

When a Knowledge Base Is Not Enough: Question Answering over Knowledge Bases with External Text Data

Denis Savenkov
Emory University
dsavenk@emory.edu

Eugene Agichtein
Emory University
eugene@mathcs.emory.edu

ABSTRACT

One of the major challenges for knowledge base question answering systems (KBQA) is to translate a natural language question to knowledge base (KB) entities and predicates. Previous systems have used a limited amount of training data to learn a lexicon that is later used for question answering. This approach does not make use of other potentially relevant text data, outside the KB, which could enrich the available information. We introduce a new system, Text2KB, that connects a KB with external text. Specifically, we revisit different phases in the KBQA process and demonstrate that text resources improve question interpretation, candidate generation and ranking. Starting with the best publicly available system, Text2KB utilizes web search results, community question answering and general text document collection data, to detect question topic entities and enrich the features of the candidates derived from the KB. Text2KB significantly improves on the initial KBQA system, and reaches the best known state of the art performance on a popular WebQuestions knowledge base question answering dataset. The results and insights developed in this work are both practically useful, and can guide future efforts on combining textual and structured KB data for question answering.

1. INTRODUCTION

It has long been recognized that searchers prefer concise and specific answers, rather than lists of document results. In particular, factual, or factoid questions, have been an active focus of research for decades due to both practical importance and relatively objective evaluation criteria. As a particularly important example, a large proportion of Web search queries are looking for entities or their attributes [19], a setting on which we focus in this work.

Two relatively separate approaches for Question Answering (QA) have emerged: text-centric, or Text-QA and knowledge base-centric, or KBQA. In the more traditional, text-QA approach, QA systems used text document collections to

retrieve passages relevant to a question and to extract candidate answers [14]. Unfortunately, an unstructured text passage does not provide explicit information about the candidate entities, which has to be inferred from the question text. The KBQA approach, which evolved from the database community, relies on large scale knowledge bases, such as dbPedia [1], Freebase [9] and WikiData [24], which store a vast amount of general knowledge about different kinds of entities. This information, encoded as [subject, predicate, object] RDF triples, can be effectively queried using structured query languages, such as SPARQL.

Both approaches need to eventually deal with natural language questions, in which information needs are expressed by the vast majority of users. While question understanding is difficult in itself, this setting is particularly challenging for KBQA systems, as it requires a translation of a text question into a structured query language. That is challenging for a number of reasons, including the complexity of a KB schema, and many differences between natural language and knowledge representations. For example, Figure 1 gives a SPARQL query that retrieves the answer to a relatively simple question “*who was the president of the Dominican Republic in 2010?*” from Freebase.

Any KBQA systems must address three challenges, namely question entity identification to anchor the query process; candidate answer entities identification; and candidate ranking. We will show that these challenges can be alleviated by the appropriate use of external textual data.

The first problem that a KBQA system faces is question entity identification. The performance of the whole system greatly depends on this stage [27], because it seeds the answer candidate search process. Question text is often quite short, may contain typos and other problems, that complicate question entity identification. Existing approaches are usually based on dictionaries that contain entity names, aliases and some other phrases, which were used to refer to entities [21]. These dictionaries are often noisy and incomplete, *e.g.*, to answer the question “*what year did tut became king?*” a system needs to detect a mention “*tut*”, which refers to the entity *Tutankhamun*. A mapping “*tut*” → *Tutankhamun* is missing in the dictionary used by one of the state of the art systems and therefore it couldn’t answer this question correctly. Instead of increasing the dictionary size we propose to use web search results to find variations of question entity names, which can be easier to link to a KB. This idea has been shown effective in entity linking in web search queries [12], one of the tracks on the Entity Recognition and Dis-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGIR '16, July 17-21, 2016, Pisa, Italy

© 2016 ACM. ISBN 978-1-4503-4069-4/16/07...\$15.00

DOI: <http://dx.doi.org/10.1145/2911451.2911536>

```

SELECT DISTINCT ?name {
  :m.027rn :government.governmental_jurisdiction.governing_officials ?gov_position .
  ?gov_position :government.government_position_held.basic_title :m.060c4 .
  ?gov_position :government.government_position_held.office_holder ?president .
  ?gov_position :government.government_position_held.from ?from_date .
  ?gov_position :government.government_position_held.to ?to_date .
  FILTER ( xsd:date(?from_date) <= "2010"^^xsd:date AND
    xsd:date(?to_date) >= "2010"^^xsd:date )
  ?president :type.object.name ?name
}

```

Figure 1: SPARQL query to retrieve the answer to the question “*who was the president of the dominican republic in 2010?*”

ambiguation Challenge 2014¹. Figure 2 presents web search results for the query “*what year did tut became king?*”, which shows that indeed many documents mention the full name of the entity.

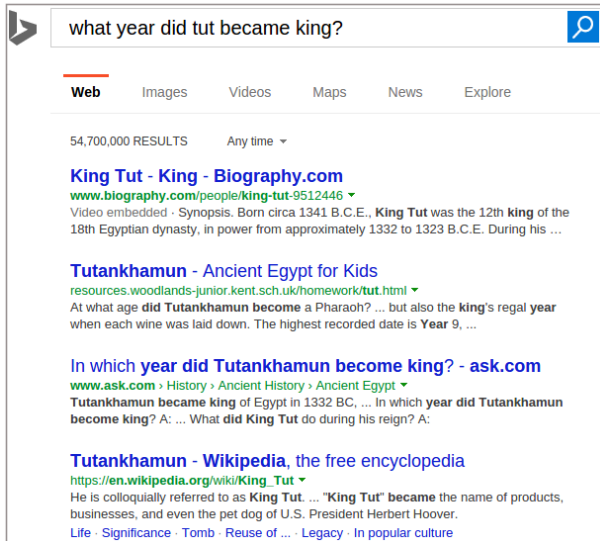


Figure 2: Search results for the question “*what year did tut became king?*”, which mention both the full name of the king and the correct answer to the question

After question entities have been identified, answer candidates need to be generated and ranked to select the best answer. A candidate query contains one or multiple KB predicates, which should somehow correspond to words and phrases in the question. Existing knowledge base question answering approaches [3, 6, 7, 8, 10, 28] rely a lexicon, learned from manually labeled training data, and supported by additional resources, such as question paraphrases [7] and weakly labeled sentences from a large text collection [29]. However, the coverage of these lexicons is limited, because such training data tends to be small compared to the number of different entities and predicates in a KB. By our estimate, in a popular WebQuestions KBQA dataset, the answers to ~5.5% of test questions (112 out of 2032) involve a predicate that does not appear in the training set. For example, an RDF triple [Bigos, food.dish.type_of_dish1, Stew] answers a test question “*what are bigos?*”, but there are no questions from the training set that are answered using the

same predicate. In addition, even if the training set contained an example targeting a particular KB predicate, the lexicon might not cover all the other possible ways the same information can be asked about. For example, questions “*who did jon gosselin cheat with?*” and “*who is the woman that john edwards had an affair with?*” are answered by the same KB predicate, but use different language. Therefore, seeing the first question in the training set might not help to answer the later question. On the other hand, traditional text-QA systems benefit from the redundancy of information on the Web, where the same facts are stated multiple times in many different ways [17]. This increases the chances of a good lexical match between a question and answer statements, which makes even some relatively simple counting-based techniques quite effective [11]. We propose to adapt these ideas from text-based question answering for KBQA. The right part of the Figure 3 shows web search results, a community question answering page, and text fragments mentioning pairs of entities, that can be useful to answer the question about John Edwards’ affair.

To summarize, our contributions are three-fold:

- A novel “hybrid” knowledge base question answering system, which uses both structured data from a knowledge base and unstructured natural language resources connected via entity links. Section 3 describes the architecture of our system, and Section 4 shows that this fusion improves the performance of a state of the art KBQA system.
- Novel data sources and techniques for knowledge base question answering. We introduce three techniques: enhancing question entity identification by analyzing web search results (Section 3.1); improving predicate matching by mining CQA data (Section 3.2); and improving candidate ranking by incorporating text-corpus statistics (Section 3.3).
- Comprehensive empirical analysis of our system on a popular WebQuestions benchmark, demonstrating that using additional text resources can improve the performance of a state-of-the-art KBQA system (Section 4). In addition, we conduct an extensive analysis of the system to identify promising directions for future improvements (Section 5).

Taken together, this work introduces a number of techniques of using external text that significantly improve the performance of the KBQA approach. More broadly, our work bridges the gap between Text-QA and KBQA worlds, demonstrating an important step forward towards combining unstructured and structured data for question answering.

¹<http://web-ngram.research.microsoft.com/ERD2014/>

2. KNOWLEDGE BASE QUESTION ANSWERING: AN OVERVIEW

In this section, we overview the existing approaches to Knowledge Base Question Answering (KBQA), as our approach, described in the next section, builds upon and extends some of these efforts.

Recently, KBQA systems have converged to two major approaches: *semantic parsing*, and *information extraction* (IE) [28]. The former focuses on question understanding, and attempts to parse the sentences into a semantic representation, *e.g.*, logical forms [6, 7, 8]. The latter, information extracting approaches [3, 30, 29] are based on identifying *topical entities* in the question, and then, using pre-defined templates for mapping the question to predicates, explore the neighborhood of these entities in a KB. Theoretically, semantic parsing-based systems would be capable of generating any required queries, and would apply to any question, seen or unseen in training, whereas the template-based approach is less likely to generalize. In practice, however, answers to most of the questions lie within two edge traversals within a KB, making the template (“information extraction”-based) approaches more effective.

Interestingly, one of the reasons for recent resurgence of interest in KBQA can be credited to creation of the WebQuestions dataset [6], which is large enough to allow both comprehensive evaluation, and training machine learning methods. Thus, the performance of KBQA systems has quickly improved, with the current state of the art systems using the information extraction approach, with sophisticated ranking and matching postprocessing [30]. Therefore, we chose to extend an existing KBQA system: Aquu [3]. It follows an information extraction approach to KBQA and achieves one of the highest scores among publicly available systems. However, as we will show, our approach is general and can be incorporated into other IE-based systems as well.

We will first describe an information extraction approach to KBQA in more detail using Aquu - our baseline system - as an example. Next, Section 3 will present our system Text2KB, which extends this approach by incorporating external text-based data on various stages of the question answering process.

2.1 The Aquu KBQA system

Recall, that a KBQA system first needs to identify question entities, which are used to initialize the “neighborhood” of potential answer entities in the KB that are connected to the question entities. For concreteness, consider a question from the Webquestions dataset “*who is the woman that john edwards had an affair with?*”. First, the system identifies a set of possible question entities. In our example, entity **John Edwards** with Freebase mid `/m/01651q` is the main question entity. However, Freebase contains millions of entities and it’s often hard to identify the topical ones (*e.g.*, entities **Woman** and **Affair** are also present in Freebase) or to disambiguate and choose between **John Edwards** a politician (`/m/01641q`), **John Edwards** an American sports car racing driver (`/m/06zs089`) and other people with the same name. There is even an entity with the name “**had an affair with**”². Aquu considers all spans of terms under certain conditions on POS tags and use a dictionary of names, aliases and anchor texts [21] to map phrases to po-

tential entities. Most recent systems, including Aquu, don’t disambiguate entities at this stage and keep a set of candidates along with some information about their popularities and mention scores.

On the next stage, SPARQL query candidates are generated by exploring the neighborhood of the question topical entities using a predefined set of query templates. Each query template has an entity, predicate and answer entity placeholders. Majority of the answers in WebQuestions dataset can be covered by just 3 templates (q_entity - question entity, a_entity - answer entity, cvt_node - Freebase mediator node, which represent tuples with more than 2 arguments):

```
SELECT DISTINCT ?a_entity {
  <q_entity> <predicate> ?a_entity .
}

SELECT DISTINCT ?a_entity {
  <q_entity> <predicate_1> ?cvt_node .
  ?cvt_node <predicate_2> ?a_entity .
}

SELECT DISTINCT ?a_entity {
  <q_entity_1> <predicate_1> ?cvt_node .
  ?cvt_node <predicate_2> <q_entity_2> .
  ?cvt_node <predicate_3> ?a_entity .
}
```

The first template retrieves a set of entities that are directly connected to the given question entity via a certain predicate. The second template accounts for the presence of a mediator node, that groups together arguments of a multi-argument relation. And the last template looks for cases, when multi-argument relations also mention another question entity, *e.g.*, **Captain Kirk** and **Star Trek** for the question “*who played captain kirk in star trek movie?*”.

Finally, each query candidate is represented with a set of features, that includes information about the popularity of question entities (mention frequency), entity linking scores, size of the answer list, entity and predicate token matches, predicate score based on question token n-grams, *etc.* The full list of features can be found in the original paper [3]. The final stage of the question answering process is filtering and ranking, using a trained random forest model. This ranking model, like the components of the other stages, is trained offline on the training subset of WebQuestions dataset.

2.2 Basic system extensions

In order to work with a high-performing KBQA system, before embarking on our drastic changes described in the next section, we analyzed and improved a number of details in the original Aquu system, as described here. The improved system will be used as a baseline in the experiments of Section 5. First, we noticed that since Aquu does not use information about answer entity Freebase types, in many cases it returns an answer that is type incompatible with the question: *e.g.*, state instead of county *etc.* Similarly to how relations are scored in Aquu, we decided to train a model to return a score, measuring a compatibility of the answer entities, using their notable types and question uni- and bigrams as features. A second extension introduced a new date range query template, which, based on the error analysis, was lacking in Aquu, but could be helpful. This template helps to solve the cases like “*what team did david beckham play for in 2011?*”, where we need to look at the ranges of dates to figure out in which range does the specified date falls.

²<http://www.freebase.com/m/0c0n01x>

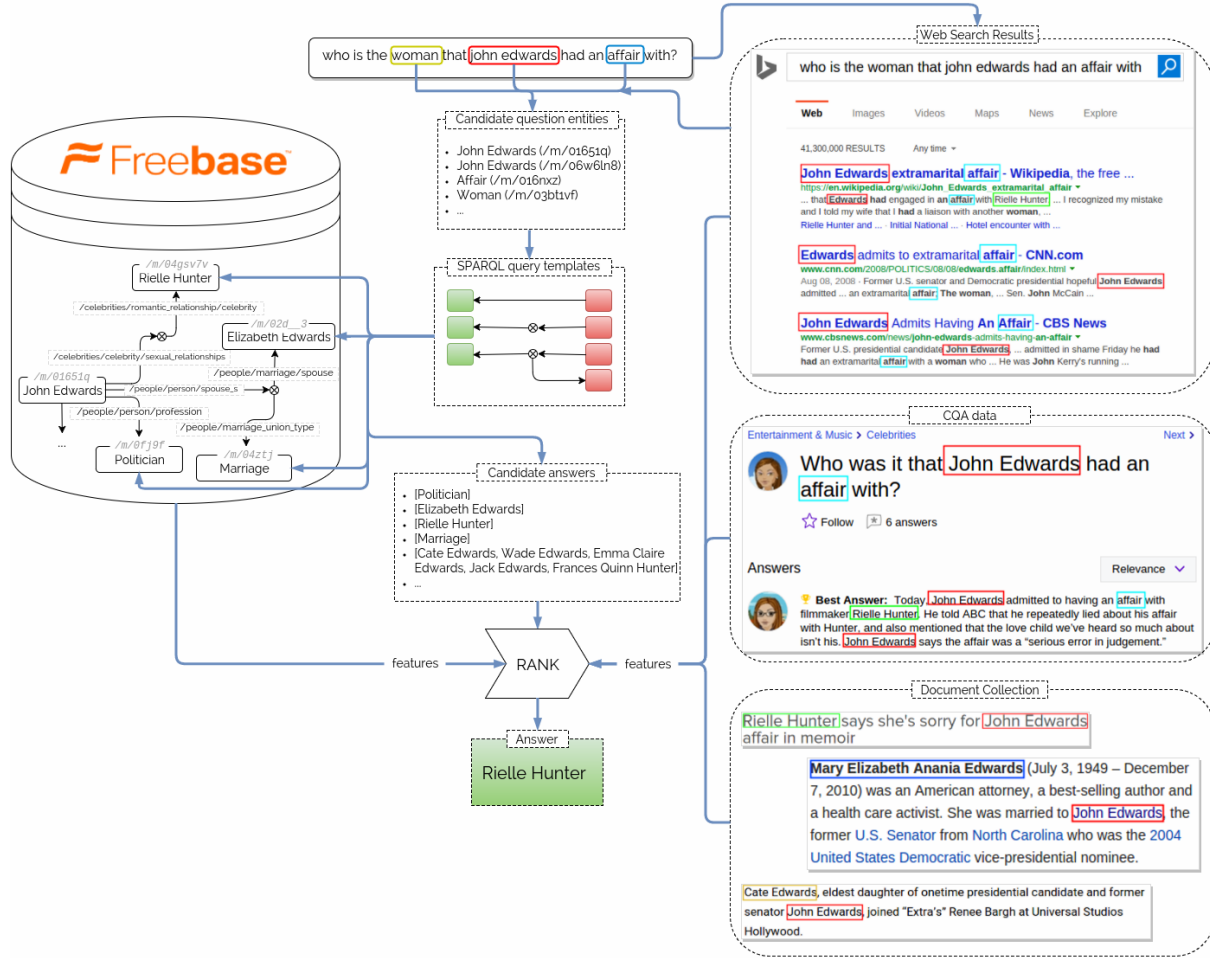


Figure 3: The architecture of our Text2KB Question Answering system

```
SELECT DISTINCT ?a_entity {
  <q_entity_1> <predicate_1> ?cvt_node .
  ?cvt_node <from_predicate> ?date_from .
  ?cvt_node <to_predicate> ?date_to .
  ?cvt_node <predicate_2> ?a_entity .
  FILTER ( <question_date> >= ?date_from AND
    <question_date> <= ?date_to )
}
```

3. TEXT2KB: INCORPORATING TEXT DATA INTO KBQA

We now introduce our system, called Text2KB, that expands upon the basic KBQA model by incorporating external textual sources throughout the QA process. The general architecture and an example use case of Text2KB is presented on Figure 3. The left part of the figure roughly corresponds to the architecture of the existing information extraction approaches to KBQA, described above. The right part introduces additional external text data sources, specifically we investigate the use of web search results, community question answering (CQA) data, and a large collection of documents with detected KB entity mentions. Recall that the main challenges in KBQA are linking topical entities in the question to the KB; identifying candidate answers in the neighborhood around the question entities; and ranking the candidates. In the rest of this section we present our approach to solving each of these challenges, by: using web

search results, CQA data, and external corpus statistics.

3.1 Web search results for KBQA

To obtain related web search results, Text2KB issues the question as a query to a commercial web search engine³, extracts top 10 search result snippets and the corresponding documents. Next, it detects KB entity mentions in both snippets and documents. This data turns out to be useful for multiple purposes, *i.e.*, question entity identification and answer candidate ranking.

Question entity identification. Question text provides only a limited context for entity disambiguation and linking; additionally, the entity name can be misspelled or an uncommon variation used. This complicates the task of entity identification, which is the foundation of whole question answering process. Fortunately, web search results help with these problems, as they usually contain multiple mentions of the same entities and provide more context for disambiguation. Text2KB uses the search result snippets to *expand* the set of detected question entities. To keep only the entities that are also mentioned in the question and avoid irrelevant entities, we use string distance. More specifically, we take names of all entities detected in the question and compute their term by term similarity with non-stopwords from the question. In this work we used Jaro-Winkler string distance

³<https://datamarket.azure.com/dataset/bing/search>

and entity was added to the list of question entities if at least one of its tokens e_t has high similarity with one of the question tokens q_t excluding stopwords (*Stop*):

$$\max_{e_t \in M \setminus Stop, q_t \in Q \setminus Stop} distance_{Jaro-Winkler}(e_t, q_t) \geq 0.8$$

Answer candidate features. Most of the information stored in knowledge bases is also present in other formats, including natural language statements, tables, *etc.* For example, on Figure 2 multiple snippets mention the date when Tutankhamun became the king. Text-QA systems usually generate answer candidates from passages extracted from retrieved documents. In our case candidates are already generated from a KB and we just need to rank them to select the best one. Text2KB uses snippets and documents to compute a set of features, which are used for answer candidate ranking. More specifically we do this following:

1. Precompute term and entity IDFs⁴. We used Google n-grams corpus to approximate terms IDF by collection frequencies and available ClueWeb Freebase entity annotations⁵ to compute entity IDFs
2. Each snippet and document is represented by two TF-IDF vectors of lowercased tokens and mentioned entities
3. In addition, vectors of all snippets and all documents are merged together to form combined token and entity vectors
4. Each answer candidate is also represented as TF-IDF vectors of terms (from entity names) and entities
5. We compute cosine similarities between answer and each snippet and document token and entity vectors. This gives us 10 similarity scores for every document for token vectors and 10 similarities for entity vectors, we take average and maximum scores as features.
6. We do the same for the combined document and use cosine similarities as features.

3.2 CQA data for Matching Questions to Predicates

Recall that a major challenge in KBQA is that natural language questions do not easily or uniquely map to entities and predicates in a KB. An established approach for this task is supervised machine learning, which requires labeled examples of questions and the corresponding answer to learn this mapping. Unfortunately, manual labeling of questions with answers is expensive, and necessarily contains only a small fraction of the different ways the same KB predicate can be inquired about using natural language questions. Researchers have proposed to use weakly supervised methods to extend the lexicon with mappings learned from *single sentence statements* mentioning entity pairs from a large corpus [29]. However, often there is a lexical gap between how information is asked about in a question and how it is expressed in a statement. On the other hand there are huge archives of questions and answers posted by real users on various community question answering websites, *e.g.*, Figure 4.

For our experiments we use 4,483,032 questions from Yahoo! Comprehensive Questions and Answers WebScope dataset⁶. Texts of each question and answer pair were run through an entity linker, that detected mentions of Freebase entities.

⁴<https://en.wikipedia.org/wiki/Tf-idf>

⁵<http://lemurproject.org/clueweb09/FACC1/>

⁶<https://webscope.sandbox.yahoo.com/catalog.php?datatype=1>

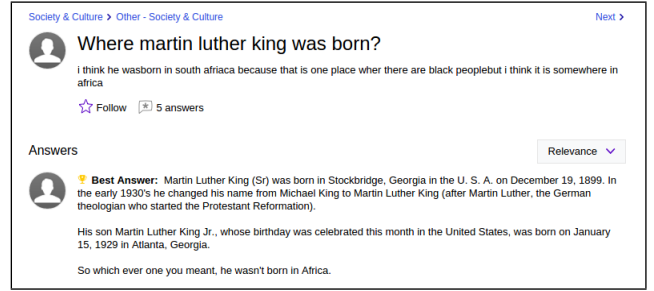


Figure 4: Example of a question and answer pair from Yahoo! Answers CQA website

Table 1: Examples of term-predicate pairs with high PMI scores, computed using distant supervision from a CQA collection

Term	Predicate	PMI score
born	people.person.date_of_birth	3.67
	people.person.date_of_death	2.73
	location.location.people_born_here	1.60
kill	people.deceased.person.cause_of_death	1.70
	book.book.characters	1.55
currency	location.country.currency_formerly_used	5.55
	location.country.currency_used	3.54
school	education.school.school_district	4.14
	people.education.institution	1.70
	sports.school_sports_team.school	1.69
illness	medicine.symptom.symptom_of	2.11
	medicine.decease.causes	1.68
	medicine.disease.treatments	1.59
win	sports.sports_team.championships	4.11
	sports.sports_league.championship	3.79

Next, similar to an idea of relation extraction from CQA data [20], we use distant supervision to label each question-answer pair with predicates between entities mentioned in the question and in the answer. As a result, we have a set of questions, annotated with KB predicates, which are, often incorrectly, assumed to answer the question. We learn the associations between question terms and predicates by computing pointwise mutual information scores⁷ (PMI) for each term-predicate pair. Examples of scores for some terms from WebQuestions dataset questions are given in Table 1.

Although noisy, the statistics look reasonable to be used for candidate ranking. In Text2KB we take candidate answer predicates and look up the PMI scores between these predicates and terms in the question. Missing pairs are given a score of 0, and minimum, average and maximum of these scores are used as features. Since this kind of data is usually sparse, we also use pretrained word2vec word embeddings⁸ to generate predicate embeddings by taking weighted average of term vectors from predicate’s PMI table. Each term’s embedding vector is weighted by its PMI value (terms with negative score are skipped). Then, we compute cosine similarities between predicate vector and question term vectors and take their minimum, average, maximum as features. Similarity between the predicate vector and average question term vector is also computed.

⁷https://en.wikipedia.org/wiki/Pointwise_mutual_information

⁸<https://code.google.com/p/word2vec/>

Table 2: Example of entity pairs along with the most popular terms mentioned around the entities

Entity 1	Entity 2	Term counts
John Edwards	Rielle Hunter	campaign, affair, mistress, child, former ...
John Edwards	Cate Edwards	daughter, former, senator, courthouse, left, greensboro, eldest ...
John Edwards	Elizabeth Edwards	wife, hunter, campaign, affair, cancer, rielle, husband ...
John Edwards	Frances Quinn Hunter	daughter, john, rielle, father, child, former, paternity...

3.3 Estimating Entity Associations

When ranking candidate answers, we are interested in estimating whether the entities in the question and the answer are related in a way asked in the question. Existing systems usually look on how candidate predicates are expressed in questions and statements. But predicate isn't the only way we can look at this. An alternative is to consider text passages, *e.g.*, sentences, that mention topical and answer entities together. For example, in the bottom right corner of Figure 3 we can see some passages that mentioned a pair of people, and the context of these mentions often expresses the nature of the relationships between the entities.

We use the ClueWeb12 corpus with existing Freebase entity annotations⁹ and compute counts of different terms that occur in the context to an entity pair mention. By an entity pair mention we mean a pair of mentions of different entities within 200 characters of each other. We take terms in between mentions and 100 character before and after mentions as the context. A small sample of this data is presented in Table 2.

First, given a set of question terms Q and an answer candidate, that starts from a question entity e_1 , we compute a language model score for every answer entity e_2 :

$$p(Q|e_1, e_2) = \prod_{t \in Q} p(t|e_1, e_2)$$

and use minimum, average and maximum as features. To address the sparsity problem, we again use embeddings, ie for each entity pair a weighted (by counts) average embedding vector of terms is computed and minimum, average and maximum cosine similarities between these vectors and question tokens vector are used as features.

3.4 Internal text data to enrich entity representation

In addition to external text data, many knowledge bases, including Freebase, contain text data as well. In Freebase, most of the entities contain a description paragraph, which often comes from the entity Wikipedia profile. These descriptions of entities in the KB were found useful for text-based question answering [22]. For completeness, we include them in our system. The aim is to improve the matching of the question text, to the unstructured description of the candidate entities. For this, each entity description is represented as a vector of tokens, and a vector of mentioned

entities. We compute cosine similarities between the question tokens and each of the candidate entity vectors, and use these scores as features for candidate ranking. In future work, we could explore incorporating any other entity profile text, such as full Wikipedia article.

4. EVALUATION

We followed the standard evaluation procedure for the WebQuestions dataset and used the original 70-30% train-test split, which results in 3,778 training and 2,032 test questions. Since each answer is potentially a list of entities a^* , the quality of an answer a is represented by F1-score:

$$f1(a^*, a) = 2 \frac{precision(a^*, a)recall(a^*, a)}{precision(a^*, a) + recall(a^*, a)}$$

where $precision(a^*, a) = \frac{|a^* \cap a|}{|a|}$ and $recall(a^*, a) = \frac{|a^* \cap a|}{|a^*|}$.

We also report average precision and recall, as well as an F1 score of average precision and recall. The results of existing approaches, our baseline and Text2KB systems is presented in Table 3.

As we can see, Text2KB significantly improves over the baseline system and reaches the current best published result - STAGG [30], and we believe that this system will also benefit from the ideas of our work, and we will explore this question in Section 5.

4.1 Ablation Study

To study effects of different components in isolation we made a series of ablation studies. For convenience, we introduce the following notations for different components of our system:

- T - notable type score model as a ranking feature
- DF - date range filter-based query template
- E - using web search result snippets for question entity identification
- W - using web search results for feature generation
- CQA - using CQA-based [question term, KB predicate] PMI scores for feature generation
- CW - features, computed from entity pairs language model, estimated on ClueWeb

In our results table we will use the notation $+<component>$ to for a system with a certain component added, and $-<component>$ when the component is removed. For example, the baseline system will be denoted as "Aqqu" according to the authors notation. The same system with additional date range filter query templates and notable types score model is denoted as "Aqqu +DF+T", which represents the same system as "Text2KB -E-W-CQA-CL". Our full system "Text2KB" can be also denoted as "Aqqu +DF+T+E+W+CQA+CL".

The first question that we are asking is what are the improvements, introduced by adding date range filter templates, notable type model, entity linking from web search results and text-based features generated from all the different sources. Results of this ablation experiment are presented in Table 4. As we can see, additional date range filters and notable types model (Text2KB -E-W-CQA-CL) are responsible for an increased recall and a drop in precision compared to the baseline model. Detecting question entities (Text2KB -W-CQA-CL) help improve both precision and recall, and therefore average F1 score by 0.096 points. An even bigger improvement is achieved by introducing all our external text-based features, and since these improvements

⁹<http://lemurproject.org/clueweb12/FACC1/>

Table 3: Performance of the Text2KB system on WebQuestions dataset compared to the existing approaches. The difference from the baseline Aqqu system is significant with p-value < 0.01

System	avg Recall	avg Precision	F1 of avg Prec and Recall	avg F1
OpenQA [16]	-	-	-	0.35
YodaQA [4]	-	-	-	0.343 ^a
SemPre [6]	0.413	0.480	0.444	0.357
Subgraph Embeddings [10]	-	-	0.432	0.392
ParaSemPre [7]	0.466	0.405	0.433	0.399
Jacana [29]	0.458	0.517	0.486	0.330
Kitt AI [27]	0.545	0.526	0.535	0.443
AgendaIL [8]	0.557	0.505	0.530	0.497
STAGG [30]	0.607	0.528	0.565	0.525
Aqqu (baseline) [3]	0.604	0.498	0.546	0.494
Our system: Text2KB	0.6354	0.5059	0.5633	0.5223

^aThe reported score is F1@1. Text-based QA systems usually return a ranked list of candidate answers and are evaluated by Precision@1 or ranking metrics. However, WebQuestions dataset contains many list answers.

Table 4: Average Recall (R), Precision (Pr), and F1 of Aqqu (baseline), Text2KB (our system), and variations of TextKB with respective components removed. * indicates significant differences at p<0.05.

System	R	Pr	F1
Aqqu (baseline)	0.604	0.498	0.494
Text2KB -E-W-CQA-CL= =Aqqu +DF+T	0.6169	0.4807	0.4987
Text2KB -W-CQA-CL	0.6272*	0.4920*	0.5083*
Text2KB -E	0.6344*	0.4966*	0.5140*

Table 5: Average Recall (R), Precision (Pr), and F1 of Text2KB variations with and without features based on web search results, CQA data and ClueWeb collection

System	R	Pr	F1
Text2KB -W	0.6327	0.4960	0.5126
Text2KB -CQA	0.6420	0.4987	0.5185
Text2KB -CL	0.6444	0.5047	0.5228
Text2KB (Web search results only)	0.6423	0.5028	0.5216
Text2KB (ClueWeb only)	0.6307	0.4978	0.5138
Text2KB (CQA only)	0.6224	0.4928	0.5077

are independent, their combination boosts the performance even more.

Now, let’s look into the relative importance of each of the data sources, we will remove or use a group of web search, cqa or clueweb-based features and see how the performance of the whole system changes. Table 5 summarizes the results of these experiments.

Features that we generate from web search results are the most effective, because even without other data sources the QA performance is almost as high as the full system. In addition, if we remove web search results based features the performance drops more than for other text data sources. Features based on ClueWeb entity pair statistics perform better than CQA-based features.

Table 6: Average Recall (R), Precision (Pr), and F1 for Text2KB (our system), STAGG and their combinations

System	R	P	F1
Our system: Text2KB	0.6354	0.5059	0.5223
STAGG [30]	0.607	0.528	0.525
Text2KB + STAGG	0.5976	0.5343	0.5320
Text2KB + STAGG (oracle)	0.7144	0.5904	0.6056

Since we used each data source to generate multiple different features for candidate ranking, it is interesting to see which particular features are more useful than others by the ranking machine learning algorithm (we used random forest). Figure 5 plots a subset of features ranked by their Gini index-based importance scores in the final answer candidate ranking model.

The figure supports the observation that web search results features are the most useful, however, other text data sources also contribute to the improvement.

In summary, Text2KB significantly outperforms the baseline system, and each of the introduced components contributes to this improvement. Web search results data turned out to be the most useful resource, and it significantly improves the quality by helping with question entity identification and candidate ranking. Next, we analyze the system performance in more detail, and investigate factors for future extension.

5. ANALYSIS AND DISCUSSION

We have shown that Text2KB outperforms the baseline. We now investigate how our system would compare to other systems on the same benchmark; then, we investigate in depth the different error modes (Section 5.1), which helps identify the areas of most substantial future improvements.

We took an existing KBQA systems and demonstrated that by combining evidence from knowledge base and external text resources we can boost the performance. A reasonable question is whether the same approach will be helpful to other systems, *e.g.*, the currently best system STAGG

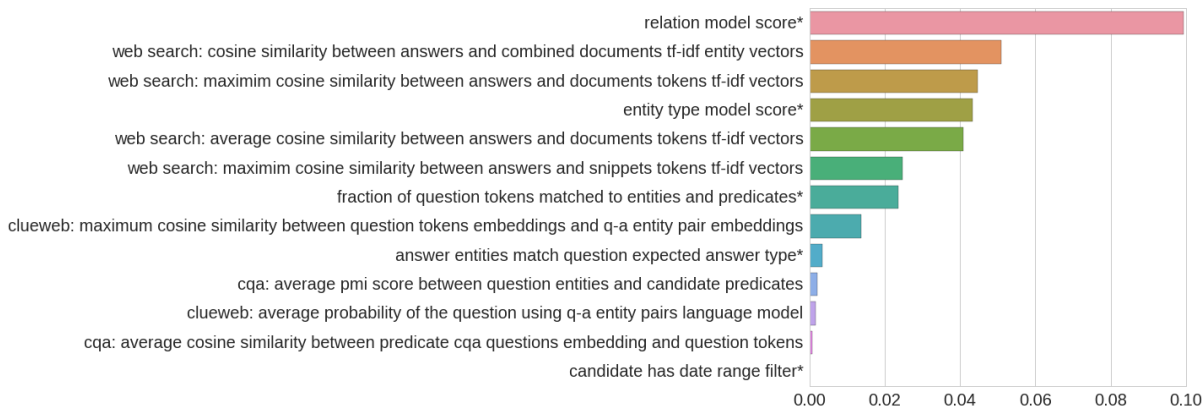


Figure 5: Importances of different text-based features for KBQA (features with * are not text-based and are provided for comparison)

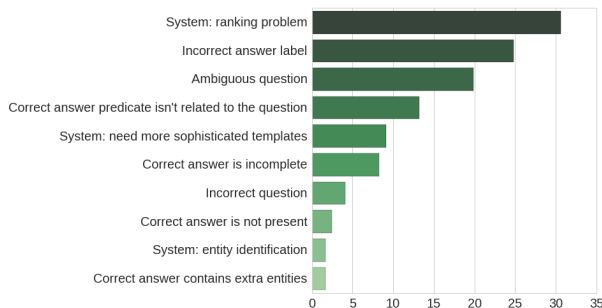


Figure 6: Distribution of problems with questions, where Text2KB returns an answer with F1<1

[30]. The differences between our baseline system Aquu and STAGG lie in the components, *i.e.*, entity linking algorithm, a set of query templates and ranking methods, therefore our approach is complementary and should be helpful. To support this claim, we made an experiment to combine answers of STAGG and Text2KB. One of the advantages of the former is its set of filters, that restricts list results to entities of certain type, gender, *etc.* Therefore, we combined answers of STAGG and Text2KB using a simple heuristic: we chose to use the answer returned by STAGG if the number of answer entities is less than in the Text2KB answer, otherwise we use the answer of our approach. Table 6 gives the results of the experiment, and as we can see the combination achieves slightly better average F1 score. Alternatively, we can look at the oracle combination of the systems, which always chooses an answer with higher F1. This experiment shows that that systems don't make exactly the same mistakes and therefore can be combined. As we can see such a combination results in a performance of 0.6056, which is much higher than either of the systems.

WebQuestions dataset is rather small as a result, answers to 112 of the test questions involve a predicate that weren't observed in the training set, which may be a problem for approaches that rely on a trained lexicon. We evaluated both systems on these questions, and indeed the performance is very low, *i.e.*, the average F1 score of Text2KB is 0.1640 compared to 0.1199 for STAGG¹⁰.

5.1 Error analysis

To get a better insights of the problems that remain, we collected 1219 questions for which Text2KB didn't return completely correct answer, *i.e.*, F1 score of the answer is less than 1. We manually looked through a couple of hundreds of these examples and grouped the problems into several clusters. The results are summarized on Figure 6.

As we can see candidate ranking is still the major problem, and it accounts for 31% of the cases. The second most popular problem is incorrect ground truth labels (almost a quarter of errors). For example: for the question *when tupac was shot?* the label says **Tupac 1994 assault** instead of **Las Vegas**. Another set of questions have incomplete or overcomplete ground truth answer list. Typical examples are questions asking for a list of movies, books, landmarks, *etc.* The ground truth answer usually contains ~ 10 entities, whereas the full list is often much larger. This seems to be an artifact of the labeling process, where the answer was selected from the Freebase entity profile page. The profile page shows only a sample of 10 entities from large lists and the others are hidden behind the "NNN values total" link. About 20% of the questions are ambiguous, *i.e.*, questions have no strict 1-1 correspondence with any of the predicates and can be answered by multiple without any obvious preferences. For example, the question *"what did hayes do?"* can be answered by profession, occupied position or some other achievements. Another problem is when there is no predicate that answers the question. For example, the question *"what do people in france like to do for fun?"* doesn't have a good match among the facts stored in Freebase. The ground truth entity **Cycling** comes from predicate related to the olympic sport competitions country participated in¹¹, which obviously isn't related to the question.

As for the system errors, there are wins and losses introduced by each of our components. Web search results helped identify the right question topical entity in a number of cases, *e.g.*, *"what did romo do?"* mentions only the last name of the Dallas Cowboys quarterback and the baseline system were unable to map it to the right entity. Web search results provides more than enough evidence that romo refers to **Tomo Romo**. However, there are a number of losses, introduced by added unrelated entities. For example, the entity **I Love Lucy** was added for the question *"what was lucille ball?"*, because the term *lucy* had high similarity with *lucille*.

¹⁰Unfortunately, the number of questions is too low to show statistical significance (p-value=0.16)

¹¹olympics.olympic_participating_country.athletes

A portion of these problems can be fixed by a better entity linking strategy, *e.g.*, [12].

An interesting example, when external text resources improved the performance is the question “*what ship did darwin sail around the world?*”. This is actually a hard question, because the ship entity is connected to the **Charles Darwin** entity through the “knownFor” predicate along with some other entities like **Natural selection**. Thus, the predicate itself isn’t related to the question, but nevertheless, the name of the ship **HMS Beagle** is mentioned multiple times in the web search results, and entity pair model computed from ClueWeb also has high scores for the terms “ship” and “world”.

There are several major reasons for the losses, introduced by features based on external text resources. Some entities often mentioned together and therefore one of them gets high values of cooccurrence features. For example, the baseline system answered the question “*when did tony romo get drafted?*” correctly, but since **Tony Romo** is often followed by **Dallas Cowboys**, Text2KB ranked the team name higher. Another common problem with our features is an artifact of entity linking, which works better for names and often skips abstract entities, like professions. For example, the correct answer to the question “*what did jesse owens won?*” is an entity with the name **Associated Press Male Athlete of the Year**, which is rarely mentioned or it’s hard to find such mentions. Some problems were introduced by a combination of components. For example, for “*where buddha come from?*” a topical entity **Buddhism** was introduced from search results, and it generated **Gautama Buddha** as one of the answer candidates. This answer was ranked the highest due to large number of mentions in the search results.

In summary, we show that ideas behind Text2KB could be integrated into other systems and improve their performance. The error analysis suggested that even though a significant number of questions in the WebQuestions dataset have incorrect or ambiguous ground truth labels, there is still a room for improvement. In particular, the future work for Text2KB will include a better strategy for entity linking using external data sources and a better context model for entity mentions in text documents, which can put more weight on entities mentioned in the context related to the question.

6. RELATED WORK

In 2011 a series of QALD (Question Answering over Linked Data) evaluation campaigns has started. You can find the most recent report in [23]. These benchmarks use dbPedia knowledge base and usually provide a training set of questions, annotated with the ground truth SPARQL queries. In QALD-3 a multilingual task has been introduced, and since QALD-4 the hybrid task is included. This task asks participants to use both structured data and free form text available in dbPedia abstracts. The formulation of the hybrid task is the most relevant to our work, but there are a couple of key differences. Questions in the hybrid track are manually created in such a way, that they can *only* be answered using a combination of RDF and free text data. Secondly, the hybrid task focuses on text data already present in a KB, whereas we are exploring external text resources. In general, because of the expensive labeling process, QALD datasets are rather small, for example, QALD-5 training set for multilingual question answering includes 300 examples and 40

examples for the hybrid task. The evaluation was performed on 50 questions for multilingual task and just 10 for hybrid. Therefore, due to the scale of datasets and slightly different focus of tasks, we didn’t attempt to evaluate our techniques on QALD benchmarks, but intend to explore it further in the future.

WebQuestions benchmark was introduced in [6]. The proposed approaches differ in the algorithms used for various components, and, what is more relevant to our work, the use of external datasets. To account for different ways a question can be formulated [7] used a dataset of question clusters from WikiAnswers to learn a question paraphrasing model. Another approach to learn term-predicate mapping is to use distant supervision [18] to label a large text corpus, such as ClueWeb [29]. In this work we build on this idea and instead of focusing on term-predicate mappings, which might be too general, consider particular entity pairs. Free-base RDF triples can automatically converted to questions using entity and predicate names [10]. Finally, many systems work with distributed vector representations for words and RDF triples and use various deep learning techniques for answer selection [10, 30]. In all of these works, external resources are used to train a lexicon for matching questions to particular KB queries. The use of external resources in this work is different, we are targeting better candidate generation and ranking by considering the actual answer entities rather than predicates used to extract them.

In general, combining different data sources, such as text documents and knowledge bases, for question answering is not a novel idea, and it has been already implemented in hybrid QA systems [5, 2]. Such systems typically have different pipelines that generate answer candidates from each of the data sources independently, and merge them to select the final answer at the end. We make a step towards integration of approaches, by incorporating text resources into different stages of knowledge base question answering process. This is similar to the work of [22], who explored the use of entity types and descriptions from a KB for text-based question answering, and [13] explored such semantic annotations for ad-hoc document retrieval. M. Yahya et al [26] explored an interesting idea of extending SPARQL triple patterns with text keywords, and proposed to use certain query relaxation techniques to improve the robustness of KBQA systems. Query relaxation constitutes in dropping certain triples patterns from SPARQL query and adding the corresponding question words as keywords to other triple patterns. These ideas are complimentary to our work, because our use of text data improves the matching between question phrases and KB concepts, whereas query relaxation replaces this problem with keyword matching and therefore can be applied on top of each other. Another KB-Text hybrid approach, proposed in [25], utilizes text resources as a post processing step for answer validation and filtering.

We should also mention OpenIE [15], which represent an interesting mixture between text and structured data. Such knowledge repositories can be queried using structured query languages, and at the same time allows keyword matching against entities and predicates [16]. In this work, we are borrowing an idea of learning about entity relationship via natural language phrases connecting them. However, since we don’t need to extract clean set of relation tuples, we can keep all kinds of phrases, mentioned around entity pairs.

7. CONCLUSIONS AND FUTURE WORK

Our work showed that unstructured text resources can be effectively utilized for knowledge base question answering to improve query understanding, candidate answer generation and ranking. We focused on three particular techniques and associated text information sources: web search results for query understanding and candidate ranking, community question answering data for candidate generation, and text fragments around entity pair mentions for ranking. Certainly, there are more resources that could be potentially adapted, *e.g.*, entity profile pages like Wikipedia, news sources, textbooks, and many others. However, we believe that the proposed approach is general enough that it could be extended and successfully incorporate these other diverse text sources.

In the future, we plan to extend our work to the more open setup, similar to the QALD hybrid task, where questions no longer have to be answered exclusively from the KB. This would require extending the described techniques, and creating new QA benchmarks.

8. REFERENCES

- [1] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, and Z. Ives. *Dbpedia: A nucleus for a web of open data*. Springer, 2007.
- [2] K. Barker. Combining structured and unstructured knowledge sources for question answering in watson. In O. Bodenreider and B. Rance, editors, *DILS*, volume 7348 of *Lecture Notes in Computer Science*, pages 53–55. Springer, 2012.
- [3] H. Bast and E. Haussmann. More accurate question answering on freebase. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, CIKM '15, pages 1431–1440, New York, NY, USA, 2015. ACM.
- [4] P. Baudiš. Systems and approaches for question answering. 2015.
- [5] P. Baudiš and J. Šedivý. Modeling of the question answering task in the yodaqa system. In *Experimental IR Meets Multilinguality, Multimodality, and Interaction*, pages 222–228. Springer, 2015.
- [6] J. Berant, A. Chou, R. Frostig, and P. Liang. Semantic parsing on freebase from question-answer pairs. In *Proceedings of EMNLP*, 2013.
- [7] J. Berant and P. Liang. Semantic parsing via paraphrasing. In *Proceedings of ACL*, volume 7, page 92, 2014.
- [8] J. Berant and P. Liang. Imitation learning of agenda-based semantic parsers. *Transactions of the Association for Computational Linguistics*, 3:545–558, 2015.
- [9] 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.
- [10] 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.
- [11] 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.
- [12] 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 and Disambiguation*, ERD '14, pages 25–30, New York, NY, USA, 2014. ACM.
- [13] J. Dalton. *Entity-based Enrichment for Information Extraction and Retrieval*. PhD thesis, University of Massachusetts Amherst, 2014.
- [14] 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.
- [15] A. Fader, S. Soderland, and O. Etzioni. Identifying relations for open information extraction. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 1535–1545. Association for Computational Linguistics, 2011.
- [16] A. Fader, L. Zettlemoyer, and O. Etzioni. Open question answering over curated and extracted knowledge bases. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '14, pages 1156–1165, New York, NY, USA, 2014. ACM.
- [17] J. Lin. An exploration of the principles underlying redundancy-based factoid question answering. *ACM Trans. Inf. Syst.*, 25(2), Apr. 2007.
- [18] M. Mintz, S. Bills, R. Snow, and D. Jurafsky. Distant supervision for relation extraction without labeled data. In *Proceedings of ACL 2009*, pages 1003–1011. Association for Computational Linguistics, 2009.
- [19] J. Pound, P. Mika, and H. Zaragoza. Ad-hoc object retrieval in the web of data. In *Proceedings of the 19th International Conference on World Wide Web*, WWW '10, pages 771–780, New York, NY, USA, 2010. ACM.
- [20] D. Savenkov, W.-L. Lu, J. Dalton, and E. Agichtein. Relation extraction from community generated question-answer pairs. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 96–102, Denver, Colorado, June 2015. Association for Computational Linguistics.
- [21] V. I. Spitzkovsky and A. X. Chang. A cross-lingual dictionary for english wikipedia concepts. In N. C. C. Chair), K. Choukri, T. Declerck, M. U. Doğan, B. Maegaard, J. Mariani, A. Moreno, J. Odijk, and S. Piperidis, editors, *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, Istanbul, Turkey, may 2012. European Language Resources Association (ELRA).
- [22] 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.

- [23] C. Unger, C. Forascu, V. Lopez, A.-C. N. Ngomo, E. Cabrio, P. Cimiano, and S. Walter. Question answering over linked data (qald-5). In L. Cappellato, N. Ferro, G. J. F. Jones, and E. SanJuan, editors, *CLEF (Working Notes)*, volume 1391 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2015.
- [24] D. Vrandečić and M. Krötzsch. Wikidata: A free collaborative knowledgebase. *Commun. ACM*, 57(10):78–85, Sept. 2014.
- [25] K. Xu, Y. Feng, S. Reddy, S. Huang, and D. Zhao. Enhancing freebase question answering using textual evidence. *arXiv preprint arXiv:1603.00957*, 2016.
- [26] 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.
- [27] X. Yao. Lean question answering over freebase from scratch. In *Proceedings of NAACL Demo*, 2015.
- [28] X. Yao, J. Berant, and B. Van Durme. Freebase qa: Information extraction or semantic parsing? *ACL 2014*, page 82, 2014.
- [29] X. Yao and B. Van Durme. Information extraction over structured data: Question answering with freebase. In *Proceedings of ACL*, 2014.
- [30] W.-t. Yih, M.-W. Chang, X. He, and J. Gao. Semantic parsing via staged query graph generation: Question answering with knowledge base. In *Association for Computational Linguistics (ACL)*, 2015.