

# Using Context Information to Enhance Simple Question Answering

Lin Li · Mengjing Zhang · Zhaohui Chao ·  
Jianwen Xiang

Received: date / Accepted: date

**Abstract** With the rapid development of knowledge bases (KBs), question answering (QA) based on KBs has become a hot research issue. In this paper, we propose two frameworks (i.e., a pipeline framework, an end-to-end framework) to focus on answering single-relation factoid questions. In both of two frameworks, we study the effect of context information on the quality of QA, such as the entity's notable type, out-degree. In the pipeline framework, it includes two cascaded steps: entity detection and relation detection. In the end-to-end framework, we combine char-level encoding and self-attention mechanisms, using weight sharing and multi-task strategies to enhance the accuracy of QA. Experimental results show that context information can get better results of simple QA whether it is the pipeline framework or the end-to-end framework. In addition, we find that the end-to-end framework achieves results competitive with state-of-the-art approaches in terms of accuracy and take much shorter time than them.

**Keywords** Question answering · Knowledge base · Context information · Self-attention mechanisms

---

Lin Li · Mengjing Zhang · Zhaohui Chao · Jianwen Xiang  
School of Computer Science and Technology, Wuhan University of Technology, Wuhan, China

Lin Li  
E-mail: cathylin@whut.edu.cn  
· Mengjing Zhang  
E-mail: zhangmengjing@whut.edu.cn  
· Zhaohui Chao  
E-mail: chaozhaohui@whut.edu.cn  
· Jianwen Xiang  
E-mail: xiangjw@gmail.com

## 1 Introduction

QA is a classic natural language processing (NLP) task, which aims at building systems that automatically answer questions formulated in natural language [1].

In recent years, several large-scale general purpose KBs have been constructed, including Freebase [2,3], its main data source is Wikipedia, and some data comes from IMDB and other websites [4], YAGO [5], DBpedia [6] and Wikidata [7], Babelnet [8]. Also, the open Chinese KBs include Zhishi.me [9], CN-DBpedia [10], Xlore [11], etc. In addition, there are some commercial KBs that are not completely open, such as Google Knowledge Graph [4, 12, 13] and the Facebook Graph. We can access the data of entities and relations through a specific API. Since Google put forward the concept of knowledge graph, the related research of KBs has reached a new level of popularity.

The father of the World Wide Web Berners-Lee [14] proposed the concept of semantic web in 1998. Different from traditional Web networks, semantic and structured descriptions are introduced in the Semantic Network, which realizes semantic-based associations between data. And it can promote the development of computers in the direction of semantic-based information exchange [15]. RDF (Resource Description Framework) is a simple and flexible data model used to describe the semantic network proposed by the World Wide Web Consortium (W3C) [16]. Most of KBs store information in the form of RDF triples (subject, predicate, object) [16, 17, 18]. For example, (*/m/02mjmr*, */people/person.place\_of\_birth*, */m/02hrh0\_*), where */m/02mjmr* is the Freebase id for Barack Obama, and */m/02hrh0\_* is the id for Honolulu. By structuring knowledge stored in this basic form, we can better organize, manage, and utilize vast amounts of knowledge. But people cannot directly understand and extract the knowledge in the knowledge graph without struct query language. Therefore by mapping the questions to the triples in the KBs, we can get the correct answer. This is a good way to use the KBs. With the development of KBs, question answering over knowledge base (KB-QA) [19] has attracted more and more attention.

KB-QA is defined as the task of retrieving the correct entity or a set of entities from a KB given a query expressed as a question in natural language [1]. For instance, in order to answer the question “*where was former U.S. President Obama born?*” we need to retrieve the entity */m/02mjmr* in Freebase to represent former U.S. president Barack Obama and the relation */people/person.place\_of\_birth* corresponds to the relation of this question. With the entity and the relation, we can form a corresponding structured query. As a result, we can obtain the right answer *Honolulu* (id: */m/02hrh0\_*). Typically, with SPARQL (an RDF query language) people can extract information expressed in natural language from KBs. Just like cross-modal retrieval [20].

The KB-QA technology can be divided into two technical routes: (1) symbol-based representations, such as traditional semantic parsing, and (2) distribution-based embedding. The former transforms the semantics of the question into a structured query, and then returns the answer by the structured query, such as  $\lambda$  paradigm [21] and DCS-Tree [22]. The corresponding semantic parsing methods are: combined category grammar (CCG) [23] and dependent combination grammar (DCS) [24]. The

method is efficient. However, in this process, experts are required to formulate dictionaries or grammar rules etc. Not only does it need a large amount of human resources, it is also difficult to migrate [25]. In the latter route, the candidate answers are first determined in the KB based on the question. The question and the candidate answers are mapped to a low-dimensional space, therefore their distributed representation (i.e., Distributed Embedding) is obtained, which is trained by the training data, therefore the question vector and its corresponding correct answer vector are in a low-dimensional space. The matching score is as high as possible. When the training of the model is completed, it can be screened according to the score of the vector expression of the candidate answer and the question, and the highest score is found as the final answer. This route can be divided into a pipeline framework and an end-to-end framework. It has strong operability and does not require any manual features. Therefore its improvement space is relatively large.

With the emergence of deep learning, the development of NLP has greatly promoted. The effect of KB-QA can be improved by combining deep learning with the above two technical routes respectively. In this paper, we mainly discuss the impact of the second route (i.e., Distribution Embedding) combined with deep learning.

Unlike traditional methods, deep learning can capture the semantic information of the text at multiple levels through learning, including words, phrases, sentences, paragraphs and chapters. Therefore the semantic gap in traditional NLP is improved or solved to some extent. Questions expressed in natural language can be answered directly. The core idea of deep learning system in KB-QA tasks is representation and matching. First, we should learn representation of both the question and the fact of KBs, which contains literature level and semantic level. Then we need to calculate the correlation between question and fact. The work of answering single-relation factoid questions was first proposed by Bordes et al. [26]. In this work, they employ Memory Networks and introduce a new dataset SimpleQuestions, which is built on Freebase. For each triple (entity1, relation, entity2), it shows the relationship between entity1 and entity2 [27]. The relation consists of entity's type and properties, defined by the Freebase ontology [17]. There are currently two main frameworks for solving KB-QA, pipeline frameworks( [28, 29, 30, 31]) and end-to-end frameworks( [1, 26, 32, 33, 34, 35, 36, 37, 38]) respectively. As in the latter case, the data processing procedure is more concise, so it is more popular than the former.

In this paper, we focus on answering single-relation factoid questions (i.e., simple questions), which can be answered by a single fact of the KBs. The task is also called simple QA. Thereby we explore two different frameworks to answer such questions, analyzing the impact of entity's context information to answer selection. We find that combining the entity's out-degree and notable type can improve the accuracy of answers. Moreover, our end-to-end framework that combines context information achieves a competitive result with state-of-the-art work of Lukovnikov et al. [34], but run much faster than it.

In our pipeline framework and end-to-end framework, there are different ways to integrate entities' type information and out-degree information.

In our pipeline framework, For the combination of out-degree information, after the entity candidate set and the relation with the highest matching degree are obtained, in order to solve the case where multiple entities in the entity candidate set

have the same relation, these entities are reordered by using the out-degree information to generate a final result. For the incorporation of type information of the entity, we calculate the matching score of the question-type pair to mine the entity type information in the question, and sort the entity candidate set and the relation candidate set to achieve final screening of the results. In the following we will show it in details.

In our end-to-end framework, type information is incorporated into two different methods. The first is to concatenate the entity type information with the entity tag and inputting it as a new tag of the entity into the multitasking end-to-end model (QA-T) to sort the answers. The second is to add a matching task between the entity type information and the question on the basis of QA-T, the entity type information is added to the original input layer, through the model we can obtain the semantic matching score between the type information and the question, then we sort the answers based on matching scores to get the final answer. In addition, for the out-degree information of the entity, it is used to sort the answers with the highest matching score, so as to further sort and filter the answers. The detailed method will be described later.

The experimental results show that the accuracy of both frameworks is improved after combining the context information. What's more, in our end-to-end framework, combining type information, using char-level embedding and self-attention mechanisms, model QA-T can get better results. As is known, the work done by Lukovnikov et al [34] using end-to-end framework gets 71.2% accuracy, and takes 48 hours. However, our work achieves 71.8% but only takes 4 hours. The result is a little better, and the time spent is greatly shortened.

The rest paper is structured as follows: We introduce our related work in Section 2. In Section 3 we present our two different frameworks in details. The experiment is reported in Section 4, and then we conclude in Section 5 with some discussions for future work.

## 2 Related Work

With the development of deep learning technology, there are two mainstream technical routes for the KB-QA task: pipeline frameworks and end-to-end frameworks.

### 2.1 Pipeline Frameworks

In pipeline frameworks, the KB-QA task is decomposed into two subtasks: entity detection and relation detection. Dai et al. [28] use the bidirectional cyclic neural network (i.e., Bi-GRU-CRF) combined with the conditional random field (i.e., CRF) to identify the questions, and cut the entity candidate sets according to the entity identification results. They also use a Bi-GRU network to encode relations in relation matching. Yin et al. [29] divide the task into two parts: entity linking and fact selection. they use the convolutional neural network based on character level (char-level) and word level (word-level) to perform entity matching and relation matching respectively, and they use the pooling method combining attention mechanism in the

relation matching process. Cui et al. [31] design a new kind of question representation: templates, and use them to improve the semantic representation of questions and better judge the type and intent of questions. Moreover, they increase the coverage of the KB by 57 times through expanding predicates in the RDF KB. Hao et al. [30] use Bi-LSTM-CRF to identify the entities in questions, then the question templates are used to correct the entity recognition results, and the relation detection is combined to improve fact selection. Multi-granular coding and multi-dimensional information are utilized in this framework.

## 2.2 End-to-end Frameworks

In end-to-end frameworks, Bordes et al. [1] present the question and fact triples in the KBs as numerical vectors of the same semantic space by embedding model, and the similarity of their vectors is measured to sort the candidate triples. Yih et al. [38] build a semantic matching model using convolutional neural networks, and propose a QA framework for single-relation questions based on semantic similarity. By measuring the entity text and entity labels, relation templates and relations in the question respectively, the similarity between the two completes the sorting of the candidate answers. Golub and He [33] employ a char-level, attention-based encoder-decoder framework for QA. The model is robust for unseen entities since it adopts char-level modeling. Hao et al. [32] present an end-to-end neural network model to represent the questions, which improves the representation of questions via cross-attention mechanisms. Lukovnikov et al. [34] present an end-to-end neural network. They merge word-level and char-level representation of questions. Wu et al. [35,36,37] employ attention mechanisms and joint learning. And they design an end-to-end network structure.

Although the above two frameworks have achieved good results on QA, they have not used context information for in-depth analysis to improve the accuracy of QA. Therefore, in this paper, we mainly discuss the impact of context information on the accuracy of QA.

## 2.3 Attention Mechanism

The attention mechanism is derived from the visual field. When the human eyes observe an image, it is preferred to quickly scan the entire image first and find the areas in the image that need to be focused, then they focus on these areas, carefully observe and analyze, and get more detailed information. The attention mechanism can be divided into Hard Attention, Soft Attention, Local Attention and Global Attention. Xu et al. [39] divide attention into Soft Attention and Hard Attention in dealing with image description generation tasks. The former refers to assigning weights to all regions of the original image. Part of the original image is assigned a corresponding weight. The latter can reduce the amount of calculations compared to the former. Relative to Global Attention, Luong et al. [40] propose Local Attention. By setting a context window, only a small part of the source language can be generated when

generating the context vector, and some irrelevant information can be filtered to reduce the amount of calculation. Bahdanau et al. [41] successfully apply the attention mechanism to the field of machine translation, greatly improving the translation effect. Wu et al. [35] present a Siamese attention architecture, and embed the attention mechanism into spatial gated recurrent units to selectively propagate relevant features and memorize their spatial dependencies through the network. Wu et al. [42] propose a deep attention-based spatially recursive model that can learn to attend to critical object parts and encode them into spatially expressive representations, which is composed of two-stream CNN layers, bilinear pooling, and spatial recursive encoding with attention, is end-to-end trainable to serve as the part detector and feature extractor whereby relevant features are localized, extracted, and encoded spatially for recognition purpose. Wu et al. [43] propose a network which is further augmented with 3D part alignment module to learn local features through Soft Attention module. These attended features are statistically aggregated to yield identity-discriminative representations.

Self attention is a special case of the attention mechanism. Different from the above, the self-attention mechanism only needs a sequence to calculate the representation. Through the interaction of the internal elements of the sequence, the structural information inside the sequence is learned to obtain better representation learning. It has achieved very good results in many NLP tasks, including machine translation [44], reading comprehension [45], statement representation, text summary [46], language comprehension [47] and other tasks. In this paper, we add the self-attention mechanism to our pipeline framework.

### 3 Our Two Frameworks

#### 3.1 Task Definition

Single-relation factoid questions (i.e., simple questions) can be answered by a single fact of the KBs. The formal definition is as follows. KB  $\{s_i, r_i, o_i\}$  is a set of triples, where  $s_i$  and  $o_i$  are the subject entity and object entity,  $r_i$  is the relation,  $(s_i, r_i, o_i)$  corresponds to one fact. For the purpose of answering question  $q$  formulated in natural language, we need to find a triple  $(s, r, o)$ , where  $s$  and  $o$  correspond to the subject and predicate in the question  $q$ , then  $o$  is the answer to the question  $q$ .

Therefore as long as we find the corresponding subject and predicate, we can turn question into a structured query to obtain the answer.

#### 3.2 Our Pipeline Framework

##### 3.2.1 Model Description

In this section, we divide the QA task into two subtasks: entity detection and relation detection.

1. We generate entity candidate set  $E$  by entity detection.

2. Based on entity candidate set  $E$ , we obtain all of the relations associated with the entity candidates.
3. Then we calculate the semantic similarity between the relation and the question by semantic matching model, and take the relation  $r$  with the highest matching score as the result to relation detection.
4. Finally we select the corresponding triple as the answer based on  $r$ .

As shown in the Figure 1, after getting entity candidate set, “/music/album/album\_content\_type” is the highest matching relation with the question. At this time, we get the entity-relation pair. According the entity-relation pair, we can find the triple (m.01hmylb, /music/album/album\_content\_type, m.06vw6v) as the fact to answer the given question.

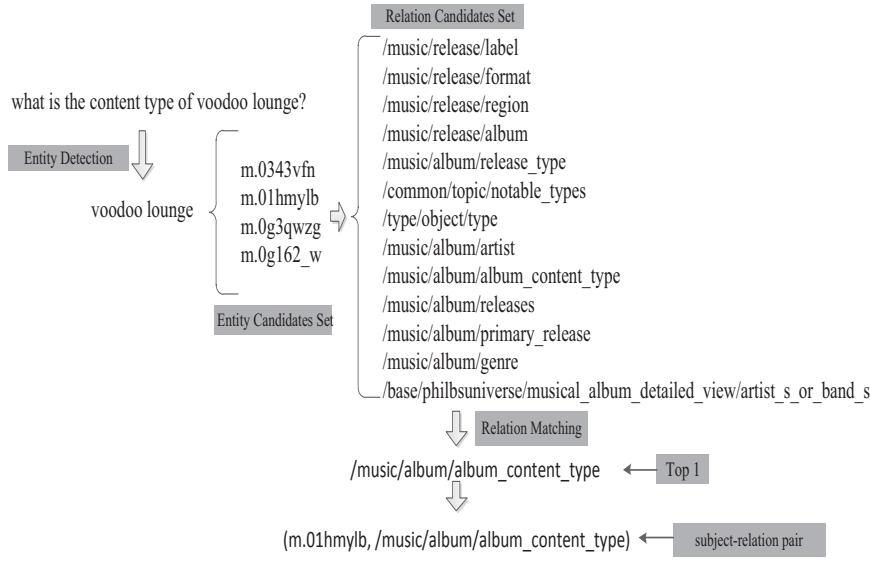


Fig. 1: Example of simple QA process: First, the entity is retrieved through entity detection according to the given question, and the entity corresponding to the fact in the KB is matched to obtain an entity candidate set. Second, according to the entities in the entity candidate set, the related relations are found, and the relation candidate set is formed. At last, according to the matching of the relation and the question, the relation with the highest matching degree is obtained, and the best entity-relation pair is obtained accordingly.

### 3.2.2 Entity Detection

Following the traditional approach, we first conduct entity recognition to mark out the words that belong to the entity in the question. Like the mainstream method, we

treat it as a sequence labeling task. We train a bidirectional LSTM [48,49] network to detect entity text in the question.

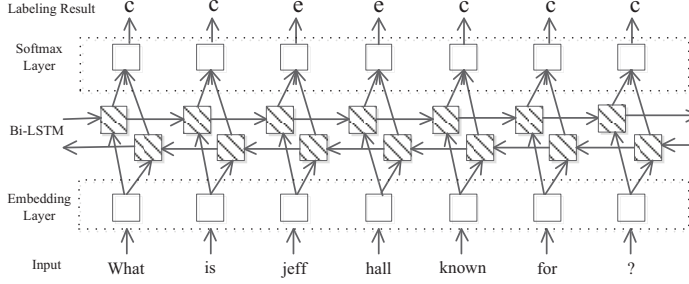


Fig. 2: The architecture of entity recognition: We use Bi-LSTM for entity recognition to extract the entity text.

As shown in Figure 2, the character “e” is the entity text, and “c” is the context text. the words that belong to the entity text are marked e, and the words that do not belong to the entity text are marked c. The fragment of entity text refers to the sequence of words corresponding to each consecutive segment e.

Because some results of entity recognition may not be completely correct, we have to employ some strategies to remedy it, therefore we borrow the method proposed by Ture et al. [50]. Based on the result of entity recognition, we obtain the fragment of entity text. Then we construct the entity candidate set through the following process.

1. We find out all entities in FB2M whose alias exactly equal to the fragment of entity text, and then form entity candidate set  $E$ . If there is no exact match entity, go to the next step.
2. The 1-gram, 2-gram and 3-gram from each fragment of entity text are extracted. If a tuple is a subset of another tuple, we should keep the long one and discard the short one. Then we can form a set of n-gram  $G$ .
3. We search the entities based on n-gram and form an entity candidate set  $E$ . Using Equation 1 to calculate the weight of the entity. Then we take the entity with the weight equal to the highest score as the result of entity detection.

$$score_i = \frac{N_i}{L_i C_i} \quad (1)$$

Where  $N_i$  is the number of words in  $G$ ,  $L_i$  is the number of words contained in the entity tag of the retrieved entity, and  $C_i$  is the number of entities in the  $E$ .

### 3.2.3 Relation Detection

In this subtask, our core mission is to measure the semantic similarity between the relations and questions. Therefore, we design a semantic matching model. We take



the matching task as a binary classification problem, matching or not matching. First we design a binary classification network. Instead of using the classification result directly, we use the value of the output layer as the matching score. As shown in Figure 3, in this process, we input the sequence text for the question and the relation, and get the word vector representation of each word through the Embedding layer. And then we get the semantic representation of the question and relation through Bi-GRU [51,52], and stitch them into a vector. As well as through a fully connected layer the final score can be obtained in the output layer. By the way, we use sigmoid as the activation function of the output layer.

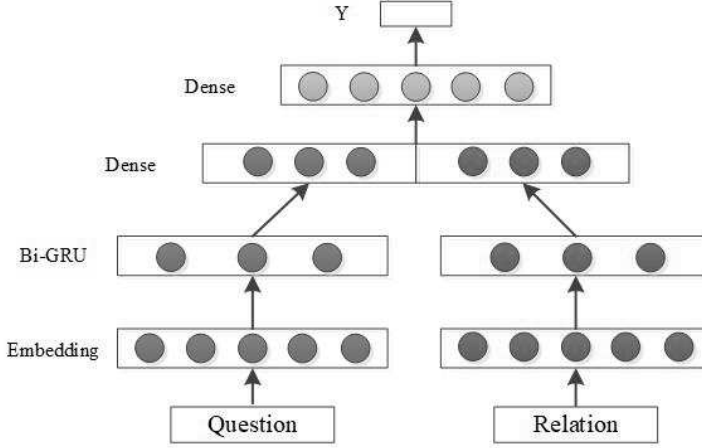


Fig. 3: The network structure of semantic matching model.

### 3.2.4 Context Information Of Pipeline Framework

In this paper, we explore to improve the resolution of the entities with the same name by entities' context information. Here context information includes the entities' out-degree and notable type.

According to whether or not the context information is combined, the selection of the entity-relation pair is divided into three different algorithms: no context information (Pipeline QA without context information i.e., P-QA) as shown in Algorithm 1, combined with out-degree information (Pipeline QA with out-degree information i.e., P-QA-Out) as shown in Algorithm 2 and combined with type information (Pipeline QA with type information i.e., P-QA-Type) as shown in Algorithm 3.

1. In Algorithm 1, after obtaining the entity candidate set  $E(e_1, e_2, \dots, e_m)$  and the relation  $r$  with the highest matching degree, the entity with the relation  $r$  can be selected as the result of the final entity detection.
2. However, sometimes there are multiple entities in  $E$  that have the same relation  $r$ , therefore here the re-sorting of the entities is performed using the out-degree information to generate the final result.

---

**Algorithm 1** Prediction algorithm without context information(P-QA)
 

---

**Input:** question  $Q$ 

     entity recognition model  $EM$ 

     relation matching model  $RM$ 
**Output:** "entity-relation" dual group  $(e, r)$ 

- 1: We use the model  $EM$  to perform entity identification on the question  $Q$ , and then the entity text  $P$  is obtained.
  - 2: The entity candidate set  $E$  is generated using entity retrieval algorithm.
  - 3: According to the entity candidate set  $E$ , we form a relation candidate set  $R$ .
  - 4: We calculate the matching score of the question  $Q$  and all relations according to the relation matching model.
  - 5: According to the relation matching result, the relation  $r$  with the highest score is selected as the result of the relation matching.
  - 6: We find the corresponding entity  $e$  from the entity candidate set  $E$  according to  $r$ .
  - 7: The prediction results i.e., entity-relation  $(e, r)$  dual group is output.
- 

In Algorithm 2, The out-degree of the entity  $e$  is the number of triples in the KB in which  $e$  is the subject. We try to rank the entity candidate set based on the out-degree. The greater the number of out-degree of an entity, the greater the scope of its association in the KB.

---

**Algorithm 2** Prediction algorithm with out-degree information(P-QA-Out)
 

---

**Input:** entity candidate set  $E(e_1, e_2, \dots, e_m)$ 

     relation  $r$ 
**Output:** "entity-relation" dual group  $(e, r)$ 

- 1: The entity candidate set  $E$  is pruned according to  $r$ , and an entity candidate subset  $E'(e'_1, e'_2, \dots, e'_k)$  with the relation  $r$  can be generated.
  - 2: We calculate the out-degree of all entities in the candidate set  $E'$ , and  $O_i(o'_1, o'_2, \dots, o'_k)$  can be generated.
  - 3: Then we sort  $E'(e'_1, e'_2, \dots, e'_k)$  in descending order according to  $O_i$ .
  - 4: The entity with the highest degree of  $e$  is selected as the result of entity detection.
  - 5: The prediction results i.e., entity-relation  $(e, r)$  dual group is output.
- 

3. In Algorithm 3, we also try to use the type information of the entity to distinguish the entities with same name. The notable type of the entity in FreeBase is a simple atomic label that indicates what the entity is notable for [53]. Freebase was acquired by Google in 2010 and officially shut down in 2016. Its data was migrated to Wikidata. Since Freebase's online API is closed, it is impossible to get the notable type information directly. The Freebase data dumps can be downloaded in an N-Triples RDF format.

We extract the notable type information from the dump files. There are 1275 kind of notable types in 2 million entities of FB2M. Then we use the same network structure like relation detection to calculate the matching score between questions and notable types. And we try to improve the accuracy of entity recognition by ranking candidate entities based on the matching scores.

Here, the matching score of the entity type information and the question is recorded as  $S_t$ , and the matching score of the relation and the question is recorded as  $S_r$ . Then the two are added to obtain the "entity-relation" binary group and the match-

ing score of question  $S$  as shown in Equation 2. The final screening of the results is achieved by scoring and sorting the entity-relation pairs.

$$S = S_t + S_r \quad (2)$$

---

**Algorithm 3** Prediction algorithm with entity type(P-QA-Type)

---

**Input:** question  $Q$

entity candidate set  $E(e_1, e_2, \dots, e_m)$

relation candidate set  $R(r_1, r_2, \dots, r_n)$

**Output:** "entity-relation" dual group  $(e, r)$

- 1: The entity type is extracted and a type list  $T(t_1, t_2, \dots, t_m)$  is generated corresponding to the entity candidate set  $E$ .
  - 2: We calculate the semantic matching score between the entity type  $T$  and the question  $Q$ ,  $S_t(s_t^1, s_t^2, \dots, s_t^m)$  is generated.
  - 3: We calculate the semantic matching score between the relation  $R$  and the question  $Q$ ,  $S_r(s_r^1, s_r^2, \dots, s_r^m)$  and entity-relation binary list  $P(p_1, p_2, \dots, p_m)$  are generated.
  - 4: According to  $S_t$  and  $S_r$ , we add them one by one according to the list  $P$  to generate  $S(s_1, s_2, \dots, s_m)$ .
  - 5: We sort the list  $P$  in descending order according to the comprehensive score  $S$ .
  - 6: The prediction results i.e., entity-relation  $(e, r)$  dual group is output.
- 

### 3.2.5 Loss Function Of Our Pipeline Framework

In our pipeline framework, we use two different loss functions.

1. In entity recognition, we use categorical\_crossentropy (i.e., Equation 3) as loss function.

$$C = -\frac{1}{n} \sum_x [y \ln a + (1 - y) \ln(1 - a)] \quad (3)$$

Where  $y$  is the expected output and  $a$  is the actual output. This function has two properties: (1) non-negative. (2) when the actual output  $a$  is close to the expected output  $y$ , the loss function is close to 0. (For instance,  $y=0, a \sim 0$ ;  $y=1, a \sim 1$ , the loss function is both close to 0.). In addition, this function also can overcome the problem that the variance cost function update weight is too slow.

2. In relation matching and in the matching of questions and type information which use Algorithm 3, The Equation 4 is as the loss function. which is a little different from Equation 3, and all characters represent the same meaning.

$$C = -\frac{1}{n} [y \ln a + (1 - y) \ln(1 - a)] \quad (4)$$

## 3.3 Our End-to-end Framework

### 3.3.1 Model Description

In this section, we explore an end-to-end framework for answering simple questions. Without entity recognition, the entity  $e$  is retrieved directly by retrieving the n-gram

tuple generated by question  $q$ . through entering the question and each fact into this model, we can get the  $score(fact)$  of each fact. And then we sort the facts by these scores, therefore we can get the final answer according to the sort result. This approach aims to build a more versatile QA model. In addition, we also study the use of self-attention mechanisms and joint representation learning, and build an end-to-end joint learning model that combines self-attention mechanisms.

As shown in Figure 4, our entire model consists of 5 major layers (or 9 small layers). The first 4 parts get the cosine similarity score of question-entity pair and question-relation pair respectively. Then we score the fact according to two different methods. We call them the automatic QA method based on weight sharing (we will call it QA-S below) and the automatic QA method based on weight sharing and multi-task (we will call it QA-T below). The shared part of the two methods in the model is further elaborated in Section 3.3.2, followed by a description of QA-S in Section 3.3.3 and a description of QA-T in Section 3.3.4.

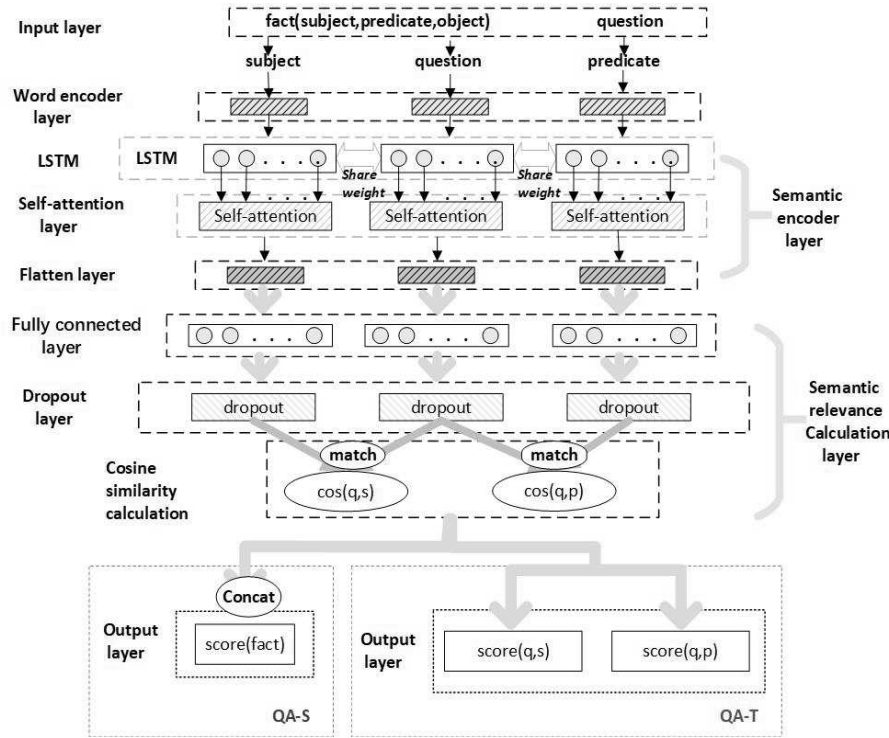


Fig. 4: Visualization of the whole model. Subject, question, predicate are sequentially respectively encoded by this model, producing their word vector representation, and then they are encoded into their semantic vector representation. Question-subject pair and question-predicate pair are scored using a cosine similarity between their representing vectors, then we can get  $score(fact)$  in two different ways.

### 3.3.2 Description Of Each Layer

1. In Input layer: We enter the fact  $f(\text{subject}, \text{predicate}, \text{object})$  and question  $q$  to the model, retrieve the subject, the predicate and send them with the question to the model for the next module to process.
2. In Word encoder layer: We encode the retrieved subject, question, and predicate to get their word vector representation, using multi-granular coding method. and the network structure is shown in Figure 5. This encoding combines word-level encoding and char-level encoding to obtain the semantic vector of the word. For the word  $w (c_1, c_2, \dots, c_n)$ , where  $c_i$  is the  $i$ th character of the word  $w$ , the corresponding word vector is  $\mathbf{v}_w$ . First, the character sequence  $(c_1, c_2, \dots, c_n)$  is input to the GRU network, with the last hidden layer state  $\mathbf{h}_i$  as the word  $w$ 's semantic representation  $\mathbf{v}_c$  of the char-level, then  $\mathbf{v}_w$  and  $\mathbf{v}_c$  are concatenated as the semantic vector  $v[\mathbf{v}_w, \mathbf{v}_c]$  of the word  $w$ .

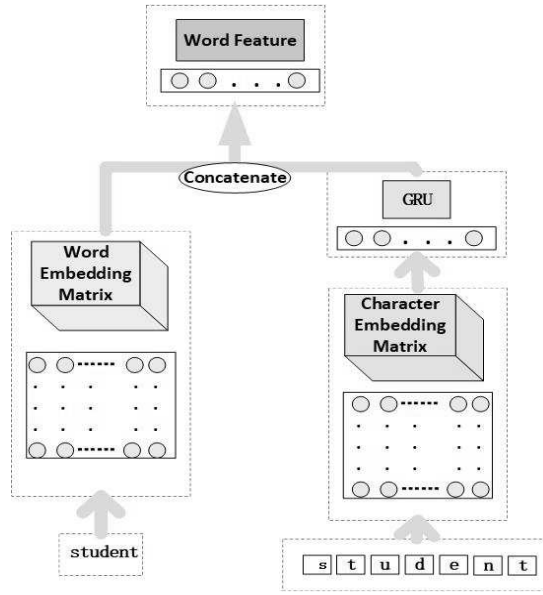


Fig. 5: The network structure of word encoder layer.

3. Semantic encoder layer: It consists of three parts: LSTM, Self-attention layer and Flatten layer.
  - (a) In LSTM: We send the semantic vector sequences of the subject, the question, and the predicate obtained by the word encoder layer into the same LSTM recurrent neural network for semantic encoding. In this part, subjects, questions and predicates share weights, which makes the processing more concise.
  - (b) In Self-attention layer: We add the self-attention mechanism to the semantic vector representation of the subject, the question and the predicate, which

can capture the long-distance interdependent features in the sentence. Since the final hidden layer state  $\mathbf{h}_i$  of LSTM is used as the semantic encoding of subjects, questions, and predicates, and the previous t-1 states have been discarded. However, the hidden state set  $h [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n]$  also corresponds to the semantic information of the current sentence. Therefore, by adding a self-attention mechanism, the model can fully calculate the internal structure of the sentence by calculating the attention between the hidden layer state sets  $h$ . And the model can generate a sentence-level embedding matrix, which makes the semantic representation vector more accurately reflect the contents of the subject, the question, and the predicate.

- (c) In Flatten layer: We compress the semantic representation matrix vectors of the subject, the question, and the predicate obtained in the previous layer into one-dimensional semantic representation vectors for convenient processing.
- 4. Semantic relevance calculation layer: It consists of Fully connected layer, Dropout layer and Cosine similarity calculation.
  - (a) In Fully connected layer: All the semantic vectors obtained are respectively made into a nonlinear transformation to obtain the probability distribution of the subject, the predicate and the question. The activation function used here is Rectified Linear Unit (i.e., ReLU) as Equation 5.

$$d = f(W_a x + b) \quad (5)$$

Where  $W_a$  is the weight coefficient of the full connection,  $b$  is the bias term, and  $f$  is the ReLU activation function as shown in Equation 6.

$$f(x) = \max(0, x) \quad (6)$$

- (b) In Dropout layer: It mainly prevents the over-fitting by temporarily discarding the neural network unit from the network according to a certain probability. The Dropout concept was proposed by Hinton et al. [54] to solve the over-fitting problem in deep image classification. Dropout refers to setting the output of each hidden neural unit to 0 with probability  $p$  as shown in Figure 6, therefore the neural unit whose output is set to 0 will not participate in the subsequent forward transmission of the network, and will not participate in the reverse. In addition Dropout can reduce the complex co-adaptation between neurons and force the network to learn more general characteristics [55].
- (c) In Cosine similarity calculation: We calculate the matching score  $score(q, s)$  between the question and the subject and the matching score  $score(q, p)$  of the question and the predicate through Equation 7 and Equation 8, respectively. The cosine similarity calculation formula is shown in Equation 9.

$$score(q, s) = \cos(\mathbf{q}, \mathbf{s}) \quad (7)$$

$$score(q, p) = \cos(\mathbf{q}, \mathbf{p}) \quad (8)$$

$$\cos(a, b) = \frac{\mathbf{a} * \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|} \quad (9)$$

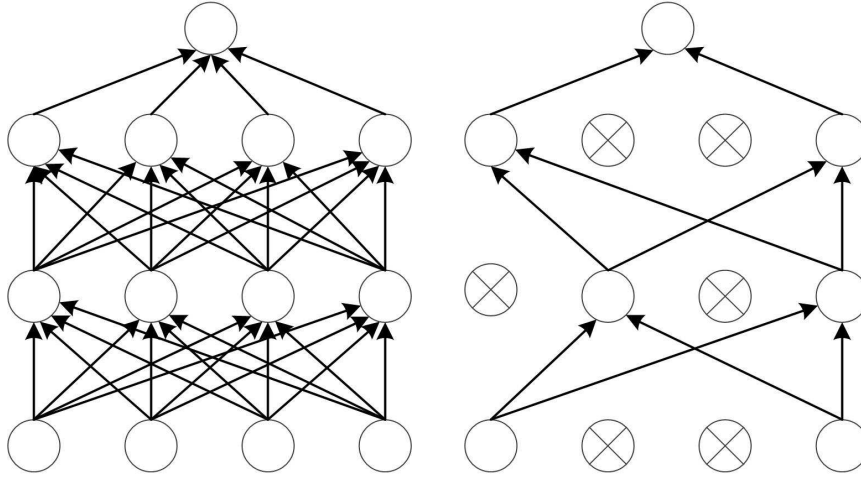


Fig. 6: Dropout layer diagram.

The above are the shared parts of the two scoring methods. Below we will describe the two different scoring methods in details.

### 3.3.3 Description Of QA-S

Through the first eight parts of the model, we can get the matching scores of the question-subject and the question-predicate pair, Based on the concatenation of them, we can get the score(fact) through Equation 10

$$score = W[S_{q,s}; S_{q,p}] \quad (10)$$

where  $W$  is the weight matrix,  $S_{q,s}$  is the matching score of question-subject pair, and  $S_{q,p}$  is the matching score of question-predicate pair.

### 3.3.4 Description Of QA-T

Multi-task learning has been successful in machine learning, such as NLP, voice recognition, and image annotation. It improves the generalization of the model by using shared representations and learning multiple tasks at the same time, and learning domain knowledge in related tasks during the training process. Rei [56] propose a sequence tagging framework with auxiliary tasks to improve the performance of the main task, which motivate the model to learn the general patterns of semantics and syntax by training the context of the predicted words. Segard et al. [57] propose a multi-level sharing model, which outputs low-level tasks in low-level networks and outputs advanced tasks in high-level networks, that is, design sharing patterns according to the semantic level of tasks. Hashimoto et al. [58] construct an end-to-end NLP framework for multi-task joint learning, which can be used to handle tasks such as part-of-speech tagging, chunking, dependency parsing, and textual entailment.

Different from QA-S, we do not concatenate the cosine similarity of the question-subject pair and the question-predicate pair, but separately, and linearly transform them according Equation 11 and Equation 12 to obtain the score of subject  $s$  and predicate  $p$  respectively.

$$score\langle q, s \rangle = w_a \cos(\mathbf{q}, \mathbf{s}) \quad (11)$$

$$score\langle q, p \rangle = w_b \cos(\mathbf{q}, \mathbf{p}) \quad (12)$$

where  $w_a, w_b$  are the linear conversion coefficients of the subject and predicate, respectively.

The above is the full introduction of our end-to-end framework. Two different methods are used to score the facts and then we will evaluate our models and methods in Section 4.

### 3.3.5 Context Information Of End-to-end Framework

Similar to Section 3.2.4, we explore the role of context information in the end-to-end framework for QA accuracy in this section.

1. We use two different methods to incorporate entity type information, as follows:
  - (a) In the end-to-end framework using method QA-T, the strategy adopted is to integrate the entity type information into the entity tag encoding process corresponding to the entity, that is, to convert the entity tag into an entity tag with entity type information. For example, for the question "**how was germany released ?**", the subject is "**germany**", its corresponding target entity is **m.017hzy7**, the entity type is "**musical recording**", then we incorporate the entity type information into the entity tag. therefore we can get the new entity label "**musical recording germany**" corresponding to **m.017hzy7**. In addition, for the two methods that do not use or use the self-attention mechanism, they are recorded as **QA-T-wt** and **QA-T-swt**, respectively.
  - (b) In the same framework with method QA-T as above, a strategy different from the above is adopted. In the Section 3.2.4, a matching model of the question and the entity type information is specifically established in order to calculate the matching score of the question and the relation. Therefore drawing on this idea, the matching of the question with the entity type information is an independent task. That is, on the output of QA-T, the matching score of the entity type information and the question is added. As an auxiliary task, it can not only improve the model to mine the entity type information in the question, but also promote the matching of the entity, the relation and the question. According to the use and non-use of the self-attention mechanism, recorded two methods as **QA-T-mwt** and **QA-T-mwst** respectively.
2. In addition to using entity type information to help distinguish entities of the same name, we can also use the out-degree information to sort the answers with the highest matching, in order to further sort and filter the answers. This section verifies the effect of the out-degree from the "Word", "Word + Self Attention" based



on the end-to-end framework with QA-S, and from the end-to-end framework with QA-T that combines the entity type information. The specific experimental results will be shown in Section 4.3.3.

### 3.3.6 Loss Function Of Our End-to-end Framework

In our end-to-end framework, We use a total of three loss functions, in which all characters have the same meaning, where  $s^+$  is the positive sample of the entity corresponding to the question  $q$ , i.e., the target entity.  $t^+$  is the entity type information of  $s^+$ ,  $s^-$  is the negative sample of the entity corresponding to  $q$ ,  $t^-$  is the entity type information of  $s^-$ ;  $p^+$  is the positive sample of the relation corresponding to  $q$ , i.e., the target relation, and  $p^-$  is the negative sample of the relation corresponding to  $q$ ,  $\gamma$  is a hyperparameter.  $S_s$ ,  $S_p$  and  $S_t$  are the matching scores corresponding to the question-subject pair, the question-predicate pair, and the question-type pair respectively.

Through the established model in Section 3.3.2, the semantic matching scores of the question-subject pair, the question-predicate pair can be separately calculated. In the training process, the positive entity sample  $s^+$  and the positive relation sample  $r^+$  of the question  $q$  are input, and a group is also input. Corresponding negative entity sample  $s^-$  and negative relation sample  $r^-$ . Under normal circumstances, the matching score of the question and the positive sample should be as much as possible than the matching score of the question and the negative sample.

(1) For the end-to-end framework with method QA-S, the loss function is shown in Equation 13. Because the matching scores of questions and subjects and predicates have been merged into a score at the output layer of the model, that is, the matching scores of the question and the "subject-predicate" correspond to  $S(q, s^+, p^+)$  and  $S(q, s^-, p^-)$ .

$$loss = \max(0, S(q, s^-, p^-) + \gamma - S(q, s^+, p^+)) \quad (13)$$

(2) For the end-to-end framework with method QA-T, the loss function is shown in Equation 14. The value of the formula consists of two parts: the entity matching loss value and the relation matching loss value.

$$loss = \sum_{(q, s^+, p^+)} (\max(0, S_s(q, s^-) - S_s(q, s^+) + \gamma) + \max(0, S_p(q, p^-) - S_p(q, p^+) + \gamma)) \quad (14)$$

The purpose of using the loss function is not to focus on the specific numerical value of the matching scores, but to focus on the difference of matching scores between the positive and negative samples and the question  $q$ , therefore strengthen the model's ability to distinguish between positive and negative samples, which can make the model learn better. The hyperparameter  $\gamma$  is defined, therefore the similarity value between the question and the positive sample can exceed the range of the similarity value between the question and the negative sample, as shown in Equation 15.

$$S_s(q, s^+) - S_s(q, s^-) \geq \gamma \quad (15)$$

(3) In addition, since in the Section 3.3.5 the matching score of the entity type and the question is added in the second method of incorporating the type information, the loss function is different from the above, and the matching score of the question and the entity type is added based on the Equation 14, as shown in Equation 16.

$$\begin{aligned} loss = \sum_{q, s^+, p^+, t^+} & (max(0, S_s(q, s^-) + \gamma - S_s(q, s^+)) + max(0, s_p(q, p^-) + \gamma - S_p(q, p^+)) \\ & + max(0, s_t(q, t^-) + \gamma - S_t(q, t^+))) \end{aligned} \quad (16)$$

## 4 Experiments

### 4.1 Dataset

we utilize SimpleQuestions dataset, which is released by Bordes et al. [26] and consists of 108442 questions written in natural language. It is constructed according corresponding facts in Freebase. The facts format as (subject, predicate, object). According to the original data division ratio, there are 75910 training data, 21687 test data, and 10845 validation data. This dataset also provides two subsets of Freebase: FB2M and FB5M. They are represented as sets of triples. We take FB2M as background KB, it includes 2 million entities and 6701 relations. In addition, we also use Freebase data dumps (22 GB compressed, 250 GB uncompressed) to extract entities' notable type information.

### 4.2 Experiment Configuration and Results Of Our Pipeline Framework

#### 4.2.1 Generation Of The Training Set

1. **Construction of the positive training set:** We construct a positive sample training set using the training data from the SimpleQuestions dataset mentioned above.
2. **Construction of the negative training set:**
  - (a) **Generation of the Entity training set  $E^-(e_1, e_2, \dots, e_{n-1})$ :** In order to generate annotation data for entity identification, it is necessary to solve the string matching of the question  $q$  with the corresponding entity  $e$  to separate the entity text from the context text. By analyzing the experimental data, the entity text in some questions  $q$  is not strictly matched with the target entity  $e$ , and there are problems such as singular and plural numbers, misspellings, etc. Therefore it is impossible to perform only an exact match between the question  $q$  and the entity  $e$ . in order to properly flag all questions, we design the following algorithm:
    - i. First, the question's word sequence  $q(w_1, w_2, \dots, w_n)$  is input to generate a 1-(n-1) tuple of  $q$ , which is denoted as  $g(g_1, g_2, \dots, g_n)$ .
    - ii. Then the following operations are performed on each entity in the KB: We calculate the edit distance  $l(l_1, l_2, \dots, l_n)$  of entity  $e$  and each element

- in the  $g$ , and the tuple  $g$  corresponding to  $l$  is taken as the matching segment.
- iii. Finally, according to the question  $q$  marked by  $g$ , all the word sequences corresponding to  $q$  are marked as ' $e$ ', and other words are marked as ' $c$ ', and all word sequences marked as ' $e$ ' constitute the Entity training set  $E^-(e_1, e_2, \dots, e_{n-1})$ .
- (b) **Generation of the Relation training set  $(q, r, tag)$ :** According to the first word of each relation, we divide all relations into 89 major categories, representing 89 domains. For example, in the domain of music, several relation examples are shown in Table 1. We generate training data in units of questions.
- For each question  $q$  and the triple  $(e, r, o)$ , we first determine the domain  $D$  according to the golden relation  $r$ .
  - Then we get all the relations  $R$  of the domain  $D$ . We generate pairs formatted in the form of  $(q, R_i, tag)$  where tag is equal to 0 or 1.
  - As illustrated in Figure 7, the relation corresponding to the question belongs to the music domain. Thus, we pair all the relations belonging to the music domain with questions and form corresponding tags. Where if relation is the target relation, then tag=1, otherwise tag=0. The Relation training set is constructed.

In addition, we copy positive cases three times in order to reduce or avoid the impact of data imbalance. Because in the construction of training data, the proportion of negative samples generated is much larger than positive samples.

Table 1: Several relation examples in music domain.

/music/live_album/concert_tour
/music/composition/compose
/music/release/label
/music/album_content_type/albums
/music/recording/producer
/music/genre/parent_genre
/music/artist/concert_tours
/music/album/compositions
.....

#### 4.2.2 Training Settings

The model word embeddings are initialized with the 300-dimensional pre-trained vectors provided by Glove [59]. We update network weights by using the Adam [60] optimizer with learning rate 0.001. The hidden layers of Bi-LSTM and Bi-GRU have size 100. In the semantic matching model, Dropout is set to 0.1.

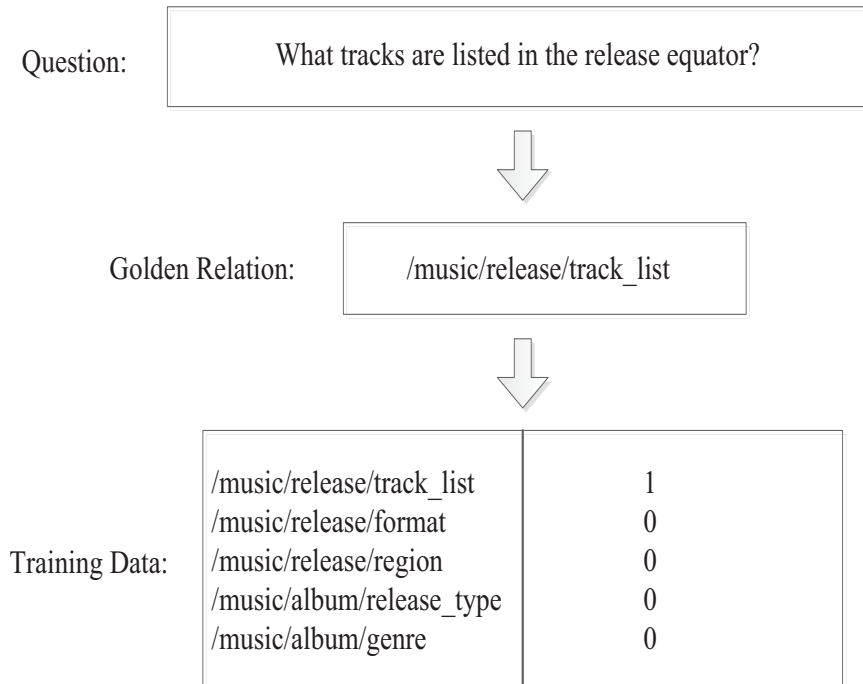


Fig. 7: Example of generating relation training data.

#### 4.2.3 Initial Preparation Work

1. All entity names are divided into words, and 1-gram, 2-gram and 3-gram are generated. Then we build an inverted index  $I$  that maps all n-grams to the entity's alias text.
2. The entity's notable type information is extracted from Freebase data dumps.
3. According to the training set, the question text should be labeled and an entity recognition training set is generated.

#### 4.2.4 Experimental Results and Discussions

Only if the entity and the relation are correctly predicted, the question  $q$  is considered being answered correctly. Therefore we use accuracy of entity-relation pair to measure the final QA results.

1. **Entity Detection:** The experiment result shows that the accuracy rate of entity recognition is 82.2%, of which 31.7% of entities cannot be uniquely identified because one name or alias may corresponds to multiple entities. In this case, the corresponding entity cannot be uniquely identified by the name or alias alone. Be-

sides, 17.8% of entities are not fully labeled correctly. This part needs to retrieve entities based on the result of entity recognition and form entity candidate sets.

2. **Relation Detection:** We generate a test data set for the model of relation detection based on the test set. For each question, we take the relations associated with its golden entity as the relation candidates. Then we use the model to find the highest matching relation. After testing, the accuracy rate of relation matching can reach 90.1%.
3. **Combination with out-degree and notable type information:** As shown in Table 2, after adding the out-degree information (method P-QA-Out), the accuracy rate increases from 65.5% to 66.6%, an increase of 0.9%, and the "Error with same label entity" has decreased by 1.1%. When the entity type information (method P-QA-Type) is added, the accuracy rate is increased by 0.7%, and the proportion of the error type "Error with same label entity" is also reduced by 0.7%.

Table 2: Comparison of results after integrating type information and out-degree in our pipeline framework. "Error with same label entity" means that the entity label in the prediction result is the same as the label of the target entity, but the entity ID is different.

Approach	Error with same label entity	Accuracy
P-QA	15.3%	65.5%
P-QA-Out	14.2%	<b>66.6%</b>
P-QA-Type	14.6%	66.2%
P-QA-Out-Type	14.3%	66.5%
P-QA-Type-Out	14.5%	66.3%

At the same time, we also explore the combination of out-degree information and entity type information. Based on the method P-QA-Out, the answers are reordered by combining the entity type information to obtain the result of the method P-QA-Out-Type. It can be seen that the final accuracy is reduced by 0.1%. In addition, in the method P-QA-Type, the answers are reordered by combining the out-degree information to obtain the result of the method P-QA-Type-Out. Compare method P-QA-Type-Out with method P-QA-Type, the final accuracy rate increased by 0.1%.

In the related work of combining entity type information in the QA, Dai et al. [28] transformed the matching problem of entity type and relation into a multi-classification task, that is, the training model classifies the question and determines the corresponding entity type in the question. On the Pipeline technology route, no strict comparisons have been found on the SimpleQuestions dataset to distinguish entities with the same name based on out-degree information.

On the Pipeline technology route, in the entity identification process, Dai et al. [28] use the Bi-GRU-CRF framework for entity identification and use Bi-GRU networks for relation matching. Yin et al. [29] use the Bi-LSTM-CRF framework for entity recognition, and then combine the character matching algorithm to classify the entity candidate set. In the fact selection process, entity label matching and relation matching are performed using a convolutional neural network (CNN) based on char-level and word-level, respectively. In this paper, we use the results of entity linking provided by Yin et al. [29], combined with the entity detection model, and the experiments are carried out according to the method in Section 3.2, and the results are shown in Table 3.

Table 3: Comparison of results after integrating entity type information and out-degree in our pipeline framework using the entity linking result from Yin et al. [29].

Approach	Error with same label entity	Accuracy
PC-QA	9.1%	70.4%
PC-QA-Out	8.3%	71.2%
PC-QA-Type	8.2%	71.3%
PC-QA-Out-Type	8.0%	71.5%
PC-QA-Type-Out	<b>7.9%</b>	<b>71.6%</b>

As can be seen from Table 3, compared with Table 2, the overall result of the experiment is greatly improved, and the effect is more obvious when using the entity type information and out-degree information. After combining them, the experimental results are found to increase by 1.1% and 1.2%, respectively. Therefore, it can be seen that entity identification and entity detection have a great influence on the results of subsequent processes, and there is a problem of error transmission. The experimental results show that the pipeline-based technical route can effectively complete the QA task. The type information and out-degree information can improve the QA results. In this process, multiple modules need to be built, and corresponding training data sets must be constructed separately. The whole process is complex and costly, and the error transmission problem needs to be solved.

### 4.3 Experiment Configuration and Results Of Our End-to-end Framework

#### 4.3.1 Generation Of The Training Set

1. **Construction of the positive training set:** We use the same method as in Section 4.2.1 to construct the positive sample training set.
2. **Construction of the negative training set:** As there are thousands of subjects and predicates in the KBs, it is impossible to use all subjects except the true sub-

ject and all predicates except the true predicate as the training data. For generating the training data set efficiently, some of the most valuable samples need to be selected to train the model. Therefore we will generate the training data set in the following three steps:

- (a) **Generation of the dictionary  $D_{rr}$ :** We calculate the edit distance  $L$  of each predicate in the Predicate set  $P$  and other predicates except it (which are expressed as  $p'$ ), then all  $p'$  are sorted in ascending order according to  $L$ , thereby the dictionary  $D_{rr}$  is generated.
- (b) **Generation of the Subject training set  $S^-(s_1, s_2, \dots, s_{n-1})$ :** For question  $q$ , the target subject is denoted as  $s$  with a label, the predicate is  $p$ , and the generated subject candidate set is  $S(s_1, s_2, \dots, s_n)$ . First, we remove the subject which has the same label with the target subject from the subject candidate set  $S$ . Then the remaining subjects constitute the Subject training set  $S^-(s_1, s_2, \dots, s_i)$  whose size is  $i$ . If  $i$  is less than 5, randomly select  $(5 - i)$  subjects from the subject candidate set  $S$ . Each time a training subject is randomly selected from  $S^-(s_1, s_2, \dots, s_5)$ .
- (c) **Generation of the Predicate training set  $P^-(p_1, p_2, \dots, p_{n-1})$ :** The true predicate  $p$  is removed from the candidate predicate set  $P$  (which is composed of all predicates related to the subject  $s$ ), then the remaining predicates constitute the Predicate training set  $P^-$  whose size is  $j$ . If  $j$  is less than 50, it is sequentially added from the dictionary  $D_{rr}$  until equal to 50. Random sampling of  $P^-$  without returning each time a negative sample is selected.

#### 4.3.2 Training Settings

We evaluate our method on the SimpleQuestions dataset which contains  $N = 21.687$  questions and the corresponding triples. For each question we follow the procedure described in Section 3.3 to find whether adding char-level encoding or adding self-attention mechanism will help the matching of question and target fact, and which scoring method can get the answer of the question more accurately.

#### 4.3.3 Experimental Results and Discussions

##### 1. Combination with char-level encoding and self-attention mechanisms:

First, we can see from the Table 4 that the multi-tasking end-to-end model based on weight sharing (QA-T) is better than the end-to-end model based on weight sharing (QA-S). The following sections will compare the two aspects of char-level encoding and self-attention mechanisms.

- (a) It can be seen that after the character information is integrated into the word encoder layer (i.e., from “Word” to “Word+Character”), it has increased by 3.3% in the model QA-S, from 56.1% to 59.4%. In the model QA-T, it increased from 67.6% to 68.8%, an increase of 1.2%.
- (b) After adding the self-attention mechanism to the semantic encoder layer (i.e., from “Word” to “Word+Self Attention”), we can see that the result of the model QA-S has improved from 56.1% to 63.7%, which improved by 7.6%;

Table 4: Comparison of results between QA-S and QA-T with char-level encoding and self-attention mechanisms. “Word” means that only word-level encoding is used in the word encoder layer, and char-level encoding is not incorporated. “Word + Self Attention” refers to the application of a self-attention mechanism at the semantic encoder level based on the “Word” coding. “Word+Character” is a multi-granular encoding method that uses word-level encoding and char-level encoding in the word encoder layer. “Word+Character+Self Attention” refers to the application of the self-attention mechanism in the semantic encoder layer based on “Word+Character”.

Approach	Accuracy	
	model QA-S	model QA-T
Word	56.1%	67.6%
Word+Self Attention	67.3%	<b>70.7%</b>
Word+Character	59.4%	68.8%
Word+Character+Self Attention	64.0%	69.8%

and the accuracy of the model QA-T has increased from 67.6% to 70.7%, an increase of 3.1%.

- (c) After adding the self-attention mechanism to the word encoder layer based on the embedding of character information (i.e., from “Word+Character” to “Word+Character+Self Attention”), the result of the model QA-S is from 59.4% to 64.0%, an increase of 4.6%; the result of the model QA-T has increased by 1.0%, from 68.8% to 69.8%.
- (d) After the multi-granularity encoding method is adopted on the basis of the semantic encoder layer adding the self-attention mechanism (i.e., from “Word + Self Attention” to “Word+Character+Self Attention”), the result of the model QA-S has increased by 0.3%. However, the result of Model QA=T has decreased by 0.9%.

From the overall comparison of model QA-S and model QA-T, we can find that the latter have a bigger improvement compared to the former. Under the same conditions, the addition of the self-attention mechanism can get improvement in both models, and the effect is more obvious in the model QA-S. The effect of the word encoder layer incorporating words’ char-level encoding is reflected in the model QA-S, Which can alleviate OOV problems and enhance the semantic representation of each word. In the model QA-T, The embedding of char-level encoding can increase the accuracy of the model by 1.2% without adding a self-attention mechanism. However, after adding the attention mechanism, the character information has a negative impact on the model’s effect, down by 0.9%.

The reason why the accuracy of the model QA-T has decreased may be that the character vector is initialized by a random vector and is continuously tuned during the training of the model. However, the data in the training data set is limited, therefore it results in a low quality of the character vector. When each element



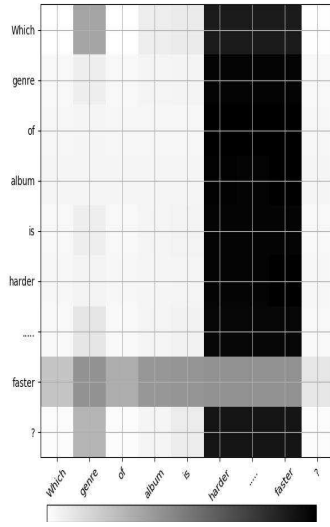


Fig. 8: Question weight distribution with example “Which genre of album is harder ..... faster?”.

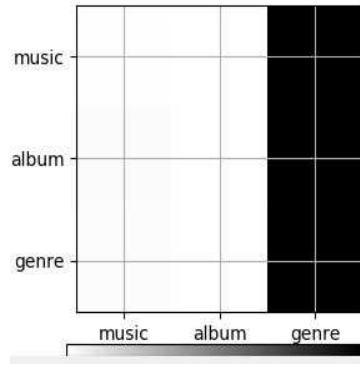


Fig. 9: Relation weight distribution with example. As can be seen in this figure, ‘genre’ has the largest weight in the relation ‘/music/album/genre’.

interacts in the self-attention mechanism, the noise is amplified, resulting in a decrease in the final result.

The experimental results show that the self-attention mechanism can better mine the internal structure of questions, entities and relations, and achieve better results without additional use of character information. The following is an example to explore the impact of the self-attention mechanism on the semantic encoder layer. For example, the question “*Which genre of album is harder ..... faster?*”, the corresponding subject is “*harder ..... faster*”, and the corresponding relation is “*music/album/genre*”. The attention weight matrix corresponding to the question and relation is shown in Figure 8 and Figure 9 respectively. The darker the color, the larger the weight.

It can be seen from Figure 8 that in the attention matrix of the question, the three words “harder ..... faster” are given a large weight, corresponding to the subject in the question. In addition, the word “genre” is given a relatively large weight in the question. It shows that the internal structure of the sentence is well studied by using the self-attention mechanism, and the subject information and predicate information are more concerned in the question.

In the attention matrix corresponding to the relation “music/album/genre” as shown in Figure 9, the weight of “genre” is large, which in turn corresponds to the predicate information “genre” in the above question. This shows that the use of the self-attention mechanism not only captures the important and discriminative

information in the relation, but also echoes the predicate information in the question.

## 2. Combination with out-degree and notable type information:

Table 5: Comparison of results after integrating entity type information.

Approach	Indistinguishable	Ambiguity	Wrong_S	Wrong_P	Accuracy
QA-T-w	3.4%	5.0%	5.5%	7.4%	67.3%
QA-T-wt	3.5%	4.0%	5.8%	7.1%	<b>68.3%</b>
QA-T-mwt	3.5%	3.9%	5.4%	7.3%	<b>68.6%</b>
QA-T-ws	3.7%	4.8%	3.1%	6.8%	70.3%
QA-T-wst	3.5%	3.8%	6.9%	6.2%	68.3%
QA-T-mwst	3.5%	3.8%	3.0%	6.8%	<b>71.5%</b>

In Table 5, QA-T-w represents the method of applying char-level encoding in the multi-task end-to-end model based on weight sharing(QA-T), and QA-T-ws represents the application of self-attention mechanism on the basis of QA-Tw. “Wrong\_S” stands for “Wrong Subject Mention” and “Wrong\_P” stands for “Wrong Predicate”.

As can be seen from Table 5, after the entity tag combines the entity type information, the accuracy of the model QA-T-wt is 1.0% higher than that of the model QA-Tw, and the case of “Ambiguity” is reduced from 5.0% to 4.0%. Therefore, when the self-attention mechanism is not added, the introduction of entity type information enriches the coding information of the entity and enhances the resolving power of the model for the case of “Ambiguity”.

After adding the self-attention mechanism, the QA-T-wst model has no improvement in the final accuracy of the model compared with the QA-T-ws model, and the error rate of the entity (“Wrong\_S”) has increased, but its “Ambiguity” situation has dropped from 4.8% to 3.8%.

From the methods QA-T-mwt and QA-T-mwst, it can be seen that the matching of the entity type and the question is a separate task, the effect of the model has been significantly improved. Method QA-T-mwt is 1.3% higher than QA-T-w. Compared with the model QA-T-ws, the accuracy rate of the model QA-T-mwst has increased from 70.3% to 71.5%, increased by 1.2%, and the “Ambiguity” error rate decreased by 1.0%, while the entity label error rate (“Wrong\_S”) fell by 0.1%.

In addition to using entity type information to help distinguish entities with the same name, we also use the out-degree information. Here the effect of it is shown in Table 6. It can be seen that when the out-degree information is added, the results of each method will have an improvement of about 0.4%, the highest result of the experiment in this paper is 71.8% on the SimpleQuestions dataset using the

QA-T-mwst method. Therefore, combining the entity information and the out-degree information can effectively improve the effect of QA.

Table 6: Comparison of the results of the model before and after adding out-degree information.

Approach	Accuracy	
	unsorted	sorted
QA-T-w	67.3%	67.6%
QA-T-ws	70.3%	70.7%
QA-T-mwst	71.5%	<b>71.8%</b>

Table 7: Comparison of end-to-end framework research results. Note: [\*] indicates that the result is a model result when FB5M is the background KB.

Approach	Accuracy
Bordes et al <sup>[26]</sup>	62.7%
Yin et al <sup>[29]</sup>	68.3%
Dai et al <sup>[28]</sup>	62.6% <sup>[*]</sup>
Golub and He al <sup>[33]</sup>	70.9%
Lukovnikov et al <sup>[34]</sup>	<b>71.2%</b>
Our approach	<b>71.8%</b>

Table 8: Comparison of model runtime and accuracy

Approach	GPU	Batch Size	Epoch	Time	Approach
Lukovnikov et al <sup>[34]</sup>	Titan X	100	50	48h	71.2%
Our Approach	GTX 1080Ti	512	40	<b>4h</b>	<b>71.8%</b>

The relevant research results of the currently known end-to-end framework on the SimpleQuestions dataset are shown in Table 7. In the end-to-end route, the work of Lukovnikov et al. [34] is currently known to be the best, and our paper uses a much less than their time to achieve a slightly better result than theirs, an increase of 0.6%.

The best result of this paper is obtained by incorporating entity type information on a multi-task end-to-end framework using word-level encoding combined with a self-attention mechanism. Lukovnikov et al. [34] also utilized character information in the encoding of questions. They use CNN to semantically model the character sequence of each word, and then the character semantic information is incorporated into the semantic representation of the corresponding word, which is used to alleviate the OOV question and also at the cost of partial time cost.

## 5 Conclusions and Future Work

In this paper, we compare the impact of the entity's out-degree, notable type information on answering single-relation factoid questions based on KBs through two different frameworks: a pipeline framework and an end-to-end framework. In the former framework, we divide this task into two subtasks: entity detection and relation detection. In the latter framework, we combine char-level encoding and self-attention mechanisms, using method of sharing weights and method of multitasking on fact selection. As the experimental results show, combining context information can improve the accuracy of QA in both frameworks. It helps to distinguish ambiguity of the entity with the same name. In addition, we find there are some ambiguities that cannot be resolved in limited context information. In practical applications, the combination of the questioner's identity information (user profile information) and more context information may solve the problem to some extent. Second, the accuracy of QA also can benefit from the adding of char-level encoding information and self-attention mechanisms. It helps achieve a more accurate semantic representation of the entity, to find the better match of fact, therefore questions can be answered with the better answer.

In future work,, we will improve the accuracy of QA from the following aspects:

1. In this paper, CNN is used to semantically encode questions, entities and relations, map them to the same semantic space, and then get the matching scores between semantic vectors by measuring the similarity of semantic vectors. However, there is no interaction between the character information of questions, entities and relations. Therefore in the subsequent research, we will try to combine more deep semantic matching models.
2. We mainly explore the use of two types of context information in this paper: entity type and out-degree information. In the future research work, we can consider introducing more context information. For example, in Freebase, the entity also has the attribute "common / topic / description", which is a holistic description of the entity. In addition, we can try to use the context information of the target entity as feedback to correct the results of QA.
3. It can be seen from the experimental process that the proportion of the negative sample has a great influence on the training result of the model within a certain range. In the follow-up study, the negative sampling method can be further explored to further improve the training efficiency and the effect of the QA model.

## References

1. Bordes, A., Weston, J., Usunier, N. Open question answering with weakly supervised embedding models. In *Joint European conference on machine learning and knowledge discovery in databases* (pp. 165-180). Springer, Berlin, Heidelberg (2014, September)
2. Bollacker, K., Cook, R., Tufts, P. Freebase: A shared database of structured general human knowledge. In *AAAI* (Vol. 7, pp. 1962-1963) (2007, July)
3. Bollacker, K., Evans, C., Paritosh, P., Sturge, T., Taylor, J. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data* (pp. 1247-1250). AcM (2008, June)
4. Li, Y. Research and Analysis of Semantic Search Technology Based on Knowledge Graph. In *2017 IEEE International Conference on Computational Science and Engineering (CSE) and IEEE International Conference on Embedded and Ubiquitous Computing (EUC)* (Vol. 1, pp. 887-890). IEEE (2017, July)
5. Suchanek, F. M., Kasneci, G., Weikum, G. Yago: a core of semantic knowledge. In *Proceedings of the 16th international conference on World Wide Web* (pp. 697-706). ACM (2007, May)
6. Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., Ives, Z. Dbpedia: A nucleus for a web of open data. In *The semantic web* (pp. 722-735). Springer, Berlin, Heidelberg. (2007).
7. Vrandečić, D., Krötzsch, M. Wikidata: a free collaborative knowledge base. (2014)
8. Navigli, R., Ponzetto, S. P. BabelNet: The automatic construction, evaluation and application of a wide-coverage multilingual semantic network. *Artificial Intelligence*, 193, 217-250 (2012)
9. Niu, X., Sun, X., Wang, H., Rong, S., Qi, G., Yu, Y. Zhishi. me-weaving chinese linking open data. In *International Semantic Web Conference* (pp. 205-220). Springer, Berlin, Heidelberg (2011, October)
10. Xu, B., Xu, Y., Liang, J., Xie, C., Liang, B., Cui, W., Xiao, Y. CN-DBpedia: a never-ending Chinese knowledge extraction system. In *International Conference on Industrial, Engineering and Other Applications of Applied Intelligent Systems* (pp. 428-438). Springer, Cham (2017, June)
11. Wang, Z., Li, J., Wang, Z., Li, S., Li, M., Zhang, D., ... Tang, J. XLORE: A Large-scale English-Chinese Bilingual Knowledge Graph. In *International semantic web conference (Posters & Demos)* (Vol. 1035, pp. 121-124) (2013, October)
12. Caifang, T., Yuan, R., Hualei, Y., Jiamin, C. Research Progress of Knowledge Graph Based on Knowledge Base Embedding. In *International Conference of Pioneering Computer Scientists, Engineers and Educators* (pp. 176-191). Springer, Singapore (2018, September)
13. Qiao, L., Yang, L., Hong, D., Yao, L., Zhiguang, Q. Knowledge graph construction techniques. *Journal of Computer Research and Development*, 53(3), 582-600 (2016)
14. Bizer, C., Heath, T., Berners-Lee, T. Linked data: The story so far. In *Semantic services, interoperability and web applications: emerging concepts* (pp. 205-227). IGI Global (2011)
15. De Virgilio, R., Del Nostro, P., Gianforme, G., Paolozzi, S. A metamodel approach to semantic web data management. In *Semantic Web Information Management* (pp. 67-91). Springer, Berlin, Heidelberg (2010)
16. Das, S., Srinivasan, J. Database technologies for RDF. In *Reasoning Web International Summer School* (pp. 205-221). Springer, Berlin, Heidelberg (2009, August)
17. Chah, N. Freebase-triples: A methodology for processing the freebase data dumps. *arXiv preprint arXiv:1712.08707* (2017)
18. Dong, X., Gabrilovich, E., Heitz, G., Horn, W., Lao, N., Murphy, K., ... Zhang, W. Knowledge vault: A web-scale approach to probabilistic knowledge fusion. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining* (pp. 601-610). ACM (2014, August)
19. Xie, Z., Zeng, Z., Zhou, G., He, T. Knowledge base question answering based on deep learning models. In *Natural Language Understanding and Intelligent Applications* (pp. 300-311). Springer, Cham (2016)
20. Wu, L., Wang, Y., & Shao, L. Cycle-consistent deep generative hashing for cross-modal retrieval. *IEEE Transactions on Image Processing*, 28(4), 1602-1612 (2019)
21. Kwiatkowski, T., Zettlemoyer, L., Goldwater, S., Steedman, M. Lexical generalization in CCG grammar induction for semantic parsing. In *Proceedings of the conference on empirical methods in natural language processing* (pp. 1512-1523). Association for Computational Linguistics (2011, July)
22. Liang, P., Jordan, M. I., Klein, D. Learning dependency-based compositional semantics. *Computational Linguistics*, 39(2), 389-446 (2013)
23. Choi, E., Kwiatkowski, T., Zettlemoyer, L. Scalable semantic parsing with partial ontologies. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th*

- International Joint Conference on Natural Language Processing (Volume 1: Long Papers) (Vol. 1, pp. 1311-1320) (2015)
24. Berant, J., Chou, A., Frostig, R., Liang, P. Semantic parsing on freebase from question-answer pairs. In Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing (pp. 1533-1544) (2013)
  25. Zhang, J., Li, W., & Ogunbona, P. Transfer learning for cross-dataset recognition: a survey. arXiv preprint arXiv:1705.04396 (2017)
  26. Bordes, A., Usunier, N., Chopra, S., Weston, J. Large-scale simple question answering with memory networks. arXiv preprint arXiv:1506.02075 (2015)
  27. Yang, B., Yih, W. T., He, X., Gao, J., Deng, L. Embedding entities and relations for learning and inference in knowledge bases. arXiv preprint arXiv:1412.6575 (2014)
  28. Dai, Z., Li, L., Xu, W. Cfo: Conditional focused neural question answering with large-scale knowledge bases. arXiv preprint arXiv:1606.01994 (2016)
  29. Yin, W., Yu, M., Xiang, B., Zhou, B., Schütze, H. Simple question answering by attentive convolutional neural network. arXiv preprint arXiv:1606.03391 (2016)
  30. Hao, Y., Liu, H., He, S., Liu, K., Zhao, J. Pattern-revising Enhanced Simple Question Answering over Knowledge Bases. In Proceedings of the 27th International Conference on Computational Linguistics (pp. 3272-3282) (2018, August)
  31. Cui, W., Xiao, Y., Wang, H., Song, Y., Hwang, S. W., Wang, W. KBQA: learning question answering over QA corpora and knowledge bases. Proceedings of the VLDB Endowment, 10(5), 565-576 (2017)
  32. Hao, Y., Zhang, Y., Liu, K., He, S., Liu, Z., Wu, H., Zhao, J. An end-to-end model for question answering over knowledge base with cross-attention combining global knowledge. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) (pp. 221-231) (2017, July)
  33. Golub, D., He, X. Character-level question answering with attention. arXiv preprint arXiv:1604.00727 (2016)
  34. Lukovnikov, D., Fischer, A., Lehmann, J., Auer, S. Neural network-based question answering over knowledge graphs on word and character level. In Proceedings of the 26th international conference on World Wide Web (pp. 1211-1220). International World Wide Web Conferences Steering Committee (2017, April)
  35. Wu, L., Wang, Y., Gao, J., Li, X. Where-and-when to look: Deep siamese attention networks for video-based person re-identification. IEEE Transactions on Multimedia (2018)
  36. Wu, L., Wang, Y., Gao, J., Li, X. Deep adaptive feature embedding with local sample distributions for person re-identification. Pattern Recognition, 73, 275-288 (2018)
  37. Wu, L., Wang, Y., Li, X., Gao, J. What-and-where to match: deep spatially multiplicative integration networks for person re-identification. Pattern Recognition, 76, 727-738 (2018)
  38. Yih, W. T., He, X., Meek, C. Semantic parsing for single-relation question answering. In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers) (Vol. 2, pp. 643-648) (2014)
  39. Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., ... Bengio, Y. Show, attend and tell: Neural image caption generation with visual attention. In International conference on machine learning (pp. 2048-2057) (2015, June)
  40. Luong, M. T., Pham, H., Manning, C. D. Effective approaches to attention-based neural machine translation. arXiv preprint arXiv:1508.04025 (2015)
  41. Bahdanau, D., Cho, K., Bengio, Y. Neural machine translation by jointly learning to align and translate. arXiv preprint arXiv:1409.0473 (2014)
  42. Wu, L., Wang, Y., Li, X., Gao, J. Deep attention-based spatially recursive networks for fine-grained visual recognition. IEEE transactions on cybernetics, (99), 1-12 (2018)
  43. Wu, L., Wang, Y., Shao, L., Wang, M. 3-d personvlad: Learning deep global representations for video-based person reidentification. IEEE transactions on neural networks and learning systems (2019)
  44. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... Polosukhin, I. Attention is all you need. In Advances in neural information processing systems (pp. 5998-6008) (2017)
  45. Cheng, J., Dong, L., Lapata, M. Long short-term memory-networks for machine reading. arXiv preprint arXiv:1601.06733 (2016)
  46. Paulus, R., Xiong, C., Socher, R. A deep reinforced model for abstractive summarization. arXiv preprint arXiv:1705.04304 (2017)
  47. Shen, T., Zhou, T., Long, G., Jiang, J., Pan, S., Zhang, C. Disan: Directional self-attention network for rnn/cnn-free language understanding. In Thirty-Second AAAI Conference on Artificial Intelligence (2018, April)

48. Hochreiter, S., Schmidhuber, J. Long short-term memory. *Neural computation*, 9(8), 1735-1780 (1997)
49. Graves, A., Schmidhuber, J. Framewise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Networks*, 18(5-6), 602-610 (2005)
50. Ture, F., Jovic, O. No Need to Pay Attention: Simple Recurrent Neural Networks Work!(for Answering" Simple" Questions). *arXiv preprint arXiv:1606.05029* (2016)
51. Chung, J., Gulcehre, C., Cho, K., Bengio, Y. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555* (2014)
52. Cho, K., Van Merriënboer, B., Bahdanau, D., Bengio, Y. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259* (2014)
53. Du, L., Kumar, A., Johnson, M., Ciaramita, M. Using entity information from a knowledge base to improve relation extraction. In *Proceedings of the Australasian Language Technology Association Workshop 2015* (pp. 31-38) (2015)
54. Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., Salakhutdinov, R. R. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580* (2012)
55. Krizhevsky, A., Sutskever, I., Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems* (pp. 1097-1105) (2012)
56. Rei, M. Semi-supervised multitask learning for sequence labeling. *arXiv preprint arXiv:1704.07156* (2017)
57. Søgaard, A., Goldberg, Y. Deep multi-task learning with low level tasks supervised at lower layers. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)* (Vol. 2, pp. 231-235) (2016)
58. Hashimoto, K., Xiong, C., Tsuruoka, Y., Socher, R. A joint many-task model: Growing a neural network for multiple nlp tasks. *arXiv preprint arXiv:1611.01587* (2016).
59. Pennington, J., Socher, R., Manning, C. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)* (pp. 1532-1543)(2014)
60. Kingma, D. P., Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*(2014)

This figure "figure-relation-weight.jpg" is available in "jpg" format from:

<http://arxiv.org/ps/1905.01995v1>



This figure "figure-whole-model.jpg" is available in "jpg" format from:

<http://arxiv.org/ps/1905.01995v1>

This figure "figure\_question\_weight.jpg" is available in "jpg" format from:

<http://arxiv.org/ps/1905.01995v1>

This figure "relation\_detection.jpg" is available in "jpg" format from:

<http://arxiv.org/ps/1905.01995v1>