

Building a Foundation System for Producing Short Answers to Factual Questions

Sameer S. Pradhan^{*}, Gabriel Illouz^{†§}, Sasha J. Blair-Goldensohn[†],
Andrew Hazen Schlaikjer[†], Valerie Krugler^{*}, Elena Filatova[†], Pablo A. Duboue[†], Hong Yu[†],
Rebecca J. Passonneau[†], Wayne Ward^{*}, Vasileios Hatzivassiloglou[†], Dan Jurafsky^{*},
Kathleen R. McKeown[†], and James H. Martin^{*}

^{*}Center for Spoken Language Research
University of Colorado
Boulder, CO 80309, USA

[†]Department of Computer Science
Columbia University
New York, NY 10027, USA

Abstract

In this paper we describe the goals of question answering research being pursued as a joint project between Columbia University and the University of Colorado at Boulder as part of ARDA’s AQUAINT program. As a foundation for targeting complex questions involving opinions, events, and paragraph-length answers, we recently built two systems for answering short factual questions. We submitted results from the two systems to TREC’s Q&A track, and the bulk of this paper describes the methods used in building each system and the results obtained. We conclude by discussing current work aiming at combining modules from the two systems in a unified, more accurate system and adding capabilities for producing complex answers in addition to short ones.

1 Introduction

The Department of Computer Science at Columbia University (CUCS) and the Center for Spoken Language Research (CSLR) at the University of Colorado at Boulder are collaborating to develop new technologies for question answering. This project is supported by the ARDA AQUAINT (Advanced QUestion Answering for INTelligence) program. The project plans to integrate robust semantics, event detection, information fusion, and summarization technologies to enable a multimedia question answering system. The goal is to develop a system capable of answering complex questions; these are questions that require interacting with the user to refine and clarify the context of the question, whose answer may be located in non-homogeneous databases of speech and text, and for which presenting the answer requires combining and summarizing information from multiple sources and over time. Generating a satisfactory answer to complex questions requires the ability to collect all relevant answers from multiple documents in different media, weigh their relative importance, and generate a coherent summary of the multiple facts and opinions reported.

In order to achieve these goals, we are developing and integrating four core technologies:

- ▷ **Semantic annotation (CSLR)** — We use a shallow, domain-independent, semantic representation and a statistical semantic parser for building this representation. The semantic representation is a basic building block for dialog management, event detection, and information fusion.

[§]Currently at LIMSI, France.

- ▷ **Context management (CSLR)** — We are developing a dialogue interface to allow the system to carry on a focused dialogue with users to answer queries. The interface will maintain context through the interaction to allow follow-up questions and will conduct clarification dialogues with the user to resolve ambiguities and refine queries.
- ▷ **Event recognition and information tracking (Columbia)** — An event is an activity with a starting and ending point, involving a fixed set of participants. We are investigating methods for identifying atomic events within each input document by extracting information about named entities that play a prominent role in the document and about the time period that the text covers. We will rely on the semantic representation of documents to allow us to identify participants and their functions.
- ▷ **Information fusion and summary generation (Columbia)** — Rather than listing a set of relevant responses to a question, we are investigating techniques to integrate summarization and language generation to produce a brief, coherent, and fluent answer. Critical to this task is the problem of selecting fragments of text from different documents that should be included in the answer and determining how to combine them, removing redundancy and integrating complementary information fluently.

As a product of the project, we envision a question answering system that will handle these complex questions and answers. Towards that goal, we built foundation systems that handle factual questions with short answers, as specified in TREC's Q&A track. Prior to the start of this project, neither site had developed a complete question answering system, so we had to build the foundation components and put them together in approximately six months. We elected to build two independent systems as a first step, one at each site, rather than attempting to integrate components across the two sites in such a short period. This gave us the added benefit that each site was able to focus more of their effort on specific components they were interested in, and avoided the need for developing complex protocols for communication between the modules across the sites. As newcomers to the TREC Q&A task, our goal was not to compete with the best systems, but rather to implement in a short amount of time baseline components and systems integrating these components so we could start answering at first simple questions. Currently, we are taking the best components from each site and integrating them in a unified system with a well-defined interface between modules. Separately from the TREC effort, we have developed components for extracting and using features that will eventually play a role in our open-ended question answering system. We are now adding some of the functionality of these components to our foundation systems. However, these advanced features, such as semantic parsing and event detection, were not applied to our TREC system and hence are not described in detail in this paper.

The paper focuses on our development of our foundation question answering systems for TREC, processing factual questions only and producing short answers. Most of the remainder of the paper (Sections 2 and 3 discusses the two architectures we developed, ways that each departs from standard question answering methodology, and our results on this year's TREC questions. We conclude with a discussion of our current integration effort and ongoing work on adding advanced components for handling additional question types to our foundation system.

2 System A — The Boulder System

The novel feature of our approach in System A is the use of shallow semantic representations to enhance potential answer identification. Most successful systems first identify a list of potential answer candidates using pure word-based metrics. Syntactic and semantic information at varying granularity is then used to re-rank those candidates [10, 9]. However, most of these semantic units are quite specific in what they label. We identify a small set of *thematic roles*—viz., *agent*, *patient*, *manner*, *degree*, *cause*, *result*, *location*, *temporal*, *force*, *goal*, *path*, *percept*, *proposition*, *source*, *state*, and *topic*—in the candidate answer sentences, using a statistical classifier [8]. The classifier is trained on the FrameNet database [2]. This is an online lexical resource for English, based on *frame semantics*. It contains hand-tagged semantic annotations of example sentences from a large text corpora, and covers thousands of lexical

items—verbs, nouns, and adjectives—representative of a wide range of semantic domains. We map the more specific *frames* and the corresponding *frame elements* to the reduced, more general set before training the classifier.

2.1 Architecture

The following sequence of actions will be taken in response to an input query:

1. Question type classification—Identify the Named Entity and Thematic Role of the expected answer type. This also defines a set of answer type patterns, and includes named entity tagging and parsing the question for thematic roles.
2. Focus identification—Identify certain salient words/phrases in the question that are very likely to be present in the answer string in one form or the other.
3. Extract a set of query words from the question, and apply semantic expansion to them.
4. Submit the query words to the *mg* (Managing Gigabytes) [14] IR engine and get back a rank-ordered set of documents.
5. Keep the top N (approximately 500) documents and prune the rest.
6. Segment documents into paragraphs and prune all but the top N' paragraphs.
7. Generate scoring features for the paragraphs, including named entity tagging and parsing of paragraphs to add thematic roles.
8. Re-rank documents based on the set of features that we compute, including answer type patterns. Some of the answer type patterns are based on the semantic labels.
9. Compute for each paragraph a confidence measure that it contains some relevant information. This includes N -Best count as one of the features.
10. Send tagged paragraphs that exceed a confidence threshold for extraction of the short answer required for TREC.

For the problem of question answering, we are more concerned with precision than recall, so we have to be careful in expanding the query words to get answers that are expressed in words quite different from the ones mentioned in the question. Semantic expansion will be performed when the system's confidence in the best candidate answer string without expansion is found to be below a certain threshold. Our mechanism for expansion is:

- a. Submit original query words to IR engine and get back a rank-ordered set of documents.
- b. Generate set of target words from top n documents based on the *idf* values of the words.
- c. Generate a set of target words from WordNet [7] synsets of original keywords.
- d. Take the intersection of the two sets and add to the keyword set.

2.2 Features

2.2.1 Answer Identification

We now discuss the features used for ranking the documents. The features are roughly ordered by decreasing salience.

- ▷ **Answer type** — In order to extract the answer to a question, the system needs to identify the expected answer type. Answers for short answer questions generally can be categorized as named entities and/or propositions. Summary information is often required for descriptive and definition questions. The system classifies the answer type by two features: 1) named entity class and 2) thematic role. Named entity class specifies one (or more) of 56 classes as the named entity class of the answer. We use 54 real named entity classes, one class representing the case where the answer is expected to be a named entity but one not known to the system, and one class for cases where the answer is not expected to be a named entity. The thematic role class identifies the thematic role in which the potential answer would tend to be found in the answering sentence.

- ▷ **Answer surface patterns** — Once the question type is known the next step is to identify candidate answers. One technique we use is based on generating a set of expected surface patterns for the answer. Sentences or snippets matching these patterns would get better scores than ones that did not. The patterns specify word and named entity-based regular expressions that are derived from a large corpus annotated with named entities; (simplified) examples include:
 1. Some common question types, e.g., [<PERSON_DESCRIPTION><PERSON_NAME>] for questions like “Who is <PERSON_NAME>?”; [<ORGANIZATION>, <CITY>, <STATE>, <COUNTRY>] for questions asking about the location or address of an organization.
 2. Likely re-phrasings of the question, e.g., [<PERSON> invented <PRODUCT>] for questions like “Who was the inventor of <PRODUCT>?”
 3. Occurrence statistics of the pattern in the corpus, e.g., [<PERSON> (<YEAR>-<YEAR>)] for birth dates of people.

- ▷ **Named entities in answer** — In the case of questions that expect a specific named entity (including the unknown named entity type) as the answer, candidates that do not contain that named entity are penalized. In the case that the answer is expected to be an unknown named entity, then candidates that contain an untagged sequence of capitalized words (a strong indicator of unknown named entities) are preferred.

- ▷ **Presence of focus word** — The presence of the focus word is an important feature for the overall score. For our purposes, a focus word is a word in the question that, or its synonym, is very likely to appear in the sentence that contains the answer. Even if all the other keywords are present in a candidate answer, but the focus word is not present, then it is not a sufficiently good candidate. The focus word can also be considered as a necessary, but not a sufficient condition for a string to qualify as a candidate answer.

- ▷ **Thematic role patterns** — While surface patterns for answers can provide valuable information when a match is found, the specific nature of the patterns and the limited occurrences of the answer string within the reformulations obtainable from the question do not always guarantee a surface pattern match. We also provide a more general set of expected answer patterns based on thematic roles. We expect that these patterns will have higher coverage than the more specific surface patterns. This feature scores sentences based on the presence of expected thematic roles and named entities existing in specific thematic roles.

Thematic role patterns serve two main purposes: Firstly, they help identify false positive answer candidates that are ranked above the correct answer candidate, and secondly, they help extract the exact answer boundary from the string, especially if the answer is not a named entity. This can be illustrated with the example in Figure 1.

- ▷ **Okapi scoring** — This is the score assigned by the information retrieval engine to the paragraph extracted in the very beginning. The formula for calculating the score [13] uses features like the number of query words present in the document, inverse document frequency of the words, the length of the document, etc. We update the Okapi score for the sentences after extracting them from paragraphs—that is, if the system expects it to be a short answer and decides on reducing the granularity below the paragraph level.

- ▷ **N-gram**—Another feature that we use is based on the length of the longest n -gram (sequence of contiguous words) in the candidate answer sentences after removing the stopwords from both the question and the answer.

- ▷ **Case match**— Documents with words in the same case tend to be more relevant than those that do not have the same case words in them, so the former are given a relatively higher score.

| |
|--|
| <p><i>Question:</i> Who assassinated President McKinley?</p> <p><i>Parse:</i> <code>[role=agent Who] [target assassinated] [role=patient [ne=person_description President] [ne=person McKinley]]?</code></p> <p><i>Keywords:</i> assassinated President McKinley</p> <p><i>Answer named entity (ne) Type:</i> Person</p> <p><i>Answer thematic role (role) Type:</i> Agent of target synonymous with “assassinated”</p> <p><i>Thematic role pattern:</i> <code>[role=agent [ne=person ANSWER] ^ [target synonym_of (assassinated)] ^ [role=patient [ne=person reference_to(President McKinley)]]</code></p> <p>This is one of possibly more than one patterns that will be applied to the answer candidates.</p> |
| <p><i>False Positives:</i></p> <p>Note: The sentence number indicates the final rank of that sentence in the returns, without using thematic role patterns.</p> <ol style="list-style-type: none"> In <code>[ne=date 1904]</code>, <code>[ne=person_description President] [ne=person Theodore Roosevelt]</code>, who had succeeded the <code>[target assassinated] [role=patient [ne=person William McKinley]]</code>, was elected to a term in his own right as he defeated <code>[ne=person_description Democrat] [ne=person Alton B. Parker]</code>. <code>[ne=person Hanna]</code>'s worst fears were realized when <code>[role=patient [ne=person_description President] [ne=person William McKinley]]</code> was <code>[target assassinated]</code>, but the country did rather well under TR's leadership anyway. <code>[ne=person Roosevelt]</code> became president after <code>[role=patient [ne=person William McKinley]]</code> was <code>[target assassinated] [role=temporal in [ne=date 1901]]</code> and served until <code>[ne=date 1909]</code>. |
| <p><i>Correct Answer:</i></p> <ol style="list-style-type: none"> <code>[role=temporal In [ne=date 1901]]</code>, <code>[role=patient [ne=person_description President] [ne=person William McKinley]]</code> was <code>[target shot] [role=agent by [ne=person_description anarchist] [ne=person Leon Czolgosz]] [role=location at the [ne=event Pan-American Exposition] in [ne=us_city Buffalo], [ne=us_state N.Y.] [ne=person McKinley]</code> died <code>[ne=date eight days later]</code>. |

Figure 1: An example where thematic role patterns help constraint the correct answer among competing candidates. All the named entities, but only roles pertaining to the *target* predicate are marked in the sentences.

2.2.2 Confidence Annotation

Once the likely answer candidates (paragraphs or sentences) are extracted for a question we need to estimate the likelihood of those being *good* answers. To do this, we have a scheme that annotates each with some level of confidence. The features that we use to estimate this confidence are:

- ▷ **N-Best count** — This is the frequency distribution of n extracted answer strings in the top N final answer candidates. We calculate prior statistics on the probability that an answer is correct if it occurs n times in the top N candidates.
- ▷ **Named entity class** — The system tends to answer questions about some named entity classes more accurately than others. We calculate prior statistics for the percent correct answers for each named entity class.

- ▷ **N-gram length** — We calculate prior statistics for the probability that an answer is correct if it contains an n -gram of length n .

These features are combined as a weighted linear combination of the individual estimates.

2.3 Current Implementation

In this section we discuss the state of the current implementation for system A, which was used to produce the first run submitted to the TREC 2002 question answering track.

2.3.1 TREC-2002 Database

The text database comprises non-stemmed, case-folded indexing of the TREC/AQUAINT corpus using a modified version of the `mg` (Managing Gigabytes) search engine [14] that incorporates the Okapi BM25 ranking scheme [13] and an expanded set of characters forming an indexed unit—so as to accommodate hyphenated words, URLs, emails etc. Each indexed document is a collection of segmented sentences forming a paragraph in the original corpus. For efficiency reasons, we use another index of the same documents pre-tagged with 29 named entities using BBN’s Identifier [3] and another rule-based tagger that tags some specific named entity sub-categories inside those identified by Identifier and some other closed-class named entities that have been frequent answer types in the previous TREC questions, bringing the total to 54.

2.3.2 Answer Identification

We currently use a rule-based question classifier that identifies the named entity type and thematic role of the expected answer to each question. Keywords are extracted from the questions by removing a small set of stopwords and some punctuation symbols that are not considered part of the indexed unit. For each query, the top N ranked documents are retrieved for processing. N is currently fixed at 2500, based on an acceptable recall level of about 85% on a sample set of past TREC questions. A set of documents are carried over from the list of documents retrieved by the IR engine, using gradually diluted boolean filters, beginning with one that is formed by AND-ing all the keywords, and then dropping the smallest length non-noun, non-adjective, non-capitalized and non-date words one-by-one, until there are no more keywords to drop, or the cumulation of filtered documents exceeds a threshold of n . The value of n is empirically set to 10. We call this the *boolean peel-off* strategy. The answer named entity type is used to filter out documents that do not contain the required named entity type. For ones that do not expect a known named entity, the ranked list is kept as is. The re-ranking of documents using the features is implemented as a series of major and minor sorts. Documents with words in the same case as the query are preferred over ones with a different case—within the group of documents containing the same number of query words, while preserving the relative Okapi ranking. Moreover, words in the same sequence as in the query are preferred over ones in a different sequence, and n -grams are promoted within the former. The paragraphs are then broken down into sentences, which are ranked using the same criterion.

2.3.3 Answer Extraction

For questions in which the answer is a specific named entity or a thematic role, if the top ranking sentence contains only one instance of that element, that instance is returned as the answer. In many cases, however, there is more than one element of the predicted type. In such cases, the system selects the element with the shortest average distance from each of the query words present in that snippet. There are penalties added for some punctuation symbols like commas, semi-colons, hyphens etc. In cases when the required answer is not a known named entity, the answer extractor tries to find a sequence of capitalized words that could be a probable named entity. In case the expected answer is not a named entity, and the thematic role cannot be identified without much ambiguity, then the system tries to find an apposition

| # wrong (W) | # unsupported (U) | # inexact (X) | # right (R) | CW Score | NIL Precision | NIL Recall |
|-------------|-------------------|---------------|-------------|----------|---------------|-------------|
| 411 | 7 | 3 | 79 | 0.226 | 0 / 9 = 0.000 | 0/46 = 0.00 |

Table 1: System A (Boulder) results; numbers are out of 500 questions.

near the question words, and extracts that as the answer. An example would be definition questions, of the style “What is X?”. Failing to get any of the above, the system replies that it does not have an answer to that particular question in the given corpus.

2.3.4 Confidence Annotation

The prior probability of correct answers for a particular question type, along with the n -best count, are used to assign a confidence measure. The system then additionally tests whether the candidate at rank two has more counts in the top ten candidates than the candidate at rank one. If so, it is promoted to rank one.

2.4 Results

The results of System A on the TREC-2002 test set are presented in Table 1. These results are consistent with what we expected based on our TREC 9/10 dev-test set. Note that we answered “no answer” or NIL only if no answer was returned by our standard algorithm, which was rarely.

3 System B — The Columbia System

3.1 Overview of System Operation

Our second foundation system, developed at Columbia University, uses a hybrid architecture that was initially dominated by the two processes of query expansion and candidate-answer scoring. This focus was a consequence of developing the original architecture in the context of web-based search for open-ended questions. By applying this original architecture to the task of providing short answers to factual questions over a specific document collection, we had the opportunity to explore strategies for performing search in parallel over open and closed collections, then combining the results. For example, as described in more detail in Section 3.5, we experimented with a strategy that interleaves information learned from the web search to re-query the TREC document set.

In this section, we give an overview of the system architecture and its evolution in the context of the TREC task. In our baseline system designed for web question-answering, we adopted the following two working assumptions:

1. Given the size and redundancy of the web, a database of paraphrase patterns with more specific patterns ordered first would be both necessary and effective for finding relevant documents (i.e., we aimed at precision rather than recall).
2. Scores for candidate answers would be a composite function of attributes of the query formulation process, and of the space of candidate answers. For example, more specific queries have higher default scores; an answer string retrieved from multiple documents is weighted higher than an answer string that appears in relatively few documents.

Finally, we assumed that the question of how best to capitalize on these two assumptions would be primarily an empirical one. Thus questions that are superficially similar might require rather distinct queries. For example, questions that contain a noun denoting a leadership position, as in “Who was the first commander of FLEET X?” and “Who was the first coach of TEAM Y?” might both benefit from query expansion rules in which the related verbs

appear (e.g., “<TERM> commanded FLEET X” and “<TERM> coached TEAM Y”), but this would not be the case for the structurally similar question “Who was the first general of FORCE Z?”. Furthermore, for domains that are well-represented on the web, more complex or more subtle scoring metrics might be required in comparison with domains that are not well-represented.

The remainder of this section documents four key modules of System B’s architecture. There is a pipeline of three modules responsible for distinct phases of the query expansion process: paraphrasing of patterns derived from the question string (Section 3.2 below); modification of queries, e.g., by inserting terms likely to occur with the answer ([1]; Section 3.3); and term prioritization (Section 3.4). At each phase of the query expansion process we modify according to the results of that phase scoring weights that guide how likely that version of the query is to lead to the correct answer. Finally, after a candidate answer set has been assembled, these weights are assembled into a single score (Section 3.5). If, as in the case of our second run submitted to TREC, candidate answer sets for a given question come from distinct document collections, the source collection is considered in ranking potential answers (Section 3.5).

We conclude this section with a brief discussion of the results obtained on the data and questions of the 2002 TREC Q&A track.

3.2 Question Paraphrasing

In order to generate queries providing high precision coverage of the answer space for a given question, custom rules were developed providing a mapping from a given question type to a set of paraphrasing patterns which would generate alternative queries. For example, the question string “Where is the Hudson River located?” may result in the generation of queries including “Hudson River is located in”, “Hudson River is located”, “Hudson River is in the”, “Hudson River is near”, and “Hudson River in”. Since we often do not have specific information about the question target, our paraphrasing rules allow changes on articles (definite, indefinite, or no article), number, and function words to maximize our coverage of the collection.

A two-level scoring scheme was implemented for these queries whereby each was scored based on the specificity of its query string as well as that of the paraphrasing pattern that generated it. Specificity here is defined by the length of the query string—or the length of the shortest possible generated query string in the case of the paraphrasing patterns—which typically leads to fewer returned results than a comparably shorter (more general) query when sent to a search engine. These query scores are used to aid the scoring and subsequent ranking of the returned results along with other factors that rate the results themselves (see Section 3.5).

In order to increase the coverage of the answer space, our paraphrasing patterns include a fall-back mechanism which would generate boolean queries based on keywords and short phrases extracted from the original question string. A separate scoring and ranking scheme was developed for choosing among the keywords, phrases, and the queries that combine them, as described in Section 3.4.

3.3 Query Modification

Knowing the type of the question can be very helpful not only for defining the type of the answer but also for better targeting of the search of the documents which might contain a potential answer. This process can be viewed as connecting the answer and question spaces. The following example illustrates the discrepancy that often exists between the words used in the question (what we call the “question space”), and the words that occur in the correct answer (the “answer space”); this discrepancy goes beyond the usual problem of dealing with synonyms or syntactic variations.

Consider the question “What is the length of the Amazon?”. We observe that the top ranked results returned by Google when this question is directly submitted as a query are not relevant to the question. Any question answering system that reorganizes information from an IR search engine is thus handicapped in this case (none of the top 20 documents contain the answer). The reason is that the question contains only two words that are not filtered as stopwords: “length” and “Amazon”. Of these, “Amazon” is highly ambiguous (the river, the mythological female

warrior, the online company, to name just the more common senses), and “length” is not likely to appear in the answer directly. This second factor represents the discrepancy between the question and answer spaces; while “length” is a very good *question* word, the corresponding good *answer* words are units of length such as “miles” and “kilometers”. If we expand the query by adding either *km* or *mi* the first hits returned by the search engine are about the Amazon river and contain the answer.

In our TREC 2002 system we used query expansion only for questions that required answers consisting of a number plus a measurement unit. All the questions that required a number as an answer were categorized under the general type NUMBER by our question classifier. We built a further classifier that translated some of the question words to subtypes of NUMBER, namely DISTANCE, WEIGHT, SPEED, TEMPERATURE, MONEY, and OTHER. The classifier was constructed by training using RIPPER [5] to produce a set of rules. For each of the questions classified into the above subtypes, the classifier returns an automatically compiled list of words expressing the appropriate measurement units. Subsequently, we expanded the queries by including the corresponding measurement units.

3.4 Query Prioritization

Query prioritization scores queries so that those with highest predicted answer precision (i.e., number of documents retrieved by a query that contain the correct answer(s) divided by the total number of documents retrieved) will be attempted first during document retrieval. This is because, first, our paraphrasing and modification mechanisms can generate thousands of queries and it is impractical to submit all of those to the search engine, and, second, because we can tell that some query rewriting mechanisms are more likely to be accurate than others.

For example, given the question “How many soldiers died in World War II?”, our system will generate many queries, using the paraphrasing (Section 3.2) and modification (Section 3.3) rules mentioned above. The generated queries include “World War II” (q_a) and “soldiers” (q_b). We want query prioritization to assign q_a a higher score since it clearly is more specific and relevant to this question, and thus is likely to have a higher predicted answer precision.

Unlike previous TREC systems [11, 6], which mainly applied heuristics for query prioritization, we empirically built our query prioritization algorithm using statistics collected over previous TREC question-answer pairs. We analyzed the relation between a query term (see Sections 3.2 and 3.3 for query term generation) and its answer precision. We considered various features of a query term, including the term’s syntactic role (e.g., noun phrase, verb phrase, and head noun), morphological features (e.g., upper case for proper noun), and inverse document frequency (IDF). We found that IDF consistently reflected answer precision. We therefore developed a query prioritization algorithm based on IDF in order to maximize the predictive power of our query scores.

The query scoring algorithm for non-paraphrase-based queries relies on two functions:

- ▷ The term-scoring function T maps terms to term scores, where terms are strings of one or more words which are part of a query. The term score of a multiword term X is the product of the IDF values of the individual words forming the term. A minimum IDF value of 1.05 is used in this case to ensure that even common terms slightly boost the score of a multiword term. In addition, a suitably high value is used for words with unknown IDF values. We use IDF values computed over a large body of recent general news text.
- ▷ The query-scoring function Q maps a set of one or more terms (a query) to a query score. The query score for query Y is the product of the term scores of the query’s component terms. Thus the score for query Y composed of terms X_1 through X_m would be

$$Q(Y) = \prod_{i=0}^m \prod_{j=0}^n \text{idf}(x_{ij})$$

where x_{ij} is the j -th word in term X_i .

As mentioned above, the query scores produced are used in document retrieval, answer selection, and answer ordering (for listing first the answers to questions where our system has the greatest confidence). For the latter purpose,

scores must be normalized. This is because answer ordering requires us to compare answer confidence across answers to different questions; since part of our confidence is determined by the score of the query leading to the answer (see Section 3.5), we must be able to make a meaningful comparison between query scores for queries produced for two different questions. Therefore, scores for each set of queries produced for a given question are normalized using division by the highest query score for a query generated for that question. Thus, the highest scoring query produced for any question will have a score of 1, while lower scoring queries will have scores between 0 and 1.

It is important to note that our system scores any paraphrase-based query above any keyword-based query, irrespective of the functions mentioned above.¹ This is because the paraphrasing rules (Section 3.2) were hand-crafted and produced with an eye toward high precision. Thus, the scores produced by this module of the system were only used to differentiate between keyword-based queries.

3.5 Combining Evidence from Multiple Sources

Even when the answer must be found in a specific collection, as in the TREC evaluation, there are benefits in combining evidence from multiple collections: the answer can be found in a larger or broader collection, in which case the smaller collection can be searched again with that specific answer in mind; and, if no answer can be found in the smaller collection, the confidence in the answer from the larger collection can help determine whether “no answer present in the collection” (NIL) should be returned.

To this end, we have designed a general mechanism using *wrappers* that interface our system to different collections and/or search engines using a common API. We have implemented wrappers for the Google, Glimpse [12], and mg [14] search engines, and for TREC-11 we ran system B using a combination of Google on the web (a broad collection) and mg on the TREC/AQUAINT collection (the collection where the answer must be found).

Our system returns a list of answers extracted from each source (search engine and collection combination), which may include the same string multiple times. Each instance of an answer has an associated confidence score, which depends on the query used to retrieve the answer (see Sections 3.2 and 3.4), the confidence score returned by the search engine (not implemented for our TREC experiments), and the confidence of the answer extractor in selecting an appropriate phrase from the returned document. These instance scores are combined for each potential answer (i.e., across all instances where the same string is returned) using the following formula for computing the cumulative score after n instances of the same string have been processed

$$\text{cumulative_score}_n = 1 - [(1 - \text{cumulative_score}_{n-1}) \times (1 - \text{instance_score}_n)]$$

which ensures that answers occurring multiple times are weighted higher, taking into account the evidence for them each time they are found.

An answer may be found in one or more sources. We employ the following algorithm for calculating a composite score based on the cumulative scores of each string proposed as a potential answer, from the web or the TREC collection. The algorithm distinguishes the following cases:

- ▷ No answers found in either collection. Since this probably represents a system failure (it is unlikely that a TREC question would not be answerable from both the web and the TREC collection), we return NIL with zero confidence.
- ▷ One or more answers found in the TREC collection, but no answers found from the Web. We return the best answer extracted from the TREC collection but depreciate its confidence by a constant factor because the coverage of the web would make it unlikely that no answers at all could be found there while the TREC collection would be able to provide one.

¹Queries augmented with the query modification techniques discussed in Section 3.3 are considered keyword-based queries in this context.

- ▷ No answers found in the TREC collection, but one or more answers found from the Web. This is our *prima facie* case for a NIL answer; however, since often we got our answer from the web using a modified query, we re-query the TREC collection using the answer identified from the web. If that step succeeds, we report the answer with a combined confidence as in the next case below; otherwise we report NIL, with confidence equal to our confidence in the web answer.
- ▷ One or more answers found in each collection, and both top ranked answers are the same. We report that answer, with reinforced confidence according to the formula

$$\text{combined_confidence} = 1 - [(1 - \text{confidence_TREC}) \times (1 - \text{confidence_web})]$$

- ▷ One or more answers found in each collection, but the top ranked answers are different. We report the TREC answer, but reduce its confidence by the formula

$$\text{combined_confidence} = \text{confidence_TREC} \times (1 - \text{confidence_web})$$

3.6 Results

For the 500 questions in the TREC-11 question, System B produced 58 correct answers, 8 unsupported, 2 inexact, and 432 wrong. The system’s unweighted precision is 11.6%, while its confidence-weighted score is 0.178, representing a significant boost over the unweighted precision. We produced a lot of NIL answers (195, or 39% of our total answers), which indicates that we were too conservative in choosing low-confidence answers found in the TREC collection. We did retrieve a third of the NIL answers present in the test set, but overall System B performed less well than System A which obtained a confidence-weighted score of 0.226 while producing very few NILs (1.8% of its total answers). On non-NIL answers, the two systems performance is closer (unweighted precision of 16.1% for system A and 14.1% for system B).

We attribute the lower performance of System B in two additional factors beyond the excessive production of NILs: First, we used a very simple extractor for selecting the answer out of the sentence where it was found, and the extractor failed to produce the correct phrase in a number of cases where we succeeded in finding the appropriate sentence. Second, our question classifier was not always successful in predicting the correct question type, producing no label for 77 (15.4%) of the questions. We performed significantly worse on those questions than in *who*, *when*, *where*, or *what* questions with well-extracted types.

In earlier years, the majority of TREC participants scored in the 15–20% range on the accuracy of the top-ranked answer. From the numbers released by NIST this year, we calculated the unweighted precision of a hypothetical system whose likelihood for getting each question right is equal to the percentage of participants who answered that question correctly. That number is 22.3%, indicating that while our two systems are below the TREC average, they are not very far behind.

4 Conclusion and Future Work

We have described two systems that handle questions with short, factual answers. These systems were developed in a very brief time frame, and represent our first entry in the TREC Q&A track. Nevertheless, early indications are that the systems score not much below average compared to the performance of all systems in this year’s evaluation, and comparable to the majority of submitted runs in previous years, when considering the top-ranked answer only.

The two systems represent for us a starting point for a question answering system, to which we are adding additional capabilities. In the months since TREC, we have worked on developing detailed XML-based protocols for communication between modules in a common architecture. The system we are building by combining elements of the two systems discussed in this paper utilizes the best-performing components of each of the current systems. The

modular architecture allows us to connect additional modules, and actual question answering is done in a distributed fashion, with most of question analysis done at Colorado and most of answer synthesis done at Columbia. Modules communicate with a central server via HTTP, while the architecture offers several communication interfaces at different levels of complexity (for example, one module may request only the high-level question type, while another may examine in detail the semantic parse). A web-based client provides a front end to the integrated system, allowing users in different locations to access the system.

Our research is moving towards questions with more complex answers, including opinions, events, definitions, and biographies. We recently completed a prototype module for processing definitional questions, and we are currently adding its functionality to our integrated system, so that for questions of the form “What is X?” we produce both a short, TREC-like, answer and an answer spanning several paragraphs. At the same time, we are continuing to tune individual components to enhance interoperability between the two sites; for example, we recently started re-focusing the semantic analysis (done at Colorado) on types of phrases that are likely to impact the processing of opinions and biographies (done at Columbia).

Acknowledgments

This work, as well as our broader research on question answering, is supported by ARDA AQUAINT contract MDA908-02-C-0008. We would like to thank Ralph Weischedel of BBN Systems and Technologies for giving us the named entity tagger (Identifinder), and Daniel Gildea for the FrameNet parser.

References

- [1] Agichtein, E., Lawrence, S., and Gravano, L., “Learning Search Engine Specific Query Transformations for Question Answering”, in *Proc. of the 10th International World-Wide Web Conference (WWW10)*, 2001.
- [2] Baker, C., Fillmore, C., and Lowe, J., “The Berkeley FrameNet project”, in *Proceedings of the COLING-ACL, Montreal, Canada, 1998*.
- [3] Bikel, D., Schwartz, R., Weischedel, R., “An Algorithm that Learns What’s in a Name”, *Machine Learning*, 34, Kluwer Academic Publishers, 1999, 211–231.
- [4] Brin, S., Page, L., “The Anatomy of a Large-Scale Hypertextual Web Search Engine”, *Computer Networks and ISDN Systems*, Vol. 30, No. 1–7, 1998, 107–117
- [5] Cohen, W. W. (1995). “Fast Effective Rule Induction”, *Proceedings of the Twelfth International Machine Learning Conference (ML-95)*.
- [6] Dumais, S., Banko, M., Brill, E., Lin, J., and Ng A. (2002). “Web question answering: Is more always better?”, in *Proceedings of SIGIR-02*, pp. 291–298.
- [7] Fellbaum, C., *editor*, “WordNet: An Electronic Lexical Database”, MIT Press, 1998.
- [8] Gildea, D., Jurafsky, D., “Automatic Labeling of Semantic Roles”, Technical Report TR-01-005, International Computer Science Institute, Berkeley, 2001
- [9] Harabagiu, S., Moldovan, D., *et al.*, “Answering complex, list and context questions with LCC’s Question-Answering Server”, *The Tenth Text REtrieval Conference (TREC-10)*, Gaithersburg, Maryland, November 13-16, 2001, pp. 355–361.

- [10] Hovy, E., Hermjakob, U., “The Use of External Knowledge of Factoid QA”, *The Tenth Text REtrieval Conference (TREC-10)*, Gaithersburg, Maryland, November 13-16, 2001, pp. 644–652.
- [11] Lee, G. G., Seo, J., Lee, S., Jung, H., Cho, B. H., Lee, C., Kwak, B. K., Cha, J., Kim, D., An, J., Kim, H., and Kim, K., “SiteQ: Engineering High Performance QA System Using Lexico-Semantic Pattern Matching and Shallow NLP”. In *Proc. TREC 2001*.
- [12] Manber, U. and Wu, S. “Glimpse: A tool to search through entire file systems”. In *Proceedings of the Winter USENIX Conference*, January 1994.
- [13] Robertson, S., Walker, S., “Okapi/Keenbow at TREC-8”, *The Eighth Text REtrieval Conference (TREC-8)*, Gaithersburg, Maryland, November 17-19, 1999, 151–162.
- [14] Witten, I., Moffat, A., Bell, T., “Managing Gigabytes: Compressing and Indexing Documents and Images”, Morgan Kaufmann Publishing, San Francisco, May 1999.