# Explainable multi-hop dense question answering using knowledge bases and text

Editor(s): Name Surname, University, Country Solicited review(s): Name Surname, University, Country Open review(s): Name Surname, University, Country

Somayeh Asadifar<sup>a</sup>, Mohsen Kahani<sup>a,\*</sup> and Saeedeh Shekarpour<sup>b</sup> <sup>a</sup>Ferdowsi University of Mashhad, Iran <sup>b</sup>Department of Computer Science, University of Dayton, Dayton, Ohio

Abstract. Much research has been conducted extracting a response from either text sources or a knowledge base (KB). The challenge becomes more complicated when the goal is to answer a question with the help of both text and KB. In these hybrid systems, we address the following challenges: i) excessive growth of search space, ii) extraction of the answer from both KB and text, iii) extracting the path to reach to the answer, and vi) the scalability in terms of the volume of documents explored. A heterogeneous graph is utilized to tackle the first challenge guided by question decomposition. The second challenge is met with the usage of the idea behind an existing text-based method, and its customization for graph development. Based on this method for multi-hop questions, an approach is proposed for the extraction of answer explanation to address the third challenge. Since the basic method uses a dense vector for scalability, the final challenge is also addressed in the proposed hybrid method. Evaluation reveals that the proposed method has the ability to extract answers in an acceptable time and volume, while offering competitive accuracy and has created a trade-off between performance and accuracy in comparison with the base methods.

Keywords: Information Retrieval, Explainable Question Answering, Question Decomposition, Multi-hop Dense Retrieval, Knowledge-Based and Corpus.

<sup>\*</sup>Corresponding author. E-mail: kahani@um.ac.ir.

#### 1. Introduction

In the past, answering open-ended questions involved defining handwritten rules [1,2], then applying machine learning techniques [3–5], and, more recently, employing deep learning methods [4,6]. This domain includes two active research areas: knowledge-based (KB) and text-based.

KB-based models fall into two general categories: semantic parsing [7–10] and information retrieval [3,11,12]. Text-based systems include two basic modules for retrieving candidates and analysing them to extract the correct answer [13]. Improving the efficiency of candidate retrieval is a current active research. Some research performed this action with term-based TF-IDF and BM25 methods [14,15].

In addition, vector space methods have been proposed to solve the problems of the previous termbased approaches [16]. They are designed to increase the efficiency of Question Answering (QA) systems in dense retrieval methods [4,5]. A new trend is processing questions that require multi-step inference (multi-hop question) [4,13,17].

Scalability and interpretability are two inseparable features in multi-hop question processing that need special attention. Addressing these issues has been successful by using dense vectors in a text-based system for multi-hop questions [5].

In recent years, hybrid systems have been employed to extract answers from both text and KB information sources [11,12]. Some hybrid systems can handle multi-hop questions, as well. For instance, GRAFT-Nets [11] and its expansion, PullNet [12], extract answers by building a heterogeneous graph from text, fact, and entity nodes.

The present research is based on the PullNet approach. However, unlike PullNet, the proposed approach does not limit the extraction of the answer to the KB form, only. In addition, the current method can interpret the answer with a smaller graph, while providing competitive accuracy with PullNet, by considering sub-questions and their execution order.

The proposed method also relies on MDR [5] to extract the best sequence of responses, which could provide answer interpretation. Step-by-step search with the help of query decomposition and simultaneous use of KB and text leads to higher speed and accuracy compared to the MDR in extracting answers.

In the current study's scenario, one text and one large KB are available, neither is enough to answer the questions that require multiple-stage inferences from both sources. The answer can also be extracted from KB entities and text fragments, and finally, an explanation is provided on how the answer is extracted.

The main contributions of this research are:

- The ability to extract an answer from both text and KB sources, if available.
- Improving the efficiency in information retrieval as well as question answering with higher accuracy and speed by using question decomposition.
- Increasing the accuracy and speed in extracting answers by using text and KB, simultaneously.
- The ability to extract answer explanation.

The rest of the paper is organized as follows. Section 2 describes the research objectives and questions that are answered in this article. Section 3 reviews related works. The proposed approach is presented in Section 4. Sections 5 provides the experimental and empirical results. In the discussion section, the findings, theoretical and practical implications of the research are reviewed. Section 7 presents the conclusion and future works. Some sample questions and answers are provided in the appendix for further explanation.

#### 2. Research Objectives

Our main goal in this research is to provide a question answering system that can answer multi-hop questions using two sources of structured and unstructured (hybrid approach).

In this study, specifically, the research objective is to address the following research questions:

- Research question 1 (RQ1): What is the effect of question decomposition on the accuracy and speed of answer extraction? (§ 5.2, Table 4 and Table 5)
- Research question 2 (RQ2): What is the relationship between the number of sub-questions and the number of answer search steps? (§ 6.1, first finding)
- Research question 3 (RQ3): How to get the sequence of intermediate answers, including text and knowledge base triples? (§ 4.1.3)
- Research question 4 (RQ4): What is the effect of simultaneous use of text and knowledge base in extracting answers? (§ 5.2, Figure 2)

#### 3. Related Works

In recent years, as information exchange has become pervasive through interfaces such as World Wide Web, large volumes of information are generated daily and made publicly available to everyone. The most important challenge that arises with this volume of information is finding the information needed.

Information retrieval (IR) methods are used as the core of many real-world applications. The goal of an IR system is to find documents containing the answer to the query. The purpose of a QA system, on the other hand, is to provide answer (not just documents). Therefore, QA is closely related to other fields, such as natural language processing (NLP) and machine learning (ML).

In early researches, open-domain Question Answering (QA) were performed to extract the answers to a question, expressed in a natural language, by a combination of manual rules and machine learning models. Recently, research has been shifted to the use of deep neural network approaches [11,12,18,19].

Also, there is a new area, called Community Question Answering (CQA), where aimed at automatically extracting the answer to the information needs of members from previous posts in a particular community on the web [20]. However, this field is different from the objective of this paper, in which open domain search is employed using large heterogeneous structured and unstructured sources.

This field has a long history of being characterized by two parts: KB-based and text-based. A new field has emerged by combining these two models. In this section, the previous works in these areas are explored in some details.

#### 3.1. KB-based

The literature in this category is based on two general methods: semantic parsing and Information Retrieval-Based.

Semantic parsing methods: These methods are based on predefined patterns or rules for converting input questions into the logical forms. The limited number of patterns has restricted the ability to answer complex questions in these systems. To remedy this problem, Abujabal et al. have developed a semiautomatic pattern learning approach [7].

Recently, systems have tended to use neural networks to increase efficiency and scalability. The question in a natural language is first transformed into an intermediate representation, such as a tree or graph. Then, the intermediated representation is converted to a logical form. Efforts have been made to increase the efficiency in identifying entities [21], generating question graph [22], and processing multihop questions [8,10].

Efforts have been made to increase the efficiency of the components of identifying entities [21] and scoring the generated graphs [22], as well as to solve the problem of answering multi-hop questions [8,10].

Performing semantic analysis using encoderdecoder models is another method that has become common in recent years. These methods, which often use a Long Short-Term Memory (LSTM) network, differ in choosing a decoder, whether a tree (Seq-to-Tree) or a sequence (Seq-to-Seq) [23]. They generally ignore the structural information and interdependence between question words by capturing only the order of the words. This problem was addressed by combining the tree structure and sequence in a graph representation using Graph-to-Seq model [9].

The methods described above, although practical, require a large amount of training data, which is costly to generate. In order to solve this challenge, efforts were made to develop methods with weaker supervision, for example, using reinforcement learning [24].

Information Retrieval-Based Methods: In these methods, first, the desired entities (i.e., topic entity) in the question are captured, then the entities are linked to the knowledge base. In the next step, subgraphs containing the desired entities are selected from the KB. Nodes in the subgraphs other than question entities are considered as answers.

In early researches [25], the classification of syntactic features of the candidate's questions and answers were employed. This not only was timeconsuming but also did not include all semantic features, as these features were defined manually. To solve the problem, research in this field turned into representation learning of question. Questions and candidate answers were represented in the vector space, and later, neural networks are used for better representation.

In addition, topic entity, the path between topic and answer entities, context (KB subgraph containing the topic entity), and answer entity are often employed to represent the candidate response [26].

For instance, the cross-attention-based neural network model captures the correlation between questions and answers and uses the above four mentioned properties to encode candidate answers. It has reported acceptable performance [3] compared to the previous approaches.

A method based on learning graph representation is presented in [11,12]. These two works utilize the text body, in addition to the KB, to find the answer, and are classified in the hybrid category. In these methods, heterogeneous graphs are extracted from knowledge base entities and body texts. Learning is performed using a convolutional Neural Network (CNN), and a classifier determines the answer. The present study is inspired by the graph representation learning method in these articles [11,12].

The trend of knowledge-based information retrieval approaches has recently tended to process multihop questions that infer answers in two ways: using memory networks and walking the path. For the former, one can refer to the work proposed by Chen et al. [27], which uses a bidirectional attentive memory network that utilizes the correlation factor to improve the representation of the question.

For the second category, Qiu et al. look at the QA issue as a matter of sequential decisions [28]. In this research, a stepwise reasoning network has been created, which is trained using reinforcement learning.

Although these methods do not use predetermined patterns, they still have difficulty processing complex questions, and most methods are challenging in creating interpretation.

Recently, KB-based systems have provided solutions for processing complex questions that include multiple entities, relationships, and constraints, often in the form of a sequence of questions. These questions can be answered by breaking into simple questions, based on question syntax and predefined templates [7,8,29] or question semantics [2]. Recent developments and challenges in complex question answering have been exhaustingly surveyed in [26]. This research examines the methods available in answering complex questions in the context of the knowledge base.

#### 3.2. Text-based

Text-based systems answer questions from existing documents during the two main operations: passage retrieval and machine reading comprehension [13]. Passage retrieval task extracts the k-top relevant documents to the question by comparing the question and the document vectors using a distance measure. Passage retrieval is a branch of information retrieval that reduces the search space to extract answer. There are many researches in this area, trying to improve the retrieval models, so that the best candidates can be extracted to help find the answer.

Initial works used term-based methods as lexical adaptations of TF-IDF and BM25 [14,15]. In these methods, retrieval is based on the bag of words concept, and the ranking function is calculated based on the term and inverse document frequencies. To address the challenge of sparse vectors with high dimensions, considering the semantic property for embedding vectors through latent semantic analysis and the concept of dense retrieval [4,16,30] can improve the performance. Recently, efforts have been made to consider a small set of question-answer pairs to create vectors via dual encoders [4,5]. Also, the inner product ranking function between the question and the passage vectors is used [5,31]. An overview of text-based methods from the perspective of information retrieval and deep learning is presented in [32].

On the other hand, the processing of complex questions has been considered in many text-based systems [4,13,17,33], primarily when the answer to the question cannot be extracted with a single text piece or when multi-step inference must be performed on several text pieces to find the answer. Such questions are called multi-hop. Concerning multi-hop questions, scalability and the ability to extract the path to the answer are two essential factors.

Derived from the sequential nature of multi-hop question answering, MDR [5] uses an iterative process and maximum inner product search technique along with dense retrieval [4] to speed up the extraction of a sequence of passages from a large pool of documents. In the first step, the most similar passage to the question is extracted, and in the following steps, a new question is generated by combining the answers of the previous steps and the initial question. The newly generated question is used to compare and find the similar passages in each step.

Although, the present research to answer the question is not just text-based and belongs to the hybrid category, it employs the MDR method to solve the challenges of scalability, interpretability, and the ability to extract answers from the text.

#### 3.3. Hybrid

Although a large volume of information is in the text format, the extraction of an answer from a text has a lot of complexities, due to the diversity of manners of expressing information in natural languages. On the other hand, in a KB, information is expressed in specific structures that make it easy to extract. However, even the largest KBs suffer from information coverage problem. Therefore, it is evident that these two sources of information can complement each other regarding the coverage and simplicity of information extraction.

In general, systems that take the advantages of these two areas operate in two main model: early fusion and late fusion [11]. In the early fusion, both sources are searched simultaneously, while in the late fusion, each source is searched separately, and later the answers are fused. The former has been shown to perform better than the latter.

Some of the models presented in this category are either primarily text-based systems, extended to be able to use a KB [34] or, vice versa [35,36]. ODQA [37] converts all available resources to text and then uses retrieval and reading tools. A few researches have extracted an answer by simultaneously using text and a KB [11,12]. These systems do not support multi-hop questions.

The first attempt to simultaneously extract an answer from text and a KB was to use key-value memory alongside universal schema [38], though not considering the rich relationships between facts and textual parts. Another approach, GRAFT-Nets [11], retrieves a response-related subgraph by creating heterogeneous graphs of entities, facts, and text pieces instead of randomly extracting them. However, the generated graphs are often huge and not scalable.

By providing a way of learning to develop nodes, PullNet [12] has shown that it can gradually produce graphs, resulting in smaller graphs. PullNet assumes that the answer can be extracted if it exists in the KB form. In addition, the explanation of how to find the answer, which is necessary in the real world, cannot be extracted. Another problem is that PullNet does not generate optimal graphs, because it creates the initial graph with all the question entities and does not consider the sequential nature of the question requiring multi-step inference. In other words, Pull-Net ignores the relationships between question entities.

As the information containing the answer to the multi-hop question cannot be obtained in a single shot [5], the present research believes that not all the question words should be considered simultaneously, but should be added step by step in the search.

The system presented in the current research falls into the category of hybrid systems. This research aims to present a system that, like PullNet, processes multi-hop questions with the help of structured and unstructured sources. Moreover, it has the ability to scale up, extract response explanation, and extract the final answer from both sources, in such a way that it can control the search space to increase the efficiency of the system.

#### 4. Approach

As mentioned, the current research belongs to the complex (often called multi-hop) question-answering systems, trying to extract answers from both text and KB sources. It is based on three fundamental axes, described as follows.

- (i) A complex question can be defined as a question that has several relations, entities and may contain various constraints, e.g., temporal, spatial, aggregation, ordinal, etc. According to this definition, in most cases, a complex question can be converted into several minor questions (sub-questions) that should be executed in a specific order to obtain the correct answer. Decomposition of the question into sub-questions has led to the improvement of system performance.
- (ii) The objective is to find answers using both text and a KB sources. Thus, an effective way to obtain the answer to a given question is a graph, as it can clearly show the relations among contexts, entities, and RDF triples. The current work constructs the proposed approach based on the PullNet [12] method, extending the graph in several steps.
- (iii) The answer to multi-hop questions, usually can be searched in sequence by finding pieces of information. In this case, the proposed approach is built on the dense retrieval multi-hop system, MDR [5], which attempts to extract the best sequence from a pool of documents. However, in our approach, the pool of documents and facts is searched.

The present study provides a solution by integrating the three considered axes to balance the accuracy and the efficiency, while extracting the answer from both sources, considering the constraints stated in the question.

Throughout the article, triple and fact, as well as document and passage, are used interchangeably. In addition, the concept of the document and the passage are the same and are considered as a single sentence.

#### 4.1. Model

The architecture of the proposed model, GraphMDR, is based on four modules. This architecture is shown in Figure 1, using a sample question for better presentation. The parts in the gray box marked with a dashed line are iteratively processed. At each



Figure 1: The architecture of the proposed model, GraphMDR

stage of iteration, Question Graph Expansion (QGE) module and Sequence Retrieval (SR) module are examined for one of the sub-questions of Question Decomposition and Constraint Recognition (QD&CR) module. These modules are described in some details, here.

### 4.1.1. Question Decomposition and Constraint Recognition (QD&CR):

First, for a given complex question, Q, containing several relationships, entities, and constraints, the sequence of sub-questions based on the order of execution,  $Q: \{sq_1, sq_2, ..., sq_n\}$ , as well as the sequence of constraint sets,  $SQC = \{sqc_1, sqc_2, ..., sqc_n\}$ , are generated, where  $sqc_i = \{c_1, ..., c_t\}$  is the constraint set of the sub-question in the *i*-th iteration,  $sq_i$ .

As GraphMDR is a hybrid method, it utilizes an open domain method [13] for question decomposition. The analysis method is performed in such a way that, with little supervision, three types of inferences, namely intersection, bridging, and comparison, (as described in [18]) are identified.

Different constraints (temporal, ordinal, aggregation, etc) are also extracted using the method presented in [8] (a KB based method), in which the answer is extracted by filtering the constraints and applying them *after* the question graph construction step. However, in GraphMDR, this activity is performed in the middle stages of finding an answer. Also, the process should not remove the constraint words of the question, contrary to [8], as in many cases, the concepts related to the constraints are in the text.

Constraints may also be employed in the Answer Extraction module. However, investigating different conditions for applying constraints to the text or performing a query over the KB is for future consideration.

#### 4.1.2. Question Graph Expansion (QGE)

To find the answer in a massive pool of passages and knowledge base triples, the current work employs a heterogeneous graph representation model utilized in previous works, PullNet [12] and GRAFT-Nets [11]. We define the graph, G as G = $\{V, E\}$ ,  $(V = V_e \cup V_d \cup V_f)$ , where nodes (V) are of one of the types: entity nodes related to the question  $(V_e)$ , document text nodes  $(V_d)$  or fact nodes containing those entities  $(V_f)$ . The edges, E, show the connection between each  $V_e$  and  $V_d$  or  $V_f$  that contain the entity  $V_e$ .

The proposed model for graph expansion is built upon PullNet, with four main differences to increase the efficiency:

1. Instead of all the entities in the given question, only the entities of sub-question  $sq_1$  are entered in the graph  $G^1$  in the first iteration. If no entity is available, the  $sq_1$  is considered as  $V_d$ . For example, we draw the reader's attention to the example of  $q_3$  from Table 6 in the Appendix A.1 section. The generated sub-question  $sq_1$  does not contain an entity, so the whole query is added as a  $V_d$  node in the graph.

**2.** In each iteration, t, in addition to extending the graph through the entities available in the graph  $G^{t-1}$ , we add the entities of sub-question  $sq_t$  to this graph (refer to the item 1 in the graph expansion process in the *t*-th iteration).

When adding passages to the graph (refer to the item 3 in the graph expansion process in the *t-th* iteration), instead of a method based on sparse vectors, dense passage and the query vectors are compared. The query and passage encoder can be implemented using any neural network.

Here, we have used two independent neural networks. The passage encoder is called  $E_p$ . By entering the sub-question  $sq_t$ , its dense vector is generated,  $sq_t$ . The answers in the previous steps are  $(s_1, \ldots, s_{t-1})$ , including passages and triples. The query at each step involves concatenating the embedding vector of answers from the previous steps,  $(s_1, \ldots, s_{t-1})$ , and the current sub-question  $sq_t$  as defined in Eq. (1). We retrieve top-most vectors similar to the embedding vector q. Maximum inner product search is used to calculate similarity.

$$q = E_p(sq_t, s_1, \dots, s_{t-1}) = sq_t \oplus S_1 \oplus \dots \oplus S_{t-1}$$
(1)

In some questions, the answer of the current sub-question depends on the answer of the previous sub-questions. Thus, not capturing the answer of the previous steps reduces the accuracy of the final answer. Example  $q_2$  and  $q_3$  in Appendix A.1 are provided to clarify the issue.

The question  $q_2$  is divided into two subquestions  $sq_1$ : "Which film is based on an opera by Giacomo Puccini?" and  $sq_2$ : "In what city did the \"Prince of tenors\" star in a film?". The extracted passage as the answer to the first question is "It is based on the 1900 opera Tosca by Giacomo Puccini, which was adapted from the 1887 play by Victorien Sardou." From "Tosca (1956 film)" Wikipedia page.

The second sub-question alone does not elicit an answer "It was made at Cinecittà in Rome", because the answer does not depend on the named entity "Prince of tenors" of the current subquestion; while it is related to the answer of the previous step. As a result, the sub-question can only be answered correctly if combined with the answer of the previous steps. As another example in the second sub-question  $sq_2$ : "the university is located in what city?" from q<sub>3</sub>, the answer "New York City" can only be obtained if  $\mathbf{sq}_2$  is combined with the answer of the previous step "Ralph Franklin Hefferline (15 February 1910 in Muncie, Indiana – 16 March 1974) was a psychology professor at Columbia University" The answer can be extracted according to the passage sp1: "Columbia University (also known as Columbia, and officially as Columbia University in the City of New York) is a private Ivy League research university in New York City.". From "Columbia University" Wikipedia page and the triple sf1:" (dbr:Columbia University, dbp:city, dbr:New York City)" from the resource "dbr:Columbia University" in DBpedia.

**3.** To find the final answer, we rely on the multihop dense retrieval method instead of training through question-answer pairs. The following section (§ 4.1.3) explains the details of using this method.

**Graph expansion process:** The graph expansion process in iteration t is as follows:

- (i) Add existing new entities in sub-question  $sq_t$  to graph  $G^t$  as a new  $V_es$  (utilize the QD&CR module).
- (ii) Select k number of  $V_es$  with high probability in the graph  $G^{t-1}$  for graph expansion.
- (iii) Retrieve top-related documents and triples related to the k-selected  $V_es$  in the last step and add new nodes ( $V_d$  or  $V_f$ ) to the graph  $G^t$ .
- (iv)Add existing new entities in  $V_d$  s as new  $V_e$ s in the graph  $G^t$ .
- (v) Add existing new entities contained in  $V_f s$  as new  $V_e s$  in the graph  $G^t$ .
- (vi)Create graph edges between each  $V_e$  and the associated  $V_d$  or the  $V_f$  in the graph  $G^t$ .

**Disjunctive and Conjunctive Sub-questions:** If there is a disjunction or conjunction between two or more sub-questions, the entities of the sub-questions are added to the graph in one step. The search for similar passages and triples for each entity added to the graph is based on a sub-question containing that entity. The example of Question  $q_3$  is one of the multi-hop questions containing conjunction in Table 8 of Appendix A.2.

#### 4.1.3. Sequence Retrieval (SR)

There are two important objectives for extracting sequences from sources containing intermediate answers: i) to create explanations for the final answer and ii) to possibly provide the final answer from the text pieces (unlike previous works, in which extracting the answer as an entity was only possible). To achieve mentioned goals, the we follow the concepts of the method introduced by [5].

In [5], the problem of finding the answer to a multi-hop question has been transformed to the problem of finding the best sequence  $P_{seq}$ : { $p_1, ..., p_n$ } of intermediate answers among the k identified sequences { $P_{seq}^1, ..., P_{seq}^k$ }. These sequences only contain documents related to the question. The probability of finding a sequence is defined as:

$$P(P_{seq}|q) = \prod_{t=1}^{n} P(p_t|q, p_1, \dots, p_{t-1})$$
(2)

$$P(s_{t_i}|sq_t, s_1, \dots, s_{t-1}) = \frac{\exp((s_{t_i}, sq_t))}{\sum_{s \in (v_e \cup v_d \in G^t)} \exp((s, sq_t))}$$
  
where  $sq_t = g(sq_t, s_1, \dots, s_{t-1})$  and  $s_{t_i} = h(s_{t_i})$  (3)

Here, we employ the key idea that multi-hop questions are often sequences of sources to extract the final answer from. The SR module examines a set of sub-questions,  $Q: \{sq_1, sq_2, ..., sq_n\}$ , generated in the QD&CR module, as well as the graph  $G^t$  in the *t*-th iteration, generated with the QGE module. Since the entities in the fact nodes are present in the graph as entity nodes, the present study chooses the intermediate answers from  $v \in v_e \cup v_d$ .

The SR module needs to retrieve quence  $S_{seq}: \{s_1, ..., s_n\}$ , where  $s_t: \{s_{t_1}, ..., s_{t_r}\}$  represents a set of nodes  $(v \in v_e \cup v_d)$  in the graph  $G^t$  that are candidate source snippets for intermediate answers in the *t*-th iteration. In addition, the k best sequences,  $\{S_{seq}^1, ..., S_{seq}^k\}$ , from several candidates are extracted. The probability of selecting a node  $(v_e \cup v_d)$  is modeled as follows, which ultimately results in the probability of sequence  $S_{seq}$ .

$$P(S_{seq}|q) = \prod_{t=1}^{n} \prod_{i=1}^{r} P(s_{t_i}|sq_t, s_1, \dots, s_{t-1})$$
(4)

In the *t-th* iteration, the maximum inner product search is performed on the dense representation of all the nodes of the entity  $(V_e)$  and the document texts  $(V_d)$  in the graph  $G^t$ . The operator,  $\langle ., . \rangle$ , is defined as the inner product between vectors  $v \in v_e \cup v_d$  and  $sq_t$  vector at each iteration.

$$P(s_{t_i}|sq_t, s_1, \dots, s_{t-1}) = \frac{\exp\left(\langle s_{t_i}, sq_t \rangle\right)}{\sum_{s \in (v_e \cup v_d \in G^t)} \exp\left(\langle s, sq_t \rangle\right)}$$
(5)

where  $sq_t = g(sq_t, s_1, ..., s_{t-1})$  and  $s_{t_i} = h(s_{t_i})$ , and two functions h(.) and g(.) are encoders that produce the dense representation of intermediate answer source  $s_{ti}$  and sub-question  $sq_t$ , respectively.

In each iteration,  $sq_t$  and the previous set of intermediate answer sources are concatenated, and vector  $q_t$  is obtained by the g(.) encoder. The proposed approach for extracting sequences from source snippets containing intermediate answers is similar to the method mention in [5], as both utilize dense retrieval. However, there are three differences: (i) for t = 1, the present study relies on sub-question  $sq_1$  instead of the whole question; (ii) at the *t-th* iteration, instead of considering the representation of the question, subquestion  $sq_t$  participates in the question representation process; and (iii) in the proposed method's t-th iteration, several candidate sources are selected to obtain an intermediate response from the entity  $(V_{\rho})$ and the document texts  $(V_d)$  nodes in graph  $G^t$  with the help of entities in previous iterations.

The evaluation in Section 5.2 shows that significant improvements are gained by considering these modifications.

#### 4.1.4. Answer Extraction (AE)

This module uses the constraints identified in the first step to extract the answer. Simple heuristics for answer extraction are employed. We consider  $s_t$  in the selected  $S_{seq}$  sequence in the *t*-th iteration. For each  $s_{t_i}$ , if  $s_{t_i}$  is an entity node, the constraint similar to the approach by Shin et al. [8] is applied. Otherwise, the sub-question  $sq_t$  is searched for the expected answer type and, based on the answer type, the answer is extracted.

#### 4.2. Model Training

The proposed model includes two different types of training:

The first training is for identifying the best entities used for expansion at each iteration. The current study follows the PullNet [12] method, in which the learning is achieved by considering the questionanswer pairs for finding the shortest paths between the question entities and the answer entity.

The second training tries to identify the best sequences of entities and documents to answer the question. It follows the MDR [5] method, with the difference that we include several positive and negative KB facts in the training process, in addition to considering the question along with the related positive and negative sentences.

Table 2: Comparison of PullNet and GraphMDR systems based on Hits@1 metric

	MetaQA			WEBOUESTIONSSP	COMPLEXWEBOUESTIONS (dev)
	(1-hop)	(2-hop)	(3-hop)		
PullNet	92.4	90.4	85.2	51.9	33.7
GraphMDR	92.4	90.4	86.1	71.9	62.3

#### 5. Experiments and Results

This section describes multi-hop question-answer datasets and baselines and reports the results of a comparison between the proposed approach and the baseline systems.

#### 5.1. Databases and Baselines

There are two main categories: simple questions (single-hop) and complex questions (multi-hop). For simple questions, there are KB-based databases, such as Wikimovie [39], and text-based databases, such as TriviaQA [40] and Squad [41]. For complex KB-based questions (multi-hop), WEBQUESTIONSSP [42], WebQuestions [36], COMPLEXWEBQUESTIONS [18], and MetaQA [6] datasets, and for text-based questions, HOTPOTQA [43] dataset have been used in the literature.

The proposed method aims to address some challenges of the PullNet methods [12] and the multi-hop dense retrieval MDR method [5]. The PullNet approach answers complex questions using a KB and text, while MDR answers complex questions using a pool of text documents. Therefore, we evaluate GraphMDR with PullNet [12] using MetaQA [6], WEBQUESTIONSSP [42], and COMPLEXWEBQUESTIONS [18] databases and with MDR method [5] using HOTPOTQA [43] database.

Evaluating with more databases and creating databases containing questions with answers in text snippets or entities are left for future work.

**MetaQA** [6]: This database is based on WikiMovies [39], and uses some Wikipedia texts to help answer questions. Also, up to 3-hop questions added to the previous single-hop collection, known as the Vanilla version. The KB includes 43k entities and 13k triplets.

The questions have been transformed into 2-hop and 3-hop questions, based on some patterns, in a sequence of 2 and 3 simple questions, respectively. Therefore, the basis for constructing these questions is the combination of sub-questions. Since the present approach is based on the decomposition of questions into simple questions, this dataset is entirely consistent with the proposed method. Table 7 in Appendix A.1.2 shows examples of 2-hop and 3-hop questions, along with the type of each question and its sub-questions.

**WEBQUESTIONSSP** [42]: This database is based on the WebQuestions database [36], where 84% of the questions are simple, and the rest are up to 2-hops questions that can be answered using Freebase alone. On the other hand, Wikipedia texts have also been used as a text corpus to provide a composite platform that requires text and a KB to respond.

**COMPLEXWEBQUESTIONS** [18]: This database has created more complex questions by adding constraints and expanding the WEBQUESTIONSSP database [42] entities.

**HOTPOTQA** [43]: This database has up to 2-hop questions, is based on Wikipedia, and provides the possibility of answer interpretation by having a set of supporting passages to reach the answer.

Because GraphMDR has the ability to search using both textual and KB sources, in order to make it possible to evaluate, the k facts of KB most similar to the supporting passages in the HOTPOTQA's database are selected using the embedding vector comparison method, and these are used alongside the supporting passages.

In addition, since GraphMDR is based on decomposing questions into sub-questions, the data set questions are decomposed, first, and the existing supporting passages for each question are separated based on the sub-questions. Examples of questions in the dataset that are enriched with supporting triples and presented based on sub-questions are shown in Table 6 of Appendix A.1.1 section.

The questions in the HOTPOTQA dataset have different numbers of supporting passages. For example, question  $\mathbf{q}_1$  from Table 6 in Appendix A.1 section, has three supporting passages. Since the proposed method is based on question decomposition, the three supporting passages are separated based on the generated sub-questions, and for each supported passage, supported triples are generated. As can be seen, the first sub-question  $\mathbf{sq}_1$  contains one supporting passage, and the second sub-question  $\mathbf{sq}_2$  contains two supporting passages.

Table 1 provides data statistics on the train, dev, and test categories of the databases used for the current paper's evaluation.

Table 1: Statistics of databases used in the evaluations

Benchmark	Train	dev	test
MetaQA	329282	39138	39093
WEBQUESTIONSSP	2848	250	1639
COMPLEXWEBQUESTIONS	27623	3518	3531
HOTPOTQA	90564	7405	7405

#### 5.2. Evaluations

As mentioned in the previous section, the two basic systems in the current work are PullNet [12] and MDR [5], of which MDR can provide responses only from text snippets, and PullNet is able to extract responses only as an entity from a KB. We need the database to contain multi-hop questions, some of whose answers are in text snippets and some in KB entities for evaluations. However, since the database is not available at this time, the proposed system, GraphMDR, is compared separately with two base systems, and experiments on such datasets are left to future work.

All experiments are conducted on a machine with a 4 Core Intel Xeon E5 CPU @ 2.00GHz with 16 GB of RAM. The FAISS<sup>1</sup> library is used to store dense vectors and calculate the best candidates. This opensource library is very effective for searching for similarities between dense vectors, because searching through vector clustering allows the system to search through the billions of dense vectors quickly. Answer selection is made using the Tersformer<sup>2</sup> framework. To identify the answer, the ELECTRA model is used, which is reported in [5] to have the best performance.

In the evaluations, whenever the difference in results is small, a *t-test* is also performed by SPSS v11.0. Results with a statistically significant difference are highlighted in the tables.

Table 2 shows the accuracy of GraphMDR and PullNet for the three provided KB-based databases. Similar to PullNet, GraphMDR provides a combination of text and a KB to answer questions in all three databases. Fifty percent of the KB is utilized and, next to the KB, is a text corpus that uses entity link tools to connect to the KB. Table 2's evaluation metric is Hits@1, which indicates the accuracy of the answer with the highest prediction.

As shown in Table 2, GraphMDR is more accurate than PullNet. The higher accuracy can be interpreted in two main factors. The first one is that in the proposed method, passage retrieval is done based on embedding vectors, while in PullNet, traditional methods (TF-IDF) are used. The second factor is the involvement of the previous steps results in retrieving the pieces of information containing the answer in each step.

The accuracy of both systems in the MetaQA dataset is much higher than the two other datasets, WEB-QUESTIONSSP and COMPLEXWEBQUESTIONS, because in MetaQA dataset, the questions are based on specific and limited patterns.

Table 3 provides comparison of GraphMDR and PullNet results using Google snippets without entity links to a KB, Wikipedia text data with existing KB links, and KB data alone for multi-hop COMPLEX-WEBQUESTIONS database.

Unlike PullNet, the GraphMDR method has the ability to extract responses from text snippets. As a result, as shown in Table 3, GraphMDR provides more accuracy for Google snippets and Wikipedia than PullNet. This result is especially in the case of Wikipedia, because the entities in the text are linked to a KB. In addition, the GraphMDR works based on embedding vector space, so it is more accurate on Hits@1 metric. In the case of Freebase, the accuracies of both methods are the same, because only KB entities are used to extract the response.

To analyze the efficiency-accuracy trade-off, the number of entities and the recalls of the two systems, PullNet and GraphMDR, using a full KB are compared. As shown in Table 4 despite the production of a smaller graph by GraphMDR, it has also reported

	Google snippets	Wikipedia	Freebase
PullNet	29.7	13.8	45.9
GraphMDR (present study)	45.2	52.2	45.9

Table 4 : Comparison of PullNet and GraphMDR based on the amount of entities/ recall of multi-hop

	MetaQA (3-hop)	COMPLEXWEBQUESTIONS
PullNet	63.3/0.98	44.1/0.68
GraphMDR (using sq)	34.3/0.98	<b>19.5</b> /0.78
GraphMDR (using q)	62.3/0.98	40.5/ <b>0.81</b>

<sup>1</sup> GitHub.com/facebookresearch/faiss

<sup>2</sup> https://huggingface.co/transformers/

Table 5: Evaluation based on retrieval performance in recall at k retrieved passages or related facts

	HotpotQA		
	R@2	R@10	R@20
MDR	65.2	77.5	80.2
GraphMDR (present study)	71.3	80.4	88.6

more accuracy.

In GraphMDR, the number of iterative steps to expand the graph and retrieve the sequence of answers is equal to the number of sub-questions (GraphMDR (steps=sq#) in Table 4). It should be noted that, as stated in the explanation of the question graph expansion method, examined in Section 4, sub-questions that are connected by disjunction or conjunction, are considered as a step.

We were interested in determining the effect of using query decomposition on the proposed method. Therefore, the results of the proposed method are reported in two experiments with (GraphMDR (using sq)) and without (GraphMDR (using q)) the use of sub-questions along with the results of PullNet (Table 4).

As can be seen, when the question is divided into sub-questions, the number of entities within the graph is significantly reduced compared to the other two cases. This observation is logical, because the graph expansion is done with the guidance of subquestions (similar to the way the human mind works). In the other two cases, the whole question is used from the beginning of the process to extract the answer, which lead to the production of a larger graph.

In the second case (GraphMDR (using q)), due to the addition of question at once, as in the PullNet method, the number of graph entities is much higher than in the previous case (GraphMDR (using sq)), while the accuracy does not increase significantly. The increase in accuracy compared to the PullNet method is due to the use of 1) dense retrieval approach, and 2) the results of the previous steps, discussed in the description of Table 3.

Table 5 compares the retrieval efficiency of related passages or facts in both GraphMDR and MDR systems. As seen, GraphMDR provides better results than MDR. This improvement in results can be explained as follows. GraphMDR, retrieves nodes in the graph that are related to the entities existing in the sub-question being processed. The current subquestion is used for comparison, so the results are closer to the supporting facts and passages in the database. Supporting facts and passages are created based on the sequential nature of answering multihop questions.

In Figure 2, the result from running two systems, MDR and GraphMDR, based on different k inputs on the HOTPOTQA database is shown (It should be noted that the MDR system was reimplemented).

Since HOTPOTQA questions are up to 2-hops, GraphMDR contains a maximum of two subquestions that include two iterations to find the answer. Consequently, there is a maximum sequence of two sets. In each set, the most appropriate doc or entity nodes are retrieved. *k* is considered as the total number of nodes suitable for recovery in sequence in a maximum of two iterations ( $\frac{k}{2}$  in each iteration).

As shown in Figure 2, for each input, k, GraphMDR offers higher accuracy in shorter execution time. These results can be interpreted in two ways: First, according to Section 4.3, MDR uses all the question words to find the sequence of texts containing the answer. In contrast, GraphMDR uses subquestion order sequences to find text snippets or entities similar to the current sub-question. This fact leads to more accurate results as well as faster speed.

Secondly, MDR only searches for answers using a sequence of text snippets. In contrast, the answer is extracted more accurately and quickly in a structured data if an answer is available in a KB. GraphMDR achieves better results by using this feature.

#### 5.3. Error Analysis

Although GraphMDR outperforms other systems, it is not error-free. Examining the errors usually can improve the system performance and show the path for future researches. For GraphMDR, three general



Figure 2: Comparison of MDR and GraphMDR based on the efficiency-performance trade-off.

categories of errors were identified as:

- (i) Syntactic errors: The first category of errors that are visible in the process of extracting the answer is the result of parsing the question. Since the decomposition method used in this research is based on dependency relations, any errors in the parser leads to the wrong result in generating the sub-questions. For example, in question  $\mathbf{q}_2$  in Table 8 of the appendix A.2 section, the question is divided into two sub-questions. The dependence of the word "head-quartered" is incorrectly recognized to the second sub-question  $\mathbf{sq}_2$ , while it should have been to the first sub-question  $\mathbf{sq}_1$ .
- (ii) Dataset preparation errors: The second category of errors is due to the method used to generate the supporting triples, which reduces the efficiency of the proposed system. To compare the proposed method with the base system, the HOTPOTQA data set is used that has supporting passages (called sp). In the current research to evaluate the proposed method, for each supporting passage, supporting triples (called sf) is also generated. The production of triples is done by comparing the embedding vector of the supporting passage with the triple vectors of the knowledge base and finding similar vectors.

However, in several cases, all or some the supporting triples produced for a passage, are not related to the answer. Since for the evaluation, the comparison between the supporting triples and the responses extracted by the proposed method is performed, the inaccuracy of the supporting triples in the data set shows a decrease in the accuracy of the proposed method in retrieving information pieces containing responses.

For example, for the  $q_3$  question, Table 8 of appendix A.2, for the sub-question  $sq_1$ , the supporting passage  $sp_1$  is generated in the HOT-POTQA dataset. We have generated the supporting triples  $sf_1$ ,  $sf_2$ , and  $sf_3$ . The correct answer to the sub-question  $sq_1$  is the first supporting triple,  $sf_1$ , and the second and third triples are incorrect. These two incorrect triples are generated only because of the similarity to the supporting passage  $sp_1$ .

Although, GraphMDR works correctly and extracts the first supporting triple in the response. not extracting the other two, shows a decrease in the accuracy. To correct these errors, the method of generating the supporting triplets should be modified. In future work, the method of generating supporting passages in the HOTPOTQA dataset should be used to generate supporting triples. In the HotpotQA dataset, some people are employed to distinguish the supporting passages from the passages pool to achieve the correct answer.

(iii) Semantic errors: The third category of errors rises from the complexity of the question concept. Answering the complex questions requires additional background knowledge. For example, in question  $q_1$  in Table 8 of the appendix A.2 section, "What distinction is held by ...?", the answer needs to find a distinction for the entity in the question among other people. The concept of "shortest person" must be extracted as the answer. In the supporting passage, a number "5 ft" is mentioned as the height and the title, "The shortest player ever to play in the National Basketball Association" for the entity. Recognizing "The shortest height" as a distinguishing feature of an entity requires additional background knowledge.

#### 6. Discussion

This section highlights the findings, and theoretical and practical implications of current research.

#### 6.1. Findings

As stated in the proposed method, to extract the answers to multi-hop questions, an iterative method, searching for the number of information snippets, including intermediate answers in each step (hop), is required. The findings of this study are as follows:

- (i) Examining the HOTPOTQA and MetaQA data sets shows that the number of steps to extract the answer is equal to the number of subquestions (disjunctive and conjunctive subquestions are also searched in one step). GraphMDR uses this feature for extracting the answer, so the accuracy of the proposed method is better than the base methods (Note the results reported in Table 2 and Table 4.). This is in contrast with the PullNet system, where there are no limits for the number of steps to find the answer.
- (ii) The second finding is that if the unstructured data contains links to the entities within the knowledge base, use of these links, along with the dense retrieval, makes it possible to extract related passages quicker. As a result, it allows the system to search for answers in a large pool of documents.

- (iii) The third finding is that using a probabilistic method, enables us to extract explanations for the answer in a hybrid system. In the real world, extracting the explanations of how to reach to the answer is essential, especially in multi-hop questions. The extracted explanations contain sequences of information (including texts and triples) containing intermediate answers.
- (iv) Using the knowledge base in a text-based state-of-the-art QA system can increase the speed and accuracy.
- (v) The last finding in this study reflects the increase in system performance in the multi-shot view in comparison to the single-shot view, considering multi-hop questions. In this study, considering the question as a set of sub-questions (looking at the question in a few shots) led to an increase in the system performance, while the keeping accuracy to be competitive. In the base systems, not using sub-questions (one-shot look at the question) has created to an increase in search space and thus reduced the speed of extracting the answer.

#### 6.2. Theoretical implications

The current research falls into the large field of information extraction and sub-branch of QA systems. Studies in this field have reached maturity in two main categories: finding answers from text and finding answers from knowledge base. Hybrid QA systems take advantage of both KB and text sources to extract answers. However, few hybrid systems could simultaneously use these two sources to extract answers (early fusion model).

Also, recent researches in hybrid QA, has tended to process multi-hop questions. To address the problem of finding answers to such questions, graph expansion methods have been used frequently. However, these systems extract passages from a small pool of documents to retrieve textual information using traditional methods. Using these methods, limits the speed of information extraction. In addition, the answer can only be extracted in the form of the entity of KB triples, while in many questions, the requested information is not available in the knowledge base and must be extracted from the text.

To solve passage retrieval, the theoretical concepts in the field of information retrieval prompted us to use the dense retrieval method. These methods have recently proven their success over traditional methods (TF-IDF). Another challenge in the state-of-the-art hybrid QA systems is the limitation in extracting explanation for the answer and finding the best sequence of answers. Considering for the problem of answering multi-hop questions in hybrid QA approaches, with the idea of the sequential nature of such questions, suggested the use of a probabilistic method in solving the challenge with a well-studied text-based QA system.

This method is based on finding the best sequence of intermediate answers in the text body to reach the final question, which in the present study has been customized to create the ability to extract the best sequence of nodes of passages and entities in the graph for selecting the best answers.

In the base systems, the multi-hop question is considered as an information unit for extracting information. However, due to the sequential nature of multihop questions, it seems logical that the information in the question should be used in a multi-step process in the answer extraction process.

In the proposed method, the idea of using subquestions is induced from what happens in a human mind when faced with a multi-hop question. Therefore, in the proposed method, the question is first decomposed into sub-questions, then we search for the final answer by using the graph of passage, entity, and triple nodes, and using the probabilistic method of finding the best sequence of answers.

#### 6.3. Practical implications

In general, in an information extraction system that aims to meet the information needs of users, several main important factors are involved: response time, the accuracy of response, scalability, and the explanation of how the response is extracted. Therefore, QA systems try to provide the existing challenges in improving these four essentials. According to the findings of this study, the efficiency of the proposed method can be expressed in the following cases according to the main factors mentioned.

- Using a knowledge base along with a text source increases the speed and accuracy in the proposed system.
- The use of sub-questions reduces the search space, which increases the speed of extracting the answer.
- Since the number of search steps according to the findings is equal to the number of subquestions, there is no need for further search, this leads to an increase in the speed. It should

be noted that disjunctive and conjunctive subquestions can be searched in one step.

- The use of dense retrieval methods, unlike conventional methods, makes it possible to extract answers from large volumes of documents. As a result, the proposed method could extract the answer in a large amount of information with acceptable speed and accuracy.
- As mentioned, determining how to reach the answer is a real-world necessity in information extraction systems. The proposed system could provide the required explanation to the user by extracting the best sequence of information pieces containing intermediate answers.

#### 7. Conclusion and Future Work

Open-domain QA systems have a long history in both text-based and knowledge-based domains. Considering the pitfalls and merits of separately extracting information from these two sources, have led the QA community to lean toward using both at the same time.

The proposed systems are still in their infancy, especially for multi-hop questions. goals are search space reduction, scalability, explainability of the answer, and the possibility of using and answer extraction answers from both sources. The aim of the current research is to solve the mentioned challenges related to the two base systems of PullNet, a hybrid QA, and MDR, a text-based QA. The results show that, by comparing the proposed method with these systems and providing the possibility of extracting responses from both textual sources and a KB, the response extraction speed rises by reducing the search space, while accuracy either remains competitive or is enhanced.

Future works includes changing the method of model training so that it is independent of the KB. In addition, considering the priority of extracting documents or similar entities, weighing them, and calculating their impacts on the accuracy can also be examined in future studies. Another plan is to examine the constraint types, especially those needing calculation, and prioritize their execution over text or a KB.

#### References

R. Das, S. Dhuliawala, M. Zaheer, L. Vilnis, I. Durugkar,
A. Krishnamurthy, A. Smola, A. McCallum, Go for a

walk and arrive at the answer: Reasoning over paths in knowledge bases using reinforcement learning, arXiv preprint arXiv:1711.05851. (2017).

- [2] W. Zheng, J.X. Yu, L. Zou, H. Cheng, Question Answering Over Knowledge Graphs : Question Understanding Via Template Decomposition, Proceedings of the VLDB Endowment 11.11 (2018) 1373-1386. doi:10.14778/3236187.3236192.
- [3] Y. Hao, Y. Zhang, K. Liu, S. He, Z. Liu, H. Wu, J. Zhao, An end-to-end model for question answering over knowledge base with cross-attention combining global knowledge, ACL 2017 - 55th Annu. Meet. Assoc. Comput. Linguist. Proc. Conf. (2017) 221–231. doi:10.18653/v1/P17-1021.
- [4] V. Karpukhin, O. Barlas, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, W. Yih, Dense Passage Retrieval for Open-Domain Question Answering, Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP) (2020) 6769-6781. doi: 10.18653/v1/2020.emnlp-main.550.
- [5] W. Xiong, X.L. Li, S. Iyer, J. Du, P. Lewis, W.Y. Wang, Y. Mehdad, W.T. Yih, S. Riedel, D. Kiela, B. Oğuz, Answering complex open-domain questions with multihop dense retrieval, International Conference on Learning Representations. (2021).
- [6] Y. Zhang, H. Dai, Z. Kozareva, A.J. Smola, L. Song, Variational Reasoning for Question Answering with Knowledge Graph, Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 32, No. 1) (2018).
- [7] A. Abujabal, M. Riedewald, M. Yahya, G. Weikum, Automated template generation for question answering over knowledge graphs, 26th Int. World Wide Web Conf. WWW 2017 (2017) 1191–1200. doi:10.1145/3038912.3052583.
- [8] S. Shin, K.H. Lee, Processing knowledge graph-based complex questions through question decomposition and recomposition, Information Sciences (2020) 234–244. doi:10.1016/j.ins.2020.02.065.
- [9] K. Xu, L. Wu,Z. Wang, M. Yu, L. Chen, V. Sheinin, Exploiting Rich Syntactic Information for Semantic Parsing with Graph-to-Sequence Model. In Proceedings of the

2018 Conference on Empirical Methods in Natural Language Processing (2018) 918-924. doi: 10.18653/v1/D18-1110.

- [10] N. Bhutani, X. Zheng, H. V. Jagadish, Learning to answer complex questions over knowledge bases with query composition, Int. Conf. Inf. Knowl. Manag. Proc. (2019) 739–748. doi:10.1145/3357384.3358033.
- [11] H. Sun, B. Dhingra, M. Zaheer, K. Mazaitis, R. Salakhutdinov, W.W. Cohen, Open domain question answering using early fusion of knowledge bases and text, Proc. 2018 Conf. Empir. Methods Nat. Lang. Process. EMNLP 2018. (2020) 4231–4242. doi:10.18653/v1/d18-1455.
- [12] H. Sun, T. Bedrax-Weiss, W.W. Cohen, PullNet: Open domain question answering with iterative retrieval on knowledge bases and text, Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing. (2019) 2380–2390.
- [13] S. Min, V. Zhong, L. Zettlemoyer, H. Hajishirzi, Multihop reading comprehension through question decomposition and rescoring, Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (2019) 6097–6109.doi:10.18653/v1/p19-1613.
- [14] T. Noraset, L. Lowphansirikul, S. Tuarob, WabiQA: A Wikipedia-Based Thai Question-Answering System, Information Processing & Management, 58(1), 102431 (2021). doi:10.1016/j.ipm.2020.102431.
- [15] W. Yang, Y. Xie, A. Lin, X. Li, L. Tan, K. Xiong, J. Lin, End-to-End Open-Domain Question Answering with BERTserini. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations) (2019) 72-77. doi:10.18653/v1/N19-4013.
- [16] Y. Luan, J. Eisenstein, K. Toutanova, M. Collins, Sparse, Dense, and Attentional Representations for Text Retrieval, transactions of the Association for Computational Linguistics 9 (2021) 329-345. doi:10.1162/tacl\_a\_00369.
- [17] L. Song, Z. Wang, M. Yu, Y. Zhang, R. Florian, D. Gildea, Exploring graph-structured passage representation for multi-hop reading comprehension with graph neural networks, arXiv:1809.02040 (2018).

- [18] A. Talmor, J. Berant, The web as a knowledge-base for answering complex questions, Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (2018) 641–651. doi:10.18653/v1/n18-1059.
- [19] M. Raison, P.E. Mazaré, R. Das, A. Bordes, Weaver: Deep Co-encoding of questions and documents for machine reading, CoRR abs/1804.10490 (2018).
- [20] Y. Wu, S. Zhao, R. Guo, A novel community answer matching approach based on phrase fusion heterogeneous information network, Information Processing & Management, 58(1) (2021). doi:10.1016/j.ipm.2020.102408.
- [21] M. Yu, W. Yin, K. S. Hasan, C. dos Santos, B. Xiang, B. Zhou, Improved Neural Relation Detection for Knowledge Base Question Answering. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) (2017) 571-581. doi:10.18653/v1/P17-1053.
- [22] K. Luo, F. Lin, X. Luo, K. Zhu, Knowledge base question answering via encoding of complex query graphs. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (2018) 2185– 2194. doi:10.18653/v1/d18-1242.
- [23] L. Dong, M. Lapata. "Language to Logical Form with Neural Attention." Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (2016) 33-43. doi: 10.18653/v1/P16-1004.
- [24] C. Liang, J. Berant, Q. Le, K. Forbus, N. Lao, Neural Symbolic Machines: Learning Semantic Parsers on Freebase with Weak Supervision. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) (2017) 23– 33. doi:10.18653/v1/P17-1003.
- [25] L. X. Yao, B. Van Durme, Information Extraction over Structured Data: Question Answering with Freebase Center for Language and Speech Processing, Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistic (2014) 956–966. doi: 10.3115/v1/P14-1090.

- [26] B. Fu, Y. Qiu, C. Tang, Y. Li, H. Yu, J. Sun, A survey on complex question answering over knowledge base: Recent advances and challenges, arxiv-2007.13069. (2020) 1–19.
- [27] Y. Chen, L. Wu, M.J. Zaki, Bidirectional Attentive Memory Networks for Question Answering over Knowledge Bases. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers) (2019) 2913–2923. doi:10.18653/v1/N19-1299.
- [28] Y. Qiu, Y. Wang, X. Jin, K. Zhang, Stepwise reasoning for multi-relation question answering over knowledge graph with weak supervision. In Proceedings of the 13th International Conference on Web Search and Data Mining (2020) 474–482. doi:10.1145/3336191.3371812.
- [29] N. Bhutani, X. Zheng, K. Qian, Y. Li, H. Jagadish, Answering Complex Questions by Combining Information from Curated and Extracted Knowledge Bases, Proceedings of the First Workshop on Natural Language Interfaces (2020) 1–10. doi:10.18653/v1/2020.nli-1.1.
- [30] D. Gillick, S. Kulkarni, L. Lansing, A. Presta, J. Baldridge, E. Ie, D. Garcia-Olano, Learning dense representations for entity retrieval, Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL) (2019) 528–537. doi:10.18653/v1/k19-1049.
- [31] G. Izacard, E. Grave, Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering, Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: . (2020) 874-880.
- [32] Z.A. Taeb, S. Momtazi, Text-based question answering from information retrieval and deep neural network perspectives: A survey, arXiv:2002.06612 (2020) 1–56.
- [33] H.K. Azad, A. Deepak, Query expansion techniques for information retrieval: A survey, Information Processing & Management 56.5 (2019): 1698-1735. doi:10.1016/j.ipm.2019.05.009.
- [34] X. Lu, R.S. Roy, Y. Wang, G. Weikum, A. Alexa, Answering Complex Questions by Joining Multi-Document Evi-

dence with Quasi Knowledge Graphs, Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval (2019) 105-114. doi:10.1145/3331184.3331252.

- [35] D. Savenkov, E. Agichtein, When a knowledge base is not enough: Question answering over knowledge bases with external text data, Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval (2016) 235–244. doi:10.1145/2911451.2911536.
- [36] K. Xu, S. Reddy, Y. Feng, S. Huang, D. Zhao, Question answering on freebase via relation extraction and textual evidence, Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) (2016) 2326–2336. doi:10.18653/v1/p16-1220.
- [37] B. Oguz, X. Chen, V. Karpukhin, S. Peshterliev, D. Okhonko, M. Schlichtkrull, S. Gupta, Y. Mehdad, S. Yih, Unified Open-Domain Question Answering with Structured and Unstructured Knowledge, arXiv:2012.14610 (2020).
- [38] R. Das, M. Zaheer, S. Reddy, A. McCallum, Question answering on knowledge bases and text using universal schema and memory networks, In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics. (2017) 358–365. doi:10.18653/v1/P17-2057.
- [39] A.H. Miller, A. Fisch, J. Dodge, A.H. Karimi, A. Bordes, J. Weston, Key-value memory networks for directly reading documents, In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing (pp. 1400-1409) (2016). 1400–1409. doi:10.18653/v1/d16-1147.
- [40] M. Joshi, E. Choi, D.S. Weld, L. Zettlemoyer, TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension, In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (2017) 1601–1611. doi:10.18653/v1/P17-1147.
- [41] P. Rajpurkar, J. Zhang, K. Lopyrev, P. Liang, SQuad: 100,000+ questions for machine comprehension of text, In Proceedings of the 2016 Conference on Empirical

Methods in Natural Language Processing (2016) 2383–2392. doi:10.18653/v1/d16-1264.

- [42] W.T. Yih, M. Richardson, C. Meek, M.W. Chang, J. Suh, The value of semantic parse labeling for knowledge base question answering. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers) (2016) 201–206. doi:10.18653/v1/p16-2033.
- [43] Z. Yang, P. Qi, S. Zhang, Y. Bengio, W.W. Cohen, R. Salakhutdinov, C.D. Manning, Hotpotqa: A dataset for diverse, explainable multi-hop question answering, Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (2018) 2369–2380. doi:10.18653/v1/d18-1259.

## Appendix A: Examples for Qualitative Analysis

In this section, examples of data sets used in this research are given, which include two subsections. The examples in the first section are used to understand the proposed model throughout the article, and the second section is devoted to the cases in which the proposed model faces challenges and are referred to in the error analysis section.

#### A.1: Running Examples throughout the program

This section provides examples of multi-hop questions from two datasets HOTPOTQA and MetaQA that are referenced throughout the article to understand the proposed method.

Table 6: Examples of HOTPOTQA datasets. sqi,spj, and sfk stand for i-th sub-question, j-th supporting passage and k-th supporting fact for question qt, respectively.

**q**<sub>1</sub>: What government position was held by the woman who portrayed Corliss Archer in the film Kiss and Tell? **answer**: Chief of Protocol.

**sq**<sub>1</sub>: Which <u>woman</u> who portrayed Corliss Archer in the film Kiss and Tell?

sp<sub>1</sub>: (Kiss and Tell (1945 film)) Kiss and Tell is a 1945 American comedy film starring then 17-year-old Shirley Temple as Corliss Archer.

sf<sub>1</sub>: not existed.

sq<sub>2</sub>: What government position was held by the <u>woman</u>?

**sp**<sub>1</sub>: (Shirley Temple) Shirley Temple Black (April 23, 1928 – February 10, 2014) was an American actress, singer, dancer, businesswoman, and diplomat who was Hollywood's number one box-office draw as a child actress from 1934 to 1938.

sf1: (dbr: Shirley Temple, dbo:wikiPageWikiLink, dbc:American film actresses).

**sp**<sub>2</sub>: (Shirley Temple) Shirley Temple As an adult, she was named United States ambassador to Ghana and to Czechoslovakia, and also served as Chief of Protocol of the United States.

sf<sub>2</sub>: (dbr: Shirley Temple, dct:subject, dbc:Chiefs of Protocol of the United States).

 $\mathbf{q}_2$ : In what city did the \"Prince of tenors\" star in a film based on an opera by Giacomo Puccini?

answer: Rome.

sq1: Which film is based on an opera by Giacomo Puccini?

sp<sub>1</sub>: (Tosca (1956 film)) It is based on the 1900 opera Tosca by Giacomo Puccini, which was adapted from the 1887 play by Victorien Sardou.

sf<sub>1</sub>: not existed.

sq<sub>2</sub>: In what city did the \"Prince of tenors\" star in a <u>film</u>?

sp1: (Tosca (1956 film)) It was made at Cinecittà in Rome.

 $\mathbf{sf_1}$ : not existed

 $\mathbf{q}_3$ : Ralph Hefferline was a psychology professor at a university that is located in what city?

answer: New York City

sq1: Ralph Hefferline was a psychology professor at what university?

**sp**<sub>1</sub>: (Ralph Hefferline) Ralph Franklin Hefferline (15 February 1910 in Muncie, Indiana – 16 March 1974) was a psychology professor at Columbia University.

sf<sub>1</sub>: not existed.

sq<sub>2</sub>: the <u>university</u> is located in what city?

sp1: (Columbia University) Columbia University (also known as Columbia, and officially as Columbia University in the

City of New York) is a private Ivy League research university in New York City.

 $sf_1$ : (dbr: Columbia University, dbp:city , dbr:New York City).

#### A.1.1: HOTPOTQA dataset

Table 6 contains examples of the HOTPOTQA data set questions that are referenced in the article for clarification.

#### A.1.2 MetaQA dataset

In this section, two examples of 2-step and 3-step questions from the data set are given. The questions were created using two and three patterns that are specified in the table as type. The sub-questions for each question are generated by the proposed approach in this paper and, as can be seen, correspond to the patterns used to construct the questions.

#### A.2. Error Cases in present approach, GraphMRD

The examples in Table 8 are selected from the HOTPOTQA dataset to interpret the challenges faced by the proposed method. The  $q_1$ ,  $q_2$ , and  $q_3$  questions are presented as examples of errors in concept complexity, question decomposition, and the method of generating supporting triples, respectively.

Table 7: Examples of 2 and 3-hop questions from the MetaQA dataset along with the fsub-questions generated by the proposed method.

	<b>q</b> <sub>1</sub> : who are the directors of the films written by [Laura Kerr]?				
	answer: H.C. Potter.				
2-hop	qtype: writer_to_movie_to_director.				
	sq <sub>1</sub> : which <u>films</u> written by [Laura Kerr]?				
	sq <sub>2</sub> : who are the directors of the <u>films</u> ?				
	q <sub>2</sub> : what types are the films directed by the director of [For Love or Money]?				
	answer: Action Comedy Western Thriller Crime.				
2 h	qtype: Movie-to-director-to-movie-to-genre.				
5-110p	sq <sub>1</sub> : director of [For Love or Money]?				
	sq <sub>2</sub> : the films directed by the director?				
	sq <sub>3</sub> : what types are the films?				

$\mathbf{q}_{1}$ : What distinction is held by the former NBA player who was a member of the Charlotte Hornets during their 1992-93 season
and was head coach for the WNBA team Charlotte Sting?
answer: shortest player ever to play in the National Basketball Association
sq1: Which NBA player who was a member of the Charlotte Hornets during their 1992-93 season and was head coach for the
WNBA team Charlotte Sting?
sp1: (1992-93_Charlotte_Hornets_season) With the addition of Mourning, along with second-year star Larry Johnson and
Muggsy Bogues, the Hornets struggled around .500 for most of the season, but won 9 of their final 12 games finishing their
season third in the Central Division with a 44-38 record, and qualified for their first ever playoff appearance.
sf1: (dbr: 1992 93 Charlotte Hornets season, dbo:wikiPageWikiLink, dbr:Muggsy Bogues).
sp2: (Muggsy Bogues) After his NBA career, he served as head coach of the now-defunct WNBA team Charlotte Sting.
sf <sub>2</sub> : (dbr:Muggsy Bogues, is dbp:coach of, dbr:2006 Charlotte Sting season).
sq <sub>2</sub> : What distinction is held by the former <u>NBA player</u> ?
sp1: (Muggsy Bogues) The shortest player ever to play in the National Basketball Association, the 5 ft Bogues played point
guard for four teams during his 14-season career in the NBA.
$\mathbf{sf}_1$ : not existed.
<b>q</b> <sub>2</sub> : Where is the company that Sachin Warrier worked for as a software engineer headquartered?
answer: Ronald Shusett.
sq <sub>1</sub> : Which company that Sachin Warrier worked for as a software engineer headquartered?
sp <sub>1</sub> : (Sachin Warrier) He was working as a software engineer in Tata Consultancy Services in Kochi.
<b>sf</b> : (dbr: Tata Consultancy Services, is dbo:wikiPageWikiLink of, dbr:Kochi).
say: Where is the company?
sp.: (Tata Consultancy Services) Tata Consultancy Services (TCS) is an Indian multinational information technology (IT)
services and consulting company, headquartered in Mumbai. Maharashtra, India and largest campus and workforce in
Chennai Tamil Nadu. India
sf.: (dbr:Tata Consultancy Service, dpp:located, dbr:Mumbai).
a : What is the name of the fight song of the university whose main commusic in Lawrence. Kanses and whose branch commuses
<b>q</b> <sub>3</sub> . What is the name of the right song of the university whose main campus is in Lawrence, Kansas and whose oranen campuses are in the Kansas City metropolitan area?
are in the Kansas City met opontal area:
answer. Kansas Song.
$\mathbf{sq}_1$ : which <u>university</u> whose main campus is in Lawrence, Kansas?
sp <sub>1</sub> : (University of Kansas) The main campus in Lawrence, one of the largest conege towns in Kansas, is on Mount Oread,
af (dhe University of Kanaga dharaity dhe Lawrence Kanaga)
$st_1$ : (dbr: University of Kansas, dbo:city, dbr:Lawrence Kansas).
st <sub>2</sub> : (dbr: Mount Oread, dbo:locatedInArea, dbr:Lawrence Kansas).
st <sub>3</sub> : (Mount Oread, dbp:elevationFt, 1037).
$\mathbf{sq}_2$ : Which <u>university</u> whose branch campuses are in the Kansas City metropolitan area?
<b>sp</b> <sub>1</sub> : (University of Kansas) Two branch campuses are in the Kansas City metropolitan area: the Edwards Campus in Over-
land Park, and the university's medical school and hospital in Kansas City.
sf <sub>1</sub> : (dbr: University of Kansas, dbo:city, dbr:Lawrence Kansas).
sf <sub>2</sub> : (dbr: University of Kansas, dbo:wikiPageWikiLink , dbr:University of Kansas Medical Center).
sq <sub>3</sub> : What is the name of the fight song of the university?
sp1: (Kansas Song) :Kansas Song (We're From Kansas) is a fight song of the University of Kansas.

 $sf_{i}: (dbr: Kansas \ Song, \ dbo: wikiPageWikiLink, \ dbr: University \ of \ Kansas).$