a quality measure.

Finally, the answer summarization module is the most difficult to evaluate, because it's result is a free form text, the quality of which needs to be manually labeled. Therefore, for this task I will refer to the wisdom of a crowd and use Amazon Mechanical Turk to label the quality of answers and compare to the top scoring passages for the same question.

4.3 Summary

The proposed research directions target various aspects of non-factoid question answering pipeline and can help improve both precision and recall of existing systems. The results of this work can help to move in the direction of systems, that produce a natural language response by analyzing

¹<http://webscope.sandbox.yahoo.com/catalog.php?datatype=1>

the information available on the web. In the next chapter, I will describe the proposed work and results on interactions between a question answering system and users.

5 Human Interaction with Question Answering Systems

Modern automatic question answering systems are still far from AI machines, that we often imagine or see in the movies. Many user information needs are still left unanswered by existing automatic question answering systems. For example, only 36% of answers of a winning approach from TREC LiveQA 2015 shared task were judged good or excellent. Therefore, in this chapter I will discuss the research on interactions between a human and a question answering system.

More specifically, Section 5.1 describes results on studying the effect of hints on user behavior for solving complex informational tasks. In Section 5.3 I propose research on using crowdsourcing for improving the performance of a question answering system. And Section 5.2 shifts the focus towards a dialogue between a user and a QA system, in particular towards clarification questions, that are often needed to understand user's question.

5.1 Search Hints for Complex Informational Tasks

SUMMARIZE SEARCH HINTS PAPER.

5.2 Clarification Questions

Today, we are observing a possible shift towards natural language interfaces, which will enable richer interaction between a human and computer, in particular question answering. Most of the existing systems are one-sided, *i.e.* they operate by returning an answer in a response to a user question. A richer model will inevitably lead to a dialogue rather than request-response kind of communication. There are many practical aspects of maintaining a dialogue for question-answering. For example, many questions that user ask allow multiple interpretations, *e.g.* "FIND GOOD EXAMPLE". In such cases a machine could ask a clarification questions to disambiguate certain aspects of the question and provide her with a reasonable response, instead of returning something non-relevant in the first place and force the user to think why did it happen.

Of course, this task is very challenging, and as a first step I propose to study what kind of clarifications do real people typically ask.

5.3 Using the Wisdom of a Crowd for Question Answering

This is our experiment for LiveQA answer crowdsourcing.

6 Summary and Discussion

In my thesis I propose several pieces of work towards improving user satisfaction with question answering systems. I plan to consider factoid and non-factoid questions, as usually classified in the community, separately, because certain techniques and data sources are most useful for one type of questions and not another.

The research I'm propose to conduct for improving factoid question answering targets a problem of combining information available in different data sources, *i.e.* structured knowledge bases and unstructured text documents. Semantic annotations of entity mentions in documents create additional connections between knowledge base entities, which should improve KB coverage and allow a system to answer more questions and with better precision. However, such an approach have certain potential limitations.

- A set of additional links for some entities is likely to be big, which makes it impossible to explore them all. Information extraction approaches on the other hand aggregate information over the whole collection, although the extraction process itself brings extra noise.
-

For non-factoid question answering I propose several improvements, targeting different stages of the question answering process. Question to query generation neural network model, trained to improve the document retrieval performance, should increase the recall of the question answering system by identifying the key phrases that needed to be search for. The answer passage scoring module should achieve a better precision by analyzing the structure of the answer origin web page, and detecting question-answer pairs and other key structural elements. Finally, the proposed direction in automatic answer summarization is a way to increase the quality of answers by combining evidence from multiple different data sources, possibly providing additional information and alternative opinions. Possible limitations of the proposed directions and approaches are:

- question to query generation model can be retrieval engine specific, which may force the model to be retrained after certain changes in the retrieval algorithm. An alternative strategy is to integrate a similar question summarization module into the retrieval engine itself.
- web page structure?...
- summarization?... The problem is that it might not work well...

Finally, I touch a user aspect of question answering, in particular user assistance with hints in case a system failed to respond to user information needs, clarification questions, which is one of the first steps in dialog-based question answering and finally using the wisdom of a crowd to improve the performance of question answering systems.

- Hints
- Clarifications
- Crowdsourcing

Bibliography

- [1] E. Agichtein, S. Lawrence, and L. Gravano. Learning search engine specific query transformations for question answering. In *Proceedings of the Tenth International World Wide Web Conference, WWW 10*, pages 169–178, 2001.
- [2] D. Ahn, V. Jijkoun, G. Mishne, K. Müller, M. de Rijke, and K. Schlobach. Using wikipedia at the trec qa track. In *The Thirteenth Text Retrieval Conference (TREC 2004)*, 2005.
- [3] A. M. N. Allam and M. H. Haggag. The question answering systems: A survey. *International Journal of Research and Reviews in Information Sciences (IJRRIS)*, 2(3), 2012.
- [4] A. Andrenucci and E. Snieders. Automated question answering: Review of the main approaches. In *null*, pages 514–519. IEEE, 2005.
- [5] I. Androutsopoulos, G. D. Ritchie, and P. Thanisch. Natural language interfaces to databases—an introduction. *Natural language engineering*, 1(01):29–81, 1995.
- [6] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, and Z. Ives. *Dbpedia: A nucleus for a web of open data*. Springer, 2007.
- [7] B. I. Aydin, Y. S. Yilmaz, Y. Li, Q. Li, J. Gao, and M. Demirbas. Crowdsourcing for multiple-choice question answering. In *AAAI*, pages 2946–2953, 2014.
- [8] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. *ICLR*, 2015.
- [9] H. Bast and E. Haussmann. More accurate question answering on freebase. In *CIKM*, 2015.
- [10] P. Baudiš. Yodaqa: a modular question answering system pipeline. In *POSTER 2015-19th International Student Conference on Electrical Engineering*, pages 1156–1165, 2015.
- [11] J. Berant, A. Chou, R. Frostig, and P. Liang. Semantic parsing on freebase from question-answer pairs. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1533–1544, 2013.
- [12] J. Berant and P. Liang. Semantic parsing via paraphrasing. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL 2014*, pages 1415–1425, 2014.
- [13] J. Berant and P. Liang. Imitation learning of agenda-based semantic parsers. *Transactions of the Association for Computational Linguistics*, 3:545–558, 2015.
- [14] D. Bernhard and I. Gurevych. Combining lexical semantic resources with question & answer archives for translation-based answer finding. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 2-Volume 2*, pages 728–736. Association for Computational Linguistics, 2009.
- [15] M. S. Bernstein, J. Teevan, S. Dumais, D. Liebling, and E. Horvitz. Direct answers for search queries in the long tail. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pages 237–246. ACM, 2012.
- [16] M. W. Bilotti, P. Ogilvie, J. Callan, and E. Nyberg. Structured retrieval for question answering. In *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 351–358. ACM, 2007.

- [17] K. Bollacker, C. Evans, P. Paritosh, T. Sturge, and J. Taylor. Freebase: A collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM International Conference on Management of Data, SIGMOD '08*, pages 1247–1250, New York, NY, USA, 2008. ACM.
- [18] A. Bordes, S. Chopra, and J. Weston. Question answering with subgraph embeddings. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014*, pages 615–620, 2014.
- [19] A. Bordes, N. Usunier, S. Chopra, and J. Weston. Large-scale simple question answering with memory networks. *arXiv preprint arXiv:1506.02075*, 2015.
- [20] L. Braunstein, O. Kurland, D. Carmel, I. Szpektor, and A. Shtok. Supporting human answers for advice-seeking questions in cqa sites. In *Advances in Information Retrieval*, pages 129–141. Springer, 2016.
- [21] E. Brill, S. Dumais, and M. Banko. An analysis of the askmsr question-answering system. In *Proceedings of the ACL-02 conference on Empirical methods in natural language processing-Volume 10*, pages 257–264. Association for Computational Linguistics, 2002.
- [22] E. Brill, J. Lin, M. Banko, S. Dumais, and A. Ng. Data-intensive question answering. In *Proceedings of TREC 2001*, January 2001.
- [23] D. Buscaldi and P. Rosso. Mining knowledge from wikipedia for the question answering task. In *Proceedings of the International Conference on Language Resources and Evaluation*, pages 727–730, 2006.
- [24] M. J. Cafarella, A. Halevy, D. Z. Wang, E. Wu, and Y. Zhang. Webtables: Exploring the power of tables on the web. *Proc. VLDB Endow.*, 1(1):538–549, Aug. 2008.
- [25] M. J. Cafarella, J. Madhavan, and A. Halevy. Web-scale extraction of structured data. *SIGMOD Rec.*, 37(4):55–61, Mar. 2009.
- [26] D. Cai, S. Yu, J.-R. Wen, and W.-Y. Ma. Vips: a visionbased page segmentation algorithm. Technical report, Microsoft technical report, MSR-TR-2003-79, 2003.
- [27] Q. Cai and A. Yates. Large-scale semantic parsing via schema matching and lexicon extension. In *ACL (1)*, pages 423–433. Citeseer, 2013.
- [28] Q. Cai and A. Yates. Large-scale semantic parsing via schema matching and lexicon extension. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics, ACL 2013*, pages 423–433, 2013.
- [29] G. Cong, L. Wang, C.-Y. Lin, Y.-I. Song, and Y. Sun. Finding question-answer pairs from online forums. In *Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 467–474. ACM, 2008.
- [30] M. Cornolti, P. Ferragina, M. Ciaramita, H. Schütze, and S. Rüd. The smaph system for query entity recognition and disambiguation. In *Proceedings of the first international workshop on Entity recognition & disambiguation*, 2014.
- [31] H. T. Dang, D. Kelly, and J. J. Lin. Overview of the trec 2007 question answering track. In *TREC*, volume 7, page 63. Citeseer, 2007.
- [32] M. De Boni and S. Manandhar. An analysis of clarification dialogue for question answering. In *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology-Volume 1*, pages 48–55. Association for Computational Linguistics, 2003.
- [33] S. Ding, G. Cong, C.-Y. Lin, and X. Zhu. Using conditional random fields to extract contexts and answers of questions from online forums. In *ACL*, volume 8, pages 710–718. Citeseer, 2008.
- [34] L. Dong, F. Wei, M. Zhou, and K. Xu. Question answering over freebase with multi-column convolutional neural networks. In *Proceedings of the 53rd Annual Meeting of the Association for Com-*

- putational Linguistics and the 7th International Joint Conference on Natural Language Processing*, volume 1, pages 260–269, 2015.
- [35] X. Dong, E. Gabrilovich, G. Heitz, W. Horn, N. Lao, K. Murphy, T. Strohmman, S. Sun, and W. Zhang. Knowledge vault: A web-scale approach to probabilistic knowledge fusion. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '14, pages 601–610, New York, NY, USA, 2014. ACM.
 - [36] H. Duan, Y. Cao, C.-Y. Lin, and Y. Yu. Searching questions by identifying question topic and question focus. In *ACL*, pages 156–164, 2008.
 - [37] O. Etzioni. Search needs a shake-up. *Nature*, 476(7358):25–26, 2011.
 - [38] O. Etzioni, M. Banko, S. Soderland, and D. S. Weld. Open information extraction from the web. *Commun. ACM*, 51(12):68–74, Dec. 2008.
 - [39] A. Fader, S. Soderland, and O. Etzioni. Identifying relations for open information extraction. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing, EMNLP 2011*, pages 1535–1545, 2011.
 - [40] A. Fader, L. Zettlemoyer, and O. Etzioni. Open question answering over curated and extracted knowledge bases. In *Proceedings of the 20th ACM International Conference on Knowledge Discovery and Data Mining*, KDD '14, pages 1156–1165, New York, NY, USA, 2014. ACM.
 - [41] A. Fader, L. S. Zettlemoyer, and O. Etzioni. Paraphrase-driven learning for open question answering. In *ACL*. Citeseer, 2013.
 - [42] D. Ferrucci, E. Brown, J. Chu-Carroll, J. Fan, D. Gondek, A. A. Kalyanpur, A. Lally, J. W. Murdock, E. Nyberg, J. Prager, et al. Building watson: An overview of the deepqa project. *AI magazine*, 31(3):59–79, 2010.
 - [43] D. Ferrucci, E. Brown, J. Chu-Carroll, J. Fan, D. Gondek, A. A. Kalyanpur, A. Lally, J. W. Murdock, E. Nyberg, J. Prager, et al. This is watson. *IBM Journal of Research and Development*, 56, 2012.
 - [44] M. J. Franklin, D. Kossmann, T. Kraska, S. Ramesh, and R. Xin. Crowddb: Answering queries with crowdsourcing. In *Proceedings of the 2011 ACM SIGMOD International Conference on Management of Data*, SIGMOD '11, pages 61–72, New York, NY, USA, 2011. ACM.
 - [45] D. Fried, P. Jansen, G. Hahn-Powell, M. Surdeanu, and P. Clark. Higher-order lexical semantic models for non-factoid answer reranking. *Transactions of the Association for Computational Linguistics*, 3:197–210, 2015.
 - [46] E. Gabrilovich, M. Ringgaard, and A. Subramanya. Facc1: Freebase annotation of cluweb corpora, 2013.
 - [47] B. F. Green Jr, A. K. Wolf, C. Chomsky, and K. Laughery. Baseball: an automatic question-answerer. In *Papers presented at the May 9-11, 1961, western joint IRE-AIEE-ACM computer conference*, pages 219–224. ACM, 1961.
 - [48] P. Gupta and V. Gupta. A survey of text question answering techniques. *International Journal of Computer Applications*, 53(4), 2012.
 - [49] R. Gupta, A. Halevy, X. Wang, S. E. Whang, and F. Wu. Biperpedia: An ontology for search applications. *Proc. VLDB Endow.*, 7(7):505–516, Mar. 2014.
 - [50] M. Heilman and N. A. Smith. Tree edit models for recognizing textual entailments, paraphrases, and answers to questions. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 1011–1019. Association for Computational Linguistics, 2010.
 - [51] L. Hirschman and R. Gaizauskas. Natural language question answering: the view from here. *natural language engineering*, 7(4):275–300, 2001.
 - [52] B. Hixon, P. Clark, and H. Hajishirzi. Learning knowledge graphs for question answering through

- conversational dialog. In *Proceedings of the the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Denver, Colorado, USA*, 2015.
- [53] E. Hovy, U. Hermjakob, and D. Ravichandran. A question/answer typology with surface text patterns. In *Proceedings of the Second International Conference on Human Language Technology Research, HLT '02*, pages 247–251, San Francisco, CA, USA, 2002. Morgan Kaufmann Publishers Inc.
 - [54] E. H. Hovy, L. Gerber, U. Hermjakob, M. Junk, and C.-Y. Lin. Question answering in webclopedia. In *TREC*, volume 52, pages 53–56, 2000.
 - [55] E. H. Hovy, U. Hermjakob, and C.-Y. Lin. The use of external knowledge of factoid qa. In *TREC*, volume 2001, pages 644–52, 2001.
 - [56] A. Ittycheriah, M. Franz, and S. Roukos. Ibm’s statistical question answering system-trec-10. In *TREC*, 2001.
 - [57] M. Iyyer, J. L. Boyd-Graber, L. M. B. Claudino, R. Socher, and H. Daumé III. A neural network for factoid question answering over paragraphs. In *EMNLP*, pages 633–644, 2014.
 - [58] J. Jeon, W. B. Croft, and J. H. Lee. Finding similar questions in large question and answer archives. In *Proceedings of the 14th ACM International Conference on Information and Knowledge Management, CIKM '05*, pages 84–90, New York, NY, USA, 2005. ACM.
 - [59] V. Jijkoun and M. de Rijke. Retrieving answers from frequently asked questions pages on the web. In *Proceedings of the 14th ACM International Conference on Information and Knowledge Management, CIKM '05*, pages 76–83, New York, NY, USA, 2005. ACM.
 - [60] V. Jijkoun, M. De Rijke, and J. Mur. Information extraction for question answering: Improving recall through syntactic patterns. In *Proceedings of the 20th international conference on Computational Linguistics*, page 1284. Association for Computational Linguistics, 2004.
 - [61] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3128–3137, 2015.
 - [62] M. Keikha, J. H. Park, W. B. Croft, and M. Sanderson. Retrieving passages and finding answers. In *Proceedings of the 2014 Australasian Document Computing Symposium*, page 81. ACM, 2014.
 - [63] O. Kolomiyets and M.-F. Moens. A survey on question answering technology from an information retrieval perspective. *Information Sciences*, 181(24):5412–5434, Dec. 2011.
 - [64] N. Kushmerick. *Wrapper induction for information extraction*. PhD thesis, University of Washington, 1997.
 - [65] C. Kwok, O. Etzioni, and D. S. Weld. Scaling question answering to the web. *ACM Transactions on Information Systems (TOIS)*, 19(3):242–262, 2001.
 - [66] B. Li, X. Si, M. R. Lyu, I. King, and E. Y. Chang. Question identification on twitter. In *Proceedings of the 20th ACM international conference on Information and knowledge management*, pages 2477–2480. ACM, 2011.
 - [67] X. Li and D. Roth. Learning question classifiers. In *19th International Conference on Computational Linguistics, COLING 2002*, 2002.
 - [68] X. Li and D. Roth. Learning question classifiers: the role of semantic information. *Natural Language Engineering*, 12(3):229–249, 2006.
 - [69] J. Lin. An exploration of the principles underlying redundancy-based factoid question answering. *ACM Transactions on Information Systems (TOIS)*, 25(2):6, 2007.
 - [70] J. Lin and B. Katz. Question answering from the web using knowledge annotation and knowledge mining techniques. In *Proceedings of the twelfth international conference on Information and knowledge management*, pages 116–123. ACM, 2003.

- [71] Y. Liu, J. Bian, and E. Agichtein. Predicting information seeker satisfaction in community question answering. In *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '08, pages 483–490, New York, NY, USA, 2008. ACM.
- [72] Y. Liu, S. Li, Y. Cao, C.-Y. Lin, D. Han, and Y. Yu. Understanding and summarizing answers in community-based question answering services. In *Proceedings of the 22Nd International Conference on Computational Linguistics - Volume 1*, COLING '08, pages 497–504, Stroudsburg, PA, USA, 2008. Association for Computational Linguistics.
- [73] F. Mahdisoltani, J. Biega, and F. Suchanek. Yago3: A knowledge base from multilingual wikipedias. In *7th Biennial Conference on Innovative Data Systems Research*. CIDR Conference, 2014.
- [74] C. Malon and B. Bai. Answer extraction by recursive parse tree descent. *ACL 2013*, page 110, 2013.
- [75] G. A. Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.
- [76] M. Mintz, S. Bills, R. Snow, and D. Jurafsky. Distant supervision for relation extraction without labeled data. In *ACL 2009, Proceedings of the 47th Annual Meeting of the Association for Computational Linguistics and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 1003–1011, 2009.
- [77] C.-S. Ong, M.-Y. Day, and W.-L. Hsu. The measurement of user satisfaction with question answering systems. *Information & Management*, 46(7):397–403, 2009.
- [78] M. Pasca and S. Harabagiu. The informative role of wordnet in open-domain question answering. In *Proceedings of NAACL-01 Workshop on WordNet and Other Lexical Resources*, pages 138–143, 2001.
- [79] J. Prager, J. Chu-Carroll, E. W. Brown, and K. Czuba. Question answering by predictive annotation. In *Advances in Open Domain Question Answering*, pages 307–347. Springer, 2006.
- [80] J. M. Prager. Open-domain question-answering. *Foundations and trends in information retrieval*, 1(2):91–231, 2006.
- [81] V. Punyakanok, D. Roth, and W.-t. Yih. Mapping dependencies trees: An application to question answering. In *Proceedings of AI&Math 2004*, pages 1–10, 2004.
- [82] S. Reddy, M. Lapata, and M. Steedman. Large-scale semantic parsing without question-answer pairs. *TACL*, 2:377–392, 2014.
- [83] A. M. Rush, S. Chopra, and J. Weston. A neural attention model for abstractive sentence summarization. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 379–389, Lisbon, Portugal, September 2015. Association for Computational Linguistics.
- [84] C. d. Santos, M. Tan, B. Xiang, and B. Zhou. Attentive pooling networks. *arXiv preprint arXiv:1602.03609*, 2016.
- [85] D. Savenkov. Ranking answers and web passages for non-factoid question answering: Emory university at trec liveqa. In *Proceedings of TREC*, 2015.
- [86] R. Sharp, P. Jansen, M. Surdeanu, and P. Clark. Spinning straw into gold: Using free text to train monolingual alignment models for non-factoid question answering. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics-Human Language Technologies (NAACL HLT)*, 2015.
- [87] D. Shen, G.-J. M. Kruijff, and D. Klakow. Exploring syntactic relation patterns for question answering. In *Natural Language Processing-IJCNLP 2005*, pages 507–518. Springer, 2005.
- [88] E. H. Shortliffe and B. G. Buchanan. A model of inexact reasoning in medicine. *Mathematical biosciences*, 23(3):351–379, 1975.

- [89] A. Shtok, G. Dror, Y. Maarek, and I. Szpektor. Learning from the past: Answering new questions with past answers. In *Proceedings of the 21st International Conference on World Wide Web, WWW '12*, pages 759–768, New York, NY, USA, 2012. ACM.
- [90] R. F. Simmons. Answering english questions by computer: A survey. *Communications of ACM*, 8(1):53–70, Jan. 1965.
- [91] R. F. Simmons. Natural language question-answering systems: 1969. *Commun. ACM*, 13(1):15–30, Jan. 1970.
- [92] P. Sondhi and C. Zhai. Mining semi-structured online knowledge bases to answer natural language questions on community qa websites. In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*, pages 341–350. ACM, 2014.
- [93] R. Soricut and E. Brill. Automatic question answering using the web: Beyond the factoid. *Information Retrieval*, 9(2):191–206, 2006.
- [94] M. M. Soubbotin and S. M. Soubbotin. Patterns of potential answer expressions as clues to the right answers. In *TREC*, 2001.
- [95] F. M. Suchanek, G. Kasneci, and G. Weikum. Yago: a core of semantic knowledge. In *Proceedings of the 16th international conference on World Wide Web*, pages 697–706. ACM, 2007.
- [96] H. Sun, H. Ma, W.-t. Yih, C.-T. Tsai, J. Liu, and M.-W. Chang. Open domain question answering via semantic enrichment. In *Proceedings of the 24th International Conference on World Wide Web, WWW '15*, pages 1045–1055, Republic and Canton of Geneva, Switzerland, 2015. International World Wide Web Conferences Steering Committee.
- [97] M. Surdeanu, M. Ciaramita, and H. Zaragoza. Learning to rank answers to non-factoid questions from web collections. *Computational Linguistics*, 37(2):351–383, 2011.
- [98] M. Tan, B. Xiang, and B. Zhou. Lstm-based deep learning models for non-factoid answer selection. *arXiv preprint arXiv:1511.04108*, 2015.
- [99] W. tau Yih, M.-W. Chang, X. He, and J. Gao. Semantic parsing via staged query graph generation: Question answering with knowledge base. In *Proceedings of the Joint Conference of the 53rd Annual Meeting of the ACL and the 7th International Joint Conference on Natural Language Processing of the AFNLP*. ACL Association for Computational Linguistics, July 2015.
- [100] C. Tsai, W.-t. Yih, and C. Burges. Web-based question answering: Revisiting askmsr. Technical report, Technical Report MSR-TR-2015-20, Microsoft Research, 2015.
- [101] C. Unger, A. Freitas, and P. Cimiano. An introduction to question answering over linked data. In *Reasoning Web. Reasoning on the Web in the Big Data Era*, pages 100–140. Springer, 2014.
- [102] E. M. Voorhees. The trec question answering track. *Natural Language Engineering*, 7(04):361–378, 2001.
- [103] D. Vrandečić and M. Krötzsch. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85, 2014.
- [104] D. Wang and E. Nyberg. Cmu oaqa at trec 2015 liveqa: Discovering the right answer with clues. In *Proceedings of TREC*, 2015.
- [105] D. Wang and E. Nyberg. A long short-term memory model for answer sentence selection in question answering. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL 2015*, pages 707–712, 2015.
- [106] M. Wang. A survey of answer extraction techniques in factoid question answering. *Computational Linguistics*, 1(1), 2006.
- [107] M. Wang and C. D. Manning. Probabilistic tree-edit models with structured latent variables for textual entailment and question answering. In *Proceedings of the 23rd International Conference on Computational Linguistics*, pages 1164–1172. Association for Computational Linguistics, 2010.

- [108] M. Wang, N. A. Smith, and T. Mitamura. What is the jeopardy model? a quasi-synchronous grammar for qa. In *EMNLP-CoNLL*, volume 7, pages 22–32, 2007.
- [109] Z. Wang and A. Ittycheriah. Faq-based question answering via word alignment. *arXiv preprint arXiv:1507.02628*, 2015.
- [110] Z. Wang, S. Yan, H. Wang, and X. Huang. Large-scale question answering with joint embedding and proof tree decoding. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, pages 1783–1786. ACM, 2015.
- [111] R. Wilensky, D. N. Chin, M. Luria, J. Martin, J. Mayfield, and D. Wu. The berkeley unix consultant project. *Computational Linguistics*, 14(4):35–84, 1988.
- [112] W. A. Woods and R. Kaplan. Lunar rocks in natural english: Explorations in natural language question answering. *Linguistic structures processing*, 5:521–569, 1977.
- [113] G. Wu and M. Lan. Leverage web-based answer retrieval and hierarchical answer selection to improve the performance of live question answering. In *Proceedings of TREC*, 2015.
- [114] K. Xu, J. Ba, R. Kiros, A. Courville, R. Salakhutdinov, R. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. *arXiv preprint arXiv:1502.03044*, 2015.
- [115] K. Xu, Y. Feng, S. Reddy, S. Huang, and D. Zhao. Enhancing freebase question answering using textual evidence. *arXiv preprint arXiv:1603.00957*, 2016.
- [116] K. Xu, S. Zhang, Y. Feng, and D. Zhao. Answering natural language questions via phrasal semantic parsing. In *Natural Language Processing and Chinese Computing*, pages 333–344. Springer, 2014.
- [117] M. Yahya, K. Berberich, S. Elbassuoni, and G. Weikum. Robust question answering over the web of linked data. In *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management*, pages 1107–1116. ACM, 2013.
- [118] J.-M. Yang, R. Cai, Y. Wang, J. Zhu, L. Zhang, and W.-Y. Ma. Incorporating site-level knowledge to extract structured data from web forums. In *Proceedings of the 18th International Conference on World Wide Web*, WWW ’09, pages 181–190, New York, NY, USA, 2009. ACM.
- [119] L. Y. Yang, Q. Ai, D. Spina, R.-C. Chen, L. Pang, W. B. Croft, J. Guo, and F. Scholer. Beyond factoid qa: Effective methods for non-factoid answer sentence retrieval. In *Proceedings of ECIR’16*, 2016.
- [120] Y. Yang, W.-t. Yih, and C. Meek. Wikiqa: A challenge dataset for open-domain question answering. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 2013–2018. Citeseer, 2015.
- [121] X. Yao. Lean question answering over freebase from scratch. In *Proceedings of NAACL Demo*, 2015.
- [122] X. Yao, J. Berant, and B. Van Durme. Freebase qa: Information extraction or semantic parsing? *ACL 2014*, page 82, 2014.
- [123] X. Yao and B. V. Durme. Information extraction over structured data: Question answering with freebase. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL 2014*, pages 956–966, 2014.
- [124] X. Yao, B. Van Durme, C. Callison-Burch, and P. Clark. Answer extraction as sequence tagging with tree edit distance. In *HLT-NAACL*, pages 858–867. Citeseer, 2013.
- [125] X. Yao, B. Van Durme, and P. Clark. Automatic coupling of answer extraction and information retrieval. In *ACL (2)*, pages 159–165, 2013.
- [126] W.-t. Yih, X. He, and C. Meek. Semantic parsing for single-relation question answering. In *ACL (2)*, pages 643–648. Citeseer, 2014.
- [127] P. Yin, N. Duan, B. Kao, J. Bao, and M. Zhou. Answering questions with complex semantic constraints on open knowledge bases. In *Proceedings of the 24th ACM International on Conference*

- on Information and Knowledge Management*, pages 1301–1310. ACM, 2015.
- [128] L. Yu, K. M. Hermann, P. Blunsom, and S. Pulman. Deep learning for answer sentence selection. *arXiv preprint arXiv:1412.1632*, 2014.
 - [129] J. M. Zelle and R. J. Mooney. Learning to parse database queries using inductive logic programming. In *Proceedings of the national conference on artificial intelligence*, pages 1050–1055, 1996.
 - [130] G. Zhou, T. He, J. Zhao, and P. Hu. Learning continuous word embedding with metadata for question retrieval in community question answering. In *Proceedings of ACL*, pages 250–259, 2015.