

Navigating Challenges and Crafting Solutions in Natural Language Processing for Effective Question-Answer Systems

Subodh Kant¹, Dr. Dinesh Yadav²

¹Research Scholar, Dept. of Computer Science

²Research Supervisor & Assistant Professor
Dept. of Computer Science
Radha Govind University, Ramgarh Jharkhand

ABSTRACT

A common problem in natural language processing is to make an automatic question-answering (QA) system that can give an appropriate answer to a question that is asked. This work gives an outline of the different approaches and techniques used to resolve this common question-answering issue. The main thought behind the QA method is to help people who want to learn more. A quick look at the different kinds of QA systems and the work that has been done so far is given in this study. It has been seen that the difference in words and meanings depending on the situation creates new problems for the question-answer system. A look at both standard and deep learning methods used to solve the study problem is attempted in order to give researchers a better idea of what they can do in this area. We suggest a structure for a question-and-answer system that uses a deep learning method. The study also talks about the system's flaws and things to think about when using it.

Keywords— Question Answer system, knowledge base, deep learning

I. INTRODUCTION

In the field of natural language processing (NLP), question-answering (QA) is a busy area of study. These days, thanks to technology like computers and the internet, we need to be able to answer user questions with the most relevant information from our knowledge base. Question-answering systems (QA systems) are the answer to this problem [1]. In order to answer a certain question, the QA system pulls data from a knowledge source. Wikipedia says, "Question answering (QA) is a field of computer science that combines information retrieval and natural language processing (NLP). Its goal is to create systems that can naturally answer questions asked by humans." The QA system is split into two main groups: open domain and closed domain. When you ask an open-domain question, it doesn't focus on any particular material.

On the other hand, when you ask a closed-domain question, it does focus on a specific source. The method looks at the problem and the situation to come up with one or more possible answers to the Question. In a closed-domain system, the information base is used to get the results.

On the other hand, web searches gave answers along with information that made the answers clear in an open-domain system. In 1999, the first papers in the TREC collections were used to test quality assurance tools that could answer true questions. A drawing example is often used to show information. In this group is the NOK method [2]. The system has to deal with a lot of problems, such as sorting questions into groups, coming up with the right queries, clearing up uncertainty, finding syntactic relatedness, and figuring out time relationships in complicated questions. Besides these problems, finding the perfect answer needs the right method for

extraction and confirmation. This paper gives a short outline of the work that has been done in the Question answering system area. The paper looks at a number of results from common machine learning methods and compares them. The study also talks about the newest ways that deep learning can be used to solve the problem of answering questions. The paper talks about the different issues that need to be fixed for the system to work properly.

The rest of the paper is organized like this: Question Answer System is explained in Section I, and linked work is shown in Section II. Question Answer System: Section III talks about how deep neural networks can be used for the QA system, Section IV talks about the knowledge base needs and how they affect the QA system, Section V shows the QA system's basic structure and block diagram, and Section VI talks about how it will be evaluated. Section VII talks about the QA method and Section VIII wraps up the study and points the way forward.

II. RELATED WORK

The QA tools are put into groups based on:

- Application domain
- User Question Types
- Source documents and analysis of users' questions
- Other Approaches

A. Application Domain

What kind of application area does the QA system have? It can be either closed (limited) or open (wide). In the General area, people can ask a wide range of questions, and the system will usually look for answers in a large collection of documents. Ontology and word knowledge were used by Kan [16] to answer. Users usually ask random questions, and the answers could be better.

To answer users' cause questions, the general area needs a big and varied information pool. The files from Wikipedia and newsgroups are enough to answer users' questions. Because the questions aren't very serious, the quality of the answers is low and depends on the person. In the Restricted domain, answers are taken from a collection of information that is specific to that domain. Because the collection is limited, the accuracy goes up. Domain theory and knowledge bases could be used to help choose the right answer. "There are many types of limited or closed domain QA systems in the literature, such as medical, patent, geographic, community-based, and time systems." If you put this whole domain together, you'd have an open-domain QA system. These systems meet the needs of advanced users who need precise answers to their questions. How happy users are based on how much they know about the subject. The answers are also very good.

B. User question type

QAs can be put into groups based on the kinds of questions people ask. Part A is fact-based questions like "which," "when," "what," "how," and "who," which need a short answer. These question-and-answer methods work well enough. These can be answered with a simple, easy-to-find database like Wikipedia without having to do a lot of complicated NLP processing. b) The answer to a list-type question is a list of facts or things. c) There is no clear answer to hypothetical questions. (d) An answer to a confirmation question is either yes or no, true or wrong. When you ask a causal question, you should give a detailed answer. It needs to be explained about the thing in Question. (f) Fuzzy questions that don't show what the user wants, so the intended answer needs to be clarified. g) Dialog questions make it hard to figure out what the user wants. h) Descriptive questions that need

you to find the word's meaning or description.

C. Source Document and Analysis of User Question

These papers are classed based on the research that has been done on them. The different groups are (1) morphological analysis: this type of analysis tries to break down words into their constituent morphemes and give each morpheme a class. (2) Syntactical analysis shows how words are put together grammatically in questions and source texts. (3) Based on the words used in a question, semantic analysis figures out what the Question might mean. (4) “The questions are analyzed in terms of their construction and meaning, and the papers are analyzed at the sentence or higher level.” (5) Expected answer type analysis needs answers that are based on the type of Question, and (6) questions need to be focused on to get the right answer.

D. Other Approaches

In addition to the methods already mentioned, the system can also be set apart by its language approach, statistical approach, or pattern-matching approach. For this kind of QA system, NLP methods and a knowledge base, also called a corpus, are used. Production rules, reasoning, patterns, semantic relatedness, and theory are used to look at the answer choices. POS tags are used to pre-process user questions, and important buzzwords are found that can be used to look for answers in the collection. In the 1960s, BASEBALL [5] and LUNAR [6] were NLP tools that let you ask questions of an organized database. They used methods to come up with standard forms that were then turned into questions.

Statistical approaches work with a lot of different kinds of data that come from the internet and online databases. To learn models that can be used in other areas, they use statistical methods such as SVM, EM, and Bayesian methods. The effects of these methods are better than those of their competitors.

The best statistical QA [7] system from IBM added a lot to the field. It used the maximum entropy method and a bag of word traits. This way, it should make things go faster and have fewer mistakes. Moschitti [8], Zhang et al. [9], Quarteroni et al. [10], and Suzuki et al. [11] all sorted questions and answers. All of them used SVM text models. Wei et al. [13] and MKQA [12] both used a changed Bayesian Classifier. In their studies, the Modified Bayesian Classifier method did a better job than the base Bayesian method.

Other rival methods need a lot of complicated processing, but the pattern-matching method only needs text patterns. QA tools like these are being used a lot these days because they can figure out text trends on their own. They didn't use words like a processor, a named-entity recognizer, a dictionary, WordNet, or anything else that is hard to understand to look for answers in text. The answer to a question can be found by looking at how similar their mirror patterns are. These patterns are made up of common sentences that mean certain things. This method works really well.

Named entity tagger was used by Greenwood et al. [15] to find trends in text, while support and confidence measures were used by Zhang et al. [14]. In order to get the right results, Cui et al. [17] used the Profile Hidden Markov Model and the Bigram model. Saxena et al. [16] used pattern matching with word expansion. Gunawerdena et al. [19] worked on it for automatic closed-domain FAQ systems. Sneiders [18] and Unger et al. [20] used the template-based method to help people answer questions.

To answer a trivia question, you use a linguistic method. It takes a deep understanding of how words work together. These systems are very good at getting answers from their information base, but they need help to handle different kinds of data. These systems need help growing because, for every new idea that is added to the knowledge base, new rules need to be made. A statistical method is often used for both factoids and other types of questions. A basic understanding of the book is enough. They can be used on a large scale because they use

statistical methods and follow guided methods for making models. Because of this, they can be used with different kinds of information. This method works for all kinds of questions, like factoids, non-factoids, meanings, names, and more. Scalability is an issue because fewer patterns need to be learned for each idea. This method works best for data sets that are small to medium in size.

Looking at things as they are now, more work needs to be done to combine language, statistical, and pattern-based methods in a way that meets the needs of all kinds of people. Most QA systems learn question answers in a low-dimensional space and choose the right answer by looking for similar features.

III. QA SYSTEMS AND DEEP NEURAL NETWORKS

In the past few years, deep neural networks have shown promise as a way to solve many AI issues. Quality assurance has also been used to create dialog systems [1] and chatbots [2] that are meant to talk like real people. But with the latest progress in deep learning, neural network models look like they could be useful for QA. DeepQA, a highly parallel software design that IBM Watson used to look at natural language material in both the clue and its library to find replies that were related to the clue. The answer is then ranked by how relevant it is, with proof for each answer. Sophia, the figure is set up to give pre-answers to a certain set of questions or sentences. Blockchain technology is used to process both the inputs and the replies. The data is then shared in a cloud network that uses deep learning.

There is a smaller learning process in these systems, but they still need a lot of training. Supporting recurrent neural networks (RNNs) with GRU and LSTM units lets them handle the longer texts needed for quality assurance. For his question-answering system, Iyyer [21] used a recurrent neural network. “Weston [22], Hochreiter [24], and Tan [23] suggested a model for a QA system memory network that uses LSTM.” For modeling questions and answer choices, Kalchbrenner [25] and Feng [26] used a model based on a convolution neural network (CNN). Authors [27], [28], [29], and [30] have used different feature-based learning and neural network techniques.

To use semantic encoding to match question-and-answer pairs, Yu et al. [31] used distributed representation learning with logistic regression. To solve the data QA system, Severyn and Moschitti [32] suggested a new deep CNN method. A measure of questions and replies was made with several convolution layers. Using frequency statistics, the two measures are compared to see what they have in common. “A three-layer, stacked bidirectional long short-term memory (LSTM) network model was used by Wang and Nyberg [33] to match the question to the answer.” Wenhui et al. [34] used attention-based self-matching networks to find the answers in the text. Deep belief networks were used by Wang et al. [35] to get results from cQA and topic datasets. A study by Ying et al. [36] [37] suggests using KABLSTM, an information-aware Attentive Bidirectional Long Short-Term Memory, to help people learn QA words better by using external information from knowledge graphs (KG). The WikiQA and TREC QA files were used to test the system. Atsushi [38] used an encoder-decoder deep neural model to try to close the vocabulary gap between searches and documents. For a multidomain FAQ system, the model learns topics from replies to questions.

The author [39] has looked into how useful it is to use modules to add background information to deep networks. They showed that this makes the system better at learning and can also be used to get information from taught deep networks. Using trust rules, they came up with a method for adding hierarchical information to deep networks while they were being trained.

The most recent models get linguistic data from external sources like WordNet and use complicated language

tools to do so. Some of the things that the system has to deal with are word order, question length, synonymy, polysemy, and lack of data.

IV. REQUIREMENTS OF KNOWLEDGE BASE

It uses three kinds of knowledge bases: a text corpus, a knowledge base corpus, and a mixed corpus.

- **Text Corpus:** CLEF and TREC are questions about facts. The TREC questions come from search engine logs, and the answers have to be taken from texts. You can get sets of question-and-answer pairs for free. There are no descriptions in these records to back up the answer given. The new file has news from journals that are important to the time. Encyclopedia knowledge doesn't matter what time it is. A lot of people use the Stanford Question Answering Dataset (SQuAD) [37] these days for deep learning. "A collection of more than 100,000 question-answer pairs about more than 500 Wikipedia pages that were gathered by the public." There are 20,120 question-answer pairs in the Customer Questions & Answers part of Amazon Product records. "There are 31,000 question-answer pairs in Microsoft Community Questions & Answers. One more set of study QA is Yahoo! Answers."
- **Collections of knowledge bases:** There are about 200 text questions in the QALD files that are used with DBpedia. That which is sent back are DBpedia URIs. There are few questions, so it's hard to use to learn. It's a set of questions and answers on Freebase that were made with Google Suggest API searches. They are sorted by Amazon Mechanical Turk [38]. That's what the Freebase URIs mean. For practice, it has 3778 cases, and for the test, it has 2032 cases. In the AAAS Project, there are 775 mixed-type questions (MCQ). You can answer 26 questions with a picture or a number, and there are 749 questions about life sciences, physical sciences, earth sciences, and what science is all about.
- **Hybrid corpus:** QALD has been working on making a hybrid QA since 2014. In order to answer questions, you need to use both DBpedia and Wikipedia triples. As of now, it has about 150 question-answer pairs to test a mixed question-answering system.

So far, different kinds of QA tools have been made using different kinds of corpora. As the size of the information collection grows, so does the accuracy of the system. Experiments with standard IR systems using texts between 10 and 100GB [40] show that the size of the corpus has a direct effect on how well the system works. Authors [41] have shown that the QA system's speed does get better up to 400–500GB, but then it hits an asymptote and starts to go down. "More study needs to be done in this area to find out if the observed asymptote is a real effect or just a problem with the system being used with the current scoring method." The answers given by the system point to a possible flaw in the way the QA system evaluates. That problem would go away if QA systems used deep learning methods, which would also make the systems more accurate.

V. METHODOLOGY AND PROPOSED SYSTEM

The QA system is made up of three main parts: question analysis, which includes question parsing, question classification, and query formulation; document analysis, which includes extracting candidate documents; and answer extraction, which finds good candidate answers and ranks them to find the best one [1]. It sorts questions into groups based on the type of Question they are, which is an important part of QA systems. AI, natural language processing, data analysis, pattern matching, and finding information are just some of the fields that are being used together in new studies. Figure 1 shows an example of a QA system's block style.

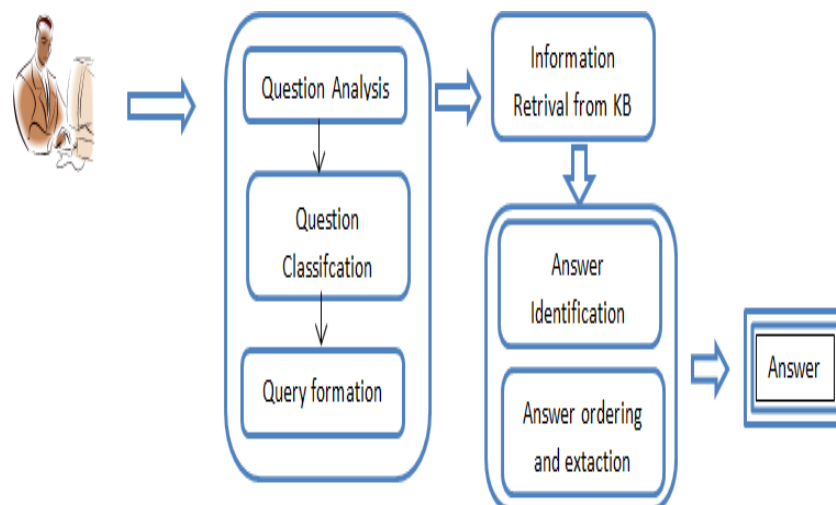


Figure 1: Basic QA system

There are three parts to how the system works: the Query Processing Phase, the Document Processing Phase, and the Answer Processing Phase. The user's Question is sent to the Query processing step as an input. The program figures out what the question is trying to do and puts it into one of several set question-type groups. This helps you figure out the pattern of answers you should get. The classification information is added to the Question to change the way the query is written. The information search engine is given the written Question. When a question is sent to the database, the machine pulls out the answers. In an open domain, this knowledge base can be the internet. In a closed domain, it can be a knowledge base that has already been set up. Based on the buzzwords in the query, the system pulls out the most likely options for the query that refers to the input question. The relevant papers that are found are sorted and shortened into outlines that should have the answer. The order of these descriptions is based on how useful they are. In order to get only the word or phrase that solves the Question, a set of rules is set up. Question processing is the part of the system that figures out what the Question is about, what kind of answer is expected, and how to rewrite the Question into multiple questions that are semantically similar. Asking the same question in different ways that have similar meanings is called query expansion, and it helps the information search system remember things better. Question asking relies on information retrieval (IR) system memory being very important. This is because if there are no right answers in a document, no further processing can be done to find an answer [3]. The last part of a question-answering system is answer extraction. This is what makes question-answering question systems different from text search systems in the normal sense. Answer extraction technology has a big impact on the question-answering system and determines the outcome.

The results are looked at based on the following:

- Relevance: The answer should have something to do with the Question
- Correctness: the answer should be based on facts
- Conciseness: the answer shouldn't have any unnecessary or useless details.
- Fullness: the answer should be whole; a partial answer shouldn't get full credit.
- Reasonability: the answer should make sense to the person who asked the question, so it's easy for them to read.
- Validation: For validation, the answer should have enough information to help the reader understand

why it was chosen as the answer to the Question.

The suggested method will use the Stanford Question Answering Dataset (SQuAD1.1), which is made up of questions that regular people have asked about a group of Wikipedia articles. Each Question's answer is a span of text from the reading section that goes with it. There are three main parts to the suggested structure. Figure 2 shows the focus layer, the decoder layer, and the embedding layer.

- **Layer for embedding:** The model's training collection is made up of questions and contexts that go with them. "These two things can be split up into separate words, which can then be turned into Word Embedding using a pre-trained vector like GloVectors." Instead of using a single hot vector for each word, word embeddings take into account the environment around them.
- **Encoder layer:** RNN uses a GRU/LSTM that works in both directions. It's the job of this encoder layer to remember what words came before and after it. It can also be used to make question vectors. A set of secret vectors is made, and then they are joined together.
- **Attention layer:** This helps you figure out the Question and what it means. Attention is found by taking the dot product of the question vector and the answer vector. If you put a softmax over the product, it will always add up to 1. Find the sum of the interest vector and the question vector.

Find possible answer vectors and rank them by how relevant they are. As an answer, the vector with the highest confidence value is chosen.

If the right answer can't be found, the information base needs to be improved. The search engine is asked to find the document, and the documents found must be used to train the model that was made in earlier steps. This would help with learning little by little.

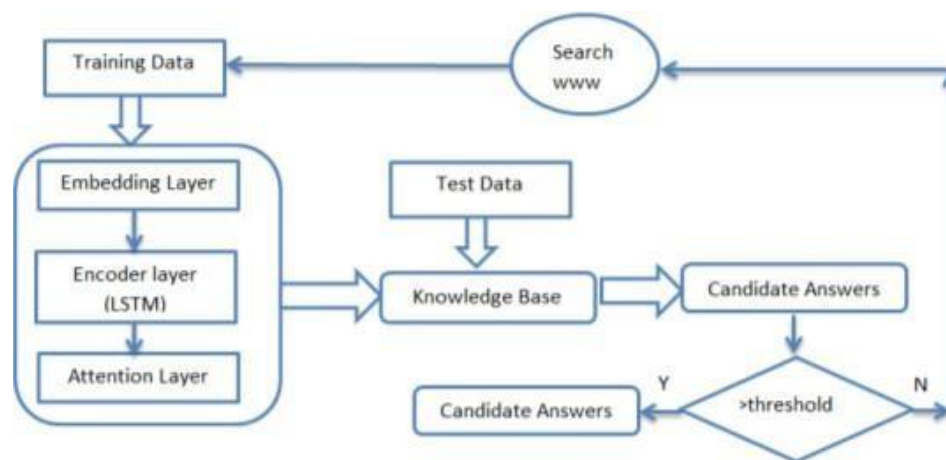


Figure 2: Block diagram of the proposed system

VI. EVALUATION MEASURES

The success of a QA system can be judged by its precision, recall, and F-measure. "In a QA system, precision (P), which is also called positive prediction value, shows how close the answer is to the question that was asked." How close a mapped answer is to a real right answer is called its accuracy. On the other hand, recall shows what percentage of relevant answers have been found out of all the relevant answers that were found. Let

P stand for precision, RA for a useful answer, and RT for the results that were retrieved. Then, precision P is found by:

Recall for a QA system is defined as

$$R = \frac{|{RA} \cap {RT}|}{|{RA}|}$$

The F1 score of test accuracy is given by

$$F1score = 2 * \frac{(P * R)}{(P + R)}$$

VII. DISCUSSION

A lot of work goes into making a good question-and-answer system, but there are still some things that need to be talked about. There is a word gap between the Question and its meaning, which makes it hard to ask a Knowledge Base with normal language. The main issues are how to match a question with the knowledge base triple and how to use its parsing method to find the named entities and their connections. There is a logical link between the Question and the possible solutions that people think they know. Because there are few word embeddings, this importance needs to be modeled and measured correctly. The type and strength of the information base are the only things that determine how proper and full the replies are. The knowledge base must be accurate, full for at least a certain area, and up to date with the most recent information. This collection needs to be updated all the time and added to the learning model so that the system is more accurate.

The information base and model that are built must allow for incremental learning. Ensemble learning methods can be used to look into the features of a question and appropriate replies in both the feature space and the sample space, which is also called the knowledge base. It costs a lot to name questions with the right replies, and there is also the problem of idea shift. Active learning should be added to a knowledge base that doesn't need set labels and can be labeled on demand. Concept drift can be stopped by changing the benchmark for appropriate replies on the fly when an old idea quickly changes into a new one. It still needs to come up with short, detailed answers to questions and come up with answers based on the type of Question and how it is graded.

VIII. CONCLUSION AND FUTURE SCOPE

The paper looks at different question-and-answer systems that use both standard and deep learning methods. Relevance, accuracy, brevity, completeness, logic, and proof are some of the most important things that determine the quality of a QA system. Machine learning the old way Statistical methods like SVM, EM, and Bayesian are often used in QA systems. A lot of the time, ontology and word knowledge are used to answer questions in community-based, medical, legal, geospatial, and open and closed application area systems. A lot of question-response systems that use deep learning choose to use LSTM and GRU along with RNN. A small number of writers have also used CNN for the QA method. Attention-based RNN and LSTM are used in the suggested system. The knowledge base is always being updated with new information, which makes the method easier to learn. The method tries to get around the problem of needing more information. Active learning could

be a way to do work in the future.

IX. REFERENCES

- [1]. L. Hirschman and R. Gaizauskas, "Natural language question answering: the view from here," *Natural Language Engineering*, vol. 7, no. 4, pp. 275-300, 2001
- [2]. Green BF, Wolf AK, Chomsky C, and Laughery K. "Baseball: An automatic question answerer", In Proceedings of Western Computing Conference, Vol. 19, 1961, pp. 219–224.
- [3]. Woods W. "Progress in Natural Language Understanding - An Application to Lunar Geology", In Proceedings of AFIPS Conference, Vol. 42, 1973, pp. 441–450.
- [4]. Ittycheriah A, Franz M, Zhu WJ, Ratnaparkhi A, and Mammone RJ, "IBM's statistical question answering system" In Proceedings of the Text Retrieval Conference TREC-9, 2000.
- [5]. Moschitti A. "Answer filtering via text categorization in question answering systems," In Proceedings of the 15th IEEE International Conference on Tools with Artificial Intelligence, 2003, pp. 241- 248.
- [6]. Zhang K, Zhao J. "A Chinese question answering system with question classification and answer clustering" in Proceedings of IEEE International Conference on Fuzzy Systems and Knowledge Discovery (FSKD), Vol.6, 2010, pp. 2692-2696.
- [7]. Han L, Yu ZT, Qiu YX, Meng XY, Guo JY and Si ST. "Research on passage retrieval using domain knowledge in Chinese question answering system," In Proceedings of IEEE International Conference on Machine Learning and Cybernetics, Vol. 5, 2008, pp. 2603-2606.
- [8]. Suzuki J, Sasaki Y, Maeda E. "SVM answer selection for open domain question answering," In Proceedings of 19th International Conference on Computational linguistics, COLING'02, Vol. 1, 2002, pp. 1-7.
- [9]. Fu J, Xu J, and Jia K. "Domain ontology-based automatic question answering," In IEEE International Conference on Computer Engineering and Technology, Vol. 2, 2009, pp. 346-349.
- [10]. Ying-wei L, Zheng-tao Y, Xiang-yan M, Wen-gang C, Cun-li M. "Question Classification Based on Incremental Modified Bayes", In Proceedings of IEEE Second International Conference on Future Generation Communication and Networking, Vol. 2, 2008, pp. 149-152.
- [11]. Zhang D and Lee WS. "Web-based pattern mining and matching approach to question answering," Proceedings of the 11th Text Retrieval Conference, 2002.
- [12]. Greenwood M. and Gaizauskas R. "Using a Named Entity Tagger to Generalise Surface Matching Text Patterns for Question Answering," In Proceedings of the Workshop on Natural Language Processing for Question Answering (EACL03), 2003, pp. 29-34.
- [13]. Saxena AK, Sambhu GV, Kaushik S, and Subramaniam LV, "Ibmirl system for question answering using pattern matching, semantic type, and semantic category recognition," Proceedings of the TREC, Vol. 2007, 2007.
- [14]. Cui H, Kan MY, and Chua TS. "Soft pattern matching models for definitional question answering," In ACM Transactions on Information Systems (TOIS), Vol. 25(2): 8, 2007.
- [15]. Schneider E. "Automated question answering using question templates that cover the conceptual model of the database," In Natural Language Processing and Information Systems, Springer Berlin Heidelberg, 2002, pp. 235-239.
- [16]. Gunawardena T, Lokuhetti M, Pathirana N, Ragel R, Deegalla S, "An automatic answering system with template matching for natural language questions" In Proceedings of 5th IEEE International Conference on Information and Automation for Sustainability (ICIAFs), 2010, pp. 353-358.
- [17]. Unger C, Bühmann L, Lehmann J, NgongaNgomo AC, Gerber D, and Cimiano P, "Template-based question answering over RDF data," Proceedings of the ACM 21st international conference on World Wide Web, 2012, pp. 639-648.
- [18]. M Iyyer, JL Boyd-Graber, LMB Claudino, R Socher, IH Daume, "A Neural Network for Factoid Question Answering over Paragraphs," In Conference on Empirical Methods on Natural Language Processing, 2014.
- [19]. J Weston, S Chopra, A Bordes, "Memory networks. " arXiv preprint arXiv:1410.3916, 2014.
- [20]. M Tan, B Xiang, B Zhou, "LSTM-based Deep Learning Models for non-factoid answer selection," arXiv preprint arXiv:1511.04108, 2015.1
- [21]. S Hochreiter, J Schmidhuber, "Long short-term memory," *Neural computation* 9.8: 1735-1780, 1997
N Kalchbrenner, E Grefenstette, P Blunsom, "A convolutional neural network for modeling sentences," In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics. June, 2014.