

Cross-Lingual Question Answering by Answer Translation

Johan Bos

Linguistic Computing Laboratory
Department of Computer Science
University of Rome “La Sapienza”
bos@di.uniroma1.it

Malvina Nissim

Laboratory for Applied Ontology
Inst. for Cognitive Science & Technology
National Research Council, Rome
malvina.nissim@loa-cnr.it

Abstract

We approach cross-lingual question answering by using a mono-lingual QA system for the source language and by translating resulting answers into the target language. As far as we are aware, this is the first cross-lingual QA system in the history of CLEF that uses this method—all other cross-lingual QA systems known to us use translation of the question or query instead. We demonstrate the feasibility of this approach by using a mono-lingual QA system for English, and translating answers and finding appropriate documents in Italian and Dutch. For factoid and definition questions, we achieve overall accuracy scores ranging from 13% (EN→NL) to 17% (EN→IT) and lenient accuracy figures from 19% (EN→NL) to 25% (EN→IT). The advantage of this strategy to cross-lingual QA is that translation of answers is easier than translating questions—the disadvantage is that answers might be missing from the source corpus and additional effort is required for finding supporting documents of the target language.

Categories and Subject Descriptors

H.3 [Information Storage and Retrieval]: H.3.1 Content Analysis and Indexing; H.3.3 Information Search and Retrieval; H.3.4 Systems and Software; H.5 [Information Interfaces and Presentation]: H.5.2 User Interfaces; I.2 [Artificial Intelligence]: I.2.1 Applications and Expert Systems; I.2.4 Knowledge Representation Formalisms and Methods; I.2.7 Natural Language Processing

General Terms

Algorithms, Design, Measurement, Performance, Experimentation

Keywords

question answering, machine translation, natural language processing

1 Introduction

Automated Question Answering (QA) is the task of providing an exact answer (instead of a document) to a question formulated in natural language. Cross-lingual QA is concerned with providing an answer in one language (the target language) to a question posed in a different language (the source language). Most systems tackle the cross-lingual problem by translating the question or query posed in the source language in the target language, and then use a QA system developed for the target language for retrieving an answer. For instance, at the last three QA@CLEF campaigns (2003–2005) there were 28 participants performing in 20 different cross-lingual tasks¹, each of them using off-the-shelf translation software (such as Babelfish, Systran, Reverso, FreeTrans, WorldLingo, Transtool) to translate the question or the keywords of the query of the source language into the target language, followed by processing the translated question using a QA system designed for the target language.

Little attention has been given to the obvious alternative approach: translating the answer, instead of the question. There are some potential advantages following this route: answers are easier to translate, due to their simpler syntactic structure. In fact, some types of answers (such as date expressions or names of persons) hardly need a translation. In addition, finding a document in the target language supporting the answer (as prescribed in the QA@CLEF exercise) is feasible: using either the translated question or keywords thereof and the translated answer, standard document retrieving tools can be used to find a document. This approach can work provided that the source and target language documents cover the same material, otherwise all bets are off.

In the context of CLEF, we tested this approach relying on the fact that the documents in the collection are from the same time period. We ran our experiments for two language pairs: EN→NL and EN→IT, using a mono-lingual QA system for English and an off-the-shelf machine translation tool for translating the answers from source into target language. In this report we briefly describe our system, show system output, evaluate the approach, and comment on the feasibility of this approach.

2 System Description

Our cross-lingual QA system that we used for our experiments at CLEF 2006 is an extension of a mono-lingual QA system for English. It deals with factoid (including list questions) and definition questions, but it uses two different streams of processing for these two types of questions. The first component in the pipeline, Question Analysis, deals with all types of questions. Then, when the question turns out to be of type factoid, Document Analysis, Answer Extraction and Answer Translation will follow. Answers to definition questions are directly searched in the target language corpora (see Figure 1 and Section 2.5). What follows is a more detailed description of each component. Examples of the different data structures of the system are shown in Figure 2 and 3.

2.1 Question Analysis

The (English) question is tokenised and parsed with a wide-coverage parser based on Combinatory Categorical Grammar (CCG). We use the parser of Clark & Curran [3]. On the basis of the output of the parser, a CCG-derivation, we build a semantic representation in the form of a Discourse Representation Structure (DRS), closely following Discourse Representation Theory [4]. This is done using the semantic construction method described in [1, 2]. The Question-DRS is the basis for generating four other sources of information required later in the question answering process: an expected answer type; a query for document retrieval; the answer cardinality; and background knowledge for finding appropriate answers.

¹To be more precise, these were BU→EN (3x), BU→FR, DE→EN (4x), DE→FR, EN→DE (3x), EN→ES (2x), EN→FR, EN→NL, EN→PT, ES→EN (2x), ES→FR, FI→EN (2x), FR→EN (8x), IT→EN (3x), IT→ES, IT→FR (2x), NL→EN, NL→FR, PT→FR (2x), and IN→EN (Information gathered from the working notes of CLEF 2003 [5], CLEF 2004 [7] and CLEF 2005 [6]).

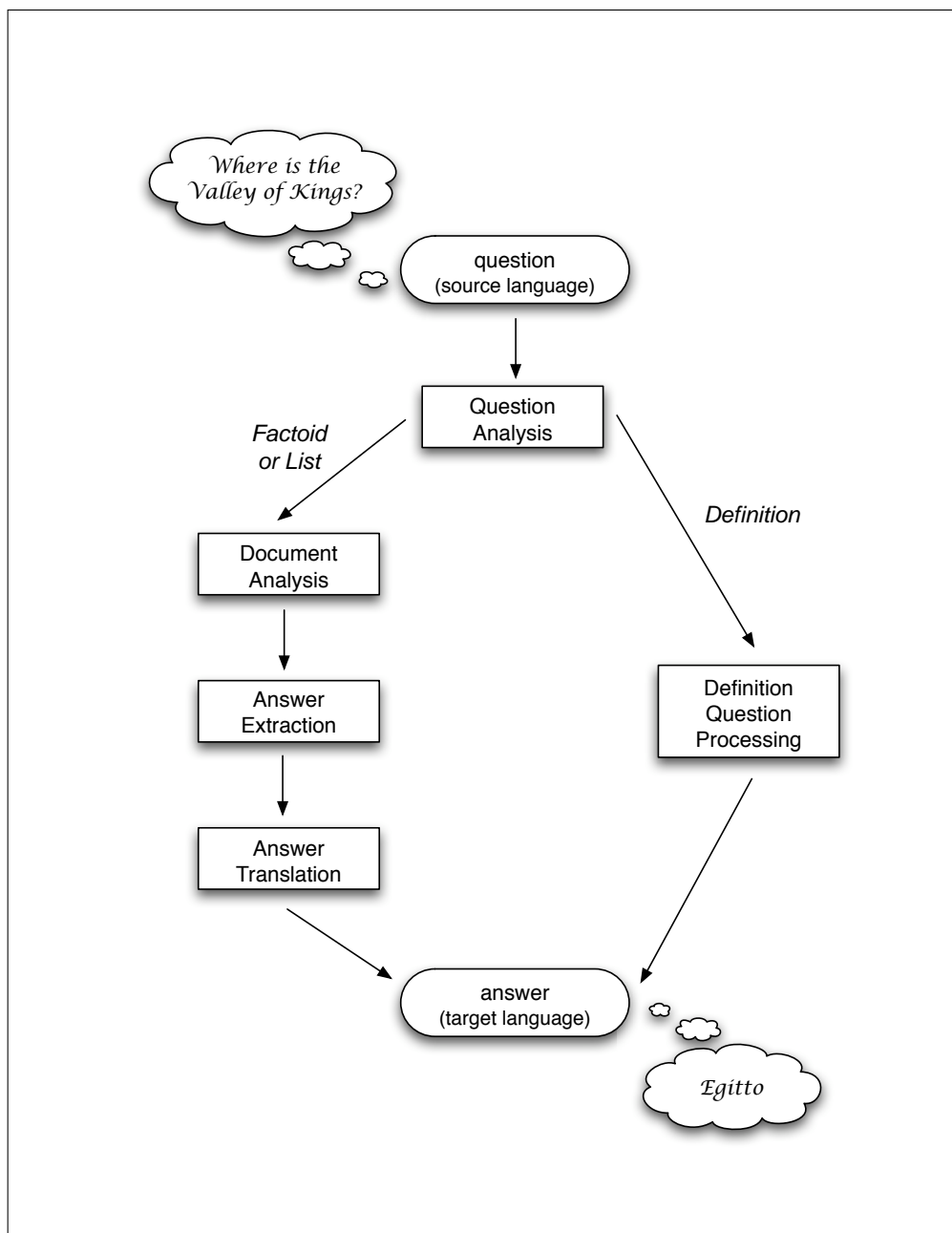


Figure 1: Simplified architecture of the QA system.

We use 14 main expected answer types which are further divided into subtypes. The main types are definition, description, attribute, numeric, measure, time, location, address, name, language, creation, instance, kind, and part. The answer cardinality denotes a range expressed by an ordered pair of two numbers, the first indicating the minimal number of answers expected, the second the maximal number of answers (or 0 if unspecified). For instance, 3–3 indicates that exactly three answers are expected, 2–0 means at least two answers. Background knowledge is a list of axioms related to the question—it is information gathered from WordNet or other lexical resources.

2.2 Document Analysis

In order to maximise the chance of finding an answer in the source language (English), we extended the English CLEF document collection with documents from the Acquaint corpus. All documents were pre-processed: sentence splitting, tokenisation, and dividing into smaller documents of two sentences each (Taking a sliding window, so each sentence will appear in two mini-documents). These mini-documents are indexed with the Indri information retrieval tools (we used version 2.2, see <http://www.lemurproject.org/indri/>). The query generated by Question Analysis (see Section 2.1) is used to retrieve the best 1,000 mini-documents, again with the use of Indri.

Using the same wide-coverage parser as for parsing the question, all retrieved documents are parsed and for each of them a Discourse Representation Structure (DRS) is generated. The parser also performs basic named entity recognition for locations, persons, and organisations. This information is used to assign the right semantic type to discourse referents in the DRS. Date expressions are normalised in the DRS.

2.3 Answer Extraction

Given the DRS of the question (the Q-DRS), and a set of DRSs of the retrieved documents (the A-DRSs), we match each A-DRS with the Q-DRS to find a potential answer. This process proceeds as follows: if the A-DRS contains a discourse referent of the expected answer type (see 2.1) matching will commence attempting to identify the semantic structure in the Q-DRS with that of the A-DRS. The result is a score between 0 and 1 indicating the amount of semantic material that could be matched. The background knowledge (such as hyponyms from WordNet) generated by the Question Analysis (see 2.1) is used to assist in the matching. All retrieved answers are re-ranked on the basis of the match-score and frequency.

2.4 Answer Translation

The answer obtained for the source language is now translated into the target language by means of the Babelfish off-the-shelf machine translation tool. However, this was not done for all types of answers: we refrained from translating answers that were person names or titles of creative works. Person names are normally not translated across the languages that we were working with. Names of creative works often are, but the machine translation software that we used did not perform well enough on some examples in our training data, so that we decided to leave these untranslated. In addition, titles of creative works often don't get a literal translation (a case in point is Woody Allen's "Annie Hall", which is translated in Italian as "Io e Annie") so a more sophisticated translation strategy is required.

Given the answer in the target language, we need to find a supporting document from the target language collection (as prescribed by the QA@TRE exercise). Hence, we also translate the original question, independent of the answer found for it. This we use to construct another query and then retrieve a document from the target language collection that contains the translated answer and as many as possible terms from the translated question. (As with the source language documents, these are indexed and retrieved using Indri.)

2.5 Definition Question Processing

We did not put much effort in dealing with definition questions. We basically adopted a simple pattern matching technique directly on the target language documents. Here is how it works.

Once a question is identified as a definition question (see Figure 1), all non-content words are removed from the question (wh-words, the copula, articles, punctuation, etc.). What is left is the topic of the definition question. For instance,

English question: “Who is Radovan Karadzic?”
Topic: “Radovan Karadzic”.

Given the topic, we search interesting information about the topic, stated in the target language. We do this using an off-line dump of the Dutch and Italian Wikipedia pages by selecting sentences containing the target. From this we generate an Indri query selecting all (one-sentence) documents containing the topic, and a combination of the interesting terms found in Wikipedia (stop words are removed). This gives us a list of documents of the target language.

The last step is some simple template matching using regular expressions to extract the relevant clauses of a sentence that could count as an answer to the question. There are only a few templates, but we use different ones for the different sub-types of definition questions (the Question Analysis component distinguishes between three sub-types of definition questions: person, organisation, and thing). Here are the patterns we used:

```
person:      /(^| )(\,|de|i|l|la|i|gli|l\') (.+) $topic /i
person:      /$topic \(\ (\d+) \)/i
organisation: /$topic \(([\^\\])+\)\)/i
thing:       /$topic , ([\^,]+) , /
thing:       /$topic \(([\^\\])+\)\)/i
```

As an example, consider the following sentences and answers in bold-face:

LASTAMPA94-041510 Il **leader** Radovan Karadzic non era reperibile.
LASTAMPA94-042168 Intanto il **leader serbo-bosniaco** Radovan Karadzic e il comandante
in capo delle forze serbe, gen.

3 Evaluation

3.1 Basic Performance

The question analysis performed fairly well. Only 2 of 200 questions of the EN→IT set, and 7 of EN→NL set could not be parsed. For 189 questions of the EN→IT set, and 177 of the EN→NL set the system determined an expected answer type. This already shows that our system had more difficulties with the EN→NL set, which is reflected in the overall scores.

Relatively many questions did not have an answer in the English collection, or at least our system failed to find one. For 41 of the 200 EN→NL questions and for 43 of the 200 EN→IT questions no English answer (correct or incorrect) was found. Answer translation introduced only few errors (see Section 3.2), but finding a supporting document proved harder than we had hoped. So several correct answers were associated with wrong documents.

We submitted four runs—two for the EN→NL task, and two for the EN→IT task. The first runs contained one answer for factoid questions and up to ten for definition and list questions. The second runs contained up to ten answers for each type of question. We used this strategy because it was unclear, at the time of submission, what kind of evaluation would be used. So the number one runs would perform better on an accuracy score, and the number two runs better on scores based on the average mean reciprocal rank (MRR).

It turned out that both accuracy and MRR were used in the evaluation. The results for definition and factoid questions are shown in Table 1, and the figures for the list questions in Table 2. As the

Input Question:	[0040] What year did Titanic sink?
Q-DRS:	<p>The Q-DRS diagram consists of a large outer box. Inside, there is a smaller box on the left containing 'x1' and 'year(x1)'. To its right is a question mark '?'. Further right is another box containing 'x2 x3' and three relations: 'named(x3,titanic,nam)', 'sink(x2) agent(x2,x3)', and 'rel(x2,x1)'.</p>
Answer Type:	[tim:yea]
Cardinality:	1-1
English Query:	#filreq(Titanic #weight(1 Titanic 3 sink))
English Answer:	1912
English Context:	[APW19981105.0654] Speaking of bad guys: Cartoon News reprints a bunch of political cartoons published after the <u>Titanic sank in 1912</u> .
A-DRS:	<p>The A-DRS diagram is a complex nested structure. It starts with a large box containing 'x0 x1 x2 x3' and relations 'speak(x1) agent(x1,x0)', 'bad(x2) guy(x2)', and 'proposition(x3) :(x2,x3)'. Inside this is a box for 'x3:' containing 'x4 x5 x6 x6 x7 x8 x9' and relations 'named(x4,cartoon_news,org)', 'bunch(x5) political(x6)', 'publish(x7) patient(x7,x6)', and 'proposition(x8) after(x7,x8)'. Inside that is a box for 'x8:' containing 'x10 x11 x12' and relations 'named(x10,titanic,nam)', 'timex(x12)=1912-XX-XX', 'sink(x11) agent(x11,x10) in(x11,x12)', 'cartoon(x6) of(x5,x6)', and 'reprint(x9) agent(x9,x4) patient(x9,x5)'. The entire structure is anchored to 'of(x1,x2)' at the bottom.</p>
Italian Answer:	1912
Italian Query:	#filreq(1912 #combine(1912 anno Titanic affondato))
Italian Context:	[AGZ.941120.0028] I cantieri di Belfast hanno confermato di aver ricevuto la commessa più prestigiosa dopo la costruzione del "Titanic": la messa in mare del "Titanic 2". Una copia perfetta del fantasmagorico transatlantico affondato nell'oceano nel 1912 durante il viaggio inaugurale è stata ordinata da un consorzio giapponese presso la ditta "Harland and Wolff" e sarà completata entro il 1999.

Figure 2: System input and output for a factoid question in the EN→IT task.

Input Question:	[0082] How many wars have been fought between India and Pakistan for the possession of Kashmir?
Q-DRS:	<div style="border: 1px solid black; padding: 10px; margin: 10px auto; width: fit-content;"> <div style="border: 1px solid black; padding: 5px; margin-bottom: 5px; width: fit-content;"> x1 x2 card(x1)=x2 quantity(x2) war(x1) </div> <div style="display: inline-block; vertical-align: middle; margin: 0 10px;">?</div> <div style="border: 1px solid black; padding: 5px; margin-bottom: 5px; width: fit-content;"> x3 x9 x7 x8 x6 x4 x5 fight(x3) patient(x3,x1) named(x9,india,loc) between(x3,x9) between(x3,x6) possession(x7) of(x7,x8) named(x8,kashmir,loc) named(x6,pakistan,loc) possession(x4) of(x4,x5) named(x5,kashmir,loc) for(x3,x7) for(x3,x4) </div> </div>
Answer Type:	[num:cou]
Cardinality:	1-1
English Query:	#filreq(Pakistan #weight(1 war 6 Kashmir 4 possession 6 India 3 fight))
English Answer:	three
English Context:	[NYT19980615.0189] Since then, the countries have fought three wars, two of them over <u>Kashmir</u> , which <u>Pakistan</u> claims was illegally grafted onto <u>India</u> at independence.
A-DRS:	<div style="border: 1px solid black; padding: 10px; margin: 10px auto; width: fit-content;"> x0 x1 x1 x2 x3 x4 x5 x6 x7 country(x0) card(x1)=3 war(x1) card(x1)=2 thing(x2) of(x1,x2) named(x3,kashmir,loc) over(x1,x3) proposition(x4) which(x1,x4) <div style="border: 1px solid black; padding: 5px; margin: 5px auto; width: fit-content;"> x8 x9 x10 x11 x12 claim(x9) nn(x8,x9) named(x8,pakistan,loc) graft(x10) patient(x10,x9) named(x11,india,loc) onto(x10,x11) independence(x12) at(x10,x12) illegally(x10) fight(x6) agent(x6,x0) patient(x6,x1) then(x7) since(x6,x7) </div> </div>
Dutch Answer:	drie
Dutch Query:	#filreq(drie #combine(oorlogen bestreden India Pakistan bezit Kashmir))
Dutch Context:	[NH19940816-0031] Tijdens de eerste van drie oorlogen tussen India en Pakistan, veroverden de Pakistanen in 1947 ongeveer eenderde deel van het voormalige grondgebied van Kashmir.

Figure 3: System input and output for a factoid question in the EN→NL task.

tables illustrate, the scores for accuracy are the same for each run, which means that providing more than one answer for a factoid question was not punished. As expected, we achieved better scores for MRR on the number two runs. It is surprising that we did so well on definition questions, given that we paid very little effort in dealing with them.

Table 1: Number of right (R), wrong (W), inexact (X), unsupported (U) answers, accuracy measure of first answers (factoid, definition, overall, and lenient), and mean reciprocal rank score (MRR) for all four submitted CLEF 2006 runs.

Run	R	W	X	U	Accuracy				MRR
					Fact. Acc.	Def. Acc.	Overall	Lenient	
EN→IT 1	32	141	4	11	15.3%	24.4%	17.0%	25.0%	0.18
EN→IT 2	32	141	4	11	15.3%	24.4%	17.0%	25.0%	0.20
EN→NL 1	25	150	6	3	11.6%	20.5%	13.4%	18.5%	0.14
EN→NL 2	25	149	7	3	11.6%	20.5%	13.4%	19.0%	0.15

Table 2: Number of right (R), wrong (W), inexact (X), unsupported (U), unassessed (I) answers, P@N score for list questions, for all submitted CLEF 2006 runs.

Run	R	W	X	U	I	P@N
EN→IT 1	12	63	1	6	0	0.10
EN→IT 2	15	70	1	7	0	0.15
EN→NL 1	6	29	0	7	24	0.18
EN→NL 2	9	40	0	5	42	0.15

3.2 Answer Translation Accuracy

Recall that depending on the kind of expected answer type, answers were translated into the target language or not. We did not translate names of persons, nor titles of creative works. Whereas this strategy worked out well for person names, it didn't for names of creative works. For instance, for question "0061 What book did Salman Rushdie write?", we found a correct answer "Satanic Verses" in the English documents, but this was not found in the Italian document collection, because the Italian translation is "Versetti Satanici". Babelfish wouldn't have helped us here either, as it translates the title into "Verses Satanic".

Table 3: Answer translation accuracy, divided over answer types (EN→NL). Quantified over the first translated answer, disregarding whether the answer was correct or not, but only for questions to which a correct expected answer type was assigned.

Expected Answer Type	Correct	Wrong	Accuracy
location	30	0	100%
numeric	9	2	82%
measure	4	2	67%
instance (names)	14	5	74%
time	15	0	100%
other	7	1	88%

Overall, answer translation introduced little errors as many answers are easy to translate (Table 3). Locations are usually correctly translated by Babelfish, and so are time expressions. Difficulties sometimes arise for numeric expressions. For instance, for the question “0084 How many times did Jackie Stewart become world champion?” we found the correct English answer “three-time”. However, this was wrongly translated in “drie-tijd”, and obviously not found in the Dutch corpus. Similarly, many answers for measurements are expressed in the English corpus using the imperial system, whereas the metric system is used in the Dutch and Italian corpora. Finally, some things are just hard to translate. Although we found a reasonably correct answer for “0136 What music does Offspring play?”, namely “rock”, this was translated in Dutch as “rots”. In itself a correct translation, but for the wrong context.

4 Conclusion

What can we say about the answer-translation strategy to cross-lingual question answering, and how can we improve it? Generally speaking, we believe it is a promising approach, because answers are easier to translate than questions—in fact, many answers don’t need to be translated, and some are relatively easy to translate, such as date expressions, locations, and numerals. The remaining translation difficulties, such as translating measure terms from the imperial to the metric system, and titles of creative works, can be dealt with by employing designated translation tools such as pre-compiled lookup-tables. It is also an advantage that the type of the answer is known—this information can serve in obtaining a better translation.

One problem that came to the surface was the omission of the answer in the source language document collection. There is only one way to deal with this situation, and that is getting a larger pool of documents. One option is to use the web for finding the answer in the source language.

Another problem that arised was finding a supporting document for the target language. The system can certainly be improved with respect to this point: it currently only takes the translated question as additional information to find a target language document. This is interesting, as the machine translated question need not be grammatically perfect to find a correct document. One way to improve this is to translate the context found for the source language as well, and use it in addition to retrieve a target language document.

References

- [1] J. Bos, S. Clark, M. Steedman, J.R. Curran, and Hockenmaier J. Wide-Coverage Semantic Representations from a CCG Parser. In *Proceedings of the 20th International Conference on Computational Linguistics (COLING '04)*, Geneva, Switzerland, 2004.
- [2] Johan Bos. Towards wide-coverage semantic interpretation. In *Proceedings of Sixth International Workshop on Computational Semantics IWCS-6*, pages 42–53, 2005.
- [3] S. Clark and J.R. Curran. Parsing the WSJ using CCG and Log-Linear Models. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL '04)*, Barcelona, Spain, 2004.
- [4] H. Kamp and U. Reyle. *From Discourse to Logic; An Introduction to Modeltheoretic Semantics of Natural Language, Formal Logic and DRT*. Kluwer, Dordrecht, 1993.
- [5] C. Peters, editor. *Results of the CLEF 2003 Cross-Language System Evaluation Campaign. Working Notes for the CLEF 2003 Workshop*, Trondheim, Norway, 21–22 August 2003.
- [6] C. Peters, editor. *Results of the CLEF 2005 Cross-Language System Evaluation Campaign. Working Notes for the CLEF 2005 Workshop*, Vienna, Austria, 21–23 September 2005.
- [7] C. Peters and F. Borri, editors. *Results of the CLEF 2004 Cross-Language System Evaluation Campaign. Working Notes for the CLEF 2004 Workshop*, Bath, UK, 15–17 September 2004.