

Semantic Attributes Model for Automatic Generation of Multiple Choice Questions

Ibrahim Eldesoky Fattoh
Faculty of Information Technology,
Misr University for Science
& Technology, Egypt

Amal Elsayed Aboutabl
Computer Science Dept., Faculty of
Computers & Information, Helwan
University, Egypt

Mohamed Hassan Haggag
Computer Science Dept., Faculty of
Computers & Information, Helwan
University, Egypt

ABSTRACT

In this research, an automatic multiple choice question generation system for evaluating semantic role labels and named entities is proposed. The selection of the informative sentence and the keyword to be asked about are based on the semantic labels and named entities that exist in the question sentence. The research introduces a novel method for the distractor selection process. Distractors are chosen based on a string similarity measure between sentences in the data set. Eight algorithms of string similarity measures are used in this research. The system is tested using a set of sentences extracted from the data set for question answering. Experimental results prove that the semantic role labeling and named entity recognition approaches can be used for keyword selection. String similarity measures have been used in generating the distractors in the process of automatic multiple choice questions generation. Combining the similarity measures of some algorithms led to enhancing the results.

General Terms

Automatic Question Generation, Natural Language Processing, Text Similarity Measures, Semantic Role Labeling.

Keywords

Automatic Multiple Choice Question Generation, String Similarity Measures, Semantic Role Labeling, Named Entity Recognition.

1. INTRODUCTION

Developing Automatic Question Generation (AQG) systems has become one of the important research issues because it requires insights from a variety of disciplines including, Artificial Intelligence (AI), Natural Language Understanding (NLU), and Natural Language Generation (NLG). There are two types of question formats; multiple choice questions which asks about a word in a given sentence where the word may be an adjective, adverb, vocabulary, etc., the second format is the entity questions systems or Text to Text QG that asks about a word or phrase corresponding to a particular entity in a given sentence. In this research, the first type of question formats is covered. The traditional multiple-choice question is made up of three components, where the sentence with a gap is defined as the question sentence, the correct choice (removed word) as the key, and the other alternative choices as the distractors [1].

Anders Celsius was born in Uppsala in

- (a) Sweden (b) France (c) Switzerland (d) Finland

The above sentence is an example of multiple choice questions, the underlined gap represents the word or phrase

that is the correct answer, and the four choices represent the true answer and three distractors. This research introduces a model for a multiple choice question generator that measures the understanding of the semantic attributes in a sentence by asking about labels extracted from the given sentence using Semantic Role Labeler (SRL) and entities extracted using Named Entity Recognizer (NER). The distractors generated for the sentence are chosen based on the string similarity between the question sentence and all other sentences in the data set. The rest of the paper is organized as follows: section 2 discusses the related work of Automatic Multiple Choice Questions (AMCQ), section 3 introduces the SRL and NER in brief, section 4 provides the different string text similarity approaches, section 5 introduces the proposed model, and section 6 shows the experimental results and evaluation, and finally section 7 introduces a conclusion and future work with some remarks.

2. RELATED WORK

In this section, a review of the previous Automatic Multiple Choice Question Generation systems is introduced. An approach for AQG for vocabulary assessment was proposed in [2]; 6 types of questions were generated: definition, synonym, antonym, hypernym, hyponym, and cloze questions. Data is retrieved from WordNet after choosing the correct sense. Concerning the distractor choice, the question generation system chooses distractors of the same part of speech and similar frequency to the correct answer. Four of the six computer-generated question types were assessed: the definition, synonym, antonym, and cloze questions. The percentage of questions generated for the four types were above 60% for 156 word list. A prototype was introduced in [3] for an automatic quiz generation system for English text to test learner comprehension of text content and English skills. They used the semantic network to represent the relationship between a vocabulary and its context. They proposed two generators for two types of questions. The first generator is for sense comprehension of adjectives; the generator extracts adjectives from the SemNet of a given text as questionnaire vocabularies and form multiple-choice cloze questions. The right answer is substituted by the synonym or a similar adjective of the applied sense of the questionnaire adjective from WordNet. The second generator is for anaphor comprehension, a learner must integrate these subnets by connecting each anaphor with its antecedents. The generator identifies the antecedent of an anaphor and form a multiple-choice cloze question by scooping the anaphor out of its sentence. The options comprise its antecedent and the distractors. An approach for multiple choice questions generation for understanding the evaluation of adjectives in a text was proposed [4]. Based on the sense association among adjectives, an adjective being examined can be usually

substituted by some other adjectives. The system was able to generate three types of questions: questions for collocations, questions for antonyms, and questions for synonyms. For a given sentence, the system extracts adjective-noun pairs that exist. Then, for each adjective-noun pair, if it is a collocation, a question is generated for it. If the original sentence has words which have negative meanings, a question is generated for antonyms. Moreover, questions are generated for synonyms or similar words. The candidates of a substitute are gathered from WordNet and filtered by web corpus searching. For evaluating the generated questions, the authors chose Far East senior high school English textbook, Book One, which contains 12 articles, as the experimental material. Experimental results have shown that the proposed answer determination approaches and question filtering strategies are effective in precision. Another automatic question generation system that can generate gap-fill questions is provided in [5]. Syntactic and lexical features are used in the process of choosing the informative sentence, determining the key, and finding the distractors. The authors introduce some features as a basis for sentence selection like its position, common tokens, contains an abbreviation and others. In key selection, part of speech tagging (POS) is used to generate a list of keys. Selecting the best key from this list depends on three parameters; number of occurrences of the key in the document, is it a word in the title, the height of the key in the syntactic tree. The distractor selection depends on features such as Dice coefficient score between the gap fill sentence and the sentence containing the distractor and others. The system was tested using two chapters of the biology book and has been evaluated manually by two biology students. The sentence selection module takes 0.7 inter evaluator agreement, the key selection takes 0.75 inter evaluator agreement, and 0.60 are useful gap fill questions which have at least one good distractor. From this literature review, it can be noted that building an automatic multiple choice question generation system comprises three steps; the first is choosing the informative sentence, the second is choosing the key word or phrase to be the right answer among the multiple choices, and the last is finding the distractors for that key word. In this research, the informative sentence selection depends on whether the sentence contains any named entities or semantic labels. Keys are chosen based on the output of semantic role labeling and named entity recognizer. Distractors selection is based on the string based similarity measures as will be explained in section 4.

3. Semantic Role Labeling (SRL) and Named Entity Recognition (NER)

Semantic role labeling describes WHO did WHAT to WHOM, WHEN, WHERE, WHY, HOW etc. for a given situation and contribute to the construction of meaning [6]. For this reason, the natural language processing community has recently experienced a growth of interest in SRL. SRL has been used in many different applications such as automatic text summarization [6] and automatic question answering [7]. Given a sentence, a semantic role labeler tries to identify the predicates (relations and actions) and the semantic entities associated with each of those predicates. The set of semantic roles used in PropBank [8] includes both predicate-specific roles whose precise meaning are determined by their predicate, and general-purpose adjunct-like modifier roles whose meaning is consistent across all predicates. The predicate specific roles are Arg0, Arg1, ..., Arg5 and ArgA. A complete list of the modifier roles as proposed in the PropBank are shown in table 1. Given a sentence such as

Anders Celsius was born in Uppsala in Sweden (1)
The SRL parse would be as seen in (2).

[Andres Celsius /A0] [born /v:] [in Uppsala /AM-Loc] [in Sweden/ AM-Loc] (2)

The relation identified in (2) is the verb (born), the predicate specific roles are (Andres celsius) identified as A0 (Arg 0), is the subject of the verb, and (in Uppsala) identified as AM Location. Also, (in Sweden) is identified semantically as AM Location which is a general purpose adjunct.

Table 1: PropBank Arguments Roles

Role	Meaning
ArgM-LOC	Location
ArgM-EXT	Extent
ArgM-DIS	Discourse connectives
ArgM-ADV	Adverbial
ArgM-NEG	Negation marker
ArgM-MOD	Modal verb
ArgM-CAU	Cause
ArgM-TMP	Temporal
ArgM-PNC	Purpose
ArgM-MNR	Manner
ArgM-DIR	Direction
ArgM-PRD	Secondary prediction

Another set of semantic attributes such as persons, organizations, locations, etc. can be recognized using named entity recognition systems. Named entity recognition is an essential task in many natural language processing applications nowadays, and is given much attention in the research community and considerable progress has been achieved in many domains, such as news wire and biomedical NER [9]. Considering the sentence in (1), the output of NER would be as in (3).

[Person Andres Celsuis] born [Loc Uppsala] in [Loc Sweden]. (3)

The entity (Andres Celsuis) is identified as person while both entities (Uppsala) and (Sweden) are identified as location. All these entities could be used as a target by replacing them with gaps, one at a time. The attributes extracted from both NER and SRL (table 2) act as the key words which we search for in the sentence.

Table 2: Keyword types (labels and entities) selected from the question sentence

Keyword Types	Source
<AM-CAUS>	SRL
<Person>	NER
<AM-LOC>	SRL
<Location>	NER
<AM-TMP>	SRL
<Date>	NER
<Time>	NER

4. TEXT SIMILARITY APPROACHES

Text similarity measures play an important role in NLP applications such as text classification, information retrieval, document clustering, short answer scoring, machine translation, text summarization and others. Finding the similarity between words is a fundamental step in finding the similarity between sentences and documents [10]. Words can be lexically similar and semantically similar. Words are lexically similar if they share a similar sequence of characters and are semantically similar in other cases such as if they have the same thing, are opposite to each other, used in the same context and one is a type of another. In this research, a set of string-based similarity algorithms are applied to measure the similarity between the question sentence and the remaining sentences existing in the knowledge base. This new methodology is proposed to select the distractors in a multiple choice question. The string metric is a metric that measures the similarity or distance between two strings. The string similarity algorithms are divided into two categories; character-based and term-based similarity algorithms. In this research, three character-based algorithms and five term-based algorithms are applied to measure the similarity between two sentences. The character-based algorithms used are Smith-Waterman [11], Damerau-Levenshtein [12, 13], and Jaro [14, 15]. The five term-based algorithms applied are N-gram, Cosine similarity, Dice's coefficient [16], Jaccard similarity [17] and Block distance [18]. These algorithms are explained and implemented in SimMetrics package [19]. Fig 1 illustrates the string based algorithms applied in this research. A survey about these algorithms and text similarity approaches exists in [20].

5. PROPOSED MODEL

The automatic multiple choice questions system proposed in this research asks for semantic roles and named entities that exist in a sentence such as attributes specified in table 2. At the beginning, a knowledge base is prepared by extracting the sentences from the used dataset, then parsing them

semantically using a semantic role labeling tool and named entity recognizer for discovering the attributes that exist in the sentence. The SENNA tool is used for both purposes [21]. The sentence that has any semantic attribute is recorded in the knowledge base and its attribute is linked with it. To generate a question, the question sentence is chosen from the knowledge base and the keyword asked for is considered the labeled word or entity word identified by SENNA tool and is substituted with a gap. The distractors for the key word asked for are selected from the other keywords for the remaining sentences in the knowledge base. To find a distractor, a string similarity measure between the question sentence and all other sentences that exist in the knowledge base is applied. Then, 3 keywords are retrieved; these keywords belong to the sentences that got the highest similarity values. The three retrieved keywords are considered to be the distractors for the question sentence. Both, Algorithm 1 an Automatic Multiple Choice Question Generation System (AMCQGS) and Fig 2 show the basic steps followed in the proposed model.

Algorithm 1(AMCQGS)

```
Begin
    Build the Knowledge Base by extracting the
    sentences which have the semantic attributes from
    the dataset
    Select a question sentence and identify the semantic
    type of the keyword by parsing it semantically
    Foreach question sentence
        Measure the similarity between the
        question sentence and all sentences in the
        knowledge base.
    End for
    Sort the obtained similarity values.
    Return the three sentences that have the highest
    similarity values
    Return three keywords of the three sentences as
    distractors and identify their types.
End
```

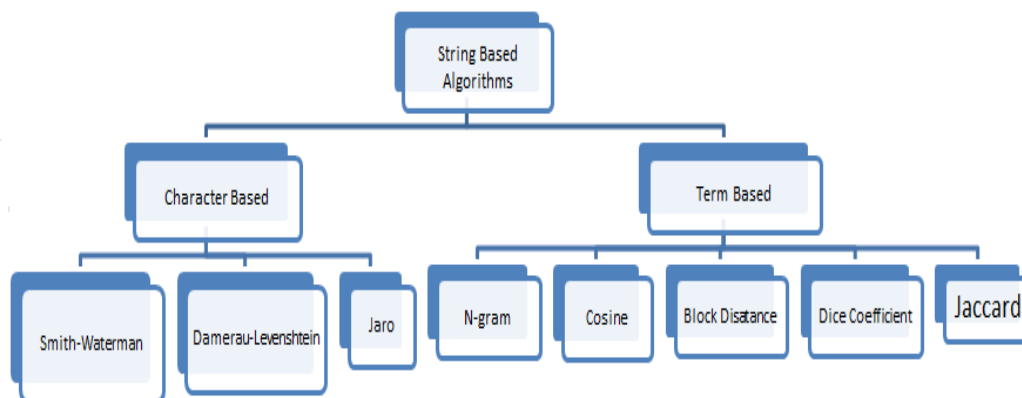


Fig 1: Text similarity algorithms applied in this research

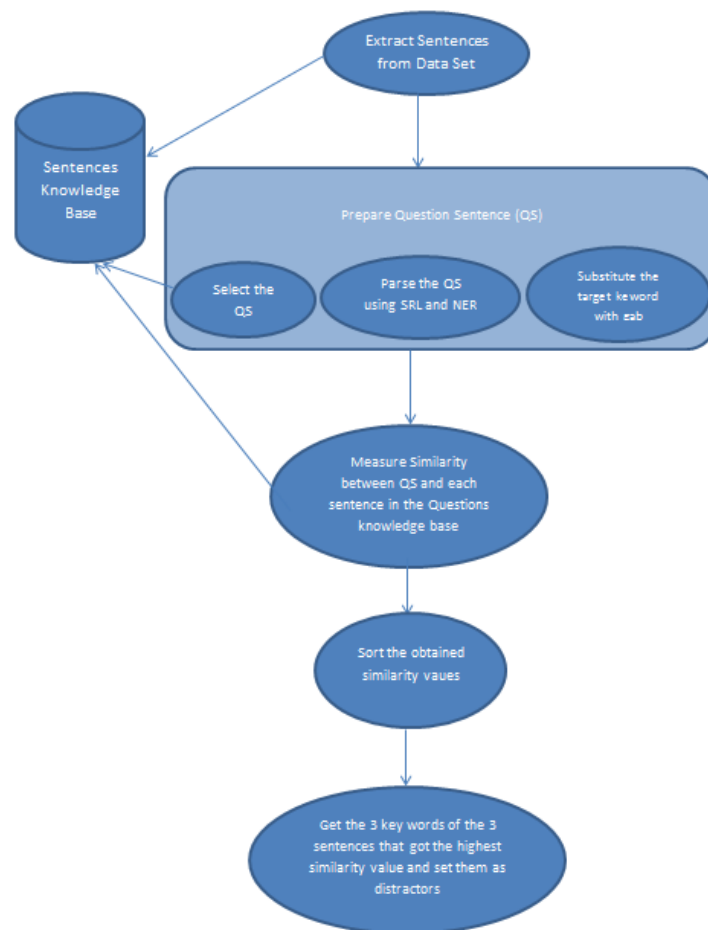


Fig 2: Flow diagram of the proposed automatic multiple choice question generation system

6. EXPERIMENTAL RESULTS AND DISCUSSION

In this section, the applied experimental results will be explained. The sentences used in this research are extracted from the TREC 2007 dataset for question answering [22]. A set of files of different domain subjects is parsed and 109 sentences are extracted to be used in testing the proposed model. The semantic attributes for these sentences are similar to types in table 2. The 109 sentences that are chosen are the sentences which yielded good result from the SENNA tool in retrieving their semantic attributes. The evaluation of both sentence selection and keyword identification depends on the output of the tool used to identify semantic attributes of a sentence. In this research out of nearly 145 parsed sentences. In this research, there were 109 sentences considered according to the keywords that are extracted from them. The distractor evaluation is the important part we tried to contribute in this research, so eight string similarity algorithms are applied trying to generate good distractors. In this research, we try to evaluate question difficulty according to the distractors generated. The question difficulties levels considered in this research are very difficult, difficult, intermediate, and easy. These levels are determined according

to the type of the generated distractor word. Each question has a true answer which is the keyword that exists in the question sentence and three distractors which are generated from the remaining sentences in the knowledge base. To evaluate the usage of the algorithms in generating the distractors, we suggest four classes for the question difficulty level. A question is very difficult if the all the generated distractors have the same type as the keyword. A question is difficult if two of the generated distractors have the same type as the keyword. A question is of intermediate difficulty if only one of the generated distractors has the same type as the keyword. A question is considered to be easy if all generated distractors are of types different from that of the key word. For more illustration, consider the question sentences in table 3.

Table 3: Examples for questions with different difficulty levels

Difficulty level	Question Sentence	Key word	Choices
Very Difficultwas the sixteenth President of the United States	Abraham Lincoln	(A) Abraham Lincoln (B) Barack Obama (C) Calvin Coolidge (D) Anders Celsius
Difficultis the sixth largest country in Europe in terms of area	Finland	(A) Abraham Lincoln (B) Finland (C) Russia (D) Switzerland
Intermediate	InSadat made a historic visit to Israel, which led to the 1979 peace treaty in exchange for the complete Israeli withdrawal from Sinai	1977	(A) Abraham Lincoln (B) 1973 (C) 1977 (D) Finland
Easy is the capital of the Republic of Austria and one of the nine states of Austria.	Vienna	(A) Abraham Lincoln (B) June 18 1953 (C) Vienna (D) 1977

According to table 3, the evaluation of the eight algorithms of string similarity is performed and their results are shown in table 4. The first column of the table shows the number of questions yielded in each class. The 45 appears in the first row for the N-gram algorithm means that the system yielded 45 questions having three distractors of the same type of the keyword asked.

Table 4: Number of sentences obtained in each class of the 8 algorithms

	N-gram	Smith	Levenshtein	Jaro	Cosine	Dice coefficient	Block Distance	Jaccard
No of very difficult questions	45	42	35	21	41	40	42	42
No of difficult questions	36	27	36	33	31	30	29	30
No of intermediate questions	21	24	26	34	19	21	20	19
No of easy questions	7	16	12	21	18	18	18	18

From table 4, it is clear that N-gram algorithm achieves the highest level of difficulty, it yielded 81 questions in the top difficult levels (very difficult and difficult), and only 28 questions for the intermediate and easy levels. Also the Jaro algorithm achieved the highest level of simplicity in the 8 algorithms; it yielded 55 questions in the intermediate and easy levels, and 54 in the difficult levels. Another measure is

introduced to measure the difficulty level of the generated questions for each algorithm using the following equation

$$\text{Difficulty level of questions} = \frac{3 * X + 2 * Y + 1 * Z}{total}$$

Where X is the number of questions that are very difficult, Y is the number of questions that are difficult, Z is the number of questions that are intermediate, and total is the total number of sentences. The overall value is divided by 3 at the end of the equation for normalizing the obtained values to get a percentage value. The value of the difficulty level of questions increases as the amount of difficult questions increases. Table 5 shows the value of the difficulty level of questions generated for each algorithm

Table 5: Difficulty level of the generated questions for the 8 algorithms

	N-gram	Smith	Levenshtein	Jaro	Cosine	Dice coefficient	Block Distance	Jaccard
Difficulty level of question	69.7%	62.4%	62.1%	48.6%	62.4%	61.5%	62.4%	62.7%

The output resulted in table 5 shows that N-gram algorithm got the highest value and the Jaro algorithm got the lowest value which proves our conclusion about both algorithms before. We assume here that a useful multiple choice question is the question having at least one good distractor. A distractor is considered to be good if it has the same type as the keyword type. Table 6 shows the percentage of good questions that are generated from each algorithm based on questions that have at least one good distractor.

Table 6: Percentage of good questions for the 8 algorithms

	N-gram	Smith	Levenshtein	Jaro	Cosine	Dice coefficient	Block Distance	Jaccard
Percentage of good questions	93.6%	85.3%	89%	80.7%	83.5%	83.5%	83.5%	83.5%

It can be noticed from table 6 that the N-gram algorithm got the highest percentage of good questions because it has the least number of easy questions. Also, the percentage value of all term-based algorithms except the N-gram is equal to 83.5%. The cause of this is that all of these algorithms have yielded the same number of easy questions as shown in table 4. The results obtained from N-gram were combined with the results obtained from both Smith and Jaccard algorithms. The cause of combining the results of these two algorithms is that they got the highest level of questions value from all 8 algorithms as shown in table 5. The similarity results obtained from N-gram is added to the similarity results yielded from Smith algorithm to perform the combination, and the same with the Jaccard algorithm. Table 7 shows the results yielded by combining the results of two different algorithms.

Table 7: Combining results of 2 different algorithms

	N-gram+ Smith	N-gram+ Jaccard
No of Very difficult questions.	42	46
No of difficult questions.	30	32
No of intermediate questions..	26	23
No of easy questions.	11	8
Difficulty level of questions	64.8%	68.8%
Percentage of good questions	89.9%	92.7%

From table 7, it is clear that the values obtained from combining both N-gram and Smith algorithms outperform the values obtained from the Smith's results only in Difficulty level of questions and Percentage of good questions. Combining N-gram's results with Jaccard's results yield an increase in both values compared to Jaccard's results. It is noticed that the N-gram algorithm performs better than when combined with either Smith or Jaccard algorithms.

7. CONCLUSION AND FUTURE WORK

This research introduces an automatic generation of multiple choice questions based on the semantic attributes in the question sentence. The semantic attributes are extracted using both semantic role labeling and named entity recognition tools. The distractor generation process introduced in this paper is based on string similarity measures between the question sentence and all other sentences that exist in the knowledge base of all sentences in the system. Eight algorithms of string based similarity are applied for all sentences and the obtained results are analyzed. A

classification is introduced to identify the question difficulty level. The eight algorithms produce promising results in the process of generating distracters. The N-gram algorithm outperformed all other algorithms in producing the highest level of questions difficulty. Combining the results of more than one algorithm is attempted and the output of this process enhances the difficulty level of some algorithms. In the future semantic similarity measures like corpus-based similarity and knowledge base similarity algorithms can be attempted. A prior classification of the sentences in the knowledge base according the key word types can be introduced to increase the level of difficulty of the generated questions.

8. REFERENCES

- [1] [1] Chen, C. Y., Liou, H. C., & Chang, J. S. (2006, July). Fast: an automatic generation system for grammar tests. In Proceedings of the COLING/ACL on Interactive presentation sessions (pp. 1-4). Association for Computational Linguistics.
- [2] Brown, J., Firshkoff, G. And Eskenazi, M. (2005) Automatic Question Generation For Vocabulary Assessment. Proceedings OfHlt/Emnlp, 819–826. Vancouver, Canada.
- [3] Sung, L., Lin, Y., And Chern, M.(2007). An Automatic Quiz Generation System For English Text. Seventh Ieee International Conference On Advanced Learning Technologies.
- [4] Lin, Y., Sung, L., AndChern, M (2007). An Automatic Multiple-Choice Question Generation Scheme For English Adjective Understanding. Workshop On Modeling, Management And Generation Of Problems/Questions In Elearning, The 15th International Conference On Computers In Education (Icce 2007), Pages 137-142, Hiroshima, Japan.
- [5] Agarwal , M. And Mannem .P. (2011). Automatic Gap-Fill Question Generation From Text Books. In Proceedings Of The 6th Workshop On Innovative Use Of Nlp For Building Educational Applications. Portland, Or, Usa. Pages 56-64.
- [6] Trandabăț, D. Using semantic roles to improve summaries.(2007). In The 13th European Workshop on Natural Language Generation (p. 164).
- [7] PizzatoL., and Molla D. (2008). Indexing on Semantic Roles for Question Answering. Coling 2008: Proceedings of the 2nd workshop on Information Retrieval for Question Answering (IR4QA). Pages 74–81. Manchester, UK.
- [8] Palmer M., Gildea D., and Kingsbury P. (2005). The proposition bank: An annotated corpus of semantic roles. Computational Linguistics, 31(1):71–106.
- [9] Tkachenko M., and Simanovisky A. (2012). Named Entity Recognition: Exploring Features. Proceedings of KONVENS 2012 (Main track: oral presentations), Vienna.
- [10] Goma, W. H. And Fahmy, A. A.(2014). Automatic Scoring for Answers to Arabic Test Questions. Computer Speech & Language 28 (4), 833-857.
- [11] Smith, F. T. & Waterman, S. M. (1981). Identification of Common Molecular Subsequences, Journal of Molecular Biology 147: 195–197.

- [12] Hall, P. A. V. & Dowling, G. R. (1980) Approximate string matching, *Comput. Surveys*, 12:381-402.
- [13] Peterson, J. L. (1980). Computer programs for detecting and correcting spelling errors, *Comm. Assoc. Comput. Mach.*, 23:676-687.
- [14] Jaro, M. A. (1989). Advances in record linkage methodology as applied to the 1985 census of Tampa Florida, *Journal of the American Statistical Society*, vol. 84, 406, pp 414-420.
- [15] Jaro, M. A. (1995). Probabilistic linkage of large public health data file, *Statistics in Medicine* 14 (5-7), 491-8.
- [16] Dice, L. (1945). Measures of the amount of ecologic association between species. *Ecology*, 26(3).
- [17] Jaccard, P. (1901). Étude comparative de la distribution florale dans une portion des Alpes et des Jura. *Bulletin de la Société Vaudoise des Sciences Naturelles* 37, 547-579.
- [18] Eugene F. K. (1987). *Taxicab Geometry*, Dover. ISBN 0-486-25202-7.
- [19] Chapman, S. (2009). *Simmetrics: a java & c#.Net library of similarity metrics*.
- [20] Gomaa W. H. And Fahmy A. A. (2013). A survey of text similarity approaches. *International Journal of Computer Applications*, 68(13), 13-18.
- [21] Collobert R., Weston J., Bottou L E., Karlen M, Kavukcuoglu K, and Kuksa P. (2011). Natural language processing (almost) from scratch. *The Journal of Machine Learning Research*, 12:2493–2537, 2011.

- [22] TREC Tracks, NIST <http://trec.nist.gov/tracks.html> visited Nov. 11, 2013.

9. AUTHOR'S PROFILE

Ibrahim Eldesoky Fattoh is currently working as a teacher assistant, Computer science department, Faculty of Information Technology, Misr University for Science and Technology (MUST), Cairo, Egypt. He is a Ph.D student, Faculty of Computer and Information, Helwan University, Egypt in the area of Automatic Question Generation under supervision of Prof. Mohammad Hassan Haggag and Assoc. Prof. Amal Elsayed Aboutabl. He finished his B. Sc. and Master degrees at Faculty of Computers and Information, Helwan University, Egypt. His master thesis was entitled "Supervised Immune System for Information Filtering". His research interests include Artificial Intelligence, Natural Language Processing, Data Mining, Text Mining, and Computational Intelligence.

Amal Elsayed Aboutabl is currently an Associate Professor at the Computer Science Department, Faculty of Computers and Information, Helwan University, Cairo, Egypt. She received her B.Sc. in Computer Science from the American University in Cairo and both of her M.Sc. and Ph.D. in Computer Science from Cairo University. She worked for IBM and ICL in Egypt for seven years. She was also a Fulbright Scholar at the Department of Computer Science, University of Virginia, USA in 2009. Her current research interests include parallel computing, performance evaluation and text processing.