# Identifying Semantically Similar questions using NLP techniques and Linked Data Principles

**Anirban Bhattacharjee**

**A Dissertation**

Presented to the University of Dublin, Trinity College

in partial fulfilment of the requirements for the degree of

**Master of Science in Computer Science (Intelligent Systems)**

Supervisor: Professor Declan O'Sullivan

August 2019

# Declaration

I, the undersigned, declare that this work has not previously been submitted as an exercise for a degree at this, or any other University, and that unless otherwise stated, is my own work.

_____

Anirban Bhattacharjee

August 13, 2019

# Permission to Lend and/or Copy

I, the undersigned, agree that Trinity College Library may lend or copy this dissertation upon request.

_____

Anirban Bhattacharjee

August 13, 2019

# Acknowledgments

There are many people who have influenced my life and supported me in every decision that I made in it. They have all, in their unique ways, made this dissertation a reality.

Firstly, I would like to thank Professor Declan O'Sullivan for his help and guidance throughout the duration of this dissertation. The insightful comments and feedback I received in our discussions helped me in not losing sight of the big picture while struggling with the details of this dissertation. He has helped shape its direction from when it was merely a rough idea in my head, to the approach that it is now.

I would like to mention my parents who deserve all the credits for supporting me to attend college and I can not thank them enough for being by my side for everything that I have ever wished to experience in my life.

Special thanks to Dr. Fabrizio Orlandi for helping me orchestrate the evaluations and his feedback regarding the same. I would also like to thank my fellow course mates and all the evaluation participants who supported me in this research.

Finally, I would like to remember Professor Séamus Lawless who has influenced me greatly during my course and will live on like a great inspiration for me throughout my life.

ANIRBAN BHATTACHARJEE

*University of Dublin, Trinity College*
*August 2019*

# Identifying Semantically Similar questions using NLP techniques and Linked Data Principles

Anirban Bhattacharjee, Master of Science in Computer Science

University of Dublin, Trinity College, 2019

Supervisor: Professor Declan O'Sullivan

In Community Question Answering (CQA) sites, despite active participation, a significant amount of questions on such sites remain unanswered due to a lot of reasons such as the question being poorly formed/worded, unavailability of any answerer or the increased inflow of questions in the same area which disinterests an answerer to answer the same question or having to redirect them multiple times to already answered questions. This research is an attempt to study if unstructured data is converted to structured data using state-of-the art natural language processing (NLP) techniques and Linked Data technologies, to what extent it could help a user in identifying semantically similar questions. One of the most contested themes in Computer Science is the ability to automatically map natural language semantics into programming languages. This research work is distinguished from other studies as we approach the problem from an ontology centred view and the idea of knowledge reuse forms the notion of this work. We evaluate our approach and open discussions on new ways to evaluate the identification of semantically similar questions. The key findings of this research demonstrate that using NLP techniques and Linked Data principles identification of semantically similar questions is viable. The proposed approach has a small but significant impact which can be leveraged for designing data models for the task of finding semantically similar questions.

# Contents

# List of Tables

# List of Figures

# Chapter 1

# Introduction

Community Question Answering(CQA) sites is a platform for submitting questions in natural language, a forum for users to respond to questions in natural language and a community built around this exchange of questions and answers. CQA forums have emerged as an effective means of information and knowledge exchange on the Web. Over the years these CQA sites have accumulated a large number of user questions and answers and their related knowledge which forms the basis of rich knowledge repositories. CQA sites differ from the traditional web search in the following ways:

1. Instead of a list of documents from the web, users receive response in natural language to their questions.

2. Based on personal experience and expertise, answers to opinion or decision based questions can be sought in CQA sites, which cannot be retrieved through traditional web search.

The different online information seeking services and their categorization are listed in Figure 1.1 presented in the work of Shah et al. [2].

| Service | Method of presenting information need | Method of presenting relevant information | Participation | Category |
|---|---|---|---|---|
| Google Web Search | Keywords | Rank list | Single-user | Web search |
| Amazon | Keywords or browsing | Rank list or recommendations | Indirect feedback from the community | Recommender system |
| Ask | Questions | Document snippets | Single-user | Automated Q&A |
| Yahoo! Answers | Questions | Answers | Direct involvement of the community | Social Q&A |

Figure 1.1: Example of different information seeking services and their categorization [2]

Shah et al. [2] states that CQA sites requires a forum for public interaction, an interface to put forward a natural language question representing an information need as opposed to keyword searching and a community in which transaction of information is based on participation level.

There are three broad categories of audience in these forums [11]:

1

1. Users who only ask questions.

2. Users who only answers questions.

3. Users who ask as well as answer questions.

The following Figure 1.2 illustrates a page from the CQA site.



Figure 1.2: Example of a page from Quora CQA site [1]

CQA sites although sounds very exciting involves the risk of providing information to the users which are of significantly low quality[12].There are also instances when the content in these sites are abusive where both questions and answers are filled with either promotional content or offensive content.

With the increasing popularity of CQA sites, it brings forward few challenges. According to relevant literature, most of the user generated content on CQA sites are redundant[13] and noisy [14][15].

## 1.1 Unanswered Questions in CQA sites

The usefulness of a CQA site depends on how effectively it can get answers for the questions posted by the users. In an anlaysis conducted by Shtok et al.[16] on Yahoo! Answers which is one of the

---

[1]https://www.quora.com/How-many-operating-systems-can-be-installed-in-one-computer

very first CQA sites on the web, they found out that 15% of the questions remains unanswered which leaves the asker unsatisfied [16]. In order to observe the state of the questions after two days from the date of their posting in CQA sites, Li et al [17] randomly tracked 3000 recently posted questions in two CQA sites namely Yahoo! Answers and Baidu Zhidao. It was found that only 17.6% of questions in Yahoo! Answers received an answer within 48 hours. In Baidu Zhidao, only 22.7% questions are resolved [17]. These observations indicated that two of the major CQA sites could not solve user questions efficiently. Most of the CQA sites in this manner suffers from the problem of "Answer Starvation"- a phenomenon where the users enjoy asking questions instead of answering them [18] which results in a large number of unanswered questions. According to [19] 12% of the questions in CQA sites are unanswered because of a duplicate question and 21.75% of questions are unanswered as it fails to attract an expert member. Research in this area [19] also signifies that duration of an unanswered question differs from site to site. Figure 1.3 is an example of unanswered questions in CQA sites for over 24 months. The duration of a question remaining unanswered in Yahoo!Answers is different than the duration of an unanswered question in Quora. Users registered in these sites may answer only a specific number of questions daily depending on their interests, determination and reputation. Although, there are active participants in CQA sites, a significant amount of questions never receives an answer due to a variety of reasons such as uninteresting questions, un-understandable questions, increased volume and variety of incoming questions which creates a situation for answerers to overlook questions and the absence of a robust mechanism to redirect questions to potential answerers.



Figure 1.3: Example of Unanswered Questions for over two years.

3

Estimating user expertise on a topic and re-directing the question to knowledge experts via inference engines as shown in Figure 1.4 could be done in a variety of ways. Research regarding this usually involves link analysis and topic modelling techniques. In relevant research literature, Jurczyk et al. [20] studied and proposed a graph modelling approach using link analysis to calculate on each expected topics the authoritative scores of users. Qu et al.[21] adopted the Probabilistic Latent Semantic Analysis (PLSA) model to capture user interactions and based on the user's history deduced correlation between answer and question, thereby improving the quality of recommending questions.



Figure 1.4: A flow chart of expert systems in CQA sites

For routing questions to most potential answerer a significant amount of research has been conducted [17] [22] [23] however the duration or the availability of the answerer to answer a recommended question is not defined. There are ways of generating answers from knowledge sources such as Wikipedia which has been discussed in relevant literature [23]. Although this approach being faster would need robust infrastructure for web scraping, data extraction and text summarization which are huge opportunities for research and are wide research fields. In this work, we research on a novel approach by using previously asked questions in CQA sites to identify semantically similar questions.

Dissatisfaction of the users could be combated by identifying semantically similar questions. It can also be used for recommending potential questions to users who are searching for an answer to a specific question and are not aware that such a question has already been answered in the past.

## 1.2  Motivation

Tim Berners-Lee, the prime inventor of the web and the semantic web initiator, the one who still leads the World Wide Web Consortium (W3C) is the first individual to publicly bring forward the need for semantics in the web at a World Wide Web (WWW) conference in the year 1994 [2]. He later proposed a road map for its architecture at a very high level in the year 1998 [3]. The semantic web vision was

---

[2]http://www.w3.org/Talks/WWW94Tim/
[3]https://www.w3.org/DesignIssues/Semantic.html

primarily popularized when an article was published in the year 2001 [24]. This article specifically attracted a lot of researchers to realize this vision. Web was mostly a large graph of web pages back in the year 2001, however, when the semantic web vision was proposed, the need to represent semantic information aroused. Similar to the link-ability of documents on the web, the semantic web's most distinguished feature is in its ability to reference specific pieces of data published by different users on the web. This web of data in contrast to the traditional web of documents can be understood and accessed by both machines as well as humans. The standards and practices suggested by the W3C to publish data on the web is referred to as the *"Linked Data"*. Machines usually understand information which is provided in the form of linked data as data models are usually specified by ontologies which enables automated reasoning. Connecting distributed pieces of information related to different kinds of datasets are favoured by linked data [25].

Our intuition to use linked data to identify semantically similar questions stems from the opportunity to make data interconnected by using the collective wisdom of the queries of real world people in CQA sites. Rather than considering every word as a random string of characters, considering words as real things in the world and connecting those things to other things to aid an user in finding questions for which they are actually looking for an answer is the reason for attempting this approach.

Semantic Web is an extension of the World Wide Web in which the meaning(semantics) of information and services is defined, making it possible to understand and cater to the needs and requests of people and machines to use the web content. "Linked data(LD) is at the heart of what Semantic Web is all about : large scale integration of, and reasoning on, data on the Web" [4]. Semantic Web is about making links such that a a machine or a human user can explore the web of data. LD refers to a set of best practices for sharing, connecting, and exposing data on the Web [26]. LD gives the power to re-use information in unexpected ways which can add value to a web of data. The information scale of unstructured data has the potential to make question answering systems over LD useful. Also, using structured resources to support information extraction for question answering or in this context finding similar questions will support the enrichment and expansion of of structured datasets from unstructured data [27].

## 1.3 Research Question

To what extent is it possible to convert natural language questions into structured data (RDF graphs) and thereby querying RDF graphs using SPARQL protocol and RDF Query Language (SPARQL) to get semantically similar questions?

---

[4]https://www.w3.org/standards/semanticweb/data

## 1.4 Research Objectives

The specific actions to be taken in order to answer the research question has been mentioned below.

- To research on an approach to analyze questions having the same meaning although worded differently using Natural Language Processing techniques and Linked Data principles.

- To transform unstructured data to structured data for better reasoning capabilities.

- To build an application to recommend semantically similar questions in response to user asked questions using structured data.

- To evaluate such an application by primarily involving human participants to identify the viability of the approach.

## 1.5 Contributions

In this dissertation - We propose a unique approach to identify semantically similar questions in CQA sites. In our approach we populate a knowledge base with questions posted in CQA sites. We create a Resource Description Framework (RDF) graph using the semantic triples generated from a natural language sentence. The nodes of the graph are either subjects or objects and the edges of the graph represent the relationship between two nodes. In this manner, unstructured natural language questions are modelled to structured sources of information in RDF data model .The triple store / RDF repository / knowledge base is then queried using SPARQL Protocol and RDF Query Language(SPARQL) to identify semantically similar questions based on the triples associated with each question. In the event of an unknown question posted by the users, the triples generated from the text of the question are used to populate the knowledge base and stored for future querying.

## 1.6 Organization of Dissertation

The rest of the dissertation is organized as follows: Chapter 2, discusses the relevant background including the state-of-the-art techniques for identifying semantically similar questions. Chapter 3, discusses our attempt to design a proof of concept application using state-of-the-art NLP techniques, RDF and SPARQL queries for identifying semantically similar questions. In Chapter 4, we discuss our implementation and the problems involved during the implementation phase. Chapter 5 states the evaluation and is a discussion on the attempt to find out how good is this approach and if this approach can be extended to make it even better than the approach it is now. In Chapter 6, we provide concluding remarks on the dissertation and an analysis of how successful was this dissertation. It ends with some final remarks, the current limitations and areas for future work.

# Chapter 2

# Background

In this chapter we discuss the history of question answering systems in general and community question answering sites in particular. The identification of similar questions from the archives of these community question answering sites has been discussed. This draws on to the state-of-the-art approaches to identify semantically similar text and the work done in relevant literature in predicting user's satisfaction depending on recommended questions. We discuss the current techniques used to extract knowledge from the CQA sites and most importantly finding user's intention for querying a CQA site. Quora's graph structure has been briefly discussed followed by a brief discussion on the current challenges in Natural language processing and linked data in the field of question answering.

## 2.1 Forerunners

Question Answering (QA) Systems or frameworks have evolved, changed and transformed a lot in the recent decades. It has kept pace with the progress of natural language processing. Natural language processing(NLP) progressed quite a lot in the 70s and the 80s, bringing about bigger undertakings in the field of question answering. One such project was the formation of Unix Consultant which was created by Robert Wilensky at U.C. Berkeley in the late 1980s. In 1990s, Boriz Katz along with his associate at Info Lab Group at the MIT Computer Science and Artificial Intelligence Laboratory built a framework called START natural language system. It was the world's first online QA system which has been running non stop since December 1993. The later part of the decade saw a massive surge in the number of search engines and by mid 2000, multiple question answering sites were built based on these search engines.

Various open source text processing projects were undertaken, few of which are Ephyra, YodaQA, PTStemmer, GATE and OpenNLP. In 2003, the Defense Advanced Research Projects Agency (DARPA) financed a Personal Assistant that Learns (PAL) program from which a spinoff named SRI

International, Siri, Inc. was formed, this venture was later acquired by Apple in the year 2010. SIRI was later incorporated in iOS (version 4), an operating system which is manufactured by Apple Inc and used specifically in mobile devices. Around the same time, research associates at IBM were developing Watson [28] to compete with the winners at Jeopardy! - an American game show. IBMs Watson ultimately triumphed over two human champions in the year 2011, a project which started in the year 2007 and opened gateways for further research in the field of deep analytics, natural language understanding and open domain question answering systems.

A few text-based QA systems were developed later, such as QANUS [29]. These QA systems retrieve answers from textual sources and the web. Structured QA frameworks like AquaLog [30] - which is an ontology - driven QA system and they retrieve the answers from pre-defined ontologies and organized data sources.

## 2.2   Community Question Answering Sites

The very first community question answering site was built by a South Korean company Naver Corporation named Knowledge iN in the year 2002. Answerbag was one such english language CQA site which was started in the year 2003 and later dissolved in the year 2015. However, Yahoo! Answers introduced by Yahoo! in the year 2005 gained much popularity and influenced CQA platforms in becoming institutionalized. Few closed domain CQA sites such as StackOverflow and StackExchange has gained immense popularity as these sites contains rich information about specific domains which is not found elsewhere in the web. CQA sites such as Quora gains active support from its community members and has become popular within a short period of time. CQA sites have in the recent years integrated with social network sites such as Facebook and has created compatible widgets to allow their users to access topics from other sites. There are various restrictions implemented in CQA sites ranging from inexperienced members not being able to choose best answers to restricting members from commenting on a specific number of questions per day.

## 2.3   Identifying questions similar to user's question from archives

Usually most of the renowned CQA sites like Quora, Yahoo! Answers maintains a huge archive of solved questions. Every-time a user asks a new question, these CQA sites tries to find this new question coming from the user in their archives in order to avoid delay in response. The concept of searching for a similar question in CQA sites is to some extent similar to passage retrieval in QA systems [31]. However, it is to be noted that matching question to question is much more difficult and constricted than matching question to passage [31]. There is a significant amount of work done in the areas of passage retrieval, for example by using dependency relations [32] which is a state of the art

technique. We have also incorporated dependency relations in our work. Parsing of natural language is a primary problem to many tasks that requires natural language processing.

The study of finding similar questions is also inspired by the work of finding information in the archives of Frequently Asked Questions(FAQ). Semantic similarities as well as statistical similarities between every question to rank FAQs has been extensively studied in FAQ Finder [33]. In this work [33]the authors used vector space models(VSM) for calculating the similarity and WordNet [34] was used to find a semantic similarity estimate between two questions. We used WordNet to retrieve synonyms for the extracted relations from each question. Lai et al. [35] developed an approach to mine FAQs from the web, however, the archives of CQA are different from the FAQ collections. The major difference between FAQs collection and CQA question archives is in the quality of both the questions and the answers. FAQ collections are usually created and maintained by private(depending on organizations) experts and the quality of questions and answers are therefore much better than CQA archives.

## 2.4 State of the Art Techniques

Two natural language sentence may mean the same but differ in the way the sentence is formed. For example, "While pursuing MSc can I start my own blog?" and "Can I start my own blog while pursuing MSc?" indicates the same meaning but differs in the way both these sentences are worded. The intuition that words with similar meaning will occur in similar context is consistent but there is a limit to how far this idea can be taken [36]. There are various approaches proposed in relevant literature for identifying similar questions, some of them has been briefly described below :

1. **Topic Modelling Approach:** This approach assumes that a particular question and their respective answer (question answer pairs) shares the same topic distribution. Researchers have developed a model to learn the latent topic space (which only emerges during the topic modelling process, therefore called latent) in these question answer pairs. This latent topic space is discovered by the topic models which helps in identifying similar questions and addressing the lexical gap problem [37] [38].

2. **Classic Term Weighting**: Research in the field of Information Retrieval related to combining document models and query models and thereby calculating the similarity based on the weights of the similar texts within documents has been proposed. [39].

3. **User Click Logs:** This approach was studied under the assumption that if two questions had similar retrieval results or similar click logs then the questions were semantically similar [40] [41].

4. **Translation Models:** Translation models usually defines the mathematical relationship between two or more languages. If two English language sentences 'e1' and 'e2' can be translated to a sentence 'f1' in second language supposedly French, then the two sentences are said to be semantically similar. Phrase Based Translation[42], Word Based Translation [43] and ClickThrough-based translation models [44] are some of the translation models which has been proposed in relevant literature.

5. **Ratio-Based Word Weighting**: A very recent approach by Bae et al.[3] studied on the assumption of similar questions having same or similar categories, a retrieval model was proposed which used ratio-based word weighting for fetching the importance of words in categories and the category importance was estimated using language modelling, translation based language modelling and category based model. Figure 2.1 represents the combination of baseline models employed to propose a ratio-based word weighting model.



Figure 2.1: Proposed approach for ratio-based word weighting [3]

6. **Deep Learning:** Question-question pair similarity is trained and tested in different variants of neural networks to predict semantically similar questions [45] [46].

## 2.5 Predicting user's satisfaction depending on recommended questions

A significant amount of work has been done [47] in studying the success of question answering and the dissatisfaction problem of information seekers in CQA. There are mainly two scenarios which keeps an asker unsatisfied and they are unsatisfactory answer and no answer. One approach for increasing user satisfaction is to submit new questions to expert users. Related research [48] defined quantitative measures of expertise to locate expert users. They presented an evidence of validating

this quantification as a measure of expertise and used it to find expert users who are best suited to answer a question.

In the field of question answering recommendation there is some relevant research conducted to find the best answers. In [49] a new topic model to simultaneously discover topical distribution of words, categories and users in question answering communities is introduced to find a ranked list of relevant answer. User activity and their related authority information along with Latent Drichlet Allocation model is used [50] to find a list of best answerers.

In relevant research[47] it was observed that 78% of the "best answers" were found to be reusable when it was put to use on future similar questions in Yahoo! Answers, however only 48% of the answers were exclusively the best.

This brings us back to our research question of identifying semantically similar questions which can be used as a measure for predicting users satisfaction. Questions which has already been answered can be reused and recommended to the users seeking information thereby enhancing user experience.

## 2.6   Extracting Knowledge from CQA sites

The ever increasing amount of data in CQA sites indicates a rich knowledge repository which resides in the form of natural language plain text. Extracting this knowledge and populating knowledge bases(KB) could be helpful in enriching knowledge repositories. There is a growing interest in automatically extracting information from these CQA sites and storing them in a structured format using KBs. Inspite of the growing interest, very little work has been done in extracting information from CQA sites for the purpose of storing them in knowledge bases and acting upon them for similar questions retrieval or question answering. Earlier research on fact extraction considering relations which are already defined between entities in text for supervised learning [51] and an unsupervised approach of extracting string of words between entities in a large corpora to cluster strings to produce relation string [52] has been studied. Another approach which has been used is of bootstrap learning [53] where in order to extract new patterns seed instances are used in an iterative manner . This approach, however, suffers from considerable drift in semantics including low precision. It is manually very costly to label large corpora to create samples for training in the case of supervised learning approaches, unsupervised or weakly supervised approaches are used in case of large datasets. One missing pattern can be detrimental for the entire model and its corresponding task of retrieving similar questions.

## 2.7 Finding user's intention for querying a CQA site.

A lot of work has already been done in finding user's intention by categorizing the question into different types. Rose et al. [54] proposed taxonomies with additional subcategories to earlier work of Broder [55] who categorized search queries in three categories namely navigational, informational and transactional. In the area of open domain question answering research, all of the highly rated systems had incorporated question taxonomies in their systems, [56],[57], [58]. Liu et al. [4] made some modifications to accommodate certain specifics of CQA services while following the work of Rose and Levinson (RL) [54]. Figure 2.2 demonstrates the taxonomy proposed by Liu et al. [4]. In this work they have retained the work of Broder's taxonomy [55] at the top levels and proposed a new Social category. The other categories such as Navigational, Informational and Transactional has been defined as it is in Broder's taxonomy, however, the Social category represented questions that didn't intend to get an answer but were just posed to start a conversation with people using CQA services. Along similar lines, Navigational category contained questions which seeks for pointers or URLs of specific websites as their information need. Transactional category contained questions which seeks for specific resources, for example, "What is the best recipe to cook baked beans?". Informational category is segregated in two subcategories *Constant* and *Dynamic*. Constant questions have a fixed set of answers for example,"Which country is the hottest in the world" while the answers for Dynamic questions are not defined e.g. "What is the temperature in Libya?". Dynamic category is further subdivided into three subcategories, *Opinion, Context Dependent* and *Open*. Opinion questions are those asking for opinions from people in CQA sites about what they think about some people, place or thing. "Is studying MSc right after college, worth it?" is an example of Opinion questions. Context-dependent questions are those questions for which the answer varies depending on the context, for example in the question, "What is the temperature in Libya?" the answer changes from minute to minute. Open questions are questions which are seeking for some facts or methods. There can be many different answers for Open questions. "Would you be able to list down all the citations and complete writing your MSc report in a month?" is one such example of open questions. Open questions are those questions which cannot be labelled in opinion or context dependent categories.

Most of the questions in CQA sites are opinion and open questions. The least number of questions are navigation questions as one would use a search engine to find the answer for a navigation question rather than choosing to use a CQA site for this purpose[4].

Figure 2.2: CQA Question Type Taxonomy [4]

## 2.8 Quora's Graph Structure

Since we have used a dataset published by Quora, it is worthwhile to briefly discuss about the internal structure of Quora. The structure of community question answering sites are a combination of questions, questions topics, answers and users. Figure 2.3 from Wang et al. [5] summarizes the relationships between Question, Topics and Users in Quora. Users can follow topics as well as other individuals. Questions in Quora are related to other questions and each question can be tagged using labels. For each question in the system, there are three participants, users who posted the question (Asker), the one who answered the question (Answerer) and the participants who voted on an answer(Voters).

Figure 2.3: Structure of questions, topics and users in Quora [5]

Quora has its own feature which has an ability to locate questions "related" to a given question. In this related question graph, nodes represents the questions and the edges represent a similarity measure as determined by Quora. Quora's internal structure is primarily dominated by three graphs.

1. User Topic Graph : Every user in Quora follows a set of topic based on which they receive updates of questions.

2. Social Graph : Users of Quora can subscribe to the activities of other users which gives rise to a social graph. Notifications are received by the users about the activities of the personalities whom they followed.

3. Question Graph : This graph is a feature used by Quora to map related questions for users to browse them quickly.

These three graphs are believed by Wang et al. [5] to be largely responsible for gaining attention of Quora users. Out of these three graphs the question graph is of primary interest for us as it is an important feature and provides Quora with an ability to locate questions "related" to a given question. The current production model for solving the problem of identifying similar questions, Quora uses a random forest model with handcrafted features, including the cosine similarity of the average of the word2vec, the number of common words, number of common topics labeled on the question and part of speech tags of the words. [59].

## 2.9 Challenges in Natural Language Processing and Linked Data

With the growing amount of structured knowledge on the web, the requirement for providing users with facilities to access to this knowledge increases steadily. However, the technical knowledge required for data modelling, forming and accessing vocabularies, and knowledge of web query languages is abstracted from the user. The main challenge for question answering over LD is the difficulty to translate user's information need in a format such that it can be used to evaluate using standard Semantic Web query processing techniques. In case a web user wants to search the linked data cloud for any information s/he would face two major obstacles. The first problem that the user will face is in the identification of a relevant dataset which contains the information that s/he is looking for, this is a very big challenge as the amount of datasets available in linked data cloud is huge. Therefore, identifying relevant dataset is not a trivial task. The second problem being the ability to formulate a query e.g. in SPARQL Protocol and RDF Query Language (SPARQL) which is the standard query language for RDF data for accessing the information from the dataset. One cannot assume that a common web user would be familiar with Semantic Web Languages or the available datasets and the structure of the linked data cloud. This is where a question answering system aims to bridge the gap between a user's information need and the structure of the data. A user's information need expressed in natural language on one side and structured queries to retrieve the information on the other side[27]. Another major challenge in dealing with natural language in general is of ambiguities. A natural language expression can be represented in such a manner that it has multiple meanings and this causes problems in mapping to multiple vocabulary element in the target dataset. As the amount of non English speaking web users creating and publishing data in their native languages are increasing, it is important to achieve a common goal that users from all nations irrespective of their language should have access to the same information. Even though, in principle, Semantic web is well suited for multilingualism as the Uniform Resource Identifiers (URIs) are language independent identifiers, a need for labels to access these URIs in various language contexts is imperative.

# Chapter 3

# Design

The previous chapters introduced the current challenges of finding similar questions in CQA sites, the attempts made in relevant literature to find a solution to this problem by introducing various techniques. It also discussed our motivation behind attempting a linked data approach to identify semantically similar questions. This chapter discusses the different categories of linked data applications and respective categorization of our approach. It illustrates the conceptual architecture at an abstract level. It discusses briefly on the different phases of the design and summarizes the effectiveness and limitations of such an approach.

## 3.1   Overview

According to "Constructivism", a term which refers to an idea that learners construct knowledge for themselves in a manner that each learner constructs meaning as s/he learns. In a work on constructivist learning theory [1] a strong assertion is made, "Constructing meaning is learning; there is no other kind"[60]. Along these lines, any form of new knowledge is learnt on the basis of what is already known by the learners. Emerging technologies such as Linked Data which advocates semantics could be made relevant for the users in this context. Such an application design to dynamically create and maintain knowledge bases is possible through a structured approach such as linked data.

We made two important design choices. Firstly the knowledge of users is captured in a knowledge base in the form of RDF triples which is the atomic unit of information in linked data. Second, the data used to query the knowledge base is also captured in an event when such knowledge is already not present in the knowledge base.

Our current system has four phases: a relationship extraction phase, an open domain ontology creation phase, a query execution phase and a phase for updating knowledge into the knowledge base.

---

[1]https://www.exploratorium.edu/education/ifi/constructivist-learning

## 3.2 Categorization of our Linked Data approach

Research related to categorization of linked data application indicates that such applications can be categorised based on dimensions that describe different technical aspects of using and representing Linked Data [6]. Figure 3.1 is a representation of the various dimensions, requirements and their specific benefits based on the categories of Linked Data applications. Applications based on Semantic web technologies if used extrinsically implies that linked data is consumed and processed by Application Program Interface (API) and for storage a relational database is required , while if used intrinsically the internal state of the application is stored in a triple-store rather than in a relational database. A linked data application can consume as well as produce linked data. Our proposed application consumes as well as produces linked data as it employs the RDF data model for internal representation of information.

There is a difference between shallow semantics and strong semantics and the important parameter lies in the kind of relationships between resources. Class properties, memberships and hierarchies can be expressed by RDF and RDF-Schema, however for increasing the richness of the semantics as well as the high level representation formalisms, variants of Web Ontology Language (OWL) is used. In terms of the expressivity of the knowledge representation techniques, the proposed application is rather constrained and mostly in RDF and RDF-Schema space. The OWL features used are mostly confined to class and property definitions. Therefore, semantic richness of the application's information space can be termed as *shallow*.

Applications using Linked Data can also be categorized in two levels *Isolated* and *Integrated*. If an application uses its own vocabulary which is different from the available vocabularies then such an application is termed as *isolated*. An application which uses, reuses and interlinks other vocabularies extensively is termed as *integrated*. An isolated application can be converted to an integrated application if such an application is published and the vocabularies are used to interlink other dataset. At the present moment, our application is in the isolated level as this application has not been published in the Linked Open Data (LOD) cloud.

Exploiting crowd intelligence for enrichment of knowledge base is also attempted by linked data applications. We consider user involvement to be primary in enriching the knowledge base and use this dimension in our proposed proof of concept application. Although an appropriate moderation process is required to prevent the publication of inappropriate material, at the present moment, application logic is used to make modifications in the predefined properties of the knowledge base.

| Dimension | Requirements | Benefits |
|---|---|---|
| *Semantic technology depth* | | |
| Extrinsic | mapping between internal information structures and RDF | standardised interaction |
| Intrinsic | sufficient query processing power | increased schema flexibility |
| *Information flow direction* | | |
| Consuming | mapping of RDF to internal information structures | wealth of additional structured information |
| Producing | mapping of internal information structures to RDF | increased information distribution |
| *Semantic richness* | | |
| Shallow | availability of structured information | pay-as-you go strategy |
| Strong | comprehensive knowledge engineering | automated reasoning |
| *Semantic integration* | | |
| Isolated | creation of own vocabularies and ontologies | simplified information governance |
| Integrated | vocabulary and identifier reuse on schema and/or instance level, co-reference and matching techniques | simplified syndication of semantic content |
| *User involvement* | | |
| Com.-oriented | provisioning of simple interaction with semantic content | exploitation of crowd intelligence |

Figure 3.1: Categories of semantic web applications [6]

## 3.3 Conceptual Architecture

According to Soni et al. [61] a conceptual architecture illustrates the important design elements of the system and it also describes the relationship among the design elements. A conceptual architecture is independent of the decisions made during implementation and it emphasizes on the interaction protocols between the design elements [61]. Architectural styles of linked data based application varies and is dependent on the constraints of the application. Our architecture reuses existing libraries and tools as components. A high level description of the components used in order to implement Semantic Web Technologies and Linked Data principles in our proposed application is discussed in the following section.

### 3.3.1 Description of Components

To discuss the component-based conceptual architecture we primarily introduce the following high-level components. Figure 3.2 represents the components and connectors of the conceptual architecture.

- Data Pre-Processing : This component is used to process data before using it for any further consumption by other components of the system.

- Data-modelling Service : This component provides a mapping from the unstructured data to a graph based data model such as RDF.

- Triple Store : This component provides a persistent storage for RDF.

- Query Engine : This component provides the ability to access the graph based data model by executing queries.

Figure 3.2: Components of conceptual architecture.

The data pre-processing task is responsible for processing the data based on anticipated problems for creating resources. Since the data we use are in natural language we use wrappers to translate the data in the form of semantic triples.We use these triples to form the data model for storing the RDF data in the triple store. We use a query engine to express queries over RDF data for finding a match in the triple store.

## 3.4 Relationship Extraction

Before persisting data in the triple store we extract the relations from a natural language question using two techniques. Firstly, we adopt the approach of extracting semantic triples using treebank parser from the work of Rusu et al. [8] and secondly we use the Open Information Extraction (OpenIE) technique [62] for relationship extraction. Once the triples has been extracted using these two approaches, they are integrated with respect to each question and used to persist data in the triple store.



Figure 3.3: Integrating triples generated from two approaches

State-of-the-art Natural Language Processing techniques were employed to extract semantic triples from the unstructured text. These triples were modelled to make the data eligible to be stored in a triple store. Figure 3.4 represents the high level flow of data from unstructured data to structured data while using a synonym engine to store the possible variants of a word along with storing the original triples generated. Once the relationship from a text are extracted in the form of triples, these words (three of them - subject, predicate and object) are fed into a synonym engine to extract their respective synonyms. These synonyms are then persisted in the triple store. Figure 3.5 presents the sequence

diagram for the event of extracting triples from an archive and persisting those triples in a triple store along with synonyms.



Figure 3.4: Flow Diagram of storing unstructured data in structured format in triple store.



Figure 3.5: Sequence diagram for storing triples in the triple store along with synonyms.

## 3.5 Synonym Engine

A synonym engine to retrieve the synonyms for all the extracted triples is employed to store semantically similar words as derived from WordNet lexical database. Semantic similarity measure has been researched in the area of information retrieval, text processing which includes text summarization, text categorization and text clustering. All these measures can be primarily grouped in four categories: feature based measures, information content based measures, path-length based measures and hybrid measures [63]. Specific "senses" of a word in a sentence for which the synonym needs to be

20

extracted is a problem and can be tackled using part-of-speech tagger which helps in determining the context. Figure3.6 presents a sample WordNet hierarchy which demonstrates that the "senses" of a word can be calculated based on the nodes and the hierarchy. In our approach the triples that are extracted from a natural language sentence are used to find the synonyms with synset types as 'noun' and 'verbs'. All the synonyms generated are further used to add description to a resource. These synonyms are also used to find a match during query execution when the searched resources are not present in the knowledge base.



Figure 3.6: A sample WordNet hierarchy [7]

## 3.6 Open Domain Modelling

In Linked Data "anybody can say anything about anybody" and everything is structured in triples. We have modelled the universe of discourse for our application in such a manner that we have used only two class definitions 'Subject' and 'Object'. 'Predicate' is used to describe the aspects of subject and is primarily used to establish a relationship between Subject and Object. Subject is the 'thing' or resource that is of interest and Object can be both a 'resource' or just a 'literal' such as a value in its raw form, for example a word or a number.

In this approach, the domain model is a set of RDF triples which includes RDFS and OWL definitions. It is not restricted to any specified domain, instances of resources represents the information items in the RDF graph. Variants of semantic triples generated are used to form knowledge in the

triple store and such an approach facilitates more possibilities of representing the same sentence in multiple set of semantic triples which are modelled in RDF graph and used for querying.

In order to perform computations over data to modify resources or to present data to the user, open domain modelling is advantageous when dealing with data from CQA sites, as the questions posted in community question answering sites are diverse and is not always constrained to any specific domain. In this approach, we choose such a modelling scheme because of the extensive diversity of questions and uncertainty of questions being posted in CQA forums.

## 3.7 RDF Triple Store

A Resource Description Framework(RDF) repository is a basis for implementing the identification of semantically similar question system. The RDF repository also referred to as the knowledge base, includes a large number of triples that assists a computer(machine) in recognition and understanding. RDF data in our approach is modeled as a directed, labeled graph G = (V,E) where V is a set of nodes and E is a set of directed edges. In this approach, we map the triples generated from a particular question to its corresponding resources. Each triple is expressed in the form of 'a subject','a predicate' and 'an object'. For example, for a given sentence: "A student studies computer science." the following are the components, a subject "student", a predicate "study", an object "computer science". Using these triples generated from the relationship extraction module a RDF repository has been created. The following are the features of the semantic graph.

- Outgoing links – The node occurred as subjects in the graph.

- Incoming links – The node occurred as objects in the graph.

- Edges of the graph — Represents the established relationship between two nodes of the graph.

We used this approach to formulate a semantic graph to contain the information of natural language in the form of triples.

## 3.8 Our Approach

Table 3.8 illustrates the comparison of our approach with other state of the art approaches as mentioned in section 2.4. Our proposed approach is a human understandable, machine readable and structured approach.

| Approaches | Human Under-standable | Machine Readable | Structured |
|---|---|---|---|
| Topic Model | No | Yes | No |
| Classic-Term Weighting | No | Yes | No |
| User Click Logs | No | Yes | No |
| Translation Models | No | Yes | No |
| Deep Learning | No | Yes | No |
| Proposed-Approach | Yes | Yes | Yes |

Table 3.1: Comparison of Proposed Approach with Other Approaches

Our proposed approach is for consuming as well as creating linked data. In the event, a user asks a natural language question, semantic triples are extracted from the question. These semantic triples are used to construct SPARQL queries to search the triple store for a matched result. The matched result is the output in this case and subsequently a validation of the matched result and the user-asked question is performed. In case, the matched result and the user-asked question are syntactically different, the user question along with the associated triples are stored in the triple store. Figure 3.7 represents a high-level flow of the process of recommending semantically similar questions and storing the related knowledge in the specified triple store. It is seen in Figure 3.7 that in case no matched questions are retrieved from the execution of generated SPARQL queries, such questions are persisted in the triple store along with associated triples. Once the persistence in triple store is accomplished by the program it finishes execution.

Figure 3.7: Flow Diagram of recommending questions and enhancing the triple store

Figure 3.8 presents a sequence diagram in the event of user asked question and depicts the interaction between modules from when a question is asked to when it is persisted to the triple store. A set of triples is generated from the user asked question and these triples are used to generate SPARQL queries to search the triple store for a similar question. This action of searching for a similar question occurs in five levels and once the match is made, the program moves on to the next set of triples which is represented by the outer loop. This process of SPARQL generation and execution for all the set of triples generated from the user asked question when concluded, the results fetched from the triple store is presented to the end user. If the user accepts the question to be semantically similar such a recommended question is linked to the user asked question. In case the user does not accept questions to be semantically similar such user asked questions are persisted in the triple store.

Sequence Diagram - User Asked Question



Figure 3.8: Sequence Diagram in an event when the user asks a question

A high-level workflow of our proposed application is presented in Figure 3.9. The relationship extraction module serves as a service for extracting semantic relations from a question archive to build a knowledge base and the same module is responsible for extracting semantic relation in the event a new question by the user is encountered by the system. After the extraction of the triples, the triples are stored in the triple store and the triples are also processed through a synonym engine for retrieving synonyms, See Figure 3.4 for each entity. The synonyms are stored in the triple store along with the extracted triples. When a new user question is encountered by the system, the relationship extraction module generates the triples which are used to construct SPARQL queries. To find a match in the knowledge base the SPARQL queries are executed on RDF data for retrieving similar questions. If a match is found, the information is presented to the user as a recommended question, in the event a match is not found, the consequent triples generated from the user question is stored in the triple store. If a recommended question and the user asked question is different and agreed to be semantically similar by the user, such information is also updated in the knowledge base for future processing and querying. If both the recommended as well as the user asked question is word to word similar, no update is made in the ontology in this case.

Figure 3.9: High Level workflow of our proposed approach.

## 3.9 Summary

We do not claim this to be the best approach to design a model for identifying semantically similar questions. However, at the current moment due to the absence of such a model for finding semantically similar sentences solely using linked data approach which facilitates both human and machine understandable paradigm, this is a proposed idea. This design can be enhanced in multiple ways, such as by using fine-tuned classifiers to rank triples and modelling the system using 5WH questions (What, Where, When, Why, Who and How) as different classes to hold related semantics and questions which can be answered with yes/no as other classes. By using named graphs to store information of specific subjects, which would enable in executing more specific SPARQL queries. Additionally other common variants of questions needs to be appropriately classified into different classes. Our research question is to what extent can natural language questions be structured into RDF graphs and with the extraction of semantic triples and usage of synonym engine to get related words it is evident that natural language can be structured in RDF graph which can then be queried using SPARQL. At the present moment it has a small but significant impact on this concept of structuring unstructured text with the help of semantic triples.

# Chapter 4

# Implementation

In the previous chapter we discussed the high level design of the application. In this chapter we will be discussing the implementation in details. We implemented our proposed approach as a JAVA application. This chapter discusses on the dataset used, the input and output considerations, how the data was pre-processed, the relationship extraction phase, mapping to triple store and finally interaction with the triple store. We finally summarize this chapter mentioning the accomplishments and the limitations of this approach.

## 4.1 Dataset

We run our experiment on Quora Question Pairs dataset [64]. This dataset is primarily related to the problem of identifying duplicate question. Identification of duplicate question has been stated as one of the challenges that arise in building a scalable knowledge-sharing platform. The dataset consists of 404,301 lines of potential pairs of duplicate questions. Each line contains IDs of the question, full text of each question and a binary value representing whether the line contains a duplicate pair of questions. Currently, Quora uses a Random Forest model to identify duplicate questions [64]. The data is in the following format:

- id - unique id of the question pair.

- qid1, qid2 - unique ids of question 1 and question 2.

- question 1, question 2 - the complete text of each question

- is_duplicate - if question 1 and question 2 are semantically similar then the value is '1', otherwise the value is stated as '0'.

A brief snapshot of the dataset is shown in Figure4.1. Labelling of the question being semantically similar is inherently subjective as it is annotated by human experts.

| id | qid1 | qid2 | question1 | question2 | is_duplicate |
|---|---|---|---|---|---|
| 0 | 1 | 2 | What is the step by step guide to invest in share market in india? | What is the step by step guide to invest in share market? | 0 |
| 1 | 3 | 4 | What is the story of Kohinoor (Koh-i-Noor) Diamond? | What would happen if the Indian government stole the Kohinoor (Koh-i-Noor) diamond back | 0 |
| 2 | 5 | 6 | How can I increase the speed of my internet connection while using a VPN? | How can Internet speed be increased by hacking through DNS? | 0 |
| 3 | 7 | 8 | Why am I mentally very lonely? How can I solve it? | Find the remainder when [math]23^{24}[/math] is divided by 24,23? | 0 |
| 4 | 9 | 10 | Which one dissolve in water quikly sugar, salt, methane and carbon di oxide? | Which fish would survive in salt water? | 0 |
| 5 | 11 | 12 | Astrology: I am a Capricorn Sun Cap moon and cap rising...what does that say about | I'm a triple Capricorn (Sun, Moon and ascendant in Capricorn) What does this say about me? | 1 |
| 6 | 13 | 14 | Should I buy tiago? | What keeps childern active and far from phone and video games? | 0 |
| 7 | 15 | 16 | How can I be a good geologist? | What should I do to be a great geologist? | 1 |
| 8 | 17 | 18 | When do you use ã,· instead of ã—? | When do you use "&" instead of "and"? | 0 |
| 9 | 19 | 20 | Motorola (company): Can I hack my Charter Motorolla DCX3400? | How do I hack Motorola DCX3400 for free internet? | 0 |
| 10 | 21 | 22 | Method to find separation of slits using fresnel biprism? | What are some of the things technicians can tell about the durability and reliability of Laptop | 0 |
| 11 | 23 | 24 | How do I read and find my YouTube comments? | How can I see all my Youtube comments? | 1 |
| 12 | 25 | 26 | What can make Physics easy to learn? | How can you make physics easy to learn? | 1 |
| 13 | 27 | 28 | What was your first sexual experience like? | What was your first sexual experience? | 1 |
| 14 | 29 | 30 | What are the laws to change your status from a student visa to a green card in the U! | What are the laws to change your status from a student visa to a green card in the US? How do | 0 |
| 15 | 31 | 32 | What would a Trump presidency mean for current international master's studen | How will a Trump presidency affect the students presently in US or planning to study in US? | 1 |
| 16 | 33 | 34 | What does manipulation mean? | What does manipulation means? | 1 |
| 17 | 35 | 36 | Why do girls want to be friends with the guy they reject? | How do guys feel after rejecting a girl? | 0 |
| 18 | 37 | 38 | Why are so many Quora users posting questions that are readily answered on Googl | Why do people ask Quora questions which can be answered easily by Google? | 1 |
| 19 | 39 | 40 | Which is the best digital marketing institution in banglore? | Which is the best digital marketing institute in Pune? | 0 |
| 20 | 41 | 42 | Why do rockets look white? | Why are rockets and boosters painted white? | 1 |
| 21 | 43 | 44 | What's causing someone to be jealous? | What can I do to avoid being jealous of someone? | 0 |
| 22 | 45 | 46 | What are the questions should not ask on Quora? | Which question should I ask on Quora? | 0 |
| 23 | 47 | 48 | How much is 30 kV in HP? | Where can I find a conversion chart for CC to horsepower? | 0 |
| 24 | 49 | 50 | What does it mean that every time I look at the clock the numbers are the same? | How many times a day do a clock's hands overlap? | 0 |
| 25 | 51 | 52 | What are some tips on making it through the job interview process at Medicines? | What are some tips on making it through the job interview process at Foundation Medicine? | 0 |
| 26 | 53 | 54 | What is web application? | What is the web application framework? | 0 |
| 27 | 55 | 56 | Does society place too much importance on sports? | How do sports contribute to the society? | 0 |
| 28 | 57 | 58 | What is best way to make money online? | What is best way to ask for money online? | 0 |
| 29 | 59 | 60 | How should I prepare for CA final law? | How one should know that he/she completely prepare for CA final exam? | 1 |
| 30 | 61 | 62 | What's one thing you would like to do better? | What's one thing you do despite knowing better? | 0 |
| 31 | 63 | 64 | What are some special cares for someone with a nose that gets stuffy during the nig | How can I keep my nose from getting stuffy at night? | 1 |
| 32 | 65 | 66 | What Game of Thrones villain would be the most likely to give you mercy? | What Game of Thrones villain would you most like to be at the mercy of? | 1 |
| 33 | 67 | 68 | Does the United States government still blacklist (employment, etc.) some United S | How is the average speed of gas molecules determined? | 0 |
| 34 | 69 | 70 | What is the best travel website in spain? | What is the best travel website? | 0 |
| 35 | 71 | 72 | Why do some people think Obama will try to take their guns away? | Has there been a gun control initiative to take away guns people already own? | 0 |

Figure 4.1: First few lines of the raw dataset

The plot representing the number of unique questions versus the number of repeated questions in Figure 4.2 clearly indicates that the dataset contains more number of unique questions. This helps us in our research as more number of unique questions create more knowledge for our knowledge base.



Figure 4.2: Number of unique questions versus repeated questions

## 4.2 Input

We select a set of 20,000 questions from the Quora Question Pairs Dataset in an iterative manner to populate the knowledge base.These questions are used for further processing which encompasses triples generation and storage in triple store.

For the entire system's framework we expect a question in natural language from a user. This question is then further processed to extract relationship information in the form of triples to interact with the knowledge base.

## 4.3 Output

For a given natural language question, after the triples extraction phase, we construct SPARQL queries which interact with the knowledge base to find a match based on the triples generated from the question. These matched questions from the knowledge-base are then presented to the user.

## 4.4 Data Preprocessing

For extraction of semantic triples, creation of resources and providing description for the resources, the following data pre-processing tasks are applied.

According to the work by Fossati et al. [65] word window and syntactic pattern are taken into consideration where a sentence below the word count of 5 is not taken into account and sentences which lack the basic pattern of Noun-Verb-Noun tags is not considered for relationship extraction. The special characters in a natural language sentence are replaced by natural language words for example, '+' has been replaced with 'plus'. This action has been performed so as to avoid creation of malformed Uniform Resource Identifiers (URIs).

## 4.5 Relationship Extraction

The main advantages of triple extraction from natural language sentences are manifold. Triple extraction provides us with a compact and simple representation of information contained in the sentence. It avoids the complexity of a full parse on a natural language sentence. The main reason for using triple extraction in our approach is because the triples contain semantic information which could be used to build a knowledge base as well as use the triples to query the knowledge base.

There are essentially two approaches for triple extraction which we have used in our implementation. Extraction using parse tree of a sentence using heuristic rules such as OpenNLP - Treebank ParseTree which has been used in numerous study for entity extraction as well as recommendation

systems [66] and the other approach is by Open Information extraction. To organize words in groups which holds grammatical sense parse tree is required. An example of a pictorial representation of a parse tree for a specified question "Who is the best tennis player in the world" as generated by Stanford Parser is shown in Figure 4.3. The meaning of all the tags are present in Appendix B. The relevant tagging and the parse representation is shown in Figure 4.4



Figure 4.3: Parse tree using Stanford Parser [2]for a question *"Who is the best tennis player in the world?"*



Figure 4.4: Tags and parse representation using Stanford Parser [4]for a question *"Who is the best tennis player in the world?"*

There are various feature which depends on the triple candidates such as length of the sentences,

number of words, context of subject, predicate and object, distance between subject, predicate and object, number of links from subject, predicate and object, depth, diameter, siblings, uncles, cousins, path to root and part of speech tagging are some of the features which are used to determine the triples that are closely relevant to a given natural language sentence.

In the relationship extraction stage the system first parses a natural language question to a set of triples. The work of Russu et al. [8] triple extraction using treebank parsers has been adopted in our work. Along with this work, we have used Stanford's Open Domain Information Extraction (OpenIE) which is a part of Stanford CoreNLP. These two methods are used for extraction of relation tuples. The major advantage of OpenIE over other open information extraction systems is that the schema for these relations which are generated need not be specified in advance. Given a natural language sentence, pre-processing of the sentence in various linguistically motivated ways are performed to produce coherent clauses. These clauses are then used to segment OpenIE triples [62]. Open IE triples has been used in a number of applications for structured relations [67] and entailment relations [68] which are critical to many natural language processing applications for semantic parsing and question answering.

### 4.5.1 Triples extraction using Treebank Parsers

While using a treebank parser a sentence (S) is represented as a parse tree with S being the root of the tree having three main children, a noun phrase (NP), a verb phrase (VP) and a full stop (.) [8].

In order to find the subject of the sentence the subtree of NP is searched using breadth first search. The very first child of NP which is a noun is selected. In the following subtrees shown in Figure 4.5 nouns can be found. The tags which starts with 'NN' in the subtree of 'NP' is considered for finding the subject in the parse tree representation of the natural language sentence.

| Subtree | The type of noun found |
|---------|------------------------|
| NN | noun, common, singular or mass |
| NNP | noun, proper, singular |
| NNPS | noun, proper, plural |
| NNS | noun, common, plural |

Figure 4.5: Selection of Subjects from the subtree[8].

For finding the predicate of the sentence, a search is performed in the VP subtree and the deepest verb descendent is considered to be the predicate. The predicates are found in the following subtrees as seen in Figure 4.6.

| Subtree | The type of verb found |
|---------|------------------------|
| VB | verb, base form |
| VBD | verb, past tense |
| VBG | verb, present participle or gerund |
| VBN | verb, past participle |
| VBP | verb, present tense, not 3rd person singular |
| VBZ | verb, present tense, 3rd person singular |

Figure 4.6: Selection of Predicates from the subtree[8].

Finally for objects, three different subtrees are searched, all the cousins of the VP subtree which contains the predicate. Subtrees which are ADJP(Adjective Phrase), PP(prepositional phrase) and NP. In NP and PP the first occurring noun is searched for while an adjective is searched in the following subtrees of ADJP shown in Figure 4.7.

| Subtree | The type of adjective found |
|---------|-----------------------------|
| JJ | adjective or numeral, ordinal |
| JJR | adjective, comparative |
| JJS | adjective, superlative |

Figure 4.7: Selection of Objects from the subtree[8].

Natural Language Toolkit(NLTK) version 3.0 has been leveraged for the purpose of using Stanford Model and Stanford Parser which has been used to carry out the task for generating the subtree and based on the labels of the subtrees the segregation of the subject, predicate and object has been executed. The results for each of the questions in the dataset has been stored in a comma-seperated values(CSV) file for further processing and integration with the triples obtained from the open information extraction approach. The triples generated by this approach is only used during the formulation of the knowledge base and not during the interaction with the triple store. In response to a user question a set of triples are generated by OpenIE approach which are used to find a similar question.

### 4.5.2 Triples extraction using Open Information Extraction

Open Information Extraction (IE) triples have been used in a number of applications for entailment graph learning [68]. Stanford Open IE is motivated from the work of finding short entailment from long sentences as described in the work of [69]. A set of clauses which can be considered on its own as syntactically and semantically correct and is entailed by the original sentence is produced. One very important project which uses Open Information Extraction is the Never Ending Language Learning(NELL) project which is in effect from the year 2010 and since then is processing the web

for evolving a knowledge base of facts and relations [70]. Figure 4.8 demonstrates two example of Part-of-Speech(POS) tagging and basic dependencies using Stanford CoreNLP.



Figure 4.8: Part of speech and basic dependency generated by Stanford CoreNLP [9]

For different NLP applications including question answering Open IE forms a part of the pipeline. Figure 4.9 represents a general pipeline for Open Information task.

**Text segmentation**

Sentence split and tokenizer.

Knowledge source

**Parsing**

Dependency parser, part-of-speech tagger, semantic role labeling or another's shallow parsers

**Extracting facts**

Rule based or machine learning method

**Another NLP task**

Ontology construction, question answering system, text summarization...

Open Information Extraction Pipeline

Figure 4.9: Pipeline of Open Information Extraction [9]

In our implementation we used Stanford CoreNLP which is written in Java. The CoreNLP

pipeline has been used in our implementation which has been inspired from the work of Manning et al. [71]. It has a lightweight framework which uses Java objects. Figure 4.10 presents the packages that were imported to facilitate the generation of relational triples. Annotators listed in Figure 4.11 are used in our implementation to create a Stanford CoreNLP pipeline. We have used the 3.9.2 version of stanford-corenlp to execute this task. In this process a CoreNLP document object is created which holds the list of sentences and every sentence is iterated over to generate relational triples.

```
import edu.stanford.nlp.ie.util.RelationTriple;
import edu.stanford.nlp.ling.CoreAnnotations;
import edu.stanford.nlp.pipeline.Annotation;
import edu.stanford.nlp.pipeline.StanfordCoreNLP;
import edu.stanford.nlp.simple.Document;
import edu.stanford.nlp.naturalli.NaturalLogicAnnotations;
import edu.stanford.nlp.util.CoreMap;
```

Figure 4.10: Packages imported for the task of extracting triples by Stanford CoreNLP OpenIE approach.

```
"annotators", "tokenize,ssplit,pos,lemma,depparse,natlog,openie")
```

Figure 4.11: Annotators used to set property for creating Stanford CoreNLP pipeline

## 4.6 RDF Data Model - Triple Store

Each entity has an associated rdfs:label which we use to store the questions as human-readable label for each entity. Figure 4.12 represents an example of the RDF graph for three independent questions having the same subject. In the representation, *base:Richard_Muller* is the subject which is associated by two predicates *base:be* and *base:do_think_of* to three questions and four objects. The rdfs:label associated with each subject is a question from which the subject was extracted. Similiarly, for each predicate and object rdfs:label is used to specifically store a related question in the form of human-readable label.

### 4.6.1 WordNet

WordNet is widely used for NLP tasks and serves as a lexical database in which the data is persisted in a related fashion as a semantic network. We use the rdfs:comment to tag the related words for each entity. Along with the questions shown in Figure4.13 the rdfs:comment is used to store the related words of *base:do_think_of* as fetched from JAVA API for WordNet Searching (JAWS) which provides

JAVA applications with the ability to retrieve data from WordNet database. JAWS was developed by Brett Spell as a student project in Southern Methodist University [5]. This API is compatible with both 2.1 and 3.0 versions of the WordNet database.It can be used with Java 1.4 and later versions. In our implementation we have used the 2.1 version of WordNet database and used this API to generate similar word forms and stored them in rdfs:comment for each entity.



Figure 4.12: RDF graph representation of a question stored in triple store [6]

---

[5]https://archive.org/stream/PracticalArtificialIntelligenceProgrammingWithJava/JavaAI3rd_djvu.txt
[6]http://www.ldf.fi/service/rdf-grapher

Figure 4.13: Predicate associated to different questions [7]

---

Figure 4.14: Object associated to a particular question: "What does Richard Muller think of Ed Witten?" [9]

The objects as seen in Figure 4.14 **base:Ed_Witten** and in Figure 4.15 **base:Michio_Kaku** is annotated with the natural language question from which it was extracted using the rdfs:label tag.



Figure 4.15: Object associated to a particular question:"What does Richard Muller think of Michio Kaku?" [11]

---

[9]http://www.ldf.fi/service/rdf-grapher
[11]http://www.ldf.fi/service/rdf-grapher

### 4.6.2 Apache Jena

Apache Jena[12] is a free open source Java framework for developing semantic web and Linked Data applications. In our implementation we have used the 3.9.0 version of Jena for processing RDF data, building knowledge base and formulating queries over RDF data. We use the OntModel interface to create an enhanced view of the Jena model to contain the ontology data.

### 4.6.3 Model

The OntClass interface which represents an ontology node representing a class description has been used in our implementation. The two classes in our implementation are the Subject and the Object and we define the Predicate as an ObjectProperty which is an interface for encapsulating properties whose range values are restricted to individuals. Figure 4.16 shown below is a representation of dynamically generating the RDF data model.

```
subjectIndividual.addProperty(predicate, objectIndividual);
```

Figure 4.16: Using Jena to form RDF data model

Figure 4.17 represents the domain and range values for the predicate. The setDomain function is for asserting a given resource to represent the class of individuals that would form the domain of the property. Similarly, the setRange function is for asserting a given resource to represent the class of individuals that would form the range of the property. In our case, the domain for predicate is set to be the subject and the range is set to be the object.

```
predicate.setDomain(subject);
predicate.setRange(object);
```

Figure 4.17: Definition of the domain and range for predicate

## 4.7 Interacting with Triple Store

Finding an answer to a query executed on a RDF dataset is a pattern matching problem. Triples which match a particular set of graph patterns is retrieved as a result. We have used TDB which is a component of Jena for RDF storage and query and is used as a RDF store in a single machine. In the event that an user asks a natural language question the relationship extraction module extracts the triples from the question and constructs SPARQL queries using the triples. There are five levels of matches which are performed to find a question based on the triples which are generated. In case the

---

[12]https://jena.apache.org/getting_started/

question asked by the user is 'new', the definition of 'new' here is a question which is not already present in the knowledge base, the question is updated in the knowledge base along with the extracted triples.

There are five levels of search which is attempted on the knowledge base to retrieve a similar question or set of similar questions. The following are the descriptions of the levels.

- *Level 1:* At this level the SPARQL query searches for the URIs of the Subject, Predicate and Object having a match for each triple that is generated from the unstructured question. In case there are no match found at this level then the application proceeds to the next level of search; Level 2. If a match is found, the system stores the required information for presentation and moves on to the next set of triples generation and tries to find a match at Level 1.

- *Level 2:* At this level the SPARQL query searches for the URIs of the Subject and Predicate with a search on the synonym set of the Objects. In case there are no match found at this level then the application proceeds to the next level of search; Level 3. If a match is found, the system stores the required information for presentation and moves on to the next set of triples generated by the system and tries to find a match at Level 1.

- *Level 3:* At this level the SPARQL query searches for the URIs of the Subject and Object with a search on the synonym set of the Predicate. In case there are no match found at this level then the application proceeds to the next level of search; Level 4. If a match is found, the system stores the required information for presentation and moves on to the next set of triples generated by the system and tries to find a match at Level 1.

- *Level 4:*At this level the SPARQL query searches for the URIs of the Predicate and Object with a search on the synonym set of the Subject. In case there are no match found at this level then the application proceeds to the next level of search; Level 5. If a match is found, the system stores the required information for presentation and moves on to the next set of triples generated by the system and tries to find a match at Level 1.

- *Level 5:* At this level the SPARQL query searches in the synonym set of the Subjects, Predicate and Object. In case when a match is found, the system stores the required information for presentation. In case there are no matches at this level, there are two tasks which executes, if the system is able to generate another set of triples from the sentence, the whole process of search commences from Level 1. If the system is unable to generate anymore triples and there are no matches found at Level 5, the triples along with the question is stored in the knowledge base.

Once all the levels of searches for all the generated triples are concluded, the users are presented with the information - the recommended question with a score. This score is not indicative of the

lexical similarity of the recommended question and the asked question. Instead, it is a measure of how many times and at which level is a match found. It is a cumulative score of the number of representational matches of the recommended questions made in the knowledge base. For scoring the levels we made an assumption that information retrieved at each levels is in descending order of importance; A match at Level 1 is more important and highly weighted than a match at Level 2 and similarly for other levels. Table 4.1 indicates the scores assigned at each level. This allocation of score is to represent a dummy confidence score for the recommended questions, it is also allocated to distinguish between two or more recommended questions and how probable is that a recommended question is semantically similar to a asked question based on the matches found at the subsequent levels.

| Levels | Score Assigned |
|---|---|
| Level 1 | 100 |
| Level 2 | 90 |
| Level 3 | 70 |
| Level 4 | 50 |
| Level 5 | 10 |

Table 4.1: Scores assigned to each Level

Figure 4.18 demonstrates the levels of matches being performed on processing through all the levels of SPARQL query execution for each triple generated. The presentation of the recommended question along with the cumulative score which is obtained by the matches made at each level is finally presented to the user. In this example, since the match for the first set of generated triples is found at level 2 and there are no subsequent matches at other levels for all the generated triples, the recommended question is presented to the user with a score of 90.

It is to be noted that when a user question is encountered by the system, the OpenIE approach of triple generation is executed to generate the triples and construct the queries based on these triples.

```
Enter your question:
Does anyone know about the best basketball player in the globe?          ◄━━━ User asked question
Subject: best_basketball_player Predicate: in    Object:      globe
No questions at Level....1
Subject: anyone Predicate: do_know_about        Object:      basketball_player
No questions at Level....1
No questions at Level....2
No questions at Level....3
No questions at Level....4
No questions at Level....5
Subject: anyone Predicate: do_know_about        Object:      best_basketball_player
No questions at Level....1
No questions at Level....2
No questions at Level....3
No questions at Level....4
No questions at Level....5
Subject: anyone Predicate: do_know_about        Object:      best_basketball_player_in_globe
No questions at Level....1
No questions at Level....2
No questions at Level....3
No questions at Level....4
No questions at Level....5
Subject: anyone Predicate: do_know_about        Object:      basketball_player_in_globe
No questions at Level....1
No questions at Level....2
No questions at Level....3
No questions at Level....4
No questions at Level....5
question: Who is the best basketball player in the world    ◄━━━ System Recommended
score: 90                                                         Question

Do you find this question similar? (y/n)
y
Updating Ontology with new knowledge!
```

Figure 4.18: Demonstration of user asked question and system recommended question

If a user accepts that the represented question is semantically similar by indicating their choice of *Yes* to the given question "Do you find this question similar? (y/n)" shown in Figure 4.18 the user asked question is linked to the system recommended question in all the common triples that these two questions shares. Additionally all the unique triples generated exclusively for the user asked question is also stored in the knowledge base. In case, the user does not agree to the output being shown (recommended question) with the input (question posed by the user), the user can indicate their response as negative. In such a scenario, the user asked question along with all the generated triples are stored in the knowledge base for future querying and reference. In this fashion, the ontology is enriched with new knowledge from the user asked questions.

In case, a question asked by the user is already present in the knowledge base and is exactly similar (word to word similar) to the question present in the knowledge base, the recommended question to the user in such a case is not updated in the knowledge base as the same question generates the same triples and such an update would mean overriding the same entities which is computationally

expensive and technically not relevant.

The update in the knowledge base is conducted by SPARQL INSERT query. Figure 4.19 presents the definition of the method and parameters used to update the knowledge base. The parameters used are the triples (subject, predicate, object) and the user asked question, which in this case is indicated by the variable *'questionString'*. All the entities are labelled with the user asked question using *rdfs:label* tag.

```
public void insertQuestions(String subject, String predicate, String object,
                            String questionString)
```

Figure 4.19: Method definition to update knowledge base

The method shown in Figure 4.19 is used for faster update, without updating the synonyms. In case the update is to be made with the synonyms an additional parameter *synonymObject* to the method as shown in Figure 4.20 is used and the synonyms are stored in *rdfs:comment* tag.

```
public void insertQuestions(String subject, String predicate, String object,
                String questionString, ExcelObject synonymObject)
```

Figure 4.20: Method definition to update knowledge base along with synonyms

## 4.8 Summary

The extraction of semantic triples, modelling them using RDF and thereby using these triples to construct queries for expressing information need over RDF data was accomplished. Unstructured data was converted to structured data and thereby using the structured data to store and query machine understandable metadata was realized in this approach. User knowledge was contained in the knowledge base with the provision of enriching the triple store. The implementation of the proposed approach suggests that it is possible to identify semantically similar natural language sentences using linked data to a considerable extent, however, the application at present is in its nascent stages and there are room for improvements. Analysis of the triples generated and ways of generating triples of better quality is a present limitation. Better engineering decisions related to this module could enrich not only the quality of the triples generated but also the quality of knowledge contained in the knowledge base. The evaluation and the findings for this proof-of-concept approach has been presented in the forthcoming chapter.

# Chapter 5

# Evaluation

For evaluating this prototype formative user-centered evaluation has been conducted. This chapter describes how the evaluation was conducted, what were the findings, which metrics were used to evaluate the system and what was the outcome of these evaluations. The chapter concludes with a critical analysis of the evaluation process and threats to validation in the process.

## 5.1 User Evaluation

Since the evaluation of semantic web technologies based approach of identifying semantically similar question is not a straightforward task as there are no available gold-standard datasets for comparison, human intervention is required and to the best of our knowledge there are no standard criteria which is defined to carry out manual evaluation.

In our case, judgement is primarily based on the ratio between the semantic similarity of the recommended question versus the asked question. In the evaluation process, the person who framed a question made the judgment as well.

## 5.2 Metrics

An overall agreement was required as there were a number of human participants. The number of human judges and the categories (whether a recommended question is semantically similar; yes or no) are considered.

The judgement of the semantic accuracy with respect to categories (Yes/No based on Semantic Similarity) by the user who asked this question (participant's judgement) has been noted.

Kappa-based metric [72] for inter-rater agreement is used. There are possibilities of bias and chance among the agreement between evaluation participants, therefore, to avoid this, we use the

Kappa [72] value to determine the inter-rater agreement between the observed data and the prior data. To obtain the value of Kappa the following formula is used.

$$K_{free} = \frac{\left[ \frac{1}{Nn(n-1)} \left( \sum_{i=1}^{N} \sum_{j=1}^{k} n_{ij}^{2} - Nn \right) \right] - \left[ \frac{1}{k} \right]}{1 - \left[ \frac{1}{k} \right]} \tag{5.1}$$

In this equation, k is the number of categories, N is the number of cases, n is the number of evaluation participants (rater), and $n_{ij}$ is an array of ratings where n_ij[i][j] evaluators assigned case i to category j.

In the following Figure 5.1 the Kappa statistics and the division of the strength of agreement is represented.

| Kappa Statistic | Strength of Agreement |
|---|---|
| <0.00 | Poor |
| 0.00–0.20 | Slight |
| 0.21–0.40 | Fair |
| 0.41–0.60 | Moderate |
| 0.61–0.80 | Substantial |
| 0.81–1.00 | Almost Perfect |

Figure 5.1: Useful benchmarks related to Kappa Statistics [10]

We also adopt the measures used in a study by Zhu et al. [1] to evaluate the performance of our proposed approach represented in Table 5.2.

| Performance Measures | | |
|---|---|---|
| Category | Measure | Description |
| Relevance Based | Success Rate | A binary measure (i.e., success or not) represents completion of each task. |
| Interaction Based | Search Time | Time taken to complete a search task |
| | Query size | The number of issued queries to complete a search task. |

Table 5.1: Measures used in evaluation [1]

We measure the success rate by the ability of the application to generate semantic triples from the

given text (user posed questions). The search time is the time taken for the application to recommend a question and the query size is the number of queries required to find a semantically similar match.

Finally we present a ratio between the following

- True Positive(TP)- Questions which are semantically similar and recommended by the system during the evaluation process.

- False Positive(FP) - Questions which are semantically not similar but still recommended by the system during evaluation process.

Since we have the values for TP and FP we can calculate the Precision of the system by the following equation.

$$Precision = \frac{TP}{TP+FP} \qquad (5.2)$$

We have considered the above-mentioned metrics to evaluate our system.

### 5.2.1 Subjects

Fifteen participants were recruited using a convenience-sampling method. The participants were sent a recruitment message for volunteering. To adjust the research setting a pilot-test was conducted. During the test, through interviews with potential participants we found that most of the participants were users of community question answering sites, who have either posted a question or posted an answer or have posted questions as well as answered questions and those who used such sites to get a solution for their query.

### 5.2.2 Tasks

The tasks designed were as such wherein users were provided with partial information, where they were shown a set of contexts based on which they were asked to frame questions. For example, they were shown the below set of semantic triples as shown in Figure 5.2 and asked to frame questions, the suggestion that there is no right or wrong answers were indicated before the test commenced. In specific situations users were asked to frame variants of the questions they posed. To reduce the variance of task complexity, each of the task sets were curated to be as similar as possible. Each task set had 4 contexts based on which the participants were asked to frame questions.

```
Context:  you     do_stop_distraction      study
          you     do_stop                  distraction
```

Figure 5.2: Example of context presented to participants for framing questions.

There were three set of contexts with each context set having four contexts. The following are the number of human participants who participated in each context set.

| Context Set | Number of Participant |
|---|---|
| Set 1 | 4 |
| Set 2 | 6 |
| Set 3 | 5 |

Table 5.2: Number of participants in each context set

### 5.2.3 Findings

We present the result of the Kappa evaluation below. Figure 5.3 indicates the rating of the users based on each context of a context set.

| Evaluation | Yes | No |
|---|---|---|
| **Context Set: 1 (4 Participants)** | | |
| Context 1: | 1 | 3 |
| Context 2: | 2 | 2 |
| Context 3: | 4 | 0 |
| Context 4: | 3 | 1 |
| **Context Set: 2 (6 Participants)** | | |
| Context 1: | 4 | 2 |
| Context 2: | 1 | 5 |
| Context 3: | 2 | 4 |
| Context 4: | 0 | 6 |
| **Context Set 3 (5 Participants)** | | |
| Context 1: | 1 | 4 |
| Context 2: | 3 | 2 |
| Context 3: | 0 | 5 |
| Context 4: | 2 | 3 |

Figure 5.3: Rating of the users based on the presented results in two categories; Yes indicates semantically similar result, No indicates results which are not semantically similar as decided by the human participants.

The overall agreement and the Kappa value is indicated in Figure 5.4. The average overall agreement obtained is 0.61 which indicates that the human participants for evaluation shares decisions quite often. The average Kappa value obtained is 0.22 which seems to be 'Fair' based on the benchmark outlined in Figure 5.1. These results indicate that the data presented to the users were coherent and it also signifies that the evaluation was not influenced by chance.

| Evaluation | Overall Agreement | Kappa |
|---|---|---|
| **Context Set: 1  (4 Participants)** | 0.583 | 0.166 |
| **Context Set: 2  (6 Participants)** | 0.650 | 0.300 |
| **Context Set 3  (5 Participants)** | 0.600 | 0.200 |

Figure 5.4: The overall agreement score and the Kappa value for each of the context set.

The evaluation for the success rate, search time and query size was measured using a system with a CPU 2.20GHz and 8GB RAM.

We calculate the success rate of each event during evaluation and measure the search time and the number of queries (Query size) required to present a result to the users. We find that in most cases the system was able to generate a result for the users except in cases when the questions formed by the users was not in the form of Noun-Verb-Noun tags representation which is required for the generation of the triples. Since, no triples were generated for such questions the system could not proceed and was not able to present data to the users. The mean search time, if all contexts are to be taken into account lies between 37.5 seconds to 214 seconds, which indicates that search time is dependent on the context and the nature of questions being formed by the users. If the word length of the question being formed is large and the semantic triples generated due to the large length of the questions is much more than for small sentences, the search time increases significantly.

There are five levels of search in our approach and the query size is representative of how many levels of validations were required to present the users with a result. For example, if the system is able to present the users with a result after search in Level 1; the query size is 1. If the system is able to present the users with a result after search in Level 5, the query size is 5. Even after, Level 5, if the system does not present the users with a result and moves on to the next set of triples and starts the search at Level 1; the query size in this case is 6( 5 queries till Level 5 of the previous search attempt and 1 query for Level 1 search for the current set of triples). Figure 5.5 represent the quantitative results of the success rate, search and query time for each context.

| Measures | Context Sets | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Context Set 1 | | | | Context Set 2 | | | | Context 3 | | | |
| | C1 | C2 | C3 | C4 | C1 | C2 | C3 | C4 | C1 | C2 | C3 | C4 |
| Success Rate | 3/4 | 4/4 | 4/4 | 4/4 | 6/6 | 6/6 | 6/6 | 5/6 | 5/5 | 5/5 | 5/5 | 5/5 |
| Search Time(Max) | 45s | 35s | 92s | 124s | 94s | 89s | 78s | 228s | 67s | 83s | 240s | 92s |
| Search Time(Min) | 30s | 21s | 70s | 88s | 77s | 63s | 42s | 192s | 43s | 72s | 188s | 74s |
| Search Time(Mean) | 37.5s | 28s | 81s | 106s | 85.5s | 76s | 60s | 210s | 55s | 77.5 s | 214s | 83s |
| Query Size(Max) | 6 | 4 | 8 | 15 | 10 | 15 | 10 | 40 | 12 | 15 | 45 | 9 |
| Query Size(Min) | 4 | 2 | 6 | 10 | 6 | 12 | 9 | 30 | 10 | 12 | 40 | 7 |
| Query Size(Mean) | 5 | 3 | 7 | 12 | 8 | 16 | 6 | 35 | 11 | 13 | 42 | 8 |

Figure 5.5: Success rate, Search Time and the Number of Queries used for each context .

Finally, the precision of the system has been calculated based on the semantic similarity judgment by the users. The 'Yes' response from the users on whether the recommended question by the system is semantically similar is considered as the 'True Positive' values; a question which is semantically similar with respect to the asked question and presented to the user. 'No' response from the users are considered as the 'False Positive' value ; a question which is not semantically similar yet recommended to the users. Figure 5.6 shows the detailed calculation for the precision of the system based on each context.

| Evaluation | Yes (TP) | No (FP) |
|---|---|---|
| **Context Set: 1  (4 Participants)** | | |
| Context 1: | 1 | 3 |
| Context 2: | 2 | 2 |
| Context 3: | 4 | 0 |
| Context 4: | 3 | 1 |
| Sum: | 10 | 6 |
| Precision: | 10/(10+6)*100=62.5% | |
| **Context Set: 2  (6 Participants)** | | |
| Context 1: | 4 | 2 |
| Context 2: | 1 | 5 |
| Context 3: | 2 | 4 |
| Context 4: | 0 | 6 |
| Sum: | 7 | 17 |
| Precision: | 7/(7+17)*100 = 29.16% | |
| **Context Set 3  (5 Participants)** | | |
| Context 1: | 1 | 4 |
| Context 2: | 3 | 2 |
| Context 3: | 0 | 5 |
| Context 4: | 2 | 3 |
| Sum: | 6 | 14 |
| Precision: | 6/(6+14)*100 = 30% | |

Figure 5.6: Precision of the system based on each context set.

49

The kind of context, the nature of user formed questions, the semantic triples generated from the user formed questions and the knowledge in the triple store are all contributing factors in determining precision of the system. Context set 1 has more precision as the nature of the user formed question generated similar semantic triples most of the times which were already present in the knowledge base. Just by having more number of similar questions in the knowledge base is not a clear indicator of good precision in our case, existence of similar semantic triples tagged to the appropriate questions and such triples getting generated during the processing of natural language questions encourages good precision. Our observations from the study suggests that if the user formed question is as such that the NLP tools employed to generate semantic triples are similar to the already stored semantic triples in the knowledge base, the precision of the recommended results increases significantly.

Observations from our evaluation indicates that our approach is human understandable. Given a choice of more words in a context for forming questions, the variance of the questions formed by the user increased. It was observed that the participants spent more time in forming questions when presented with more options (words in a context). Few participants formed multiple questions from the same context easily when presented with more options. Restricting the words in a context gave less freedom to the participants to form questions and the questions they formed usually generated the semantic triples which were present in the knowledge base. In situations when the system could not generate semantic triples for the user posed questions, the users were asked to form variants of the question that they posed and it occurred that such formation of new variants of questions were not able to generate a positive recommendation.

## 5.3 Critical Analysis

It is expensive and unfeasible in terms of time and effort required by human participants to evaluate a large number of instances. Moreover, a fixed domain facilitates the recruitment of better human participants, since, in our proposed approach we have made an open domain attempt, pre-defined domain specific rules as well as evaluation by human-experts in specific domains is absent. Strengths and weaknesses of the system from several perspectives were examined and the experimental design and the results represented in this study can serve as a useful guideline for further studies on the same topic of identifying semantically similar question. However, some aspects which could be influencing the agreement or disagreement among participants are as follows:

- The abstract context presented to the participants and based on the understanding of the context or the core idea, variance in the formulation of the question by the participants and thereby influencing the judgement whether such a question is semantically similar against the recommended questions.

- The reliance on Natural language processing (NLP) tools for relationship extraction and the interpretation of the results from such tools by the human participants was an influencing factor. For example, a series of semantic triples generated from a given question incited confusion and influenced their decision of marking a specified case as negative.

- Traditional knowledge of Information Retrieval concepts such as n-grams and stop words influenced the decision of the judgement. For example, prepositions in our approach is used as a predicate to establish a relationship between a subject and an object, however, it is considered as good candidates for stop words in Information Retrieval, these unfamiliar understanding by the participants impacted the evaluation.

Even though the count of participants were less and there were issues in judgment by the participants our proposed approach to identify semantically similar questions by using NLP tools and Linked Data principles holds significant promise.

## 5.4  Summary

The outcome of the evaluation indicates that identifying semantically similar unstructured text with the help of linked data approach is a possibility. If both domain-specific or domain independent data modelling is done based on the given problem statement, this approach is extensible and maintainable. There are various limitations of this approach, however, the primary user centered evaluation confirms that this approach is human understandable and with the help of knowledge experts could be fine-tuned in the event of semantic drift within the knowledge base making the system extendable and relevant. We believe based on the pilot methodology and the preliminary results obtained from our evaluation, the proposed approach is viable. This approach could with multiple layers of analysis-each layer focusing on enriching the knowledge base and recommending semantically similar questions holds promise. It facilitates not only identifying semantically similar questions from unstructured text but also in creation of a reusable, extendable and human-understandable knowledge base which could be reasoned upon and used for gaining insights.

# Chapter 6

# Conclusion

The study conducted to answer the research question as mentioned in section 1.3 yielded interesting and insightful findings. In this chapter we discuss the possibilities and limitations of this approach, recommendations for future work and finally we summarize stating how successful was this dissertation.

## 6.1  Possibilities

In the domain of question answering and with respect to the specific case of our work of identifying similar questions, linked data approach has a great potential and has shown significant promise during the course of our research. It is the opinion of the author that with additional refinements, investigation and experimentation of this approach a web-scale application can be realized which would aid in dynamic creation and enrichment of knowledge base which could be exploited to gain reasonable insights for general and specific purposes. This approach has sufficient potential in the area of knowledge harvesting using a structured approach. The growing advancement in the area of NLP techniques for relationship extraction and the infusion of results from such techniques to model data in the format as specified by W3C standards thereby using such data to gain information cannot be ignored or undermined. This approach also facilitates the avenue of data monitoring which allows in maintaining consistent standard of data quality.

## 6.2  Limitations

We would like to discuss on two types of limitations, primarily the engineering limitations which can be fine-tuned with more time and effort and the specific limitations of this approach. In our current approach the generation of triples and the corresponding synonyms generated from those

triples is highly dependent on the NLP tools and the WordNet API that has been used to gather data for construction of the knowledge base. Research related to the quality of triples generation and comparison between various state-of-the art techniques to choose the best technique or customize a current technique would help in increasing the quality of the generated triples. Our current approach does not focus on the ranking of triples based on the semantic similarity of the sentence from which such triples are generated. This is an engineering limitation and can be solved with state-of-the art algorithms for classification and clustering. Another important feature of spelling correction before using generated triples to model a knowledge base so as to decrease the amount of noise in the knowledge base could be introduced as such a feature does not exist within the scope of our work. A more thorough analysis of WordNet and to retrieve only meaningful synonyms and the related meronym, hyponym and thereby building a customized synonym engine would facilitate in creating robust building blocks for the knowledge base. Monitoring, maintaining and updating the triple store is required by the knowledge experts from time to time.

## 6.3 Future Work

Limitations of the present approach and the possible solution has been discussed in the previous section. In this section we would like to provide few recommendations for future work. There are various options by which the present work could be expanded based on the current findings. The design of the knowledge base from a technical perspective could be further developed as a five star linked data application which would facilitate inter-linking of data as well as providing richer context for the users. Such design developments, would encourage in using more external vocabularies as well as publishing data on the Linked Open Data cloud. Given that the findings of the current study emphasizes on the different ways in which users forms questions depending on what context is shown to them. It would be interesting to investigate as well as to compare the contexts shown to the users (and the user posed questions based on these contexts) to derive richer semantics which could then be thoroughly analyzed to be included in the knowledge base. For future research, criteria for selection of participants having knowledge in the fields of cognitive sciences, human reasoning or people involved in the study of human behaviour and mental processes could be recruited as they would be able to provide richer variety of contexts. Inclusion of such participants would greatly help further research on this proposed structured approach of identifying semantically similar questions using linked data. In depth comparison between this approach and the state of the art techniques discussed in this dissertation in terms of user control and user scrutability could be studied extensively to aid in creation of benchmarks and gold standards for evaluation.

## 6.4 Summary

An end-to-end system for the identification of semantically similar natural language questions has been developed during the course of this research. The results of this approach were user evaluated and the findings demonstrates that such an approach is human as well as machine understandable. This approach also indicates that it is more dependent on the human-understandable parameters rather than just machine-understandable parameters.

The fact that data from CQA sites could be processed using NLP techniques, modelled in a structured format such as RDF and used for querying information using SPARQL is a significant accomplishment of this work. It also provides a meaningful contribution in the area of community question answering using linked data.

The purpose of this study was to explore the possibilities of identifying semantically similar questions using NLP techniques and linked data principles. By creating an end-to-end system for identification of semantically similar natural language questions and thereby evaluating such a system involving primarily human participants, it was envisaged that such a system is viable and human understandable (although abstractly at the present moment). The findings show that although there were differences in the opinions of a whether a question being semantically similar between participants, the system expressed results which showed signs of interest in the human participants. The possibility of such an approach is a testament to the positive response by user participants on the structured approach versus unstructured and non-human understandable approaches of doing the same task by other state-of-the-art techniques.

# Appendix A

# Appendix

## A.1 User Evaluation Data

The below images represents the raw data of user evaluation conducted.

**Context Set 1:**

1) You      do_bake      chocolate_cake
   You      do_bake_chocolate_cake_without      egg
   You      can_bake_chocolate_cake_without      egg
   You      can_bake      chocolate_cake

2) Someone_drive      in      busy_road
   Busy_road      in      Dublin

3) I      can_make      paper_boat
   I      can_make_paper_boat_for      kid

4) You      do_find      peace
   You      be_in      when_India
   You      be_in      India
   You      do_find_peace_in      India
   You      do_find_peace_in      when_India

**Context Set 2:**

1) Many_best_player        in        globe
   You     can_tell        many_best_player_in_globe
   You     can_tell        how_many_best_player_in_globe


2) You     do_stop_distraction     study
   You     do_stop        distraction


3) You     can_stop_distraction    study
   You     can_stop        distraction


4) Human_being        study_for        exam
   Being        study_for        exam
   Human_being        do_stop_distraction     study_for_exam
   Being        do_stop_distraction     study
   Being        do_stop_distraction     study_for_exam


**Context Set 3:**

1) You        would_prefer_tea_in        morning
   You        would_prefer        tea
   You        would_prefer_coffee_in        morning
   You        would_prefer        coffee


2) I        can_find        information
   I        can_find        information_about_other_program
   I        can_find        information_about_program
   I        can_find        information_about_other_graduate_program
   I        can_find        information_about_graduate_program


3) You        will_eradicate_poverty_from        you_country
   You        will_eradicate        poverty


4) Someone     can_bring_stability_to        life
   Someone     can_bring        stability
   Someone     can_bring_stability_to        chaotic_life

Content Set 1 :

Context 1 :

Participant 1 : How can you bake a chocolate cake without eggs?

Recommended : Can you bake a chocolate cake without eggs.

Yes/No    :    No . [ Present in the knowledge base - user question]
            x

Participant 2 : Can chocolate cakes be baked without eggs ?

Recommended : Triples generated; no knowledge in the knowledge base ;
            nothing recommended :

Yes/No  :   No → knowledge stored in triple store .
            x

Participant 3 : I don't have eggs but can I still bake a
            chocolate cake ?

Recommended : How can you bake a chocolate cake without eggs?

Yes/No :   Yes

Participant 4 : Can one bake chocolate cakes without eggs ?

Recommended : Triples not generated; therefore no questions were
            recommended . x

---

Content Set 1

Context 2 :   Someone - drive in busy road ; busy - road in Dublin

Participant 1 : How can someone drive in a busy road in Dublin ? ,

Recommended :  Can someone drive in a busy road in Dublin ? ,

Yes/No :   Yes . ✓

Participant 2 : How busy are the roads in Dublin for driving ?

Recommended : Is it true that the roads in Dublin are too busy
                                    to drive ?

Yes/No :   No . x

Participant 3 : Is it possible to drive in a busy road in Dublin ?

Yes/No :   No .   x

Recommended : Is it possible to be an outgoing introvert .
                ( generated from triple . it be possible)

Participant 4 : Are the roads in Dublin too busy for someone
            to drive .

Recommended :  Is it true that the roads in Dublin are too
            busy to drive .

Yes/No :   Yes . ✓

---

Content Set 1 :

Context 3 :

Participant 1 : How can I make paper boat for kids ?

Recommended : How can I make paper boats for kids?

Yes/No :   Yes ✓

Participant 2 : Can I make a paper boat for kids, if so, how?

Recommended : How can I make paper boat for kids

Yes/No  :   Yes ✓

Participant 3 : How can I make paper boats for kids

Recommended : How can I make paper boats for kids .

Yes/No :   Yes ✓

Participant 4 : For kids, how I can make a paper boat ?

Recommended : How can I make paper boats for kids .

Yes/No :   Yes ✓

---

Content Set 1 :

Context 4 :

Participant 1 : When in India where will you go to find peace ?

Recommended: Where do you find peace when you are in India?

Yes/No  :   Yes . ✓

Participant 2 : Where will you go to find peace when you are
            in India ?

Recommended : When in India where will you go to find peace

Yes/No  :   Yes . ✓

Participant 3 : Do you find peace in India ?

Recommended : Where will you go to find peace when you are
            in India ?

Yes/No :   No . x

Participant 4 : Where will you try finding peace when you are
            in India?

Recommended: Where will you go to find peace when you are
            in India ?

Yes/No :   Yes ✓

Content Set 2:
Content 1:

Participant 1 : How many best players are there in the globe?
Recommended: Can you tell me how many best players are there in the globe?
Yes/No :  Yes. ✓

Participant 2 : Who are the best players in tennis around the globe?
Recommended : 1st instance - no knowledge in the ontology; knowledge updated - variant question not accepted.
Yes/No :  No : ✗

Participant 3 : How many best players are in the world?
Recommended: Can you tell me how many best players are there in the globe?
Yes/No :  Yes.

Participant 4 : Who is the best player in the globe?
Recommended : Can you tell me how many best players are there in the globe.
Yes/No :  No : ✗

Participant 5 : How many best players in globe?
Recommended : Can you tell me how many best players are in globe.
Yes/No :  Yes. ✓

Participant 6 : In this world, how many best players are there?
Recommended : Can you tell me how many best players are in globe.
Yes/No :  Yes.

Content Set 2:
Content 2:

Participant 1 : How do you stop distractions while studying?
Recommended : How on earth do you stop distractions?
Yes/No :  No. ✗

Participant 2 : How can you stop distractions during study?
Recommended: How can you stop distractions while studying?
Yes/No :  Yes ✓

Participant 3 : Do distractions stop studying?
Recommended : Triples generated but no questions recommended ✗

Participant 4 : How do you stop getting distracted while studying?
Recommended : What does it mean if you have a dream getting a snake bite.
Yes/No :  No (Irrelevant recommendation!) ✗

Participant 5 : Can you stop distractions for me?
Recommended : How can you stop distractions while studying?
Yes/No :  No. ✗

Participant 6 : Can distractions be stopped entirely?
Recommended : Triples generated but no question has been recommended. ✗

Asked variants of questions for not recommended question. Results retrieved were not accepted as semantically similar by participants; therefore marked as 'No'.

Content Set 2:
Content 3:

Participant 1 : How can you stop distractions during study?
Recommended : How can you stop distractions while studying.
Yes/No :  Yes ✓

Participant 2 : Is it possible to distract a student?
Recommended : Is it possible to be an outgoing introvert?
Yes/No :  No ✗

Participant 3 : How distractions can stop study?
Recommended : How can you stop distractions while studying.
Yes/No :  No ✗

Participant 4 : How can you stop distractions in study?
Recommended : How can you stop distractions while studying.
Yes/No :  Yes ✓

Participant 5 : Can you completely stop distractions?
Recommended : How can you stop distractions while studying.
Yes/No :  No. ✗

Participant 6 : Can I stop distractions, like you do during studying?
Recommended : How do you get noticed
Yes/No :  No ✗

Content Set 2:
Content 4:

Participant 1 : Can human beings study?
Recommended : Unable to generate triples.
— Variant questions asked, answer retrieved. None not accepted to be semantically similar; therefore case is marked as 'No' ✗

Participant 2 : Why do human beings get distracted while studying for exams?
Recommended : How do human beings stop distractions while studying for exams?
Yes/No —  No ✗

Participant 3 : How distractions stope human beings from studying for exams?
Recommended : How do human beings stop distractions while studying for exams?
Yes/No —  No ✗

Participant 4 : Why does human beings need to study for the exam?
Recommended — Why do you need to study for exam?
Yes/No —  No ✗

Participant 5 : How can people stop getting distracted while studying for an exam?
Recommended — How can you stop distraction during study.
Yes/No —  No ✗

Participant 6 : Can human beings study for exams?
Recommended — How can human beings study for exams?
Yes/No —  No ✗

## Content set 3: Content 1

Participant 1 : In the morning, what would you prefer, tea or coffee?

Recommended : Which drink would you prefer tea or coffee in the morning?

Yes/No : Yes.

Participant 2 : What is your preferred drink in the morning?

Recommended : Triples generated; but no question was recommended.

Participant 3 : Is tea better than coffee in the morning?

Recommended : Triples generated; but no question was recommended; as such triples didn't exist in the knowledge base.

Participant 4 : Would you prefer coffee in the morning?

Recommended : Which drink would you prefer tea or coffee in the morning?

Yes/No : No — (Rationale given by the participant; the information of 'tea' was not contained in the question asked!)

Participant 5 : Which beverages would you like to have in the morning?

Recommended : Triples generated; but no question was recommended; as the generated triples were not present in the knowledge base.

Variants of questions asked for questions for which no-triples were generated; marked as 'No' by evaluators

## Content set 3: Content 2

Participant 1 : How can I find information about graduate programs?

Recommended : Is there a place from where I can find information about graduate programs?

Yes/No : Yes. ✓

Participant 2 : Can you find information about other graduate programs?

Recommended : Where can I find information about other graduate programs?

Yes/No : Yes. ✓

Participant 3 : Can I find information about graduate programs?

Recommended : Is there a place from where I can find information about graduate programs?

Yes/No : Yes ✓

Participant 4 : Can I find information about other programs?

Recommended : Can I find information about graduate programs?

Yes/No : No ✗

Participant 5 : If I want to know about graduate programs where should I be looking for it?

Recommended : I want to know reviews about point group Barcelona internship.

Yes/No : No ✗

## Content set 3: Content 3

Participant 1 : Will you be able to eradicate poverty from your country?

Recommended : Is it possible to eradicate poverty, if so how?

Yes/No : No. ✗

Participant 2 : How will you eradicate poverty?

Recommended : How will you eradicate poverty from your country?

Yes/No : No. ✗

Participant 3 : What are the ways in which you can eradicate poverty from your country?

Yes/No : Triples generated but no questions retrieved ✗

Participant 4 : Can you eradicate poverty from your country?

Recommended : What are the ways in which you can eradicate poverty from your country?

Yes/No : No. ✗

Participant 5 : How can you eradicate poverty?

Recommended : What are the ways in which you can eradicate poverty from your country?

Yes/No : No. ✗

## Content set 3: Content 4

Participant 1 : What will you suggest someone to bring stability to a chaotic life?

Recommended : How can someone help bring stability to a chaotic life?

Yes/No : Yes. ✓

Participant 2 : How often does someone attempt to bring stability to a chaotic life?

Recommended : How can someone help in bringing stability to a chaotic life?

Yes/No : No ✗ [Persisted in knowledge base - user question]

Participant 3 : Is it true that someone can bring stability to a chaotic life?

Recommended : How can someone help in bringing stability to a chaotic life?

Yes/No : No ✗ [Persisted in knowledge base - user question]

Participant 4 : How to bring stability to someone's chaotic life.

Recommended : Can someone bring stability to someone's chaotic life?

Yes/No : No. ✗

Participant 5 : Can someone try to bring stability to a chaotic life?

Recommended : How often does someone attempt to bring stability to a chaotic life?

Yes/No : Yes ✓ [Persisted by a user question previously]

# Appendix B

# Appendix

## B.1 Definition of tags

Description of some of the Treebank Constituent Tags [1].

**Clause Level**

**S** - simple declarative clause, i.e. one that is not introduced by a (possible empty) subordinating conjunction or a *wh*-word and that does not exhibit subject-verb inversion.
**SBAR** - Clause introduced by a (possibly empty) subordinating conjunction.
**SBARQ** - Direct question introduced by a *wh*-word or a *wh*-phrase. Indirect questions and relative clauses should be bracketed as SBAR, not SBARQ.
**SINV** - Inverted declarative sentence, i.e. one in which the subject follows the tensed verb or modal.
**SQ** - Inverted yes/no question, or main clause of a *wh*-question, following the *wh*-phrase in SBARQ.

**Phrase Level**

**ADJP** - Adjective Phrase.
**ADVP** - Adverb Phrase.
**CONJP** - Conjunction Phrase.
**FRAG** - Fragment.
**INTJ** - Interjection. Corresponds approximately to the part-of-speech tag UH.
**LST** - List marker. Includes surrounding punctuation.
**NAC** - Not a Constituent; used to show the scope of certain prenominal modifiers within an NP.
**NP** - Noun Phrase.
**NX** - Used within certain complex NPs to mark the head of the NP. Corresponds very roughly to N-bar level but used quite differently.
**PP** - Prepositional Phrase.
**PRN** - Parenthetical.
**PRT** - Particle. Category for words that should be tagged RP.
**QP** - Quantifier Phrase (i.e. complex measure/amount phrase); used within NP.
**RRC** - Reduced Relative Clause.
**UCP** - Unlike Coordinated Phrase.
**VP** - Vereb Phrase.
**WHADJP** - *Wh*-adjective Phrase. Adjectival phrase containing a *wh*-adverb, as in *how hot*.
**WHAVP** - *Wh*-adverb Phrase. Introduces a clause with an NP gap. May be null (containing the 0 complementizer) or lexical, containing a *wh*-adverb such as *how* or *why*.
**WHNP** - *Wh*-noun Phrase. Introduces a clause with an NP gap. May be null (containing the 0 complementizer) or lexical, containing some *wh*-word, e.g. *who, which book, whose daughter, none of which*, or *how many leopards*.
**WHPP** - *Wh*-prepositional Phrase. Prepositional phrase containing a *wh*-noun phrase (such as *of which* or *by whose authority*) that either introduces a PP gap or is contained by a WHNP.
**X** - Unknown, uncertain, or unbracketable. X is often used for bracketing typos and in bracketing *the...the*-constructions.

---

[1]http://www.surdeanu.info/mihai/teaching/ista555-fall13/readings/PennTreebankConstituents.html

## Word level

**CC** - Coordinating conjunction
**CD** - Cardinal number
**DT** - Determiner
**EX** - Existential there
**FW** - Foreign word
**IN** - Preposition or subordinating conjunction
**JJ** - Adjective
**JJR** - Adjective, comparative
**JJS** - Adjective, superlative
**LS** - List item marker
**MD** - Modal
**NN** - Noun, singular or mass
**NNS** - Noun, plural
**NNP** - Proper noun, singular
**NNPS** - Proper noun, plural
**PDT** - Predeterminer
**POS** - Possessive ending
**PRP** - Personal pronoun
**PRP$** - Possessive pronoun (prolog version PRP-S)
**RB** - Adverb
**RBR** - Adverb, comparative
**RBS** - Adverb, superlative
**RP** - Particle
**SYM** - Symbol
**TO** - to
**UH** - Interjection
**VB** - Verb, base form
**VBD** - Verb, past tense
**VBG** - Verb, gerund or present participle
**VBN** - Verb, past participle
**VBP** - Verb, non-3rd person singular present
**VBZ** - Verb, 3rd person singular present
**WDT** - Wh-determiner
**WP** - Wh-pronoun
**WP$** - Possessive wh-pronoun (prolog version WP-S)
**WRB** - Wh-adverb

# Bibliography

[1] Y. Zhu, M. C. Kim, and E. Yan, "Evaluating interactive bibliographic information retrieval systems: A user-centered approach," *Proceedings of the Association for Information Science and Technology*, vol. 55, no. 1, pp. 628–637, 2018.

[2] C. Shah, S. Oh, and J. S. Oh, "Research agenda for social qa," *Library Information Science Research*, vol. 31, no. 4, pp. 205 – 209, 2009.

[3] K. Bae and Y. Ko, "Improving question retrieval in community question answering service using dependency relations and question classification," *Journal of the Association for Information Science and Technology*, 2019.

[4] Y. Liu, S. Li, Y. Cao, C.-Y. Lin, D. Han, and Y. Yu, "Understanding and summarizing answers in community-based question answering services," in *Proceedings of the 22Nd International Conference on Computational Linguistics - Volume 1*, COLING '08, (Stroudsburg, PA, USA), pp. 497–504, Association for Computational Linguistics, 2008.

[5] G. Wang, K. Gill, M. Mohanlal, H. Zheng, and B. Y. Zhao, "Wisdom in the social crowd: an analysis of quora," in *WWW*, 2013.

[6] M. R. Martin and S. Auer, "Categorisation of semantic web applications," 2010.

[7] T. Wei, Y. Lu, H. Chang, Q. Zhou, and X. Bao, "A semantic approach for text clustering using wordnet and lexical chains," *Expert Systems with Applications*, vol. 42, no. 4, pp. 2264 – 2275, 2015.

[8] D. Rusu, L. Dali, B. Fortuna, M. Grobelnik, and D. Mladení, "Triplet extraction from sentences," 2007.

[9] R. Glauber and D. B. Claro, "A systematic mapping study on open information extraction," *Expert Syst. Appl.*, vol. 112, pp. 372–387, 2018.

[10] J. R. Landis and G. G. Koch, "The measurement of observer agreement for categorical data," *Biometrics*, vol. 33, pp. 159–174, 1977.

[11] L. A. Adamic, J. Zhang, E. Bakshy, and M. S. Ackerman, "Knowledge sharing and yahoo answers: Everyone knows something," in *Proceedings of the 17th International Conference on World Wide Web*, WWW '08, (New York, NY, USA), pp. 665–674, ACM, 2008.

[12] A. Baltadzhieva and G. Chrupa, "Question quality in community question answering forums: A survey," *SIGKDD Explor. Newsl.*, vol. 17, pp. 8–13, Sept. 2015.

[13] J. Jeon, W. B. Croft, J. H. Lee, and S. Park, "A framework to predict the quality of answers with non-textual features," in *SIGIR*, 2006.

[14] X. Xue, J. Jeon, and W. B. Croft, "Retrieval models for question and answer archives," in *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '08, (New York, NY, USA), pp. 475–482, ACM, 2008.

[15] E. Agichtein, C. Castillo, D. Donato, A. Gionis, and G. Mishne, "Finding high-quality content in social media," in *Proceedings of the 2008 International Conference on Web Search and Data Mining*, WSDM '08, (New York, NY, USA), pp. 183–194, ACM, 2008.

[16] A. Shtok, G. Dror, Y. Maarek, and I. Szpektor, "Learning from the past: Answering new questions with past answers," in *Proceedings of the 21st International Conference on World Wide Web*, WWW '12, (New York, NY, USA), pp. 759–768, ACM, 2012.

[17] B. Li and I. King, "Routing questions to appropriate answerers in community question answering services," in *Proceedings of the 19th ACM International Conference on Information and Knowledge Management*, CIKM '10, (New York, NY, USA), pp. 1585–1588, ACM, 2010.

[18] L. Chen, "Understanding and exploiting user intent in community question answering," in *PhD thesis, Birkbeck, University of London*, 2014.

[19] M. Asaduzzaman, A. S. Mashiyat, C. K. Roy, and K. A. Schneider, "Answering questions about unanswered questions of stack overflow," in *2013 10th Working Conference on Mining Software Repositories (MSR)*, pp. 97–100, May 2013.

[20] P. Jurczyk and E. Agichtein, "Discovering authorities in question answer communities by using link analysis," in *Proceedings of the Sixteenth ACM Conference on Conference on Information and Knowledge Management*, CIKM '07, (New York, NY, USA), pp. 919–922, ACM, 2007.

[21] M. Qu, G. Qiu, X. He, C. Zhang, H. Wu, J. Bu, and C. Chen, "Probabilistic question recommendation for question answering communities," in *Proceedings of the 18th International Conference on World Wide Web*, WWW '09, (New York, NY, USA), pp. 1229–1230, ACM, 2009.

[22] D. Horowitz and S. D. Kamvar, "The anatomy of a large-scale social search engine," in *Proceedings of the 19th International Conference on World Wide Web*, WWW '10, (New York, NY, USA), pp. 431–440, ACM, 2010.

[23] E. Mendes Rodrigues and N. Milic-Frayling, "Socializing or knowledge sharing?: Characterizing social intent in community question answering," in *Proceedings of the 18th ACM Conference on Information and Knowledge Management*, CIKM '09, (New York, NY, USA), pp. 1127–1136, ACM, 2009.

[24] T. Berners-Lee, J. Hendler, and O. Lassila, "The semantic web," *Scientific American*, pp. 34–43, 2001.

[25] V. Rodríguez-Doncel, C. Santos, P. Casanovas, and A. Gómez-Pérez, "Legal aspects of linked data – the european framework," *Computer Law Security Review*, vol. 32, no. 6, pp. 799 – 813, 2016.

[26] C. Bizer, T. Heath, and T. Berners-Lee, "Linked data - the story so far," *Int. J. Semantic Web Inf. Syst.*, vol. 5, pp. 1–22, 2009.

[27] C. Unger, A. Freitas, and P. Cimiano, "An introduction to question answering over linked data," in *Reasoning Web*, 2014.

[28] D. A. Ferrucci, "Introduction to "this is watson"," *IBM Journal of Research and Development*, vol. 56, pp. 1:1–1:15, May 2012.

[29] J.-P. Ng and M.-Y. Kan, "Qanus: An open-source question-answering platform," *ArXiv*, vol. abs/1501.00311, 2015.

[30] V. Lopez, M. Pasin, and E. Motta, "Aqualog: An ontology-portable question answering system for the semantic web," in *The Semantic Web: Research and Applications* (A. Gómez-Pérez and J. Euzenat, eds.), (Berlin, Heidelberg), pp. 546–562, Springer Berlin Heidelberg, 2005.

[31] K. Wang, Z. Ming, and T.-S. Chua, "A syntactic tree matching approach to finding similar questions in community-based qa services," in *Proceedings of the 32Nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '09, (New York, NY, USA), pp. 187–194, ACM, 2009.

[32] Soojong Lim, Euisuk Jung, and Myoung-Gil Jang, "Dependency relation analysis using case frame for encyclopedia question-answering system," in *30th Annual Conference of IEEE Industrial Electronics Society, 2004. IECON 2004*, vol. 3, pp. 3120–3124 Vol. 3, Nov 2004.

[33] Robin D. Burke, Kristian J. Hammond, Vladimir Kulyukin, Steven L. Lytinen, Noriko Tomuro, and Scott Schoenberg, "Question answering from frequently asked question files: Experiences with the faq finder system," in *AI Magazine, Vol 18 No 2: Summer 1997*, AI Magazine, 1997.

[34] C. Fellbaum, "Wordnet: An electronic lexical database.," in *Cambridge, MA*, MIT Press, 1998.

[35] Y.-S. Lai, K.-A. Fung, and C.-H. Wu, "Faq mining via list detection," in *Proceedings of the 2002 Conference on Multilingual Summarization and Question Answering - Volume 19*, MultiSumQA '02, (Stroudsburg, PA, USA), pp. 1–7, Association for Computational Linguistics, 2002.

[36] A. J. Gill, L.-C. Umr, and R. M. French, "Level of representation and semantic distance: Rating author personality from texts," 2006.

[37] K. Zhang, W. Wu, H. Wu, Z. Li, and M. Zhou, "Question retrieval with high quality answers in community question answering," in *CIKM*, 2014.

[38] Z. Ji, F. Xu, B. Wang, and B. He, "Question-answer topic model for question retrieval in community question answering," in *CIKM*, 2012.

[39] J. D. Lafferty and C. Zhai, "Document language models, query models, and risk minimization for information retrieval," in *SIGIR*, 2001.

[40] J.-R. Wen, J.-Y. Nie, and H. Zhang, "Query clustering using user logs," *ACM Trans. Inf. Syst.*, vol. 20, pp. 59–81, 2002.

[41] A. Tombros, R. Villa, and C. J. van Rijsbergen, "The effectiveness of query-specific hierarchic clustering in information retrieval," *Inf. Process. Manage.*, vol. 38, pp. 559–582, 2002.

[42] G. Zhou, L. Cai, J. Zhao, and K. Liu, "Phrase-based translation model for question retrieval in community question answer archives," in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1*, HLT '11, (Stroudsburg, PA, USA), pp. 653–662, Association for Computational Linguistics, 2011.

[43] K. Liu, L. Xu, and J. Zhao, "Opinion target extraction using word-based translation model," in *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, EMNLP-CoNLL '12, (Stroudsburg, PA, USA), pp. 1346–1356, Association for Computational Linguistics, 2012.

[44] J. Gao, X. He, and J.-Y. Nie, "Clickthrough-based translation models for web search: From word models to phrase models," in *Proceedings of the 19th ACM International Conference on Information and Knowledge Management*, CIKM '10, (New York, NY, USA), pp. 1139–1148, ACM, 2010.

[45] W. Zhang, Y. Li, and S. Wang, "Learning document representation via topic-enhanced lstm model," *Knowl.-Based Syst.*, vol. 174, pp. 194–204, 2019.

[46] M. Yang, W. Tu, Q. Qu, W. Zhou, Q. Liu, and J. Zhu, "Advanced community question answering by leveraging external knowledge and multi-task learning," *Knowl.-Based Syst.*, vol. 171, pp. 106–119, 2019.

[47] Y. Liu, J. Bian, and E. Agichtein, "Predicting information seeker satisfaction in community question answering," in *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '08, (New York, NY, USA), pp. 483–490, ACM, 2008.

[48] A. Mockus and J. D. Herbsleb, "Expertise browser: A quantitative approach to identifying expertise," in *In 2002 International Conference on Software Engineering*, pp. 503–512, ACM Press, 2002.

[49] J. Guo, S. Xu, S. Bao, and Y. Yu, "Tapping on the potential of q&
a community by recommending answer providers," in *Proceedings of the 17th ACM Conference on Information and Knowledge Management*, CIKM '08, (New York, NY, USA), pp. 921–930, ACM, 2008.

[50] M. Liu, Y. Liu, and Q. Yang, "Predicting best answerers for new questions in community question answering," in *Proceedings of the 11th International Conference on Web-age Information Management*, WAIM'10, (Berlin, Heidelberg), pp. 127–138, Springer-Verlag, 2010.

[51] G. Zhou, L. Qian, and J. Fan, "Tree kernel-based semantic relation extraction with rich syntactic and semantic information," *Inf. Sci.*, vol. 180, pp. 1313–1325, 2010.

[52] O. Etzioni, M. Banko, S. Soderland, and D. S. Weld, "Open information extraction from the web," *Commun. ACM*, vol. 51, pp. 68–74, 2008.

[53] S. Brin, "Extracting patterns and relations from the world wide web," in *WebDB*, 1998.

[54] D. E. Rose and D. Levinson, "Understanding user goals in web search," in *WWW*, 2004.

[55] A. Broder, "A taxonomy of web search," *SIGIR Forum*, vol. 36, pp. 3–10, Sept. 2002.

[56] D. Moldovan, S. Harabagiu, M. Pasca, R. Mihalcea, R. Girju, R. Goodrum, and V. Rus, "The structure and performance of an open-domain question answering system," in *Proceedings of the 38th Annual Meeting on Association for Computational Linguistics*, ACL '00, (Stroudsburg, PA, USA), pp. 563–570, Association for Computational Linguistics, 2000.

[57] E. Hovy, L. Gerber, U. Hermjakob, C.-Y. Lin, and D. Ravichandran, "Toward semantics-based answer pinpointing," in *Proceedings of the First International Conference on Human Language Technology Research*, HLT '01, (Stroudsburg, PA, USA), pp. 1–7, Association for Computational Linguistics, 2001.

[58] S.-Y. Yang, C.-L. Hsu, D. liang Lee, and L. Y. Deng, "Faq-master: an ontological multi-agent system for web faq services," 2008.

[59] N. Dandekar, "Semantic question matching with deep learning," 2017.

[60] P. G. E. Hein, "Constructivist learning theory," *CECA (International Committee of Museum Educators) Conference*, 1991.

[61] D. Soni, R. L. Nord, and C. Hofmeister, "Software architecture in industrial applications," *1995 17th International Conference on Software Engineering*, pp. 196–196, 1995.

[62] G. Angeli, M. J. J. Premkumar, and C. D. Manning, "Leveraging linguistic structure for open domain information extraction," in *ACL*, 2015.

[63] L. Meng, R. Huang, and J. Gu, "An effective algorithm for semantic similarity metric of word pairs," *International Journal of Multimedia and Ubiquitous Engineering*, vol. 8, no. 2, pp. 1–12, 2013.

[64] S. Iyer, N. Dandekar, and K. Csernai, "First quora dataset release: Question pairs," in *Quora*, 2017.

[65] M. Fossati, E. Dorigatti, and C. Giuliano, "N-ary relation extraction for simultaneous t-box and a-box knowledge base augmentation," *Semantic Web*, vol. 9, pp. 413–439, 2017.

[66] N. Jaywant, S. Shetty, and V. Musale, "Digital identity based recommendation system using social media," in *2016 2nd International Conference on Applied and Theoretical Computing and Communication Technology (iCATccT)*, pp. 288–292, July 2016.

[67] S. Riedel, L. Yao, A. McCallum, and B. M. Marlin, "Relation extraction with matrix factorization and universal schemas," in *HLT-NAACL*, 2013.

[68] M. J. Hosseini, N. Chambers, S. Reddy, X. R. Holt, S. B. Cohen, M. Johnson, and M. Steedman, "Learning typed entailment graphs with global soft constraints," *Transactions of the Association for Computational Linguistics*, vol. 6, pp. 703–717, 2018.

[69] L. Romano, M. Kouylekov, I. Szpektor, I. Dagan, and A. Lavelli, "Investigating a generic paraphrase-based approach for relation extraction," in *EACL*, 2006.

[70] A. Carlson, J. Betteridge, B. Kisiel, B. Settles, E. R. Hruschka, and T. M. Mitchell, "Toward an architecture for never-ending language learning," in *AAAI*, 2010.

[71] C. D. Manning, M. Surdeanu, J. Bauer, J. Finkel, S. J. Bethard, and D. McClosky, "The Stanford CoreNLP natural language processing toolkit," in *Association for Computational Linguistics (ACL) System Demonstrations*, pp. 55–60, 2014.

[72] J. J. Randolph, "Free-marginal multirater kappa (multirater k[free]): An alternative to fleiss' fixed-marginal multirater kappa.," 2005.