

QUESTION-ANSWERING SYSTEMS AS EFFICIENT SOURCES OF TERMINOLOGICAL INFORMATION: AN EVALUATION

Abstract

Question-answering systems (or QA Systems) stand as a new alternative for Information Retrieval Systems. Most users often need to retrieve specific information about a factual question in order to obtain a whole document. We conducted a study to evaluate the efficiency of QA systems as terminological sources for physicians, specialized translators, and users in general. To this end we analyzed the performance of one open-domain QA system, START, and one restricted-domain QA system, MedQA. The research entailed a collection of two hundred definitional questions (*What is...?*), either general or specialized, from WebMed. We studied the sources that QA systems used to retrieve the answers, and later applied different evaluation measures to mark the quality of answers. Both QA systems were determined to be appropriate for the retrieval of terminology, proving reliable as sources and supplying correct answers.

Keywords

Question-answering systems, Evaluation of QA systems, definitional questions, sources of health information, MedQA, Start.

Introduction

Question-answering systems (heretofore QA Systems) can be viewed as a new alternative to the more familiar Information Retrieval Systems. These systems try to offer detailed, understandable answers to factual questions, in order to retrieve a collection of documents related to a particular search.¹ In recent years, the development of QA systems has been encouraged and furthered through the TREC meetings (*Text REtrieval Conference*)² – mainly since TREC-8³. This Conference has proven to be an important international forum, putting together and improving research efforts behind the different aspects of information retrieval. The QA systems try to make retrieval easier through the short-answer question models.³⁻⁴ Accordingly, users do not have to read the full text of documents such as a web page, an article of scientific journal, etc., in order to arrive at the information they need because the QA system shows the correct answer by means of a number, a noun, a short phrase or a concise extract of text.

The questions used in QA systems can be expressed using interrogative adverbs (*who, what, which, how, when, where*), or in imperative form (tell me, show, list...). Once the question is provided, the QA systems extract natural language answers, to be comprehended in a natural way for humans.⁵ QA systems follow three main steps: first of all, the systems retrieve the documents to obtain relevant sentences about the search term; they retrieve and select the sentences; and finally, they choose non-redundant definition sentences from the overall results of sentence retrieval, to

1
2
3
4
5
6 delimit the response.⁶⁻⁷ The objective pursued is for the systems to retrieve only the
7
8 correct information to answer the users' questions.⁸ Evaluation is one of the most
9
10 important dimensions in QA systems, as the process of assessing, comparing and
11
12 ranking is key to monitoring progress in the field.⁹ The main component of these
13
14 systems consists of measuring modules, which analyze the tagged sentences in the
15
16 documents selected, and compare them with the question in order to find the most
17
18 similar sentence.¹⁰⁻¹¹ Generally speaking, QA systems feature very simple and user-
19
20 friendly interfaces, and rely on methods of linguistic analysis and natural language
21
22 processing in the different phases of operation. For example, when dealing with the
23
24 questions posed by the users, they identify their component parts and then
25
26 determine the kind of answer anticipated.¹² The ones that allow users to query in
27
28 different languages are known as multilingual QA systems.
29
30
31
32
33
34
35

36 Although these systems have enriched the possibilities of information retrieval
37
38 tools, QA systems are also hampered by certain restrictions. To date, most of them
39
40 are not open-domain systems or general systems,¹³ but rather restricted-domain
41
42 systems specializing in a particular field. Moreover, all of these QA systems are
43
44 based on prototypes; that is, they are available as demos, and only in a few cases
45
46 have been marketed. A further problem can be found in the design of these systems:
47
48 what is needed is a more interactive QA procedure that allows for real feedback
49
50 between questions and answers, and user communication with the system on a
51
52 conversational level.
53
54
55
56
57
58
59
60

1
2
3
4
5
6 While not many QA systems are available on the Internet, we do have some open-
7
8 domain QA systems such as START¹⁴, developed at the Massachusetts Institute of
9
10 Technology, it is a very atypical which includes calls to OMNIBASE, a system that
11
12 integrates heterogeneous data sources using an *object-property-value model* (Katz et al.,
13
14 2002); NSIR¹⁵, developed by the University of Michigan; or Qualim¹⁶, financed by
15
16 Microsoft; in addition to some restricted-domain QA systems including MedQA¹⁷,
17
18 developed by Columbia University. In the case of NSIR and Qualim, answers are
19
20 constructed on the basis of information provided by Google¹⁸ and Wikipedia¹⁹,
21
22 respectively. Although START also retrieves information from Wikipedia, it uses
23
24 other specialized sources such as directories, databases, dictionaries, or
25
26 encyclopaedias. Meanwhile, MedQA retrieves information from the medical
27
28 database Medline, specialized dictionaries, Wikipedia and certain search engines like
29
30 Google.

31
32
33
34
35
36
37
38
39 In the Web setting, the overload of information may be perceived more acutely
40
41 than in other contexts. When users pose a given question by means of search engine
42
43 tools (including directories or metasearchers), the systems tend to retrieve an
44
45 excessive number of web pages, many of which are not relevant or useful in light of
46
47 the users' needs. Professionals in different areas claim that QA systems constitute a
48
49 good method to obtain specialized information in quick and efficient manner.²⁰⁻²²
50
51

52
53
54 In a study by Ely,²³ participating physicians spent on the average less than two
55
56 minutes looking for information to resolve clinical queries, although many of their
57
58
59
60

1
2
3
4
5
6 questions remained unanswered. Regarding this point, some researchers have
7
8 shown that the physicians trust QA systems as search methods for specialized
9
10 information retrieval.^{21,24} The general public increasingly consults knowledge
11
12 resources like the Web as well: before or after seeing a doctor, for themselves or for
13
14 relatives, to obtain information about the nature of a disease, the indications and
15
16 contraindications of a treatment, etc.¹²
17
18
19

20
21 While researchers have looked into various aspects of QA systems in recent years,
22
23 one facet that is widely overlooked is the formal evaluation of this tool and the
24
25 results it supplies. Indeed, no study to date has focused specifically on the
26
27 information sources from which responses are derived. This is the main aim of our
28
29 line of research. Ideally, QA systems should create coherent definitions in a dynamic
30
31 way, and ones that contain and summarize the most descriptive information
32
33 contained in a document collection, in view of the specific term or focus of the user
34
35 query.^{12,25}
36
37
38
39
40

41 Our objective led us to use definition-type questions in order to evaluate two QA
42
43 systems and determine the different sources behind the retrieval of medical
44
45 information. In the sections below we describe the questions used, the QA systems
46
47 analyzed and the measures of evaluation applied. Finally, we show our results and
48
49 briefly expound some conclusions.
50
51
52
53
54
55
56
57
58
59
60

Methodology

We took a sample of two hundred definitional questions about different medical issues as the basis of this study. The questions were obtained from the webpage WebMD²⁶, a US health portal providing valuable health information and support, tools for managing health problems, and specialized background on a number of illnesses. It was created by health specialists who aspire to explain, briefly yet credibly, in-depth medical information, reference material, and online community programs.

The collection of two hundred questions was created using the expression “*What is...?*” (i.e. what is irritable bowel syndrome?) in the internal search engine of the website; and in turn, WebMD provided a list of some 6000 responses in their characteristic question-answer format. We chose around 250 factual questions about different health issues, specified in Table 1. It was not our intention to evaluate the coverage of the databases sources of QA systems START and MedQA, but merely to appraise how they work and what sources they retrieve data from. This led us to finally choose 200 questions to be answered by both systems.

Authors Ely and colleagues²⁷ suggest a classification of five hierarchical categories to categorize medical questions. Firstly they distinguish between clinical and non-clinical questions: the clinical questions are further divided into general and specific; general questions are divided into evidence and no-evidence; and in turn, the evidence questions are divided into intervention and no-intervention. While not all

1
2
3
4
5
6 these categories were appropriate for our purposes, they served as the foundation
7
8 for our classification (Table 1).
9

10
11 Table 1. Categories of reference of definitional questions.
12

13 START, a QA system allowing users to pose questions about various health
14 issues, can respond to even very specialized questions within the area of health
15 care.²⁸ It has a dynamic yet easy interface, and responds quickly. Information is
16 retrieved from a very wide list of sources, such as *World Book*, *The World Factbook*
17 *2008*, *START KB*, *Internet Public Library*, and many others.
18
19
20
21
22
23
24
25

26 Meanwhile, MedQA¹⁷ is a specialized QA system that analyses thousands of
27 documents to arrive at a coherent response. Because it works specifically in the area
28 of health care, its sources are more specialized.²⁴ It also has a user-friendly interface,
29 but it is slower than START. It retrieves information from a wide array of sources,
30 including *Wikipedia*, *Medline* or *Medline Plus*.
31
32
33
34
35
36
37
38

39 After presenting the questions to both QA systems, we analyzed and evaluated
40 the answers obtained, and identified the source or sources used by the system.
41 Answers were marked as: incorrect (0 points), inexact (1 point) or correct (2 points),
42 according to the guidelines of CLEF (Cross Language Evaluation Forum)²⁹. To be
43 judged as correct, the answer had to respond accurately to the question asked, not
44 use more than 100 words in its response, and not contain irrelevant information. All
45 the questions that were answered correctly yet did not fulfil these criteria were
46 considered inexact. Likewise, we recorded the response time and the partial or total
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6 repetitions of information by the systems. The mark obtained by each question was
7
8 the baseline for application of further evaluation measures, explained below.³⁰
9

10
11 *Mean Reciprocal Rank* (MRR) is a statistical tool for evaluating any process that
12 produces a list of possible answers to a query. The reciprocal rank of a query
13 response is the multiplicative inverse of the rank of the first correct answer (for
14 example, if a question gets the correct answer in the 1st place, it will receive a score of
15 1, it would be $\frac{1}{2}$ if it is in the 2nd place, $\frac{1}{3}$ in the 3rd place...). If the answer is not
16 found, a score of 0 is assigned. MRR can be used with several correct answers, but it
17 only takes into account the first correct answer found.
18
19
20
21
22
23
24
25
26
27

28
29 *Total Reciprocal Rank* (TRR) is useful when there is more than one correct answer
30 to a question. In these cases, it is not sufficient to consider the first correct answer in
31 evaluations; instead, TRR takes into consideration all the correct answers and assigns
32 a weight to each according to its ranking in the list provided by the system. For
33 example, if the QA system provides two correct answers (the first and the third
34 ones), the TRR will be $\frac{1}{1} + \frac{1}{3}$.
35
36
37
38
39
40
41
42
43

44 *First Hit Success* (FHS) assigns 1 if the first answer returned by the system is
45 correct, and 0 if it is not. This measure, then, only accepts the first questions in the
46 list of results.
47
48
49
50

51 We used the measurement of "precision" in the evaluation of information
52 retrieval. It is understood as the capacity of system to retrieve documents or answers
53
54
55
56
57
58
59
60

(in the case of QA systems) relevant to the query and well ranked (in the case of systems ranking the results).

$$\textit{precision} = \frac{\textit{Number of relevant documents retrieved}}{\textit{Total number of documents retrieved}}$$

Results

After posing 200 questions in our QA systems, we identified the sources used by them to obtain the answers. START provided answers to the medical questions from six sources, appearing in this order: *Wikipedia*, *American Medical Association* (the only specialized source used by START), *The Internet Movie Data Bases*, *Webopedia.com*, *Yahoo* and *Merriam Webster Dictionary*.

Very briefly, *Wikipedia* is a widely used online encyclopaedia able to offer information about different issues in several languages. The website *American Medical Association*³¹ offers useful information about health for patients and physicians. *The Internet Movie Database (IMBD)*³² is an American movie site, available in some languages, with data about movies, series and actors from all over the world. *Yahoo*³³ is a directory that categorizes web pages under different subjects. *Webopedia.com*³⁴ is an online computer dictionary and internet search engine for internet terms and technical support. And finally, the *Merriam-Webster Dictionary*³⁵ is a free dictionary and thesaurus more strictly speaking, with definitions, etymology, pronunciation, etc. for each entry.

1
2
3
4
5
6 The source that offered more answers was *Wikipedia*, with a total of 182. Second
7
8 was *Merriam Webster Dictionary* with 84 answers –although 31 of these answers
9
10 repeated exactly the same information, for which reason we rejected them. *American*
11
12 *Medical Association*, the only specialized source, gave 36 answers. The sources
13
14 providing the fewest answers were *The Internet Movie Data Base (IMDB)*, *Yahoo* and
15
16 *Webopedia.com*, with 5, 2 and 1 answers obtained, respectively.
17
18
19

20
21 Table 2. Sources used by START
22

23
24 In evaluating the quality of the results by the START sources (Table 3), *Wikipedia*
25
26 was found to be the source giving more correct answers (104), with 42 answers that
27
28 were inexact and 36 others that were incorrect. Some of the inexact answers pointed
29
30 to an intermediating “window” of sorts with several options related with the query.
31
32 The general dictionary *Merriam-Webster Dictionary* offered 45 correct answers, 7
33
34 inexact ones and only one incorrect answer. The *American Medical Association*
35
36 supplied just one correct answer and 35 inexact answers. The only response obtained
37
38 through *Webopedia.com* was considered correct, whereas all the answers of *IMDB* and
39
40 *Yahoo* were incorrect.
41
42
43
44
45

46 Table 3. Answers provided by START
47

48
49 The number of answers retrieved by MedQA was higher than for START, and
50
51 most sources were of a specialized nature. *Medline*³⁶ answered all the questions. This
52
53 bibliographic database created by the *U.S. National Library of Medicine* includes
54
55 citations and specialized articles from approximately 5000 selected journals, from
56
57 1966 to the present.
58
59
60

Table 4. Sources used by MedQA

The *Dictionary of Cancer Terms*³⁷ (created by the U.S. National Institute of Cancer) and *Wikipedia* offered 192 and 191 answers, respectively. *Google* is appraised by previous authors as one of the best sources for answering definitional questions;²⁴ this search engine offered 174 answers in our experience, though 34 were rejected as repetitions. *Dorland's Illustrated Medical Dictionary*,³⁸ another non-free dictionary for health issues, gave 143 answers. *Medline Plus*,³⁹ as a multilingual medical portal with information about medication, disease and other health issues, features a medical encyclopaedia, tutorials and videos for patients; it gave us 105 answers. The multilingual glossary of *Technical and Popular Medical Terms*,⁴⁰ set up by *The European Commission* and executed by *Heymans Institute of Pharmacology* and *Mercator School*, provided 29 results. The *National Immunization Program Glossary*⁴¹ of the U.S. *Department of Health & Human Services* supplied just 3 answers.

The two QA systems evaluated here gave similar figures for repeated answers (31 repetitions in START and 34 in MedQA). In START, all the repetitions were exactly identical, and came from the same sources (*Merriam-Webster Dictionary*). In MedQA, the repetitions offered more or less the same answer, but their sources were different (*Wikipedia* and *Google*). Although a question may harvest different yet equally valid answers at a given time, when the same answer is repeated, users tend to feel confused, and the list of results increases unnecessarily. This is why we “penalized” the QA systems by not considering these answers as valid.

Table 5. Answers shown by MedQA

1
2
3
4
5
6 As we see in Table 5, there were five sources providing more correct answers than
7
8 inexact or incorrect ones: these were *Medline Plus*, *Wikipedia*, *Google*, *Technical and*
9
10 *Popular Medical Terms* and *National Immunization Program Glossary*. The only source
11
12 supplying a majority of inexact answers was *Dorland's Illustrated Medical Dictionary*,
13
14 which remitted irrelevant information about *Dorland's* itself (copyright, edition and
15
16 other non-pertinent information) in most responses. *Medline* and the *Dictionary of*
17
18 *Cancer Terms* gave more incorrect answers, and this dictionary sometimes offered
19
20 irrelevant or incorrect information. *Medline* is a bibliographical database, and it
21
22 rarely showed definitions about specific terms, but instead supplied extracts from
23
24 studies (or abstracts) by health specialists or other researchers. Thus, we may infer
25
26 that the questions were not expressed in the best possible terms. This is due to
27
28 MedQA was specifically designed and evaluated on definitional question-
29
30 answering.
31
32
33
34
35
36
37
38

39 Calculation of the time of response (time elapsing before appearance of results on
40
41 the screen) for each question led us to some interesting findings. The values obtained
42
43 were quite different for the two systems: the average response time for START was 2
44
45 to 4 seconds, while MedQA was considerably slower –with a minimum of 10
46
47 seconds and a maximum of 135 seconds. Overall, nearly 50% of the queries were
48
49 solved in a period between 26 and 35 seconds (Figure 1). During the wait, MedQA
50
51 tells users that operations are underway at that moment –first of all, the system looks
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6 over *Google*, then in *Medline*, and finally, it removes all the redundant answers to
7
8 generate the coherent ones.
9

10
11 Figure 1. Analysis of frequencies according to the response time in MedQA
12

13
14 In identifying the sources used by the two systems, we applied specific measures
15
16 for the evaluation of information retrieval. Table 6 indicates that the average number
17
18 of answers retrieved for each question is considerably higher with MedQA (5.2) than
19
20 with START (1.41). Moreover, MedQA gave, on the average, more correct responses
21
22 per question, 2.17, as compared with the 0.94 of START. This finding comes to
23
24 confirm that the more specialized system offers a more adequate coverage by subject
25
26 for the sort of query collection used here; and aside from the greater yield of
27
28 responses provided by MedQA, the average offerings of incorrect and inexact
29
30 responses are also greater under this system (1.93 and 1.08, respectively) than with
31
32 the general-domain system START (0.22 incorrect and 0.25 inexact ones).
33
34
35
36
37

38
39 Table 6. Measures for evaluating the quality of answers
40

41
42 As we explained in the section on Methods, MRR calculates the inverse value of
43
44 the first correct answer, whereas FHS simply evaluates if the first answer was correct
45
46 or not. The two measures show us, in this case, that MedQA ranks their results more
47
48 adequately, because the first correct answer tends to appear in the first place of the
49
50 list (more frequently than with START). This proves very important, as no algorithm
51
52 is involved in the ranking process. These systems, then, maintain the ranking of
53
54 answers as determined by the source they came from. In terms of user-friendliness,
55
56
57
58
59
60

1
2
3
4
5
6 FHS might be a somewhat more realistic or convenient measure, because users
7
8 usually focus on the first answer retrieved.
9

10
11 The measure TRR is lower in MedQA, however. This figure takes into account not
12
13 just the first one but all the correct responses supplied by the system, and weights
14
15 the value of the correct response in light of its placement within the list of results.
16
17 Since MedQA provides a greater amount of results, the correct responses in the
18
19 lower positions of the ranking receive less weight, and the TRR drops with respect to
20
21 that of the START, which consistently yielded fewer responses.
22
23
24

25
26 Finally, we assessed the precision of the two systems. The value obtained for
27
28 START precision was higher (67% relevant responses) than for MedQA (42%). The
29
30 percentages increased if the inexact answers were also included as relevant (84%
31
32 with START and 67% for MedQA) Therefore, we may affirm that the more
33
34 specialized system produces a greater degree of documental noise –that is, that the
35
36 correct responses are accompanied by numerous incorrect and/or inexact one.
37
38
39
40
41
42
43

44 Discussion

45
46
47 The results obtained by presenting 200 questions to the two separate systems
48
49 analysed here, START and MedQA, allowed us to subsequently evaluate their
50
51 effectiveness and their use of different information sources. Despite certain
52
53 limitations on the part of both systems (a lack of accessibility for the general public,
54
55 and insufficient development in some specific areas), we were able to confirm that
56
57
58
59
60

1
2
3
4
5
6 both are very useful in the retrieval of valid definitional health-care information,
7
8
9 with responses from both proving coherent and precise to an acceptable degree.
10
11 However, as one might expect, the answers supplied by MedQA were more reliable
12
13 that those of START in the sense that they came from specialized clinical or academic
14
15 sources, most of them showing links to research articles addressing the matter at
16
17 hand.
18
19

20
21 Another interesting finding is that the responses do not appear under a truly
22
23 representative ranking of relevance, but rather, with both systems, results are shown
24
25 in a pre-established order according to source of the information. The systems give
26
27 priority in the display of results to those sources that consistently provide answers
28
29 (like *Wikipedia* or *Google*), regardless of the criteria of reliability and credibility that
30
31 should be demanded of scientific information. Notwithstanding, we did observe that
32
33 MedQA always makes use of *Medline* in responding to queries, which can be
34
35 interpreted as a sign of reliability (yet not necessarily of precision).
36
37
38
39
40

41 Results are encouraging in that they point to the potential of this type of tool in
42
43 the more general realm of information access, as they may be a good, reliable and
44
45 reasonably precise alternative on occasions, alleviating informational overload. They
46
47 are able to provide concrete results quickly and easily. Recent studies^{9,42} have
48
49 explored various possible means of enhancing the performance of such QA systems,
50
51 for instance through the incorporation of ontology, which would heighten the
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6 quality of the answers obtained by structuring, inter-relating and formalizing all
7
8 relevant information from the thematic domain of reference.
9

10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60

- 1 Jackson, P. & Schilder, F. Natural Language Processing: Overview". In:
2
3 Brown, K. (eds). *Encyclopedia of Language & Linguistics*, 2nd. Ed. Amsterdam:
4 Elsevier Press., 2005: 503–518.
- 5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
2 Access to Text REtrieval Conference (TREC). Available from:
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
<http://trec.nist.gov/>
- 3 Voorhees, E.M. The TREC 8 Question Answering Track Report. In: Voorhees,
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
E.M. & Harman, D.K. (eds). *Proceedings of the 8th Text REtrieval Conference*,
vol. 500-246 in NIST Special Publication, NIST, Gaithersburg, Md, 1999:
107–130.
- 4 Blair-Goldensohn, S.B. & Schlaikjer, A.H. Answering Definitional Questions:
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
A Hybrid Approach. *New Directions In Question Answering*, 2004, 4: 47-58.
- 5 Costa, L.F. & Santos, D. Question Answering Systems: a partial answer.
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
SINTEF, Oslo, 2007.
- 6 Cui, H., Kan, M.Y., Chua, T.S. & Xiao, J. A Comparative Study on Sentence
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
Retrieval for Definitional Question Answering. *SIGIR Workshop on Information
retrieval for Question Answering*, Sheffield, 2004.

- 1
2
3
4
5
6 7 Mollá, D. & Vicedo, J.L. *Question-Answering in Restricted Domains*. Menlo
7
8 Park, California, AAAI Press, 2005.
9
- 10
11 8 Tsur, O. *Definitional Question-Answering Using Trainable Text Classifiers*.
12
13 PhD Thesis. *Institute of Logic Language and Computation (ILLC), University of*
14
15 *Amsterdam*, 2003.
16
17
- 18
19 9 Sing, G.O., Ardil, C., Wong, W. & Sahib, S. Response Quality Evaluation in
20
21 Heterogeneous Question Answering System: A Black-box Approach.
22
23 *Proceedings of World Academy of Science, Engineering and Technology*, 2005: 9.
24
25
- 26
27 10 Alfonseca, E., De Boni, M., Jara, J.L. & Manandhar, S. A prototype Question
28
29 Answering system using syntactic and semantic information for answer
30
31 retrieval. In: *Proceedings of the 10th Text Retrieval Conference (TREC-10)*.
32
33 Gaithersburg, 2002.
34
35
- 36
37 11 Jacquemart, P. & Zweigenbaum, P. Towards a Medical Question-Answering
38
39 System: a Feasibility Study. In: Beux, P. L. & Baud, R. (eds). *Proceedings of*
40
41 *Medical Informatics Europe (MIE '03)*, vol. 95 of *Studies in Health Technology and*
42
43 *Informatics*, San Palo, California, 2003: 463–468.
44
45
- 46
47 12 Zweigenbaum, P. Question answering in biomedicine. In: De Rijke, M. &
48
49 Webber, B. (eds). *Proceedings Workshop on Natural Language Processing for*
50
51 *Question Answering*, Budapest: ACL, EACL 2003: 1–4.
52
53
- 54
55 13 Frank, A., Kirefer, H.U., Xu, F., Uszkoreit, H., Crysmann, B., Jörg, B. &
56
57 Schäfer, U. Question answering from structured knowledge sources. *Journal of*
58
59
60

- 1
2
3
4
5
6 *Applied Logic, Special Issue on Questions and Answers: Theoretical and Applied*
7
8 *Perspectives*, 2006: 5, 20–48
9
- 10
11 14 Access to START (Natural Language Question Answering System). Available
12
13 from: <http://start.csail.mit.edu/>
14
15
- 16 15 Access to NSIR (Question Answering System). Available from:
17
18 <http://tangra.si.umich.edu/clair/NSIR/html/nsir.cgi/>
19
20
- 21 16 Access to QuaLiM (Question Answering Demo). Available from:
22
23 <http://demos.inf.ed.ac.uk:8080/qualim/>
24
25
- 26 17 Access to MedQA. Available from:
27
28 <http://monkey.ims.uwm.edu:8080/MedQA/>
29
30
- 31 18 Access to Google. Available from: <http://www.google.com/>
32
33
- 34 19 Access to Wikipedia. Available from: <http://www.wikipedia.org/>
35
36
- 37 20 Crouch, D., Saurí, R. & Fowler, A. AQUAINT Pilot Knowledge-Based
38
39 Evaluation: Annotation Guidelines. Tech. rep., *Palo Alto Research Center*, 2005.
40
- 41 21 Lee, M.; Cimino, J., Zhu, H.R., Sable, C., Shanker, V., Ely, J. & Yu, H. Beyond
42
43 Information Retrieval – Medical Question Answering. *AMIA*. Washington
44
45 DC, 2006.
46
47
48
- 49 22 Yu, H., Lee, M, Kaufman, D., Ely, J., Osheroff, J.A., Hripcsak, G. & Cimino, J.
50
51 Development, implementation, and a cognitive evaluation of a definitional
52
53 question answering system for physicians. *Journal of Biomedicine Informatics*,
54
55 2007: 4, 236-251.
56
57
58
59
60

- 1
2
3
4
5
6 23 Ely, J.W.; Osheroff, P.N.; Ebell, M.; Bergus, G.; Barcey, L.; Chambliss, M. And
7 E. Evans, 1999. "Analysis of questions asked by family doctors regarding
8 patient care". *British Medical Journal*, 1999: 319, 358–361.
9
10
11 24 Yu, H. & Kaufman, D. A cognitive evaluation of four online search engines
12 for answering definitional questions posed by physicians. *Pacific Symposium*
13 *on Biocomputing*, 2007: 12, 328-339.
14
15
16 25 Blair-Goldensohn, S.B., McKeow, K.R. & Schlaikjer, A.H. A hybrid Approach
17 for QA Track Definitional Questions. *Proceedings of TREC*, 2003: 336-343.
18
19 26 Access to WebMD. Available from: <http://www.webmd.com/>
20
21 27 Ely, J.W.; Osheroff, J., Gorman, P.N., Ebell, M.H., Chambliss, M.L., Pifer, E.A.
22 & Stavri, P.Z. A taxonomy of generic clinical questions: classification study.
23 *BMJ*, 2000: 321, 429–432.
24
25
26 28 Katz, B., Felshin, S., Yuret, D., Ibrahim, A., Lin, J., Martion, G., McFarland,
27 A.J. & Temelkuran, B. Omnibase: Uniform Access to Heterogeneous Data for
28 Question Answering. In: *Proceedings of the 7th International Workshop on*
29 *Applications of Natural Language to Information Systems (NLDB 2002)*. 2002: 230–
30 234.
31
32 29 Access to Cross Language Evaluation Forum (CLEF). Available from:
33 <http://www.clef-campaign.org/>
34
35 30 Raved, D.R., Qi, H., Wu, H. & Fan, W. *Evaluating Web-based Question*
36 *Answering Systems*. Technical Report, University of Michigan, 2001.
37
38 31 Access to American Medical Association (AMA). Available from:
39 <http://www.ama-assn.org/>
40
41 32 Access to Internet Movie Database (IMDb). Available from:
42 <http://www.imdb.com/>
43
44 33 Access to Yahoo. Available from: <http://www.yahoo.com/>
45
46 34 Access to Webopedia. Available from: <http://www.webopedia.com/>
47
48
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3
4
5
6 35 Access to Merriam-Webster. Available from: [http://www.merriam-](http://www.merriam-webster.com/)
7 [webster.com/](http://www.merriam-webster.com/)
8
9
10 36 Access to Medline. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/>
11
12 37 Access to Dictionary of Cancer Terms. Available from:
13 <http://www.cancer.gov/dictionary/>
14
15 38 Access to Dorland's Illustrated Medical Dictionary. Available from:
16 <http://www.dorlands.com/wsearch.jsp/>
17
18 39 Access to MedlinePlus. Available from: <http://medlineplus.gov/>
19
20 40 Access to Glossary of Technical and Popular Medical Terms. Available from:
21 <http://users.ugent.be/~rvdstich/eugloss/welcome.html/>
22
23 41 Access to National Immunization Program Glossary. Available from:
24 <http://www.cdc.gov/vaccines/about/terms.htm/>
25
26 42 Buitelaar, P., Cimiano, P., Frank, P., Hartung, M. & Racioppa, S. Ontology-
27 based information extraction and integration from heterogeneous data
28 sources. *Int. J. Human-Computer Studies*, 2008: 66, 759 – 788.
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 1

Question Number	PAIN	INFLAMMATION	DISEASE	SYNDROME	INFECTION	TREATMENT	OTHERS
	8	16	97	11	10	38	15

Table 1. Categories of reference of definitional questions.

For Peer Review

Table 2

Sources	Answers obtained
<i>Wikipedia</i>	182
<i>Merriam Webster Dictionary</i>	84 (31 repetitions)
<i>American Medical Association</i>	36
<i>IMDB</i>	5
<i>Yahoo</i>	2
<i>Webopedia.com</i>	1
Total	310

Table 2. Sources used by START

For Peer Review

Table 3

Source	Correct	Inexact	Incorrect
<i>Wikipedia</i>	104	42	36
<i>Merriam- Webster Dictionary</i>	45	7	1
<i>American Medical Association</i>	1	35	0
<i>Webopedia.com</i>	1	0	0
<i>Yahoo</i>	0	0	2
<i>IMDB</i>	0	0	5
Total	151	84	44

Table 3. Answers provided by START

For Peer Review

Table 4

Sources	Answer obtained
<i>Medline</i>	200
<i>Dictionary of Cancer Terms</i>	192
<i>Wikipedia</i>	191
<i>Google</i>	174 (34 repetitions)
<i>Dorland's Illustrated Medical Dictionary</i>	143
<i>Medline Plus</i>	105
<i>Technical and Popular Medical Terms</i>	29
<i>National Immunization Program Glossary</i>	3
Total	1037

Table 4. Sources used by MedQA

Table 5

Source	Correct	Inexact	Incorrect
Google	122	26	26
Wikipedia	117	31	43
Medline Plus	95	1	9
Dictionary of Cancer Terms	51	0	140
Technical and Popular Medical Terms	21	3	5
Dorland's Illustrated Medical Dictionary	14	94	35
Medline	12	61	127
National Immunization Program Glossary	2	0	1
Total	434	216	386

Table 5. Answers shown by MedQA

Table 6

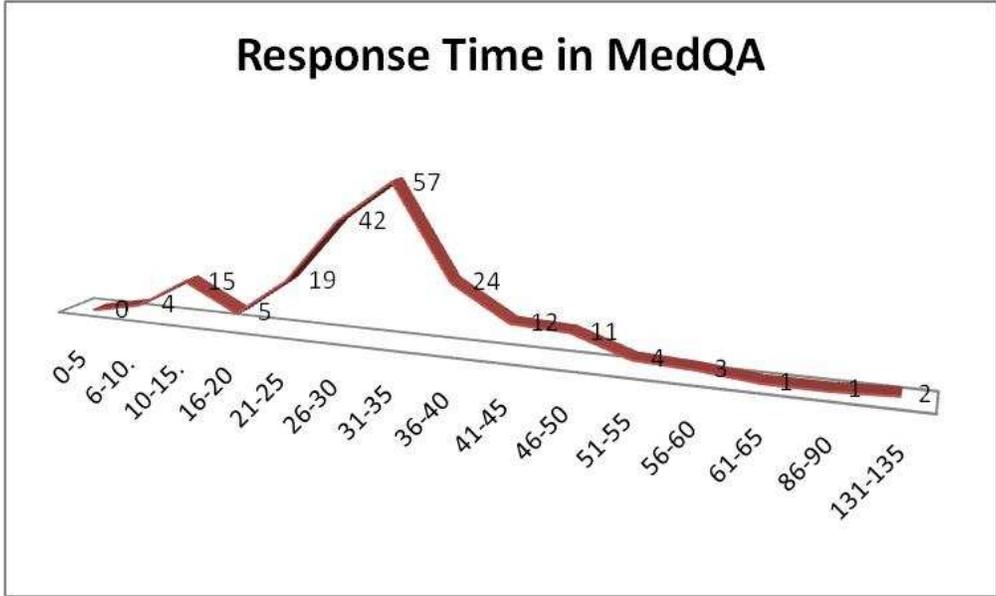
	Average answers retrieved per question	Average correct answers per question	Average incorrect answers per question	Average inexact answers per question	MRR	FHR	TRR	Precision (1)	Precision (2)
MedQA	5'18	2'17	1'93	1'08	0,86	0,75	0,40	42%	63%
START	1'41	0'94	0'22	0'25	0,60	0,61	0,59	67%	84%

(1) Taking only correct responses into account

(2) Taking both correct and inexact responses into account

Table 6. Measures for evaluating the quality of answers

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



Analysis of frequencies according to the response time in MedQA
199x120mm (96 x 96 DPI)

Peer Review