

# Recurrent One-Hop Predictions for Reasoning over Knowledge Graphs

Wenpeng Yin<sup>1</sup>, Yadollah Yaghoobzadeh<sup>2</sup>, Hinrich Schütze<sup>3</sup>

<sup>1</sup>Dept. of Computer Science, University of Pennsylvania, USA

<sup>2</sup>Microsoft Research, Montreal, Canada

<sup>3</sup>CIS, LMU Munich, Germany

wenpeng@seas.upenn.edu

## Abstract

Large scale knowledge graphs (KGs) such as Freebase are generally incomplete. Reasoning over multi-hop (mh) KG paths is thus an important capability that is needed for question answering or other NLP tasks that require knowledge about the world. mh-KG reasoning includes diverse scenarios, e.g., given a head entity and a relation path, predict the tail entity; or given two entities connected by some relation paths, predict the unknown relation between them. We present *ROPs*, recurrent one-hop predictors, that predict entities at each step of mh-KB paths by using recurrent neural networks and vector representations of entities and relations, with two benefits: (i) modeling mh-paths of arbitrary lengths while updating the entity and relation representations by the training signal at each step; (ii) handling different types of mh-KG reasoning in a unified framework. Our models show state-of-the-art for two important multi-hop KG reasoning tasks: Knowledge Base Completion and Path Query Answering.<sup>1</sup>

## 1 Introduction

Natural language understanding (NLU) is impossible without knowledge about the world. Large scale knowledge graphs (KGs) such as Freebase (Bollacker et al., 2008) are structures that store world knowledge. Unfortunately, KGs suffer from incomplete coverage (Min et al., 2013); e.g., Freebase contains Brandon Lee, but not his ethnicity.

The knowledge in KGs needs to be expanded to cover more facts; reasoning is one way to do so. For example, we could infer that (*Microsoft*, ?, *United States*) instantiates “CountryOfHQ” given the facts (*Microsoft*, *IsBasedIn*, *Seattle*) and (*Seattle*, *LocatedIn*, *United States*); or we could infer Brandon Lee’s ethnicity from his parents’ ethnicity, i.e., answering the query (*Brandon Lee*, *Ethnicity*, ?) by facts (*Brandon Lee*, *Father*, *Bruce Lee*) and (*Bruce Lee*, *Ethnicity*, *Chinese*). We refer to the two reasoning examples as “knowledge base completion (KBC)” and “path query answering (PQA)”, respectively.

The most successful approach for modeling KGs is the *embedding approach*. It embeds KG elements (entities and relations) into low-dimensional dense vectors; controlling the dimensionality of the vector space forces generalization to new facts (Nickel et al., 2011).

In this work, we are mainly interested in three issues. (i) Compared to modeling one-hop KG paths, a bigger challenge is how to model multi-hop paths, e.g., the path query (*U.S.A*, *president*→*spouse*→*born\_in*, ?) for the question “Where was the first lady of the United States born?” (ii) How can we address different KG reasoning problems driven by multi-hop paths in a universal paradigm rather than via different systems? (iii) How can we combine specific multi-hop KG reasoning tasks with generic KG representation learning, so that KG representation learning can either stand alone or be incorporated into diverse multi-hop KG-related NLU problems. We get inspiration from following two types of work.

First, *prior work in KG reasoning*. Guu et al. (2015) extend one-hop reasoning regimes such as TransE (Bordes et al., 2013) to multi-hop PQA. However, these basic one-hop models do not encode the relation

<sup>1</sup><https://github.com/yinwenpeng/KBPath>

This work is licensed under a Creative Commons Attribution 4.0 International License. License details: <http://creativecommons.org/licenses/by/4.0/>.

order when used in compositional training schemes. For example, path query  $q_1 = (h, r_1 \rightarrow r_2 \rightarrow \dots \rightarrow r_k, ?)$  will be encoded into the same embedding as path query  $q_2 = (h, r_2 \rightarrow r_1 \rightarrow \dots \rightarrow r_k, ?)$ , resulting in (often incorrect) prediction of the same tail entity. Instead, the relation order should influence the prediction. This limitation in modeling multi-hop relation paths motivates the *RNN approach*: using recurrent neural networks (RNN (Elman, 1990)) to model relation paths (Neelakantan et al., 2015). Das et al. (2017) further extend this approach by incorporating entity information and apply it to multi-hop KBC. Intermediate entities should influence the reasoning decision. For example, given two paths with the same relation sequence: (Donald Trump, *child*, Ivanka Trump, *mother*, Ivana Trump) and (Donald Trump, *child*, Barron Trump, *mother*, Melania Trump), even though both paths have the relation sequence [*child*, *mother*], the relation between (Donald Trump, Melania Trump) is “spouse” while it does not hold between (Donald Trump, Ivana Trump) due to the intermediate entities: “Ivanka Trump” vs. “Barron Trump”. Similarly, paths (JFK, *located\_in*, NYC, *located\_in*, NY) and (Yankee Stadium, *located\_in*, NYC, *located\_in*, NY) would predict the same score for target relation “airport\_serves\_place” if we do not consider that Yankee Stadium is not an airport (Das et al., 2017).

Second, *sequence labeling tasks* such as POS tagging, chunking and NER have been successfully addressed by RNNs (Huang et al., 2015; Lample et al., 2016). These approaches model the mechanism in a structure of form “input<sub>1</sub>, tag<sub>1</sub>, input<sub>2</sub>, tag<sub>2</sub>, input<sub>3</sub>, tag<sub>3</sub>, ..., input<sub>t</sub>, tag<sub>t</sub>”, which resembles the structure of multi-hop KG paths.

Inspired by this prior work, we propose *ROP*, Recurrent One-hop Predictor. Given a head entity, ROP encodes a multi-hop sequence of relations and predicts a sequence of entities using an RNN. More formally, given relation sequence “ $r_1, r_2, \dots, r_t$ ” and the head entity  $e_h$ , ROP predicts the sequence  $e_1, e_2, \dots, e_t$ , thus generating a complete KG path  $e_h, r_1, e_1, r_2, e_2, \dots, r_t, e_t$ . Intuitively, our model memorizes history of path context; given a new relation, it predicts the next entity, then the memory is updated, and the process – given new relation, predicting new entity, updating memory – keeps going. Grouping step-wise updates in a chain gives our model two advantages. (i) A better vector space representation of KG entities and relations, with training signals either from in-path entities or from labels of reasoning tasks or from both. (ii) A unified approach for two different reasoning tasks (mh-KBC and mh-PQA) and state-of-the-art in each.

In summary, our contributions are: (i) ROP, a novel RNN sequence modeling of multi-hop KG paths that updates entity and relation embeddings by training signal at each step and leads to better KG embeddings; (ii) unified framework for solving different multi-hop reasoning tasks over KGs; (iii) showing the importance of modeling within-path entities in mh-PQA; (iv) state-of-the-art results for both mh-KBC and mh-PQA; (v) releasing an enhanced version of a mh-PQA dataset by adding within-path entities.

§2 introduces the mh-KBC and mh-PQA tasks. §3 discusses related work. §4 presents three ROP architectures and §5 evaluates them. §6 concludes.

## 2 Multi-Hop Path Reasoning Tasks

We first give background on the two KG reasoning tasks we address in this work.

**Knowledge Base Completion (mh-KBC).** In mh-KBC, the goal is to predict new relations between entities using existing path connections. For example,  $(A, \textit{LivesIn}, B)$  is implied, with some probability, from  $(A, \textit{CEO}, X)$  and  $(X, \textit{HQIn}, B)$ . This kind of reasoning lets us infer new or missing facts from KGs.

Figure 1(a) shows two paths between *Microsoft* and *United States*: (*Microsoft*, *IsBasedIn*, *Seattle*, *IsLocatedIn*, *Washington*, *IsLocatedIn*, *United States*) (blue) and (*Microsoft*, *CEO*, *Satya Nadella*, *PlaceOfBirth*, *India*, *LargestTradingPartner*, *United States*) (red). The task is then to predict the direct relation that connects *Microsoft* and *United States*; i.e., *CountryOfHQ* in this case. There can exist multiple long paths between two entities; the example shows that the target relation may only be inferrable from one path. The difficulty of finding the most informative path makes this task challenging.

**Path Query Answering (mh-PQA).** In mh-PQA, the goal is to *predict missing properties* of an entity, such as the earlier mentioned ethnicity of Brandon Lee. More generally, given an entity and relations of interest, predict what the target entity is, as Figure 1(b) shows. This task corresponds to answering

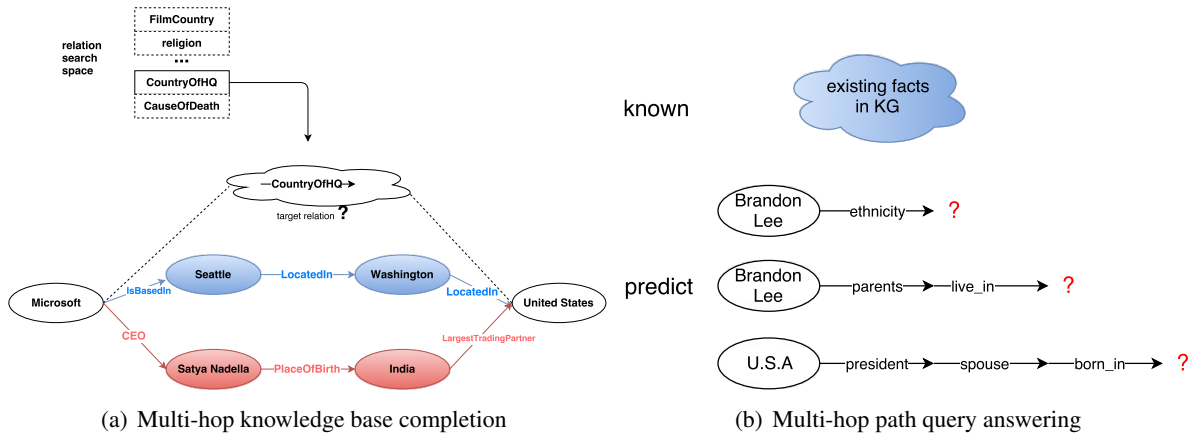


Figure 1: Multi-Hop Path Reasoning Tasks

compositional natural questions. For example, the question “Where do Brandon Lee’s parents live?” can be formulated by the path query `brandon_lee/parents/live_in`. mh-PQA tries to find answers to the path queries and hence compositional questions. Unfortunately, KGs often have missing facts (edges), which makes mh-PQA a non-trivial problem.

A path query  $q_t$  consists of an initial anchor entity,  $e_h$ , followed by a sequence of  $t$  relations to be traversed,  $p = (r_1, \dots, r_t)$ . Following (Guu et al., 2015), the answer or denotation of the query is the set of all entities that can be reached from  $e_h$  by traversing  $p$ .

### 3 Related Work

Here we focus on the multi-hop path reasoning literature. Some work (Neelakantan et al., 2015; Guu et al., 2015; Lin et al., 2015; Lin et al., 2016; Shen et al., 2016) does some composition over relation paths. Given relation path  $p = (r_1, \dots, r_t)$ , the composition operation is *add* ( $\mathbf{p} = \mathbf{r}_1 + \dots + \mathbf{r}_t$ ), *multiplication* ( $\mathbf{p} = \mathbf{r}_1 \cdot \dots \cdot \mathbf{r}_t$ ) or an *RNN step*:  $\mathbf{p}_i = \text{RNN}(\mathbf{p}_{i-1}, \mathbf{r}_i)$ , where  $\mathbf{p}_i$  is the accumulated relation information up to step  $i$ . Some work explores compositional encoding of long paths (Lin et al., 2016; Shen et al., 2016), but still performs reasoning in one-hop scenario. Neelakantan et al. (2015) use RNNs to model multi-hop paths.

Das et al. (2017) extend the RNN approaches by leveraging within-path entities into the encoding of inputs along with relations. We also include within-path entities, but we do not give them as inputs; instead, *we force our RNN to predict them as outputs and do updates at each step in the path*. This supports representation learning for KG entities and relations even without task-specific annotations. Toutanova et al. (2016) propose a dynamic programming algorithm to model both relation types and intermediate entities in the compositional path representations and test on WordNet and a biomedical KG. These two works address mh-KBC; for mh-PQA, there is no prior work on using within-path entities<sup>2</sup>, including Das et al. (2017), in which the system uses within-path entities as input, while those entities are not available for testing. So Das et al. (2017) use RNN for mh-PQA to encode the relation sequence, but it does not incorporate the intermediate entities involved.

### 4 Recurrent One-Hop Prediction

We propose three ROPs, recurrent one-hop predictors, to model paths such as  $p = (e_h, r_1, e_1, \dots, r_i, e_i, \dots, r_t, e_t)$ . Entities and relations appear alternately in  $p$ ;  $e_h$  and  $e_t$  are head and tail entities;  $t$  steps (hops) connect  $e_h$  and  $e_t$ ; each step is a single fact triple  $(e_{i-1}, r_i, e_i)$ . We stipulate  $\mathbf{e}_i \in \mathbb{R}^{d_e}$  and  $\mathbf{r}_i \in \mathbb{R}^{d_r}$ . We allow  $d_e \neq d_r$ .

<sup>2</sup>Guu et al. (2015) attempt to measure how severe the cascading errors along the path are by reconstructing the intermediate entities along a path. Their operation only generates an evaluation score for the path representation. We instead turn this into a training objective to fine-tune the path representations.

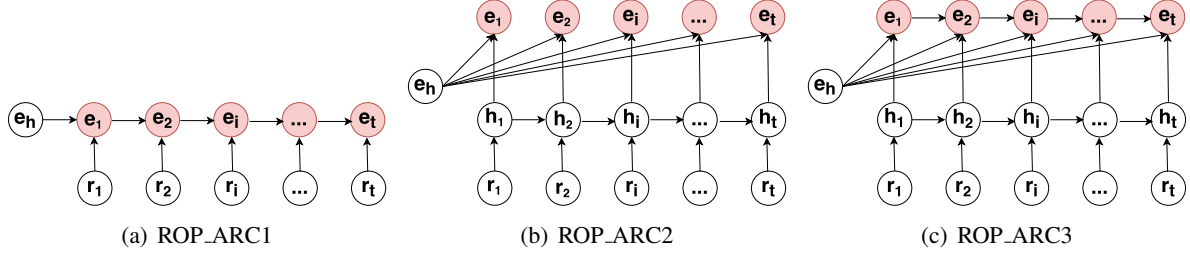


Figure 2: Three recurrent one-hop predictors

Paths are encoded by GRUs (Cho et al., 2014):

$$\mathbf{g}_z = \sigma(\mathbf{x}_i \mathbf{U}^z + \mathbf{h}_{i-1} \mathbf{W}^z) \quad (1)$$

$$\mathbf{g}_r = \sigma(\mathbf{x}_i \mathbf{U}^r + \mathbf{h}_{i-1} \mathbf{W}^r) \quad (2)$$

$$\hat{\mathbf{h}}_i = \tanh(\mathbf{x}_i \mathbf{U}^q + (\mathbf{h}_{i-1} \circ \mathbf{g}_r) \mathbf{W}^q) \quad (3)$$

$$\mathbf{h}_i = (1 - \mathbf{g}_z) \circ \hat{\mathbf{h}}_i + \mathbf{g}_z \circ \mathbf{h}_{i-1} \quad (4)$$

where  $x$  is the input sequence with  $x_i$  at position  $i$ ,  $h$  is the output sequence with  $h_i$  at position  $i$ .  $g_z$  and  $g_r$  are gates. All  $\mathbf{U}$ s and  $\mathbf{W}$ s are parameters. In following, we define the whole Eqs. 1–4 as a single GRU step as:

$$h_i = \text{GRU}(h_{i-1}, x_i) \quad (5)$$

Thus, we interpret each GRU step as a composition function of two objects:  $h_{i-1}$  and  $x_i$ .

We now introduce the three architectures ROP\_ARC1, ROP\_ARC2 and ROP\_ARC3 that encode the context “ $e_h, r_1, e_1, \dots, r_i$ ” in different ways to predict entity  $e_i$ .

**ROP\_ARC1** (Figure 2(a)) models KG paths as:

$$\hat{\mathbf{e}}_0 = \mathbf{e}_h \quad (6)$$

$$\hat{\mathbf{e}}_i = \text{GRU}(\hat{\mathbf{e}}_{i-1}, \mathbf{r}_i) \quad (7)$$

$\hat{\mathbf{e}}_0$  is initialized to the embedding of the *true* head entity  $\mathbf{e}_h$ . At position  $i, i > 0$ ,  $\hat{\mathbf{e}}_i$  is the *predicted* entity embedding.

ROP\_ARC1 is essentially a recurrent process with a pre-set starting state; the key is to use the head entity embedding  $\mathbf{e}_h$  as the initialization of the hidden state of GRU. We hope this start point guides where the path goes and what state to reach at each position. As a result, *relations lie in the input space and entities lie in the hidden space*. All *predicted* entities  $\hat{\mathbf{e}}_i$  will be compared with the gold intermediate entities  $\mathbf{e}_i$ , then the loss (red in Figure 2) is used to train the system.

ROP\_ARC1 only encodes head entity  $\mathbf{e}_h$  at the starting hidden state, possibly far from the tail entity  $\mathbf{e}_t$  for long paths. Thus,  $\mathbf{e}_h$  cannot provide effective guidance for prediction of  $\mathbf{e}_t$ . This motivates ROP\_ARC2, a modification of ROP\_ARC1.

In **ROP\_ARC2**, we want the head entity  $\mathbf{e}_h$  to participate in the entity prediction at each step more directly and effectively. To this end, ROP\_ARC2 first encodes the relation sequence as standard GRU:

$$\mathbf{h}_0 = \mathbf{0} \quad (8)$$

$$\mathbf{h}_i = \text{GRU}(\mathbf{h}_{i-1}, \mathbf{r}_i)$$

$\mathbf{h}_0$  is initialized to  $\mathbf{0}$ .  $\mathbf{h}_i, i > 0$ , contains the information of the relation sequence  $r_1, r_2, \dots, r_i$  independent of the head entity  $\mathbf{e}_h$ . To make sure the path reaches the correct state, ROP\_arc2 then composes each  $\mathbf{h}_i$  with  $\mathbf{e}_h$  to predict  $\mathbf{e}_i$  as:

$$\hat{\mathbf{e}}_i = \text{COMP}_1(\mathbf{e}_h, \mathbf{h}_i) \quad (9)$$

As composition function  $\text{COMP}_1$  we use ADD (addition, as in TransE) or GRU (Eq 5).

In ROP\_ARC2, head entity  $e_h$  directly participates in entity prediction at each step. Unfortunately, there is often another issue – head entity  $e_h$  may not match the hidden state  $h_i$  if they are far away from each other –  $h_i$  mainly encodes some latest relation inputs that are less related to the head entity. Besides, there is a second information source, in addition to  $e_h$ , that is clearly relevant for accurate prediction: the preceding predicted entity  $\hat{e}_{i-1}$ . But ROP\_ARC2 does not make it available. Our solution is ROP\_ARC3. This architecture predicts the next entity  $e_i$  using both  $e_h$  and  $\hat{e}_{i-1}$ , based on the intuition that  $e_i$  is directly related to its predecessor  $e_{i-1}$ .

**ROP\_ARC3** combines the benefits of ROP\_ARC1 and ROP\_ARC2. It encodes the relation sequence as in Eq 8 in ROP\_ARC2, but composes head entity  $e_h$  as well as the *predicted* entity  $\hat{e}_{i-1}$  with  $h_i$  to predict the entity  $e_i$  as:

$$\hat{e}_i = \text{COMP}_2(e_h, \hat{e}_{i-1}, \mathbf{h}_i) \quad (10)$$

As composition function  $\text{COMP}_2$ , we use ADD (addition) or an extended GRU step, defined as:

$$\mathbf{g}_{zh} = \sigma(\mathbf{h}_i \mathbf{U}^{zh} + \mathbf{e}_h \mathbf{W}^{zh}) \quad (11)$$

$$\mathbf{g}_{rh} = \sigma(\mathbf{h}_i \mathbf{U}^{rh} + \mathbf{e}_h \mathbf{W}^{rh}) \quad (12)$$

$$\mathbf{g}_{zp} = \sigma(\mathbf{h}_i \mathbf{U}^{zp} + \hat{e}_{i-1} \mathbf{W}^{zp}) \quad (13)$$

$$\mathbf{g}_{rp} = \sigma(\mathbf{h}_i \mathbf{U}^{rp} + \hat{e}_{i-1} \mathbf{W}^{rp}) \quad (14)$$

$$\hat{\mathbf{h}}_i = \tanh(\mathbf{h}_i \mathbf{U}^q + (\mathbf{e}_h \circ \mathbf{g}_{rh}) \mathbf{W}^{qh} + (\hat{e}_{i-1} \circ \mathbf{g}_{rp}) \mathbf{W}^{qp}) \quad (15)$$

$$\hat{e}_i = (1 - \mathbf{g}_{zh} - \mathbf{g}_{zp}) \circ \hat{\mathbf{h}}_i + \mathbf{g}_{zh} \circ \mathbf{e}_h + \mathbf{g}_{zp} \circ \hat{e}_{i-1} \quad (16)$$

where super/subscript  $h$  refers to *head* entity and  $p$  refers to *prior* predicted entity  $\hat{e}_{i-1}$ .

In following work, we define Eqs. 11–16 as eGRU (extended GRU) step as:

$$\hat{e}_i = \text{eGRU}(e_h, \hat{e}_{i-1}, \mathbf{h}_i) \quad (17)$$

We extend GRU into eGRU, so that one architecture can compose three objects: a hidden state  $h_i$  and two entity states  $e_h$  and  $\hat{e}_{i-1}$ .

**Training.** For the gold entity sequence  $e_1, e_2, \dots, e_t$ , we define the **loss function** as the margin-based ranking criterion used in TransE (Bordes et al., 2013):

$$l_{\text{seq}} = \sum_i \max(0, \alpha + s(\mathbf{e}_i^-, \hat{\mathbf{e}}_i) - s(\mathbf{e}_i, \hat{\mathbf{e}}_i)) \quad (18)$$

where  $\alpha$  is the margin,  $e_i^-$  a negative sample for entity  $e_i$  and  $s()$  a similarity function.

**Discussion.** The three ROP models *differ* in three aspects.

(i) ROP\_ARC1 has only one composition process, i.e., GRU, to encode from head entity  $e_h$  to  $r_i$ ; ROP\_ARC2 and ROP\_ARC3 each have two composition processes, one composes relation sequence  $r_1, \dots, r_i$  into hidden state  $\mathbf{h}_i$ , the other composes entities ( $e_h$  in ROP\_ARC2,  $e_h$  and  $\hat{e}_{i-1}$  in ROP\_ARC3) with  $\mathbf{h}_i$  to predict  $e_i$ .

Figure 2 shows that ROP\_ARC1 uses the GRU hidden states as outputs to compare with gold entities. ROP\_ARC2 and ROP\_ARC3 instead compose their hidden states with head entity or preceding predicted entity to generate a new output space, then compare with gold entities.

(ii) ROP\_ARC2 only uses the head entity  $e_h$  along with the current hidden state  $\mathbf{h}_i$  to predict the next entity  $e_i$  whereas ROP\_ARC1 and ROP\_ARC3 in addition use the preceding predicted entity  $\hat{e}_{i-1}$ . Hence, ROP\_ARC3 roughly uses the combined information of ROP\_ARC1 and ROP\_ARC2 to do the prediction.

(iii) GRUs like ROP\_ARC2&ARC3 often zero-initialize first hidden state  $\mathbf{h}_0$ . Setting  $\mathbf{h}_0$  to the head entity in ROP\_ARC1 supports prediction of subsequent entities.

The ROP models are *similar* in three aspects.

(i) They model entities and relations in different spaces: relations are in the input space, entities are in hidden or output spaces. Our motivation is similar to SE, TransH and TransR (cf. §3).

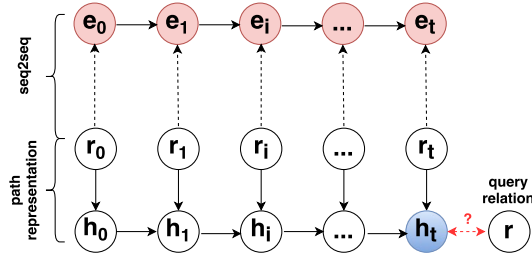


Figure 3: Architecture for mh-KBC task

(ii) The gating mechanism enables flexible compositions between entities and relations based on path context. In contrast, *one-hop KG embedding approaches model compositions statically as shared parameters and consider no context.*

(iii) Unlike Neelakantan et al. (2015), Guu et al. (2015) and Lin et al. (2015) (who also do some composition over relation paths), we finetune entity/relation embeddings in each step of the path. Our intuition is that many more training signals coming from each step of the sequence enable better learning of entity/relation embeddings. Das et al. (2017), as prior work that incorporates entities in the paths, only update the entity embeddings in the path once, when reaching the end of the path. We test this effect in our experiments.

## 5 Experiments

### 5.1 Knowledge Base Completion (mh-KBC)

**Dataset.** We use Das et al. (2017)’s dataset, in which there are 46 query relations, each has *train*, *dev* and *test* files containing positive and negative entity pairs (head entity, tail entity). Each entity pair is also provided with multiple multi-hop paths connecting them.

This is a binary classification task for each query relation: does the query relation hold between a pair of entities? The classifier builds a ranked list of entity pairs for corresponding query relation. Evaluation: mean average precision (MAP) across all 46 relations.

**Task setup.** We use ROP for mh-KBC. Figure 3 shows that we extend the basic ROP architecture (§4) by a GRU layer (bottom layer in Figure 3), denoted as  $GRU_r$ , that learns the representation of the relation sequence. Finally, we use  $h_t$  (the last hidden state of  $GRU_r$ , shown in blue in Figure 3) to match the query relation  $r$ .<sup>3</sup>

Thus, the training loss of this task has two parts: one comes from the loss of ROP training ( $l_{seq}$ , Eq 18); the other is the prediction loss of query relation, denoted as  $l_{pred}$ . Our preliminary experiments always showed worse performance of ADD, so the first loss  $l_{seq}$  here only considers (e)GRU. The second loss is as in (Neelakantan et al., 2015); we show that our ROP part boosts the system as it enables relation paths to encode entity information with multiple updating – as opposed to a single update per path as in (Das et al., 2017). Similar to  $l_{seq}$ , we define:

$$l_{pred} = \sum_{j,i} \max(0, \beta + s(\mathbf{h}_j^-, \mathbf{r}) - s(\mathbf{h}_i, \mathbf{r})) \quad (19)$$

where  $\mathbf{r}$  is the embedding of a query relation in Figure 3,  $\mathbf{h}_i$  (resp.  $\mathbf{h}_j^-$ ) is the representation of positive (resp. negative) path example  $i$  (resp.  $j$ ), shown as  $h_t$  in Figure 3. In testing, all paths are ranked in terms of the given query relation.

The target relation is often not inferable from all paths between head and tail (see Figure 1(a)); it can be entailed by a single path or a subset of paths. Hence, given the representation of a group of paths, how to match them with the representation of the target relation is a key problem. We use  $\max$  (over all paths) because we found it works well in selecting the best path for the target relation. See last paragraph of §5.1 for discussion.

<sup>3</sup>Using the last hidden state in ROP (which participates in predicting the tail entity) to predict the query relation performed much worse. We suspect the finetuning by tail entity makes it not indicative for query relations.

Model	MAP
PRA (Lao et al., 2011)	64.43
PRA+ Bigram (Neelakantan et al., 2015)	64.93
RNN-path (Neelakantan et al., 2015)	68.43
RNN-path-entity (Das et al., 2017)	71.74
RNN-path-types (Das et al., 2017)	73.26
ROP_ARC1	74.23
ROP_ARC2	74.46
ROP_ARC3	<b>76.16</b>

Table 1: Results of mh-KBC task

We use AdaGrad (Duchi et al., 2011) with learning rate 0.1. Relation embedding dimension is 200, margins  $\alpha$  and  $\beta$  are 0.5 (Eqs. 18–19), entity embedding dimension 200, batch size 20, negative sampling size 4. Longer paths are truncated to 8, the first 30 paths are kept for each entity pair.

**Results.** Table 1 compares our ROP systems to five baselines. RNN-path (Neelakantan et al., 2015) composes the relations occurring in a path using a vanilla RNN. It ignores all information about within-path entities and trains separate models per relation. RNN-path-entity (Das et al., 2017) models the path entities and improves the results. RNN-path-type (Das et al., 2017) further improves the result by representing entities with the sum of their type embeddings. Thus, RNN-path-type uses extra information, the entity types. Apart from these baselines, it is also feasible to compare with the baselines based on composition of one-hop triple-based embedding models. However, the performance of these baselines is very poor for this task (Neelakantan et al., 2015) and therefore, we do not include them in our comparison.

The performance order of our architectures for mh-KBC is  $ROP\_ARC3 > ROP\_ARC2 > ROP\_ARC1$ . All our ROP systems are superior to the state-of-the-art. In particular, they are superior to RNN-path-type (Das et al., 2017), the state-of-the-art, even though we do not need and do not use information about the types of entities. Comparing ROP performance to RNN-path-entity is more fair, and this makes it even clearer that our modeling is effective. Das et al. (2017) encode entities along with relations into the path representation explicitly. Their motivation is that entity incorporation prevents prediction of the same target relation when relation sequences are the same, but path entities are different. Our ROP models achieve this implicitly, as the relation sequence can predict the entity sequence; this makes sure the representation of the relation sequence is specific to the entity sequence. As a result, the same goal can be achieved as (Das et al., 2017).

In (Das et al., 2017), an entity and a relation *are concatenated as a new unit* in the path; both entity and relation representations are trained based on the given gold relation as label. ROP predicts intermediate entities and updates the entity/relation embeddings and other parameters in each step; the training signals can come from the in-path target entities and from the reasoning task specific annotations. Thus, ROP makes use of the in-path structures to learn good quality entity/relation embeddings, which are further employed and finetuned to solve the reasoning problem.

This can also explain why max (over all paths) works well for us while Das et al. (2017) found Log-SumExp (over all paths) works better. As their system solely relies on target relations as training signals, max selection prevents all other paths from generating gradients to update, so max tends to select randomly in the initial stage. However, in our system, the representations of entities and relations get rich training signals at each path hop, so that the path representations are more reliable. Hence, max can select the most informative path and avoid misleading paths to support the claim of target relation. In a similar vein, max pooling is widely found more effective than mean/sum pooling for classification as this task mostly relies on the dominant features.

## 5.2 Path Query Answering (mh-PQA)

**Dataset.** We use the mh-PQA dataset released by Guu et al. (2015) and refer to it as *BaseKGP*<sup>4</sup>. BaseKGP contains paths like  $e_h, r_1, r_2, \dots, r_t, e_t$ , where  $e_h$  and  $e_t$  are head and tail entities, connected

<sup>4</sup>Das et al. (2017) use a dataset based on the WordNet, however they conclude that the dataset is not an ideal test bed for mh-PQA due to some limitations: it is fairly small, with very short paths, few unseen paths during test time, and only one path between an entity pair. Therefore, we experiment on BaseKGP.

	#paths	#entities	#relations
train	6,266,058	75,043	26
dev	27,163	41,010	26
test	109,557	96,858	26

Table 2: Statistics of EnhancedKGP for mh-PQA

methods		MQ	H@10
BaseKGP (baselines)	Comp-Bilinear	83.5	42.1
	Comp-Bilinear-Diag	84.8	38.6
	Comp-TransE	88.0	50.5
BaseKGP (our model)	ROP_ARC1	88.3	52.7
	ROP_ARC2 (ADD)	89.3	54.3
	ROP_ARC2 (GRU)	89.6	54.9
EnhancedKGP (our model)	ROP_ARC1	89.4	54.2
	ROP_ARC2 (ADD)	90.3	55.5
	ROP_ARC2 (GRU)	90.5	56.3
	ROP_ARC3 (ADD)	90.3	55.8
	ROP_ARC3 (eGRU)	<b>90.7</b>	<b>56.7</b>

Table 3: Results of mh-PQA

by relation sequence  $r_1, \dots, r_t$ . There are 6,266,058/27,163/109,557 paths in train/dev/test.

BaseKGP is based on a Freebase subset released by Socher et al. (2013). The original subset contains a collection of Freebase triplets in form of (head, relation, tail). Guu et al. (2015) generated paths by traversing the triplet space. Paths in BaseKGP of form  $e_h, r_1, r_2, \dots, r_t, e_t$  do not contain intermediate entities. We create *EnhancedKGP* by enhancing each BaseKGP path in *train* as follows. We search entities at each step of a path by traversing the subset (Socher et al., 2013) until reaching the tail entity  $e_t$ . When there are multiple entity choices at a step, we randomly choose one. EnhancedKGP *train* has the same size as BaseKGP *train*, except that paths are filled by intermediate entities. Table 2 gives statistics. We will release EnhancedKGP, the first dataset for mh-PQA that includes within-path entities.

**Task setup.** We tune parameters on dev. We sample 10 negative entities for each ground truth entity in the path and use ranking loss (Eq 18) (with  $\alpha=0.3$ ,  $s(\cdot) = \text{cosine similarity}$ ). For testing, we ignore intermediate predicted entities and only output the tail entity. We update parameters – relation embeddings, entity embeddings (both dimension 300) and GRU parameters – using AdaGrad with learning rate 0.01 and mini-batch size 50.

We compare ROP with the three compositional training schemes *Bilinear*, *Bilinear-Diag* and *TransE* in (Guu et al., 2015). The compositional training of TransE, denoted as *Comp-TransE* in this work, is the state-of-the-art in mh-PQA. For ROP\_ARC2, we report the performance of using ADD and GRU for COMP<sub>1</sub> in Eq 9. Similarly, for ROP\_ARC3, we report the performance of using ADD and eGRU for COMP<sub>2</sub> in Eq 10.

In addition, to better investigate ROP models, we run them on both BaseKGP and EnhancedKGP. Note that since there are no intermediate entities in BaseKGP, ROP\_ARC3 is reduced to ROP\_ARC2. Hence, we report results on BaseKGP only for ROP\_ARC1 and ROP\_ARC2.

Two benchmark metrics are reported for this task (Guu et al., 2015): *hits at 10* (H@10), percentage of ground truth tail entities ranked in the top 10 of all retrieved; and *mean quantile* (MQ), normalized version of mean rank.

**Results.** Table 3 gives results for mh-PQA: top baselines on BaseKGP (block 1), ROP results on BaseKGP (block 2) and ROP results on EnhancedKGP (block 3). Overall, our ROP models all get new state-of-the-art by large margins (2–4% H@10 on BaseKGP, 4–6% H@10 on EnhancedKGP). The improvements of the second block over the first are evidence that recurrent neural networks are an appropriate paradigm in this task. The third block shows that encoding intermediate entities in an appropriate



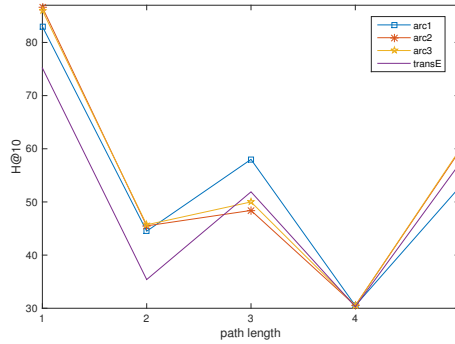


Figure 4: H@10 vs. path lengths on test set

way, like ROP does, can give a further boost.

Comparing ROP\_ARC2 and ROP\_ARC1, we realize that forwarding head entity  $e_h$  *directly* to predict each intermediate entity  $e_i$  is better than putting  $e_h$  as the start state of the GRU. We suspect this is due to the fact that the latest state of GRU is gradually less influenced by  $e_h$  when the following relation path is getting longer and longer. In ROP\_ARC2, no matter how long the relation sequence  $r_1, r_2, \dots, r_i$  is, we always compose its whole representation  $\mathbf{h}_i$  with the representation of  $e_h$  so that  $e_h$  can participate in the prediction of  $e_i$  more directly.

ROP\_ARC3 further improves results by composing not only head entity  $e_h$ , but also the preceding predicted entity  $\hat{e}_{i-1}$  with  $h_i$  to predict entity  $e_i$ . The result suggests that  $\hat{e}_{i-1}$  provides critical information in this prediction process. This may be due to the property of RNN that latest inputs tend to be remembered better than old inputs, so the latest hidden state  $\mathbf{h}_i$  is heavily influenced by the adjacent relations of position  $i$ ; employing the preceding predicted entity  $\hat{e}_{i-1}$  hence can help the prediction of  $e_i$ .

For composition, GRU/eGRU always surpass ADD – presumably due to the higher expressibility of (e)GRU’s gating mechanism.

**Performance vs. path lengths.** Figure 4 graphs H@10 on different path lengths for our three ROP systems and the Comp-TransE baseline. We expected that H@10 should decrease gradually. However, Figure 4 shows that even-length paths are harder than odd-length ones. This surprising phenomenon is *due to a large number of inverse relations, as most inverse relations in BaseKGP are 1-to-N connections.*

The percentages of inverse relations in paths of length 1, 2, 3, 4, 5 are 0.0, 49.7, 38.3, 49.5 and 42.9, respectively. Spearman correlation coefficients between these percentages and system performance are always higher than 0.9. Thus, the more inverse relations, the worse performance. Unfortunately, to date there is no good way to model pairs of a relation  $\mathbf{r}$  and its inverse  $\ast\mathbf{r}$ , e.g. “gender” and “\*gender”, as separate, yet at the same time as systematically related. Guu et al. (2015)’s TransE implementation enforces  $\ast\mathbf{r} + \mathbf{r} = \mathbf{0}$ . Figure 4 shows this is not as effective as expected, perhaps because it fails for 1-to-N projections. ROP treats inverse relations as independent, but we did not observe any worse performance. Better modeling of inverse relations is an interesting challenge for future work.

## 6 Conclusion

This work presented three ROP architectures for multi-hop KG reasoning. ROP models a KG path of arbitrary length as a pair of a relation sequence and an entity sequence, using the former to predict the latter by encoding and decoding step by step. Our neural sequential modeling of KG paths enables better learning of entity/relation embeddings because there are more training signals at each multi-hop path. ROPs showed state-of-the-art in two representative KG reasoning tasks, multi-hop KBC and multi-hop PQA. Incorporating knowledge from textual sources by initializing the entity embeddings with distributional representation of entities (Yaghoobzadeh and Schütze, 2015) could improve our results further, which we will explore in the future.

## Acknowledgement

We gratefully acknowledge the support of the European Research Council for this research (ERC Advanced Grant NonSequeToR, # 740516).

## References

- Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of SIGMOD*, pages 1247–1250.
- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. In *Proceedings of NIPS*, pages 2787–2795.
- Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder-decoder approaches. *Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*.
- Rajarshi Das, Arvind Neelakantan, David Belanger, and Andrew McCallum. 2017. Chains of reasoning over entities, relations, and text using recurrent neural networks. In *Proceedings of EACL*, pages 132–141.
- John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *JMLR*, 12:2121–2159.
- Jeffrey L. Elman. 1990. Finding structure in time. *Cognitive Science*, 14(2):179–211.
- Kelvin Guu, John Miller, and Percy Liang. 2015. Traversing knowledge graphs in vector space. In *Proceedings of EMNLP*, pages 318–327.
- Zhiheng Huang, Wei Xu, and Kai Yu. 2015. Bidirectional LSTM-CRF models for sequence tagging. *CoRR*, abs/1508.01991.
- Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016. Neural architectures for named entity recognition. In *Proceedings of NAACL*, pages 260–270.
- Ni Lao, Tom Mitchell, and William W Cohen. 2011. Random walk inference and learning in a large scale knowledge base. In *Proceedings of EMNLP*, pages 529–539.
- Yankai Lin, Zhiyuan Liu, Huan-Bo Luan, Maosong Sun, Siwei Rao, and Song Liu. 2015. Modeling relation paths for representation learning of knowledge bases. In *Proceedings of EMNLP*, pages 705–714.
- Xixun Lin, Yanchun Liang, and Renchu Guan. 2016. Compositional learning of relation paths embedding for knowledge base completion. *CoRR*, abs/1611.07232.
- Bonan Min, Ralph Grishman, Li Wan, Chang Wang, and David Gondek. 2013. Distant supervision for relation extraction with an incomplete knowledge base. In *Proceedings of HLT-NAACL*, pages 777–782.
- Arvind Neelakantan, Benjamin Roth, and Andrew McCallum. 2015. Compositional vector space models for knowledge base completion. In *Proceedings of ACL*, pages 156–166.
- Maximilian Nickel, Volker Tresp, and Hans-Peter Kriegel. 2011. A three-way model for collective learning on multi-relational data. In *Proceedings of ICML*, pages 809–816.
- Yelong Shen, Po-Sen Huang, Ming-Wei Chang, and Jianfeng Gao. 2016. Implicit reasoner: Modeling large-scale structured relationships with shared memory. *CoRR*, abs/1611.04642.
- Richard Socher, Danqi Chen, Christopher D Manning, and Andrew Ng. 2013. Reasoning with neural tensor networks for knowledge base completion. In *Proceedings of NIPS*, pages 926–934.
- Kristina Toutanova, Victoria Lin, Wen-tau Yih, Hoifung Poon, and Chris Quirk. 2016. Compositional learning of embeddings for relation paths in knowledge base and text. In *Proceedings of ACL*.
- Yadollah Yaghoobzadeh and Hinrich Schütze. 2015. Corpus-level fine-grained entity typing using contextual information. In *Proceedings of EMNLP*, pages 715–725.