

# QA4IE: A Question Answering based Framework for Information Extraction

Lin Qiu<sup>1</sup>, Hao Zhou<sup>1</sup>, Yanru Qu<sup>1</sup>, Weinan Zhang<sup>1</sup>, Suoheng Li<sup>2</sup>,  
Shu Rong<sup>2</sup>, Dongyu Ru<sup>1</sup>, Lihua Qian<sup>1</sup>, Kewei Tu<sup>3</sup>, and Yong Yu<sup>1</sup>

<sup>1</sup> Shanghai Jiao Tong University

{lqiu, kevinqu, maxru, qianlihua, yyu}@apex.sjtu.edu.cn

{zhou1998, wnzhang}@sjtu.edu.cn

<sup>2</sup> Yitu Tech

{suoheng.li, shu.rong}@yitu-inc.com

<sup>3</sup> ShanghaiTech University

tukw@shanghaitech.edu.cn

**Abstract.** Information Extraction (IE) refers to automatically extracting structured relation tuples from unstructured texts. Common IE solutions, including Relation Extraction (RE) and open IE systems, can hardly handle cross-sentence tuples, and are severely restricted by limited relation types as well as informal relation specifications (e.g., free-text based relation tuples). In order to overcome these weaknesses, we propose a novel IE framework named QA4IE, which leverages the flexible question answering (QA) approaches to produce high quality relation triples across sentences. Based on the framework, we develop a large IE benchmark with high quality human evaluation. This benchmark contains 293K documents, 2M golden relation triples, and 636 relation types. We compare our system with some IE baselines on our benchmark and the results show that our system achieves great improvements.

## 1 Introduction and Background

Information Extraction (IE), which refers to extracting structured information (i.e., relation tuples) from unstructured text, is the key problem in making use of large-scale texts. High quality extracted relation tuples can be used in various downstream applications such as Knowledge Base Population [16], Knowledge Graph Acquisition [22], and Natural Language Understanding. However, existing IE systems still cannot produce high-quality relation tuples to effectively support downstream applications.

### 1.1 Previous IE Systems

Most of previous IE systems can be divided into Relation Extraction (RE) based systems [27,51] and Open IE systems [3,8,36].

Early work on RE decomposes the problem into Named Entity Recognition (NER) and relation classification. With the recent development of neural networks (NN), NN based NER models [18,26] and relation classification models [48] show better performance than previous handcrafted feature based methods. The recently proposed

RE systems [33,52] try to jointly perform entity recognition and relation extraction to improve the performance. One limitation of existing RE benchmarks, e.g., NYT [34], Wiki-KBP [23] and BioInfer [31], is that they only involve 24, 19 and 94 relation types respectively comparing with thousands of relation types in knowledge bases such as DBpedia [4,6]. Besides, existing RE systems can only extract relation tuples from a single sentence while the cross-sentence information is ignored. Therefore, existing RE based systems are not powerful enough to support downstream applications in terms of performance or scalability.

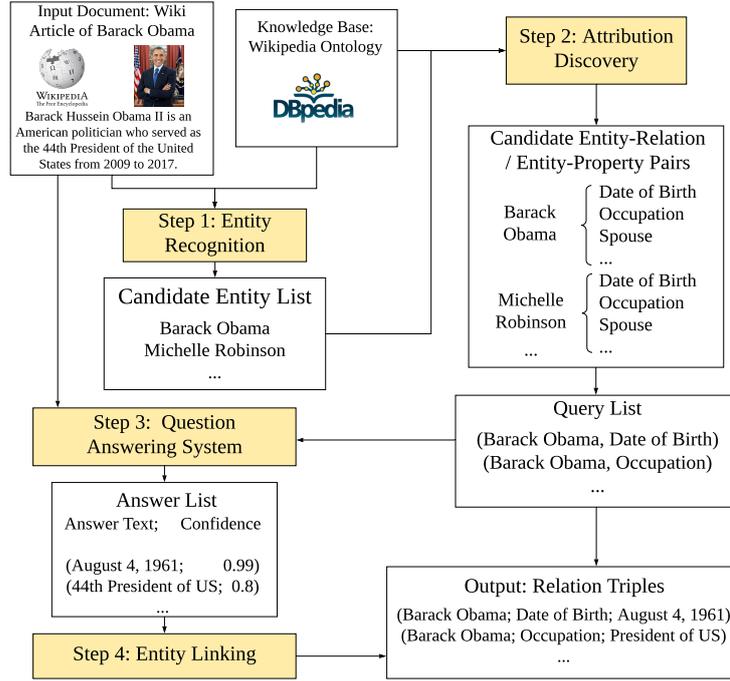
On the other hand, early work on Open IE is mainly based on bootstrapping and pattern learning methods [2]. Recent work incorporates lexical features and sentence parsing results to automatically build a large number of pattern templates, based on which the systems can extract relation tuples from an input sentence [3,8,36]. An obvious weakness is that the extracted relations are formed by free texts which means they may be polysemous or synonymous and thus cannot be directly used without disambiguation and aggregation. The extracted free-text relations also bring extra manual evaluation cost, and how to automatically evaluate different Open IE systems fairly is an open problem. Stanovsky and Dagan [41] try to solve this problem by creating an Open IE benchmark with the help of QA-SRL annotations [10]. Nevertheless, the benchmark only involves 10K golden relation tuples. Hence, Open IE in its current form cannot provide a satisfactory solution to high-quality IE that supports downstream applications.

There are some recently proposed IE approaches which try to incorporate Question Answering (QA) techniques into IE. Levy et al. [21] propose to reduce the RE problem to answering simple reading comprehension questions. They build a question template for each relation type, and by asking questions with a relevant sentence and the first entity given, they can obtain relation triples from the sentence corresponding to the relation type and the first entity. Roth et al. [35] further improve the model performance on a similar problem setting. However, these approaches focus on sentence level relation argument extractions and do not provide a full-stack solution to general IE. In particular, they do not provide a solution to extract the first entity and its corresponding relation types before applying QA. Besides, sentence level relation extraction ignores the information across sentences such as coreference and inference between sentences, which greatly reduces the information extracted from the documents.

## 1.2 QA4IE Framework

To overcome the above weaknesses of existing IE systems, we propose a novel IE framework named QA4IE to perform document level general IE with the help of state-of-the-art approaches in Question Answering (QA) and Machine Reading Comprehension (MRC) area.

The input of QA4IE is a document  $D$  with an existing knowledge base  $K$  and the output is a set of relation triples  $R = \{e_i, r_{ij}, e_j\}$  in  $D$  where  $e_i$  and  $e_j$  are two individual entities and  $r_{ij}$  is their relation. We ignore the adverbials and only consider the entity pairs and their relations as in standard RE settings. Note that we process the entire document as a whole instead of processing individual sentences separately as in



**Fig. 1.** An overview of our QA4IE Framework.

previous systems. As shown in Figure 1, our QA4IE framework consists of four key steps:

1. Recognize all the candidate entities in the input document  $D$  according to the knowledge base  $K$ . These entities serve as the first entity  $e_i$  in the relation triples  $R$ .
2. For each candidate entity  $e_i$ , discover the potential relations/properties as  $r_{ij}$  from the knowledge base  $K$ .
3. Given a candidate entity-relation or entity-property pair  $\{e_i, r_{ij}\}$  as a query, find the corresponding entity or value  $e_j$  in the input document  $D$  using a QA system. The query here can be directly formed by the word sequence of  $\{e_i, r_{ij}\}$ , or built from templates as in [21].
4. Since the results of step 3 are formed by free texts in the input document  $D$ , we need to link the results to the knowledge base  $K$ .

This framework determines each of the three elements in relation triples step by step. Step 1 is equivalent to named entity recognition (NER), and state-of-the-art NER systems [25,26] can achieve over 0.91 F1-score on CoNLL'03 [43], a well-known NER benchmark. For attribution discovery in step 2, we can take advantage of existing knowledge base ontologies such as Wikipedia Ontology to obtain a candidate relation/property list according to NER results in step 1. Besides, there is also some existing work on attribution discovery [20,49] and ontology construction [9] that can be used to solve the problem in step 2. The most difficult part in our framework is step 3 in which we need to find the entity (or value)  $e_j$  in document  $D$  according to the previous entity-relation (or entity-property) pair  $\{e_i, r_{ij}\}$ . Inspired by recent success in QA and

MRC [37,46,47], we propose to solve step 3 in the setting of SQuAD [32] which is a very popular QA task. The problem setting of SQuAD is that given a document  $\tilde{D}$  and a question  $q$ , output a segment of text  $a$  in  $\tilde{D}$  as the answer to the question. In our framework, we assign the input document  $D$  as  $\tilde{D}$  and the entity-relation (or entity-property) pair  $\{e_i, r_{ij}\}$  as  $q$ , and then we can get the answer  $a$  with a QA model. Finally in step 4, since the QA model can only produce answers formed by input free texts, we need to link the answer  $a$  to an entity  $e_j$  in the knowledge base  $K$ , and the entity  $e_j$  will form the target relation triple as  $\{e_i, r_{ij}, e_j\}$ . Existing Entity Linking (EL) systems [28,38] directly solve this problem especially when we have high quality QA results from step 3.

As mentioned above, step 1, 2 and 4 in the QA4IE framework can be solved by existing work. Therefore, in this paper, we mainly focus on step 3. According to the recent progress in QA and MRC, deep neural networks are very good at solving this kind of problem with a large-scale dataset to train the network. However, all previous IE benchmarks [41] are too small to train neural network models typically used in QA, and thus we need to build a large benchmark. Inspired by WikiReading [12], a recent large-scale QA benchmark over Wikipedia, we find that the articles in Wikipedia together with the high quality triples in knowledge bases such as Wikidata [45] and DBpedia can form the supervision we need. Therefore, we build a large scale benchmark named QA4IE benchmark which consists of 293K Wikipedia articles and 2M golden relation triples with 636 different relation types.

Recent success on QA and MRC is mainly attributed to advanced deep learning architectures such as attention-based and memory-augmented neural networks [5,42] and the availability of large-scale datasets [11,13] especially SQuAD. The differences between step 3 and SQuAD can be summarized as follows. First, the answer to the question in SQuAD is restricted to a continuous segment of the input text, but in QA4IE, we remove this constraint which may reduce the number of target relation triples. Second, in existing QA and MRC benchmarks, the input documents are not very long and the questions may be complex and difficult to understand by the model, while in QA4IE, the input documents may be longer but the questions formed by entity-relation (or entity-property) pair are much simpler. Therefore, in our model, we incorporate Pointer Networks [44] to adapt to the answers formed by any words within the document in any order as well as Self-Matching Networks [47] to enhance the ability on modeling longer input documents.

### 1.3 Contributions

The contributions of this paper are as follows:

1. We propose a novel IE framework named QA4IE to overcome the weaknesses of existing IE systems. As we discussed above, the problem of step 1, 2 and 4 can be solved by existing work and we propose to solve the problem of step 3 with QA models.
2. To train a high quality neural network QA model, we build a large IE benchmark in QA style named QA4IE benchmark which consists of 293K Wikipedia articles and 2 million golden relation triples with 636 different relation types.

3. To adapt QA models to the IE problem, we propose an approach that enhances existing QA models with Pointer Networks and Self-Matching Networks.
4. We compare our model with IE baselines on our QA4IE benchmark and achieve a great improvement over previous baselines.
5. We open source our code and benchmark for repeatable experiments and further study of IE.<sup>4</sup>

## 2 QA4IE Benchmark Construction

This section briefly presents the construction pipeline of QA4IE benchmark to solve the problem of step 3 as in our framework (Figure 1). Existing largest IE benchmark [41] is created with the help of QA-SRL annotations [10] which consists of 3.2K sentences and 10K golden extractions. Following this idea, we study recent large-scale QA and MRC datasets and find that WikiReading [12] creates a large-scale QA dataset based on Wikipedia articles and WikiData relation triples [45]. However, we observe about 11% of QA pairs with errors such as wrong answer locations or mismatch between answer string and answer words. Besides, there are over 50% of QA pairs with the answer involving words out of the input text or containing multiple answers. We consider these cases out of the problem scope of this paper and only focus on the information within the input text.

Therefore, we choose to build the benchmark referring the implementation of WikiReading based on Wikipedia articles and golden triples from Wikidata and DBpedia [4,6]. Specifically, we build our QA4IE benchmark in the following steps.

**Dump and Preprocessing.** We dump the English Wikipedia articles with Wikidata knowledge base and match each article with its corresponding relation triples according to its title. After cleaning data by removing low frequency tokens and special characters, we obtain over 4M articles and 18M triples with over 800 relation types.

**Clipping.** We discard the triples with multiple entities (or values) for  $e_j$  (account for about 6%, e.g., a book may have multiple authors). Besides, we discard the triples with any word in  $e_j$  out of the corresponding article (account for about 50%). After this step, we obtain about 3.5M articles and 9M triples with 636 relation types.

**Incorporating DBpedia.** Unlike WikiData, DBpedia is constructed automatically without human verification. Relations and properties in DBpedia are coarse and noisy. Thus we fix the existing 636 relation types in WikiData and build a projection from DBpedia relations to these 636 relation types. We manually find 148 relations which can be projected to a WikiData relation out of 2064 DBpedia relations. Then we gather all the DBpedia triples with the first entity is corresponding to one of the above 3.5M articles and the relation is one of the projected 148 relations. After the same clipping process as above and removing the repetitive triples, we obtain 394K additional triples in 302K existing Wikipedia articles.

**Distillation.** Since our benchmark is for IE, we prefer the articles with more golden triples involved by assuming that Wikipedia articles with more annotated triples are more informative and better annotated. Therefore, we figure out the distribution of the

---

<sup>4</sup> Our source code and benchmark datasets can be found at <https://github.com/SJTU-lqiu/QA4IE>

**Table 1.** Detailed Statistics of QA4IE Benchmark.

		S	M	L	Total
SPAN	# Docs	52898	29352	65124	147374
	# Triples	342361	195944	457509	995814
SEQ	# Docs	52559	29188	64385	146132
	# Triples	341820	196138	457033	994991
	# Seq-triples	46521	27176	57507	131204
	%Seq-triples	13.61	13.86	12.58	13.19

**Table 2.** Comparison between existing IE benchmarks and QA benchmarks. The first two are IE benchmarks and the rest four are QA benchmarks.

Dataset	Source	#Docs	#Triples/Queries	Remarks
<b>QA4IE Benchmark</b>	Wikipedia / WikiData / DBpedia	293K	2M	automatical generation
Open IE [41]	WSJ / Wikipedia	3.2K	10K	generated from QA-SRL annotations
Zero-Shot Benchmark [21]	Wikipedia / WikiData	N/A	30M	sentence level docs, only 120 relation types
WikiReading [12]	Wikipedia / WikiData	4.7M	18.58M	11% errors, 50% out of document answers
SQuAD [32]	Wikipedia	536	100K	crowdsourced, span answers only
CNN/Daily Mail [11]	CNN / Daily Mail	300K	1.4M	semi-synthetic cloze-style query
CBT [13]	Children’s Book	688K	688K	semi-synthetic cloze-style query

number of golden triples in articles and decide to discard the articles with less than 6 golden triples (account for about 80%). After this step, we obtain about 200K articles and 1.4M triples with 636 relation types.

**Query and Answer Assignment.** For each golden triple  $\{e_i, r_{ij}, e_j\}$ , we assign the relation/property  $r_{ij}$  as the query and the entity  $e_j$  as the answer because the Wikipedia article and its corresponding golden triples are all about the same entity  $e_i$  which is unnecessary in the queries. Besides, we find the location of each  $e_j$  in the corresponding article as the answer location. As we discussed in Section 1, we do not restrict  $e_j$  to a continuous segment in the article as required in SQuAD. Thus we first try to detect a matched span for each  $e_j$  and assign this span as the answer location. Then for each of the rest  $e_j$  which has no matched span, we search a matched sub-sequence in the article and assign the index sequence as the answer location. We name them **span-triples** and **seq-triples** respectively. Note that each triple will have an answer location because we have discarded the triples with unseen words in  $e_j$  and if we can find multiple answer locations, all of them will be assigned as ground truths.

**Dataset Splitting.** For comparing the performance on span-triples and seq-triples, we set up two different datasets named QA4IE-SPAN and QA4IE-SEQ. In QA4IE-SPAN, only articles with all span-triples are involved, while in QA4IE-SEQ, articles with seq-triples are also involved. For studying the influence of the article length as longer articles are normally more difficult to model by LSTMs, we split the articles according to the article length. We name the set of articles with lengths shorter than 400 as S, lengths between 400 and 700 as M, lengths greater than 700 as L. Therefore we obtain 6 different datasets named QA4IE-SPAN-S/M/L and QA4IE-SEQ-S/M/L. A 5/1/5 splitting of train/dev/test sets is performed. The detailed statistics of QA4IE benchmark are provided in Table 1.

We further compare our QA4IE benchmark with some existing IE and QA benchmarks in Table 2. One can observe that QA4IE benchmark is much larger than previous IE and QA benchmarks except for WikiReading and Zero-Shot Benchmark.

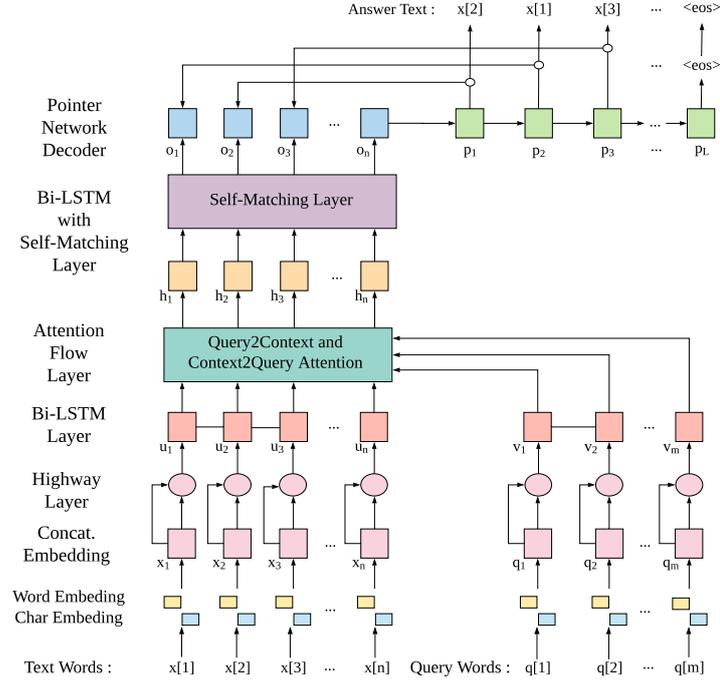


Fig. 2. An overview of our QA model.

However, as we mentioned at the beginning of Section 2, WikiReading is problematic for IE settings. Besides, Zero-Shot Benchmark is a sentence-level dataset and we have described the disadvantage of ignoring information across sentences at Section 1.1. Thus to our best knowledge, QA4IE benchmark is the largest document level IE benchmark and it can be easily extended if we change our distillation strategy.

### 3 Question Answering Model

In this section, we describe our Question Answering model for IE. The model overview is illustrated in Figure 2. The input of our model are the words in the input text  $x[1], \dots, x[n]$  and query  $q[1], \dots, q[n]$ . We concatenate pre-trained word embeddings from GloVe [30] and character embeddings trained by CharCNN [17] to represent input words. The  $2d$ -dimension embedding vectors of input text  $x_1, \dots, x_n$  and query  $q_1, \dots, q_n$  are then fed into a Highway Layer [40] to improve the capability of word embeddings and character embeddings as

$$\begin{aligned}
 g_t &= \text{sigmoid}(W_g x_t + b_g) \\
 s_t &= \text{relu}(W_x x_t + b_x) \\
 u_t &= g_t \odot s_t + (1 - g_t) \odot x_t.
 \end{aligned} \tag{1}$$

Here  $W_g, W_x \in \mathbb{R}^{d \times 2d}$  and  $b_g, b_x \in \mathbb{R}^d$  are trainable weights,  $u_t$  is a  $d$ -dimension vector. The function  $\text{relu}$  is the rectified linear units [19] and  $\odot$  is element-wise multiply over two vectors. The same Highway Layer is applied to  $q_t$  and produces  $v_t$ .

Next,  $u_t$  and  $v_t$  are fed into a Bi-Directional Long Short-Term Memory Network (BiLSTM) [14] respectively in order to model the temporal interactions between sequence words:

$$\begin{aligned} u'_t &= \text{BiLSTM}(u'_{t-1}, u_t) \\ v'_t &= \text{BiLSTM}(v'_{t-1}, v_t). \end{aligned} \quad (2)$$

Here we obtain  $\mathbf{U} = [u'_1, \dots, u'_n] \in \mathbb{R}^{2d \times n}$  and  $\mathbf{V} = [v'_1, \dots, v'_m] \in \mathbb{R}^{2d \times m}$ . Then we feed  $\mathbf{U}$  and  $\mathbf{V}$  into the attention flow layer [37] to model the interactions between the input text and query. We obtain the  $8d$ -dimension query-aware context embedding vectors  $h_1, \dots, h_n$  as the result.

After modeling interactions between the input text and queries, we need to enhance the interactions within the input text words themselves especially for the longer text in IE settings. Therefore, we introduce Self-Matching Layer [47] in our model as

$$\begin{aligned} o_t &= \text{BiLSTM}(o_{t-1}, [h_t, c_t]) \\ s_j^t &= w^T \tanh(W_h h_j + \tilde{W}_h h_t) \\ \alpha_i^t &= \exp(s_i^t) / \sum_{j=1}^n \exp(s_j^t) \\ c_t &= \sum_{i=1}^n \alpha_i^t h_i. \end{aligned} \quad (3)$$

Here  $W_h, \tilde{W}_h \in \mathbb{R}^{d \times 8d}$  and  $w \in \mathbb{R}^d$  are trainable weights,  $[h, c]$  is vector concatenation across row. Besides,  $\alpha_i^t$  is the attention weight from the  $t^{\text{th}}$  word to the  $i^{\text{th}}$  word and  $c_t$  is the enhanced contextual embeddings over the  $t^{\text{th}}$  word in the input text. We obtain the  $2d$ -dimension query-aware and self-enhanced embeddings of input text after this step. Finally we feed the embeddings  $\mathbf{O} = [o_1, \dots, o_n]$  into a Pointer Network [44] to decode the answer sequence as

$$\begin{aligned} p_t &= \text{LSTM}(p_{t-1}, c_t) \\ s_j^t &= w^T \tanh(W_o o_j + W_p p_{t-1}) \\ \beta_i^t &= \exp(s_i^t) / \sum_{j=1}^n \exp(s_j^t) \\ c_t &= \sum_{i=1}^n \beta_i^t o_i. \end{aligned} \quad (4)$$

The initial state of LSTM  $p_0$  is  $o_n$ . We can then model the probability of the  $t^{\text{th}}$  token  $a^t$  by

$$\begin{aligned} P(a^t | a^1, \dots, a^{t-1}, \mathbf{O}) &= (\beta_1^t, \beta_2^t, \dots, \beta_n^t, \beta_{n+1}^t) \\ P(a_i^t) &\triangleq P(a^t = i | a^1, \dots, a^{t-1}, \mathbf{O}) = \beta_i^t. \end{aligned} \quad (5)$$

Here  $\beta_{n+1}^t$  denotes the probability of generating the ‘‘eos’’ symbol since the decoder also needs to determine when to stop. Therefore, the probability of generating the answer sequence  $\mathbf{a}$  is as follows

$$P(\mathbf{a} | \mathbf{O}) = \prod_t P(a^t | a^1, \dots, a^{t-1}, \mathbf{O}). \quad (6)$$

Given the supervision of answer sequence  $\mathbf{y} = (y_1, \dots, y_L)$ , we can write down the loss function of our model as

$$L(\theta) = - \sum_{t=1}^L \log P(a_{y_t}^t). \quad (7)$$

To train our model, we minimize the loss function  $L(\theta)$  based on training examples.

**Table 3.** Comparison of QA models on SQuAD datasets. We only include the single model results on the dev set from published papers.

	Dev Set
<i>Span Model</i>	EM / F1
LR Baseline [32]	40.0 / 51.0
Match-LSTM [46]	64.1 / 73.9
BiDAF [37]	67.7 / 77.3
R-Net [47]	71.1 / 79.5
MEMEN [29]	71.0 / 80.4
M-Reader+RL [15]	72.1 / 81.6
SAN [24]	<b>76.2 / 84.1</b>
<i>Sequence Model</i>	
Match-LSTM (Seq) [46]	54.4 / 68.2
<b>Our Model</b>	<b>61.7 / 72.5</b>

## 4 Experiments

### 4.1 Experimental Setup

We build our QA4IE benchmark following the steps described in Section 2. In experiments, we train and evaluate our QA models on the corresponding train and test sets while the hyper-parameters are tuned on dev sets. In order to make our experiments more informative, we also evaluate our model on SQuAD dataset [32].

The preprocessing of our QA4IE benchmark and SQuAD dataset are all performed with the open source code from [37]. We use 100 1D filters with width 5 to construct the CharCNN in our char embedding layer. We set the hidden size  $d = 100$  for all the hidden states in our model. The optimizer we use is the AdaDelta optimizer [50] with an initial learning rate of 2. A dropout [39] rate of 0.2 is applied in all the CNN, LSTM and linear transformation layers in our model during training. For SQuAD dataset and our small sized QA4IE-SPAN/SEQ-S datasets, we set the max length of input texts as 400 and a mini-batch size of 20. For middle sized (and large sized) QA4IE datasets, we set the max length as 700 (800) and batch size as 7 (5). We introduce an early stopping in training process after 10 epochs. Our model is trained on a GTX 1080 Ti GPU and it takes about 14 hours on small sized QA4IE datasets. We implement our model with TensorFlow [1] and optimize the computational expensive LSTM layers with LSTMBlockFusedCell<sup>5</sup>.

### 4.2 Results in QA Settings

We first perform experiments in QA settings to evaluate our QA model on both SQuAD dataset and QA4IE benchmark. Since our goal is to solve IE, not QA, the motivation of this part of experiments is to evaluate the performance of our model and make a comparison between QA4IE benchmark and existing datasets. Two metrics are introduced in the SQuAD dataset: Exact Match (EM) and F1-score. EM measures the percentage that the model prediction matches one of the ground truth answers exactly while F1-score measures the overlap between the prediction and ground truth answers. Our QA4IE benchmark also adopts these two metrics.

<sup>5</sup> [https://www.tensorflow.org/api\\_docs/python/tf/contrib/rnn/LSTMBlockFusedCell](https://www.tensorflow.org/api_docs/python/tf/contrib/rnn/LSTMBlockFusedCell)

**Table 4.** Comparison of QA models on 6 datasets of our QA4IE benchmark. The BiDAF model cannot work on our SEQ datasets thus the results are N/A.

Model	SPAN-S	SPAN-M	SPAN-L	SEQ-S	SEQ-M	SEQ-L
	EM / F1					
BiDAF [37]	88.89 / 90.89	82.37 / 85.04	68.00 / 70.29	N/A	N/A	N/A
Match-LSTM [46]	85.88 / 88.21	79.19 / 82.05	66.87 / 70.44	89.60 / 91.95	83.57 / 87.40	62.64 / 68.98
<b>Our Model</b>	<b>91.53 / 93.19</b>	<b>86.04 / 88.65</b>	<b>70.86 / 74.51</b>	<b>91.20 / 93.04</b>	<b>85.52 / 88.43</b>	<b>71.96 / 76.11</b>

**Table 5.** Model ablation on QA4IE-SEQ-S. The first line is our original model and each of the following lines is the original model with a component ablated.

	EM / F1	Training hours
<b>Our Original Model</b>	91.20 / 93.04	14
– Char Embedding	89.78 / 91.76	14
– Highway	90.04 / 91.97	14
– Self Matching	89.55 / 91.60	10
– LSTMBlockFusedCell		90

Table 3 presents the results of our QA model on SQuAD dataset. Our model outperforms the previous sequence model but is not competitive with span models because it is designed to produce sequence answers in IE settings while baseline span models are designed to produce span answers for SQuAD dataset.

The comparison between our QA model and two baseline QA models on our QA4IE benchmark is shown in Table 4. For training of both baseline QA models,<sup>6</sup> we use the same configuration of max input length as our model and tune the rest of hyperparameters on dev sets. Our model outperforms these two baselines on all 6 datasets. The performance is good on S and M datasets but worse for longer documents. As we mentioned in Section 4.1, we set the max input length as 800 and ignore the rest words on L datasets. Actually, there are 11% of queries with no answers in the first 800 words in our benchmark. Processing longer documents is a tough problem [7] and we leave this to our future work.

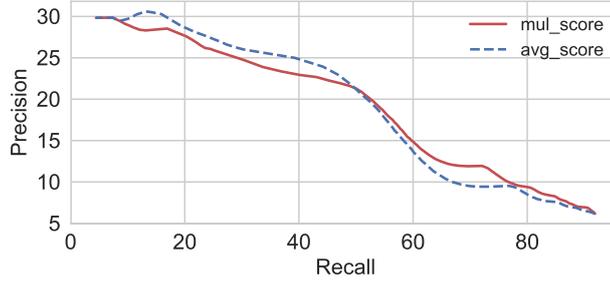
To study the improvement of each component in our model, we present model ablation study results in Table 5. We do not involve Attention Flow Layer and Pointer Network Decoder as they cannot be replaced by other architectures with the model still working. We can observe that the first three components can effectively improve the performance but Self Matching Layer makes the training more computationally expensive by 40%. Besides, the LSTMBlockFusedCell works effectively and accelerates the training process by 6 times without influencing the performance.

### 4.3 Results in IE Settings

In this subsection, we put our QA model in the entire pipeline of our QA4IE framework (Figure 1) and evaluate the framework in IE settings. Existing IE systems are all free-text based Open IE systems, so we need to manually evaluate the free-text based results in order to compare our model with the baselines. Therefore, we conduct experiments

<sup>6</sup> The code of BiDAF is from <https://github.com/allenai/bi-att-flow>.

The code of Match-LSTM is from <https://github.com/fuhuamosi/MatchLstm>.



**Fig. 3.** Precision-recall curves with two confidence scores on the dev set of QA4IE-SPAN-S.

on a small dataset, the dev set of QA4IE-SPAN-S which consists of 4393 documents and 28501 ground truth queries.

Our QA4IE benchmark is based on Wikipedia articles and all the ground truth triples of each article have the same first entity (i.e. the title of the article). Thus, we can directly use the title of the article as the first entity of each triple without performing step 1 (entity recognition) in our framework. Besides, all the ground truth triples in our benchmark are from knowledge base where they are disambiguated and aggregated in the first place, and therefore step 4 (entity linking) is very simple and we do not evaluate it in our experiments.

A major difference between QA settings and IE settings is that in QA settings, each query corresponds to an answer, while in the QA4IE framework, the QA model take a candidate entity-relation (or entity-property) pair as the query and it needs to tell whether an answer to the query can be found in the input text. We can consider the IE settings here as performing step 2 and then step 3 in the QA4IE framework.

In step 2, we need to build a candidate query list for each article in the dataset. Instead of incorporating existing ontology or knowledge base, we use a simple but effective way to build the candidate query list of an article. Since we have a ground truth query list with labeled answers of each article, we can add all the neighboring queries of each ground truth query into the query list. The neighboring queries are defined as two queries that co-occur in the same ground truth query list of any articles in the dataset. We transform the dev set of QA4IE-SPAN-S above by adding neighboring queries into the query list. After this step, the number of queries grows to 426336, and only 28501 of them are ground truth queries labeled with an answer.

In step 3, we require our QA model to output a confidence score along with the answer to each candidate query. Our QA model produces no answer to a query when the confidence score is less than a threshold  $\delta$  or the output is an “eos” symbol. For the answers with a confidence score  $\geq \delta$ , we evaluate them by the EM measurement with ground truth answers and count the true positive samples in order to calculate the precision and recall under the threshold  $\delta$ . Specifically, we try two confidence scores calculated as follows:

$$\text{Score}_{\text{mul}} = \prod_{t=1}^L P(a_{i_t}^t), \quad \text{Score}_{\text{avg}} = \sum_{t=1}^L P(a_{i_t}^t)/L, \quad (8)$$

where  $(a_{i_1}^1, \dots, a_{i_L}^L)$  is the answer sequence and  $P(a_i^t)$  is defined in Eq. (5).  $\text{Score}_{\text{mul}}$  is equivalent to the training loss in Eq. (7) and  $\text{Score}_{\text{avg}}$  takes the answer length into account.

**Table 6.** Results of three Open IE baselines on the dev set of QA4IE-SPAN-S.

	Open IE 4	Stanford IE	ClauseIE
#Extracted Triples	32309	120147	75078
#After Filtering	487	467	554
#True Positive	403	301	133

The precision-recall curves of our framework based on the two confidence scores are plotted in Figure 3. We can observe that the EM rate we achieve in QA settings is actually the best recall (91.87) in this curve (by setting  $\delta = 0$ ). The best F1-scores of the two curves are 29.97 (precision = 21.61, recall = 48.85,  $\delta = 0.91$ ) for  $\text{Score}_{\text{mul}}$  and 31.05 (precision = 23.93, recall = 44.21,  $\delta = 0.97$ ) for  $\text{Score}_{\text{avg}}$ .  $\text{Score}_{\text{avg}}$  is better than  $\text{Score}_{\text{mul}}$ , which suggests that the answer length should be taken into account.

We then evaluate existing IE systems on the dev set of QA4IE-SPAN-S and empirically compare them with our framework. Note that while [21] is closely related to our work, we cannot fairly compare our framework with [21] because their systems are in the sentence level and require additional negative samples for training. [35] is also related to our work, but their dataset and code have not been published yet. Therefore, we choose to evaluate three popular Open IE systems, Open IE 4 [36], Stanford IE [3] and ClauseIE [8].

Since Open IE systems take a single sentence as input and output a set of free-text based triples, we need to find the sentences involving ground truth answers and feed the sentences into the Open IE systems. In the dev set of QA4IE-SPAN-S, there are 28501 queries with 44449 answer locations labeled in the 4393 documents. By feeding the 44449 sentences into the Open IE systems, we obtain a set of extracted triples from each sentence. We calculate the number of true positive samples by first filtering out triples with less than 20% words overlapping with ground truth answers and then asking two human annotators to verify the remaining triples independently. Since in the experiments, our framework is given the ground-truth first entity of each triple (the title of the corresponding Wikipedia article) while the baseline systems do not have this information, we ask our human annotators to ignore the mistakes on the first entities when evaluating triples produced by the baseline systems to offset this disadvantage. For example, the 3rd case of ClauseIE and the 4th case of Open IE 4 in Table 7 are all labeled as correct by our annotators even though the first entities are pronouns. The two human annotators reached an agreement on 191 out of 195 randomly selected cases.

The evaluation results of the three Open IE baselines are shown in Table 6. We can observe that most of the extracted triples are not related to ground truths and the precision and recall are all very low (around 1%) although we have already helped the baseline systems locate the sentences containing ground truth answers.

#### 4.4 Case Study

In this subsection, we perform case studies of IE settings in Table 7 to better understand the models and benchmarks. The baseline Open IE systems produce triples by analyzing the subjects, predicates and objects in input sentences, and thus our annotators lower the bar of accepting triples. However, the analysis on semantic roles and parsing

**Table 7.** Case study of three Open IE baselines and our framework on dev set of QA4IE-SPAN-S, the results of baselines are judged by two human annotators while the results of our framework are measured by Exact Match with ground truth. The triples in red indicate the wrong cases.

Input Sentence	Ground Truth Triple	Open IE 4	Stanford IE	ClauseIE	Ours
Dieter Kesten was born on 9 June 1914 at Gelsenkirchen.	(Dieter Kesten; date of birth; 9 June 1914)	(Dieter Kesten; was born; on 9 June 1914 at Gelsenkirchen)	(Dieter Kesten; was born on; 9 June 1914)	(Dieter Kesten; was born; on 9 June 1914)	(Dieter Kesten; date of birth; 9 June 1914)
Hamilton died on 2 March 1625 at Whitehall, London, from a fever and was buried in the family mausoleum at Hamilton, on 2 September of that year.	(James Hamilton; date of death; 2 March 1625)	(Hamilton; died; on 2 March 1625 at Whitehall)	(Hamilton; died on; 2 September) (Hamilton; died on; 2 March 1625)	(Hamilton; died on; 2)	(James Hamilton; date of death; 2 March 1625)
She attended Texas A&M University, where she swam for the Texas A&M Aggies swimming and diving team in National Collegiate Athletic Association (NCAA) competition from 2011 to 2014.	(Breeja Larson; member of sports team; Texas A&M Aggies)	(She; attended; Texas A&M University)	(She; attended; M University)	(She; attended; Texas A&M University) (she; swam; for the Texas A&M Aggies swimming and diving team)	(Breeja Larson; member of sports team; Texas A&M Aggies)
His grave and memorial are at Balbeggie Churchyard, St. Martin's, near Perth, Scotland.	(John Simpson; place of death; St. Martin's)	(His grave and memorial; are; at Balbeggie Churchyard, St. Martin's, near Perth)	(Perth; near Churchyard is; St. Martin's)	(Balbeggie Churchyard near Perth Scotland; is; St. Martin's)	(John Simpson; place of death; Balbeggie Churchyard)
He served in the British Army and was wounded in World War I.	(William Dobbie; conflict; World War I)	(He; was wounded; in World War I)	(He; was wounded in; World War I)	(He; was wounded; in World War I)	(William Dobbie; conflict; World War I)

trees cannot work very well on complicated input sentences like the 2nd and the 3rd cases. Besides, the baseline systems can hardly solve the last two cases which require inference on input sentences.

Our framework works very well on this dataset with the QA measurements  $EM = 91.87$  and  $F1 = 93.53$  and the IE measurements can be found in Figure 3. Most of the error cases are the fourth case which is acceptable by human annotators. Note that our framework takes the whole document as the input while the baseline systems take the individual sentence as the input, which means the experiment setting is much more difficult for our framework.

#### 4.5 Human Evaluation on QA4IE Benchmark

Finally, we perform a human evaluation on our QA4IE benchmark to verify the reliability of former experiments. The evaluation metrics are as follows:

**Triple Accuracy** is to check whether each ground truth triple is accurate (one cannot find conflicts between the ground truth triple and the corresponding article) because the ground truth triples from WikiData and DBpedia may be incorrect or incomplete.

**Contextual Consistency** is to check whether the context of each answer location is consistent with the corresponding ground truth triple (one can infer from the context to obtain the ground truth triple) because we keep all matched answer locations as ground truths but some of them may be irrelevant with the corresponding triple.

**Triple Consistency** is to check whether there is at least one answer location that is contextually consistent for each ground truth triple. It can be calculated by counting the results of Contextual Consistency.

We randomly sample 25 articles respectively from the 6 datasets (in total of 1002 ground truth triples with 2691 labeled answer locations) and let two human annotators label the Triple Accuracy for each ground truth triple and the Contextual Consistency for each answer location. The two human annotators reached an agreement on 131 of 132 randomly selected Triple Accuracy cases and on 229 of 234 randomly selected Contextual Consistency cases. The human evaluation results are shown in Table 8. We can find that the Triple Accuracy and the Triple Consistency is acceptable while

**Table 8.** Human evaluation on QA4IE benchmark.

	SPAN-S	SPAN-M	SPAN-L	SEQ-S	SEQ-M	SEQ-L	Total
Triple Accuracy	98.8%	96.9%	98.1%	97.1%	96.2%	97.8%	97.5%
	161 / 163	154 / 159	159 / 162	170 / 175	152 / 158	181 / 185	977 / 1002
Contextual Consistency	78.6%	65.1%	70.3%	75.4%	73.9%	82.4%	74.6%
	195 / 248	239 / 367	494 / 703	230 / 305	264 / 357	586 / 711	2008 / 2691
Triple Consistency	93.3%	87.4%	91.4%	92.0%	92.4%	92.4%	91.5%
	152 / 163	139 / 159	148 / 162	161 / 175	146 / 158	171 / 185	917 / 1002

the Contextual Consistency still needs to be improved. The Contextual Consistency problem is a weakness of distant supervision, and we leave this to our future work.

## 5 Conclusion

In this paper, we propose a novel QA based IE framework named QA4IE to address the weaknesses of previous IE solutions. In our framework (Figure 1), we divide the complicated IE problem into four steps and show that the step 1, 2 and 4 can be solved well enough by existing work. For the most difficult step 3, we transform it to a QA problem and solve it with our QA model. To train this QA model, we construct a large IE benchmark named QA4IE benchmark that consists of 293K documents and 2 million golden relation triples with 636 different relation types. To our best knowledge, our QA4IE benchmark is the largest document level IE benchmark. We compare our system with existing best IE baseline systems on our QA4IE benchmark and the results show that our system achieves a great improvement over baseline systems.

For the future work, we plan to solve the triples with multiple entities as the second entity, which is excluded from problem scope in this paper. Besides, processing longer documents and improving the quality of our benchmark are all challenging problems as we mentioned previously. We hope this work can provide new thoughts for the area of information extraction.

## Acknowledgements

W. Zhang is the corresponding author of this paper. The work done by SJTU is sponsored by National Natural Science Foundation of China (61632017, 61702327, 61772333) and Shanghai Sailing Program (17YF1428200).

## References

1. Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al.: Tensorflow: A system for large-scale machine learning. In: OSDI. vol. 16, pp. 265–283 (2016)
2. Agichtein, E., Gravano, L.: Snowball: Extracting relations from large plain-text collections. In: Proceedings of the fifth ACM conference on Digital libraries. pp. 85–94. ACM (2000)
3. Angeli, G., Premkumar, M.J.J., Manning, C.D.: Leveraging linguistic structure for open domain information extraction. In: ACL. vol. 1, pp. 344–354 (2015)
4. 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 (2007)

5. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. In: International Conference on Learning Representations (ICLR) (2015)
6. Bizer, C., Lehmann, J., Kobilarov, G., Auer, S., Becker, C., Cyganiak, R., Hellmann, S.: Dbpedia-a crystallization point for the web of data. *Web Semantics: science, services and agents on the world wide web* 7(3), 154–165 (2009)
7. Chen, D., Fisch, A., Weston, J., Bordes, A.: Reading wikipedia to answer open-domain questions. In: ACL. vol. 1, pp. 1870–1879 (2017)
8. Del Corro, L., Gemulla, R.: Clausie: clause-based open information extraction. In: Proceedings of international conference on World Wide Web. pp. 355–366 (2013)
9. Gupta, R., Halevy, A., Wang, X., Whang, S.E., Wu, F.: Biperpedia: An ontology for search applications. *Proceedings of the VLDB Endowment* 7(7), 505–516 (2014)
10. He, L., Lewis, M., Zettlemoyer, L.: Question-answer driven semantic role labeling: Using natural language to annotate natural language. In: EMNLP. pp. 643–653 (2015)
11. Hermann, K.M., Kocisky, T., Grefenstette, E., Espeholt, L., Kay, W., Suleyman, M., Blunsom, P.: Teaching machines to read and comprehend. In: Advances in Neural Information Processing Systems. pp. 1693–1701 (2015)
12. Hewlett, D., Lacoste, A., Jones, L., Polosukhin, I., Fandrianto, A., Han, J., Kelcey, M., Berthelot, D.: Wikireading: A novel large-scale language understanding task over wikipedia. In: ACL. vol. 1, pp. 1535–1545 (2016)
13. Hill, F., Bordes, A., Chopra, S., Weston, J.: The goldilocks principle: Reading children’s books with explicit memory representations. *arXiv preprint arXiv:1511.02301* (2015)
14. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural computation* (1997)
15. Hu, M., Peng, Y., Qiu, X.: Reinforced mnemonic reader for machine comprehension. *CoRR*, abs/1705.02798 (2017)
16. Ji, H., Grishman, R.: Knowledge base population: Successful approaches and challenges. In: ACL. pp. 1148–1158 (2011)
17. Kim, Y., Jernite, Y., Sontag, D., Rush, A.M.: Character-aware neural language models. In: AAAI. pp. 2741–2749 (2016)
18. Lample, G., Ballesteros, M., Subramanian, S., Kawakami, K., Dyer, C.: Neural architectures for named entity recognition. In: Proceedings of NAACL-HLT. pp. 260–270 (2016)
19. Le, Q.V., Jaitly, N., Hinton, G.E.: A simple way to initialize recurrent networks of rectified linear units. *arXiv preprint arXiv:1504.00941* (2015)
20. Lee, T., Wang, Z., Wang, H., Hwang, S.w.: Attribute extraction and scoring: A probabilistic approach. In: 29th International Conference on Data Engineering. pp. 194–205 (2013)
21. Levy, O., Seo, M., Choi, E., Zettlemoyer, L.: Zero-shot relation extraction via reading comprehension. In: CoNLL. pp. 333–342 (2017)
22. Lin, Y., Liu, Z., Sun, M., Liu, Y., Zhu, X.: Learning entity and relation embeddings for knowledge graph completion. In: AAAI. vol. 15, pp. 2181–2187 (2015)
23. Ling, X., Weld, D.S.: Fine-grained entity recognition. In: AAAI (2012)
24. Liu, X., Shen, Y., Duh, K., Gao, J.: Stochastic answer networks for machine reading comprehension. *arXiv preprint arXiv:1712.03556* (2017)
25. Luo, G., Huang, X., Lin, C.Y., Nie, Z.: Joint entity recognition and disambiguation. In: EMNLP. pp. 879–888 (2015)
26. Ma, X., Hovy, E.: End-to-end sequence labeling via bi-directional lstm-cnns-crf. In: ACL. vol. 1, pp. 1064–1074 (2016)
27. Mintz, M., Bills, S., Snow, R., Jurafsky, D.: Distant supervision for relation extraction without labeled data. In: ACL. pp. 1003–1011 (2009)
28. Moro, A., Raganato, A., Navigli, R.: Entity linking meets word sense disambiguation: a unified approach. *TACL* 2, 231–244 (2014)
29. Pan, B., Li, H., Zhao, Z., Cao, B., Cai, D., He, X.: Memen: Multi-layer embedding with memory networks for machine comprehension. *arXiv preprint arXiv:1707.09098* (2017)

30. Pennington, J., Socher, R., Manning, C.: Glove: Global vectors for word representation. In: EMNLP. pp. 1532–1543 (2014)
31. Pyysalo, S., Ginter, F., Heimonen, J., Björne, J., Boberg, J., Järvinen, J., Salakoski, T.: Bioinfer: a corpus for information extraction in the biomedical domain. *BMC bioinformatics* **8**(1), 50 (2007)
32. Rajpurkar, P., Zhang, J., Lopyrev, K., Liang, P.: Squad: 100,000+ questions for machine comprehension of text. In: EMNLP. pp. 2383–2392 (2016)
33. Ren, X., Wu, Z., He, W., Qu, M., Voss, C.R., Ji, H., Abdelzaher, T.F., Han, J.: Cotype: Joint extraction of typed entities and relations with knowledge bases. In: Proceedings of the 26th International Conference on World Wide Web. pp. 1015–1024 (2017)
34. Riedel, S., Yao, L., McCallum, A.: Modeling relations and their mentions without labeled text. In: Joint European Conference on Machine Learning and Knowledge Discovery in Databases. pp. 148–163. Springer (2010)
35. Roth, B., Conforti, C., Poerner, N., Karn, S., Schütze, H.: Neural architectures for open-type relation argument extraction. arXiv preprint arXiv:1803.01707 (2018)
36. Schmitz, M., Bart, R., Soderland, S., Etzioni, O., et al.: Open language learning for information extraction. In: EMNLP. pp. 523–534 (2012)
37. Seo, M., Kembhavi, A., Farhadi, A., Hajishirzi, H.: Bidirectional attention flow for machine comprehension. arXiv preprint arXiv:1611.01603 (2016)
38. Shen, W., Wang, J., Han, J.: Entity linking with a knowledge base: Issues, techniques, and solutions. *IEEE Transactions on Knowledge and Data Engineering* **27**(2), 443–460 (2015)
39. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research* **15**(1), 1929–1958 (2014)
40. Srivastava, R.K., Greff, K., Schmidhuber, J.: Highway networks. arXiv preprint arXiv:1505.00387 (2015)
41. Stanovsky, G., Dagan, I.: Creating a large benchmark for open information extraction. In: EMNLP. pp. 2300–2305 (2016)
42. Sukhbaatar, S., Weston, J., Fergus, R., et al.: End-to-end memory networks. In: Advances in neural information processing systems. pp. 2440–2448 (2015)
43. Tjong Kim Sang, E.F., De Meulder, F.: Introduction to the conll-2003 shared task: Language-independent named entity recognition. In: Proceedings of NAACL-HLT. pp. 142–147 (2003)
44. Vinyals, O., Fortunato, M., Jaitly, N.: Pointer networks. In: Advances in Neural Information Processing Systems. pp. 2692–2700 (2015)
45. Vrandečić, D., Krötzsch, M.: Wikidata: a free collaborative knowledgebase. *Communications of the ACM* **57**(10), 78–85 (2014)
46. Wang, S., Jiang, J.: Machine comprehension using match-lstm and answer pointer. arXiv preprint arXiv:1608.07905 (2016)
47. Wang, W., Yang, N., Wei, F., Chang, B., Zhou, M.: Gated self-matching networks for reading comprehension and question answering. In: ACL. vol. 1, pp. 189–198 (2017)
48. Xu, K., Feng, Y., Huang, S., Zhao, D.: Semantic relation classification via convolutional neural networks with simple negative sampling. In: EMNLP. pp. 536–540 (2015)
49. Yahya, M., Whang, S., Gupta, R., Halevy, A.: Renoun: Fact extraction for nominal attributes. In: EMNLP. pp. 325–335 (2014)
50. Zeiler, M.D.: Adadelata: an adaptive learning rate method. arXiv preprint arXiv:1212.5701 (2012)
51. Zeng, D., Liu, K., Chen, Y., Zhao, J.: Distant supervision for relation extraction via piecewise convolutional neural networks. In: EMNLP. pp. 1753–1762 (2015)
52. Zheng, S., Wang, F., Bao, H., Hao, Y., Zhou, P., Xu, B.: Joint extraction of entities and relations based on a novel tagging scheme. In: ACL. vol. 1, pp. 1227–1236 (2017)