

N-gram vs. Keyword-based Passage Retrieval for Question Answering

Davide Buscaldi and José Manuel Gomez and Paolo Rosso and Emilio Sanchis

Dpto. de Sistemas Informáticos y Computación (DSIC),
Universidad Politécnica de Valencia, Spain,
{dbuscaldi,jogomez,proso,esanchis}@dsic.upv.es

Abstract. In this paper we describe the participation of the Universidad Politécnica de Valencia to the 2006 edition, which was focused on the comparison between a Passage Retrieval engine (JIRS) specifically aimed to the Question Answering task and a standard, general use search engine such as Lucene. JIRS is based on n -grams, Lucene on keywords. We participated in three monolingual tasks: Spanish, Italian and French. The obtained results show that JIRS is able to return high quality passages, especially in Spanish.

1 Introduction

Most of the Question Answering (QA) systems that are currently used in the CLEF and TREC¹ QA exercises are based on common keyword-based Passage Retrieval (PR) methods [3,1,8,5]. QUASAR is a mono/cross-lingual QA system built around the JIRS PR system [6], that is based on n -grams instead of keywords. JIRS is specifically oriented to the QA task; it can also be considered as a language-independent PR system, because it does not use any knowledge about the lexicon and the syntax of the language during question and passage processing phases. The crucial step performed by JIRS is the comparison between the n -grams of the question and those of the passages retrieved by means of typical keyword-based search. The main objective of our participation to QA@CLEF2006 was the comparison of JIRS with a classical PR engine (in this case, Lucene²). In order to do this, we actually implemented two versions of QUASAR, which differ only for the PR system used. With regard to the improvements over last year's system, our efforts were focused on the Question Analysis module, which in contrast to the one used in 2005 does not use Support Vector Machines in order to classify the questions. Moreover, all modules are now better integrated in the complete system.

The 2006 CLEF QA task introduced some challenges with respect to the previous edition: list questions, the lack of a label distinguishing “definition” questions from other ones, and the introduction of another kind of “definitions”,

¹ <http://trec.nist.gov>

² <http://lucene.apache.org>

that we named *object definitions*. This forced us to change slightly the class ontology we used in 2005.

In the next section, we describe the structure and the building blocks of our QA system. In section 3 we discuss the results of QUASAR in the 2006 CLEF QA task.

2 Architecture of QUASAR

The architecture of QUASAR is shown in Fig.1.

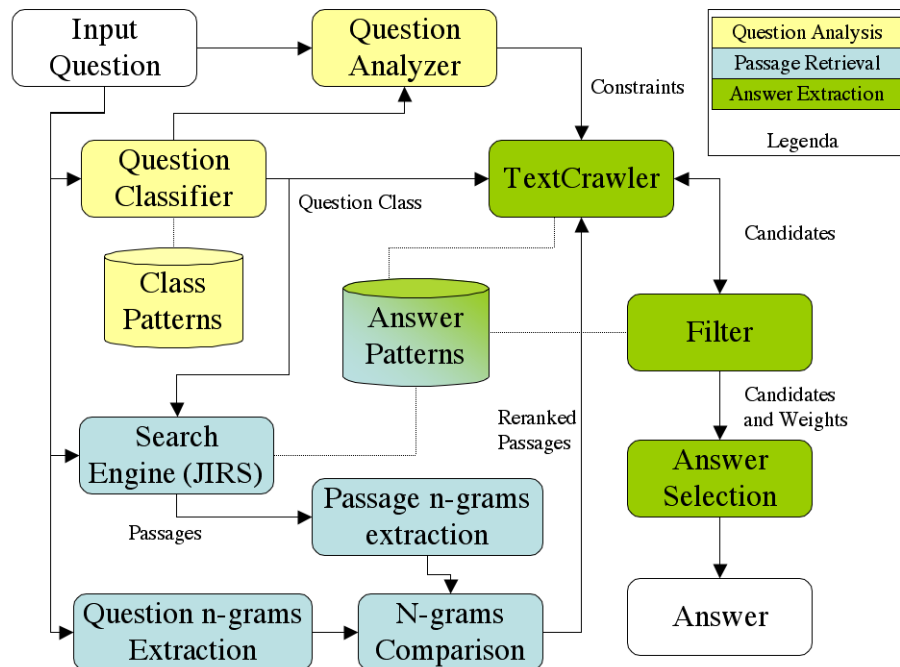


Fig. 1. Diagram of the QA system

Given a user question, this will be handed over to the *Question Analysis* module, which is composed by a *Question Analyzer* that extracts some constraints to be used in the answer extraction phase, and by a *Question Classifier* that determines the class of the input question. At the same time, the question is passed to the *Passage Retrieval* module, which generates the passages used by the *Answer Extraction* (AE) module together with the information collected in the question analysis phase in order to extract the final answer.

2.1 Question Analysis Module

This module obtains both the expected answer type (or *class*) and some constraints from the question. Question classification is a crucial step of the processing since the Answer Extraction module uses a different strategy depending on the expected answer type; as reported by Moldovan et al. [4], errors in this phase account for the 36.4% of the total number of errors in Question Answering.

The different answer types that can be treated by our system are shown in Table 1. We introduced the “FIRSTNAME” subcategory for “NAME” type questions, because we defined a pattern for this kind of questions in the AE module. We also specialized the “DEFINITION” class into three subcategories: “PERSON”, “ORGANIZATION” and “OBJECT”, which was introduced this year (e.g.: *What is a router?*). With respect to CLEF 2005, the Question Classifier does not use a SVM classifier.

Table 1. QC pattern classification categories.

L0	L1	L2
NAME	ACRONYM PERSON TITLE FIRSTNAME LOCATION	COUNTRY CITY GEOGRAPHICAL
DEFINITION	PERSON ORGANIZATION OBJECT	
DATE	DAY MONTH YEAR WEEKDAY	
QUANTITY	MONEY DIMENSION AGE	

Each category is defined by one or more patterns written as regular expressions. The questions that do not match any defined pattern are labeled with *OTHER*. If a question matches more than one pattern, it is assigned the label of the longest matching pattern (i.e., we consider longest patterns to be less generic than shorter ones).

The Question Analyzer has the purpose of identifying the constraints to be used in the AE phase. These constraints are made by sequences of words extracted from the POS-tagged query by means of POS patterns and rules. For instance, any sequence of nouns (such as *ozone hole*) is considered as a relevant

pattern. The POS-taggers used were the SVMtool³ for English and Spanish, and the TreeTagger⁴ for Italian and French.

There are two classes of constraints: a *target* constraint, which is the word of the question that should appear closest to the answer string in a passage, and zero or more *contextual* constraints, keeping the information that has to be included in the retrieved passage in order to have a chance of success in extracting the correct answer. For example, in the following question: “*Dónde se celebraron los Juegos Olímpicos de Invierno de 1994?*” (*Where did the Winter Olympic games of 1994 take place?*) *celebraron* is the target constraint, while *Juegos Olímpicos de Invierno* and *1994* are the contextual constraints. There is always only one target constraint for each question, but the number of contextual constraint is not fixed. For instance, in “*Quién es Neil Armstrong?*” the target constraint is *Neil Armstrong* but there are no contextual constraints.

2.2 The JIRS Passage Retrieval Module

The passages containing the relevant terms are retrieved by JIRS using a classical keyword-based IR system. This year, the module was modified in order to rank better the passages which contain an answer pattern matching the question type. Therefore, this module is not as language-independent as in 2005 because it uses informations from the Question Classifier and the patterns used in the Answer Extraction phase.

Sets of unigrams, bigrams, ..., n -grams are extracted from the passages and from the user question. In both cases, n is the number of question terms. These n -gram sets are compared in order to obtain the weight of each passage, which is proportional to the size of the question n -grams found in the passage.

For instance, if the question is “*What is the capital of Croatia?*” and the system retrieves the following two passages: “*...Tudjman, the president of Croatia, met Eltsin during his visit to Moscow, the capital of Russia...*”, and “*...they discussed the situation in Zagreb, the capital of Croatia...*”. The second passage must have more importance because it contains the 4-gram “*the capital of Croatia*”, whereas the first one contains the 3-gram “*the capital of*” and the 1-gram “*Croatia*”. This example also shows the advantage of considering n -grams instead of keywords: the two passages contains the same question keywords, but only one of them contains the right answer.

In order to calculate the weight of n -grams of every passage, the greatest n -gram in the passage is identified and it is assigned a weight equal to the sum of all its term weights. Subsequently, smaller n -grams are searched. The weight of every term is determined by means of formula (1):

$$w_k = 1 - \frac{\log(n_k)}{1 + \log(N)} . \quad (1)$$

³ <http://www.lsi.upc.edu/nlp/SVMTool/>

⁴ <http://www.ims.uni-stuttgart.de/projekte/complex/TreeTagger/DecisionTreeTagger.html>

Where n_k is the number of passages in which the term appears, and N is the number of passages. We make the assumption that stopwords occur in every passage (i.e., $n_k = N$ for stopwords). Therefore, if the term appears once in the passage collection, its weight will be equal to 1 (the greatest weight).

Sometimes a term unrelated to the question can obtain a greater weight than those assigned to the Named Entities (NEs), such as names of persons, organizations and places, dates. The NEs are the most important terms of the question and it does not make sense to return passages which do not contain them. Therefore, NEs are given a greater weight than the other question terms, in order to force its presence in the first ranked passages. NEs are recognized by simple rules, such as capitalization, or by checking if they are a number. Once all the terms have been weighted, the sum is normalized.

JIRS can be obtained at the following URL: <http://jirs.dsic.upv.es>.

2.3 Answer Extraction

The input of this module is constituted by the n passages returned by the PR module and the constraints (including the expected type of the answer) obtained through the *Question Analysis* module. A *TextCrawler* is instantiated for each of the n passages with a set of patterns for the expected type of the answer and a pre-processed version of the passage text. For CLEF 2006, we corrected some errors in the patterns and we also introduced new ones.

Some patterns can be used for all languages; for instance, when looking for proper names, the pattern is the same for all languages. The pre-processing of passage text consists in separating all the punctuation characters from the words and in stripping off the annotations of the passage. It is important to keep the punctuation symbols because we observed that they usually offer important clues for the individuation of the answer (this is true especially for *definition* questions): for instance, it is more frequent to observe a passage containing “*The president of Italy, Giorgio Napolitano*” than one containing “*The president of Italy IS Giorgio Napolitano*” ; moreover, movie and book titles are often put between apices.

The positions of the passages in which occur the constraints are marked before passing them to the TextCrawlers. A difference with 2005 is that now we do not use the Levenshtein-based spell-checker to compare strings in this phase now.

The TextCrawler begins its work by searching all the passage’s substrings matching the expected answer pattern. Then a weight is assigned to each found substring s , depending on the positions of s with respect to the constraints, if s does not include any of the constraint words. If in the passage are present both the target constraint and one or more of the contextual constraints, then the product of the weights obtained for every constraint is used; otherwise, it is used only the weight obtained for the constraints found in the passage.

The *Filter* module takes advantage of a mini knowledge base in order to discard the candidate answers which do not match with an allowed pattern or that do match with a forbidden pattern. For instance, a list of country names

in the four languages has been included in the knowledge base in order to filter country names when looking for countries. When the Filter module rejects a candidate, the TextCrawler provide it with the next best-weighted candidate, if there is one.

Finally, when all TextCrawlers end their analysis of the text, the *Answer Selection* module selects the answer to be returned by the system. The following strategies apply:

- Simple voting (SV): The returned answer corresponds to the candidate that occurs most frequently as passage candidate.
- Weighted voting (WV): Each vote is multiplied for the weight assigned to the candidate by the TextCrawler and for the passage weight as returned by the PR module.
- Maximum weight (MW): The candidate with the highest weight and occurring in the best ranked passage is returned.
- Double voting (DV): As simple voting, but taking into account the second best candidates of each passage.
- Top (TOP): The candidate elected by the best weighted passage is returned.

We used the Confidence Weighted Score (CWS) to select the answer to be returned to the system, relying on the fact that in 2005 our system was the one returning the best values for CWS [7]. For each candidate answer we calculated the CWS by dividing the number of strategies giving the same answer by the total number of strategies (5), multiplied for other measures depending on the number of returned passages (n_p/N , where N is the maximum number of passages that can be returned by the PR module and n_p is the number of passages actually returned) and the averaged passage weight. The final answer returned by the system is the one with the best CWS. Our system always return only one answer (or NIL), although 2006 rules allowed to return more answers per question. The weighting of NIL answers is slightly different, since it is obtained as $1 - n_p/N$ if $n_p > 0$, 0 elsewhere.

The snippet for answer justification is obtained from the portion of text surrounding the first occurrence of the answer string. The snippet size is always 300 characters (150 before and 150 after the answer) + the number of characters of the answer string.

3 Experiments and Results

We submitted two runs for each of the following monolingual task: Spanish, Italian and French. The first runs (labelled *upv.061*) use the system with JIRS as PR engine, whereas for the other runs we used Lucene, adapted to the QA task with the implementation of a weighting scheme that privileges long passages and is similar to the word-overlap scheme of the MITRE system [2]. In Table 2 we show the overall accuracy obtained in all the runs.

With respect to 2005, the overall accuracy increased by $\sim 3\%$ in Spanish and Italian, and by $\sim 7\%$ in French. We suppose that the improvement in French

Table 2. Accuracy results for the submitted runs. Overall: overall accuracy, factoid: accuracy over factoid questions; definition: accuracy over definition questions; nil: precision over nil questions (correctly answered nil/times returned nil); CWS: confidence-weighted score.

task	run	overall	factoid	definition	nil	CWS
es-es	upv_061	36.84%	34.25%	47.62%	0.33	0.225
	upv_062	30.00%	27.40%	40.48%	0.32	0.148
it-it	upv_061	28.19%	28.47%	26.83%	0.23	0.123
	upv_062	28.19%	27.78%	29.27%	0.23	0.132
fr-fr	upv_061	31.58%	31.08%	33.33%	0.36	0.163
	upv_062	24.74%	26.35%	19.05%	0.18	0.108

is due to the fact that the target collection was larger this year. Spanish is still the language in which we obtain the best results, even if we are not sure about the reason: a possibility is that this can be due to the better quality of the POS-tagger used in the analysis phase for the Spanish language.

We obtained an improvement over the 2005 system in factoid questions, but also worse results in definition ones, probably because of the introduction of the *object* definitions.

The JIRS-based systems performed better than the Lucene-based ones in Spanish and French, whereas in Italian they obtained almost the same results. The difference in the CWS values obtained in both Spanish and French is consistent and weights in favour of JIRS. This prove that the quality of passages returned by JIRS for these two languages is considerably better.

We measured also the inter-agreement of the two systems, by counting the number of times that the two systems returned the same source document divided by the number of times that they returned the same answer (we called this measure *Agreement on Answer* or *AoA*), and when the answer was the right one (in this case we call it *AoRA*).

Table 3. Inter-agreement between the two systems, calculated by means of the AoA and AoRA measures.

task	AoA	AoRA	Collection size
es-es	61.33%	42.52%	1086MB
it-it	71.92%	46.68%	170MB
fr-fr	53.47%	29.81%	487MB

As it can be observed in Table 3, the best agreement has been obtained in Italian, as one would expect due to the smaller size of the collection.

4 Conclusions and Further Work

We obtained a slight improvement over the results of our 2005 system. This is consistent with the small amount of modifications introduced, principally because of the new rules defined for the 2006 CLEF QA task. The most interesting result is that JIRS demonstrated to be more effective for the QA task than a standard keyword-based search engine such as Lucene in two languages over three. Our further works on the QUASAR system will concern the implementation of a specialized strategy for definition questions, and probably a major revision of the Answer Extraction module.

Acknowledgments

We would like to thank the TIN2006-15265-C06-04 research project for partially supporting this work. This paper is a revised version of the work titled “The UPV at QA@CLEF 2006” included in the CLEF 2006 Working Notes.

References

1. Lili Aunimo, Reeta Kuuskoski, and Juha Makkonen. Cross-language question answering at the university of helsinki. In *Workshop of the Cross-Lingual Evaluation Forum (CLEF 2004)*, Bath, UK, 2004.
2. Marc Light, Gideon S. Mann, Ellen Riloff, and Eric Breck. Analyses for elucidating current question answering technology. *Nat. Lang. Eng.*, 7(4):325–342, 2001.
3. Bernardo Magnini, Matteo Negri, Roberto Prevete, and Hristo Tanev. Multilingual question/answering: the DIOGENE system. In *The 10th Text REtrieval Conference*, 2001.
4. Dan Moldovan, Marius Pasca, Sanda Harabagiu, and Mihai Surdeanu. Performance issues and error analysis in an open-domain question answering system. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, New York, USA, 2003.
5. Gunther Neumann and Bogdan Sacaleanu. Experiments on robust nl question interpretation and multi-layered document annotation for a cross-language question/answering system. In *Workshop of the Cross-Lingual Evaluation Forum (CLEF 2004)*, Bath, UK, 2004.
6. José Manuel Gómez Soriano, Davide Buscaldi, Empar Bisbal, Paolo Rosso, and Emilio Sanchis. Quasar: The question answering system of the universidad politécnica de valencia. In *Lecture Notes in Computer Science (LNCS)*, volume 4022, Vienna, Austria, 2006. Springer Verlag.
7. Alessandro Vallin, Danilo Giampiccolo, Lili Aunimo, Christelle Ayache, Petya Osenova, Anselmo Peas, Maarten de Rijke, Bogdan Sacaleanu, Diana Santos, and Richard Sutcliffe. Overview of the clef 2005 multilingual question answering track. In *CLEF 2005 Proceedings*, 2005.
8. José L. Vicedo, Ruben Izquierdo, Fernando Llopis, and Rafael Munoz. Question answering in spanish. In *Workshop of the Cross-Lingual Evaluation Forum (CLEF 2003)*, Trondheim, Norway, 2003.