

# A Hybrid Question Answering System for Multiple Choice Question (MCQ)

Pinaki Bhaskar<sup>1</sup>, Somnath Banerjee<sup>1</sup>, Partha Pakray<sup>1</sup>, Samadrita Banerjee<sup>2</sup>,  
Sivaji Bandyopadhyay<sup>1</sup>, and Alexander Gelbukh<sup>3</sup>

<sup>1</sup> Department of Computer Science and Engineering,  
Jadavpur University, Kolkata – 700032, India

<sup>2</sup> School of Cognitive Science,  
Jadavpur University, Kolkata – 700032, India

<sup>3</sup> Center for Computing Research,  
National Polytechnic Institute, Mexico City, Mexico  
{pinaki.bhaskar, s.banerjee1980, parthapakray, samadrita.banerjee}@gmail.com,  
sivaji\_cse\_ju@yahoo.com, gelbukh@gelbukh.com

**Abstract.** The article presents the experiments carried out as part of the participation in the main task (English dataset) of QA4MRE@CLEF 2013. In the developed system, we first combine the question Q and each candidate answer option A to form (Q, A) pair. Each pair has been considered a Hypothesis (H). We have used Morphological Expansion to rebuild the H. Then, each H has been verified by assigning a matching score. Stop words and interrogative words are removed from each H and query words are identified to retrieve the most relevant sentences from the associated document using Lucene. Relevant sentences are retrieved from the associated document based on the TF-IDF of the matching query words along with n-gram overlap of the sentence with the H. Each retrieved sentence defines the Text T. Each T-H pair is assigned a ranking score that works on textual entailment principle. The inference weight i.e., matching score has automatically been assigned to each answer options based on their inference matching. Each sentence in the associated document has contributed an inference score to each H. The candidate answer option that receives the highest inference score has been identified as the most relevant option and selected as the answer to the given question.

**Keywords:** Question Answering technique, QA4MRE Data Sets, Named Entity, Textual Entailment, Machine Reading.

## Introduction

Machine Reading is currently one of the most difficult and challenging task of Artificial Intelligence. Machine Reading involves not only parsing of text but also constructing a coherent internal model of the world that the text is describing using extensive

background knowledge to fill in the gaps and resolve ambiguities (Schank and Abelson, 1977). The main objective of QA4MRE [3] is to develop a methodology for evaluating Machine Reading systems through Question Answering and Reading Comprehension Tests. Machine Reading task obtains an in-depth understanding of just one or a small number of texts. The task focuses on the reading of single documents and identification of the correct answer to a question from a set of possible answer options. The identification of the correct answer requires various kinds of inference and the consideration of previously acquired background knowledge. Ad-hoc collections of background knowledge have been provided for each of the topics in all the languages involved in the exercise so that all participating systems work on the same background knowledge. Texts have been included from a diverse range of sources, e.g. newspapers, newswire, web, blogs, Wikipedia entries.

Answer Validation (AV) is the task of deciding for given a question and an answer from a QA system, whether the answer is correct or not and it was defined as a problem of RTE in order to promote a deeper analysis in Question Answering [3]. Answer Validation Exercise (AVE) is a task introduced in the QA@CLEF competition. AVE task is aimed at developing systems that decide whether the answer of a Question Answering system is correct or not. There were three AVE competitions AVE 2006 [4], AVE 2007 [5] and AVE 2008 [6]. AVE systems receive a set of triplets (Question, Answer and Supporting Text) and return a judgment of “SELECTED”, “VALIDATED” or “REJECTED” for each triplet.

Section 2 describes the task; Section 3 describes the corpus statistics; Section 4 describes the system architecture. The experiments carried out on test data sets are discussed in Section 5 along with the results. The conclusions are drawn in Section 6.

## 2 Task Description

In contrast to text mining (or text harvesting, sometimes also called macro-reading), where the system reads and combines evidence from hundreds or even thousands of texts, Machine Reading is the task of obtaining an in-depth understanding of just one, or a small number, of texts.

As in the previous campaign, the task focuses on the reading of single documents and the identification of the answers to a set of questions about information that is stated or implied in the text. Systems should be able to use knowledge obtained automatically from given texts to answer a set of questions posed for single documents at a time. Questions are in the form of multiple choice, where a significant portion of questions have no correct answer among the given alternatives proposed. While the principal answer is to be found among the facts contained in the test documents provided, systems may use knowledge from additional given texts (the ‘Background Corpus’) to assist them with answering the questions. Some questions will also test a

system's ability to understand certain propositional aspects of meaning such as modality and negation.

Participating systems will be required to answer the questions of test data. Test Questions are in the form of multiple choices: for each question, 5 possible answers are given. The system has to focus on testing the comprehension of single document. Though the direct and immediate answer is always present in the test document, but to recognize that it is the answer, systems may need some background knowledge and various kinds of textual inferences may be needed, e.g., lexical (acronymy, synonymy, hyperonymy), syntactic (nominalization / verbalization, causative, paraphrase, active/passive), discourse (coreference, anaphora ellipsis), etc. There will always be one and only one correct option. Systems will also have the chance to leave some questions unanswered if they are not confident about the correctness of their response. The system is not required to answer every question, as the C@1 measure is used for evaluation. Therefore, there are three possibilities:

- To submit an answer and ask for it to be evaluated,
- Not to submit an answer,
- To submit an answer and ask for it not to be evaluated.

### 3 Corpus Statistics

In addition to the test document, systems are provided with a collection of additional texts on the same topic, from which they may acquire the reading capabilities and draw the knowledge, if needed, to overcome any knowledge gaps in the source text. The 2013 background collections are based on but not identical to the 2012 collections. Texts are drawn from many sources: newspapers, newswire, web pages, blogs and Wikipedia entries. Thus the kind of knowledge provided is generic with respect to each topic, containing for example the most common classes and instances, frequent assertions, and general relations between these assertions such as causality, etc.

The 2013 test set will be composed of 4 topics, namely “**Aids**”, “**Climate change**” and “**Music and Society**” and “**Alzheimer**”. Each topic includes 4 reading tests. Each reading test will consist of one document, accompanied by 15 to 20 questions, each with a set of five answer options per question. So, for each language task, there will be in total:

- - 16 test documents (4 documents for each of the four topics)
- - 240/320 questions (15/20 questions for each document) with
- - 1200/1600 choices/options (5 for each question)

Test documents, questions, and options are made available in Arabic, Bulgarian, English, Romanian, and Spanish. These materials will be exactly the same in all

languages, created using parallel translations. We have worked only with English language data. The Background Collections (one for each topic) are comparable (but not identical) topic-related collections created in all the different languages.

#### 4 Machine Reading System Architecture

The architecture of machine reading system is described in Figure 1. Proposed architecture is made up of four main modules along with knowledgebase. Each of these modules is now being described in subsequent subsections.

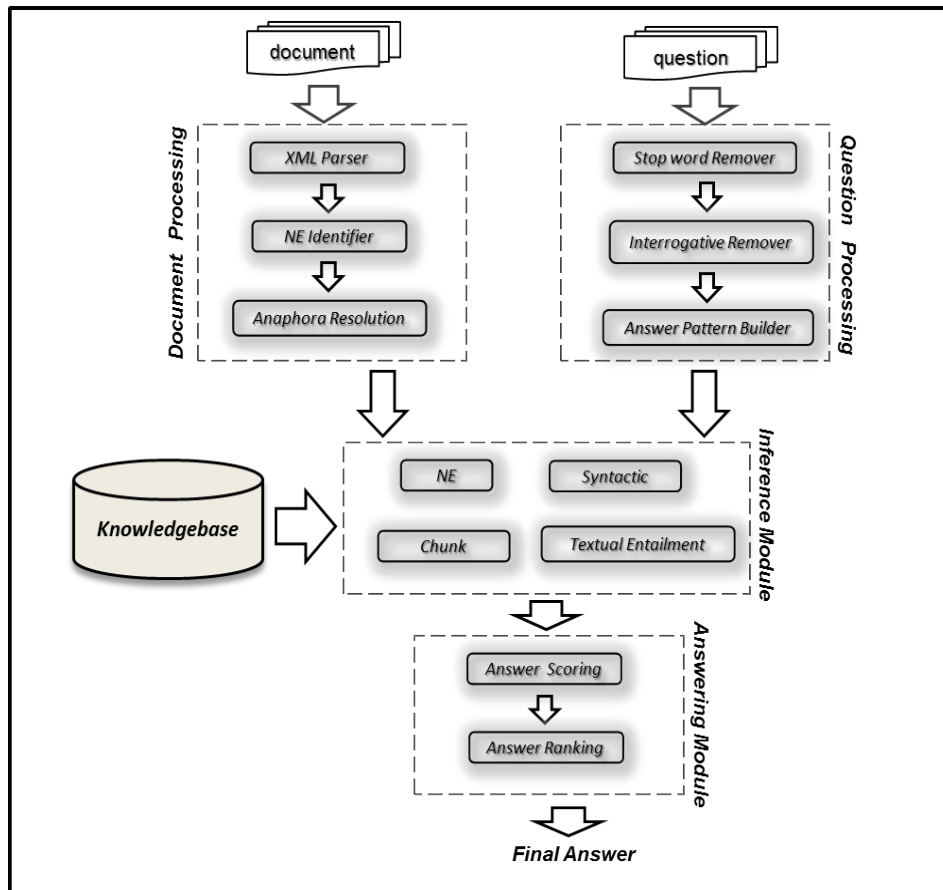


Fig. 1: System Architecture

#### 4.1 Document Processing Module

Document processing module consists of three sub-modules: XML Parser, Named Entity (NE) Identification and Anaphora Resolution.

**XML parser.** The given XML corpus has been parsed using XML parser. The XML parser extracts the document and associated questions. After parsing, the documents and the associated questions are extracted from the given XML documents and stored in the system.

**Named Entity (NE) Identification.** For each question, system must identify the correct answer among the proposed alternative answer options. Each generated answer pattern corresponding to a question is compared with each sentence in the document to assign an inference score. The score assignment module requires that the named entities in each sentence and in each answer pattern are identified. The CRF-based Stanford Named Entity Tagger<sup>1</sup> (NE Tagger) has been used to identify and mark the named entities in the documents and queries. The tagged documents and queries are passed to the lexical inference sub-module.

**Anaphora Resolution.** It has been observed that resolving the anaphors in the sentences in the documents improves the inference score of the sentence with respect to each associated answer option. To resolve the anaphora BART<sup>2</sup> (Beautiful Anaphora Resolution Toolkit) has been used in the present task. BART performs automatic co-reference resolution, including all necessary preprocessing steps.

#### 4.2 Question Processing Module

This module responsible for generating the answer pattern and deciding answers or not answers to a question. This module consists of - *Stop word Remover*, *Interrogative Remover*, and *Answer Pattern Builder*. Stop words and interrogatives have been removed from question text by *Stop word Remover*, *Interrogative Remover* sub-modules respectively to build query terms (QT). Then, an answer pattern is built by (QT, OPTION<sub>T</sub>) pair; where OPTION<sub>T</sub> refers the T-th answer option. Each answer pattern is considered as a hypothesis H. So, five hypotheses have been built for each question.

Then, Porter stemmer (Porter, 1980) has been applied to the hypothesis H. *Morphological Expansion* has been applied to the hypothesis. [13] show that in some cases performing stemming decreases the overall IR performance compared with a

---

<sup>1</sup> <http://nlp.stanford.edu/ner/index.shtml>

<sup>2</sup> <http://www.bart-coref.org/>

simple bag-of-words approach. They also show that both recall and the ranking of relevant documents is negatively affected by stemming (they evaluate the ranking of relevant documents using a weighted recall approach referred to as total document reciprocal rank). They also suggest that morphological variations may be added to query to improve performance. For example the question “What lays blue eggs” would, under the three different approaches, be converted to:

**Bag-of-Words:** blue  $\wedge$  eggs  $\wedge$  lays

**Stemming:** blue  $\wedge$  egg  $\wedge$  lai

**Morphological Expansion:** blue  $\wedge$  (eggs  $\vee$  egg)  $\wedge$  (lays  $\vee$  laying  $\vee$  lay  $\vee$  laid)

So, if the hypothesis H contains the words  $Q_1, Q_2, \dots, Q_N$ , then after morphological expansion the hypothesis H may be :

$$(Q_1 \vee M_{11} \vee M_{12} \vee \dots \vee M_{1K}) \wedge (Q_2 \vee M_{21} \vee M_{22} \vee \dots \vee M_{2P}) \wedge \dots \wedge (Q_N \vee M_{N1} \vee M_{N2} \vee \dots \vee M_{NQ})$$

Where,  $(M_{11}, M_{12}, \dots, M_{1K}), (M_{21}, M_{22}, \dots, M_{2P}), \dots, (M_{31}, M_{32}, \dots, M_{NQ})$  are the morphological variants of  $Q_1, Q_2, \dots, Q_N$  respectively.

### 4.3 Inference Module

This module assigns inference score to each hypothesis. Each hypothesis has been considered as a query in this module.

**Answer Validation.** The corpus is in XML format. All the XML test data has been parsed before indexing using our XML Parser. The XML Parser extracts the sentences from the document. After parsing the documents, they are indexed using Lucene, an open source full text search tool.

*Query Word Identification and Sentence Retrieval,* After indexing has been done, the queries have to be processed to retrieve relevant sentences from the associated documents. Each answer pattern or query is processed to identify the query words for submission to Lucene. Each hypothesis has been submitted to Lucene after removing *stop words* (using the stop word list<sup>3</sup>). The remaining words are identified as the query words. Query words may appear in inflected forms in the question. For English,

---

<sup>3</sup> <http://members.unine.ch/jacques.savoy/clef/>

standard Porter Stemming algorithm<sup>4</sup> has been used to stem the query words. After searching using Lucene, a set of sentences in ranked order are retrieved.

First of all, all query words are fired with AND operator. If at least one sentence is retrieved using the query with AND operator then the query is removed from the query list and need not be searched again. The rest of the queries are fired again with OR operator. OR searching retrieves at least one sentence for each query. Now, the top ranked relevant ten sentences for each query are considered for further processing. In case of AND search only the top ranked sentence is considered. Sentence retrieval is the most crucial part of this system. We take only the top ranked relevant sentences assuming that these are the most relevant sentences in the associated document for the question from which the query has been generated.

Each retrieved sentence is considered as the Text (T) and is paired with each generated hypothesis (H). Each T-H pair identified for each answer option corresponding to a question is now assigned a score based on the NER module, Textual Entailment module, Chunking module, Syntactic Similarity module and Question Type module.

*NER Module.* It is based on the detection and matching of Named Entities (NEs) [9] in the Retrieved Sentence (T) - generated Hypothesis (H) pair. Once the NEs of the hypothesis and the text have been detected, the next step is to determine the number of NEs in the hypothesis that match in the corresponding retrieved sentence. The measure NE\_Match is defined as  $NE\_Match = \frac{\text{number of common NEs between T and H}}{\text{Number of NEs in Hypothesis}}$ .

If the value of NE\_Match is 1, i.e., 100% of the NEs in the hypothesis match in the text, then the T-H pair is considered as an entailment. The T-H pair is assigned the value "1", otherwise, the pair is assigned the value "0".

*Textual Entailment Module (TE).* This TE module [8] is based on three types of matching, i.e., WordNet based Unigram Match and Bigram Match and Skip-bigram Match.

a) WordNet based Unigram Match: In this method, the various unigrams in the hypothesis for each Retrieved Sentence (T) - generated Hypothesis (H) pair are checked for their presence in the retrieved text. WordNet synsets are identified for each of the unmatched unigrams in the hypothesis. If any synset for the H unigram match with any synset of a word in the T then the hypothesis unigram is considered as a successful WordNet based unigram match. If the value of Wordnet\_Unigram\_Match is 0.75 or more, i.e., 75% or more unigrams in the H match either directly or through WordNet synonyms, then the T-H pair is considered as an entailment. The T-H pair is then assigned the value "1", otherwise, the pair is assigned the value "0".

---

<sup>4</sup> <http://tartarus.org/~martin/PorterStemmer/java.txt>

b) **Bigram Match:** Each bigram in the hypothesis is searched for a match in the corresponding text part. The measure *Bigram\_Match* is calculated as the fraction of the hypothesis bigrams that match in the corresponding text, i.e.,  $\text{Bigram\_Match} = (\text{Total number of matched bigrams in a T-H pair} / \text{Number of hypothesis bigrams})$ . If the value of *Bigram\_Match* is 0.5 or more, i.e., 50% or more bigrams in the H match in the corresponding T, then the T-H pair is considered as an entailment. The T-H pair is then assigned the value “1”, otherwise, the pair is assigned the value “0”.

c) **Skip-grams:** A skip-gram is any combination of *n* words in the order as they appear in a sentence, allowing arbitrary gaps. In the present work, only 1-skip-bigrams are considered where 1-skip-bigrams are bigrams with one word gap between two words in a sentence. The measure *1-skip\_bigram\_Match* is defined as

$$1\_skip\_bigram\_Match = skip\_gram(T,H) / n,$$

where *skip\_gram(T,H)* refers to the number of common 1-skip-bigrams (pair of words in order with one word gap) found in T and H and *n* is the number of 1-skip-bigrams in the hypothesis H. If the value of *1-skip\_bigram\_Match* is 0.5 or more, then the T-H pair is considered as an entailment. The text-hypothesis pair is then assigned the value “1”, otherwise, the pair is assigned the value “0”.

*Chunk Module.* The question sentences are pre-processed using Stanford dependency parser. The words along with their part of speech (POS) information are passed through a Conditional Random Field (CRF) based chunker [11] to extract phrase level chunks of the questions. A rule-based module is developed to identify the chunk boundaries. The question-retrieved text pairs that achieve the maximum weight are identified and the corresponding answers are tagged as “1”. The question-retrieved text pair that receives a zero weight is tagged as “0”.

*Syntactic Similarity Module.* This module is based on the Stanford dependency parser [9], which normalizes data from the corpus of text and hypothesis pairs, accomplishes the dependency analysis and creates appropriate structures.

**Matching Module.** After dependency relations are identified for both the retrieved sentence and the hypothesis in each pair, the hypothesis relations are compared with the retrieved text relations. The different features that are compared are noted below. In all the comparisons, a matching score of 1 is considered when the complete dependency relations along with all of its arguments match in both the retrieved sentence and the hypothesis. In case of a partial match for a dependency relation, a *matching score* of 0.5 is assumed.

- a. **Subject-Verb Comparison:** The system compares hypothesis subject and verb with retrieved sentence subject and verb that are identified through the *nsubj* and *nsubjpass* dependency relations. A matching score of 1 is assigned in case of a complete match. Otherwise, the system considers the following matching process.



- b. WordNet Based Subject-Verb Comparison: If the corresponding hypothesis and sentence subjects do match in the subject-verb comparison, but the verbs do not match, then the WordNet distance between the hypothesis and the sentence is compared. If the value of the WordNet distance is less than 0.5, indicating a closeness of the corresponding verbs, then a match is considered and a *matching score* of 0.5 is assigned. Otherwise, the subject-subject comparison process is applied.
- c. Subject-Subject Comparison: The system compares hypothesis subject with sentence subject. If a match is found, a score of 0.5 is assigned to the match.
- d. Object-Verb Comparison: The system compares hypothesis object and verb with retrieved sentence object and verb that are identified through *dobj* dependency relation. In case of a match, a *matching score* of 0.5 is assigned.
- e. WordNet Based Object-Verb Comparison: The system compares hypothesis object with text object. If a match is found then the verb corresponding to the hypothesis object with retrieved sentence object's verb is compared. If the two verbs do not match then the WordNet distance between the two verbs is calculated. If the value of WordNet distance is below 0.5 then a *matching score* of 0.5 is assigned.
- f. Cross Subject-Object Comparison: The system compares hypothesis subject and verb with retrieved sentence object and verb or hypothesis object and verb with retrieved sentence subject and verb. In case of a match, a *matching score* of 0.5 is assigned.
- g. Number Comparison: The system compares numbers along with units in the hypothesis with similar numbers along with units in the retrieved sentence. Units are first compared and if they match then the corresponding numbers are compared. In case of a match, a *matching score* of 1 is assigned.
- h. Noun Comparison: The system compares hypothesis noun words with retrieved sentence noun words that are identified through *nn* dependency relation. In case of a match, a matching score of 1 is assigned.
- i. Prepositional Phrase Comparison: The system compares the prepositional dependency relations in the hypothesis with the corresponding relations in the retrieved sentence and then checks for the noun words that are arguments of the relation. In case of a match, a *matching score* of 1 is assigned.
- j. Determiner Comparison: The system compares the determiner in the hypothesis and in the retrieved sentence that are identified through *det* relation. In case of a match, a *matching score* of 1 is assigned.
- k. Other relation Comparison: Besides the above relations that are compared, all other remaining relations are compared verbatim in the hypothesis and in the retrieved sentence. In case of a match, a *matching score* of 1 is assigned.

API for WordNet Searching RiWordnet<sup>5</sup> provides Java applications with the ability to retrieve data from the WordNet database.

Each of the matches through the above comparisons is assigned some weight.

---

<sup>5</sup> <http://www.rednoise.org/rita/wordnet/documentation/index.htm>

### 4.3 Answering Module

In this module, we have got the weight from Named Entity Recognition (NER) Module , Textual Entailment (TE) Module, Question Type Analysis Module , Chunk Boundary and Syntactic Similarity Module.

Each sentence in the associated document is assigned an inference score with respect to each  $(QT, OPTION_T)$  pair. Each question has five answer options and the task is to identify the best answer to the question from an associated document. Each question in the system is identified as the (question, document) pair represented as  $\{q_i, d\_id\}$  where  $i=1\dots5$ . There are 5 questions corresponding to each document. Each answer option is represented in the system as  $\{d\_id, q\_id_i, a\_id_j\}$ , where,  $d\_id$ =document id,  $q\_id_i$ =  $i$  th query, where  $i=1\dots5$ ,  $a\_id_j$ =  $j$  th answer option, where  $j=1\dots5$ .

Each query frame is defined in the system as  $(DOC, QT, OPTION_T)$  where,

$DOC$  = Give Document to be used for verifying answer options

$QT$  = Query Term, is a list of words after removing the stop words and interrogative word from the given question.

$OPTION_T$  = T-th answer option

Now, for each given answer option a score is calculated and the answer option with highest score is taken as correct answer for the given query. The algorithm *SelectAnswerOption* describes the option selection procedure.

## 5 Evaluation

The objective of the reading perspective evaluation is to offer information about the performance of a system “understanding” the meaning of each single document. The main measure used in this evaluation campaign is  $c@1$ , which is defined in equation 1.

$$c @ 1 = \frac{1}{n} (n_R + n_U \frac{n_R}{n}) \quad (1)$$

where,  $n_R$ : the number of correctly answered questions,  $n_U$ : number of unanswered questions and  $n$ : the total number of questions

Afterwards, these  $c@1$  scores can be aggregated at topic and global levels in order to obtain the following values:

- Median, average and standard deviation of  $c@1$  scores at test level, grouped by topic,
- Overall median, average and standard deviation of  $c@1$  values at test level.

**Table 1.** Algorithm SelectAnswerOption (Answer Set)

<b>Algorithm <i>SelectAnswerOption</i>(DOC, QT, OPTION<sub>T</sub>)</b>
<b>Step 1:</b> [ <i>Initialization</i> ] correct_option= $\infty$ // not answered
<b>Step 2:</b> [ <i>Calculate score for each sentence</i> ] For each sentence $S_i \in \text{DOC}$ and answer option $q_j \in Q$ Where, $j=1 \dots 5$ $A_{ji} = \text{AnswerScore}(S_i, QT, \text{OPTION})$ End For
<b>Step 3:</b> [ <i>Applying Matching Score(<math>M_{score}</math>)</i> ] For each answer option $AQ_j \in AQ$ $AQ_j = M_{score}(AQ_j)$ End For
<b>Step 4:</b> [ <i>Select the answer option</i> ] correct_option= index of maximum $AQ = \{ AQ_1, AQ_2, AQ_3, AQ_4, AQ_5 \}$
END

The median  $c@1$  has been provided under the consideration that it can be more informative at reading level than average values. This is because median is less affected by outliers than average, and therefore, it offers more information about the ability of a system to understand a text.

This approach allows us to evaluate systems in a similar way to the manner new language learners are graded.

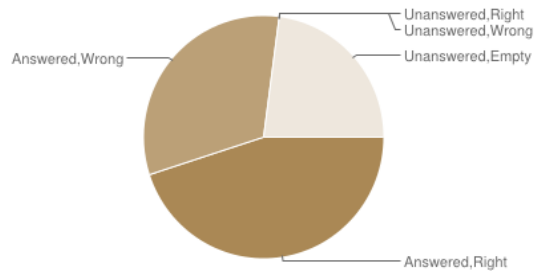
Thus, we can consider that a system passes a test from this evaluation perspective if it achieves a score equal or higher than 0.5. In the case of obtaining an overall average  $c@1$  higher than 0.5, we say that the system passes this evaluation perspective.

The dataset was composed of a total of 284 questions of which 240 are main questions and 44 are auxiliary questions. The difference between main and auxiliary questions resides in the presence of an inference. In fact an auxiliary question is just a duplicate of a main question minus the inference. The idea is that the simpler versions (auxiliary) could be added to main questions: if a system gets the difficult version wrong and the easy version right, it could be that it could not perform the required inference.

## 5.1 Evaluation on the main questions

### A) Evaluation at question-answering level

- Number of questions ANSWERED: **185**
- Number of questions UNANSWERED: **55**
  
- Number of questions ANSWERED with RIGHT candidate answer : **108**
- Number of questions ANSWERED with WRONG candidate answer : **77**
- Number of questions UNANSWERED with RIGHT candidate answer : **0**
- Number of questions UNANSWERED with WRONG candidate answer : **0**
- Number of questions UNANSWERED with EMPTY candidate : **55**



**Fig. 2:** Pie Chart Representation of Evaluation at QA level (main)

Accuracy (*answered with judgment=correct*) calculated over all questions:

Overall **accuracy** =  $108/240 = 0.45$

Proportion of answers correctly discarded:  $0/55 = 0.00$

**Table 2.** Overall c@1 per topic

Topic	n	n <sub>R</sub>	n <sub>U</sub>	c@1
AIDS	60	35	10	0.68
Music and Society	60	28	26	0.67
Climate Change	60	15	07	0.28
Alzheimer	60	30	12	0.60

**Overall c@1 measure** =  $(108+55(108/240))/240 = 0.55$

**B) Evaluation at reading-test level**

Median: 0.47 - Average: 0.46 - Standard Deviation: 0.21 -calculated over c@1 of all 16 reading tests

**Topic  $t\_id = '1'$  - Alzheimer**

Median: 0.70 - Average: 0.68 - Standard Deviation: 0.17 -calculated over the c@1 of the four reading tests

**Table 3.** Overall c@1 for Alzheimer

Reading ID(r_id)	n	n <sub>R</sub>	n <sub>U</sub>	c@1
1	15	10	03	0.80
2	15	11	02	0.83
3	15	08	02	0.60
4	15	06	03	0.48

**Topic  $t\_id = '2'$  - Music and society**

Median: 0.48 - Average: 0.47 - Standard Deviation: 0.12 -calculated over the c@1 of the four reading tests

**Table 4.** Overall c@1 for Music and society

Reading ID(r_id)	n	n <sub>R</sub>	n <sub>U</sub>	c@1
5	20	07	05	0.44
6	19	04	10	0.32
7	20	09	06	0.59
8	19	08	05	0.53

**Topic  $t\_id = '3'$  - Climate Change**

Median: 0.20 - Average: 0.22 - Standard Deviation: 0.11 -calculated over the c@1 of the four reading tests

**Table 5.** Overall c@1 for Climate Change

Reading ID(r_id)	n	n <sub>R</sub>	n <sub>U</sub>	c@1
09	18	06	02	0.37
10	18	03	03	0.19
11	18	02	01	0.12
12	20	04	01	0.21

**Topic  $t\_id = '4'$  - AIDS**

Median: 0.51 - Average: 0.48 - Standard Deviation: 0.20 -calculated over the  $c@1$  of the four reading tests

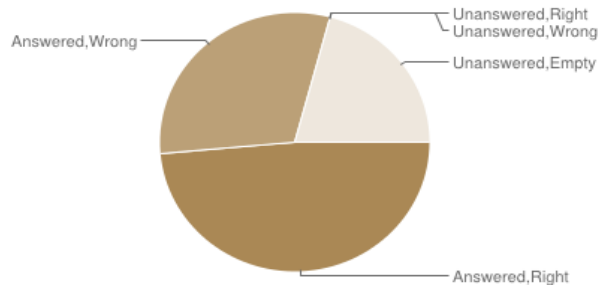
**Table 6.** Overall  $c@1$  for AIDS

Reading ID( $r\_id$ )	$n$	$n_R$	$n_U$	$c@1$
13	18	03	05	0.21
14	18	08	04	0.54
15	18	11	02	0.68
16	18	08	01	0.47

## 5.2 Evaluation on all questions (main + auxiliary)

### A) Evaluation at question-answering level

- Number of questions ANSWERED : **225**
- Number of questions UNANSWERED : **59**
  
- Number of questions ANSWERED with RIGHT candidate answer : **138**
- Number of questions ANSWERED with WRONG candidate answer : **87**
- Number of questions UNANSWERED with RIGHT candidate answer : **0**
- Number of questions UNANSWERED with WRONG candidate answer : **0**
- Number of questions UNANSWERED with EMPTY candidate : **59**



**Fig. 3:** Pie Chart Representation of Evaluation at QA level (main+auxiliary)

Accuracy (*answered with judgment=correct*) calculated over all questions:  
Overall accuracy =  $138/284 = 0.49$

Proportion of answers correctly discarded:  $0/59 = 0.00$

**Table 7.** Overall c@1 per topic

Topic	n	n <sub>R</sub>	n <sub>U</sub>	c@1
AIDS	60	35	10	0.68
Music and Society	60	42	28	0.73
Climate Change	60	20	09	0.30
Alzheimer	60	41	12	0.66

**Overall c@1 measure** =  $(138+59(138/284))/284 = 0.59$

**B) Evaluation at reading-test level**

Median: 0.62 - Average: 0.59 - Standard Deviation: 0.22 -calculated over c@1 of all 16 reading tests

**Topic t\_id = '1' - Alzheimer**

Median: 0.70 - Average: 0.68 - Standard Deviation: 0.17 -calculated over the c@1 of the four reading tests

**Table 8.** Overall c@1 for Alzheimer

Reading ID(r_id)	n	n <sub>R</sub>	n <sub>U</sub>	c@1
1	15	10	03	0.80
2	15	11	02	0.83
3	15	08	02	0.60
4	15	06	03	0.48

**Topic t\_id = '2' - Music and society**

Median: 0.72 - Average: 0.72 - Standard Deviation: 0.14 -calculated over the c@1 of the four reading tests

**Table 9.** Overall c@1 for Music and society

Reading ID(r_id)	n	n <sub>R</sub>	n <sub>U</sub>	c@1
5	20	10	06	0.65
6	19	07	10	0.56
7	20	13	07	0.88
8	19	12	05	0.80

**Topic  $t\_id = '3'$  - Climate Change**

*Median: 0.31 - Average: 0.30 - Standard Deviation: 0.10 -calculated over the  $c@1$  of the four reading tests*

**Table 10.** Overall  $c@1$  for Climate Change

Reading ID( $r\_id$ )	$n$	$n_R$	$n_U$	$c@1$
09	18	06	03	0.39
10	18	04	03	0.26
11	18	03	02	0.19
12	20	07	01	0.37

**Topic  $t\_id = '4'$  - AIDS**

*Median: 0.66 - Average: 0.65 - Standard Deviation: 0.18 -calculated over the  $c@1$  of the four reading tests*

**Table 11.** Overall  $c@1$  for AIDS

Reading ID( $r\_id$ )	$n$	$n_R$	$n_U$	$c@1$
13	18	06	05	0.43
14	18	10	04	0.68
15	18	14	02	0.86
16	18	11	01	0.65

## 6 Conclusion

The question answering system has been developed as part of the participation in the QA4MRE track as part of the CLEF 2013 evaluation campaign. The overall system has been evaluated using the evaluation metrics provided as part of the QA4MRE 2013 track. It has been observed from evaluation results that our proposed model works very well on the topics- “Aids”, “Music and Society” and “Alzheimer”. And the system performance decrease to handle “Climate change” documents and questions. As we have prepared the domain base knowledgebase for all the domains except Climate change domain. This may be one of the reasons for poor results of this domain. Hence it’s proved that domain knowledgebase has a strong effect on each of our system. But, the overall evaluation results are satisfactory in terms of  $c@1$ . Future works will be motivated towards improving the performance of the system and introducing domain knowledgebase for “Climate Change” domain.



**Acknowledgements.** We acknowledge the support of the support of the Department of Electronics and Information Technology (DeitY), Ministry of Communications & Information Technology (MCIT), Government of India funded project “Development of Cross Lingual Information Access (CLIA) System Phase II”.

## References

1. Anselmo Peñas, Pamela Forner, Richard Sutcliffe, Álvaro Rodrigo, Corina Forăscu, Iñaki Alegria, Danilo Giampiccolo, Nicolas Moreau, Petya Osenova.: Overview of ResPubliQA 2009: Question Answering Evaluation over European Legislation. In Working Notes for the CLEF 2009 Workshop, 30 September-2 October, 2009, Corfu, Greece.
2. Anselmo Peñas, Pamela Forner, Álvaro Rodrigo, Richard Sutcliffe, Corina Forăscu and Cristina Mota.: Overview of ResPubliQA 2010: Question Answering Evaluation over European Legislation. In Working Notes for the CLEF 2010 Workshop, Padua, Italy, 20-23 September 2010.
3. Anselmo Peñas, Eduard Hovy, Pamela Forner, Álvaro Rodrigo, Richard Sutcliffe, Corina Forascu, Caroline Sporleder. Overview of QA4MRE at CLEF 2011: Question Answering for Machine Reading Evaluation, Working Notes of CLEF 2011. (2011)
4. Peñas, A.,Rodrigo, Á. , Sama, V., Verdejo, F.: Overview of the answer validation exercise 2006. Working Notes of CLEF 2006. (2006)
5. Peñas, A., Rodrigo, Á, Verdejo, F.: Overview of the Answer Validation Exercise 2007. Working Notes of CLEF 2007. (2007)
6. Rodrigo, Á., Peñas, A., Verdejo, F.: Overview of the answer validation exercise 2008. Working Notes of CLEF 2008. (2008).
7. Partha Pakray, Pinaki Bhaskar, Santanu Pal, Dipankar Das, Sivaji Bandyopadhyay and Alexander Gelbukh: JU\_CSE\_TE: System Description QA@CLEF 2010 – ResPubliQA. CLEF 2010 Workshop on Multiple Language Question Answering (MLQA 2010).
8. Pakray, P., Gelbukh, A., Bandyopadhyay, S.: Answer Validation using Textual Entailment. 12th CILing, Lecture Notes in Computer Science, 2011, Volume 6609/2011, 353-364, DOI: 10.1007/978-3-642-19437-5\_29. (2011)
9. E. Briscoe, J. Carroll, and R. Watson.: The Second Release of the RASP System. In Proceedings of the COLING/ACL 2006 Interactive Presentation Sessions.
10. Marie-Catherine de Marneffe, Bill MacCartney, and Christopher D. Manning.: Generating Typed Dependency Parses from Phrase Structure Parses. In 5th International Conference on Language Resources and Evaluation (LREC) (2006)
11. Xuan-Hieu Phan.: CRFChunker: CRF English Phrase Chunker. PACLIC 2006. (2006)
12. P. Pakray, P. Bhaskar, S. Banerjee, B. Pal, A. Gelbukh, S. Bandyopadhyay: A Hybrid Question Answering System based on Information Retrieval and Answer Validation, In: the proceedings of Question Answering for Machine Reading Evaluation (QA4MRE) at CLEF 2011, Amsterdam. (2011)
13. M. W. Bilotti, B.Katz, and J.Lin.: What Works Better for Question Answering: Stemming or Morphological Query Expansion. In Proceedings of the SIGIR Workshop on Information Retrieval for Question Answering (IR4QA), pages 1–7,Sheffield, UK, July 29. (2004)

14. P. Bhaskar, P. Pakray, S. Banerjee, S. Banerjee, S. Bandyopadhyay, A. Gelbukh.: Question Answering System for QA4MRE@CLEF2012. In: the proceedings of Question Answering for Machine Reading Evaluation (QA4MRE) at CLEF 2012, Rome, Italy. (2012)