

# IMPLICIT REASONET: IMPLICITLY MODELING LARGE-SCALE STRUCTURED RELATIONSHIPS WITH SHARED MEMORY

Yelong Shen\*, Po-Sen Huang\*, Ming-Wei Chang, Jianfeng Gao

Microsoft Research, Redmond, WA, USA

{yeshen, pshuang, minchang, jfgao}@microsoft.com

## ABSTRACT

The knowledge base completion task has been as an important open research problem. To recover missing relationships in a knowledge base, a model needs to perform inference based on the relationships observed during training. In this paper, we propose a new neural model, Implicit Reasonet (IRN), to alleviate the needs to inject structured relationships between instances manually. Our model is designed to perform large-scale inference *implicitly* by using a multi-step procedure. The inference is guided by a *search controller* to access the memory that is *shared* across instances. Our proposed model outperforms all previous approaches on the popular FB15k and WN18 benchmarks by more than 5.7% and 1.0%, respectively.

## 1 INTRODUCTION

Knowledge bases such as WordNet (Fellbaum, 1998), Freebase (Bollacker et al., 2008), or Yago (Suchanek et al., 2007) contain many real-world facts expressed as triples, e.g. (*Bill Gates*, *FounderOf*, *Microsoft*). These knowledge bases are useful for many downstream applications such as question answering (Berant et al., 2013; Yih et al., 2015) and information extraction (Mintz et al., 2009). However, despite their formidable size, many important facts are still missing in the knowledge bases. For example, West et al. (2014) showed that 21% of the 100K most frequent PERSON entities have no known nationality in a recent version of Freebase. Thus, the link prediction or the knowledge base completion task (KBC) has become an important open research problem since Nickel et al. (2011). Formally, we seek to infer unknown relations based on the observed triples.

Neural-network based methods have been very popular for solving the KBC task. Following Bordes et al. (2013), one of the most popular approaches for the KBC task is to learn a vector-space representation of entities and relations during training, and then apply simple operations to infer the missing relations at test time. However, several recent papers demonstrate limitations of approaches relying upon vector-space models alone. By themselves, these models cannot capture the structured relationships between multiple triples adequately (Guu et al., 2015; Toutanova et al., 2016; Lin et al., 2015a). These papers each propose new ways to inject structured information at training time.

Existing neural KBC models, despite their strong empirical performance, often use very simple prediction functions. Typically, the probability of two entities that have a previously unknown relationship is proportional to a linear or bi-linear function of their corresponding vector or matrix representations. However, to make a correct prediction for a complicated instance, a multi-step search or inference procedure might be needed. For example, assume that we want to fill in the missing relation for the triple (*Obama*, *NATIONALITY*, ?), a multi-step search procedure can make use of the evidence in the observed triples such as (*Obama*, *BORNIN*, *Hawaii*) and (*Hawaii*, *PARTOF*, *U.S.A*). Unfortunately, due to the size of knowledge bases, performing search with neural models on top of all observed triples can be very costly, or heavy constraints needed to be applied to inference procedures and models.

In this paper, we take a different approach from prior work for the KBC task by addressing the challenges of performing large-scale inference through the design of *search controller* and *shared*

\*These authors contributed equally to this work.

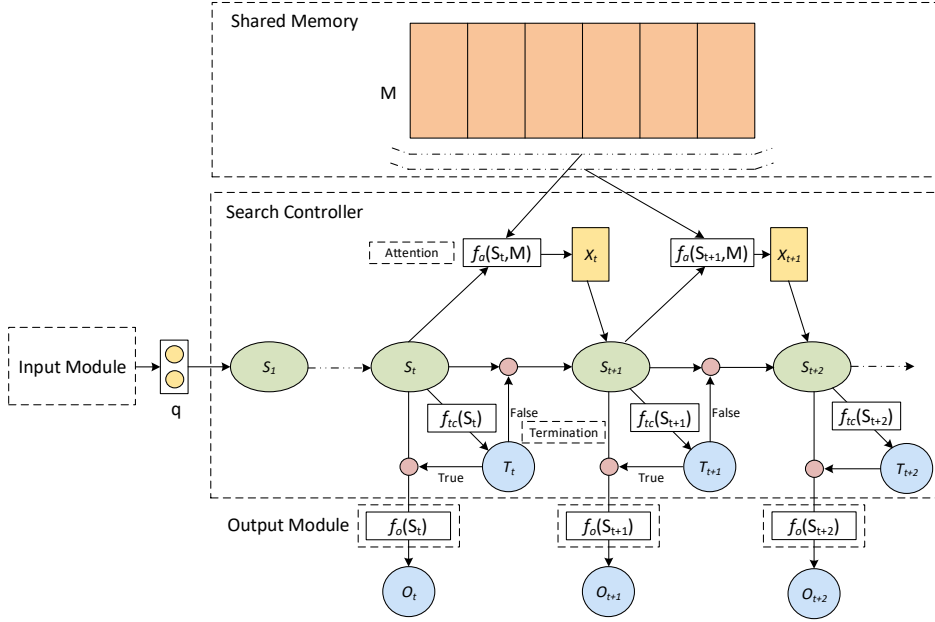


Figure 1: An IRN Architecture.

*memories*. Our inference procedure is controlled by the *search controller*, which only operates on the shared memories instead of on the set of observed triples. After receiving the representation from the *input module*, the search controller will interact with the *shared memory* and check the *termination gate* several times until the model feels confident enough to output the prediction. Intuitively, the shared memory is designed to store key information about the overall structures it learned during training, and hence the search controller only needs to access the shared memory instead of the observed triples.

There are several advantages of using IRNs. First, the cost of inference can be controlled because the search controller only needs to access the shared memory. Second, all the modules, including the search controller and memory, are jointly trained, and hence alleviate the needs to inject structured relationships between instances manually. Finally, we can easily extend IRNs to other tasks that require modeling structured relationships between instances, e.g. link prediction for social network, by switching the input and output modules.

The main contributions of our paper are as follows:

- We propose Implicit Reasoning Networks (IRNs), which use a shared memory guided by a search controller to model large-scale structured relationships implicitly.
- IRNs surpass prior state-of-the-art approach by 5.7% and 1.0% on the popular FB15k and WN18 benchmarks, respectively.
- We analyze the behavior of IRNs with a newly proposed generation task, shortest path synthesis. We show that IRNs outperform a standard sequence-to-sequence model, and are capable of performing meaningful multi-step inference.

## 2 REASONET FOR IMPLICIT INFERENCE

In this section, we first describe the general architecture of IRNs in a way that is agnostic to the KBC task. In an IRN, there are four main components, including an input component, an output component, a shared memory, and a search controller, as shown in Figure 1. The descriptions of these modules are as follow:

Input/Output Module: These two modules are task-dependent ones. The input module takes a query and converts the query into a vector representation  $q$ . The output module is a function  $f_o$ , which

converts the hidden state received from the search controller ( $s$ ) into an output  $O$ . We optimize the whole model using the output prediction  $O$  with respect to a ground-truth target using a task-specified loss function.

**Shared Memory:** The shared memory is denoted as  $M$ . It consists of a list of memory vectors,  $M = \{m_i\}_{i=1\dots I}$ , where  $m_i$  is a fixed dimensional vector. The memory vectors are randomly initialized and automatically updated through back-propagation. The shared memory component is shared across all instances.

**Search Controller:** The search controller is a recurrent neural network and controls the search process by keeping internal state sequences to track the current search process and history. The search controller uses an attention mechanism to fetch information from relevant memory vectors in  $M$ , and decides if the model should output the prediction or continue to generate the next possible output.

- **Internal State:** The internal state of the search controller is denoted as  $S$ , which is a vector representation of the search process. The initial state  $s_1$  is usually the vector representation of the input vector  $q$ . The internal state at  $t$ -th time step is represented by  $s_t$ . The sequence of internal states is modeled by an RNN:  $s_{t+1} = \text{RNN}(s_t, x_t; \theta_s)$ .
- **Attention to memory:** The attention vector  $x_t$  at  $t$ -th time step is generated based on the current internal state  $s_t$  and the shared memory  $M$ :  $x_t = f_{att}(s_t, M; \theta_x)$ . Specifically, the attention score  $a_{t,i}$  on a memory vector  $m_i$  given a state  $s_t$  is computed as  $a_{t,i} = \text{softmax}_{i=1,\dots,|M|} \lambda \cos(W_1 m_i, W_2 s_t)$ , where  $\lambda$  is set to 10 in our experiments and the weight matrices  $W_1$  and  $W_2$  are learned during training. The attention vector  $x_t$  can be written as  $x_t = f_{att}(s_t, M; \theta_x) = \sum_i^{|M|} a_{t,i} m_i$ .
- **Termination Control:** The terminate gate produces a stochastic random variable according to the current internal state,  $t_t \sim p(\cdot | f_{tc}(s_t; \theta_{tc}))$ .  $t_t$  is a binary random variable. If  $t_t$  is true, the IRN will finish the search process, and the output module will execute at time step  $t$ ; otherwise the IRN will generate the next attention vector  $x_{t+1}$  and feed into the state network to update the next internal state  $s_{t+1}$ . In our experiments, the termination variable is modeled by a logistical regression:  $f_{tc}(s_t; \theta_{tc}) = \text{sigmoid}(W_{tc} s_t + b_{tc})$ , where the weight matrix  $W_{tc}$  and bias vector  $b_{tc}$  are learned during training.

Compared IRNs to Memory Networks (MemNN) (Weston et al., 2014; Sukhbaatar et al., 2015) and Neural Turing Machines (NTM) (Graves et al., 2014; 2016), the biggest difference between our model and the existing frameworks is the search controller and the use of the shared memory. The search controller allows us to dynamically perform a multi-step inference depending on the complexity of the instance. MemNN and NTM explicitly store instance-specific inputs in a working memory, i.e., the memory vectors are reset for every instance. In contrast, in IRNs, the shared memory is a long-term memory and is updated per instance in training. Intuitively, the shared memory in IRNs, after training using task-specific objectives, contains the information that is useful to accomplish the task.

## 2.1 STOCHASTIC INFERENCE PROCESS

The inference process of an IRN is as follows. First, the model converts a task-dependent input to a vector representation through the input module. Then, the model uses the input representation to initialize the search controller. In every time step, the search controller determines whether the process is finished by sampling from the distribution according to the terminate gate. If the outcome is termination, the output module will generate a task-dependent prediction given the search controller states. If the outcome is continuation, the search controller will move on to the next time step, and create an attention vector based on the current search controller state and the shared memory. Intuitively, we design whole process by mimicking a search procedure that iteratively finds its target through a structure and output its prediction when a satisfied answer is found. The detailed inference process is described in Algorithm 1.

The inference process of an IRN is considered as a Partially Observable Markov Decision Process (POMDP) (Kaelbling et al., 1998) in the reinforcement learning (RL) literature. The IRN produces the output vector  $o_T$  at the  $T$ -th step, which implies termination gate variables  $t_{1:T} = (t_1 = 0, t_2 = 0, \dots, t_{T-1} = 0, t_T = 1)$ , and then takes prediction action  $p_T$  according to the probability distribution

**Algorithm 1:** Stochastic Inference Process in an IRN**Input :** Shared memory  $M$ ; Input vector  $q$ ; Maximum step  $T_{\max}$ **Output :** Output vector  $o$ 

- 1 Define  $s_1 = q$ ;  $t = 1$ ;
- 2 Sample  $t_t$  from the distribution  $p(\cdot|f_{tc}(s_t; \theta_{tc}))$ ;
- 3 if  $t_t$  is false, go to Step 4; otherwise Step 7;
- 4 Generate an attention vector  $x_t = f_{att}(s_t, M; \theta_x)$ ;
- 5 Update the internal state  $s_{t+1} = \text{RNN}(s_t, x_t; \theta_s)$ ;
- 6 Set  $t = t + 1$ ; if  $t < T_{\max}$  go to Step 2; otherwise Step 7;
- 7 Generate output  $o_t = f_o(s_t; \theta_o)$ ;
- 8 Return  $o = o_t$ ;

given  $o_T$ . Therefore, the IRN learns a stochastic policy  $\pi((t_{1:T}, p_T)|q; \theta)$  with parameters  $\theta$  to get a distribution over termination actions, and over prediction actions. The termination step  $T$  varies from instance to instance. The parameters of the IRN  $\theta$  are given by the parameters of the embedding matrices  $W$  for the input/output module, the shared memory  $M$ , the attention network  $\theta_x$ , the search controller RNN network  $\theta_s$ , the output generation network  $\theta_o$ , and the termination gate network  $\theta_{tc}$ . The parameters  $\theta = \{W, M, \theta_x, \theta_s, \theta_o, \theta_{tc}\}$  are trained by maximizing the total reward that the IRN could expect when interacting with the environment. The expected reward for an instance is defined as:

$$J(\theta) = \mathbb{E}_{\pi(t_{1:T}, p_T; \theta)} \left[ \sum_{t=1}^T r_t \right]$$

The reward can only be received at the final termination step when a prediction action  $p_T$  is performed. The rewards on intermediate steps are zeros,  $\{r_t = 0\}_{t=1 \dots T-1}$ .

We employ the approach from our previous work (Shen et al., 2016), REINFORCE (Williams, 1992) based Contrastive Reward method, to maximize the expected reward. The gradient of  $J$  can be written as:

$$\nabla_{\theta} J(\theta) = \sum_{(t_{1:T}, a_T) \in \mathbb{A}^{\dagger}} \pi(t_{1:T}, p_T; \theta) \left[ \nabla_{\theta} \log \pi(t_{1:T}, p_T; \theta) \left( \frac{r_T}{b^i} - 1 \right) \right]$$

where  $\mathbb{A}^{\dagger}$  is all the possible episodes, the baseline  $b^i = \sum_{(t_{1:T}, a_T) \in \mathbb{A}^{\dagger}} \pi(t_{1:T}, a_T; \theta) r_T$  is the expected reward on the  $|\mathbb{A}^{\dagger}|$  episodes for the  $i$ -th training instance.

### 3 APPLYING IRNS TO KNOWLEDGE BASE COMPLETION

The goal of KBC tasks (Bordes et al., 2013) is to predict a head or a tail entity given the relation type and the other entity, i.e. predicting  $h$  given  $(?, r, t)$  or predicting  $t$  given  $(h, r, ?)$ , where  $?$  denotes the missing entity. For a KBC task, the input to our model is a head or tail entity and a relation. The task-dependent input module first extracts the embedding vectors for the entity and relation from an embedding matrix. We then represent the query vector  $q$  for an IRN as the concatenation of the two vectors. For the task dependent output module, we use a nonlinear projection to project the search controller state into an output vector  $o$ :  $f_o(s_t; \theta_o) = \tanh(W_o s_t + b_o)$ , where the  $W_o$  and  $b_o$  are the weight matrix and bias vector, respectively. We define the ground truth target entity embedding as  $y$ , and use the  $L_1$  distance measure between the output  $o$  and target entity  $y$ , namely  $d(o, y) = |o - y|_1$ . We sample a set of incorrect entity embeddings  $N = \{y_i^-\}_{i=1}^{|N|}$  as negative examples. The probability of selecting a prediction  $\hat{y} \in D$  can be approximated as

$$p(\hat{y}|o) = \frac{\exp(-\gamma d(o, \hat{y}))}{\sum_{y_k \in D} \exp(-\gamma d(o, y_k))}$$

where  $D = N \cup \{y\}$ . We set  $|N|$  and  $\gamma$  to 20 and 5, respectively, for the experiments on FB15k and WN18 datasets. The IRN performs a prediction action  $p_T$  on selecting  $\hat{y}$  with probability  $p(\hat{y}|o)$ .

We define the reward of the prediction action as one if the ground truth entity is selected, and zero otherwise.

## 4 EXPERIMENTAL RESULTS

In this section, we evaluate the performance of our model on the benchmark FB15k and WN18 datasets for KBC tasks (Bordes et al., 2013). These datasets contain multi-relations between head and tail entities. Given a head entity and a relation, the model produces a ranked list of the entities according to the score of the entity being the tail entity of this triple. To evaluate the ranking, we report **mean rank (MR)**, the mean of rank of the correct entity across the test examples, and **hits@10**, the proportion of correct entities ranked in the top-10 predictions. Lower MR or higher hits@10 indicates a better prediction performance.

We use the same hyper-parameters of our model for both FB15k and WN18 data sets. Entity embeddings (which are not shared between input and output modules) and relation embedding are both 100-dimensions. There are 64 memory vectors with 200 dimensions each, initialized by random vectors with unit  $L_2$ -norm. We use single-layer GRU with 200 cells as the search controller. We set the maximum inference step of the IRN to 5. We randomly initialize all model parameters, and use SGD as the training algorithm with mini-batch size of 64. We set the learning rate to a constant number, 0.01. To prevent the model from learning a trivial solution by increasing entity embeddings norms, we follow Bordes et al. (2013) to enforce the  $L_2$ -norm of the entity embeddings as 1. We use hits@10 as the validation metric for the IRN. Following the work (Lin et al., 2015a), we add reverse relations into the training triplet set to increase the training data.

Following Nguyen et al. (2016), we divide the results of previous work into two groups. The first group contains the models that directly optimize a scoring function for the triples in a knowledge base without using extra information. The second group of models make uses of additional information from multi-step relations. For example, RTransE (García-Durán et al., 2015) and PTransE (Lin et al., 2015a) models are extensions of the TransE (Bordes et al., 2013) model by explicitly exploring multi-step relations in the knowledge base to regularize the trained embeddings. The NLFeat model (Toutanova et al., 2015) is a log-linear model that makes use of simple node and link features.

The results are in Table 1. According to the table, our model significantly outperforms previous baselines, regardless of whether previous approaches use additional information or not. Specifically, on FB15k, the MR of our model surpasses all previous results by 20, and our hit@10 outperforms others by 5%. On WN18, the IRN obtains the highest hit@10 while maintaining similar MR results compared to previous work.<sup>1</sup>

We analyze hits@10 results on FB15k with respect to the relation categories. Following the evaluation in Bordes et al. (2013), we evaluate the performance in four types of relation: 1-1 if a head entity can appear with at most one tail entity, 1-Many if a head entity can appear with many tail entities, Many-1 if multiple heads can appear with the same tail entity, and Many-Many if multiple head entities can appear with multiple tail entities. The detailed results are shown in Table 2. The IRN significantly improves the hits@10 results in the Many-1 category on predicting the head entity (18%), the 1-Many category on predicting the tail entity (16%), and the Many-Many category (over 10% in average).

To analyze the behavior of IRNs, we pick some examples for the tail entity prediction in Table 3. Interestingly, we observed that the model can gradually increase the ranking score of the correct tail entity during the inference process.

## 5 ANALYSIS: APPLYING IRNs TO A SHORTEST PATH SYNTHESIS TASK

We construct a synthetic task, shortest path synthesis, to evaluate the inference capability over a shared memory. The motivations of applying our model to this task are as follows. First, we want to test if IRNs can work on another task that also requires performing inference to get the correct

<sup>1</sup>Nguyen et al. (2016) reported two results on WN18, where the first one is obtained by choosing to optimize hits@10 on the validation set, and second one is obtained by choosing to optimize MR on the validation set. We list both of them in Table 1.

Table 1: The knowledge base completion (link prediction) results on WN18 and FB15k.

Model	Additional Information	WN18		FB15k	
		MR	Hits@10 (%)	MR	Hits@10 (%)
SE (Bordes et al., 2011)	NO	985	80.5	162	39.8
Unstructured (Bordes et al., 2014)	NO	304	38.2	979	6.3
TransE (Bordes et al., 2013)	NO	251	89.2	125	47.1
TransH (Wang et al., 2014)	NO	303	86.7	87	64.4
TransR (Lin et al., 2015b)	NO	225	92.0	77	68.7
CTransR (Lin et al., 2015b)	NO	218	92.3	75	70.2
KG2E (He et al., 2015)	NO	348	93.2	59	74.0
TransD (Ji et al., 2015)	NO	212	92.2	91	77.3
TATEC (García-Durán et al., 2015)	NO	-	-	58	76.7
NTN (Socher et al., 2013)	NO	-	66.1	-	41.4
DISTMULT (Yang et al., 2014)	NO	-	94.2	-	57.7
STransE (Nguyen et al., 2016)	NO	244 ( <b>206</b> )	94.7 (93)	69	79.7
RTransE (García-Durán et al., 2015)	Path	-	-	50	76.2
PTransE (Lin et al., 2015a)	Path	-	-	58	84.6
NLFeat (Toutanova et al., 2015)	Node + Link Features	-	94.3	-	87.0
<b>IRN</b>		249	<b>95.3</b>	<b>38</b>	<b>92.7</b>

Table 2: Hits@10 (%) in the relation category on FB15k. (M stands for Many)

Model	Predicting head $h$				Predicting tail $t$			
	1-1	1-M	M-1	M-M	1-1	1-M	M-1	M-M
SE (Bordes et al., 2011)	35.6	62.6	17.2	37.5	34.9	14.6	68.3	41.3
Unstructured (Bordes et al., 2014)	34.5	2.5	6.1	6.6	34.3	4.2	1.9	6.6
TransE (Bordes et al., 2013)	43.7	65.7	18.2	47.2	43.7	19.7	66.7	50.0
TransH (Wang et al., 2014)	66.8	87.6	28.7	64.5	65.5	39.8	83.3	67.2
TransR (Lin et al., 2015b)	78.8	89.2	34.1	69.2	79.2	37.4	90.4	72.1
CTransR (Lin et al., 2015b)	81.5	89.0	34.7	71.2	80.8	38.6	90.1	73.8
KG2E (He et al., 2015)	<b>92.3</b>	94.6	66.0	69.6	<b>92.6</b>	67.9	94.4	73.4
TransD (Ji et al., 2015)	86.1	95.5	39.8	78.5	85.4	50.6	94.4	81.2
TATEC (García-Durán et al., 2015)	79.3	93.2	42.3	77.2	78.5	51.5	92.7	80.7
STransE (Nguyen et al., 2016)	82.8	94.2	50.4	80.1	82.4	56.9	93.4	83.1
PTransE (Lin et al., 2015a)	91.0	92.8	60.9	83.8	91.2	74.0	88.9	86.4
<b>IRN</b>	87.2	<b>96.1</b>	<b>84.8</b>	<b>92.9</b>	86.9	<b>90.5</b>	<b>95.3</b>	<b>94.1</b>

answer. Second, we select the *sequence generation* task so that we are able to analyze the inference capability of IRNs in details.

In the shortest path synthesis task, as illustrated in Figure 2, a training instance consists of a start node and an end node (e.g.,  $215 \rightsquigarrow 493$ ) of an underlying weighted directed graph that is unknown to models. The output of each instance is the shortest path between the given start and end nodes of the underlying graph (e.g.,  $215 \rightarrow 101 \rightarrow 493$ ). Specifically, models can only observe the start-end node pairs as input and their shortest path as output. The whole graph is unknown to the models and the edge weights are not revealed in the training data. At test time, a path sequence is considered correct if it connects the start node and the end node of the underlying graph, and the cost of the predicted path is the same as the optimal path.

Note that the task is very difficult and *cannot* be solved by dynamic programming algorithms since the weights on the edges are not revealed to the algorithms or the models. To recover some of the shortest paths at the test time, the model needs to infer the correct path from the observed instances. For example, assume that we observe two instances in the training data, “ $A \rightsquigarrow D: A \rightarrow B \rightarrow G \rightarrow D$ ” and “ $B \rightsquigarrow E: B \rightarrow C \rightarrow E$ ”. In order to answer the shortest path between  $A$  and  $E$ , the model needs to infer that “ $A \rightarrow B \rightarrow C \rightarrow E$ ” is a possible path between  $A$  and  $E$ . If there are multiple possible paths, the model has to decide which one is the shortest one using statistical information.

Table 3: Test examples in FB15k dataset, given a head entity and a relation, the IRN predicts the tail entity with multiple search steps.

<b>Input:</b> (Dean Koontz, /PEOPLE/PERSON/PROFESSION)						
<b>Target:</b> Film Producer						
Step	Termination Prob.	Rank	Predict top-3 entities			
1	0.018	9	Author	TV. Director	Songwriter	
2	0.052	7	Actor	Singer	Songwriter	
3	0.095	4	Actor	Singer	Songwriter	
4	0.132	4	Actor	Singer	Songwriter	
5	0.702	3	Actor	Singer	<b>Film Producer</b>	

<b>Input:</b> (1964 Summer Olympics, /OLYMPICS/OLYMPIC_GAMES/SPORTS)						
<b>Target:</b> Judo						
Step	Termination Prob.	Rank	Predict top-3 entities			
1	0.018	498	Rugby union	Golf	Curling	
2	0.079	4	Alpine skiing	Tennis	Archery	
3	0.164	3	Alpine skiing	Tennis	<b>Judo</b>	
4	0.215	1	<b>Judo</b>	Alpine skiing	Archery	
5	0.524	1	<b>Judo</b>	Alpine skiing	Archery	

In the experiments, we construct a graph with 500 nodes and we randomly assign two nodes to form an edge. We split 20,000 instances for training, 10,000 instances for validation, and 10,000 instances for testing. We create the training and testing instances carefully so that the model needs to perform inference to recover the correct path. We describe the details of the graph and data construction parts in the appendix section. A sub-graph of the data is shown in Figure 2.

For the settings of the IRN, we switch the output module to a GRU decoder for a sequence generation task. We assign reward  $r_T = 1$  if all the prediction symbols are correct and 0 otherwise. We use a 64-dimensional embedding vector for input symbols, a GRU controller with 128 cells, and a GRU decoder with 128 cells. We set the maximum inference step  $T_{\max}$  to 5.

We compare the IRN with two baseline approaches: dynamic programming without edge-weight information and a standard sequence-to-sequence model (Sutskever et al., 2014). Without knowing the weights from the data, dynamic programming cannot perfectly infer the correct path and recovers only 589 correct paths during test time. The sequence-to-sequence model with a similar parameter size recovers 904 correct paths. The IRN outperforms both baselines and recovers 1,319 paths. 76.9% of the predicted paths from IRN are *valid* paths, where a path is valid if the path connects the start and end node nodes of the underlying graph. In contrast, only 69.1% of the predicted paths from the sequence-to-sequence model are valid.

To further understand the inference process of the IRN, Figure 2 shows the inference process of a test instance. Interestingly, to make the correct prediction on this instance, the model has to perform a fairly complicated inference.<sup>2</sup> We observe that the model cannot find a connected path in the first three steps. Finally, the model finds a valid path at the forth step and predict the correct shortest path sequence at the fifth step.

## 6 RELATED WORK

**Link Prediction and Knowledge Base Completion** Given that  $r$  is a relation,  $h$  is the head entity, and  $t$  is the tail entity, most of the embedding models for link prediction focus on finding the scoring function  $f_r(h, t)$  that represents the implausibility of a triple. (Bordes et al., 2011; 2014; 2013; Wang et al., 2014; Ji et al., 2015; Nguyen et al., 2016). In most studies, the scoring function  $f_r(h, t)$  is often linear or bi-linear. For example, in TransE (Bordes et al., 2013), the function is implemented as  $f_r(h, t) = \|\mathbf{h} + \mathbf{r} - \mathbf{t}\|$ , where  $\mathbf{h}$ ,  $\mathbf{r}$  and  $\mathbf{t}$  are the corresponding vector representations.

<sup>2</sup> In the example, to find the right path, the model needs to search over observed instances “215  $\rightsquigarrow$  448: 215  $\rightarrow$  101  $\rightarrow$  448” and “76  $\rightsquigarrow$  493: 76  $\rightarrow$  308  $\rightarrow$  101  $\rightarrow$  493”, and to figure out the distance of “140  $\rightarrow$  493” is longer than “101  $\rightarrow$  493” (there are four shortest paths between 101  $\rightarrow$  493 and three shortest paths between 140  $\rightarrow$  493 in the training set).

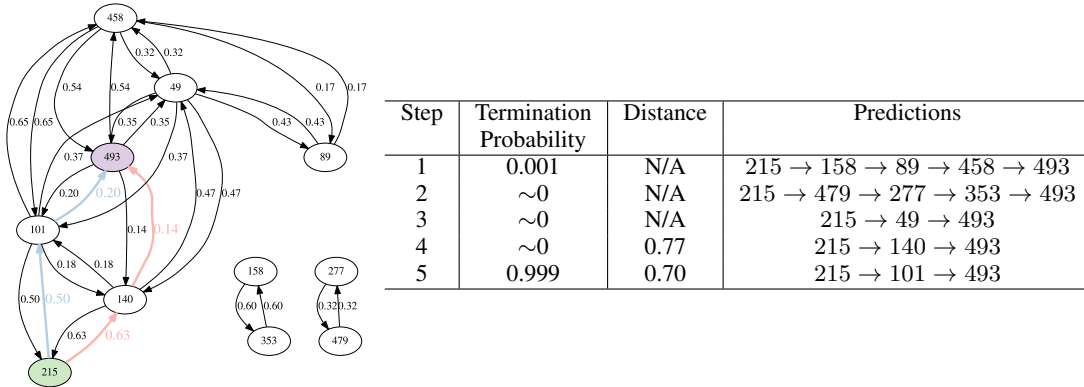


Figure 2: An example of the shortest path synthesis dataset, given an input “215  $\rightsquigarrow$  493” (Answer: 215 → 101 → 493). Note that we only show the nodes that are related to this example here. The corresponding termination probability and prediction results are shown in the table. The model terminates at step 5.

Recently, different studies (Guu et al., 2015; Lin et al., 2015a; Toutanova et al., 2016) demonstrate the importance for models to also learn from multi-step relations. Learning from multi-step relations injects the structured relationships between triples into the model. However, this also poses a technical challenge of considering exponential numbers of multi-step relationships. Prior approaches address this issue by designing a path-mining algorithm (Lin et al., 2015a) or considering all possible paths using a dynamic programming algorithm with the restriction of using linear or bi-linear models only (Toutanova et al., 2016). Toutanova & Chen (2015) shows the effectiveness of using simple node and link features that encode structured information on FB15k and WN18. In our work, the IRN outperforms prior results and shows that similar information can be captured by the model without explicitly designing features.

Knowledge base completion is more than just link prediction. Studies such as (Riedel et al., 2013) show that incorporating textual information can further improve the knowledge base completion tasks. It would be interesting to incorporate the information outside the knowledge bases in our model in the future.

**Neural Frameworks** Sequence-to-sequence models (Sutskever et al., 2014; Cho et al., 2014) have shown to be successful in many applications such as machine translation and conversation modeling (Sordoni et al., 2015). While sequence-to-sequence models are powerful, recent work has shown that the necessity of incorporating an external memory to perform inference in simple algorithmic tasks (Graves et al., 2014; 2016).

We are inspired by our previous work (Shen et al., 2016) for using a search controller module to perform iteratively search over a given paragraph until a satisfied answer is found. In contrast, in this work, IRNs use shared memory along with the search controller to perform inference over large-scale structured relationships.

## 7 CONCLUSION

In this paper, we propose Implicit Reasoning Networks (IRNs), which perform inference over a shared memory that models large-scale structured relationships implicitly. The inference process is guided by a search controller to access the memory that is shared across instances. We demonstrate and analyze the multi-step inference capability of IRNs in the knowledge base completion tasks and a shortest path synthesis task. Our model, without using any explicit knowledge base information, outperforms all prior approaches on the popular FB15k and WN18 benchmarks by more than 5.7% and 1.0%, respectively. For future work, we aim to further develop IRNs for modeling the relationships from unstructured data in natural language, and exploit ways to generate human understandable reasoning interpretation from the shared memory.

For future work, we aim to further extend IRNs in two ways. First, inspired from Ribeiro et al. (2016), we would like to develop techniques to exploit ways to generate human understandable reasoning



interpretation from the shared memory. Second, we want to apply IRNs to infer the relationships in unstructured data such as natural language. One important application is to capture the common sense knowledge. For example, when the model receives a natural language query “are rabbits animals?”, the model can infer the answer implicitly in the shared memory without performing inference directly on top of huge amount of observed sentences such as “all mammals are animals” and “rabbits are animals”. We believe the ability to perform inference implicitly is crucial for modeling large-scale structured relationships.

#### ACKNOWLEDGMENTS

We thank Scott Wen-Tau Yih and Zachary Lipton for their thoughtful feedback and discussions.

#### REFERENCES

- Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. Semantic parsing on Freebase from question-answer pairs. In *Proceedings of EMNLP*, 2013.
- Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. Freebase: A collaboratively created graph database for structuring human knowledge. In *Proceedings of SIGMOD-08*, pp. 1247–1250, 2008.
- Antoine Bordes, Jason Weston, Ronan Collobert, and Yoshua Bengio. Learning structured embeddings of knowledge bases. In *Proceedings of the Twenty-Fifth AAAI Conference on Artificial Intelligence*, pp. 301–306, 2011.
- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In *Advances in Neural Information Processing Systems*, pp. 2787–2795, 2013.
- Antoine Bordes, Xavier Glorot, Jason Weston, and Yoshua Bengio. A semantic matching energy function for learning with multi-relational data. *Machine Learning*, 94(2):233–259, 2014.
- Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- C. Fellbaum. *WordNet: An Electronic Lexical Database*. MIT Press, 1998.
- Alberto García-Durán, Antoine Bordes, and Nicolas Usunier. Composing relationships with translations. In *EMNLP*, pp. 286–290, 2015.
- Alberto García-Durán, Antoine Bordes, Nicolas Usunier, and Yves Grandvalet. Combining two and three-way embeddings models for link prediction in knowledge bases. *CoRR*, abs/1506.00999, 2015.
- Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. *arXiv preprint arXiv:1410.5401*, 2014.
- Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka Grabska-Barwińska, Sergio Gómez Colmenarejo, Edward Grefenstette, Tiago Ramalho, John Agapiou, et al. Hybrid computing using a neural network with dynamic external memory. *Nature*, 2016.
- Kelvin Guu, John Miller, and Percy Liang. Traversing knowledge graphs in vector space. *arXiv preprint arXiv:1506.01094*, 2015.
- Shizhu He, Kang Liu, Guoliang Ji, and Jun Zhao. Learning to represent knowledge graphs with gaussian embedding. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, pp. 623–632, 2015.
- Guoliang Ji, Shizhu He, Liheng Xu, Kang Liu, and Jun Zhao. Knowledge graph embedding via dynamic mapping matrix. In *ACL*, pp. 687–696, 2015.
- Leslie Pack Kaelbling, Michael L. Littman, and Anthony R. Cassandra. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101:99–134, 1998.
- Yankai Lin, Zhiyuan Liu, Huanbo Luan, Maosong Sun, Siwei Rao, and Song Liu. Modeling relation paths for representation learning of knowledge bases. In *Proceedings of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2015a.

- Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. Learning entity and relation embeddings for knowledge graph completion. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, AAAI'15, pp. 2181–2187, 2015b.
- Mike Mintz, Steven Bills, Rion Snow, and Daniel Jurafsky. Distant supervision for relation extraction without labeled data. In *Proceedings of ACL-IJCNLP-09*, pp. 1003–1011, 2009.
- Dat Quoc Nguyen, Kairit Sirts, Lizhen Qu, and Mark Johnson. STransE: a novel embedding model of entities and relationships in knowledge bases. In *NAACL*, pp. 460–466, 2016.
- Maximilian Nickel, Volker Tresp, and Hans-Peter Kriegel. A three-way model for collective learning on multi-relational data. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp. 809–816, 2011.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier. In *KDD*, 2016.
- Sebastian Riedel, Limin Yao, Andrew McCallum, and Benjamin M. Marlin. Relation extraction with matrix factorization and universal schemas. In *HLT-NAACL*, pp. 74–84, 2013.
- Yelong Shen, Po-Sen Huang, Jianfeng Gao, and Weizhu Chen. Reasonet: Learning to stop reading in machine comprehension. *CoRR*, abs/1609.05284, 2016.
- Richard Socher, Danqi Chen, Christopher D. Manning, and Andrew Y. Ng. Reasoning With Neural Tensor Networks For Knowledge Base Completion. In *Advances in Neural Information Processing Systems*, 2013.
- Alessandro Sordani, Michel Galley, Michael Auli, Chris Brockett, Yangfeng Ji, Margaret Mitchell, Jian-Yun Nie, Jianfeng Gao, and Bill Dolan. A neural network approach to context-sensitive generation of conversational responses. *arXiv preprint arXiv:1506.06714*, 2015.
- F. M. Suchanek, G. Kasneci, and G. Weikum. Yago: A Core of Semantic Knowledge. In *WWW*, 2007.
- Sainbayar Sukhbaatar, Jason Weston, Rob Fergus, et al. End-to-end memory networks. In *Advances in Neural Information Processing Systems*, pp. 2440–2448, 2015.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, pp. 3104–3112, 2014.
- Kristina Toutanova and Danqi Chen. Observed versus latent features for knowledge base and text inference. In *Proceedings of the 3rd Workshop on Continuous Vector Space Models and their Compositionality*, pp. 57–66, 2015.
- Kristina Toutanova, Danqi Chen, Patrick Pantel, Hoifung Poon, Pallavi Choudhury, and Michael Gamon. Representing text for joint embedding of text and knowledge bases. In *EMNLP*, 2015.
- Kristina Toutanova, Xi Victoria Lin, Scott Wen tau Yih, Hoifung Poon, and Chris Quirk. Compositional learning of embeddings for relation paths in knowledge bases and text. In *ACL*, 2016.
- Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. Knowledge graph embedding by translating on hyperplanes. In *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence*, pp. 1112–1119, 2014.
- Robert West, Evgeniy Gabrilovich, Kevin Murphy, Shaohua Sun, Rahul Gupta, and Dekang Lin. Knowledge base completion via search-based question answering. In *Proceedings of the 23rd international conference on World Wide Web*, pp. 515–526. ACM, 2014.
- Jason Weston, Sumit Chopra, and Antoine Bordes. Memory networks. *arXiv preprint arXiv:1410.3916*, 2014.
- Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8:229–256, 1992.
- Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding entities and relations for learning and inference in knowledge bases. *CoRR*, abs/1412.6575, 2014.
- Wen-tau Yih, Ming-Wei Chang, Xiaodong He, and Jianfeng Gao. Semantic parsing via staged query graph generation: Question answering with knowledge base. In *Proc. of ACL*, 2015.

## A DETAILS OF THE GRAPH CONSTRUCTION FOR THE SHORTEST PATH SYNTHESIS TASK

We construct the underlying graph as follows: on a three-dimensional unit-sphere, we randomly generate a set of nodes. For each node, we connect its  $K$ -nearest neighbors and use the euclidean distance between two nodes to construct a graph. We randomly sample two nodes and compute its shortest path if it is connected between these two nodes. Given the fact that all the sub-paths within a shortest path are shortest paths, we incrementally create the dataset and remove the instances which are a sub-path of previously selected paths or are super-set of previous selected paths. In this case, all the shortest paths can not be answered through directly copying from another instance. In addition, all the weights in the graph are hidden and not shown in the training data, which increases the difficulty of the tasks. We set  $k = 50$  as a default value.