# Relational Message Passing for Knowledge Graph Completion

Hongwei Wang
Stanford University
Stanford, California, United States
hongweiw@cs.stanford.edu

Hongyu Ren
Stanford University
Stanford, California, United States
hyren@cs.stanford.edu

Jure Leskovec
Stanford University
Stanford, California, United States
jure@cs.stanford.edu

## ABSTRACT

Knowledge graph completion aims to predict missing relations between entities in a knowledge graph. In this work, we propose a *relational message passing* method for knowledge graph completion. Different from existing embedding-based methods, relational message passing only considers edge features (i.e., relation types) without entity IDs in the knowledge graph, and passes relational messages among edges iteratively to aggregate neighborhood information. Specifically, two kinds of neighborhood topology are modeled for a given entity pair under the relational message passing framework: (1) *Relational context*, which captures the relation types of edges adjacent to the given entity pair; (2) *Relational paths*, which characterize the relative position between the given two entities in the knowledge graph. The two message passing modules are combined together for relation prediction. Experimental results on knowledge graph benchmarks as well as our newly proposed dataset show that, our method PATHCON outperforms state-of-the-art knowledge graph completion methods by a large margin. PATHCON is also shown applicable to inductive settings where entities are not seen in training stage, and it is able to provide interpretable explanations for the predicted results. The code and all datasets are available at https://github.com/hwwang55/PathCon.

## CCS CONCEPTS

• **Computing methodologies → Semantic networks**; **Statistical relational learning**; • **Mathematics of computing →** *Graph algorithms*.
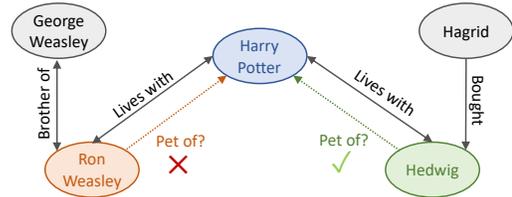
## KEYWORDS

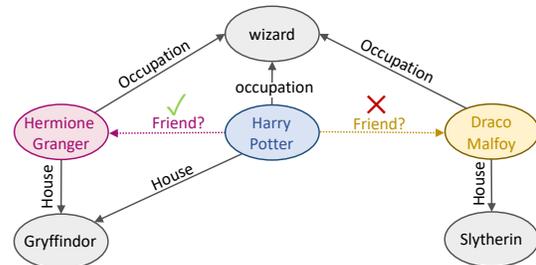Knowledge graph completion; message passing; graph neural networks

(a) Consider we aim to predict whether Ron Weasley or Hedwig is a Pet of Harry Potter. Both entities have the same relational path (Lives with) to Harry Potter but they have distinct relational context: Ron Weasley has {Brother of, Lives with}, while Hedwig has {Bought, Lives with}. Capturing the relational context of entities allows our model to make a distinction between Ron Weasley, who is a person, and Hedwig, which is an owl.



(b) Two head entities Hermione Granger and Draco Malfoy have the same relational context {Occupation, House}, but different relational paths to the tail entity Harry Potter {(House, House), (Occupation, Occupation)} vs. {(Occupation, Occupation)}, which allows our model to predict friendship between Harry Potter and Hermione Granger vs. Draco Malfoy.

Figure 1: (a) Relational context of an entity and (b) relational paths between entities. Our model is able to capture both.

## 1 INTRODUCTION

Knowledge graphs (KGs) store structured information of real-world entities and facts. A KG usually consists of a collection of triplets. Each triplet $(h, r, t)$ indicates that head entity $h$ is related to tail entity $t$ through relationship type $r$. Nonetheless, KGs are often incomplete and noisy. To address this issue, researchers have proposed a number of KG completion methods to predict missing links/relations in KGs [3, 9, 12, 13, 19, 22, 24, 35, 36, 39, 40].
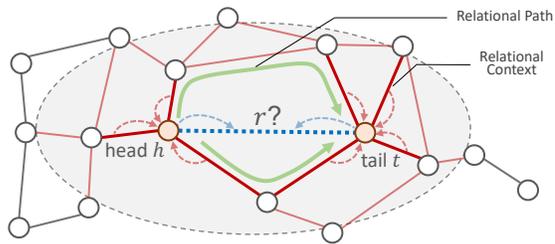
In general, relation types are not uniformly distributed over a KG but spatially correlated with each other. For example, the neighboring relations of "graduated from" in the KG are more likely to be "person.birthplace" and "university.location" rather than "movie.language". Therefore, for a given entity pair $(h, t)$, characterizing the relation types of neighboring links of $h$ and $t$ will provide valuable information when inferring the relation type between $h$ and $t$. Inspired by recent success of graph neural networks [11, 15, 34], we propose using *message passing* to capture the neighborhood structure for a given entity pair. However, traditional message passing methods usually assume that messages

are associated with nodes and messages are passed from nodes to nodes iteratively, which are not suitable for KGs where edge features (relation types) are more important.

**Relational message passing**. To address the above limitation, we propose *relational message passing* for KG completion. Unlike traditional node-based message passing, relational message passing *only* considers edge features (relation types), and passes messages of an edge directly to its neighboring edges. Note that since relational message passing only models relations rather than entities, it brings three additional benefits compared with existing knowledge graph embedding methods [3, 13, 22, 24, 35, 39]: (1) it is *inductive*, since it can handle entities that do not appear in the training data during inference stage; (2) it is *storage-efficient*, since it does not calculate embeddings of entities; and (3) it is *explainable*, since it is able to provide explainability for predicted results by modeling the correlation strength among relation types. However, a potential issue of relational message passing is that its computational complexity is significantly higher than node-based message passing (Theorem 2). To solve this issue, we propose *alternate relational message passing* that passes relational messages between nodes and edges *alternately* over the KG. We prove that alternate message passing scheme greatly improves time efficiency and achieves *the same order of computational complexity* as traditional node-based message passing (Theorem 1 and 3).

**Relational context and relational paths**. Under the alternate relational message passing framework, we explore two kinds of local subgraph topology for a given entity pair $(h, t)$ (see Figure 1 for an illustrating example): (1) *Relational context*. It is important to capture the neighboring relations of a given entity in the KG, because neighboring relations provide us with valuable information about what is the nature or the "type" of the given entity (Figure 1a). Many entities in KGs are not typed or are very loosely typed, so being able to learn about the entity and its context in the KG is valuable. We design a multi-layer relational message passing scheme to aggregate information from multi-hop neighboring edges of $(h, t)$. (2) *Relational paths*. Note that modeling only relational context is not able to identify the relative position of $(h, t)$. It is also important to capture the set of relational paths between $(h, t)$ (Figure 1b). Here different paths of connections between the entities reveal the nature of their relationship and help with the prediction. Therefore, we calculate all relational paths connecting $h$ and $t$ in the KG and pass relational messages along these paths. Finally, we use an attention mechanism to selectively aggregate representations of different relational paths, then combine the above two modules together for relation prediction.

**Experiments**. We conduct extensive experiments on five well-known KGs as well as a new KG proposed by us, *DDB14 dataset*. Experimental results demonstrate that our proposed model PathCon (short for relational PATHs and CONtext) significantly outperforms state-of-the-art KG completion methods, for example, the absolute Hit@1 gain over the best baseline is 16.7% and 6.3% on WN18RR and NELL995, respectively. Our ablation studies show the effectiveness of our approach and demonstrate the importance of relational context as well as relational paths. Our method is also shown to maintain strong performance in inductive KG completion, and it



**Figure 2: An example of PathCon considering both the relational context within 2 hops of the head and the tail entities (denoted by red edges) and relational paths of length up to 3 relations that connect head to tail (denoted by green arrows). Context and paths are captured based on relation types (not entities) they contain. By combining the context and paths PathCon predicts the probability of relation $r$.**

provides high explainability by identifying important relational context and relation paths for a given predicted relation.

**Contributions**. Our key contributions are listed as follows:

- We propose *alternate relational message passing* framework for KG completion, which is *inductive*, *storage-efficient*, *explainable*, and *computationally efficient* compared with existing embedding-based methods;
- Under the proposed framework, we explore two kinds of subgraph topology: *relational context* and *relational paths*, and show that they are critical to relation prediction;
- We propose a new KG dataset DDB14 (Disease Database with 14 relation types) that is suitable for KG-related research.

## 2 PROBLEM FORMULATION

Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be an instance of a knowledge graph, where $\mathcal{V}$ is the set of nodes and $\mathcal{E}$ is the set of edges. Each edge $e$ has a relation type $r \in \mathcal{R}$. Our goal is to predict missing relations in $\mathcal{G}$, i.e., given an entity pair $(h, t)$, we aim to predict the relation of the edge between them.[1] Specifically, we aim to model the distribution over relation types given a pair of entities $(h, t)$: $p(r|h, t)$. This is equivalent to modeling the following term

$$p(r|h, t) \propto p(h, t|r) \cdot p(r) \tag{1}$$

according to Bayes' theorem. In Eq. (1), $p(r)$ is the prior distribution over relation types and serves as the regularization of the model. Then the first term can be further decomposed to

$$p(h, t|r) = \frac{1}{2}\Big(p(h|r) \cdot p(t|h, r) + p(t|r) \cdot p(h|t, r)\Big). \tag{2}$$

Eq. (2) sets up the guideline for designing our model. The term $p(h|r)$ or $p(t|r)$ measures the likelihood of an entity given a particular relation. Since our model does not consider the identity of entities, we use an entity's *local relational subgraph* instead to represent the entity itself, i.e., $p(C(h)|r)$ and $p(C(t)|r)$ where $C(\cdot)$

[1]Some of the related work formulates this problem as predicting the missing tail (head) entity given a head (tail) entity and a relation. The two problems are actually reducible to each other: Given a model $\Phi(\cdot|h, t)$ that outputs the distribution over relation types for an entity pair $(h, t)$, we can then build a model $\Gamma(\cdot|h, r) = \text{SOFTMAX}_t (\Phi(r|h, t))$ that outputs the distribution over tail entities given $h$ and $r$, and vice versa. Since the two problems are equivalent, we only focus on relation prediction in this work.

| Symbol | Description |
|---|---|
| $h, t$ | Head entity and tail entity |
| $r$ | Relation type |
| $s_e^i$ | Hidden state of edge $e$ at iteration $i$ |
| $m_v^i$ | Message of node $v$ at iteration $i$ |
| $\mathcal{N}(e)$ | Endpoint nodes of edge $e$ |
| $\mathcal{N}(v)$ | Neighbor edges of node $v$ |
| $s_{(h,t)}$ | Context representation of the entity pair $(h,t)$ |
| $s_{h \rightarrow t}$ | Path representation of all paths from $h$ to $t$ |
| $\alpha_P$ | Attention weight of path $P$ |
| $\mathcal{P}_{h \rightarrow t}$ | Set of paths from $h$ to $t$ |

Table 1: Notation used in this paper.

denotes the local relational subgraph of an entity. This is also known as *relational context* for $h$ and $t$. The term $p(t|h,r)$ or $p(h|t,r)$ in Eq. (2) measures the likelihood of how $t$ can be reached from $h$ or the other way around given that there is a relation $r$ between them. This inspires us to model the *relational paths* between $h$ and $t$ in the KG. In the following we show how to model the two factors in our method and how they contribute to relation prediction.

## 3 OUR APPROACH

In this section, we first introduce the relational message passing framework, then present two modules of the proposed PATHCON: relational context message passing and relational path message passing. Notations used in this paper are listed in Table 1.

### 3.1 Relational Message Passing Framework

**Traditional node-based message passing**. We first briefly review traditional node-based message passing method for general graphs. Assume that each node $v$ is with feature $x_v$. Then the message passing runs for multiple timesteps over the graph, during which the hidden state $s_v^i$ of each node $v$ in iteration $i$ is updated by

$$m_v^i = A\left(\left\{s_u^i\right\}_{u \in \mathcal{N}(v)}\right), \tag{3}$$

$$s_v^{i+1} = U\left(s_v^i, m_v^i\right), \tag{4}$$

where $m_v^i$ is the message received by node $v$ in iteration $i$, $\mathcal{N}(v)$ denotes the set of neighbor nodes of $v$ in the graph, $A(\cdot)$ is message aggregation function, and $U(\cdot)$ is node update function. The initial hidden state $s_v^0 = x_v$.

The above framework, though popular for general graphs and has derived many variants such as GCN [15], GraphSAGE [11], and GIN [34], faces the following challenges when applied to knowledge graphs: (1) Unlike general graphs, in most KGs, edges have features (relation types) but nodes don't, which makes node-based message passing less natural for KGs. Though node features can be set as their identities (i.e., one-hot vectors), this will lead to another two issues: (2) Modeling identity of nodes cannot manage previously unseen nodes during inference and fails in inductive settings. (3) In real-world KGs, the number of entities are typically much larger than the number of relation types, which requires large memory for storing entity embeddings.

**Relational message passing**. To address the above problems, a natural thought is to perform message passing over edges instead

of nodes:

$$m_e^i = A\left(\left\{s_{e'}^i\right\}_{e' \in \mathcal{N}(e)}\right), \tag{5}$$

$$s_e^{i+1} = U\left(s_e^i, m_e^i\right), \tag{6}$$

where $\mathcal{N}(e)$ denotes the set of neighbor edges of $e$ (i.e., edges that share at lease one common end-point with $e$) in the graph, and $s_e^0 = x_e$ is the initial edge feature of $e$, i.e., the relation type. Therefore, Eqs. (5) and (6) are called *relational message passing*.

Relational message passing avoids the drawbacks of node-based message passing, however, it brings a new issue of computational efficiency when passing messages. To see this, we analyze the computational complexity of the two message passing schemes (proofs are given in Appendix A and B):

THEOREM 1 (COMPLEXITY OF NODE-BASED MESSAGE PASSING). *Consider a graph with $N$ nodes and $M$ edges. The expected cost of node-based message passing (Eqs. (3) and (4)) in each iteration is $2M + 2N$.*

THEOREM 2 (COMPLEXITY OF RELATIONAL MESSAGE PASSING). *Consider a graph with $N$ nodes and $M$ edges. The expected cost of relational message passing (Eqs. (5) and (6)) in each iteration is $N \cdot \text{Var}[d] + \frac{4M^2}{N}$, where $\text{Var}[d]$ is the variance of node degrees in the graph.*

**Alternate relational message passing**. According to the above theorems, the complexity of relational message passing is much higher than node-based message passing, especially in real-world graphs where node distribution follows the power law distribution whose variance ($\text{Var}[d]$) is extremely large due to the long tail. To reduce the redundant computation in relational message passing and improve its computational efficiency, we propose the following message passing scheme for KGs:

$$m_v^i = A_1\left(\left\{s_e^i\right\}_{e \in \mathcal{N}(v)}\right), \tag{7}$$

$$m_e^i = A_2\left(m_v^i, m_u^i\right), \ v, u \in \mathcal{N}(e), \tag{8}$$

$$s_e^{i+1} = U\left(s_e^i, m_e^i\right). \tag{9}$$

We decompose edge aggregation in Eq. (5) into two steps as Eqs. (7) and (8). In Eq. (7), for each node $v$, we aggregate all the edges that $v$ connects to by an aggregation function $A_1(\cdot)$ and get message $m_v^i$, where $\mathcal{N}(v)$ denotes the set of neighbor edges for node $v$. Then in Eq. (8), we obtain message $m_e^i$ of edge $e$ by aggregating messages from its two end-points $v$ and $u$ using function $A_2(\cdot)$, where $\mathcal{N}(e)$ denotes the set of neighbor nodes for edge $e$. The hidden state of edge $e$ is finally updated using the message $m_e^i$ as in Eq. (9).

An intuitive understanding of alternate relational message passing is that nodes here serve as "distribution centers" that collect and temporarily store the messages from their neighbor edges, then propagate the aggregated messages back to each of their neighbor edges. Therefore, we call Eqs. (7)-(9) *alternate relational message passing*, as messages are passed alternately between nodes and edges.

The complexity of alternate relational message passing is given as follows (proof is given in Appendix C):

THEOREM 3 (COMPLEXITY OF ALTERNATE RELATIONAL MESSAGE PASSING). *Consider a graph with $N$ nodes and $M$ edges. The expected cost of alternate relational message passing (Eqs. (7)-(9)) in each iteration is $6M$.*

From Theorem 3 it is clear to see that alternate relational message passing greatly reduces the time overhead and achieves the same order of complexity as node-based message passing.

**Remarks.** We present the following two remarks to provide more insight on the proposed framework:

REMARK 1 (RELATIONSHIP WITH BELIEF PROPAGATION). *Alternate relational message passing is conceptually related to belief propagation (BP) [37], which is also a message-passing algorithm that passes messages between nodes and edges. But note that they are significantly different in: (1) application fields. BP is used to calculate the marginal distribution of unobserved variables in a graphical model, while our method aims to predict the edge type in KGs; (2) the purpose of using edge-node alternate message passing. BP uses this because of the special structure of factor graphs, while we use this to reduce the computational overhead.*

REMARK 2 (UTILIZING NODE FEATURES). *Though our proposed framework is claimed to only use edge features, it can be easily extended to the case where node features are present and assumed to be important, by additionally including the feature vector of node $v$ in Eq. (7), i.e., $m_v^i = A_1\left(\left\{s_e^i\right\}_{e \in \mathcal{N}(v)}, x_v\right)$, where $x_v$ is the feature of node $v$. As long as node features do not contain node identities, our proposed framework is still inductive. We do not empirically study the performance of our method on node-feature-aware cases, because node features are unavailable for all datasets used in this paper. We leave the exploration of this extension to future work.*

## 3.2 Relational Context

For a KG triplet $(h, r, t)$, relational context of $h$ and $t$ is usually highly correlated with $r$. For example, if $r$ is "graduated from", it's reasonable to guess that the surrounding relations of $h$ are "person.birthplace", "person.gender", etc., and the surrounding relations of $t$ are "institution.location", "university.founder", "university.president", etc. Therefore, the context of $h$ and $t$ will provide valuable clues when identifying the relation type of the edge between them, and here we use the proposed message passing method to learn from relational context.

Denote $s_e^i$ as the hidden state of edge $e$ in iteration $i$, and $m_v^i$ as the message stored at node $v$ in iteration $i$. We instantiate the alternate relational message passing in Eqs. (7)-(9) to learn the representation of each edge:

$$m_v^i = \sum_{e \in \mathcal{N}(v)} s_e^i, \tag{10}$$

$$s_e^{i+1} = \sigma\left(\left[m_v^i, m_u^i, s_e^i\right] \cdot W^i + b^i\right), \ v, u \in \mathcal{N}(e), \tag{11}$$

where $[\cdot]$ is the concatenation function, $W^i$, $b^i$, and $\sigma(\cdot)$ are the learnable transformation matrix, bias, and nonlinear activation function, respectively.[2] $s_e^0 = x_e$ is initial feature of edge $e$, which

can be taken as the one-hot identity vector of the relation type that $e$ belongs to.[3]

Relational context message passing in Eqs. (10) and (11) are repeated for $K$ times. The final message $m_h^{K-1}$ and $m_t^{K-1}$ are taken as the representation for head $h$ and tail $t$, respectively. We also give an illustrative example of relational context message passing in Figure 2, where the red/pink edges denote the first-order/second-order contextual relations.

## 3.3 Relational Paths

We follow the discussion in Section 2 and discuss how to model the term $p(t|h, r)$ or $p(h|t, r)$. Note that we do not consider node/edge identity in relational context message passing, which leads to a potential issue that our model is not able to identify the relative position between $h$ and $t$ in the KG. For example, suppose for a given entity pair $(h, t)$, $h$ is surrounded by "person.birthplace", "person.gender", etc., and $t$ is surrounded by "institution.location", "university.founder", "university.president", etc. Then it can be inferred that $h$ is probably a person and $t$ is probably a university, and there should be a relation "graduated_from" between them because such a pattern appears frequently in the training data. However, the person may have no relationship with the university and they are far from each other in the KG. The reason why such false positive case happens is that relational context message passing can only detect the "type" of $h$ and $t$, but is not aware of their relative position in the KG.

To solve this problem, we propose to explore the connectivity pattern between $h$ and $t$, which are represented by the paths connecting them in the KG. Specifically, a raw path from $h$ to $t$ in a KG is a sequence of entities and edges: $h(v_0) \xrightarrow{e_0} v_1 \xrightarrow{e_1} v_2 \cdots v_{L-1} \xrightarrow{e_{L-1}} t(v_L)$, in which two entities $v_i$ and $v_{i+1}$ are connected by edge $e_i$, and each entity in the path is unique.[4] The corresponding relational path $P$ is the sequence of relation types of all edges in the given raw path, i.e., $P = (r_{e_0}, r_{e_1}, ..., r_{e_{L-1}})$, where $r_{e_i}$ is the relation type of edge $e_i$. Note that we do not use the identity of nodes when modeling relational paths, which is the same as for relational context.

Denote $\mathcal{P}_{h \to t}$ as the set of all relational paths from $h$ to $t$ in the KG. Our next step is to define and calculate the representation of relational paths. In PATHCON, we assign an independent embedding vector $s_P$ for each relational path $P \in \mathcal{P}_{h \to t}$.[5] A potential concern here is that the number of different paths increases exponentially with the path length (there are $|r|^k$ $k$-hop paths), however, in practice we observe that in real-world KGs most paths actually do not occur (e.g., only 3.2% of all possible paths of length 2 occur in FB15K dataset), and the number of different paths is actually quite manageable for relatively small values of $k$ ($k \le 4$).

An illustrative example of relational paths is shown in Figure 2, where the two green arrows denote the relational paths from head entity $h$ to tail entity $t$.

---

[2]We shall discuss other implementations of Eqs. (7)-(9) in Section 3.6 and examine their performance in experiments.

[3]In cases where relation types have names, initial features can also be bag-of-words (BOW) or sentence embeddings learned by language models like BERT [7]. We shall investigate the performance of different initial feature types in experiments.

[4]Entities in a path are required to be unique because a loop within a path does not provide additional semantics thus should be cut off from the path.

[5]Other methods for calculating path representations are also possible. We shall discuss them in Section 3.6.

## 3.4 Combining Relational Context and Paths

For relational context, we use massage passing scheme to calculate the final message $m_h^{K-1}$ and $m_t^{K-1}$ for $h$ and $t$, which summarizes their context information, respectively. $m_h^{K-1}$ and $m_t^{K-1}$ are further combined together for calculating the context of $(h, t)$ pair:

$$s_{(h,t)} = \sigma\left(\left[m_h^{K-1}, m_t^{K-1}\right] \cdot W^{K-1} + b^{K-1}\right), \tag{12}$$

where $s_{(h,t)}$ denotes the context representation of the entity pair $(h, t)$. Note here that Eq. (12) should only take messages of $h$ and $t$ as input without their connecting edge $r$, since the ground truth relation $r$ should be treated unobserved in the training stage.

For relational paths, note that there may be a number of relational paths for a given $(h, t)$ pair, but not all paths are logically related to the predicted relation $r$, and the importance of each path also varies. In PATHCON, since we have already known the context $s_{(h,t)}$ for $(h, t)$ pair and it can be seen as prior information for paths between $h$ and $t$, we can calculate the importance scores of paths based on $s_{(h,t)}$. Therefore, we first calculate the attention weight of each path $P$ with respect to the context $s_{(h,t)}$:

$$\alpha_P = \frac{\exp\left(s_P^\top s_{(h,t)}\right)}{\sum_{P \in \mathcal{P}_{h \to t}} \exp\left(s_P^\top s_{(h,t)}\right)}, \tag{13}$$

where $\mathcal{P}_{h \to t}$ is the set of all paths from $t$ to $t$. Then the attention weights are used to average representations of all paths:

$$s_{h \to t} = \sum_{P \in \mathcal{P}_{h \to t}} \alpha_P s_P, \tag{14}$$

where $s_{h \to t}$ is the aggregated representation of relational paths for $(h, t)$. In this way, the context information $s_{(h,t)}$ is used to assist in identifying the most important relational paths.

Given the relational context representation $s_{(h,t)}$ and the relational path representation $s_{h \to t}$, we can predict relations by first adding the two representation together and then taking softmax as follows:

$$p(r|h, t) = \text{SOFTMAX}\left(s_{(h,t)} + s_{h \to t}\right). \tag{15}$$

Our model can be trained by minimizing the loss between predictions and ground truths over the training triplets:

$$\min \mathcal{L} = \sum_{(h,r,t) \in \mathcal{D}} J\big(p(r|h, t), r\big), \tag{16}$$

where $\mathcal{D}$ is the training set and $J(\cdot)$ is the cross-entropy loss.

It is worth noticing that the context representation $s_{(h,t)}$ plays two roles in the model: It directly contributes to the predicted relation distribution, and it also helps determine the importance of relational paths with respect to the predicted relation.

## 3.5 Discussion on Model Explainability

Since PATHCON only models relations without entities, it is able to capture pure relationship among different relation types thus can naturally be used to explain for predictions. The explainability of PATHCON is two-fold:

On the one hand, modeling relational context captures the correlation between contextual relations and the predicted relation, which can be used to indicate important neighbor edges for the given relation. For example, "institution.location", "university.founder",

and "university.president" can be identified as important contextual relations for "graduated from".

On the other hand, modeling relational paths captures the correlation between paths and the predicted relation, which can indicate important relational paths for the given relation. For example, ("schoolmate of", "graduated from") can be identified as an important relational path for "graduated from".

It is interesting to see that the explainability provided by relational paths is also connected to first-logic logical rules with the following form:

$$B_1(h, x_1) \land B_2(x_1, x_2) \land \cdots \land B_L(x_{L-1}, t) \Rightarrow r(h, t),$$

where $\bigwedge B_i$ is the conjunction of relations in a path and $r(h, t)$ is the predicted relation. The above example of relational path can therefore be written as the following rule:

$$(h, \text{ schoolmate of, } x) \land (x, \text{ graduated from, } t)$$
$$\Rightarrow (h, \text{ graduated from, } t).$$

Therefore, PATHCON can also be used to learn logical rules from KGs just as prior work [9, 12, 19, 36, 40].

## 3.6 Design Alternatives

Next we discuss several design alternatives for PATHCON. In our ablation experiments we will compare PATHCON with the following alternative implementations.

When modeling relational context, we propose two alternatives for context aggregator, instead of the Concatenation context aggregator in Eqs. (11) and (12):

**Mean context aggregator**. It takes the element-wise mean of the input vectors, followed by a nonlinear transformation function:

$$s_e^{i+1} = \sigma\left(\frac{1}{3}\big(m_v^i + m_u^i + s_e^i\big)W + b\right), \ v, u \in \mathcal{N}(e), \tag{17}$$

The output of Mean context aggregator is invariant to the permutation of its two input nodes, indicating that it treats the head and the tail equally in a triplet.

**Cross context aggregator**. It is inspired by combinatorial features in recommender systems [31], which measure the interaction of unit features (e.g., AND(gender=female, language=English)). Note that Mean and Concatenation context aggregator simply transform messages from two input nodes separately and add them up together, without modeling the interaction between them that might be useful for link prediction. In Cross context aggregator, we first calculate all element-level pairwise interactions between messages from the head and the tail:

$$m_v^i m_u^{i\top} = \begin{bmatrix} m_v^{i\,(1)} m_u^{i\,(1)} & \cdots & m_v^{i\,(1)} m_u^{i\,(d)} \\ \cdots & & \cdots \\ m_v^{i\,(d)} m_u^{i\,(1)} & \cdots & m_v^{i\,(d)} m_u^{i\,(d)} \end{bmatrix}, \tag{18}$$

where we use superscript with parentheses to indicate the element index and $d$ is the dimension of $m_v^i$ and $m_u^i$. Then we summarize all interactions together via flattening the interaction matrix to a vector then multiplied by a transformation matrix:

$$s_e^{i+1} = \sigma\left(\text{flatten}\big(m_v^i m_u^{i\top}\big)W_1^i + s_e^i W_2^i + b^i\right), \ v, u \in \mathcal{N}(e). \tag{19}$$

It is worth noting that Cross context aggregator preserves the order of input nodes.

**Learning path representation with RNN**. When modeling relational paths, recurrent neural network (RNN) can be used to learn the representation of relational path $P = (r_1, r_2, ...)$:

$$s_P = \text{RNN}(r_1, r_2, ...), \tag{20}$$

instead of directly assigning an embedding vector to $P$. The advantage of RNN against path embedding is that its number of parameters is fixed and does not depend on the number of relational paths. Another potential benefit is that RNN can hopefully capture the similarity among different relational paths.

**Mean path aggregator**. When calculating the final representation of relational paths for $(h, t)$ pair, we can also simply average all the representations of paths from $h$ to $t$ instead of the Attention path aggregator in Eqs. (13) and (14):

$$s_{h \to t} = \sum_{P \in \mathcal{P}_{h \to t}} s_P. \tag{21}$$

Mean path aggregator can be used in the case where representation of relational context is unavailable, since it does not require attention weights as input.

## 4 EXPERIMENTS

In this section, we evaluate the proposed PATHCON model, and present its performance on six KG datasets.

### 4.1 Experimental Setup

**Datasets**. We conduct experiments on five standard KG benchmarks: **FB15K**, **FB15K-237**, **WN18**, **WN18RR**, **NELL995**, and one KG dataset proposed by us: **DDB14**.

**FB15K** [4] is from Freebase [2], a large-scale KG of general human knowledge. **FB15k-237** [23] is a subset of FB15K where inverse relations are removed. **WN18** [4] contains conceptual-semantic and lexical relations among English words from WordNet [18]. **WN18RR** [6] is a subset of WN18 where inverse relations are removed. **NELL995** [33] is extracted from the 995th iteration of the NELL system [5] containing general knowledge.

In addition, we present a new dataset **DDB14** that is suitable for KG-related tasks. DDB14 is collected from Disease Database[6], which is a medical database containing terminologies and concepts such as diseases, symptoms, drugs, as well as their relationships. We randomly sample two subsets of 4,000 triplets from the original one as validation set and test set, respectively.

The statistics of the six datasets are summarized in Table 2. We also calculate and present the mean and variance of node degree distribution (i.e., $\mathbb{E}[d]$ and $\text{Var}[d]$) for each KG. It is clear that $\text{Var}[d]$ is large for all KGs, which empirically demonstrates that the complexity of relational message passing is fairly high, thus alternate relational message passing is necessary for real graphs.

**Baselines**. We compare PATHCON with several state-of-the-art models, including **TransE** [3], **ComplEx** [24], **DistMult** [35], **RotatE** [22], **SimplE** [13], **QuatE** [39], and **DRUM** [19]. The first six

---

[6] http://www.diseasedatabase.com

| | FB15K | FB15K-237 | WN18 | WN18RR | NELL995 | DDB14 |
|---|---|---|---|---|---|---|
| #nodes | 14,951 | 14,541 | 40,943 | 40,943 | 63,917 | 9,203 |
| #relations | 1,345 | 237 | 18 | 11 | 198 | 14 |
| #training | 483,142 | 272,115 | 141,442 | 86,835 | 137,465 | 36,561 |
| #validation | 50,000 | 17,535 | 5,000 | 3,034 | 5,000 | 4,000 |
| #test | 59,071 | 20,466 | 5,000 | 3,134 | 5,000 | 4,000 |
| $\mathbb{E}[d]$ | 64.6 | 37.4 | 6.9 | 4.2 | 4.3 | 7.9 |
| $\text{Var}[d]$ | 32,441.8 | 12,336.0 | 236.4 | 64.3 | 750.6 | 978.8 |

**Table 2: Statistics of all datasets. $\mathbb{E}[d]$ and $\text{Var}[d]$ are mean and variance of the node degree distribution, respectively.**

models are embedding-based methods, while DRUM only uses relational paths to make prediction. The implementation details of baselines (as well as our method) is provided in Appendix D.

We also conduct extensive ablation study and propose two reduced versions of our model, **CON** and **PATH**, which only use relational context and relational paths, respectively, to test the performance of the two components separately.

The number of parameters of each model on DDB14 are shown in Table 3. The result demonstrates that PATHCON is much more storage-efficient than embedding-based methods, since it does not need to calculate and store entity embeddings.

| Method | TransE | ComplEx | DisMult | RotatE | SimplE | QuatE | PATHCON |
|---|---|---|---|---|---|---|---|
| #param. | 3.7M | 7.4M | 3.7M | 7.4M | 7.4M | 14.7M | 0.06M |

**Table 3: Number of parameters of all models on DDB14.**

**Evaluation Protocol**. We evaluate all methods on relation prediction, i.e., for a given entity pair $(h, t)$ in the test set, we rank the ground-truth relation type $r$ against all other candidate relation types. It is worth noticing that most baselines are originally designed for head/tail prediction, therefore, their negative sampling strategy is to corrupt the head or the tail for a true triple $(h, r, t)$, i.e., replacing $h$ or $t$ with a randomly sampled entity $h'$ or $t'$ from KGs, and using $(h', r, t)$ or $(h, r, t')$ as the negative sample. In relation prediction, since the task is to predict the missing relation for a given pair $(h, t)$, we modify the negative sampling strategy accordingly by corrupting the relation $r$ of each true triplet $(h, r, t)$, and use $(h, r', t)$ as the negative sample where $r'$ is randomly sampled from the set of relation types. This new negative sampling strategy can indeed improve the performance of baselines in relation prediction.

We use **MRR** (mean reciprocal rank) and **Hit@1, 3** (hit ratio with cut-off values of 1 and 3) as evaluation metrics.

### 4.2 Main Results

**Comparison with baselines**. The results of relation prediction on all datasets are reported in Table 4. In general, our method outperforms all baselines on all datasets. Specifically, the absolute Hit@1 gain of PATHCON against the best baseline in relation prediction task are 0.2%, 0.6%, 0.9%, 16.7%, 6.3%, and 1.8% in the six datasets, respectively. The improvement is rather significant for WN18RR and NELL995, which are exactly the two most sparse KGs according to the average node degree shown in Table 2. This empirically demonstrates that PATHCON maintains great performance for sparse KGs, and this is probably because PATHCON has much fewer parameters

| | FB15K | | | FB15K-237 | | | WN18 | | | WN18RR | | | NELL995 | | | DDB14 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | MRR | Hit@1 | Hit@3 | MRR | Hit@1 | Hit@3 | MRR | Hit@1 | Hit@3 | MRR | Hit@1 | Hit@3 | MRR | Hit@1 | Hit@3 | MRR | Hit@1 | Hit@3 |
| TransE | 0.962 | 0.940 | 0.982 | 0.966 | 0.946 | 0.984 | 0.971 | 0.955 | 0.984 | 0.784 | 0.669 | 0.870 | _0.841_ | _0.781_ | _0.889_ | _0.966_ | _0.948_ | 0.980 |
| ComplEx | 0.901 | 0.844 | 0.952 | 0.924 | 0.879 | 0.970 | _0.985_ | _0.979_ | _0.991_ | 0.840 | 0.777 | 0.880 | 0.703 | 0.625 | 0.765 | 0.953 | 0.931 | 0.968 |
| DistMult | 0.661 | 0.439 | 0.868 | 0.875 | 0.806 | 0.936 | 0.786 | 0.584 | 0.987 | 0.847 | _0.787_ | 0.891 | 0.634 | 0.524 | 0.720 | 0.927 | 0.886 | 0.961 |
| RotatE | 0.979 | 0.967 | 0.986 | 0.970 | 0.951 | 0.980 | 0.984 | _0.979_ | 0.986 | 0.799 | 0.735 | 0.823 | 0.729 | 0.691 | 0.756 | 0.953 | 0.934 | 0.964 |
| SimplE | _0.983_ | _0.972_ | _0.991_ | 0.971 | 0.955 | 0.987 | 0.972 | 0.964 | 0.976 | 0.730 | 0.659 | 0.755 | 0.716 | 0.671 | 0.748 | 0.924 | 0.892 | 0.948 |
| QuatE | _0.983_ | _0.972_ | _0.991_ | _0.974_ | _0.958_ | _0.988_ | 0.981 | 0.975 | 0.983 | 0.823 | 0.767 | 0.852 | 0.752 | 0.706 | 0.783 | 0.946 | 0.922 | 0.962 |
| DRUM | 0.945 | 0.945 | 0.978 | 0.959 | 0.905 | 0.958 | 0.969 | 0.956 | 0.980 | _0.854_ | 0.778 | _0.912_ | 0.715 | 0.640 | 0.740 | 0.958 | 0.930 | _0.987_ |
| Con | 0.962 ± 0.000 | 0.934 ± 0.000 | 0.988 ± 0.000 | 0.978 ± 0.000 | 0.961 ± 0.001 | **0.995** ± 0.000 | 0.960 ± 0.002 | 0.927 ± 0.005 | 0.992 ± 0.001 | 0.943 ± 0.002 | 0.894 ± 0.004 | 0.993 ± 0.003 | 0.875 ± 0.003 | 0.815 ± 0.004 | 0.928 ± 0.003 | 0.977 ± 0.000 | 0.961 ± 0.001 | 0.994 ± 0.001 |
| Path | 0.937 ± 0.001 | 0.918 ± 0.001 | 0.951 ± 0.001 | 0.972 ± 0.001 | 0.957 ± 0.001 | 0.986 ± 0.001 | 0.981 ± 0.000 | 0.971 ± 0.005 | 0.989 ± 0.001 | 0.933 ± 0.001 | 0.897 ± 0.001 | 0.961 ± 0.001 | 0.737 ± 0.001 | 0.685 ± 0.002 | 0.764 ± 0.002 | 0.969 ± 0.001 | 0.948 ± 0.001 | 0.991 ± 0.000 |
| PathCon | **0.984** ± 0.001 | **0.974** ± 0.002 | **0.995** ± 0.001 | **0.979** ± 0.000 | **0.964** ± 0.001 | 0.994 ± 0.001 | **0.993** ± 0.001 | **0.988** ± 0.001 | **0.998** ± 0.000 | **0.974** ± 0.001 | **0.954** ± 0.002 | **0.994** ± 0.000 | **0.896** ± 0.001 | **0.844** ± 0.004 | **0.941** ± 0.004 | **0.980** ± 0.000 | **0.966** ± 0.001 | **0.995** ± 0.000 |

**Table 4: Results of relation prediction on all datasets. Best results are highlighted in bold, and best results of baselines are highlighted with underlines.**

than baselines and is less prone to overfitting. In contrast, performance gain of PathCon on FB15K is less significant, which may be because the density of FB15K is very high so that it is much easier for baselines to handle.

In addition, the results also demonstrate the stability of PathCon as we observe that most of the standard deviations are quite small.

Results in Tables 4 also show that, in many cases Con or Path can already beat most baselines. Combining relational context and relational paths together usually leads to even better performance.

**Inductive KG completion**. We also examine the performance of our method in inductive KG completion. We randomly sample a subset of nodes that appears in the test set, then remove these nodes along with their associated edges from the training set. The remaining training set is used to train the models, and we add back the removed edges during evaluation. The evaluation transforms from fully conductive to fully inductive when the ratio of removed nodes increases from 0 to 1. The results of PathCon, DistMult, and RotatE on relation prediction task are plotted in Figure 3. We observe that the performance of our method decreases slightly in fully inductive setting (from 0.954 to 0.922), while DistMult and RotatE fall to a "randomly guessing" level. This is because the two baselines are embedding-based models that rely on modeling node identity, while our method does not consider node identity thus being naturally generalizable to inductive KG completion.

## 4.3 Model Variants

**The number of context hops and maximum path length**. We investigate the sensitivity of our model to the number of context hops and maximum path length. We vary the two numbers from 0 to 4 (0 means the corresponding module is not used), and report the results of all combinations (without (0, 0)) on WN18RR in Figure 4. It is clear to see that increasing the number of context hops and maximum path length can significantly improve the result when they are small, which demonstrates that including more neighbor edges or counting longer paths does benefit the performance. But the marginal benefit is diminishing with the increase of layer numbers. Similar trend is observed on other datasets too.

**Context aggregators**. We study how different implementations of context aggregator affect the model performance. The results of Mean, Concat, and Cross context aggregator on four datasets are shown in Figure 5 (results on FB15K and WN18 are omitted as they are similar to FB15K-237 and WN18RR, respectively). The results show that Mean performs worst on all datasets, which indicates the importance of node orders when aggregating features from nodes to edges. It is also interesting to notice that the performance comparison between Concat and Cross varies on different datasets: Concat is better than Cross on NELL995 and is worse than Cross on WN18RR, while their performance is on par on FB15K-237 and DDB14. However, note that a significant defect of Cross is that it has much more parameters than Concat, which requires more running time and memory resource.

**Path representation types and path aggregators**. We implement four combinations of path representation types and path aggregators: Embedding+Mean, Embedding+Attention, RNN+Mean, and RNN+Attention, of which the results are presented in Figure 6. Different from context aggregators, results on the six datasets are similar for path representation types and path aggregators, so we only report the results on WN18RR. We find that Embedding is consistently better than RNN, which is probably because the length of relational paths are generally short (no more than 4 in our experiments), so RNN can hardly demonstrate its strength in modeling sequences. The results also show that Attention aggregator performs slightly better than Mean aggregator. This demonstrates that the contextual information of head and tail entities indeed helps identify the importance of relational paths.

**Initial edge features**. Here we examine three types of initial edge features: identity, BOW, and BERT embedding of relation types. We choose to test on NELL995 because its relation names consist of relatively more English words thus are semantically meaningful (e.g., "organization.headquartered.in.state.or.province"). The results are reported in Figure 7, which shows that BOW features are slightly better than identity, but BERT embeddings perform significantly worse than the other two. We attribute this finding to that BERT embeddings are better at identifying semantic relationship among relation types, but our model aims to learn the mapping from BERT embeddings of context/paths to the identity of predicted relation
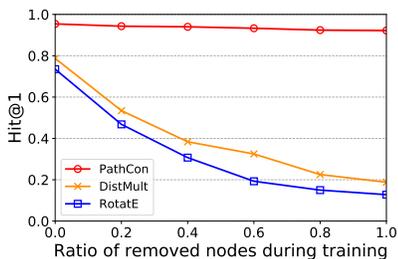
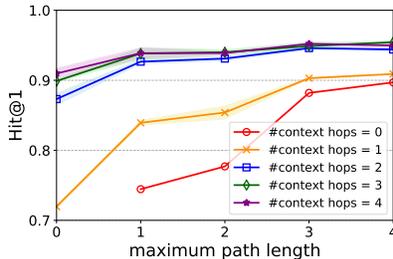Figure 3: Results of inductive KG completion on WN18RR.



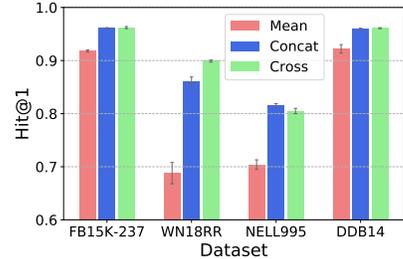Figure 4: Results of PATHCON with different hops/length on WN18RR.



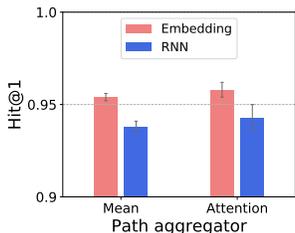Figure 5: Results of CON with different context aggregators.



Figure 6: Results of PATH-CON with different path representation types and path aggregators on WN18RR.
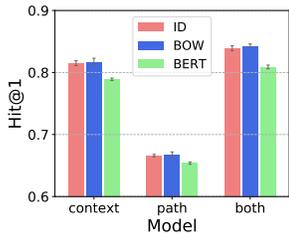


Figure 7: Results of CON, PATH, and PATHCON with different initial features of relations on NELL995.

types. In other words, BERT may perform better if the predicted relation types are also represented by BERT embeddings, so that this mapping is learned within the embedding space. We leave the exploration as future work.

### 4.4 Case Study on Model Explainabilty

We choose FB15K-237 and DDB14 as the datasets to show the explainability of PATHCON. The number of context hops is set to 1 and the maximum path length is set to 2. When training is completed, we choose three relations from each dataset and list the most important relational context/paths to them based on the transformation matrix of the context/path aggregator. The results are presented in Table 5, from which we find that most of the identified context/paths are logically meaningful. For example, "education campus of" can be inferred by "education institution in", and "is associated with" is found to be a transitive relation. In addition, more visualized results and discussion on DDB14 dataset are included in Appendix E.

## 5 RELATED WORK

### 5.1 Knowledge Graph Completion

KGs provide external information for a variety of downstream tasks such as recommender systems [27–29] and semantic analysis [26]. Most existing methods of KG completion are based on embeddings, which normally assign an embedding vector to each entity and relation in the continuous embedding space and train the embeddings based on the observed facts. One line of KG embedding methods is *translation-based*, which treat entities as points in a continuous

space and each relation translates the entity point. The objective is that the translated head entity should be close to the tail entity in real space [3], complex space [22], or quaternion space [39], which have shown capability to handle multiple relation patterns and achieve state-of-the-art result. Another line of work is *multi-linear* or *bilinear models*, where they calculate the semantic similarity by matrix or vector dot product in real [35] or complex space [24]. Besides, several embedding-based methods explore the architecture design that goes beyond point vectors [6, 21]. However, these embedding-based models fail to predict links in inductive setting, neither can they discover any rules that explain the prediction.

### 5.2 Graph Neural Networks

Existing GNNs generally follow the idea of neural message passing [10] that consists of two procedures: propagation and aggregation. Under this framework, several GNNs are proposed that take inspiration from convolutional neural networks [8, 11, 15, 25], recurrent neural networks [17], and recursive neural networks [1]. However, these methods use node-based message passing, while we propose passing messages based on edges in this work.

There are two GNN models conceptually connected to our idea of identifying relative position of nodes in a graph. DEGNN [16] captures the distance between the node set whose representation is to be learned and each node in the graph, which is used as extra node attributes or as controllers of message aggregation in GNNs. SEAL [38] labels nodes with their distance to two nodes $a$ and $b$ when predicting link existence between $(a, b)$. In contrast, we use relational paths to indicate the relative position of two nodes.

Researchers also tried to apply GNNs to knowledge graphs. For example, Schlichtkrull *et al.* [20] use GNNs to model the entities and relations in KGs, however, they are limited in that they did not consider the relational paths and cannot predict in inductive settings. Wang *et al.* [30, 32] use GNNs to learn entity embeddings in KGs, but their purpose is to use the learned embeddings to enhance the performance of recommender systems rather than KG completion.

## 6 CONCLUSION AND FUTURE WORK

We propose PATHCON for KG completion. PATHCON considers two types of subgraph structure in KGs, i.e., contextual relations of the head/tail entity and relational paths between head and tail entity. We show that both relational context and relational paths are critical to relation prediction, and they can be combined further to achieve

| | Predicted relation | Important relational context | Important relational paths |
|---|---|---|---|
| FB15K-237 | award winner | award honored for, award nominee | (award nominated for), (award winner, award category) |
| | film written by | film release region | (film edited by), (film crewmember) |
| | education campus of | education major field of study | (education institution in) |
| DDB14 | may cause | may cause, belongs to the drug family of | (is a risk factor for), (see also, may cause) |
| | is associated with | is associated with, is a risk factor for | (is associated with, is associated with) |
| | may be allelic with | may be allelic with, belong(s) to the category of | (may cause, may cause), (may be allelic with, may be allelic with) |

**Table 5: Examples of important context/paths identified by PathCon on FB15K-237 and DDB14.**

state-of-the-art performance. Moreover, PathCon is also shown to be inductive, storage-efficient, and explainable.

We point out four directions for future work. First, as we discussed in Remark 2, it is worth studying the empirical performance of PathCon on node-feature-aware KGs. Second, as we discussed in Section 4.3, designing a model that can better take advantage of pre-trained word embeddings is a promising direction; Third, it is worth investigating why RNN does not perform well, and whether we can model relational paths better; Last, it is interesting to study if the context representation and path representation can be assembled in a more principled way.

## REFERENCES

[1] Monica Bianchini, Marco Gori, and Franco Scarselli. 2001. Processing directed acyclic graphs with recursive neural networks. *IEEE Transactions on Neural Networks* 12, 6 (2001), 1464–1470.

[2] Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *SIGMOD*. 1247–1250.

[3] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. In *NeurIPS*. 2787–2795.

[4] Antoine Bordes, Jason Weston, Ronan Collobert, and Yoshua Bengio. 2011. Learning structured embeddings of knowledge bases. In *AAAI*. 301–306.

[5] Andrew Carlson, Justin Betteridge, Bryan Kisiel, Burr Settles, Estevam R Hruschka, and Tom M Mitchell. 2010. Toward an architecture for never-ending language learning. In *AAAI*. 1306–1313.

[6] Tim Dettmers, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. 2018. Convolutional 2d knowledge graph embeddings. In *AAAI*. 1811–1818.

[7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: pre-training of deep bidirectional transformers for language understanding. *arXiv preprint* (2018).

[8] David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P Adams. 2015. Convolutional networks on graphs for learning molecular fingerprints. In *NeurIPS*. 2224–2232.

[9] Luis Galárraga, Christina Teflioudi, Katja Hose, and Fabian M Suchanek. 2015. Fast rule mining in ontological knowledge bases with amie+. *The VLDB Journal* 24, 6 (2015), 707–730.

[10] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. 2017. Neural message passing for quantum chemistry. In *ICML*. 1263–1272.

[11] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *NeurIPS*. 1024–1034.

[12] Vinh Thinh Ho, Daria Stepanova, Mohamed H Gad-Elrab, Evgeny Kharlamov, and Gerhard Weikum. 2018. Rule learning from knowledge graphs guided by embedding models. In *International Semantic Web Conference*. Springer, 72–90.

[13] Seyed Mehran Kazemi and David Poole. 2018. Simple embedding for link prediction in knowledge graphs. In *NeurIPS*. 4284–4295.

[14] Diederik P Kingma and Jimmy Ba. 2015. Adam: a method for stochastic optimization. In *ICLR*.

[15] Thomas N Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *ICLR*.

[16] Pan Li, Yanbang Wang, Hongwei Wang, and Jure Leskovec. 2020. Distance Encoding: Design Provably More Powerful Neural Networks for Graph Representation Learning. In *NeurIPS*.

[17] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. 2016. Gated graph sequence neural networks. In *ICLR*.

[18] George A Miller. 1995. Wordnet: a lexical database for English. *Commun. ACM* 38, 11 (1995), 39–41.

[19] Ali Sadeghian, Mohammadreza Armandpour, Patrick Ding, and Daisy Zhe Wang. 2019. Drum: end-to-end differentiable rule mining on knowledge graphs. In *NeurIPS*. 15321–15331.

[20] Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. 2018. Modeling relational data with graph convolutional networks. In *European Semantic Web Conference*. Springer, 593–607.

[21] Richard Socher, Danqi Chen, Christopher D Manning, and Andrew Ng. 2013. Reasoning with neural tensor networks for knowledge base completion. In *NeurIPS*. 926–934.

[22] Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. 2019. Rotate: knowledge graph embedding by relational rotation in complex space. In *ICLR*.

[23] Kristina Toutanova and Danqi Chen. 2015. Observed versus latent features for knowledge base and text inference. In *Workshop on Continuous Vector Space Models and their Compositionality*. 57–66.

[24] Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. 2016. Complex embeddings for simple link prediction. In *ICML*.

[25] Hongwei Wang and Jure Leskovec. 2020. Unifying graph convolutional neural networks and label propagation. *arXiv preprint* (2020).

[26] Hongwei Wang, Fuzheng Zhang, Min Hou, Xing Xie, Minyi Guo, and Qi Liu. 2018. Shine: Signed heterogeneous information network embedding for sentiment link prediction. In *WSDM*. 592–600.

[27] Hongwei Wang, Fuzheng Zhang, Jialin Wang, Miao Zhao, Wenjie Li, Xing Xie, and Minyi Guo. 2018. Ripplenet: Propagating user preferences on the knowledge graph for recommender systems. In *CIKM*. 417–426.

[28] Hongwei Wang, Fuzheng Zhang, Jialin Wang, Miao Zhao, Wenjie Li, Xing Xie, and Minyi Guo. 2019. Exploring high-order user preference on the knowledge graph for recommender systems. *ACM Transactions on Information Systems* 37, 3 (2019), 1–26.

[29] Hongwei Wang, Fuzheng Zhang, Xing Xie, and Minyi Guo. 2018. DKN: Deep knowledge-aware network for news recommendation. In *WWW*. 1835–1844.

[30] Hongwei Wang, Fuzheng Zhang, Mengdi Zhang, Jure Leskovec, Miao Zhao, Wenjie Li, and Zhongyuan Wang. 2019. Knowledge-aware graph neural networks with label smoothness regularization for recommender systems. In *KDD*. 968–977.

[31] Hongwei Wang, Fuzheng Zhang, Miao Zhao, Wenjie Li, Xing Xie, and Minyi Guo. 2019. Multi-task feature learning for knowledge graph enhanced recommendation. In *WWW*. 2000–2010.

[32] Hongwei Wang, Miao Zhao, Xing Xie, Wenjie Li, and Minyi Guo. 2019. Knowledge graph convolutional networks for recommender systems. In *WWW*.

[33] Wenhan Xiong, Thien Hoang, and William Yang Wang. 2017. Deeppath: a reinforcement learning method for knowledge graph reasoning. In *EMNLP*. 564–573.

[34] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2019. How powerful are graph neural networks?. In *ICLR*.

[35] Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. 2015. Embedding entities and relations for learning and inference in knowledge bases. In *ICLR*.

[36] Fan Yang, Zhilin Yang, and William W Cohen. 2017. Differentiable learning of logical rules for knowledge base reasoning. In *NeurIPS*. 2319–2328.

[37] Jonathan S Yedidia, William T Freeman, and Yair Weiss. 2003. Understanding belief propagation and its generalizations. *Exploring artificial intelligence in the new millennium* 8 (2003), 236–239.

[38] Muhan Zhang and Yixin Chen. 2018. Link prediction based on graph neural networks. In *NeurIPS*. 5165–5175.

[39] Shuai Zhang, Yi Tay, Lina Yao, and Qi Liu. 2019. Quaternion knowledge graph embeddings. In *NeurIPS*. 2731–2741.

[40] Wen Zhang, Bibek Paudel, Liang Wang, Jiaoyan Chen, Hai Zhu, Wei Zhang, Abraham Bernstein, and Huajun Chen. 2019. Iteratively learning embeddings and rules for knowledge graph reasoning. In *WWW*. 2366–2377.

# APPENDIX

## A  Proof of Theorem 1

PROOF. In each iteration of node-based message passing:

The aggregation (Eq. (3)) is performed for $N$ times, and each aggregation takes $\mathbb{E}[d] = \frac{2M}{N}$ elements as input in expectation, where $\mathbb{E}[d]$ is the expected node degree. Therefore, the expected cost of aggregation in each iteration is $N \cdot \mathbb{E}[d] = 2M$;

The update (Eq. (4)) is performed for $N$ times, and each update takes 2 elements as input. Therefore, the cost of update in each iteration is $2N$.

In conclusion, the expected cost of node-based message passing in each iteration is $2M + 2N$. □

## B  Proof of Theorem 2

For relational message passing, it actually passes messages on the *line graph* of the original graph. The line graph of a given graph $\mathcal{G}$, denoted by $L(\mathcal{G})$, is a graph such that each node of $L(\mathcal{G})$ represents an edge of $\mathcal{G}$, and two nodes of $L(\mathcal{G})$ are adjacent if and only if their corresponding edges share a common endpoint in $\mathcal{G}$. We show by the following lemma that the line graph is much *larger* and *denser* than the original graph:

LEMMA 1. *The number of nodes in line graph $L(\mathcal{G})$ is $M$, and the expected node degree of $L(\mathcal{G})$ is*

$$\mathbb{E}_{L(\mathcal{G})}[d] = \frac{N \cdot \mathrm{Var}_{\mathcal{G}}[d]}{M} + \frac{4M}{N} - 2, \qquad (22)$$

*where $\mathrm{Var}_{\mathcal{G}}[d]$ is the variance of node degrees in $\mathcal{G}$.*

PROOF. It is clear that the number of nodes in line graph $L(\mathcal{G})$ is $M$ because each node in $L(\mathcal{G})$ corresponds to an edge in $\mathcal{G}$. We now prove that the expected node degree of $L(\mathcal{G})$ is $\mathbb{E}_{L(\mathcal{G})}[d] = \frac{N \cdot \mathrm{Var}_{\mathcal{G}}[d]}{M} + \frac{4M}{N} - 2$.

Let's first count the number of edges in $L(\mathcal{G})$. According to the definition of line graph, each edge in $L(\mathcal{G})$ corresponds to an unordered pair of edges in $\mathcal{G}$ connecting to a same node; On the other hand, each unordered pair of edges in $\mathcal{G}$ that connect to a same node also determines an edge in $L(\mathcal{G})$. Therefore, the number of edges in $L(\mathcal{G})$ equals the number of all unordered pairs of edges connecting to a same node:

$$\# \ edges \ in \ L(\mathcal{G}) = \sum_i \binom{d_i}{2} = \sum_i \frac{d_i(d_i - 1)}{2} = \frac{1}{2}\sum_i d_i^2 - M,$$

where $d_i$ is the degree of node $v_i$ in $\mathcal{G}$ and $M = 2\sum_i d_i$ is the number of edges. Then the the expected node degree of $L(\mathcal{G})$ is

$$\begin{aligned}
\mathbb{E}_{L(\mathcal{G})}[d] &= 2 \cdot \frac{\# \ edges \ in \ L(\mathcal{G})}{\# \ nodes \ in \ L(\mathcal{G})} = \frac{\sum_i d_i^2 - 2M}{M} \\
&= \frac{N \cdot \mathbb{E}_{\mathcal{G}}[d^2]}{M} - 2 = \frac{N \left( \mathrm{Var}_{\mathcal{G}}[d] + \mathbb{E}_{\mathcal{G}}^2[d] \right)}{M} - 2 \\
&= \frac{N \cdot \mathrm{Var}_{\mathcal{G}}[d] + N \left( \frac{2M}{N} \right)^2}{M} - 2 \\
&= \frac{N \cdot \mathrm{Var}_{\mathcal{G}}[d]}{M} + \frac{4M}{N} - 2.
\end{aligned}$$

□

From Lemma 1 it is clear to see that $\mathbb{E}_{L(\mathcal{G})}[d]$ is at least twice of $\mathbb{E}_{\mathcal{G}}[d] = \frac{2M}{N}$, i.e. the expected node degree of the original graph $\mathcal{G}$, since $\mathrm{Var}_{\mathcal{G}}[d] \geq 0$ ($-2$ is omitted). Unfortunately, in real-world graphs (including KGs), node degrees vary significantly, and they typically follow the power law distribution whose variance is extremely large due to the long tail (this is empirically justified in Table 2, as we can see that $\mathrm{Var}_{\mathcal{G}}[d]$ is quite large for all KGs). This means that $\mathbb{E}_{L(\mathcal{G})}[d] \gg \mathbb{E}_{\mathcal{G}}[d]$ in practice. On the other hand, the number of nodes in $L(\mathcal{G})$ (which is $M$) is also far larger than the number of nodes in $\mathcal{G}$ (which is $N$). Therefore, $L(\mathcal{G})$ is generally much larger and denser than its original graph $\mathcal{G}$. Based on Lemma 1, Theorem 2 is proven as follows:

PROOF. In each iteration of relational message passing:

The aggregation (Eq. (5)) is performed for $M$ times, and each aggregation takes $\mathbb{E}_{L(\mathcal{G})}[d] = \frac{N \cdot \mathrm{Var}_{\mathcal{G}}[d]}{M} + \frac{4M}{N} - 2$ elements as input in expectation. So the expected cost of aggregation in each iteration is $M \cdot \mathbb{E}_{L(\mathcal{G})}[d] = N \cdot \mathrm{Var}_{\mathcal{G}}[d] + \frac{4M^2}{N} - 2M$;

The update ((Eq. (6))) is performed for $M$ times, and each update takes 2 elements as input. Therefore, the cost of update in each iteration is $2M$.

In conclusion, the expected cost of relational message passing in each iteration is $N \cdot \mathrm{Var}_{\mathcal{G}}[d] + \frac{4M^2}{N}$. □

## C  Proof of Theorem 3

PROOF. In each iteration of alternate relational message passing:

The edge-to-node aggregation operation (Eq. (7)) is performed for $N$ times, and each aggregation takes $\mathbb{E}[d] = \frac{2M}{N}$ elements as input in expectation. Therefore, the expected cost of edge-to-node aggregation in each iteration is $N \cdot \mathbb{E}[d] = 2M$;

The node-to-edge aggregation (Eq. (8)) is performed for $M$ times, and each aggregation takes 2 elements as input. So the cost of node-to-edge aggregation in each iteration is $2M$;

The update (Eq. (9)) is performed for $M$ times, and each update takes 2 elements as input. Therefore, the cost of update in each iteration is $2M$.

In conclusion, the expected cost of alternate relational message passing in each iteration is $6M$. □

## D  Implementation Details

**Baselines**. The implementation code of TransE, DistMult, ComplEx, and RotatE comes from https://github.com/DeepGraphLearning/KnowledgeGraphEmbedding; the implementation code of SimplE is at https://github.com/baharefatemi/SimplE; the implementation code of QuatE is at https://github.com/cheungdaven/QuatE, and we use QuatE[2] (QuatE without type constraints) here; the implementation code of DRUM is at https://github.com/alisadeghian/DRUM. For fair comparison, the embedding dimension for all the baselines are set to 400. We train each baseline for 1,000 epochs, and report the test result when the result on validation set is optimal.

**Our method**. Our proposed method is implemented in TensorFlow and trained on single GPU. We use Adam [14] as the optimizer with learning rate of 0.005. L2 regularization is used to prevent overfitting and the weight of L2 loss term is $10^{-7}$. Batch size is 128, the number of epochs is 20, and the dimension of all hidden
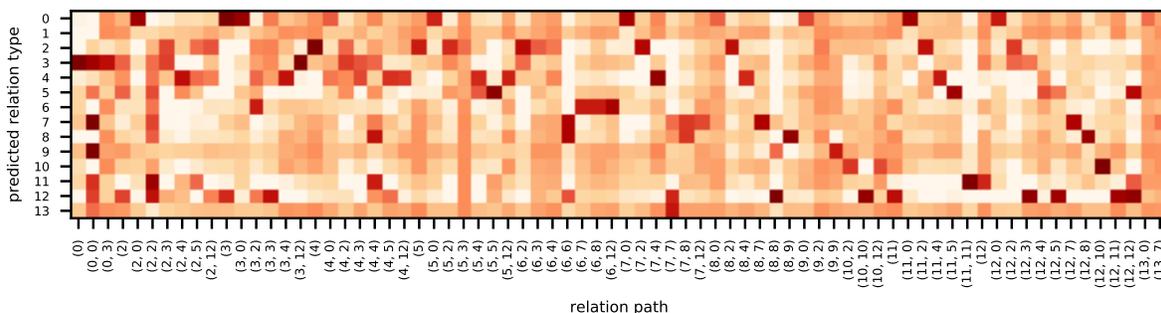
Figure 8: The learned correlation between all relational paths with length ≤ 2 and the predicted relations on DDB14.

| | FB15K | FB15K-237 | WN18 | WN18RR | NELL995 | DDB14 |
|---|---|---|---|---|---|---|
| #context hops | 2 | 2 | 3 | 3 | 2 | 3 |
| Max. path len. | 2 | 3 | 3 | 4 | 3 | 4 |

Table 6: Dataset-specific hyper-parameter settings: the number of context hops and the maximum path length.

states is 64. Initial relation features are set as their identities, while BOW/BERT features are studied in Section 4.3. The above settings are determined by optimizing the classification accuracy on the validation set of WN18RR, and kept unchanged for all datasets.

During experiments we find that performance of different number of context hops and the maximum path length largely depends on datasets, so these hyper-parameters are tuned separately for each dataset. We present their default settings in Table 6, and search spaces of hyper-parameters as follows:

- Dimension of hidden states: $\{8, 16, 32, 64\}$;
- Weight of L2 loss term: $\{10^{-8}, 10^{-7}, 10^{-6}, 10^{-5}\}$;
- Learning rate: $\{0.001, 0.005, 0.01, 0.05, 0.1\}$;
- The number of context hops: $\{1, 2, 3, 4\}$;
- Maximum path length: $\{1, 2, 3, 4\}$.

Each experiment of PathCon is repeated for three times. We report average performance and standard deviation as the results.

## E  More Results of Explainability on DDB14

After training on DDB14, we print out the transformation matrix of the context aggregator and the path aggregator in PathCon, and the results are shown as heat maps in Figures 9 and 8, respectively. The degree of darkness of an entry in Figure 9 (Figure 8) denotes the strength of correlation between the existence of a contextual relation (a relational path) and a predicted relation. Relation IDs as well as their meanings are listed as follows for readers' reference:

| | |
|---|---|
| 0: belong(s) to the category of | 7: interacts with |
| 1: is a category subset of | 8: belongs to the drug family of |
| 2: may cause | 9: belongs to drug super-family |
| 3: is a subtype of | 10: is a vector for |
| 4: is a risk factor for | 11: may be allelic with |
| 5: is associated with | 12: see also |
| 6: may contraindicate | 13: is an ingredient of |

Figure 9 shows that most of large values are distributed along the diagonal. This is in accordance with our intuition, for example, if we want to predict the relation for pair $(h, ?, t)$ and we observe
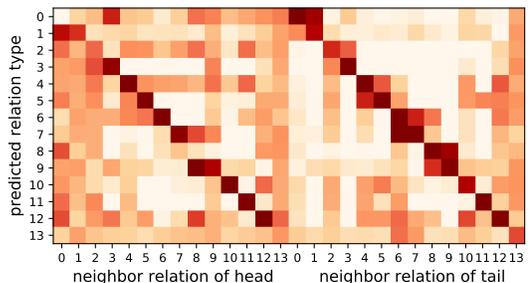


Figure 9: The learned correlation between the contextual relations of head/tail and the predicted relations on DDB14.

that $h$ appears in another triplet $(h, \text{is a risk factor for}, t')$, then we know that the type of $h$ is risk factor and it is likely to be a risk factor of other entities in the KG. Therefore, "?" are more likely to be "is a risk factor for" than "belongs to the drug family of" since $h$ is not a drug. In addition, we also find some large values that are not in the diagonal, e.g., (belongs to the drug family of, belongs to the drug super-family) and (may contraindicate, interacts with).

We also have some interesting findings from Figure 8. First, we find that many rules from Figure 8 is with the form:

$$(a, \text{see also}, b) \land (b, \text{R}, c) \Rightarrow (a, \text{R}, c),$$

where R is a relation type in the KG. These rules are indeed meaningful because $(a, \text{see also}, b)$ means $a$ and $b$ are equivalent thus can interchange with each other.

We also find PathCon learns rules that show the relation type is transitive, for example:

$(a, \text{is associated with}, b) \land (b, \text{is associated with}, c)$
$\Rightarrow (a, \text{is associated with}, c)$;

$(a, \text{may be allelic with}, b) \land (b, \text{may be allelic with}, c)$
$\Rightarrow (a, \text{may be allelic with}, c)$.

Other interesting rules learned by PathCon include:

$(a, \text{belong(s) to the category of}, b) \Rightarrow (a, \text{is a subtype of}, b)$;

$(a, \text{is a risk factor for}, b) \Rightarrow (a, \text{may cause}, b)$;

$(a, \text{may cause}, c) \land (b, \text{may cause}, c) \Rightarrow (a, \text{may be allelic with}, b)$;

$(a, \text{is a risk factor for}, c) \land (b, \text{is a risk factor for}, c)$
$\Rightarrow (a, \text{may be allelic with}, b)$.