

Improving Embedded Knowledge Graph Multi-hop Question Answering by Introducing Relational Chain Reasoning

Weiqliang Jin¹, Biao Zhao¹, Hang Yu^{2*}, Xi Tao², Ruiping Yin^{3,4} and Guizhong Liu¹

¹School of Information and Communications Engineering, Xi'an Jiaotong University, Xi'an, 710049, Shaanxi, China.

²School of Computer Engineering and Science, Shanghai University, Baoshan, 200444, Shanghai, China.

³Information Faculty of Computer School, Beijing University of Technology, Chaoyang, 100124, Beijing, China.

⁴Engineering Research Center of Intelligence Perception and Autonomous Control, Ministry of Education, 100124, Beijing, China.

*Corresponding author(s). E-mail(s): yuhang@shu.edu.cn;
Contributing authors: weiqliangjin@stu.xjtu.edu.cn;
biaozhao@xjtu.edu.cn; 20721546@shu.edu.cn;
yinruiping@bjut.edu.cn; liugz@xjtu.edu.cn;

Abstract

Knowledge Graph Question Answering (KGQA) aims to answer user-questions from a knowledge graph (KG) by identifying the reasoning relations between topic entity and answer. As a complex branch task of KGQA, multi-hop KGQA requires reasoning over the multi-hop relational chain preserved in KG to arrive at the right answer. Despite recent successes, the existing works on answering multi-hop complex questions still face the following challenges: i) The absence of an explicit relational chain order reflected in user-question stems from a misunderstanding of a user's intentions. ii) Incorrectly capturing relational types on weak supervision of which dataset lacks intermediate reasoning chain annotations due to expensive labeling cost. iii) Failing to consider implicit relations between the topic entity and the

answer implied in structured KG because of limited neighborhoods size constraint in subgraph retrieval-based algorithms. To address these issues in multi-hop KGQA, we propose a novel model herein, namely Relational Chain based Embedded KGQA (Rce-KGQA), which simultaneously utilizes the explicit relational chain revealed in natural language question and the implicit relational chain stored in structured KG. Our extensive empirical study on three open-domain benchmarks proves that our method significantly outperforms the state-of-the-art counterparts like GraftNet, PullNet and EmbedKGQA. Comprehensive ablation experiments also verify the effectiveness of our method on the multi-hop KGQA task. We have made our model's source code available at github: <https://github.com/albert-jin/Rce-KGQA>.

Keywords: Data Mining and Search, Question Answering, Knowledge Graph based Multi-hop QA, Neural Semantic Parsing, Knowledge Graph Embedding

1 Introduction

Knowledge Base Question Answering (KBQA) [1] is an attractive service mining and analytics method that has attracted extensive attention from academic and industrial circles in recent years. Given a natural language question, the KBQA system aims to answer the correct target entities from a given knowledge base (KB) [2]. It relies on certain capabilities including capturing rich semantic information to understand natural language questions clearly and seek correct answers in large scale structured knowledge databases accurately. Knowledge Graph Question Answering (KGQA) [3, 4] is a popular research branch of KBQA which uses a knowledge graph (KG) as its knowledge source [2, 5] and uses factoid triples stored in KG to answer natural language questions. Thanks to KG's unique data structure and its efficient querying capability, users can benefit from a more efficient acquisition of the substantial and valuable KG knowledge, and gain excellent customer experience.

Early works [6, 7] on KGQA focus on answering a simple question, where only a single relation between the topic entity and the answer are involved. For example, in the question "What films did [Martin Lawrence] act in?", as depicted in Fig.1, there only exists a single relation 'starred_actors' between the topic entity 'Martin Lawrence' and the answer. The final answer only relies on just a single KG fact (Martin Lawrence, starred_actors.reverse, Black Knight). To solve simple question tasks, most traditional methods [8, 9] create diverse pre-defined manual templates and then utilize these templates to map unstructured questions into structured logical forms. Unfortunately, these pre-defined templates and hand-crafted syntactic rules are both labor-intensive and expensive. Moreover, such approaches require crowd workers to be familiar with linguistic and specific domain expert knowledge. Due to the dependency on large-scale fixed rules and manual templates, these methods cannot handle complex questions which require multiple relations inferences.

To make KGQA more applicable in realistic application scenarios, researchers have shifted their attentions from simple questions to complex ones. Knowledge graph multi-hop question answering is a challenging task which aims to seek answers which is multiple hops away from the topic entity in the knowledge graph. For example, the question “Who directed the films which [Martin Lawrence] acted in?” is a complex multi-hop question which requires a relational chain (starred_actors, directed_by) which has multiple relationships to arrive at the corresponding answers. This task is a relatively complex task compared with its simple counterparts due to the multi-hop relations retrieval procedure [5] which requires more than one KG fact in the inference.

Previous approaches [10–13] for handling complex question answering constructed a specialized pipeline consisting of multiple semantic parsing and knowledge graph answer-retrieving modules to complete the KGQA task. However, most have several drawbacks and encounter various challenges. We summarize these challenges as follows:

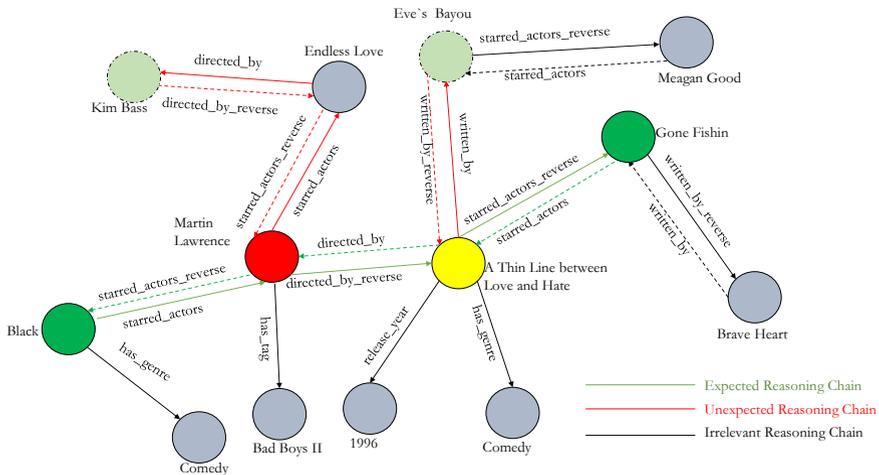


Fig. 1 The Freebase subgraph entered on the topic entity [Martin Lawrence] of the example questions. The red, yellow, green, green with an imaginary line, and grey circles denote the topic, intermediate, expected, unexpected and irrelevant entity nodes, respectively. The green, red and grey colored edges indicate the correct, incorrect and irrelevant reasoning chains, respectively.

Unexpected Relation Type Recognition. As the first step of KGQA, a semantic parser of most existing methods performs with poor accuracy in recognizing the correct relational type implied in questions, which hinders downstream answering reasoning. For example, as shown in Fig.1, let us consider questions where the topic entity and answer are connected by a multiple-hop reasoning chain, e.g., “Who acted in the movies directed by the director [Martin Lawrence]?”. To answer this type of question, the related two facts (Martin Lawrence, directed_by_reverse) and (A Tine Line,

starred_actors_reverse, Gone Fishin) help derive the answers within the neighborhood of the topic entity [Martin Lawrence]. Typically, these methods are prone to encounter incorrect relational reasoning ($AB \rightarrow AC$) when we mistake the unexpected relational chain (starred_actors_reverse, written_by) for the expected relational chain (starred_actors_reverse, directed_by_reverse). Thus, it is necessary to optimize relational semantics parsing for more accurate user intention recognition.

Unexpected Relation Order Recognition. Semantic parsing-based [14–17] methods mostly do not effectively capitalize on the correlation information of relationship order and direction from user-question expression. They become more susceptible to incorrect understanding when the questions are complicated from both semantic and syntactic aspects. The accuracy rate of parsing syntactics can be dramatically decreased by those with long-distance dependency. More especially, tracing back to the above question example and Fig.1, in addition to the correct chain (with green arrows), the spurious multi-hop chain (with red arrows) from entity node [Martin Lawrence] to [Kim Bass] can lead to incorrect reasoning results when the semantic parser module fails to parse such semantics as, reversing ($AB \rightarrow BA$) or shuffling ($ABC \rightarrow BCA$), the correct order of the relational chain. In short, we need an accurately capture of longer ordered-relational mappings implied in user language expressions to reach correct answers.

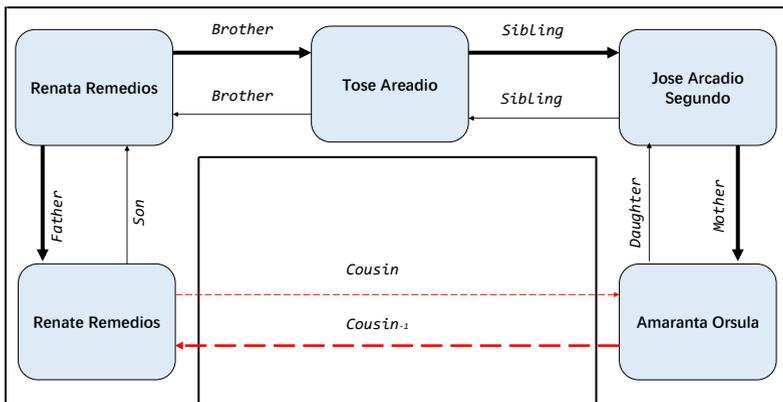


Fig. 2 A structured knowledge graph related to the user-question, “Who is [Renate Remedios]’s cousin?”. The illustration typically introduces implicit relations between the topic entity and the answer that hides in KG.

Implicit Relation Reasoning & Subgraph Neighborhood Constraint. Most mainstream KGQA methods [18, 19] cannot indirectly capture knowledge of implicit relational chains for reasoning due to the limiting constraint of neighborhood size constraint. All the answers are provided by

retrieving the extracted question subgraph. Let us consider the question “Who is [Renate Remedio]’s cousin?”. As depicted in Fig.2, the corresponding knowledge graph has no direct relational chain between the topic entity [Renate Remedio] and answer [Amaranta Orsule]. In other words, for future successive KGQA solutions, it is important to be able to discover the implicit factoid knowledge (Remedio \leftarrow Cousin \rightarrow Orsule) in the incomplete KG by the explicit relational chain (Remedio \leftarrow Father - Siblings - Mother \rightarrow Orsule), similar to the KG-based link-prediction task [20]. Furthermore, most existing methods labor under the undesirable constraint of answer detection from a pre-specified localized KG sub-graph neighborhood. For example, the state-of-the-art (SOTA) method GraftNet [21] whose answer is restricted to being a subset of the entities present in a localized KG sub-graph neighborhood, only reports a recall around 0.55 on an incomplete KG where only half of the original triples are presented.

To alleviate these limitations and challenges for the multi-hop KGQA task, our paper introduces a novel architecture, namely Relational Chain based Embedded KGQA, which supports the integration to learn the explicit relational mappings implied in the user’s expression and the implicit knowledge from the structured KG, similar to link prediction. Our proposed approach counters these limitations by simultaneously using the explicit semantic relational chain described in the question and the implicit relational chain between the structured KG nodes. We use the knowledge graph embedding and construct the *Answer Filtering Module* to calculate the mutual relationship between the topic entity and answer. Motivated by the previous work Embed-KGQA [22], we show how our model leverages an end-to-end neural network that employs the KG entity and relation embeddings to provide complex questions with answers from the KG. Since our model replaces the traditional pipeline procedure of generating and retrieving a localized subgraph at intermediate reasoning steps, it helps to decrease memory costs efficiently and obtain computational efficiency. To obtain a more competitive performance in large-scale KG, we apply an extra reasoning procedure called the *Relational Chain Reasoning Module* to prune the candidate entities ranked by the *Answer Filtering Module*. We apply a Siamese architecture [23] based on the *long short-term memory* (LSTM) [24] and transformer *RoBERTa* [25] to learn the semantic similarity between the relational chain of the problem description and the KG factoid relational chain. It also leverages the external supervised signal of the relational chain from the training sample. By calculating the semantic similarity between question semantics and candidate entity retrieval chain, we can further determine the final answer more accurately. The internal construction details of our model are introduced in Sec. 4.

We summarize the contributions of this paper as follows:

1. We propose a novel approach namely Relational Chain based Embedded KGQA which includes two main modules: the *Answer Filtering Module* and the *Relational Chain Reasoning Module*. Moving away from previous studies, our model simultaneously takes advantage of the knowledge graph embedding

and training with weak supervision by predicting the intermediate relational chain implied in KG to perform the multi-hop KGQA task.

2. We introduce the *Answer Filtering Module*, a knowledge graph embedding based end-to-end network for preliminary answer filtering. This module can address the problem of inadequacy that is a factor in the missing links in the incomplete knowledge graph thanks to its capability of capturing implicit KG relationships. Furthermore, we consider all entities as candidate answers in this step, so our model won't suffer from the out-of-reach issues brought by limited subgraph neighborhood constraint.

3. Our proposed *Relational Chain Reasoning Module* can help capture the multi-hop relations surrounding the topic node to support the results more accurately. We apply the Siamese network [23] to calculate the vector representation-based semantic similarity score between the user's question and KG structured knowledge. To the best of our knowledge, our proposed sub-module is the first to consider the question relational direction and order information by using the Siamese network.

4. Our experimental results on three widely adopted KGQA benchmarks demonstrate our method's competitive capability compared with most SOTA methods (average 1.2% absolute improvement across *hit@1* evaluation metric). Furthermore, using an extensive ablation study, we demonstrate the superiority and effectiveness of our proposed model for the multi-hop KGQA task.

The rest of the paper is organized as follows: We first provide a thorough review of the related KGQA works in Sec. 2. Next, we introduce the preliminary knowledge about the KGQA task in Sec. 3. Following the internal structure of our model, we then explicate our two features: *Answer Filtering Module* in Sec. 4.2 and *Relational Chain Reasoning Module* in Sec. 4.3. Sec. 5 describes the experimental details on three open-domain datasets. Finally, in Sec. 6, we conclude our contributions to this work and suggest several promising innovations for our Relational Chain based Embedded KGQA in the future.

2 Related Work

Our work is closely related to the Multi-hop Knowledge Graph Question Answering, Knowledge Graph Embedding, Siamese Network, and Pretrained Language Model.

2.1 Multi-hop KGQA

Multi-hop Knowledge Base Question Answering comprises two mainstream branches: Information Retrieval-based (IR) and Semantic Parsing-based (SP-based). The most popular methods fall into these two categories.

2.1.1 Semantic Parsing Methods

SP-based approaches follow a parse-then-execute procedure. These methods [22, 26–29] can be summarized as the following steps: (1) *question semantic understanding*: parsing relations and subjects involved in complex questions, (2) *logical formula construction*: decode the subgraph into an executable logic form such as high-level programming languages or structured queries such as SPARQL, (3) *KG-based positioning and querying*: search from KGs using the query language and provide query results as the final answer. Owing to its intermediate procedure of generating expressive logic forms, SP-based methods are more interpretable than their IR-based methods counterparts. Nonetheless, for most existing SP-based methods, more relations in complex questions indicate a larger search space of potential logic forms for parsing, which will dramatically increase the computational cost.

Yu et al. [30] pointed out that KG relation detection is a core component and entity linking is a key step in KGQA tasks. To improve the recognition accuracy of both sub-tasks, they proposed the Hierarchical Residual BiLSTM (HR-BiLSTM) to encode question descriptions and word-level and phrase-level relationship path. The new HR-BiLSTM module calculates the similarity scores for all the questions and textual relationships, which integrates for these two components entity linking and relationship path identification into a single step and enhances each other. When in inference, the model only selects the highly-scored (relations, topic entity) pairs as correct answers from candidates.

Miller et al. [31] proposed an ideal domain-specific KGQA framework, called Key Value-Memory networks (KV-MemNN), which has proved to be effective to support answer reasoning over specific domain multi-source knowledge like textual documents and structured KG. It performs QA by employing a widely used long-term memory mechanism to reason on a key-value structured memory network. They defined three operations, i.e., key hashing, the model first fetches all KG triples relevant to given questions and then stores their topic entities and relationships in the key slot, tail entity in the value slot; key addressing, the model assigns each memory unit with a normalized relevance weight by the dot product operation as the relevance probability between the question and each key representations in the memory; finally, value reading, where the model reads the values of all addressed memories by taking their weighted sums of all values and relevance weights, and use the outputs to represent intermediate reasoning results, which is then used to update the question representation. To obtain the final prediction over all candidate answers, the model repeats the key addressing and value reading steps in the Ranking component several times.

However, the KV-MemNN obviously presents the following challenges: 1) It often fails to precisely update multi-relation question queries during multiple memory reading. 2) It reads the memory repeatedly since they can not well determine accurately when to stop. 3) It focuses more on memory facts understanding rather than the properly questions understanding, so it does

not perform as well as expected when applied to the scenario where its questions are complicated and associated with complex constraints, such as an open-domain KGQA task. 4) It selects the candidate with the highest similarity score as the only answer in default. So, it conducts inefficiently when the questions contain more than one answer.

To solve these challenges, Xu et al. [32] proposed an interpretable mechanism to enable a basic KV-MemNN model to work for complex questions, which yielded state-of-the-art performances on three benchmarks. Enhanced KV-MemNN introduced a novel **STOP** strategy into multi-hop memory reading to generate a flexible number of queries and introduce a new query updating method, which considers the already-addressed keys in previous hops as well as the value representations that avoids repeated or invalid memory readings. For multi-constraint questions, the model considers the value representation of each hop by accumulating all the value representations of both current and previous hops to address each relevant constraint at different hops.

In addition to the above representative methods, many knowledge base question answering approaches based on Graph Neural Network (GNN) [33] and Graph Convolutional Network (GCN) [34] have been proposed in recent years, approaches such as Graph Convolutional Network-based Multi-Relation Question Answering system (QAGCN) [35] and Case-Based Reasoning SUB-Graph model (CBR-SUBG) [36]. As the name suggests, QAGCN is a simple but effective model that leverages attentional graph convolutional networks that can perform multi-step reasoning during the encoding of knowledge graphs. Able to leverage highly-efficient embedding computations, the model's significant advantage is that it can essentially simplify complex reasoning mechanisms. CBR-SUBG is a semiparametric model for weakly-supervised KGQA that retrieves similar queries and utilizes the similarities in graph structure of local subgraphs to answer a query. It contains a parametric component comprising a graph neural network (GNN). Through experiments, it performs competitively with state-of-the-art KGQA models on multiple benchmarks. Due to its capacity for reasoning pattern identification, the method CBR-SUBG can also provide interpretable paths for returned answers, which could bring slightly better interpretability.

2.1.2 Information Retrieval Methods

IR-based approaches typically include a series of procedures as follows: question-specific graph extraction, question semantics representation, extracted graph-based reasoning and answer candidates ranking. Given a complex question description, these methods [37, 38] first construct a question-specific subgraph which includes all question-related entity nodes and relation edges from KGs without generating an executable logic formula, This is followed by employing a question representation module to encode user-question tokens as low-dimensional vectors. Secondly, an extracted-graph based reasoning module conducts a semantic matching algorithm to aggregate the center entity's neighborhoods' information from the question-specific subgraph. At

the end of the reasoning, they rank all the entities' scores in the subgraph by applying an answer-ranking module to predict the top-ranked entities as the final answers. Based on feature representation technology, IR-based approaches can be divided into feature engineering-based approaches and representation learning-based approaches.

IR-based feature engineering approaches [39] rely on manually defined and extracted features, which are time-consuming and cannot detect the whole question semantics. To solve these problems, representation learning *IR-based* methods convert questions and related entities into distributed vector representations in the same dimension space and treat KGQA tasks as semantic matching between distributed representations of questions and candidate answers [2].

Sun et al. [21] propose a integrated framework namely *GRAFT-Net*, which is adopted an knowledge fusion strategy, where the answers are selected from a heterogeneous question-specific subgraph constructed from the KG and textual documents based on the given questions. The subgraph contains three factors: entity nodes, sentence nodes and a special type of edges which indicates the mutual relations between entity and sentence nodes. During answer detection, the convolution neural network *GRAFT-Net* spreads central entity node feature to neighboring nodes in several iterations and determines whether an entity node is an answer or not.

However, the automatically constructed subgraph in *GRAFT-Net* relies heavily on heuristic rules and can lead to serious error cascading and bring incorrect reasoning. Thus, soon after proposing the *GRAFT-Net* [21], Sun et al. [10] presented a learned iterative process for topic-entity-centric graph construction. The improved method, called *Pull-Net*, where the "pull" classifier is weakly supervised so that only QA pairs are used for supervision. It first selects seed entity nodes by *GRAFT-Net* and a novel classification model at each step. Then, more and more extra valuable entities and sentences are introduced into the current graph through several pre-defined operational iterations, with the final answer determined by the same procedure as *GRAFT-Net* [21]. Experimentally, *PullNet* improves dramatically over prior state-of-the-art methods [21, 38, 39] even under weakly supervised signals and incomplete KGs.

A significant challenge in multi-hop Knowledge Base Question Answering is the lack of supervision signals at intermediate steps. To address this challenge, He et al. [27] propose an elaborate teacher-student framework by adapting the generic Neural State Machine (NSM) [40] as the student network, while the teacher network aims to learn intermediate supervision signals to improve the student network. The extensive evaluation results with three benchmark datasets show that their proposed model is superior to previous methods in terms of effectiveness for the multi-hop KGQA task. Moreover, other detailed experiments prove that their approach is more flexible to extend itself to other neural architectures or learning strategies on graphs.

Apart from these traditional subgraph-generation methods, researchers also try to incorporate the KG embedding mechanism as extra information into

entity and relation representations to alleviate the incomplete KG sparsity problems. Inspired by relationship completion and missing link prediction tasks in the KGs, Saxena et al. [22] propose a novel framework, named *EmbedKGQA*, which leverages the pre-trained KG embeddings to enrich the learned entity and relation representations. Extensive comparative experiments on multiple benchmarks show that EmbedKGQA is particularly effective in performing multi-hop KGQA over sparse KGs.

2.2 Knowledge Graph Embedding

Knowledge Graph Embedding [28] is to embed a KG's factoid triples knowledge including all entities and relations into continuous and low-dimensional embedding representation space, such that the original entities and relations are well preserved in the vectors. Representative KG embedding models exploit the distance-based scoring function $f(\cdot)$ which is used to measure the plausibility of a triple (*topic entity, predicate, and tail entity*) as the distance between the head and tail entity such as TransE [41] and its extensions (e.g. TransH [42]), DistMult [43] and ComplEx [44]. In short, a typical KG embedding technique generally consists of three steps: (1) representing entities and relations, (2) defining a scoring function, and (3) learning entity and relation representations. Thanks to its ability to simplify the manipulation while preserving the KG inherent structure, it can benefit a variety of downstream tasks to take the entire KG into consideration, such as entity alignment [45], relation prediction [20] and even KGQA work [22]. The effectiveness of knowledge graph embedding in various real-world NLP tasks [46, 47] motivates us to explore its potential advantages in the KGQA task.

2.3 Siamese Network

The *Siamese network* [23] is a semantic textual similarity metric that is built on top of a feature representation network such as CNN [48] and RNN [23]. Given an example input pair, the Siamese network first maps these two inputs into sequences of the word embedding vector using large-scale pretrained embeddings like Glove [49], then passes these vectors through the representation extractor's forward procedure and get the semantic vector representations, respectively. Finally, the Siamese network applies ℓ^1 norm (Manhattan distance) or ℓ^2 (Euclidean distance) norm as the distance measurement function to calculate the similarity between these two representations. Furthermore, the long short-term memory (LSTM) is superior to the original RNNs for learning long-range dependencies because its memory cell units can capture rich features across lengthy language token sequences. Therefore, in this work, we employ the Siamese adaptation of the bidirectional LSTM network to learn the semantic relational chain.

2.4 Pretrained Language Model: RoBERTa

Bidirectional Encoder Representations from Transformers [50], or BERT, is a revolutionary self-supervised pretraining technique that learns to predict intentionally hidden (masked) sections of text. Crucially, the representations learned by BERT have been shown to generalize well to downstream tasks and, when BERT was first released in 2018, it achieved state-of-the-art results on many NLP benchmark datasets.

RoBERTa [25], as is proposed by Liu et al., is built on BERT’s language-masking strategy and modifies key hyperparameters in BERT, and it can be regarded as a heavily optimized version of BERT. It includes removing BERT’s next-sentence pretraining objective, and trains with much larger mini-batches and learning rates. It was also trained on an order of magnitude more data than BERT, for a longer amount of time. This allows RoBERTa representations to generalize even better to downstream tasks compared to BERT. The improved model RoBERTa achieves state-of-the-art results on GLUE, RACE and SQuAD benchmarks, without multi-task finetuning for GLUE or additional data for SQuAD.

3 Preliminaries

In this section, we formally introduce the preliminary knowledge [2, 5] on the multi-hop KGQA task formulation and its related definitions. Before the formulaic description, all the summarized pre-defined notations for the KGQA task are given as follows: we denote a KG as $\mathcal{G}(\varepsilon, \mathcal{R})$ in which ε, \mathcal{R} respectively denote the entities and relation set, and we use (h, ℓ, t) to represent a factoid triple in KG. We use an uppercase and lowercase letter to denote a matrix (e.g. \mathbf{W}) and a vector (e.g. \mathbf{v}). The ℓ^n norm of a vector is denoted as $\|\mathbf{p}\|_n$.

Definition 1 (Multi-hop question) [2, 5, 22] If a natural language question involves more than one predicate between the topic entity and answer, then we believe the answer is multiple hops away from the topic entity in the KG. Thus, we identify this as a multi-hop question. For example, let us consider the multi-hop question: “When did the film production company announce which actor also directed the movie [Cast a Deadly Spell]?”, which consists of several predicates which correspond to the KG relational links: *release_year*, *starred_actors*, *directed_by* respectively.

Definition 2 (Knowledge Graph Embedding) [51] The KG embedding algorithm [28, 51] aims to map all the KG components including entity and relation to a low-dimension and continuous vector space. Given a KG consisting of n entities and m relations, we firstly initialize the values of h, ℓ and t randomly. Then, a scoring function $f_\ell(h, t)$ which we defined measures the relation of a fact triple (h, ℓ, t) . Finally, the embedding algorithm utilizes a margin-based ranking criterion to optimize the embedding distribution that maximizes the overall plausibility of factoid triples (h, ℓ, t) and to minimize the plausibility of spurious triples (h', ℓ', t') simultaneously.

Definition 3 (Multi-hop KGQA task) [2, 5] The multi-hop question was introduced in Definition 1. In this section, we define a knowledge graph (KG) as \mathcal{G} . \mathcal{G} is a directed graph whose nodes represent entities and edges represent relations, and each triple in the KG represents an atomic realistic fact, such as (Joseph Robinette Biden, president_of, USA).

Formally, given a complex natural language question in the format of a sequence of tokens $\Pi = w_1, w_2, \dots, w_l$ and the available KG \mathcal{G} , the KGQA task first links the topic entity w_i, \dots, w_j to the KG \mathcal{G} . The subject mentioned in a question is also named as a topic entity. Then, it identifies the most possible KG relations which are related to the user’s question. Using these two steps, the goal of KGQA is to determine the factual answer with triples stored in KG, denoted by the set \mathcal{A}_q , to query q from the candidate entities E by leveraging the topic entity and related relations in KG. Specifically, we focus on solving complex question answering, termed the *multi-hop* KGQA task, where the answer is multiple hops away from the topic entity in a knowledge graph, which means these questions require more than one KG triple.

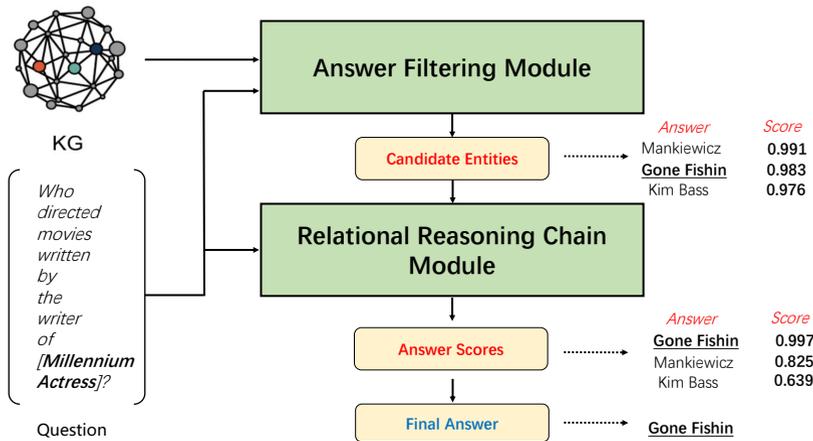


Fig. 3 This figure shows our overall pipeline architecture for the multi-hop KGQA task. The green rectangles denote our two sub-modules, the solid arrows and dashed arrows indicate the information flowing through our model and the intermediate results. The next figures [5, 6] also illustrate this by using the typical user’s complex question “Who acted in the movies directed by the director [Martin Lawrence]?”

4 Our Proposed Model

Our *Relational Chain based Embedded KGQA* is a two-stage pipeline model which consists of two components: *Answer Filtering Module* and *Relational Chain Reasoning Module*.

4.1 Overview

As illustrated in Fig.3, given a real-world question and an available KG, the *Answer Filtering Module* first jointly leverages topic entity embedding and question representation to score all possible candidate entities in this KG to provide a set of pruned candidate answers for this question. However, the entity nodes in a KG are often on a scale as large as a million, hence it could be noisy and inaccurate when comparing the topic entity with all other entities $\hat{\mathbf{t}}$. To make the learning more efficient and accurate, we do not directly select the top-1 scoring entity from the sorted entities as our final answer. Instead, we introduce the extra module *Relational Chain Reasoning Module* to take the relation type and order of semantic relational chain into consideration for a higher *hit@1* accuracy result.

Before being fed into the next stage, we transform these intermediate candidate entities to their shortest relational chains which point to the question's topic entity by retrieving them in KG and mapping these ordered chains to sequences of embedding which correspond with our embedded KG. The *Relational Chain Reasoning Module* receives the intermediate results generated by the last step. Then, it simultaneously utilizes the relational chain sequences and user-question to measure the mutual similarity score through our Siamese network. Taking the question relational chain reasoning details into consideration can help increase the accuracy of answer prediction compared with the first stage, the *Answer Filtering Module*. Finally, after sorting the scored candidates, we choose the entity which has the highest similarity score as the final answer. Fig.4 formally illustrates the algorithm of how our method works and predicts the final answer for a given multi-hop question, where ℓ denotes the question semantic representation, and ϕ denotes the **Complex** scorer.

Algorithm 1 Stepwise Answer Reasoning for Rce-KGQA

Require: Knowledge Graph $\mathcal{G}(\varepsilon, \mathcal{R})$; L length Question q , denote $\{\mathbf{w}_j\}_{j=1}^L$; question topic entity \mathbf{h} in KG;

Ensure: Predicted Answer $e_{ans} = \underset{t' \in \mathcal{A}}{\operatorname{argmax}} \phi(h, \ell, t')$

- 1: Through **Complex**, we obtain KG relation embeddings \mathbf{W}_r and entity embeddings \mathbf{W}_e where $e \in \varepsilon$ and $r \in \mathcal{R}$.
 - 2: **for all** t' such that $t' \neq \mathbf{h} \cap t' \in \varepsilon$ **do**
 - 3: $\ell = \operatorname{Attn}(\operatorname{BiLSTM}(\{\mathbf{w}_j\}_{j=1}^L))$
 - 4: $\operatorname{Score}(t') = \phi(h, \ell, t')$
 - 5: **end for**
 - 6: Select top-N ($N \in \{5, 10, 15\}$) scoring tail entities from last step, denotes $\{t''\}_{i=1}^N$
 - 7: **for all** $t'' \in \{t''\}_{i=1}^N$ **do**
 - 8: Retrieval its shortest relational chain $r_i^* = \{\mathbf{r}_j\}_{j=1}^M$, where M is the chain length.
 - 9: **end for**
 - 10: Collect all candidates and its relational chain as a set $\{[t'', r_i^*]\}_{i=1}^N$.
 - 11: **for all** $[t'', r_i^*] \in \{[t'', r_i^*]\}_{i=1}^N$ **do**
 - 12: $\mathcal{V}_q = f_{c2}(\operatorname{dropout}(\operatorname{relu}(f_{c1}(\operatorname{RoBERTa}(\{\mathbf{w}_j\}_{j=1}^L))))))$
 - 13: $\mathcal{V}_{r_i^*} = \operatorname{Attn}(\operatorname{BiLSTM}(r_i^*))$
 - 14: $\operatorname{Score}(\mathcal{V}_q, \mathcal{V}_{r_i^*}) = \exp\left(-\left\|\mathcal{V}_q^{(2)} - \mathcal{V}_{r_i^*}^{(2)}\right\|_2\right)$
 - 15: **end for**
 - 16: Select the highest scoring entity $\underset{t''}{\operatorname{argmax}} \operatorname{Score}(\mathcal{V}_q, \mathcal{V}_{r_i^*})$ from $\{[t'', r_i^*]\}_{i=1}^N$, denote as e_{ans} .
-

Fig. 4 The formulaic illustration of our Rce-KGQA model.

4.2 Answer Filtering Module

As the first step of our model, our *Answer Filtering Module* aims to filter an entity set from all KG entities as candidate answers via three steps. These three operational steps are illustrated in Fig.5 and relate respectively to three sub-modules: Graph Embedding Generator, Question Semantic Parser and Answer Scorer. We introduce each sub-module followed by the system processing order.

4.2.1 Graph Embedding Generator

Traditional solutions could not handle many scenario problems such as *Implicit Relation Reasoning* and *Subgraph Neighborhood Constraint*. Inspired by the competitive performance of previous work EmbedKGQA [22], we observe that the global relation knowledge and structure information preserved in KG embedding could potentially be used to resolve these issues efficiently and improve the overall accuracy of question answering.

In this work, our used KG is also embedded in continuous low-dimensional vector space to obtain all sparse representations for all entities and relations that existed in KG so we can simplify computations on the KG. We apply the *Complex Embeddings* (ComplEx) [44] approach to embed relations and entities in complex vector space. Compared with traditional KG embedding methods like *TransE* [41] and its extensions, semantic matching models such as Holographic Embeddings [52], ComplEx and RESCAL [53] have shown that they can generally yield better results. All KG embeddings are initialized randomly from uniform distributions. Generally, the hyper-parameters about entity and relation embedding dimension are not less than 100 and, in this paper, we set embedding dimension at 200 which follows previous similar works.

In each training step, positive facts which present correct real-world relational triples are sampled from all factoid triples existing in KG. Negative facts which present fake factoid relational triples are generated from negative sampling in the negative example generation step, where it randomly replaces the tail entity with an incorrect entity or replaces the relation with incorrect relation.

Considering the samples count ratio of positive to negative ones, Trouillon et al. [44] further investigated the influence of different numbers of negative samples for each positive one. Their work demonstrates that generating more negatives usually leads to better performance, and around fifty negatives per positive example is an appropriate trade-off between reasoning accuracy and training cost. So, our implementations also follow this prior setting.

Given $h, t \in \varepsilon$ and $\ell \in \mathcal{R}$, this embedding approach would provide $v_h, v_\ell, v_t \in \mathbf{C}^d$ for each relation triple (h', ℓ', t') , and the scoring function is defined as follows:

$$\begin{aligned} \phi(h', \ell', t') &= \text{Re}(\langle v_h, v_r, \bar{v}_t \rangle) \\ &= \text{Re} \left(\sum_{k=1}^d v_h^{(k)} v_r^{(k)} \bar{v}_t^{(k)} \right) \end{aligned} \quad (1)$$

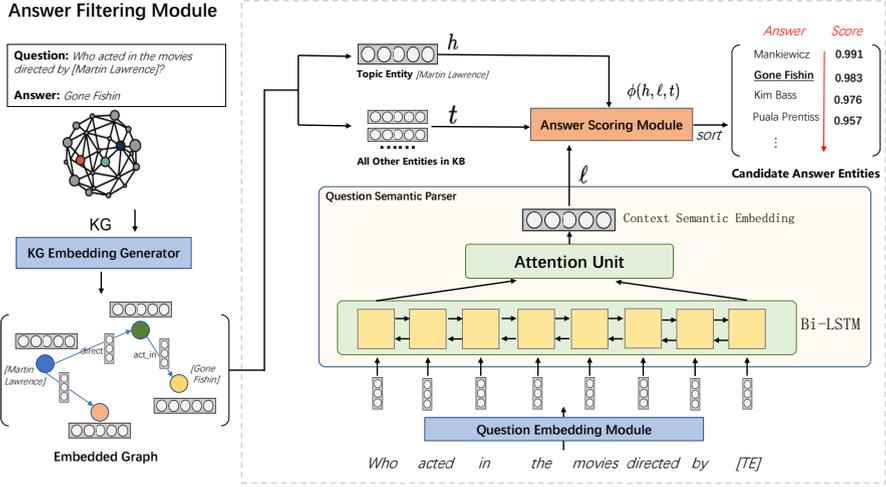


Fig. 5 Structure illustration of the Answer Filtering Module.

$$\phi(h', \ell', t') > 0 \quad \forall a \in \mathcal{A} \quad (2)$$

$$\phi(h', \bar{\ell}', \bar{t}') < 0 \quad \forall \phi(h', \bar{\ell}', \bar{t}') \notin \mathcal{A} \quad (3)$$

where the $Re(\cdot)$ means taking the real part of a complex value, the \bar{v}_t denotes the conjugate of v_t , the $\bar{\ell}', \bar{t}'$ is the random replaced wrong relation and wrong tail entity of the ℓ', t' , and the \mathcal{A} means the set including all real-world knowledge triples.

Optimization Eq.1 aims to minimize the values for all false triples less than 0, Eq.3 and maximizes the values for all true triples greater than 0, Eq.2. It can be easily carried out by stochastic gradient descent (SGD) or Adam optimizer at each training iteration. Lastly, the original structure and relation information in the KG are preserved in these learned vectors, which helps efficient completion of the downstream procedures.

4.2.2 Question Semantic Parser

In this section, we introduce the *Question Semantic Parser*, which consists of a recurrent neural network (bidirectional-LSTM) and extra self-attention operation, to help represent the question's meanings. During the inference procedure, the *Question Semantic Parser* takes a question as the input and provides a predicted vector $\hat{\ell}$ as this question's relationship representation between the topic and answer in KG.

As shown in Fig.5, we build this sub-module based on the hierarchical neural network. Firstly, we encode the L length question $\{t_j\}$, for $j = 1, \dots, L$ into a sequence of word embedding vectors $\{v_j\}$ through our word embedding layer whose parameters are learnable during the training procedure. The word

embedding dimension is consistent with the recurrent network’s hidden dimension. Then we employ a single layer bidirectional LSTM to learn a forward hidden state sequence $(\vec{h}_1, \vec{h}_2, \dots, \vec{h}_L)$ and a backward hidden state sequence $(\overleftarrow{h}_1, \overleftarrow{h}_2, \dots, \overleftarrow{h}_L)$. Compared to general RNN, bidirectional LSTM is a special RNN which mainly solves the gradient disappearance and gradient explosion challenges and captures better long-distance semantics during long sequence modeling. Our LSTM component uses 256 as its dimension of hidden representations h_t and memory cells c_t . It is well known that the performance of LSTMs depends crucially on their initialization and offers a strong starting point to facilitate model convergence, so we initialize our bidirectional LSTM weights with Xavier [54] initialization which is markedly superior to other initialization methods such as Gaussian, Uniform and Kaiming initialization[55].

Taking the forward step as an example, the next state \vec{h}_j based on last state \vec{h}_{j-1} is computed via the following operations.

$$\mathbf{f}_j = \sigma \left(\mathbf{W}_{xf} \mathbf{x}_j + \mathbf{W}_{hf} \overrightarrow{\mathbf{h}}_{j-1} + \mathbf{b}_f \right) \quad (4)$$

$$\mathbf{i}_j = \sigma \left(\mathbf{W}_{xi} \mathbf{x}_j + \mathbf{W}_{hi} \overrightarrow{\mathbf{h}}_{j-1} + \mathbf{b}_i \right) \quad (5)$$

$$\mathbf{o}_j = \sigma \left(\mathbf{W}_{xo} \mathbf{x}_j + \mathbf{W}_{ho} \overrightarrow{\mathbf{h}}_{j-1} + \mathbf{b}_o \right) \quad (6)$$

$$\mathbf{c}_j = \mathbf{f}_j \circ \mathbf{c}_{j-1} + \mathbf{i}_j \tanh \left(\mathbf{W}_{xc} \mathbf{x}_j + \mathbf{W}_{hc} \overrightarrow{\mathbf{h}}_{j-1} + \mathbf{b}_c \right) \quad (7)$$

$$\overrightarrow{\mathbf{h}}_j = \mathbf{o}_j \circ \tanh(\mathbf{c}_j) \quad (8)$$

The variables f_j, i_j, o_j in the above equations are the input, forget and output gate’s activation vectors respectively, where \mathbf{c}_{j-1} and \mathbf{c}_j are the cell state vectors in the $j - 1$ time and j time, \tanh and σ are the hyperbolic tangent and sigmoid functions. Eq.4 denotes the forget gate operation, which aims to control whether to forget the hidden cell state’s part information of the last moment with a certain probability. Eq.5 denotes the input gate operation, which is responsible for processing the current sequence input. Eq.6 denotes the output gate operation, which determines the degree to which information is updated and output. Eqs. 7 and 8 are the steps to update the old cell state, which is determined jointly by the state of the previous sequence, this sequence’s current input and the activation function. After the information flows through LSTM, we concatenate the forward \vec{h}_j and backward \overleftarrow{h}_j and obtain the combined features $\mathbf{h}_j = [\vec{h}_j; \overleftarrow{h}_j]$. After that, the last hidden state is considered to be the question semantic representation.

Different word tokens make different contributions to the relationship semantic recognition. For example, words which are prepositions and articles are more irrelevant for discovering question semantics than relational demonstrators. Thus, after LSTM we apply the self-attention mechanism to capture

more valuable features. The attention operation details are shown in Eq.9 and Eq.10. Given an LSTM hidden representation, a full connect layer, activation function \tanh and softmax operation will jointly generate the attention weight α_j at first. Then, as shown in Eq.10, the final attention vector representation s_j which is the semantic representation of the language question is aggregated by the weighted sum operation of \mathbf{h}_j and a_j .

$$\alpha_j = \frac{\exp(a_j)}{\sum_{i=1}^L \exp(a_i)} \quad \text{where} \quad a_j = \tanh(\mathbf{w}_a^\top [\mathbf{h}_j] + b) \quad (9)$$

$$s_j = \sum_i \alpha_{ij} \mathbf{h}_{ij} \quad (10)$$

All the weight matrices, weight vector \mathbf{W} , and bias terms are calculated based on the training data, i.e. LSTM gate unit weight matrix $\{\mathbf{W}_f, \mathbf{W}_i, \mathbf{W}_o\}$ and attention weight matrix \mathbf{w}_a^\top . In this way, we obtain the rich relationship semantics implied in natural language questions for answer reasoning.

4.2.3 Answer Scorer

Like the ComplEx [44] scoring function which depicted in Eqs. 1, 2, and 3, as shown in Eq.11, we learn an answer ranker $Rank(t)$ for each candidate \mathbf{t} , namely *Answer Scorer* to score the (topic entity, relationship semantic) pair against all possible KG entities $\mathbf{t} \in \varepsilon$ by maximizing the probabilities of positive samples $t \in \mathcal{A}$ and minimizing the negative sample $t' \notin \mathcal{A}$, where the \mathcal{A} means the set including all real-world knowledge triples.

$$Rank(t) = \begin{cases} \max(\phi(h, \ell, t)), & \forall t \in \mathcal{A} \\ \min(\phi(h, \ell, t')), & \forall t' \notin \mathcal{A} \end{cases} \quad (11)$$

Since our *Question Semantic Parser* is designed to fit realistic relationship features, all the pretrained KG entity embeddings are frozen during the model convergence procedure.

Instead of simply selecting the entity with the highest score due to its low accuracy but high recall performance, we conduct a rough filtering by selecting top-n, where $n \in \{5, 10, 15\}$ to obtain the intermediate scored candidate entities that have a high answer recall rate. For the inference, the *Answer Scoring Module* gives each candidate a plausibility score to indicate its answer confidence and filter out the top-n scored intermediate result which is fed into the next step, the *Relational Chain Reasoning* module.

4.3 Relational Chain Reasoning Module

Our available KG often contains a large number of entities and has enormous factoid triples, and it could be inaccurate when comparing all candidate embedding representations against with each other. Specifically, after training the *Answer Filtering Module*, we obtain all the scored entities for each training sample. During the prediction result analysis, we observe that the model performance on *hit@5* outperforms that on *hit@1* metric, which could be due to

the influence of a large number of noisy entities number in the large-scale KG. Furthermore, since the answer is given as the only ground-truth information, a major challenge for multi-hop KGQA is that it usually lacks intermediate reasoning supervision signals.

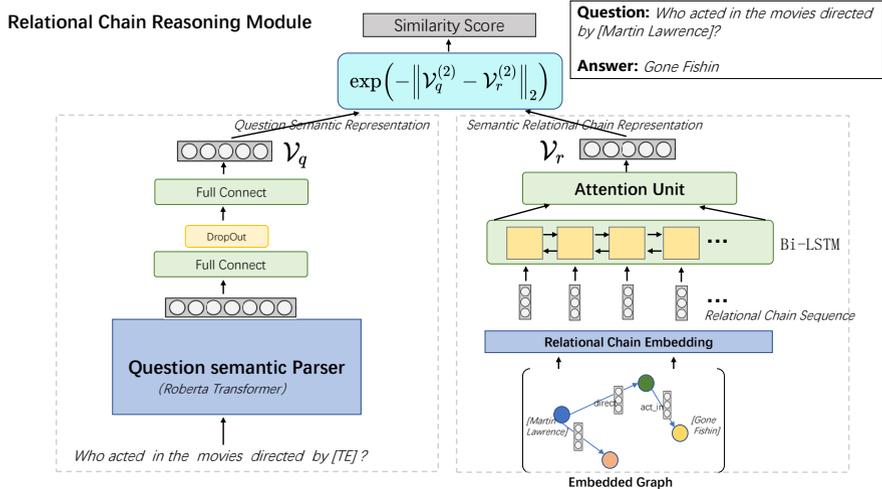


Fig. 6 Structure of the designed Relational Chain Reasoning Module.

To tackle these issues, we propose an extra component for our KGQA work, termed the *Relational Chain Reasoning Module*. As the final and important step of our approach, it aims to improve the reasoning accuracy through considering the reasoning chain order and its relational type under a weak supervised situation. Its training procedure is irrelevant to the *Answer Filtering Module*, but its training dataset is constructed from the *Answer Filtering Module*'s prediction results.

Formally, we use the trained *Answer Filtering Module* to obtain the sorted scored entities $\{[e_{it}; s_{it}]\}_{i=1}^n$, $\{[e_{iv}; s_{iv}]\}_{i=1}^n$ and $\{[e_{ie}; s_{ie}]\}_{i=1}^n$ for the training, validating, and testing datasets, respectively. Then, these scored results are truncated by only reserving top- n ($n \in \{5, 10, 15\}$) candidates. Next, looking at the rough filtered results in each experimental sample, if it belongs to the correct answers, we construct a corresponding positive sample $[\mathbf{q}; \{\mathbf{r}\}_{i=1}^c]$, in which \mathbf{q} denotes the question tokens and $\{\mathbf{r}\}_{i=1}^c$ is the c length shortest path searched by the graph retrieval algorithm, as do the negative samples that are not inside in the correct answers.

4.3.1 Siamese Network Based Similarity Scoring

As shown in Fig.6, we introduce our novel sub-module *Siamese Network Based Similarity Scoring*, which is essentially a Siamese network [23]. The module's two feature detectors aim to extract the question Eq.12 and the relational

chain semantic representations Eq.13. They are constructed by a pretrained transformer **RoBERTa** [25] to generate question vector \mathcal{V}_q from question $q = \{\mathbf{w}_i\}_{i=1}^N$ and a single-layer bidirectional **LSTM** to generate relational chain vector \mathcal{V}_r from relational chain $r^* = \{\mathbf{r}_j\}_{j=1}^M$. RoBERTa is a revolutionary self-supervised pretraining technique that learns to predict intentionally hidden sections of language text. Moreover, the representations learned by RoBERTa has been shown it can generalize outcomes superior even to many NLP downstream tasks compared to original BERT [50]. Structure details will be described in the next successive two subsection. After information passes through these two encoded networks, the semantic feature similarities Eq14 in the vector representation space are subsequently used to infer the semantic similarity between the question and relational chain from the topic entity to answer.

$$\mathcal{V}_q = W_2^\top (\text{drop}(W_1^\top (\text{RoBERTa}(q)) + b_1)) + b_2 \quad (12)$$

$$\mathcal{V}_r = \text{Attn}(\text{BiLSTM}(r^*)) \quad (13)$$

$$\text{Score}(\mathcal{V}_q, \mathcal{V}_r) = \exp\left(-\left\|\mathcal{V}_q^{(2)} - \mathcal{V}_r^{(2)}\right\|_2\right) \quad (14)$$

Here, we explain why we use a semantic similarity score between relation chain \mathcal{V}_r and question representations \mathcal{V}_q to determine the right answer in detail. Taking the above question “Who acted in the movies directed by the director [Martin Lawrence]?” as an example, we denote the natural language tokens as q . Inspired by literature in recent years such as [2, 36], we intuitively suppose that the implied question semantic is similar and its representation is closer to the relational chain representation in vector space, “*directed_by_reverse* \rightarrow *starred_actors_reverse*”, which is correct in both order and type. It is inconsistent and far away from the relational chains which have the wrong type or order, such as “*starred_actors_reverse* \rightarrow *directed_by_reverse*” and “*directed_by_reverse* \rightarrow *written_by_reverse*”.

As shown in Eq.14, we train the similarity scorer using the stochastic gradient descent (SGD) backward propagation algorithm under the mean-squared-error (MSE) loss function. Then we endow our training criterion with the ℓ^2 (Euclidean distance) norm metric to avoid the model parameter distribution being highly warped. Regarding the predicted relatedness labels to lie in $[0, 1]$, for the positive sample we maximize the prediction value as close as possible to 1, and the negative sample as close as possible to 0. Finally, after sorting the scored candidates, we choose the entity which has the highest similarity score as the final answer.

4.3.2 Question Semantic Representation

As depicted in Fig.6 left, we use and finetune a standard version of **RoBERTa** to obtain the hidden state \mathcal{V}_q lying in start token [*CLS*] for our question encoders. Note that we reformat question q through replacing the topic entity’s

mention in the question with a token “NE”. And we respectively supplement two special characters ‘[CLS]’ and ‘[SEP]’ before and after the reformatted question, which follows the BERT default configuration. This aforesaid operations can help our model better distinguish the topic entity and other question words mentioned. Afterwards, we link the question’s topic entity mention to the KG node through matching with standard KG entity literal representations. We regard it as the semantic signal for answer reasoning and then adopt two full connect layers, neuron activation function ReLU [56] and a dropout layer to get better feature learning capability. After the above step, a vector representation is generated by the last full-connected layer, and we use the vector to compare space distance similarity with the output of another feature encoder, *Relational Chain Representation Module*.

4.3.3 Relational Chain Representation

As depicted in Fig.6 right, we provide the Relational Chain Representation module which consists of a single-layer bidirectional LSTM and a self-attention layer. Given the KG relational chain embedding sequence as input, this module learns and captures the relevant and necessary semantic information for answering reasoning. It has a similar structure to the *Question Semantic Parser* in the *Answer Filtering Module*, but the token embeddings used are not initialized randomly in this case. Instead, we apply the pretrained KG relation embedded representations existing in Sec 4.2.1 to embed our relational chain.

Formally, as depicted in Eq.13, for relational chain $\mathbf{r}^* = \{\mathbf{r}_j\}_{j=1}^M$, where M denotes the chain length, each relationship representation $\mathbf{r}_j (j = 1, 2, \dots, M)$ is initialized with the pretrained KG relational embeddings. Next, we feed the embedded vectors \mathbf{r}^* into a single-layer bidirectional LSTM network to obtain a series of output states h_1, h_2, \dots, h_M , where the j -th relation h_j denotes $[\vec{\mathbf{h}}_j; \overleftarrow{\mathbf{h}}_j]$, a combined vector of forward LSTM state output $\vec{\mathbf{h}}_j$ and backward LSTM state output $\overleftarrow{\mathbf{h}}_j$. Then the new question representation $q = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_M]$ can be transformed through a self-attention operation, which is similar to *Answer filtering Module* counterpart, shown in Eq.9, 10. Through this forward propagation, that is similar to *Question Semantic Parser*, we can obtain a question semantic vector \mathcal{V}_q which has the same dimension as relational chain vector \mathcal{V}_r for the following similarity computation step.

5 Experiment

In this section, we evaluate our proposed Rce-KGQA against competitive baselines on three benchmark datasets to investigate whether our model can outperform other methods on reasoning over the weakly supervised signal and incomplete knowledge graph. And we also append extensive ablation experiments and a case study to carefully verify and vividly demonstrate that the necessity, superiority, and meaning about our ideas in this work. The datasets as well as the pytorch implementation of our model are publicly available at <https://github.com/albert-jin/Rce-KGQA>.

5.1 Datasets and Evaluation Metric

Here, we first describe the three benchmark datasets we use in this work and then give a brief introduction of the metric *hit@1* we used for model evaluation.

MetaQA [57] is a large KGQA dataset which provides an original version *Vanilla* and two variations. In this paper, we use the original ones because they are designed manually. This dataset includes up to 75w QA pairs which are merely 2-hop questions. The questions’ literal descriptions are generated by cross-language translation, *English*→*French*→*English*. This relies on a large-scale movie domain which contains 9 relationship types, 43234 entities and up to 135k factoid triples. Here, we use the dataset edition which is generated from Apoorv et al. [22].

WebQuestionsSP-tiny [58] dataset is a relatively small dataset with a total of 4736 QA pairs. This QA dataset’s available KG is a subset of Freebase that contains all the facts within 2-hops of any entity mentioned in the questions of the original WebQuestionsSP, which have more than 188w entities and 1000 relationship types. Following [22], for all topic entities labeled in the original Freebase, He et al. [27] construct a subgraph containing other KG entities close to them by the PageRank-Nibble algorithm (PRN) [59]. In this way, theoretically, the pruned KG is likely to contain the corresponding answer entity for close to whole questions. In this work, we use the same train/dev/test splits as GraftNet [21].

Complex WebQuestionsSP (Complex-WebQSP)[60] is a more complex multi-hop reasoning dataset and its questions require up to a 4-hops relational path during answer reasoning. There are four types of question: conjunction(about 45%), composition(about 45%), comparative (about 5%), and superlative (about 5%). Our used edition of the Complex-WebQSP dataset is obtained from He et al. [27].

Table 1 Statistics for dataset MetaQA, WebQuestionsSP-tiny and Complex WebQuestionsSP. MetaQA contains three subsets of different complex question relational chain length, **1/2/3** hop MetaQA. WebQuestionsSP-tiny dataset requires up to 2-hop reasoning from knowledge base. Meanwhile, the Complex-WebQSP dataset requires up to 4-hops of reasoning on the KG.

Datasets	Train	Dev	Test
MetaQA 1-hop	208970	19947	9992
MetaQA 2-hop	231844	14872	14872
MetaQA 3-hop	227060	14274	14274
WebQSP-tiny	2848	250	1639
Complex-WebQSP	27639	3519	3531

Metric *hit@1* is a standard assessment for measuring the ratio in all validation samples that the highest scored entity belongs to the correct answers. In brief, if the QA system provides the user with a single entity and this entity

is right, we then determine that this prediction is correct. This evaluating indicator is popular and publicly recognized and has been used in many recent KGQA works [22, 26, 37].

5.2 Experiment Setting

As we know, hyperparameter choices have a significant impact on the model’s final performance [61–63]. Our optimal model hyperparameter configuration is summarized as follows. All the LSTM modules we used as feature encoders have a single layer with a hidden dimension of 256. All the dropout layers randomly drop about 30% of their features for inputs during the training step but they do not drop any features during testing. We apply Xavier initialization to each network layer’s training parameters in our model. Our applied pretrained transformer *RoBERTa* is a PyTorch-implemented configuration, which uses the BERT-base architecture [50, 64, 65], consisting of 12 layers, 768-d hidden size and 12 attention heads for efficient training and inference. Roberta-base encoder [25, 66] contains up to about 125M parameters. The Answer Filtering Module is trained for up to 200 epochs with a batch size of 128, and the relational chain reasoning module is trained for up to 120 epochs with a batch size of 32 on three benchmarks. For every 10 training epochs, we adopt the early-stopping strategy by evaluating *hit@1* on the test set to avoid overfitting. During model convergence, the stochastic gradient descent (SGD) optimizer with initial learning rate $lr = 1e-5$ is adopted. We used different random seeds to validate our best-configured model independently 5 times and report the average validation performances of our model in the next sections.

5.3 Compared Methods

In our experiment, the state-of-the-art methods for comparison are described as follows:

- **EmbedKGQA** [22] is a KG embedding driving method for multi-hop KGQA which matches the pretrained entity embeddings with question embeddings generated from the transformer.
- **SRN** [11] is an RL-based multi-hop question answering model which conducts the QA task by extending the inference chains on a KG.
- **KVMem** [31] The Key-Value Memory Network first attempts to conduct QA over incomplete KGs by augmenting it with text. It uses a memory table which stores the KG facts encoded into key-value pairs to retrieve a question-specific subgraph for reasoning.
- **GraftNet** [21] is a question description-based semantic sub-graph driving method that uses a variational graph CNN to perform QA tasks over question-specific subgraphs containing KG facts, entities and discourses from textual corpora.
- **PullNet** [10] improves GraftNet on the retrieval subgraph by introducing the graph retrieval module which utilizes shortest path from the topic entity to answer as the additional supervised signal.

- NSM_s [27] is a series of teacher-student learning approaches implemented as based on the Neural State Machine [40]. NSM , NSM_{+p} , and NSM_{+h} are three model variants which effectively employ the intermediate supervision signals. Specifically, NSM do not use the teacher network, NSM_{+p} use the teacher network with parallel reasoning, and NSM_{+h} use the teacher network with hybrid reasoning.

5.4 Main Results

In this section, we compare our model with the state-of-the-art baseline methods on three benchmarks, and the following questions are answered:

Q1. How effective and accurate is the performance of our model compared with other SOTA models?

Q2. Can our model really identify the implicit relations and the indirectly linked answer when there is no direct relational chain from topic entity to answer?

Q3. How many parameters does our *Rce-KGQA* contain, and does our model have a certain high execution efficiency?

5.4.1 Answer Reasoning on MetaQA

As illustrated in Table 2, the overall experimental results on the *MetaQA* test set clearly demonstrate that our proposed KGQA architecture significantly outperforms state-of-the-art methods on *hit@1* metric. Specifically, according to the performance of our model’s result on 1-hop MetaQA (identical to WikiMovies), compared to the other methods, Row 1 ~5, we observe that our method’s *hit@1* accuracy achieves much higher performance up to **98.3%**, which is an increase of **1.3%** compared to **GraftNet** [21] and **PullNet** [10], an increase of **0.8%** compared to **EmbedKGQA** [22] and an increase of **1.1%** compared to NSM_{+h} [27].

For the evaluation of multi-relation questions which require at least two hops of inference to find the answers, *hit@1* results on 2-hop and 3-hop MetaQA also show better performance than most competitive state-of-the-art baselines. Although our model test results on 2-hop did not achieve the best score, it is nevertheless comparable to other SOTA models. It is worth noting that the retrieval-and-reason process of **PullNet**, which can simultaneously extract answers from both corpora and KGs, is good at reasoning answers over large-scale KGs such as MetaQA. In contrast, we can see that our model *Rce-KGQA*, which takes relation chain reasoning into consideration, does not drop significantly in performance but remains almost unchanged when the relational chain hop increases. We think the reason may be that the baseline models only consider the question shadow semantic representations, and inevitably introduce noise and incorrect retrieving path over KG. On the other hand, since our model focuses on relational chain order and relation type, it is less sensitive to the hop size and shows robustness on complex multi-constraint queries over KG.

In summary, these results have shown their effectiveness and superiority when considering the question semantic and its relational chain into reasoning, which largely improves the KGQA performance.

Model	1-hop MetaQA	2-hop MetaQA	3-hop MetaQA
EmbedKGQA	<u>97.5</u>	98.8	94.8
SRN	97.0	95.1	75.2
KVMem	96.2	82.7	48.9
GraftNet	97.0	94.8	77.7
PullNet	97.0	99.9	91.4
NSM	97.1	99.9	98.9
NSM+p	97.3	99.9	98.9
NSM+h	97.2	99.9	98.9
Our Model	98.3	<u>99.7</u>	<u>97.9</u>

Table 2 Effectiveness comparisons on three subsets of MetaQA. The first group of results was taken from papers on recent methods. The values are reported using hits@1. The number in **bold** and underlined number denote the best and second-best methods, respectively. This figure corresponds to Sec. 5.4.1.

However, we also consistently find that our proposed Rce-KGQA’s two-stage pipeline mechanism could bring deviation cascade propagation between the coarse-grained answer filtering procedure [5] and the fine-grain answer selecting procedure [6]. As is shown in Table 2, and the comparison results from the columns of 2/3-hop MetaQA indicated. Our approach performs poorly when compared with other KGQA models such as **PullNet** and **NSMs**. They all consistently achieved very high performance with the *hit@1* metric of 99.9 percentage on 2-hop MetaQA and 98.9 percentage *hit@1* on the 3-hop MetaQA dataset. Correspondingly, our Rce-KGQA underperforms with about 0.2 percentage points behind on the 2-hop MetaQA and about 1.0 percentage points behind on the 3-hop MetaQA dataset.

In general, although our two-stage pipeline solution outperforms many baselines such as GraftNet and PullNet [10, 21], our designed Rce-KGQA’s separate architecture determines the fact that the quality of the final answer provided by *Relational Chain Reasoning Module* completely depends on the quality of the candidate entities provided by *Answer Filtering Module*. The feature of pipeline architecture like this is the principal reason that inevitably causes the cascaded error propagation, which would bring down the overall performance of the question answering service. We think it is crucial to enhance our model *Rce-KGQA* by integrating these two separated modules [5, 6] into one joint module, which we leave for future work.

In addition to the embedded graph vector’s parameter volumes, our two-stage modules respectively contain 46M parameters (mainly owned by *BiLSTM-Attn* block) and 197M parameters (mainly owned by *Roberta* encoder and *BiLSTM-Attn* block). Through our experiments, we observe that the total inference time fluctuates between 1.9 seconds and 2.7 seconds, in which the *Answer Filtering Module* contributes about 0.6s and the *Relational Chain Reasoning Module* contributes about 1.4s. It is conceivable that if the encoder were

switched to RoBERTa using the BERT-large architecture [25, 66], it would mean more heavy parameters, including up to 355M parameters and additional training optimization procedures and more inference time-consuming. Based on the experimental reproducibility and the consideration of the model lightweight, we chose the basic configuration of the *Roberta* encoder.

5.4.2 Experiments on WebQSP-tiny and Complex-WebQSP

WebQuestionsSP-tiny [22, 27, 36] is a relatively small dataset for training but relies on a large-scale KG (Freebase) whose entities' count is greater than 10 million. Table 3 presents the evaluation results on the *WebQuestionsSP-tiny* validation dataset, from which we can observe that our KGQA system still performs better than other state-of-the-art counterparts, **EmbedKGQA** (has 3.8% lower hits-at-one than our model) and **PullNet** (has 2.3% lower hits-at-one than our model).

Model	WebQuestionsSP-tiny	Complex-WebQSP
KVMem	46.7	21.1
GraftNet	66.4	32.8
EmbedKGQA	66.6	-
PullNet	68.1	45.9
NSM	68.7	47.6
NSM+p	<u>73.9</u>	<u>48.3</u>
NSM+h	74.3	48.8
Our Model	70.4	<u>48.3</u>

Table 3 Experiment results (% Hits@1) compared with SOTA methods on the WebQuestionsSP-tiny and Complex WebQuestionsSP validation datasets. All QA pairs in WebQuestionsSP-tiny are 2-hop relational questions. We copy the results for KV-Mem, GraftNet, EmbedKGQA, PullNet and NSM from [10, 21, 22, 27, 31], respectively. The best score is in **bold** and the second-best score is underlined.

Specifically, the last row shows that our full model achieves accuracy of up to 70.4% *hit@1*, which improves a large margin to other prior models. A possible explanation is that the filtering model equips the extra relational chain module with better reasoning perception, leveraging KG and question implicit features more efficiently, and emphasizing the order of relational triples selection to help our model make a correct decision. Even in large-scale KGs along with small training dataset situations like *WebQuestionsSP-tiny*, our *Rce-KGQA* solution can still be robust and helpful for handling realistic QA applications.

Complex-WebQSP [27] is a derivative KGQA dataset edition which is generated from *WebQuestionsSP-tiny* by extending the question entities or adding constraints to answers. As its name indicates, most of the questions this dataset included require up to 4-hops of relational chain reasoning from the topic entity to the corresponding answers.

The third column of Table 3 reports the *hit@1* metric statistics on the Complex-WebQSP benchmark. Our model outperforms competitively with

state-of-the-art KGQA baselines on such a complex multi-hop question answering scenario. More specifically, among most baselines (KVMem \sim NSM), our Rce-KGQA significantly surpasses other baselines, and respectively achieves 27.2%, 15.5%, 3.3%, and 0.7% absolute gains over these baselines (KVMem \sim NSM) in terms of the overall metric *hit@1*. This establishes the fact that our *Rce-KGQA* is better than previous approaches in terms of answering the questions with long-distance relational dependency. The **NSM_{+h}** achieves the best performance on the two adopted benchmarks. The **NSM_{+p}** and our *Rce-KGQA* both gain the second best performance on the *Complex-WebQSP* benchmark, proving both our model's competitive capability and the importance of the added teacher network. This is an important observation and advancement when it comes to handling such complex question answering tasks since our proposed approach is robust and efficient in dealing with these complex questions in multi-relational-hop answer reasoning situations.

5.5 Answer Reasoning for implicit relationship discovery

As shown in Table 4, we verify our method's ability to discover missing implicit relationships through comparison experiments. The KG which MetaQA uses has no missing link during the reasoning path because the QA question pairs are constructed upon this KG. However, to make it become a realistic setting, we simulate an incomplete KG by randomly removing half (with probability = 0.5) of the factoid triples from it. We call this pruned setting **half** and we call the full KG setting **full** in the text.

The experiments show our method's implicit relation discovering capability substantially outperforms other state-of-the-art methods over incomplete KGs. The amount of improvement is significant, with an increase of 43.7% compared to **KVMem** in *hit@1*. Furthermore, our competitive model also delivers an average 1.7% *hit@1* rate on 2-hop MetaQA half setting and performs well on 3-hop MetaQA half setting while **PullNet** still achieves the highest *hit@1* score.

Hence many baseline methods such as **GraftNet**, **PullNet** require constructed question-specific subgraphs, indicating they lack the capability to recall the answer nodes out of their generated subgraph and cannot perform well in real QA scenarios. Fortunately our model, which exploits the KG link prediction properties, does not limit its capability due to this constraint. Although those complex questions in *WebQuestionsSP-tiny* could be easily covered by hand-crafted rules, as many have been, our model is not suitable for such pre-defined rules. We think it is crucial to use more advanced reasoning capabilities to enhance our model *Rce-KGQA* correctly, a task which we leave for future work.

5.6 Answer Filtering Result Analysis

To further examine whether our proposed enhancement to the extra module *Relational Chain Reasoning Module* with advanced and obvious improvements,

Models	2-hop MetaQA		3-hop MetaQA	
	full	half	full	half
KVMem	82.7	48.4	48.9	37.6
GraftNet	94.8	69.5	77.7	66.4
PullNet	99.9	90.4	91.4	85.2
EmbedKGQA	98.8	<u>91.8</u>	<u>94.8</u>	70.3
Our Model	<u>99.7</u>	92.1	97.9	<u>84.7</u>

Table 4 Experimental results about reasoning on incomplete KG (*hit@1* as a percentage). We consider two different KG settings, **full** and **half**. **Full** denotes the complete KG and **half** denotes a KG subset whose 50% factoid triples are randomly removed.

we analyze the reasoning performance of first module *Answer Filtering Module* and show the answer distributions with prediction in Table 5. In Table 5, as we can see that, if we regard the scored candidate entities provided by this module as the final answer, the *hit@1* accuracy rate drops markedly compared to *hit@5* and *hit@10*. This observation indicates that the model which does not consider relational chain order and relation type achieves very poor performance and proves the necessity and superiority of our proposed module, the *Relational Chain Reasoning Module*.

Furthermore, from Table 5, we can clearly observe that almost right answers of our used datasets are collectively distributed in the top-5 of our model’s predictions. Due to the high recall rate of our first module, the *Answer Filtering Module*, we think we can only rely on a few top-scoring candidates (such as top-15, top-10, or even top-5) to further filter the final answer more accurately. So, during our model training and inference experiments, we tried several experimental schemes and considered the number of top-scoring candidates specifically, how to select candidates to automatically generate positive or negative samples and how to cut the top-N candidates for the sub-module *Relational Chain Reasoning Module* inference. The related experimental details are shown in Sec. 5.6.

Dataset	<i>Hit@1_{-r}</i>	<i>Hit@5_{-r}</i>	<i>Hit@10_{-r}</i>
2-hop MetaQA	0.861	0.995	0.999
3-hop MetaQA	0.858	0.984	0.997

Table 5 Our *Answer Filtering Module* answer reasoning performance on three different metrics *Hit@1, 5, 10*. Model performance on *Hit@5* and *Hit@10* accuracy highly outperform over *Hit@1* accuracy.

5.7 Candidates Filtering Strategy Analysis

Our QA solution *Rce-KGQA* is a complex pipeline system, and we inevitably must choose some crucial hyper-parameters to acquire an optimal model. These pre-defined parameters include full connected/LSTM layer number, dropout rate, learning rate, and so on. As illustrated in Sec. 5.6 and Table 5, different selection modes about top-scoring candidates during training and inference

have a huge impact on the final model performances. We now further investigate what influence would our model experience in different candidates selection mode.

Firstly, from Table 5 we can observe that our first step ‘answering filtering’ consistently achieves high recall performances on *hit@5*. Now we manually choose and cut-difference sorted candidate answers and conduct our comparison experiments. Specifically, from Sec. 4.3 we know that the dataset for *Relational Chain Reasoning Module* training is dynamically constructed by the last step. And the data construction follows the selection of top- N sorted scoring entities in which number N simultaneously decides the proportion of positive/negative samples and the intermediate results selection during the inference step. For example, in question Q , during model evaluation, we firstly obtain the intermediate scoring candidate answers $\{[A_i; S_i]\}_{i=1}^N$. In the next step, we should use the *Relational Chain Reasoning Module* to provide all filtered candidates with a more precise score to reach a final answer. Selecting how many top scoring entities to conduct the further step precisely becomes our focus research point.

We choose four strategies which include top- $\{5, 10, 15, 20\}$; then, from the corresponding experiments, we receive the following results which is shown in Table 6. From experimental comparison results we clearly find that in top-5 selection strategy, our model achieves the highest *hit@1* performances and achieves the second highest performances in the top-10 strategy, in both 2-hop MetaQA and 3-hop MetaQA datasets. These phenomena can come from two aspects. First, the recall rate of the module *Relational Chain Reasoning Module* is good enough for the next fine-grained screening and the positive/negative samples proportion should be in a suitable extent. In addition, more candidates in the model inference step could introduce more noise entities which could affect model answer judgment and decrease the overall model performances.

Policy	top-5 pick	top-10 pick	top-15 pick	top-20 pick
2-hop MetaQA	0.997	<u>0.994</u>	0.989	0.985
3-hop MetaQA	0.979	<u>0.971</u>	0.967	0.967

Table 6 Our model final performance statistics about impacts of the four selection strategies: choose top-5, choose top-10, choose top-15 and choose top-20. The values reported are *hit@1*. Bold and underlined fonts denote the best and the second-best selection strategies.

5.8 Ablation Study

To better understand and gain a deep insight into our model design, we also perform ablation experiments to investigate systematically the impact and contributions of different components. **RceKGQA_{-r}**, **RceKGQA_{-a}** and **RceKGQA_{-b}** are variants of our full model, RceKGQA. Note that, for our ablated experiments, we remove one component each time. Here we briefly introduce these variants for the ablated experiments.

- **RceKGQA_{-r}** removes the *Relational Chain Reasoning Module* and the highest scoring entity is provided by the *Answer Filtering Module* as the final answer.
- **RceKGQA_{-a}** removes all the self-attention operations from the model.
- **RceKGQA_{-b}** replaces RoBERTa with LSTM in the question semantic representation part of the *Relational Chain Reasoning Module*.
- **RceKGQA** is the full model introduced in this paper.

In this section, the following questions are answered:

Q1. How much does the *Relational Chain Reasoning Module* help our model’s reasoning accuracy?

Q2. Can the attention mechanism really help increase our model’s overall performance?

Q3. Is the effectiveness of our method due to the use of **RoBERTa** in the *Relational Chain Reasoning Module*?

According to the result comparison between RceKGQA and its variant RceKGQA_{-r}, we can clearly conclude that removing the *Relational Chain Reasoning Module* from the proposed model has a huge impact on the results.

The performance gap between RceKGQA_{-r} which is shown in Row **2** and our full model as shown in Row **1** indicates that the semantic relational chain factor plays a pivotal role in answer reasoning, which incorporates relational chain order and relationship type to provide more accurate answers for the terminal user.

As shown in Row **3**, when our full model is compared with RceKGQA_{-a}, we can see an average 3.5% performance drop across the *hit@1* metric if the self-attention components are removed, demonstrating the importance of the self-attention mechanism used in our method, as it effectively helps with the final answer prediction. A possible reason is that the self-attention mechanism can distinguish the most relevant and interesting signals from noise information, and help our model better understand the question and relational chain semantic.

RceKGQA_{-b}, which is shown in Row **4**, only loses 1.6% *hit@1* accuracy compared with our full model, which replaces RoBERTa with BiLSTM, and demonstrates that using the transformer is not a major factor in increasing the overall model performance.

The above-ablated statistics confirm that all three components introduced for handling the multi-hop relation question answering contribute to the overall model performance.

5.9 Case Study

The major novelty of our approach lies in the introduced *relational chain reasoning* network. Here, we present a case study to demonstrate its contribution to improving the overall model architecture.

As shown in Fig. 7, given the question, “*In which years were movies released which starred actors who appeared in the movie [Thunderbolt]?*”, the

Model	1-hop MetaQA	2-hop MetaQA	3-hop MetaQA
RceKGQA	98.3	99.7	97.9
RceKGQA _{-r}	85.8	86.1	84.8
RceKGQA _{-a}	<u>96.1</u>	95.9	93.4
RceKGQA _{-b}	95.9	<u>98.2</u>	<u>95.6</u>

Table 7 Ablation study statistical results for RceKGQA and its three variants (*Hit@1* by percentage). Compared with our full model, suffix _{-r} denotes RceKGQA whose *Relational Chain Reasoning Module* is removed, suffix _{-a} denotes the RceKGQA variant whose attention operation is dropped and suffix _{-b} denotes the RceKGQA variant whose question encoder RoBERTa is replaced with BiLSTM.

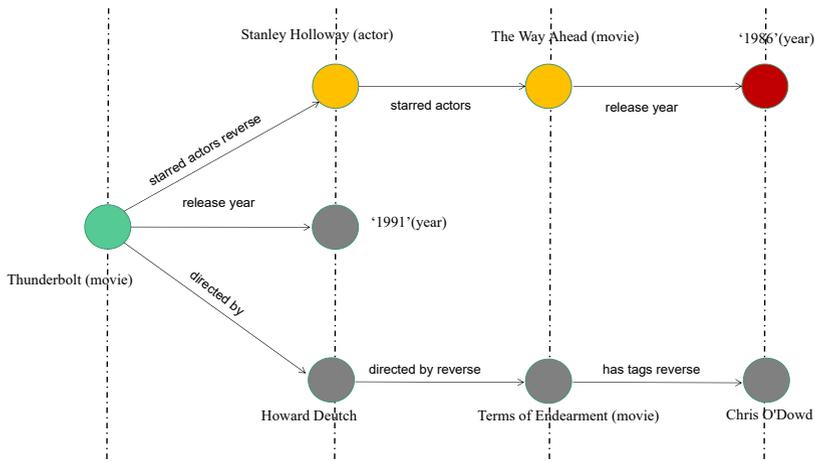


Fig. 7 Case analysis from the 3-hop MetaQA dataset. We use green, red, yellow and grey circles to denote the topic KG nodes, correct answer, intermediate nodes and irrelevant nodes, respectively. The orange and red coloured circles denote the actual reasoning intermediate nodes and answer nodes. The colour darkness indicates the relevant degree of an entity by a method.

right reasoning relational chain preserved in KG is **Thunderbolt(movie)-starred_actors_reverse** → **Stanley Holloway(actor)-starred_actors** → **The Way Ahead(movie)-release_year** → **'1986'(year)**. When ignoring the relational chain feature factor and only using the high-scored entity generated by the *Answer Filtering Module* as the final answer, the network mistakenly selects a wrong reasoning path **Thunderbolt(movie)-release_year** → **'1991'(year)** for the aforesaid question with a very high probability of **0.96** as the answer. Its attention only focuses on the relationship: *release_year*, which ignores the repeated relations of *starred_actors* and *starred_actors_reverse*. In comparison, the complete model which considers the relational chain factor and utilizes the *relational chain reasoning* network for fine-grained selection can easily and correctly provide the right answer **'1986'(year)** with a high probability of **0.99** from KG.

This example shows that our *relational chain reasoning* network indeed provides very useful supervision signals of relational chain recognition at intermediate steps to improve our model’s overall QA performance.

6 Conclusion and Future Perspective

In this work, we introduce an elaborate KG embedding-based pipeline approach for the multi-hop KGQA task, termed Relational Chain-based Embedded KGQA. Novel techniques are proposed to effectively utilize QA relational chain parsing to identify the semantics more accurately and leverage the structure information preserved in KG embedding to reason the implicit answer indirectly. Our comprehensive empirical results on three benchmarks demonstrate that our method outperforms many of its state-of-the-art counterparts. The experimental comparison between our approach and its ablated variants also verifies that the proposed model components contribute to the answer reasoning result. We believe KGQA will continue to be an attractive and promising research direction with realistic industrial and domestic scenarios, such as Intelligent Recommendation, Smart Personal Assistant, Big Data Mining Services, and Automatic Customer Services.

In the future, we plan to study the following major problems: (i) To support real-world dynamic application scenarios, the KGQA application is always updated quickly and inevitably accumulates new and immense external knowledge in real time. How can we augment our available KG’s knowledge reserve automatically and incrementally to expand our system’s knowledge coverage? (ii) This model is trained on relatively small QA datasets under weak supervision without external prior knowledge. How can we introduce external knowledge such as knowledge from web pages and other open-domain KGs to improve our question answering system’s performance?

Following the universal solution patterns of the KGQA task, the method presented in this paper assumes that “Our model will always choose an optimal answer from the KG”. Therefore, the method based on this assumption is definitely not suitable for the case where the answer does not exist in the KG. In future work, we will add research on the sub-task of “Detecting whether the answer exists in the KG”. Moreover, our proposed Rce-KGQA is essentially a pipeline mechanism, which could bring deviation propagation and poor performance. In future work, we will also enhance our *Rce-KGQA* by integrating the separated *Answer Filtering Module* and *Relational Chain Reasoning Module* together. Concretely, the major factor, which hinders the joint modelling, is the multiple step-by-step answer retrieval due to the KG’s traditional structured storage pattern. Inspired by Fabio et al. [67] who prove that the huge-volume PLMs have surprising knowledge storing capabilities, we will try to infuse the KG knowledge into the Transformer-based PLMs, which can theoretically solve the end-to-end fashion modelling difficulty well from the root.

Acknowledgement

This work was partially supported by the Shanghai Yangfan Program (Project Code: 22YF1413600), the Major Research Plan of National Natural Science Foundation of China (Project Code: 92167102), and the Shaanxi Province Key Industrial Chain Projects (Project Code: NO.2018ZDCXL-GY-04-03-02). The authors would like to thank Guizhong Liu and Ruiping Yin for providing helpful discussions and comments.

References

- [1] Xiao, J., Kalia, A.K., Vukovic, M.: Juno: An intelligent chat service for its service automation. In: *Service-Oriented Computing – ICSOC 2018 Workshops*, pp. 486–490. Springer, Cham (2019)
- [2] Bin, F., Yunqi, Q., Chengguang, T., Yang, L., Haiyang, Y., Jian, S.: A survey on complex question answering over knowledge base: Recent advances and challenges. *CoRR* (2020) [arXiv:2007.13069](https://arxiv.org/abs/2007.13069)
- [3] Hao, Y., Zhang, Y., Liu, K., Zhao, J.: An end-to-end model for question answering over knowledge base with cross-attention combining global knowledge, pp. 221–231. *Association for Computational Linguistics*, (2017)
- [4] Michael, P., Luke, Z.: Simplequestions nearly solved: A new upperbound and baseline approach. *CoRR* (2018)
- [5] Yunshi, L., Gaole, H., Jinhao, J., Jing, J., Wayne, X., Jirong, W.: A survey on complex knowledge base question answering: Methods, challenges and solutions. *CoRR* (2021) [arXiv:2105.11644](https://arxiv.org/abs/2105.11644)
- [6] Lan, Y., Wang, S., Jiang, J.: Knowledge base question answering with a matching-aggregation model and question-specific contextual relations. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* **27**(10), 1629–1638 (2019)
- [7] Yunshi, L., Shuohang, W., Jing, J.: Knowledge base question answering with topic units. In: *IJCAI* (2019)
- [8] Bast, H., Haussmann, E.: More accurate question answering on free-base. In: *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, pp. 1431–1440 (2015)
- [9] Abujabal, A., Yahya, M., Riedewald, M.: Automated template generation for question answering over knowledge graphs. In: *Proceedings of the 26th International Conference on World Wide Web*, pp. 1191–1200. *International Conference on World Wide Web*, (2017)

- [10] Haitian, S., Tania, B.-W., William, W.C.: Pullnet: Open domain question answering with iterative retrieval on knowledge bases and text, pp. 474–482. EMNLP, (2019)
- [11] Yunqi, Q., Yuanzhuo, Wang, X.J., Kun, Z.: Stepwise reasoning for multi-relation question answering over knowledge graph with weak supervision, pp. 474–482. WSDM, (2020)
- [12] Dong, L., Wei, F., Zhou, M., Xu, K.: Question answering over Freebase with multi-column convolutional neural networks, pp. 260–269. Association for Computational Linguistics, Beijing, China (2015)
- [13] Yu, H., Lu, J., Zhang, G.: An online robust support vector regression for data streams. *IEEE Transactions on Knowledge and Data Engineering*, 1–1 (2020). <https://doi.org/10.1109/TKDE.2020.2979967>
- [14] Boris, G.: Question-answering system for teaching autistic children to reason about mental states. Technical report (2000)
- [15] Gao, S., Chen, X., Ren, Z., Zhao, D., Yan, R.: Meaningful answer generation of e-commerce question-answering. *ACM Trans. Inf. Syst.* **39**(2) (2021). <https://doi.org/10.1145/3432689>
- [16] Yu, H., Lu, J., Zhang, G.: Continuous support vector regression for nonstationary streaming data. *IEEE Transactions on Cybernetics*, 1–14 (2020). <https://doi.org/10.1109/TCYB.2020.3015266>
- [17] Xia, N., Yu, H., Wang, Y., Xuan, J., Luo, X.: Dafs: a domain aware few shot generative model for event detection. *Machine Learning* **11**(12) (2022). <https://doi.org/10.1007/s10994-022-06198-5>
- [18] Huang, H., Wei, X., Nie, L., Mao, X., Xu, X.-S.: From question to text: Question-oriented feature attention for answer selection. *ACM Trans. Inf. Syst.* **37**(1) (2018). <https://doi.org/10.1145/3233771>
- [19] Molino, P., Aiello, L.M., Lops, P.: Social question answering: Textual, user, and network features for best answer prediction. *ACM Trans. Inf. Syst.* **35**(1) (2016). <https://doi.org/10.1145/2948063>
- [20] Deepak, N., Jatin, C., Charu, S., Manohar, K.: Learning attention-based embeddings for relation prediction in knowledge graphs. In: Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, pp. 4710–4723 (2019)
- [21] Haitian, S., Bhuwan, D., Manzil, Z., Kathryn, M., Ruslan, S., Cohen, W.W.: Open domain question answering using early fusion of knowledge bases and text (2018)

- [22] Saxena, A., Tripathi, A., Talukdar, P.: Improving multi-hop question answering over knowledge graphs using knowledge base embeddings, pp. 4498–4507 (2020)
- [23] Mueller, J., Thyagarajan, A.: Siamese recurrent architectures for learning sentence similarity. In: AAAI, pp. 2786–2792. Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, (2016)
- [24] Gao, J., Yu, H., Zhang, S.: Joint event causality extraction using dual-channel enhanced neural network. *Knowledge-Based Systems* **258**, 109935 (2022). <https://doi.org/10.1016/j.knosys.2022.109935>
- [25] Yinhan, L., Myle, O., Naman, G., Jingfei, D.: A robustly optimized bert pretraining approach: Roberta. PMLR, (2019)
- [26] LAN, Y., Jing, J.: Query graph generation for answering multi-hop complex questions from knowledge bases (2020)
- [27] He, G., Lan, Y., Jiang, J., Zhao, W.X., Wen, J.-R.: Improving multi-hop knowledge base question answering by learning intermediate supervision signals. In: Proceedings of the 14th ACM International Conference on Web Search and Data Mining. WSDM '21, pp. 553–561. Association for Computing Machinery, New York, NY, USA (2021). <https://doi.org/10.1145/3437963.3441753>
- [28] Xiao, H., Jingyuan, Z., Dingcheng, L., Ping, L.: Knowledge graph embedding based question answering, pp. 105–113. Proceedings of the 13th ACM International Conference on Web Search and Data Mining, (2019)
- [29] Chen, Y., Subburathinam, A., Chen, C.-H., Zaki, M.J.: Personalized food recommendation as constrained question answering over a large-scale food knowledge graph. Proceedings of the 14th ACM International Conference on Web Search and Data Mining (2021)
- [30] Mo, Y., Wenpeng, Y., Kazi, S.H., Bowen, Z.: Improved neural relation detection for knowledge base question answering. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, pp. 571–581. Association for Computational Linguistics, (2017)
- [31] Alexander, H.M., Adam, F., Jesse, D., Amir-Hossein, K.: Key-value memory networks for directly reading documents, pp. 249–256. EMNLP, (2016)
- [32] Kun, X., Yuxuan, L., Yansong, F., Zhiguo, W.: Enhancing key-value memory neural networks for knowledge based question answering. In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics, pp. 2937–2947. Association for

- Computational Linguistics, (2019)
- [33] Bill, Yuchen, L., Xinyue, C., Jamin, C., Xiang, R.: Kagnet: Knowledge-aware graph networks for commonsense reasoning, pp. 2829–2839. Association for Computational Linguistics, (2019)
 - [34] Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. In: 5th International Conference on Learning Representations, pp. 486–490. ICLR, April 24-26 (2017)
 - [35] Wang, R., Rossetto, L., Cochez, M., Bernstein, A.: QAGCN: A graph convolutional network-based multi-relation question answering system. arXiv (2022). <https://doi.org/10.48550/ARXIV.2206.01818>
 - [36] Das, R., Godbole, A., Naik, A., Tower, E., Zaheer, M., Hajishirzi, H., Jia, R., Mccallum, A.: Knowledge base question answering by case-based reasoning over subgraphs. In: Chaudhuri, K., Jegelka, S., Song, L., Szepesvari, C., Niu, G., Sabato, S. (eds.) Proceedings of the 39th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 162, pp. 4777–4793. PMLR, (2022)
 - [37] Zi-Yuan, C., Chih-Hung, C., Lun-Wei, K.: Uhop: An unrestricted-hop relation extraction framework for knowledge-based question answering. In: NAACL (2019)
 - [38] Jain, S.: Question answering over knowledge-base using factual memory networks. In: NAACL (2016)
 - [39] Yao, X., Van, D.: Information extraction over structured data: question answering with freebase (2014)
 - [40] Drew, A.H., Christopher, D.M.: Learning by abstraction: The neural state machine. In: NeurIPS, pp. 5901–5914 (2019)
 - [41] A.Bordes, N.Usunier, A.Garcla-Duran, J.Weston, O.Yakhnenko: Translating embeddings for modeling multi-relational data. In: Proc. Adv. Neural Inf. Process, pp. 2787–2795 (2013)
 - [42] Wang, A., Zhang, J., Feng, J., Chen, Z.: Knowledge graph embedding by translating on hyperplanes. In: 28th AAAI, pp. 1112–1119 (2014)
 - [43] Min-Chul, Y., Do-Gil, L., HaeChang, R.: Knowledge-based question answering using the semantic embedding space. In: Expert Systems with Applications, pp. 9086–9104 (2015)
 - [44] Trouillon, T., Welbl, J., Riedel, S., Gaussier, E., Bouchard, G.: Complex embeddings for simple link prediction. In: International Conference on

Machine Learning, pp. 2071–2080 (2016)

- [45] Sun, Z., Wang, C., Hu, W., Chen, M., Dai, J., Zhang, W., Qu, Y.: Knowledge graph alignment network with gated multi-hop neighborhood aggregation **34**, 222–229 (2020)
- [46] Tingting, J., Hao, W., Xiangfeng, L., Xie, S., Jingchao, W.: Mifas: Multi-source heterogeneous information fusion with adaptive importance sampling for link prediction. (2021). <https://doi.org/10.1111/exsy.12888>
- [47] Afzal, A., Sading, M., Hussain, M., Ali, M., Lee, S., Khattak, A.: Knowledge-based reasoning and recommendation framework for intelligent decision making. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, pp. 571–581. Expert Systems, (2018). <https://doi.org/10.1111/exsy.12242>
- [48] Chopra, S., LeCun, Y.: Learning a similarity metric discriminatively with application to face verification. In: IEEE Computer Society Conference Computer Vision and Pattern Recognition, pp. 539–546 (2005)
- [49] Jeffrey, P., Richard, S., Christopher, M.: Global vectors for word representation. In: In EMNLP, pp. 1532–1543 (2014)
- [50] Ashish, V., Noam, S., Niki, P., Jakob, U., Llion, J., Gomez, A.N., Lukasz, K., Illia, P.: Attention Is All You Need (2017)
- [51] Wang, Q., Mao, Z., Wang, B., Guo, L.: Knowledge graph embedding: A survey of approaches and applications. IEEE Transactions on Knowledge and Data Engineering **29**(12), 2724–2743 (2017)
- [52] Nickel, M., Rosasco, L., Poggio, T.: Holographic embeddings of knowledge graphs. In: in Proc. 30th AAAI Conf 2016, pp. 1955–1961 (2016)
- [53] M, N., Tresp, V., Kriegel, H.-P.: A three-way model for collective learning on multi-relational data. In: in Proc. 28th Int Conf, pp. 809–816 (2011)
- [54] Glorot, Bengio, Y.: Understanding the difficulty of training deep feed-forward neural networks. In: Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, vol. 9, pp. 249–256. PMLR, (2010)
- [55] He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification (2015)
- [56] Xavier, G., Antoine, B., Yoshua, B.: Deep sparse rectifier neural networks. In: Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, vol. 15, pp. 315–323 (2011)

- [57] Yuyu, Z., Hanjun, D., Zornitsa, Kozareva, Alexander, J. S., Le, S.: Variational reasoning for question answering with knowledge graph. In: In Thirty-Second AAAI Conference on Artificial Intelligence, pp. 2787–2795 (2018)
- [58] Wentau, Y., Matthew, R., Christopher, M., Ming-Wei, C., Jina, S.: The value of semantic parse labeling for knowledge base question answering. In: In ACL, pp. 2787–2795 (2016)
- [59] Reid, A., Fan, R.K.C., Kevin, J.L.: Local graph partitioning using pagerank vectors. In: FOCS (2006)
- [60] Alon, T., Jonathan, B.: The web as a knowledge-base for answering complex questions. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, pp. 641–651. In NAACL-HLT, (2018)
- [61] Ding, H., Huang, S., Jin, W., Shan, Y., Yu, H.: A novel cascade model for end-to-end aspect-based social comment sentiment analysis. *Electronics* **11**(12) (2022). <https://doi.org/10.3390/electronics11121810>
- [62] Jin, W., Yu, H., Luo, X.: Cvt-assd: Convolutional vision-transformer based attentive single shot multibox detector. In: 2021 IEEE 33rd International Conference on Tools with Artificial Intelligence (ICTAI), pp. 736–744 (2021). <https://doi.org/10.1109/ICTAI52525.2021.00117>
- [63] Zhao, Z., Yu, H., Luo, X., Gao, J., Xu, X., Shengming, G.: Ia-icgc: Integrating prior knowledge via intra-event association and inter-event causality for chinese causal event extraction. In: Pimenidis, E., Angelov, P., Jayne, C., Papaleonidas, A., Aydin, M. (eds.) *Artificial Neural Networks and Machine Learning – ICANN 2022*, pp. 519–531. Springer, Cham (2022)
- [64] Tao, Q., Luo, X., Wang, H., Xu, R.: Enhancing relation extraction using syntactic indicators and sentential contexts. In: 2019 IEEE 31st International Conference on Tools with Artificial Intelligence (ICTAI), pp. 1574–1580 (2019). <https://doi.org/10.1109/ICTAI.2019.00227>
- [65] Gu, H., Yu, H., Luo, X.: Dbgare: Across-within dual bipartite graph attention for enhancing distantly supervised relation extraction. In: Memmi, G., Yang, B., Kong, L., Zhang, T., Qiu, M. (eds.) *Knowledge Science, Engineering and Management*, pp. 400–412. Springer, Cham (2022)
- [66] Ott, M., Edunov, S., Baevski, A., Fan, A., Gross, S., Ng, N., Grangier, D., Auli, M.: fairseq: A fast, extensible toolkit for sequence modeling. In: *Proceedings of NAACL-HLT 2019: Demonstrations* (2019)

- [67] Petroni, F., Rocktäschel, T., Riedel, S., Lewis, P., Bakhtin, A., Wu, Y., Miller, A.: Language models as knowledge bases? In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 2463–2473. Association for Computational Linguistics, Hong Kong, China (2019). <https://doi.org/10.18653/v1/D19-1250>. <https://aclanthology.org/D19-1250>