

Integrating Web-based and Corpus-based Techniques for Question Answering

Boris Katz, Jimmy Lin, Daniel Loreto, Wesley Hildebrandt,
Matthew Bilotti, Sue Felshin, Aaron Fernandes, Gregory Marton, Federico Mora
MIT Computer Science and Artificial Intelligence Laboratory
Cambridge, MA 02139

1 Introduction

MIT CSAIL's entry in this year's TREC Question Answering track focused on integrating Web-based techniques with more traditional strategies based on document retrieval and named-entity detection. We believe that achieving high performance in the question answering task requires a combination of multiple strategies designed to capitalize on different characteristics of various resources.

The system we deployed for the TREC evaluation last year relied exclusively on the World Wide Web to answer factoid questions (Lin et al., 2002). The advantages that the Web offers are well known and have been exploited by previous systems (Brill et al., 2001; Clarke et al., 2001; Dumais et al., 2002). The immense amount of freely available unstructured text provides data redundancy, which can be leveraged with simple pattern matching techniques involving the expected answer formulations. In many ways, we can utilize huge quantities of data to overcome many thorny problems in natural language processing such as lexical ambiguity and paraphrases. Furthermore, Web search engines such as Google provide a convenient front-end for accessing and filtering enormous amounts of Web data. We have identified this class of techniques as the *knowledge mining* approach to question answering (Lin and Katz, 2003).

In addition to viewing the Web as a repository of unstructured documents, we can also take advantage of structured and semistructured sources available on the Web using *knowledge annotation* techniques (Katz, 1997; Lin and Katz, 2003). Through empirical analysis of real world natural language questions, we have noticed that large classes of commonly occurring queries can be parameterized and captured using a simple object-property-value data model (Katz et al., 2002). Furthermore, such a data model is easy to impose on Web resources through a framework of wrapper scripts. These techniques allow our system to view the Web as if it were a "virtual database" and use knowledge contained therein to answer user questions.

While the Web is undeniably a useful resource

for question answering, it is not without drawbacks. Useful knowledge on the Web is often drowned out by the sheer amount of irrelevant material, and statistical techniques are often insufficient to separate right answers from wrong ones. Overcoming these obstacles will require addressing many outstanding issues in computational linguistics: anaphora resolution, paraphrase normalization, temporal reference calculation, and lexical disambiguation, just to name a few. Furthermore, the setup of the TREC evaluations necessitates an extra step in the question answering process for systems that extract answers from external sources, typically known as *answer projection*. For every Web-derived answer, a system must find a supporting document from the AQUAINT corpus, even if the corpus was not used in the answer extraction process.

This year's main task included definition and list questions in addition to factoid questions. Although Web-based techniques have proven effective in handling factoid questions, they are less applicable to tackling definition and list questions. The data-driven approach implicitly assumes that each natural language question has a unique answer. Since a single answer instance is sufficient, algorithms were designed to trade recall for precision. For list and definition questions, however, a more balanced approach is required, since multiple answers are not only desired, but necessary. We believe that the best strategy is to integrate Web-based approaches with more traditional question answering techniques driven by document retrieval and named-entity detection. Corpus- and Web-based strategies should play complementary roles in an overall question answering framework.

2 List Questions

For answering list questions, our system employs a traditional pipeline architecture with distinct stages for document retrieval, passage retrieval, answer extraction, and duplicate removal (see Figure 1). The general idea is to successively narrow down the AQUAINT corpus, first to a candidate list of documents, then to manageable-sized passages, and fi-

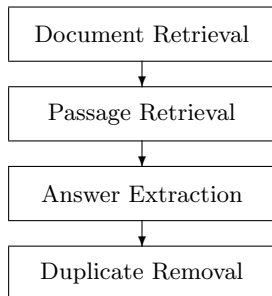


Figure 1: Architecture for answering list questions.

nally employ knowledge of fixed lists to extract relevant answers. The following subsections describe this process in greater detail.

2.1 Document Retrieval

In response to a natural language question, our document retriever provides a set of candidate documents that are likely to contain the answer; these documents serve as the input to additional processing modules. As such, the importance of document retrieval cannot be overstated: if no relevant documents are retrieved, any amount of additional processing would be useless.

For our document retriever, we relied on Lucene, a freely available open-source IR engine.¹ Lucene supports a weighted boolean query language, although the system performs ranked retrieval using a standard *tf.idf* model. We have previously discovered that for the purposes of passage retrieval, Lucene performs on par with state-of-the-art probabilistic systems based on the Okapi weighting (Tellex et al., 2003).

An often effective way to boost document retrieval recall is to employ query expansion techniques. In our TREC entry this year, we implemented two separate query generators that take advantage of linguistic resources to expand query terms. Lucene provides a structured query interface that gives us the ability to fine-tune our query expansion algorithms. In the following subsections, we describe these two techniques in greater detail.

2.1.1 Method 1

Our first query generator improves on a simple bag-of-words query by taking inflectional and derivational morphology into account: queries are a conjunction of disjuncts, where each disjunct contains morphological variants of a single term. Base query terms are extracted from the natural language question by removing all stopwords. Assuming we have three query terms, A , B , and C , arranged in increas-

ing *idf*, our first query method would generate the following queries:

$$\begin{aligned}
 &A \wedge B \wedge C \\
 &e(A) \wedge e(B) \wedge e(C) \\
 &e(B) \wedge e(C) \\
 &e(C) \\
 &e(A) \wedge e(B) \\
 &e(B) \\
 &e(A)
 \end{aligned}$$

where

$$e(x) = x \vee \text{inflect}(x)^{0.75} \vee \text{derive}(x)^{0.50}$$

where $\text{inflect}(x)$ and $\text{derive}(x)$ represent a disjunct of inflectional and derivational morphological forms of x , respectively. The first query is simply a conjunction of all non-stopwords from the question. The second query is a conjunction where each of the conjoined elements is a disjunct of the morphological expansions of a query term. Inflectional variants are generated with the assistance of WordNet (to handle irregular forms). Derivational variants are generated by a version of CELEX that we manually annotated. Using Lucene’s query weighting mechanism, inflected forms are given a weight of 0.75, and derivational forms a weight of 0.5. To generate subsequent queries, the system successively drops disjuncts starting with the disjunct associated with the lowest *idf* term until all disjuncts have been dropped—this has the effective of query relaxation. After that, the highest *idf* disjunct is dropped, and the generator starts a fresh cycle of successively dropping the lowest *idf* disjuncts.

Our document retriever is given a target hit list size, and successively executes queries from the query generator until the target number of documents has been found. This ensures that downstream modules will always be given a consistently-sized set of documents to process.

2.1.2 Method 2

Our second query generation algorithm takes advantage of named-entity recognition technology and other lexical resources to chunk natural language questions so that query terms are not broken across constituent boundaries. To identify relevant named entities, we use Sepia (Marton, 2003), an information extraction system based on Combinatory Categorical Grammar (CCG). In particular, personal names are recognized so that inappropriate queries are never generated; for example, a name such as “John Fitzgerald Kennedy” can produce legitimate queries involving “John F. Kennedy”, “John Kennedy”, and “Kennedy”, but never “John Fitzgerald” or simply “John”. For certain classes of named-entity types, we have encoded a set of heuristic rules that generates the acceptable variants. Our

¹jakarta.apache.org/lucene/docs/index.html

query generator takes advantage of Lucene’s ability to execute phrase queries to ensure that the best matching documents are returned.

Our second query generator also leverages WordNet to identify multi-word expressions that should not be separated in the query process. Multi-token collocations such as “hot dog” should never be broken down into `hot` and `dog`, since the meaning of hot dog cannot be compositionally derived from the individual words. Because these multi-word expressions cannot be predicted syntactically (e.g., compare “hot dog” with “fast car”), one practical solution is to employ a fixed list of such lexical items. If a query term is neither a recognized entity nor a multi-word expression, our second query generator expands the term with inflectional and derivational variants using the same technique as the first method.

We found that our first query generation method traded off precision for recall with its elaborate term dropping strategy—often, the first few queries are too restrictive, and because of this, most of the documents are retrieved by overly general queries. The result is often a hit list that has been “padded” with irrelevant documents; it appears that loose queries with few terms aren’t precise enough to retrieve good candidate documents. As an alternative, we implemented a slightly different strategy for our second query generator. It drops query disjuncts in order of increasing *idf* until no terms remain, and then stops. As a simple example, if the query has three (non-stopword) terms, *A*, *B*, and *C*, arranged in increasing *idf*, our second query generator would produce the following queries:

$$\begin{aligned} e(A) \wedge e(B) \wedge e(C) \\ e(B) \wedge e(C) \\ e(C) \end{aligned}$$

where $e(x)$ represents the expansions of an individual query term, as described in this section.

2.2 Passage Retrieval

The next stage in the processing pipeline for answering list questions is passage retrieval, which attempts to narrow down the set of candidate documents to a set of candidate passages, which are sentences in our architecture.

In a separate study of passage retrieval algorithms (Tellex et al., 2003), we determined that IBM’s passage scoring method (Ittycheriah et al., 2000; Ittycheriah et al., 2001) produced the most accurate results. To determine the best passage (sentence in our case), our system breaks each candidate document into sentences and scores each one based on the IBM algorithm.

The IBM passage retrieval algorithm computes a series of distance measures for each passage. The

“matching words measure” sums the *idf* values of words that appear in both the query and the passage. The “thesaurus match measure” sums the *idf* values of words in the query whose WordNet synonyms appear in the passage. The “mis-match words measure” sums the *idf* values of words that appear in the query and not in the passage. The “dispersion measure” counts the number of words in the passage between matching query terms, and the “cluster words measure” counts the number of words that occur adjacently in both the question and the passage. These various measures are linearly combined to give the final score for a passage.

We modified the IBM passage scoring algorithm to take into account linguistic knowledge provided by our query generator. The modified algorithm includes scores for matching hyponyms, inflectional variants, derivational variants, and antonyms (negative weight). In addition, our modified algorithm takes advantage of multi-word expressions tokenized from the question, that is, occurrences of “hot” and “dog” within a passage will not match “hot dog”.

One of our goals is to determine the effects of additional linguistic knowledge on performance, and for our TREC submissions, we set up a matrix experiment with two query generators and two passage retrievers (the original IBM method and our modified algorithm). The results will be discussed later in Section 5.

2.3 Answer Extraction

The first step of the answer extraction process is to determine the question focus—the word or phrase in the question that is used to identify the ontological type of the entity we are looking for (i.e., the target type). For this, we enlisted the parser of the START question answering system (Katz, 1997). In addition, we have also constructed a mapping from question focus to target type. Consider a question such as “List journalists that have won the Pulitzer Prize more than once?”: START would recognize *journalist* as the question focus, and PERSON as the target type (since we don’t have a specific category for journalists in our ontology).

Separately, we have compiled offline a large knowledge base of entities, mostly in the form of fixed lists, that correspond to the various target types. For example, we have gathered lists of U.S. states, major U.S. cities, major world cities, countries, person names, etc. If the target type is among one of these categories for which we have a fixed list, our answer extractor simply extracts instances of the target type from the top ranking passages collected by the previous stage.

As an example, consider the following question:

In which U.S. states have there been fatalities caused by snow avalanches? (q2183)

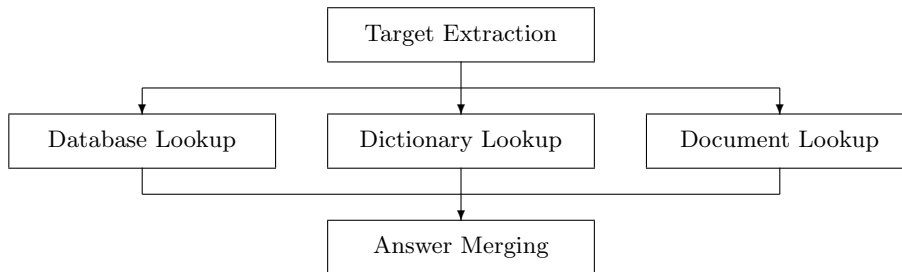


Figure 2: Architecture for answering definition questions.

Our system correctly identifies the question focus as “U.S. state” (corresponding to the target type US STATE) and extracts all instances of U.S. states from top ranking passages. Since the passage retrieval algorithm returns passages that already have occurrences of terms from the question, instances of the target type are likely to be the correct answer.

If the target type is not in our knowledge base, we employ two backoff procedures. Occasionally, answers to list questions have the question focus directly embedded in them (e.g., “littleneck clam” is a type of clam), and in the absence of any additional knowledge, noun phrases containing the question focus are extracted as answer instances. Finally, if no noun phrases containing the question focus can be found, our answer extraction module simply picks the noun phrase closest to the question focus in each of the passages and returns that as the answer.

After collecting all the answer candidates, we discard ones with query terms in them. Noun phrases containing keywords from the query typically repeat some aspect of the original user question and make little sense as answers. This heuristic has worked well in our previous question answering system (Lin and Katz, 2003).

2.4 Duplicate Removal

Answer instances extracted from the previous stage typically contain duplicates, which our system removes using a thresholded edit-distance measure. Finally, the system computes the number of answer instances to return based on a relative thresholding scheme. Each answer candidate is given a score equal to the score of the passage from which it was extracted, and all candidate answers below 10% of the maximum score are discarded. The remaining instances are returned as the final answers.

3 Definition Questions

Our architecture for answering definition questions is shown in Figure 2. The target extraction module

first analyzes the natural language question to determine the unknown term. Once the target term has been found, three parallel techniques are employed to retrieve relevant nuggets that “define” the term: lookup in a database of relational information created from the AQUAINT corpus, lookup in a Web dictionary followed by answer projection, and lookup directly in the AQUAINT corpus with information retrieval techniques. Answers from the three different sources are merged to produce the final system output. The following subsections briefly describe each of these techniques; please refer to our forthcoming paper (Hildebrandt et al., 2004) for more details.

3.1 Target Extraction

We have developed a pattern-based parser to analyze definition questions and extract the target term using simple regular expressions. If the natural language question does not fit any of our patterns, the parser heuristically extracts the last sequence of capitalized words in the question as the target. Our simple definition target extractor was tested on definition-style questions from the previous TREC evaluations and performed quite well on those training questions.

3.2 Database Lookup

The use of surface patterns for answer extraction has proven to be an effective strategy for question answering. Typically, surface patterns are applied to a candidate set of documents that have been returned by traditional document retrieval systems. While this strategy may be effective for factoid questions, it generally suffers from low recall. In the case of factoid questions, where only one instance of an answer is necessary, recall is not a primary concern. However, definition questions require a system to find as many relevant nuggets as possible, making recall very important.

Instead of using surface patterns post-retrieval, we

copular pattern: A **fractal** *is* a pattern that is irregular, but self-similar at all size scales
 appositive pattern: The **Aga Khan**, Spiritual Leader of the Ismaili Muslims
 occupation pattern: *steel magnate* **Andrew Carnegie**
 verb pattern: **Althea Gibson** *became* the first black tennis player to win a Wimbledon singles title
 parenthesis pattern: **Alice Rivlin** (director of the Office of Management and Budget)

Figure 3: Sample nuggets extracted from the AQUAINT corpus using surface patterns. The target terms are in bold, the nuggets underlined, and the pattern landmarks in italics.

employ an alternative strategy: by applying a set of surface patterns offline, we are able to “precompile” from the AQUAINT corpus knowledge nuggets about every entity mentioned within it. In essence, we have automatically constructed an immense relational knowledge base, which, for each entity, contains all the nuggets distilled from every article within the corpus. Once this database has been constructed, the task of answering definition questions becomes a simple database lookup.

Our surface patterns operate both at the word level and the part-of-speech level. We utilize patterns over part-of-speech tags to perform rudimentary chunking, such as marking the boundaries of noun phrases. Our system uses a total of thirteen patterns, some of which are described below (Figure 3 shows several examples):

- **Copular pattern.** Copular constructions often provide a definition of the target term. However, the pattern is a bit more complex than finding the verb *be* and its inflectional variants; in order to filter out spurious nuggets (e.g., the progressive tense), our system throws out all definitional nuggets that do not begin with a determiner; this ensures that we only get “NP₁ be NP₂” patterns, where either NP₁ or NP₂ can be the nugget.
- **Appositive pattern.** Commas typically provide strong evidence for the presence of an appositive. With the assistance of part-of-speech tags, identifying “NP₁, NP₂” patterns is relatively straightforward. Most often, NP₁ is the target term and NP₂ is the nugget, but occasionally the positions are swapped. Thus, we index both NPs as the target term.
- **Occupation pattern.** Common nouns preceding proper nouns typically provide some relevant information such as occupation, affiliation, etc. In order to boost the precision of this pattern, our system discards all common noun phrases that do not contain an “occupation” such as *actor*, *spokesman*, *leader*, etc. We mined this list of occupations from WordNet and Web resources.

- **Verb pattern.** By statistically analyzing a corpus of biographies of famous people, we were able to compile a list of verbs that are commonly used to describe people and their accomplishments, including *became*, *founded*, *invented*, etc. This list of verbs is employed to extract “NP₁ verb NP₂” patterns, where NP₁ is the target term, and NP₂ is the nugget.
- **Parenthesis pattern.** Parenthetical expressions following noun phrases typically provide some interesting nuggets about the preceding noun phrase; for persons, it often contains birth/death years or occupation/affiliation.

Typically, our patterns identify short nuggets on the order of a few dozen characters. In answering definition questions, we decided to return responses that include additional context. To accomplish this, we simply expand all nuggets around their center point to encompass one hundred characters. We found that this technique enhances the readability of our responses: many nuggets seem odd and out of place without context and surrounding text is often necessary for disambiguation. Furthermore, returning a longer answer means that our responses sometimes contain additional relevant nuggets that are not part of the nugget matched by the original pattern. In definition questions, these “fortuitous” nuggets serve as an additional source of answers.

One drawback to our knowledge base of nuggets is the tremendous amount of redundancy contained within it. Because we compiled all patterns from all entities within the entire AQUAINT corpus, common nuggets are often repeated. In order to deal with this, we employed a simple heuristic to remove duplicate information: if any two responses share more than sixty percent of their keywords, one of them is randomly thrown out.

3.3 Dictionary Lookup

Another component of our system for answering definition questions utilizes an existing Web-based dictionary for nuggets. Obviously, such an approach cannot be applied directly, because all nuggets must originate from the AQUAINT corpus. To address

this issue, we developed answer projection techniques to “map” dictionary definitions back onto AQUAINT documents. The mapping component is based on the idea that if you already know the answer, it is much easier to find relevant nuggets in the corpus.

Given the target term, our dictionary wrapper goes online to the Merriam-Webster website and fetches the term’s definition. Keywords from the definition are used as the query to the Lucene document retriever. Once a set of candidate documents has been returned, we break each document into sentences and score each sentence based on its keyword overlap with the dictionary definition. The sentences with the highest scores are retained and, if necessary, shortened to one hundred characters centered around the target term.

3.4 Document Lookup

As a last resort (i.e., if no answers are found with the first two techniques), our system employs standard document retrieval to extract relevant nuggets. The target term is used as a Lucene query to gather a set of candidate documents. These documents are chunked into separate sentences and those sentences containing the target term are retained as responses. As before, these sentences are shortened appropriately if needed.

3.5 Answer Merging

The input to the answer merging stage is a series of one hundred character responses from each of the sources: database, dictionary, and corpus. The responses are arranged according to an ad-hoc priority scale we developed based on the accuracy of each approach. For example, we found that verb patterns generally return very good nuggets, and copular constructions are often less accurate. The priority of dictionary answers lies somewhere between the best and worst patterns, ordered such that some dictionary responses (if any) would always be returned in the final answer. Responses extracted directly from document lookup are used only if the two other methods return no answers: document lookup is considered a strict back-off method used only as a last resort. See (Hildebrandt et al., 2004) for a more detailed description of our ordering algorithm.

Finally, the answer merging stage of our system also decides the number of one hundred character responses to return. Since the length penalty for returning long answers is not very steep, we decided to return longer answers in hopes of including more relevant nuggets. Given n responses, we calculated the final number of responses to return as:

$$\begin{array}{ll} n & \text{if } n \leq 10 \\ n + \sqrt{n - 10} & \text{if } n > 10 \end{array}$$

This strategy ensures that our system will always return a generous number of nuggets, and has proven to work well empirically.

4 Factoid Questions

Our system for answering factoid questions was largely unchanged from last year. We employed the Aranea question answering system (Lin et al., 2002; Lin and Katz, 2003), which embraces two different views of the World Wide Web: as a heterogeneous collection of unstructured documents and as a source of carefully crafted and organized knowledge about specific topics.

Aranea’s approach is primarily motivated by an observation that the distribution of user queries qualitatively obeys Zipf’s Law—a small fraction of question types accounts for a significant fraction of all question instances. Large classes of commonly-occurring questions translate naturally into database queries and are handled by Aranea using a technique we call *knowledge annotation*, which allows our system to access semistructured and heterogeneous data as if it were a uniform database. In addition, we have discovered that a simple object–property–value data model captures the content of both Web resources and natural language questions (Katz et al., 2002). To take advantage of these observations, we have built a framework of site-specific wrappers that provide uniform access to knowledge contained in a variety of Web resources. These wrappers are connected to natural language questions through parameterized schemata.

As with all Zipf curves, there is a broad tail where individual instances are either unique or account for an insignificant fraction of total questions. To answer questions that cannot be easily classified into common categories or grouped by simple patterns, Aranea employs what we call redundancy-based *knowledge mining* techniques. Knowledge mining leverages the massive amounts of information available on the Web to overcome many thorny problems associated with natural language processing. The insight is simple: the more data available, the greater the chance that the answer to a natural language question is stated as a reformulation of that question. In such cases, simple pattern matching techniques suffice to accurately extract answers.

The setup of the TREC evaluation requires each answer to be supported with a document from the AQUAINT corpus. Since Aranea does not use the AQUAINT corpus in the question answering process, Web-based answers must then be “projected” back onto AQUAINT documents. Answer projection is accomplished in a two-step process: first, a set of candidate documents is gathered; then, a modified passage retrieval algorithm scans the documents to pick the best document. For obtaining the

set of candidate documents, we tried three different approaches: using the NIST-supplied PRISE documents, using documents generated by our first query generation algorithm (Section 2.1.1), and using documents generated by our second query generation algorithm (Section 2.1.2).

After a set of candidate documents has been gathered, the answer projection module applies a modified window-based passage retrieval algorithm to score the documents. Each 140-byte window is given a score equal to the number of times keywords from both the question and candidate answer appear, with the restriction that at least one keyword from the question must appear in the passage. The score of a document is simply the score of the highest scoring passage. The highest scoring document is paired with the Web-derived candidate answer as the final response unit.

5 Results

A summary of our results at this year’s TREC evaluation is shown in Table 1. Out of twenty-five groups, we ranked sixth in factoid questions, third in list questions, and eighth in definition questions. Our final weighted score ranked us sixth out of all the participating groups. For factoid questions, the query generation algorithms used for answer projection in each of the runs are shown in Table 2. For list questions, the query generator and passage scoring algorithms used for each of the runs are also shown in Table 2. For definition questions, all three submissions were exactly the same.

5.1 List Results

In this section, we discuss our results for answering list questions with respect to query generation, passage retrieval, and the question focus/target type.

5.1.1 Query Generation

Our second query generation method performed slightly better than our first query generation method. In particular, tokenization of multi-token expressions had the biggest positive impact on performance. Consider the following question:

What countries have had school bus accidents that resulted in fatalities? (q2180)

The second query generation algorithm correctly identified “school bus” as a collocation and thus never broke up the expression into “school” and “bus”.

5.1.2 Passage Retrieval

In general, the modified IBM passage scoring algorithm performed slightly worse than the original IBM algorithm. However, since they returned exactly the same responses most of the time, it is difficult to determine if the score differences are

above the margin of error inherent in human judgments. In retrospect, we believe that our modified IBM algorithm was too lax in matching various forms of expansions (too high a score was given to variants). It is a well-known result that uncontrolled expansion of lexical-semantic relations (e.g., synonyms and hyponyms) results in lower performance (Voorhees, 1994). It has likewise been shown that inflectional and derivational expansion does not significantly increase performance. However, these previous experiments were focused solely on document retrieval, using queries that were typically much longer than TREC-style natural language questions. For the question answering task, we believe that linguistically-motivated query expansions will have a positive impact on performance. While our experiments have not yet shown a significant overall positive effect, we attribute this to implementational deficiencies in our overall system, rather than conceptual shortcomings.

As an illustrative example, we present a case where matching of expanded terms did increase performance:

What countries still have royalty? (q2250)

The original IBM passage retrieval algorithm did not return any correct answers, whereas the modified version returned four correct answers. The performance increase can be directly attributed to looser query matching of expanded terms. The original algorithm found passages related to the economic sense of royalty, whereas the modified algorithm retrieved passages with the correct sense related to monarchy.

5.1.3 Question Focus and Target Type

Our strategy for answering list questions crucially depends on correctly identifying the question focus and the associated target type (the ontological type of entity sought after). For a few questions, our system was unable to correctly determine the question focus, resulting in a score of zero for those questions. To address this shortcoming, we will improve START’s ability to recognize question focus.

Although identifying the question focus helps in answering a question, care is needed to map the focus word into a corresponding target type (the specific ontological category). Consider the following questions:

List the names of cell phone manufacturers. (q2096)

Name recipients of funds given by the various foundations of Bill and Melinda Gates. (q2291)

Our system correctly identified “manufacturer” as the question focus in the first question, but chose the

Task	MITCSAIL03a	MITCSAIL03b	MITCSAIL03c	best	median	best rank
Factoid	0.293	0.295	0.291	0.7	0.149	6th
List	0.13	0.118	0.134	0.396	0.053	3rd
Definition	0.309	0.282	0.282	0.555	0.189	8th
weighted total	0.256	0.248	0.250	0.559	0.135	6th

Table 1: Summary of MIT CSAIL submissions.

	MITCSAIL03a	MITCSAIL03b	MITCSAIL03c
List questions:			
Query generator	method 1	method 2	method 2
Passage retriever	IBM	IBM	modified IBM
Factoid questions:			
Answer projection	PRISE	method 1	method 2

Table 2: Variations in each of the TREC runs.

wrong sense for the target type. The term was on our list of professions, so the system incorrectly looked for personal names. The second question demonstrates that not all focus terms, even when correctly identified, are useful. “Recipients” are so general that they can be anything: people, companies, organizations, and even countries.

Not surprisingly, our system performed well for questions whose target type had corresponding fixed lists in our knowledge base. Since we had exhaustive lists for entities like cities, countries, and U.S. presidents, all our answers were at least of the correct type. However, since our system ignored syntactic relations within the passage, it often overgenerated wrong answers. Consider the following question:

What countries have won the men’s World Cup for soccer? (q2346)

Since our system returned all countries found near the relevant keywords, most of the answers were countries that played in the World Cup, not winners of it. As a result, we obtained high recall, but poor precision, on this question. This is certainly a case where the use of syntactic relations can dramatically improve question answering performance (Katz and Lin, 2003).

Our backoff method of looking for the question focus in candidate answers worked for the following question:

What grapes are used in making wine? (q1940)

The system extracted correct answers like “Chardonnay Grapes”. However, the same technique didn’t work when the question focus was “team” or “food” because journalists typically do not write “X team” or “Y food”.

5.2 Definition Results

Although the responses were identical in each of our three submitted runs for definition questions, the scores were not; that is, given the same exact answer string, assessors came up with different judgments some of the time. Out of the 317 responses we submitted for the 50 definition questions, there were 19 responses which were not judged the same over all three runs. However, 7 of these were cases where assessors found the same nugget in different responses for a question. In addition, there are clear instances where an answer nugget is in one of our responses and the assessors missed it, even when the nugget was present word for word. Voorhees’ analysis (2003) of the definition results indicates that the margin of judging error was 0.043, i.e., scores for pairs of identical runs differed by as much as 0.043 (F-measure). Furthermore, due to the small testset size, a score difference of at least 0.1 in F-measure is required in order for two evaluation results to be considered statistically different (at 95% confidence).

Target term extraction was the single biggest source of error in answering definition questions. If the target term is not correctly identified, then all subsequent modules have little chance of providing relevant nuggets.

We did not anticipate the presence of stopwords in names. Consider the following questions:

What is Bausch & Lomb? (q1917)
 Who is Vlad the Impaler? (q1933)
 Who is Akbar the Great? (q1955)

Our naive pattern-based target extractor identified “Lomb”, “Impaler”, and “Great” as the target terms for the above questions, respectively. Fortunately, “Impaler” is such a rare word that we actually returned nuggets concerning “Vlad the Im-

	MITCSAIL03a		MITCSAIL03b		MITCSAIL03c	
	PRISE		Method 1		Method 2	
Right	121	29.30%	122	29.54%	120	29.06%
Inexact	18	4.36%	15	3.63%	15	3.63%
Unsupported	26	6.30%	21	5.08%	21	5.08%
Wrong	248	60.05%	255	61.74%	257	62.23%
Total	413		413		413	

Table 3: Detailed analysis of factoid questions.

paler”. Similarly, “Lomb” so frequently co-occurs with “Bausch & Lomb” that our system was able to provide relevant nuggets. However, since “Great” is a very common word, our definitions for “Akbar the Great” were meaningless.

The system’s inability to parse certain forms of names is related to our simple assumption that the final consecutive sequence of capitalized words in a question is the target. However, this turned out to be an incorrect assumption:

Who was Abraham in the Old Testament?
(q1972)
What is ETA in Spain? (q1987)
What is Friends of the Earth? (q2222)

Our pattern-based target extractor marked “Old Testament”, “Spain”, and “Earth” as the targets for those questions, respectively. The inability to correctly identify the target term resulted in the system’s failure to return relevant nuggets.

Another problem our target extractor encountered is apposition. Take the following example:

What is the medical condition shingles?
(q2348)

The target extractor incorrectly identified “medical condition shingles” as the target term. As a result, our system did not identify a single relevant nugget. To better extract target terms for definition questions, we will employ START and Sepia in the future, which we were unable to utilize for definition questions this year for technical reasons.

5.3 Factoid Results

Table 3 shows a detailed analysis of factoid questions. As in previous years, answer projection appears to be the Achilles’ heel in our Web-based question answering strategy, as shown by the relatively large fraction of unsupported and inexact answers (in comparison to typical results of other teams). Furthermore, it does not appear that any of our more advanced query generation algorithms had any significant impact of the final score of factoid questions.

6 Conclusion

The focus of our research this year was to integrate Web- and corpus-based question answering techniques under a unified framework. This falls under our general research agenda of employing linguistic techniques, at the lexical, morphological, syntactic, and semantic levels, in conjunction with statistical techniques when appropriate. Although our TREC experiments have yet to show significant benefits from linguistically-informed processing techniques, we believe that high performance in the question answering task can only be achieved through fusion of multiple strategies and multiple resources.

References

- Eric Brill, Jimmy Lin, Michele Banko, Susan Dumais, and Andrew Ng. 2001. Data-intensive question answering. In *Proceedings of the Tenth Text REtrieval Conference (TREC 2001)*.
- Charles Clarke, Gordon Cormack, Thomas Lynam, C.M. Li, and Greg McLearn. 2001. Web reinforced question answering (MultiText experiments for TREC 2001). In *Proceedings of the Tenth Text REtrieval Conference (TREC 2001)*.
- Susan Dumais, Michele Banko, Eric Brill, Jimmy Lin, and Andrew Ng. 2002. Web question answering: Is more always better? In *Proceedings of the 25th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2002)*.
- Wesley Hildebrandt, Boris Katz, and Jimmy Lin. 2004. Answering definition questions with multiple knowledge sources. In *Proceedings of the 2004 Human Language Technology Conference and the North American Chapter of the Association for Computational Linguistics Annual Meeting (HLT-NAACL 2004)*.
- Abraham Ittycheriah, Martin Franz, Wei-Jing Zhu, and Adwait Ratnaparkhi. 2000. IBM’s statistical question answering system. In *Proceedings of the Eighth Text REtrieval Conference (TREC-9)*.
- Abraham Ittycheriah, Martin Franz, and Salim Roukos. 2001. IBM’s statistical question answering system—TREC-10. In *Proceedings of the Tenth Text REtrieval Conference (TREC 2001)*.

- Boris Katz and Jimmy Lin. 2003. Selectively using relations to improve precision in question answering. In *Proceedings of the EACL 2003 Workshop on Natural Language Processing for Question Answering*.
- Boris Katz, Sue Felshin, Deniz Yuret, Ali Ibrahim, Jimmy Lin, Gregory Marton, Alton Jerome McFarland, and Baris Temelkuran. 2002. Omnibase: Uniform access to heterogeneous data for question answering. In *Proceedings of the 7th International Workshop on Applications of Natural Language to Information Systems (NLDB 2002)*.
- Boris Katz. 1997. Annotating the World Wide Web using natural language. In *Proceedings of the 5th RIAO Conference on Computer Assisted Information Searching on the Internet (RIAO '97)*.
- Jimmy Lin and Boris Katz. 2003. Question answering from the Web using knowledge annotation and knowledge mining techniques. In *Proceedings of Twelfth International Conference on Information and Knowledge Management (CIKM 2003)*.
- Jimmy Lin, Aaron Fernandes, Boris Katz, Gregory Marton, and Stefanie Tellex. 2002. Extracting answers from the Web using knowledge annotation and knowledge mining techniques. In *Proceedings of the Eleventh Text REtrieval Conference (TREC 2002)*.
- Gregory A. Marton. 2003. Sepia: Semantic parsing for named entities. Master's thesis, Massachusetts Institute of Technology.
- Stefanie Tellex, Boris Katz, Jimmy Lin, Gregory Marton, and Aaron Fernandes. 2003. Quantitative evaluation of passage retrieval algorithms for question answering. In *Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2003)*.
- Ellen M. Voorhees. 1994. Query expansion using lexical-semantic relations. In *Proceedings of the 17th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR-1994)*.
- Ellen M. Voorhees. 2003. Overview of the TREC 2003 question answering track. In *Proceedings of the Twelfth Text REtrieval Conference (TREC 2003)*.