**Volume 16, No. 2
April 2012**

Juncal Gutierrez-Artacho

**Professor María-Dolores Olvera-Lobo** holds a Ph.D. in Documentation and is professor in the Department of Library and Information Science at the University of Granada, and teaches as well in the School of Translation and Interpretation. Her scholarly books and articles have concentrated on Retrieval Information, Information Science and Translation Teaching and Learning. She is a member of the associated unit of the SCImago research group, under Spain´s Higher Council for Scientific Research (CSIC). She has been a visiting professor at many institutions all over the world.

Professor Olvera-Lobo can be reached at molvera@ugr.es.

Professor **Juncal Gutierrez-Artacho** graduated in 2008 from the University of Granada (UGR) with a degree in Translation and Interpreting. Currently she is undertaking her Doctoral Studies in the University of Granada and she is teaching Translation at the Pablo Olavide University (Seville, Spain). She is doing her dissertation about the analysis and evaluation of mono-lingual and multi-lingual Question-answering systems. Her scholarly articles have concentrated on Retrieval Information and Translation Teaching and Learning. Presently, she is a member of some research groups and some teaching innovation projects. She is also a professional translator.

Professor Gutierrez-artacho can be reached at juncalgutierrez@ugr.es.

## Front Page

Select one of the previous 59 issues.

Select an issue: ▼

🎾 **Index 1997-2012**

🎾 **TJ Interactive:**
**Translation Journal Blog**

**Translator Profiles**

🎾 Planning and Passion
by Helen Eby

**The Profession**

🎾 The Bottom Line

# Language Resources for Translation in Multilingual Question Answering Systems

by María-Dolores Olvera-Lobo and Juncal Gutiérrez-Artacho
CSIC, Unidad Asociada Grupo SCImago, Madrid, Spain and Departament of Library and Information Science, Universidad de Granada, Spain
Departament of Translation and Interpreting, Universidad de Granada, Spain and Departament of Philology and Translation, Universidad Pablo de Olavide, Seville, Spain

### Abstract

In the field of Information Retrieval monolingual and multilingual tools are being created that can greatly assist specialists in their work; as well as helping other users find a wide variety of information. Multilingual tools are evolving but several years of study and research are still needed to improve implementations. One of the main difficulties facing these tools is the task of translating queries made by users and the documentary sources found in response (Diekema, 2003). Given the current expansion in research, development, and the creation of multilingual IR systems, it was considered worthwhile analysing and evaluating the resources used by one type of these systems, the multilingual Question Answering systems.

Our article is primarily intended as a general purpose analysis and aims to encompass translation in the study of multilingual QA systems. The second general aim is to identify and analyse the linguistic resources and tools found in these systems. Specific objectives include identifying the main types of language resources and tools useful in the multilingual IR processes associated with multilingual QA systems, and establishing how much use is made of these tools by multilingual QA systems.

### Keywords

Information Retrieval, Question Answering Systems, Multilingual Question Answering Systems, Evaluation, Linguistic Resources, Translation

### Introduction

*I*n the field of information retrieval (hereafter IR) monolingual and multilingual tools are being created that can greatly assist specialists in their work, as well as help other users find a wide variety of information. Multilingual tools are evolving, but several years of study and research are still needed to improve implementations. One of the main difficulties facing these tools is the task of translating queries made by users and the documentary sources found in response (Diekema, 2003). Given the current expansion in research, development, and the creation of multilingual IR systems, it was considered worthwhile to analyze and evaluate the resources used by one type of these systems, the multilingual question answering systems (hereafter QA systems).

Monolingual and multilingual tools are being created that can greatly assist specialists in their work.

Although research in this area began just over a decade ago, QA systems remain largely unknown outside the field of IR. In this context, a study from the perspective of translation may offer a different approach focused on the problem of translation and linguistic resources. Researchers currently working in the field of multilingual QA systems

are searching for new methods to optimize the efficiency of IR without using too many resources for language problems. However, a system cannot easily retrieve relevant information for the user without optimal use of translation resources. For this reason, translation is crucial in this environment and enables problems to be analyzed from a fresh point of view. Any progress made in solving problems of multilingual communication can be added to existing information retrieval systems.

Recent advances in IR and Web globalization mean that multilingual search systems have been developed in which translation and language resources are as important as the documentary and computer tools. This type of system has opened a new research field that examines the most effective methods for IR and investigates which resources are required for a correct translation.

This article is primarily intended as a general purpose analysis and aims to encompass translation in the study of multilingual QA systems. The second general aim is to identify and analyze the linguistic resources and tools found in these systems. Specific objectives include identifying the main types of language resources and tools useful in the multilingual IR processes associated with multilingual QA systems, and establishing how much use is made of these tools by multilingual QA systems.

### Information retrieval

Traditionally, IR is understood as a fully automatic process that responds to a user query by examining a collection of documents and returns a sorted document list that should be relevant to the user requirements as expressed in the query. An information retrieval system is a system used to store items of information that need to be processed, searched, retrieved, and disseminated to users (Salton and Mc Gill, 1983).

When a need for information arises, a process called the "search strategy" is set in motion, which leads to the supply of documents by the system (Belkin and Croft, 1987). The process must entail (Chowdhury, 1999): *a)* definition of the informational need; *b)* selection of the information sources to be used; *c)* translation of the user query expressed in natural language into the indexing language of the information source, if necessary; *d)* translation of the expression from the indexing language to the query language of each information system; *e)* implementation of expressions obtained from the query language; *f)* assessment of results by the user and the redefinition of the query expressions if the results are not relevant; and *g)* selecting and obtaining the documents that respond to the user's needs.

An optimal IR system retrieves *all* the relevant documents (implying an exhaustive search, i.e. a high recall) and *only* the relevant documents (implying perfect accuracy, that is to say, a high precision) (Baeza-Yates and Ribeiro-Neto, 1999). This traditional model involves many implied restrictions: *a)* the assumption that users want full-text documents, rather than answers, and that the query will satisfied with these documents; *b)* that the process is direct and unidirectional rather than interactive; and finally, *c)* that the query and the document share the same language.

Multilingual IR or CLIR (*cross-lingual information retrieval*) involves at least two languages in this process. In a multilingual environment such as the Web, most IR systems (search engines) are limited to finding documents in the language of the query; or alternatively, they include machine translation systems, which are only useful once the documents are located and do not effectively cross the language barrier. Given a particular query, CLIR systems run on a collection of multilingual documents and retrieve relevant information regardless of the language used in the query (Grefenstette, 1998). Within the area of multilingual IR, the object of our study is multilingual QA systems. These systems are opening a new field of research that is becoming increasingly important within CLIR. Traditionally, CLIR is described as the problem of finding documents that the user cannot read (Oard, 2001).

One of the earlier works in CLIR was conducted by Salton (1970) and compared the effectiveness of English and German queries with that of queries obtained using a bilingual thesaurus for retrieving documents in both languages. Salton (1970) empirically showed that CLIR using a hand-crafted bilingual thesaurus is comparable with monolingual information retrieval in performance. Usually with CLIR a multilingual thesaurus of some sort is created to hold a list of descriptors for each document in a collection and the semantic relations between them, and each term in the thesaurus must be translated for each language involved (Fluhr, 1996). The descriptors can be added to the thesaurus manually or automatically (if the system can learn from previous indexing which terms are likely to be important) (López Ostenero, 2002).

These circumstances have fueled academic interest in multilingual IR, or CLIR and the techniques of natural language processing. Although Salton (1970) is considered the "father" of the earliest research initiatives concerning CLIR, the first Workshop geared specifically to CLIR topics was held in Zurich and it was organized by the Association for Computing Machinery (ACM) during the Special Interest Group on Information Retrieval, SIGIR-96 conference (Grefenstette, 1998). Nowadays, there are three important international forums about the evaluation of IR systems focusing on techniques and proceedings related with CLIR: Text Retrieval Conference (TREC) [1], the Cross-Language Evaluation Forum (CLEF)[2], the NII Text Collection for IR Systems (NTSIR)[3] or the Language Resources and Evaluation Conference (LREC) [4](Olvera-Lobo, 2009).

The TREC, co-sponsored by the National Institute of Standards and Technology (NIST) and U.S. Department of Defense, was started in 1992 (Vorhees, 1999) as part of the TIPSTER Text program. TREC was the first conference to address the issue of information retrieval and has been held annually. Its purpose was to support research within the information retrieval community by providing the infrastructure necessary for large-scale evaluation of text retrieval methodologies. TREC has successfully met its dual goals of improving the state-of-the-art in information retrieval and of facilitating technology transfer. Retrieval system effectiveness approximately doubled in the first six years of TREC. Most of today's commercial search engines include technology first developed in TREC.

One year later a new workshop is developed, called NII Test Collection for IR Systems (NTCIR). The NTCIR Workshop is a series of evaluation workshops designed to enhance research in Information Access (IA) technologies including information retrieval, question answering, text summarization, extraction, etc. The main aims of this conference are: a) to encourage research in IA technologies by providing large-scale test collections reusable for experiments and a common evaluation infrastructure allowing cross-system comparisons; b) to provide a forum for research groups interested in cross-system comparison and to exchange research ideas in an informal atmosphere; c) to investigate evaluation methods of IA techniques and methods for constructing a large-scale data set reusable for experiments.

In NTCIR it is presented works in Japanese, Chinese, Korean, and English, although multilingual papers in other languages are accepted. Nine workshops have been held since 1998, and they are divided into different areas or disciplines.

In 2000, a new conference, CLEF, was created, the most important European forum for the evaluation of multilingual and multimedia retrieval systems. CLEF has been developed to promote R&D in multilingual information access by *a)* developing an infrastructure for the testing, tuning and evaluation of information retrieval systems operating on European languages in both monolingual and multilingual contexts, and *b)* creating test-suites of reusable data which can be employed by system developers for benchmarking purposes. The final objective is to boost and encourage information retrieval technologies development in Europe in order to guarantee its competitiveness on a global sphere.

**Question answering systems**

Information overload is felt more strongly on the Web than elsewhere. All too often a query made with a web search tool (search engine or meta-search engine) results in the retrieval of too many pages, many of which are useless or irrelevant to the user. Therefore, professionals from various areas are beginning to recognize the usefulness of other types of systems, such as QA systems, for quickly and effectively finding specialist information (Crouch, et al. 2005; Lee, et al. 2006; Yu, et al. 2007).

QA systems are an evolutionary improvement on IR systems. As an alternative to traditional IR systems, they give correct and understandable answers to factual questions, rather than just offering a list of documents related to the search (Jackson and Schilder, 2005). The benefit is that users do not have to read whole documents to find the desired information. QA systems have attracted major attention since the TREC-8 (Text REtrieval) conference. QA systems were intended to stimulate human language behavior, therefore being able to answer natural language questions. As a result, the machine attempts to understand the question and thus responds by answering. In other words, understanding natural language by the machine is an essential process in the development of QA systems (Belkin and Vickery, 1985; Sultan, 2006).

Although there are various templates for making queries in QA systems, most of these systems understand questions expressed with interrogative particles

(*who, what, where, why, when, and how*); while some understand the imperative form (*tell me*). When a query is entered into the interface, the system proceeds to analyze the question by separating the word or keywords. The system then locates and extracts one or several answers from different sources of information, depending on the specialist area of the question (Olvera-Lobo and Gutiérrez-Artacho, 2010). Subsequently, the system evaluates and eliminates redundant information, or information that does not respond correctly to the question, and submits one or more prepared responses to the user (Cui, et al. 2004; Tsur, 2003).

These systems usually have a simple interface, where users can enter their queries, while some offer a list of recent queries to help users understand how the system works. QA systems handle these queries by applying algorithms and methods of linguistic analysis, as well as by using natural language processing to identify the components and determine the expected response (Zweigenbaum, 2005). This analysis usually uses a variety of standard questions in which certain words are replaced by labels accepted by the system.

QA systems may be general domain and so answer questions from diverse fields like START[5] Natural Language Question Answering System, or NSIR Question Answering System[6]. Alternatively, they may be domain-specific and focus on a specialized area, in the same way as HONQA[7] (Health On the Net Foundation). Domain-specific systems use specific linguistic resources that enable more precise answers to be given (Olvera-Lobo and Gutiérrez-Artacho, 2011a).

Another key aspect of these systems is that the system-user relationship is two-way. Establishing an interaction helps QA systems find better answers, and in turn, the QA system helps users find answers more quickly. However, it remains necessary to deepen the interactive design of these systems and enable true feedback between questions and answers, so that users communicate with the system in a conversational manner.

While the development of QA systems represents progress, the systems nevertheless suffer from restrictions. Many were only developed as prototypes, or demonstration versions, and few were marketed. Some researchers have designed and created systems that were presented and discussed at various forums and conferences. However, because the usefulness of the systems was limited to very specific contexts, or because of problems of implementation, only a few of these systems were later developed for end users (Olvera-Lobo and Robinson-García, 2009).

The aim of QA systems is to find exact and correct answer for users' questions. Analyzing the question, searching, and choosing an answer are three important items in a QA system. In addition to user interaction, various QA systems contain at least three following parts: question processing, document processing, and answer processing (Kangavari et al., 2008). The question reformulation or processing module tries to identify various ways of expressing an answer given to a natural language question. This reformulation is often used in Question Answering systems to retrieve answers in a large document collection (Verona & Motta, 2004). The objective of this step is to obtain features from the question that could be helpful in the following steps. All the information obtained by this module is given to the following steps of the system.

The document processing module depends on the architecture of the QA systems. It is used to perform a first selection of paragraphs that are considered relevant to the input question.

And finally, the answer processing module consists of two main components: answer extraction and answer validation. First, candidate answers are extracted from documents which are retrieved by the search engine in the answer extraction module. After that, the module validates answers by filtering and ranking candidate answers, and the final system suggests answers for approval by user voting (Kangavari et al., 2008). The mission of this step is to eliminate possible incorrect paragraphs contained in the list returned by the QA system. Thus, there are more possibilities of giving a correct answer at the end of the process (Rodrigo, et al. 2010).

QA systems are classified into two types according to the source of the question answer, that is, whether it is from a structured source or free text (Guthrie, 2006). The former systems return answers drawn from a database, and there are considered to be the earliest QA systems. The next generation of QA systems evolved to extract answers from the plain machine-readable text that is found in a collection of documents or on the Web.

**Multilingual QA system**

The emergence of multi-language QA systems is a relatively recent development, and many of the existing systems are still in an experimental stage. In multilingual QA systems, the language of the question may differ from the language of the retrieved document. However, QA systems differ from other CLIR systems because they do not retrieve whole documents and instead respond to queries with a short answer. According to Aceves Pérez (2008), QA systems are a set of coordinated monolingual systems in which each extracts responses from a collection of separate monolingual documents. Normally, multilingual QA systems are similar to monolingual QA systems, the main difference being the incorporation of a translation module and/or linguistic tool for multilingual retrieval (Olvera-Lobo and Gutiérrez-Artacho, 2011c). The most basic method of multilingual QA systems contains two steps. In the first step, the original question is translated to the target language by machine translation. In the second step, the translated question is processed by a QA system and the system returns an answer in the target language. More advanced methods can incorporate natural language processing tools and text mining.

Usually QA systems that deal with multiple languages rely on a translation module. The user enters his specific query, generally including some interrogative adverb (How? When? Where?) in a given natural source language. This question is translated by an automatic translator. In the stage of query analysis, the QA system examines the user´s question and determines what type of information is being requested. The classification of the questions is a key for the system, since this information will be utilized in the search stage, and in the selection and extraction of the potential responses (García Cumbreras et al., 2005). The resulting search expression will be, then, the *input*, or the formulation of the query to be used by the search engine of the system for comparing and matching it with the documents in the database. Once the documents that are relevant to the query are located, the system breaks them up into sections, selects the excerpts that include the candidate responses, and selects a final response. This response, along with its location in the corresponding document, is finally delivered to the user (Olvera-Lobo and García-Santiago, 2010).

Besides the user's interface, the multilingual QA systems architecture also contains translation software. Nowadays, machine translation has a number of different facets and views. All of them have in common that this translation must be carried out by software in a more or less automatic way. The rate and quality of this translation can vary. But even the most sophisticated IR systems cannot yet produce translations on a large scale that need absolutely no revision by a human. IR systems also have restrictions about the nature of the texts that they can translate better (Olvera-Lobo and García-Santiago, 2010; García-Santiago and Olvera-Lobo, 2010).

Most of the work in the field of multilingual QA systems has been carried out by artificial intelligence and computer specialists, so that the translation problems have not been given priority. However, translation is crucial in CLIR because queries and documents do not always share the same language. The main translation problems identified are: lexical ambiguity, lack of translation coverage, multi-modal lexemes, and errors in lexical resources (Diekema, 2003). There are some researches about the different linguistic resources and tools used by the multilingual QA systems (Olvera-Lobo and Gutierrez-Artacho, 2011b). These previous works have showed that the automatic translators remain the most popular option, despite the fact that the authors of the papers recognise the resulting problems of ambiguity.

**Method section**

An empirical methodology was adopted for this study and the collection of data about the tools and linguistic resources employed by these systems; as well as their use and implementation.

The first stage of the study focused on identifying the major conferences, meetings, and forums that address multilingual QA systems. The aim was to find, analyze, and compare the different types of linguistic resources. Although a growing number of IR conference are held each year, not all include a section devoted exclusively to QA systems, and even fewer tackle multilingual aspects. However, we identified several conferences and forums that mainly focus on research into multilingual QA systems.

In total, some 315 papers published between 2000 and 2011 at the above, and other, conferences were reviewed. No papers from 2011 were included because some of the conferences that were held are yet to publish their proceedings. Over 75% of the published articles were presented at CLEF. NTCIR had the second highest number of analysed articles, and TREC was in fourth place (see Figure 1).

Figure 1. Papers analyzed by Conference

For the studied period, the years with the largest number of papers published on multilingual QA systems were 2005 and 2008. A growing level of interest peaked in 2008; and from 2007 interest began shifting to other types of QA systems such as image, voice, and expertise domains.

Figure 2. Papers analyzed by year

We explored the resources used by multilingual QA systems showed in different conferences. We did an important documentary observation phase, because of analyzing and evaluating how many different linguistic resources and tools are used by these systems, so, it enabled us to monitor the progress made by these developers.

We studied the subject discussed in each paper – including the language resources and tools used. Although all the papers discussed the linguistic aspects of multilingual QA systems, only some tackled this as the main theme.

In a second phase we explored the resources used by existing multilingual QA systems. For some systems, it was relatively easy to obtain information because the linguistic resources were freely accessible and developers provided all the relevant literature. However, these were the exceptions. Most of the systems were partially developed prototypes and access was not available. For this reason, the documentary observation phase of our study was so important because it enabled us to monitor the progress made by these developers.

**Results and discussion**

Five main types of linguistic resources used in multilingual QA systems were identified following an analysis of the literature. The main resource types were databases, corpora, dictionaries, ontologies, and thesauri.

There were also two types of linguistic tools used by these systems, namely, automatic translators and computational grammars. These resources and tools, along with their various types and subtypes, do not run in the same way and use differing methods of processing information. Other methods for solving the problems of multilingual communication were also used—such as translation and transliteration—and these tools play an important role in several of the systems. Sometimes, a single resource was insufficient and several resources were used together to achieve better results.

Previous works (Diekema, 2003) identified four major sources of translation in CLIR – ontologies, bilingual dictionaries, automatic translators, and corpora. This study shows that CLIR has grown in popularity in recent years and that some resources are often used. Following an analysis of the literature and after identifying the resources and tools used by multilingual QA systems, a classification was made dividing these resources into two large groups: linguistic resources and linguistic tools.

Recent research and advances made in multilingual QA systems relate mainly to the more effective incorporation of new language resources, the creation of faster and more efficient systems, and the production of more transparent results. However, there remains an unsolved challenge: translation.

In analyzing the literature, it was found that the resource most used by multilingual QA systems was corpora (mostly parallel), followed by machine translation, and Wikipedia (see Figure 3). We can realize that the traditional trilogy of resources in multi-lingua QA systems (dictionaries, machine translation and corpora) has changed (Nguyen, et al. 2009).

Figure 3. Resources used in the papers reviewed

Corpora were used in 94 of the 315 papers reviewed. The apparent popularity of corpora is explained by the fact that many variants of corpora are included within the heading. The most surprising aspect of this resource was its nearly steady growth in recent years and the peak in 2008. We saw a significant decline in use in 2006, but this may be partially attributed to the fact that only 29 papers on multilingual QA systems were found for the year. Linguistic corpora are very useful resources for specialized domains. This is because the information received by users will be complete and correct when a translation is made or reviewed by professional translators. Existing corpora can be made available on the web in several languages, thus solving two of the main problems raised earlier, computational cost and storage.

Automatic translators were used in 92 of the 315 papers reviewed. This tool is often incorporated individually or in combination with other linguistic resources to offer better coverage. Although most authors confirm the problems of ambiguity and the poor quality of texts, they continue to prefer this tool because it is one of the cheapest and easiest to incorporate into systems. Automatic translation usually gives better results in general domain QA systems than in specific domains. This is because automatic translators cannot identify and correctly translate certain specialized terms. Nor can this tool be recommended for systems that use non-Western languages, or more than two languages. In fact, few automatic translators are effective in these tasks.

However, the use of automatic translators has declined in recent years. However, their presence declines substantially over the next three years (2001, 2002, and 2003). The number of multilingual QA systems using automatic translation rose again after 2006, yet not individually as in earlier years, but in combination or in support of other language resources. Automatic translators have continued to be used in the most recent years, but in a smaller number of systems.

The third most commonly used resource was Wikipedia. Wikipedia was used on 50 occasions. This is one of the most innovative resources and is growing rapidly in popularity. It was first incorporated in 2005, and its presence grew substantially the following year. The use of databases, which were used as often as Wikipedia, was very irregular, being entirely absent during some years. The results obtained by incorporating Wikipedia into such systems are unclear; some researchers claim it is one of the most useful resources, while others stress that there are many errors that are difficult to resolve.

The fourth and fifth most commonly used resources were ontologies and dictionaries, with 40 and 31 appearances respectively. Dictionaries, together with automatic translators and corpora, are the resources traditionally used by these systems, and so a similar trend is found for all three resources. However, grammar and ambiguity problems have recently reduced their popularity, so that only 5 of the 79 systems studied over the past four years used this resource.

Very different behavior is seen with the fifth application, ontologies. This resource was not used in the early years, but from the year 2004 has begun to slowly gain acceptance in multilingual QA systems. Ontologies offer many advantages and especially in specialized domain systems. Most systems are composed of texts that have been completely translated into various working languages, and so relationships are easily established. Another advantage is that there are many research teams working closely with multilingual ontologies and studying the various relationships that can be established between terms—and this existing body of work ensures a quality final product.

Other five linguistic resources were used to solve the translation problem: computational grammars, web pages, thesauri, databases and EuroWordNet.


**Conclusion**

This study has analyzed the main publications over the past nine years—from 2000 to 2010. Literature from 2011 was not included because some of the conferences that were held are yet to publish their proceedings and did not include real data regarding the situation. In total, 315 papers presented at major conferences (CLEF, TREC, NTSIR, among others) were studied and as much data as possible was extracted for an overview of the situation.

Seven of the most used resources were identified and studied: databases, dictionaries, corpora, ontologies, thesauri, Wikipedia and EuroWordNet. The second group in our study consisted of two linguistic tools: computational grammars and automatic translators. The inclusion of grammars in multilingual QA systems is relatively recent, and so the above classifications have not been taken into account.

After defining and describing these resources and tools, we considered each of the systems and techniques presented in the 315 papers. We found that corpora, in particular, parallel corpora, remain the most popular option. The second one is machine translation, despite the fact that the authors of the papers recognize the resulting problems of ambiguity. The low computational cost and ease of storage are two of the main advantages of these two resources. In our opinion, these resources can be adequate for IR when combined with others. However, there have been some changes in the use and incorporation of these resources and tools. The three most popular traditional resources (machine translation, dictionaries, and corpora) are gradually leaving a widening gap for others, such as ontologies and the free

encyclopaedia Wikipedia. In addition, other approaches such as computational grammars are slowly attracting more researchers who are experienced in handling the results they produce.

A comparison of the evolution and use of different resources and tools shows that trends favor the traditionally more popular tools (automatic translators, dictionaries and corpora). However, ontologies and Wikipedia show trends that match, or nearly match, that of the traditional resources. The remaining tools are timidly growing in popularity and have promising futures. However, the trends for each combination of tools in multilingual QA systems were not studied exhaustively.

There is a growing trend toward the translation of documents, although the option to translate queries remains the most widely one used by researchers. We believe the combination of both approaches is the most useful route and offers the best results – even handling more than two languages without difficulty.

## References

Aceves Pérez, R. M. (2008) *Búsqueda de Respuestas en Fuentes Documentales Multilingües,* PhD Thesis, Instituto Nacional de Astrofísica, Óptica y Electrónica, Mexico.

Baeza-Yates, R. and Ribeiro-Neto, B. (1999) *Modern information retrieval*. New York: ACM Press

Belkin, N.J. and Croft, W. B. (1987) Retrieval techniques. *Annual Review of Information Science and Technology*, **22**, 109–146.

Chowdhury, G. G. (1999) *Introduction to modern information retrieval*. London: Library Association.

Crouch, D., Saurí, R., and Fowler, A. (2005) AQUAINT Pilot Knowledge-Based Evaluation: Annotation Guidelines. *Palo Alto Research Center*, at: www2.parc.com/isl/groups/nltt/papers/aquaint_kb_pilot_evaluation_guide.pdf

Cui, H., Kan, M. Y., Chua, T.S. and Xiao, J. (2004) A Comparative Study on Sentence Retrieval for Definitional Question Answering, *SIGIR Workshop on Information retrieval for Question Answering* (IR4QA), Sheffield.

Diekema, A. R. (2003) *Translation Events in Cross-Language Information Retrieval: Lexical ambiguity, lexical holes, vocabulary mismatch, and correct translations.* PhD Thesis. University of Syracuse.

Fluhr, C. (1996). Multilingual Information Retrieval. *Survey of the State of the Art in Human Language Technology,* at: cslu.cse.ogi.edu/HLTsurvey/ch8node7.html

García Cumbreras, M. Á., Ureña López, L.A., Martínez Santiago, F., and Montejo Raez, A. (2005) Búsqueda de respuestas multilingüe: clasificación de preguntas en español basadas en aprendizaje, *Procesamiento del lenguaje natural*, 34, 31–40, at: www.sepln.org/revistaSEPLN/revista/34/03.pdf.

García-Santiago, L., and Olvera-Lobo, M.D. (2010) Automatic Web Translators as Part of a Multilingual Question-Answering (QA) System: Translation of Questions, *Journal of Documentation,* **66**, 434–455.

Grefenstette, G. (1998) Cross-Language Information Retrieval. *Kluwer academic publishers*, **1**.

Jackson, P., and Schilder, F. (2005) Natural Language Processing: Overview. In Brown (ed.), *Encyclopedia of Language & Linguistics,* **2** , 503–518. Amsterdam, Elsevier Press.

Kangavari, M.R., Ghandchi, S., Golpour, M. (2008) A New Model for Question Answering Systems, *World Academy of Science, Engineering and Technology,* **42**.

Lee, M., Cimino, J., Zhu, H.R., Sable, C., Shanker, V., Ely, J., and Yu, H (2006) Beyond Information Retrieval –Medical Question Answering. *AMIA.* Washington DC, USA.

López-Ostenero, F. (2002) Un Sistema Interactivo para la Búsqueda de Información en Idiomas Desconocidos por el Usuario. PhD Thesis, Departamento de Lenguajes y Sistemas Informáticos, Universidad Nacional de Educación a Distancia, Spain.

Nguyen, D., Overwijk, A., Hauff, C., Trieschnigg, D.R.B., Hiemstra, D. and de Jong, F.M.G. (2009). WikiTranslate: Query Translation for Cross-lingual

Information Retrieval using only Wikipedia. In Peters, et al. (eds.): *CLEF 2008, LNCS 5706*, 58–65.

Oard, D.W. (2001). Interactive cross-language information retrieval. *SIGIR FORUM,* **35** (1), 1–3.

Olvera-Lobo, M.D. (2009) Cross-language Information Retrieval on the Web. In Cruz-cunha, M. M. (ed.)**,**Oliveira, E. F. (ed.)**,**Tavares, A. J. V. (ed.) y Ferreira, L. G. (ed.), *Handbook Of Research On Social Dimensions Of Semantic Technologies And Web Services*, Chapter XXXIV, 704–719.

Olvera-Lobo, M.D., and Robinson-García, N. (2009) Tratamiento lingüístico de las preguntas en los sistemas de búsqueda de respuestas. *El profesional de la documentación,* **18** (2): 180–187.

Olvera-Lobo, M. D., and García-Santiago, L. (2010) Analysis of errors in the automatic translation of questions for translingual QA systems, *Journal of Documentation*, **66** (3), 434–455.

Olvera-Lobo, M.D.; and Gutiérrez-Artacho, J. (2010). Question-Answering Systems as Efficient Source of Terminological Information: Evaluation. *Health Information and Library Journal.* Vol. 27 (4): 268 – 276.

Olvera-Lobo, M.D.; and Gutiérrez-Artacho, J. (2011a). Evaluation of Open- vs. Restricted- Domain Question Answering Systems in the Biomedical Field. *Journal of Information Science.* Vol. 37 (2): 152-162.

Olvera-Lobo, M.D.; and Gutiérrez-Artacho, J. (2011b). Language resources used in multilingual Question Answering Systems. *Online Information Preview.* Vol. 35 (4): 543 – 557

Olvera-Lobo, M.D.; and Gutiérrez-Artacho, J. (2011c). Multilingual Question-Answering System in biomedical domain on the Web: an evaluation. *Lecture Notes in Computer Science.* Vol. 6941

Rodrigo, A., Pérez-Iglesias, J., Peñas, A., Garrido, G., Araujo, L. (2010) A Question Answering System based on Information Retrieval and Validation, *CLEF (Notebook Papers/LABs/Workshops) 2010.*

Salton, G. (1970) Automatic Processing of Foreign Language Documents. *Journal of American Society for Information Sciences*, 21:187–194.

Salton, G. and Mc Gill, M.J. (1983) *Introduction to Modern Information Retrieval*. New York: Mc Graw-Hill Computer Series.

Sultan, M. (2006). Multiple Choice Question Answering. PhD Thesis. University of Sheffield.

Tsur, O. (2003). *Definitional Question-Answering Using Trainable Text Classifiers.* PhD Thesis. University of Amsterdam.

Verona, M.V. and Motta, E. (2004). *AQUA, A Knowledge-Based Architecture for a Question Answering System*, Tech Report Kmi-o4-15, Knowledge media institute, Milton Keynes, England.

Voorhees, E.M. (1999) The TREC 8 Question Answering Track Report". In Voorhees and Harman (eds.), *Proceedings of the 8th Text REtrieval Conference,* 107–130.

Yu, H. and Kaufman, D. (2007) A cognitive evaluation of four online search engines for answering definitional questions posed by physicians. *Pacific Symposium on Biocomputing,* **12**, 328–339.

Zweigenbaum, P. (2005). Question answering in biomedicine. In Rijke and Webber (eds.), *Proceedings Workshop on Natural Language Processing for Question Answering,* EACL 2003, 1–4. ACL, Budapest.

---

[1] Available at: trec.nist.gov/

[2] Available at: www.clef-campaign.org/

[3] Available at: research.nii.ac.jp/ntcir/index-en.html

[4] Available at: www.lrec-conf.org/

[5] Available at: start.csail.mit.edu/

[6] Available at: tangra.si.umich.edu/clair/NSIR/html/nsir.cgi

[7] Available at: services.hon.ch/cgi-bin/QA10/qa.pl