

Combining Rules and Embeddings via Neuro-Symbolic AI for Knowledge Base Completion

Prithviraj Sen, Breno W. S. R. de Carvalho, Ibrahim Abdelaziz, Pavan Kapanipathi,
Francois Luus, Salim Roukos, Alexander Gray

IBM Research

Abstract

Recent interest in Knowledge Base Completion (KBC) has led to a plethora of approaches based on reinforcement learning, inductive logic programming and graph embeddings. In particular, rule-based KBC has led to interpretable rules while being comparable in performance with graph embeddings. Even within rule-based KBC, there exist different approaches that lead to rules of varying quality and previous work has not always been precise in highlighting these differences. Another issue that plagues most rule-based KBC is the non-uniformity of relation paths: some relation sequences occur in very few paths while others appear very frequently. In this paper, we show that not all rule-based KBC models are the same and propose two distinct approaches that learn in one case: 1) a mixture of relations and the other 2) a mixture of paths. When implemented on top of neuro-symbolic AI, which learns rules by extending Boolean logic to real-valued logic, the latter model leads to superior KBC accuracy outperforming state-of-the-art rule-based KBC by 2-10% in terms of mean reciprocal rank. Furthermore, to address the non-uniformity of relation paths, we combine rule-based KBC with graph embeddings thus improving our results even further and achieving the best of both worlds.

1 Introduction

A number of approaches have been proposed for knowledge base completion (KBC), a popular task that addresses the inherent incompleteness of Knowledge Graphs (KG) (Bollacker et al. 2008; Rebele et al. 2016). Compared to embeddings-based techniques (Sun et al. 2019a; Lacroix, Usunier, and Obozinski 2018), rule learning techniques for KBC can handle cold-start challenges for new entities whose embeddings are not available (inductive setting) (Yang, Yang, and Cohen 2017a; Sadeghian et al. 2019a; Das et al. 2018), and learn multi-hop, human-interpretable, first order logic rules that enables complex reasoning over KGs. For example, the following rule (from Freebase) infers a person’s nationality given her/his place of birth and its corresponding country:

$$\forall P, N \exists L: \text{nationality}(P, N) \leftarrow \text{bornIn}(P, L) \wedge \text{partOf}(L, N)$$

There exist two kinds of prevalent rule learning approaches based on their mechanism to select relations. The first approach, denoted *Chain of Mixtures* (CM), represents each relation hop in the body of the rule, towards the right of the implication \leftarrow symbol, as a mixture of relations. With

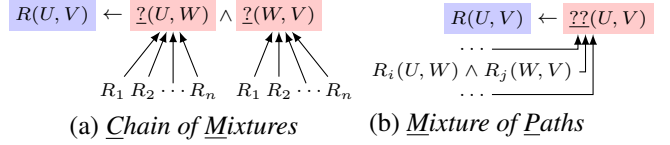


Figure 1: Rule-based KBC approaches: (a) CM, (b) MP.

close ties to NeuralLP (Yang, Yang, and Cohen 2017a), DRUM (Sadeghian et al. 2019a), CM is depicted in Figure 1 (a) where we learn a rule with two relations in the body. The second approach, denoted *Mixture of Paths* (MP), learns relation paths with close ties to MINERVA (Das et al. 2018), RNNLogic (Qu et al. 2020). Figure 1 (b) depicts MP where the body of the rule is defined as a mixture of all possible length 2 relation paths in the KG. Crucially, both approaches, which rely on recurrent neural networks (RNN), use increasingly complex training algorithms. More precisely, NeuralLP and DRUM generate mixture probabilities using a long short-term memory (LSTM) (Hochreiter and Schmidhuber 1997) and bi-directional LSTM, respectively. RNNLogic, the latest in the line of MP-based works, samples *sets* of rules from a multinomial distribution defined by its RNN’s latent vectors and learning this RNN requires sampling from intractable posteriors. While RNN training has improved significantly (Pascanu, Mikolov, and Bengio 2013), we are still faced with issues when these are used to learn long-range dependencies (Trinh et al. 2018). Suffices to say that it is unclear whether the current RNN-based rule-based KBC models are in their simplest form.

Instead of following the current trend of devising increasingly complex rule-based KBC approaches, we ask whether it is possible to devise an approach without resorting to RNNs? To this end, we propose to utilize Logical Neural Networks (LNN) (Riegel et al. 2020), a member of Neuro-Symbolic AI family (NeSy), to devise both a CM-based (*LNN-CM*) and a MP-based approach (*LNN-MP*). LNN is an extension of Boolean logic to the real-valued domain and is differentiable thus enabling us to learn rules end-to-end via gradient-based optimization. Another advantage of LNN is that while other members of NeSy harbor tenuous connections to Boolean logic’s semantics, e.g., Dong et al. (2019), LNN shares strong ties thus making the learned rules fully interpretable. More

importantly, our *RNN-less* approach is conceptually simpler to understand and reason about than recent rule-based KBC approaches while still performing at par, if not better.

One shortcoming of MP (path-based rule learning) is that it suffers from sparsity. Since not all relation paths are equally prevalent in the KG, our goal is to guide the learning process towards more prevalent, effective relation paths as opposed to paths that appear less frequently. We show how to use pre-trained knowledge graph embeddings (KGE) to overcome such issues and learn more effective rules. While RNNLogic (Qu et al. 2021) scores paths using RotatE (Sun et al. 2019a), a specific KGE, we show how to use any of the vast array of KGEs available in the literature. Our contributions are:

- We propose to implement with logical neural networks (LNN) two simple approaches for rule-based KBC. LNN is a differentiable framework for real-valued logic that has strong connections to Boolean logic semantics.
- To address the non-uniform distribution of relation paths that is likely to exist in any KG, we propose a simple approach that combines rule-based KBC with any KGE in order to reap the complementary benefits of both worlds.
- Our experiments on 4 different KBC benchmarks show that we: (a) outperform other rule-based baselines by 2%-10% (mean reciprocal rank); (b) are comparable to state-of-the-art rule with embedding-based technique RNNLogic.

2 Preliminaries: Logical Neural Networks

We begin with a brief overview of *Logical Neural Networks* (LNN) (Riegel et al. 2020), a differentiable extension of Boolean logic that can learn rules end-to-end while still maintaining strong connections to Boolean logic semantics. In particular, LNN extends propositional Boolean operators with learnable parameters that allows a better fit to data. We next review operators such as LNN- \wedge , useful for modeling \wedge in the example rule (Section 1), and LNN-pred, useful for expressing mixtures used in both CM and MP (Figure 1).

2.1 Propositional LNN Operators

To address the non-differentiability of classical Boolean logic, previous work has resorted to t -norms from fuzzy logic which are differentiable but lack parameters and thus cannot adapt to data. For instance, NeuralLP (Yang, Yang, and Cohen 2017a) uses product t -norm defined as $x \wedge y \equiv xy$ instead of Boolean conjunction. In contrast, LNN conjunction includes parameters that can not only express conjunctive semantics over real-valued logic but also fit the data better. The following LNN- \wedge operator extends the Łukasiewicz t -norm, defined as $x \wedge y \equiv \max(0, x + y - 1)$, with parameters:

$$\text{LNN-}\wedge(\mathbf{x}; \beta, \mathbf{w}) \equiv \max[0, \min\{1, \beta - \mathbf{w}^\top(\mathbf{1} - \mathbf{x})\}]$$

$$\text{subject to: } \beta \mathbf{1} - \alpha \mathbf{w} \leq (1 - \alpha) \mathbf{1} \quad (1)$$

$$\beta - (1 - \alpha) \mathbf{1}^\top \mathbf{w} \geq \alpha \quad (2)$$

$$\mathbf{w} \geq \mathbf{0}$$

where \mathbf{w} denotes weights and β denotes bias, both learnable parameters, \mathbf{x} denotes a vector of inputs (in the case of *binary* conjunction $|\mathbf{x}| = 2$), $\mathbf{0}$ ($\mathbf{1}$) denote a vector of 0s (1s), and

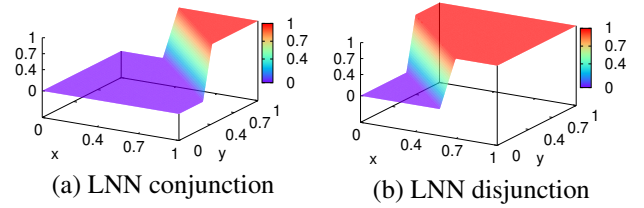


Figure 2: (a) LNN- \wedge and (b) LNN- \vee learned with $\alpha = 0.7$.

α denotes a hyperparameter. Recall that, a (binary) Boolean conjunction operator returns *true* (1) if both inputs are *true* (1) and *false* (0) otherwise. Intuitively, output is “high” if both inputs are “high”, otherwise the output is “low”. There are 4 steps to understanding how LNN extends this behavior to real-valued logic. 1) Given variable x whose truth value lies in the range $[0, 1]$, LNN utilizes $\alpha > \frac{1}{2}$ to define “low” as $x \in [0, 1 - \alpha]$ and “high” as $x \in [\alpha, 1]$. 2) Note that, the non-negativity constraint enforced on \mathbf{w} in LNN- \wedge ensures that it is a monotonically increasing function with respect to input \mathbf{x} . 3) Given its monotonicity, we can ensure that LNN- \wedge returns a “high” value when all entries in \mathbf{x} are “high” by simply constraining the output to be “high” when $\mathbf{x} = \alpha \mathbf{1}$ (Equation 2). 4) Similarly, to ensure that the output is “low” if any input entry is “low”, i.e. $1 - \alpha$, Equation 1 constrains the output to be low when all but one entries in \mathbf{x} is 1. This constrained formulation of LNN- \wedge ensures that conjunctive semantics is never lost during the learning process while still providing a better fit.

Figure 2 (a) shows a learned (binary) LNN- \wedge ($\alpha = 0.7$). Notice how, the output (z -axis) is close to 0 when either input x or y is “low” ($\in [0, 0.3]$) and jumps to 1 when both are “high” ($\in [0.7, 1]$), thus precisely capturing the semantics of logical conjunction. Note that, \max and \min have subgradients available thus making LNN- \wedge amenable to training via gradient-based optimization. Just to contrast with LNN disjunction (\vee) that is defined as $1 - \text{LNN-}\wedge(1 - \mathbf{x}; \beta, \mathbf{w})$, Figure 2 (b) shows a learned LNN- \vee ($\alpha = 0.7$). In contrast to LNN- \wedge , in this case, output is “low” when both $x, y \in [0, 0.3]$ and jumps to 1 when either of them are “high” ($\in [0.7, 1]$) thus capturing semantics of Boolean disjunction.

Another operator we will find useful is LNN-pred. In the next section, we show how to use this to express mixtures, either over relations (CM) or relation paths (MP). LNN-pred is a simple operator with one non-negative weight per input:

$$\text{LNN-pred}(\mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \mathbf{x} \text{ subject to } \mathbf{w} \geq \mathbf{0}, \mathbf{w}^\top \mathbf{1} = 1$$

where \mathbf{x} denotes (a vector of) inputs and \mathbf{w} denotes non-negative, learnable parameters constrained to sum to 1.

2.2 Training LNNs

While LNN operators are amenable to gradient-based learning, one complication we are yet to address is implementing *constrained* optimization to learn LNN parameters. Fortunately, there exist approaches that can convert any system of linear constraints (including equalities and inequalities) into a sequence of differentiable operations such that we can sample parameters directly from the feasible set (Frerix,

Nießner, and Cremer 2020) which is what we use to train all our LNNs. We also refer the interested reader to Riegel et al. (2020) that describes additional LNN training algorithms.

3 Learning LNN Chain Rules for KBC

3.1 Notation and Problem Definition

Let $\mathcal{G} = \langle \mathcal{V}, \mathcal{R}, \mathcal{E} \rangle$ denote a knowledge graph (KG) such that edge $e \in \mathcal{E}$ comprises a triple $\langle h, r, t \rangle$ where $h, t \in \mathcal{V}$ denote source and destination vertices, and $r \in \mathcal{R}$ denotes a relation. In this work, we focus on predicting destinations, i.e., given query $\langle h, r, ? \rangle$ predict the answer t . Following previous work (Yang, Yang, and Cohen 2017a), we learn *chain rules* in first-order logic (also called open path rules) that chains together multiple relations from \mathcal{R} to model the given relation:

$$\begin{aligned} \forall X_0, X_m \exists X_1, \dots, X_{m-1} : \\ r_0(X_0, X_m) \leftarrow r_1(X_0, X_1) \wedge \dots \wedge r_m(X_{m-1}, X_m) \\ r_i \in \mathcal{R} \forall i = 0, \dots, m \end{aligned} \quad (3)$$

where the *head* r_0 denotes the relation being modeled, r_1, \dots, r_m form the relations in the *body*, and X_0, \dots, X_m denote logical constants that can take values from \mathcal{V} . Note that, relations can repeat within the body, i.e., $r_i = r_j, i \neq j$. Furthermore, we allow r_0 to appear within the body which leads to a recursive rule. Chain rules are closely connected to multi-hop paths. Essentially, the above rule claims that $\langle u, r_0, v \rangle$ exists *if* there exists an m -length path $p = u \xrightarrow{r_1} \dots \xrightarrow{r_i} \dots \xrightarrow{r_m} v$ connecting $X_0 = u$ and $X_m = v$. Given such a multi-hop path p , we refer to the sequence of relations r_1, \dots, r_m as its *relation path*. Furthermore, given an m -length relation path $\mathbf{r} = r_1, \dots, r_m$, $\mathcal{P}_{\mathbf{r}}(u, v)$ denotes the *set* of all paths connecting u to v in \mathcal{G} via relation path \mathbf{r} .

Our goal is to learn to predict destinations for each $r \in \mathcal{R}$, however, in interest of keeping the approach simple we pose each relation-specific learning task in isolation. Given that standard KBC benchmarks (such as WN18RR) consist of sparse graphs, following previous work (Yang, Yang, and Cohen 2017a), we too introduce inverse relations. In other words, for each $r \in \mathcal{R}$ we introduce a new relation r^{-1} by adding for each $\langle h, r, t \rangle \in \mathcal{E}$ a new triple $\langle t, r^{-1}, h \rangle$ to \mathcal{G} . We refer to the augmented set of relations, including the inverse relations, as \mathcal{R}^+ . Learning to predict destinations for a given relation $r \in \mathcal{R}^+$ is essentially a binary classification task where $\langle h, r, t \rangle \in \mathcal{E}$ and $\langle h, r, t \rangle \notin \mathcal{E}, \forall h, t \in \mathcal{V}$, denote positive and negative examples, respectively. The approaches we present next differ in how they score a triple $\langle h, r, t \rangle$ which in turn depends on paths $\mathcal{P}_{\mathbf{r}}(u, v)$. Since we are interested in using such models to score unseen triples, we need to remove the triple from \mathcal{G} during training because, as mentioned earlier, we use triples from \mathcal{G} as positive examples. More precisely, when scoring a triple $\langle h, r, t \rangle$ during training we always compute $\mathcal{P}_{\mathbf{r}}(u, v)$ not from \mathcal{E} but from $\mathcal{E} \setminus \{ \langle h, r, t \rangle, \langle t, r^{-1}, h \rangle \}$, i.e., we always remove 2 edges: the triple itself and its corresponding inverse triple.

3.2 Chains of LNN-pred Mixtures

Given \mathcal{G} , a relation to model $r \in \mathcal{R}$, and a user-defined rule-length m , our first approach for learning chain rules uses m

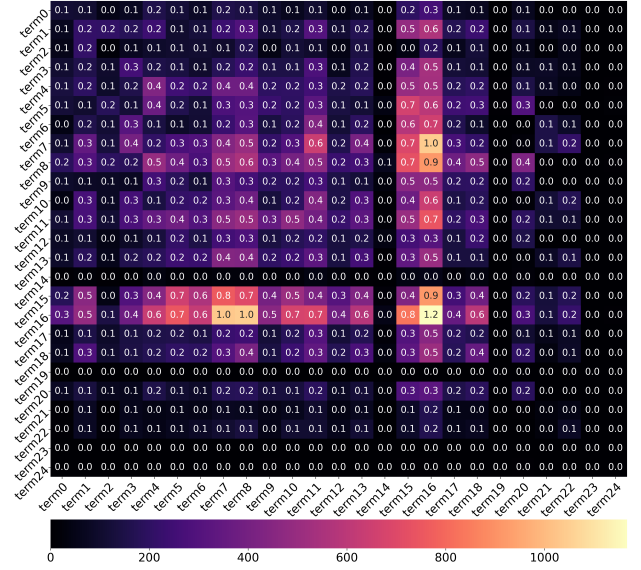


Figure 3: Kinship path counts (in thousands).

LNN-pred operators and one m -ary LNN- \wedge operator. Score for triple $\langle u, r, v \rangle$ is defined as:

$$\sum_{\mathbf{r} = r_1, \dots, r_m} \sum_{p \in \mathcal{P}_{\mathbf{r}}(u, v)} \text{LNN-}\wedge \{ \text{LNN-pred}(e_{r_1}), \dots, \text{LNN-pred}(e_{r_m}) \}$$

where $e_r \in \{0, 1\}^{|\mathcal{R}|}$ is a one-hot encoding whose r^{th} entry is 1 with 0s everywhere else. To be clear, the above model consists of $m + 1 + m|\mathcal{R}|$ parameters: $m + 1$ for LNN- \wedge (\mathbf{w}, β) and m LNN-pred operators each consisting of $|\mathcal{R}|$ -sized \mathbf{w} parameters. Due to one-hot encodings and the fact that the inner summation term does not depend on path p , the above expression can be simplified to:

$$\sum_{\mathbf{r} = r_1, \dots, r_m} |\mathcal{P}_{\mathbf{r}}(u, v)| \text{LNN-}\wedge (w_{r_1}^1, \dots, w_{r_m}^m) \quad (4)$$

where w_r^i denotes the r^{th} weight of the i^{th} LNN-pred.

This model is closely related to NeuralLP (Yang, Yang, and Cohen 2017a) that also chains together mixtures of relations (see Equation 5 in Yang, Yang, and Cohen). Differences between the above model and NeuralLP include the use of 1) LNN operators instead of product t -norm and, 2) a rule generator. NeuralLP uses an RNN-based rule generator to generate the relations present in the body of the chain rule, we instead sum over all relation paths. NeuralLP was among the first NeSy-based KBC approaches on top of which later approaches are built, e.g. DRUM (Sadeghian et al. 2019b).

3.3 Mixture of Relation Paths

One of the drawbacks of the previous model is that it treats each LNN-pred independently. To get around this, we propose our second approach which consists of one LNN-pred *across all relation paths*. Given \mathcal{G}, r, m , score for $\langle u, r, v \rangle$ is:

$$\sum_{\mathbf{r} = r_1, \dots, r_m} |\mathcal{P}_{\mathbf{r}}(u, v)| \text{LNN-pred}(e_{\mathbf{r}})$$

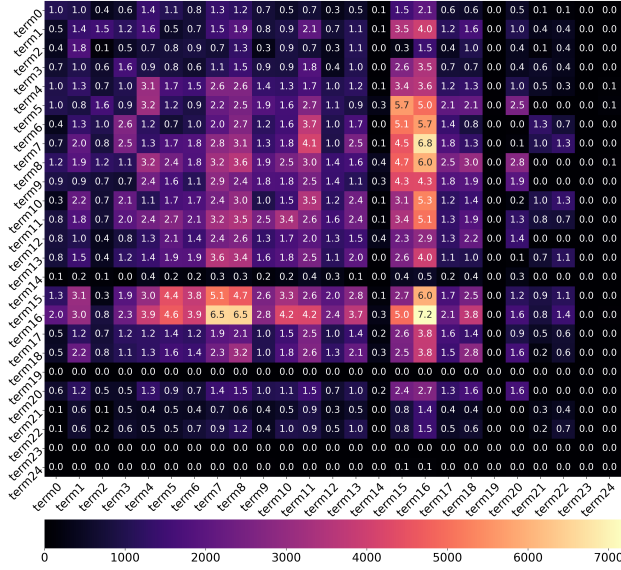


Figure 4: Kinship CP-N3 path scores ($\times 1000$).

where $\mathbf{e}_r \in \{0, 1\}^{|\mathcal{R}|^m}$ is a one-hot encoding with 1 in its \mathbf{r}^{th} position and 0 everywhere else. One way to define a unique index for relation path $\mathbf{r} = r_1, \dots, r_m$ is $\sum_{i=1, \dots, m} r_i |\mathcal{R}|^i$.

While recent approaches (Das et al. 2018; Qu et al. 2021) have followed a similar model, the idea was proposed in the path ranking algorithm (PRA) (Lao, Mitchell, and Cohen 2011). Instead of path counts $|\mathcal{P}_r(u, v)|$, PRA computes the random walk probability of arriving from u to v . Another difference lies in the parameterization. While LNN-pred constrains its parameters to sum to 1, PRA uses elastic net regularization instead. Note that, PRA has been shown to perform quite poorly on KBC benchmarks, e.g. see Qu et al. (2021).

3.4 Handling Sparsity with Graph Embeddings

One of the drawbacks of Mixture of Relation Paths (MP) is that *only paths with relation path \mathbf{r} contribute to the estimation of the weight parameter for \mathbf{r} (w_r)*. In contrast, any path that contains r_i in the i^{th} edge of its relation path (a much larger set of paths) may contribute to the estimation of $w_{r_i}^i$ in Equation 4. This is problematic because in most knowledge graphs, there is a stark difference in path counts for various relation paths. For example, Figure 3 shows the counts for all length 2 relation paths in Kinship and barring 4 ($\text{term16}(U, W) \wedge \text{term7}(W, V)$, $\text{term16}(U, W) \wedge \text{term8}(W, V)$, $\text{term16}(U, W) \wedge \text{term16}(W, V)$, $\text{term7}(U, W) \wedge \text{term16}(W, V)$) all relation path counts lie below 1000. In fact, a vast majority lie in single digits. This implies that for a rare relation path \mathbf{r} , estimated w_r may lack in statistical strength and thus may be unreliable.

To address sparsity of paths, we utilize knowledge graph embeddings (KGE). The literature is rife with approaches that embed \mathcal{V} and \mathcal{R} into a low-dimensional hyperspace by learning distributed representations or *embeddings*. Our assumption is that such techniques score *more prevalent* paths *higher* than paths that are less frequent. Let $\sigma(p)$ denote

the score of path p given such a pre-trained KGE. Score of $\langle u, r, v \rangle$ under the modified MP model is now given by:

$$\sum_{\mathbf{r}=r_1, \dots, r_m} \text{LNN-pred}(\mathbf{e}_r) \sum_{p \in \mathcal{P}_r(u, v)} \sigma(p)$$

Essentially, the goal is to bias the learning process so that relation paths corresponding to larger $\sigma(p)$ are assigned larger weights. Figure 4 shows the (sum of) path scores for all length 2 relation paths in Kinship measured via CP-N3 embeddings (Lacroix, Usunier, and Obozinski 2018) where the cell values are much larger, e.g. while there are only 1200 $\text{term16}(U, W) \wedge \text{term16}(W, V)$ paths its total path score is 7200. Scoring paths using pre-trained KGE was only introduced recently in RNNLogic (Qu et al. 2020) which used RotatE (Sun et al. 2019a) specifically. We next show how to utilize a much larger class of KGEs for the same purpose.

There are at least 2 kinds of KGEs available that rely on a: 1) similarity measure of the triple $\langle h, r, t \rangle$, e.g., CP-N3 (Lacroix, Usunier, and Obozinski 2018), and 2) distance measure used to contrast t 's embedding with some function of h and r 's embeddings, e.g., TransE (Bordes et al. 2013a), RotatE (Sun et al. 2019a). We describe $\sigma(p)$ for both cases:

$$\begin{aligned} \text{similarity-based: } \sigma(p) &= \sum_{\langle h, r, t \rangle \in p} \frac{1}{1 + \exp\{-\text{sim}(h, r, t)\}} \\ \text{distance-based: } \sigma(p) &= \sum_{\langle h, r, t \rangle \in p} \frac{\exp\{2(\delta - d(h, r, t))\} - 1}{\exp\{2(\delta - d(h, r, t))\} + 1} \end{aligned}$$

where δ denotes the margin parameter used by the underlying distance-based KGE to convert distances into similarities. For both of these, we break the path into a series of edges, use the underlying KGE to compute similarity $\text{sim}()$ or distance $d()$ for each triple (as the case may be) and aggregate across all triples in the path. Based on extensive experimentation, we recommend `sigmoid` and `tanh` as the non-linear activation for similarity-based and distance-based KGE, respectively.

3.5 Training Algorithm

Among various possible training algorithms, based on extensive experimentation, we have found the following scheme to perform reliably. In each iteration, we sample uniformly at random a mini-batch of positive triples B^+ from $\langle h, r, t \rangle \in \mathcal{E}$ and negative triples B^- from $\langle h, r, t \rangle \notin \mathcal{E}, \forall h, t \in \mathcal{V}$, such that $|B^+| = |B^-|$ to minimize the following loss:

$$\sum_{\langle h, r, t \rangle \in B^+} \sum_{\langle h', r, t' \rangle \in B^-} \max\{0, \text{score}(h', t') - \text{score}(h, t) + \gamma\}$$

where γ denotes the margin hyperparameter.

4 Experimental Setup

In this section, we describe our experimental setup. We introduce the datasets used and compare against state-of-the-art baselines. We primarily validate our LNN-based rule-learning approaches and their combination with KGE.

Datasets: To evaluate our approach, we experiment on standard KBC benchmarks, viz. Unified Medical Language System (UMLS) (Kok and Domingos 2007), Kinship (Kok and

Domingos 2007), WN18RR (Dettmers et al. 2018), and FB15K-237 (Toutanova and Chen 2015). Table 1 provides dataset statistics. We use standard train/validation/test splits for all datasets. Note that, RNNLogic defined its own splits for Kinship and UMLS, we report results on these too for a fair comparison (numbers in parenthesis in Table 1).

Baselines: We compare our approaches against state-of-the-art KBC approaches categorized as follows:

- *Rule-learning approaches:* Neural-LP (Yang, Yang, and Cohen 2017b), DRUM (Sadeghian et al. 2019b), RNNLogic (rules only) (Qu et al. 2021) and CTP (Minervini et al. 2020).
- *Rule-learning approaches that utilize KGE:* We compare against RNNLogic (w/ embd.) (Qu et al. 2021) which uses RotatE (Sun et al. 2019b), MINERVA (Das et al. 2018) and MultiHopKG (Lin, Socher, and Xiong 2018).
- *Embedding-based approaches:* We report results from ComplEx-N3, the best KGE we know of, and CP-N3, the KGE used by our methods, (both from Lacroix, Usunier, and Obozinski (2018)), RotatE (Sun et al. 2019b), used in RNNLogic, and Complex (Trouillon et al. 2016).

Metrics: Given an unseen query $\langle h, r, ? \rangle$ we compute *filtered* ranks (Bordes et al. 2013a) for destination vertices after removing the destinations that form edges with h and r present in the train, test and validation sets. Based on Sun et al. (2020)’s suggestions, our definitions of mean reciprocal rank (MRR) and Hits@K satisfy two properties: 1) They assign a larger value to destinations ranked lower, and 2) If destinations t_1, \dots, t_m share the same rank $n + 1$ then each of them is assigned an average of ranks $n + 1, \dots, n + m$:

$$\text{MRR}(t_i) = \frac{1}{m} \sum_{r=n+1}^{n+m} \frac{1}{r}, \text{ Hits@K}(t_i) = \frac{1}{m} \sum_{r=n+1}^{n+m} \delta(r \leq K)$$

where $\delta(\cdot)$ denotes the Dirac delta function. We include inverse triples and report averages across the test set.

Implementation: We evaluate rules learned with *Chain of Mixtures* (LNN-CM), *Mixture of Paths* (LNN-MP) and their combinations with KGE. We use Adagrad (Duchi, Hazan, and Singer 2011) with step size $\in \{0.1, 1.0\}$, margin $\gamma \in \{0.1, 0.5, 1.0, 2.0\}$ and batch size $|B^+| = |B^-| = 8$. We use the validation set to perform hyperparameter tuning and learn rules of length up to 4 for FB15K-237, 5 for WN18RR, and 3 for Kinship, UMLS. We combine our rule-learning approach with pre-trained CP-N3 embeddings of dimension $4K$ for FB15K-237, $3K$ for WN18RR, and $8K$ for Kinship, UMLS.

5 Results and Discussion

Table 2 shows all our results comparing LNN-CM and LNN-MP on standard splits of all datasets. We compare against all the baselines including KGE-based approaches (EMBD.), approaches that only learn rules (RULES), and approaches that learn rules with KGE (w/ EMBD.). For the last category, we indicate the KGE used in the name of the approach.

LNN-CM vs. NeuralLP, DRUM: Recall that LNN-CM is closely related to NeuralLP and DRUM. Table 2 shows that, in terms of MRR, LNN-CM outperforms DRUM and NeuralLP on UMLS and FB15K-237 while being comparable to DRUM on Kinship and WN18RR. LNN-CM achieves this

Dataset	Train	Valid	Test	Relations	Entities
Kinship	8544 (3206)	1068 (2137)	1074 (5343)	50	104
UMLS	5216 (1959)	652 (1306)	661 (3264)	92	135
WN18RR	86835	3034	3134	11	40943
FB15K-237	272155	17535	20466	237	14541

Table 1: Dataset Statistics. Numbers in parenthesis reflects splits from RNNLogic (Qu et al. 2021).

without using an RNN-based rule generator and is thus conceptually much simpler. This is clear evidence that an RNN is not necessary for KBC. DRUM and NeuralLP’s results in Table 2 are significantly lower than what was reported in Yang, Yang, and Cohen (2017a); Sadeghian et al. (2019a). As explained in Sun et al. (2020), this is likely due to the carefully defined MRR and Hits@K used in our evaluation, not using which can lead to overly-optimistic and unfair results¹. **LNN-MP vs. LNN-CM:** Table 2 shows that LNN-MP (RULES) consistently outperforms LNN-CP. While previous comparisons against NeuralLP and/or DRUM (Qu et al. 2021; Lin, Socher, and Xiong 2018; Das et al. 2018) also hinted at this result, a surprising finding is that the margin of difference depends on the dataset. In particular, on FB15K-237 LNN-CM’s MRR comes within 1% of state-of-the-art RNNLogic (RULES)’s. Across all standard splits, LNN-MP (RULES) outperforms all rule-learning methods. On the smaller Kinship and UMLS datasets, LNN-MP outperforms CTP, and on the larger WN18RR and FB15K-237 it outperforms RNNLogic. In Table 3, we also report LNN-MP’s results on RNNLogic’s splits which include a much smaller training set for Kinship and UMLS (as shown in Table 1). Given the substantial increases resulting from the switch to standard splits (+14.5% and +8.5% in MRR on Kinship and UMLS, respectively), it seems that the simpler LNN-MP (RULES) exploits the additional training data much more effectively. On the other hand, RNNLogic (RULES) hardly shows any improvement on Kinship and in fact, deteriorates on UMLS. A possible reason could be that the inexact training algorithm employed by RNNLogic (based on expectation-maximization and ELBO bound) fails to leverage the additional training data.

Learning Rules with Embeddings: Since LNN-MP outperforms LNN-CM (rules only), we combine it with CP-N3, one of the best KGE available, to learn rules with embeddings. In comparison to LNN-MP (RULES), LNN-MP w/ CP-N3 shows consistent improvements ranging from +1.2% MRR (on WN18RR, Table 2) to +9.2% MRR (on Kinship standard splits, Table 2). In comparison to RNNLogic w/ RotatE, LNN-

¹Investigating more, Table 4 in Sadeghian et al. reports DRUM’s Hits@10 on WN18RR as 0.568 with rule length up to 2. However, 56% of WN18RR’s validation set triples require paths longer than 2 to connect the source with the destination (test set fraction is similar). This implies that one can only hope to achieve a maximum Hits@10 of 0.44 when restricted to rules of length 2 learned via *any rule-based KBC* technique which is < 0.568 and a clear contradiction.

		Kinship			UMLS			WN18RR			FB15K-237		
		MRR	Hits@10	Hits@3	MRR	Hits@10	Hits@3	MRR	Hits@10	Hits@3	MRR	Hits@10	Hits@3
EMBD.	CP-N3	88.9	98.7	94.8	93.3	99.7	98.9	47.0*	54.0*	—	36.0*	54.0*	—
	Complex-N3	89.2	98.6	94.7	95.9	99.7	98.9	48.0*	57.0*	—	37.0*	56.0*	—
	RotatE	75.5	97.3	86.5	86.1	99.5	97.0	47.6*	57.1*	49.2*	33.8*	53.3*	37.5*
	Complex	83.7	98.4	91.6	94.5	99.7	98.0	44.0	49.6	44.9	31.8	49.1	34.7
RULES	NeuralLP	48.8	89.1	63.0	55.3	93.0	75.4	33.7	50.2	43.3	25.8	48.5	31.4
	DRUM	40.0	86.1	48.2	61.8	97.9	91.2	34.8	52.1	44.6	25.8	49.1	31.5
	CTP	70.3	93.9	79.7	80.1	97.0	91.0	—	—	—	—	—	—
	RNNLogic	64.5	91.1	72.9	71.0	91.1	82.1	45.5*	53.1*	47.5*	28.8*	44.5*	31.5*
	LNN-CM (Ours)	39.1	68.7	43.6	78.7	95.1	88.9	36.6	49.2	39.2	28.0	43.5	30.4
	LNN-MP (Ours)	81.9	98.4	89.3	90.0	99.4	98.3	47.3	55.5	49.7	30.7	47.0	34.2
w/ EMBD.	RNNLogic w/ RotatE	—	—	—	—	—	—	48.3*	55.8*	49.7*	34.4*	53.0*	38.0*
	LNN-MP (Ours) w/ CP-N3	91.1	99.2	96.5	94.5	100	99.2	48.5	56.1	50.2	35.1	53.0	39.1

Table 2: Results of LNN-CM and LNN-MP in comparison to other state-of-the-art approaches on standard splits of Kinship, UMLS, WN18RR, and FB15K-237. Bold font denotes best within each category of KBC approaches. * denotes results copied from original papers. CTP did not scale to larger datasets. ComplEx-N3 does not report Hits@3 for WN18RR and FB15K-237. RNNLogic does not report Kinship and UMLS results on standard splits w/ RotatE.

	Kinship			UMLS		
	MRR	Hits@10	Hits@3	MRR	Hits@10	Hits@3
CP-N3	60.2	92.2	70.0	76.6	95.0	85.6
ComplEx-N3	60.5*	92.1*	71.0*	79.1*	95.7*	87.3*
RotatE	65.1*	93.2*	75.5*	74.4*	93.9*	82.2*
Complex	54.4	89.0	63.8	74.0	92.5	81.1
RNNLogic	63.9	92.4	73.1	74.5	92.4	83.3
LNN-MP (Ours)	67.4	95.0	77.0	81.5	96.6	91.8
RNNLogic w/ RotatE	72.2*	94.9*	81.4*	84.2*	96.5*	89.1*
LNN-MP w/ CP-N3 (Ours)	72.2	96.9	81.6	87.6	98.0	94.7

Table 3: Results on RNNLogic’s splits for Kinship and UMLS. * denotes results copied from Qu et al. (2021).

MP w/ CP-N3 shows small but consistent improvements on UMLS (RNNLogic’s splits, Table 3), WN18RR, FB15K-237 (Table 2), and comparable results on Kinship (RNNLogic’s splits, Table 3). This is encouraging due to the simplicity of LNN-MP which adds to its appeal. Table 4 shows that LNN-MP w/ CP-N3 also outperforms MINERVA and MultiHopKG, two more MP-based approaches that utilize embeddings. Note that, following MINERVA and MultiHopKG, for this comparison we leave out inverse triples from the test set.

5.1 Qualitative Analysis and Discussion

Learned rules can be extracted from LNN-MP by sorting the relation paths in descending order of w_r . Table 5 presents some of the learned rules appearing in the top-10.

Rules for FB15K-237: Rule 1 in Table 5 describes a learned rule that infers the language a person speaks by exploiting knowledge of the language spoken in her/his country of nationality. In terms of relation paths, this looks like:

$$P(\text{person}) \xrightarrow{\text{nationality}} N(\text{nation}) \xrightarrow{\text{spoken_in}} L(\text{language})$$

	WN18RR			FB15K-237		
	MRR	Hits@10	Hits@3	MRR	Hits@10	Hits@3
MINERVA*	44.8	51.3	45.6	29.3	45.6	32.9
MultiHopKG*	47.2	54.2	—	40.7	56.4	—
RNNLogic*	47.8	55.3	50.3	37.7	54.9	41.2
RNNLogic* w/ RotatE	50.6	59.2	52.3	44.3	64.0	48.9
LNN-MP (Ours)	51.6	58.4	53.5	40.0	57.4	44.4
LNN-MP w/ CP-N3 (Ours)	51.9	59.0	53.9	44.8	63.6	49.5

Table 4: Evaluating on direct triples (excluding inverses).

Similarly, Rule 2 uses the film_country relation instead of nationality to infer the language used in a film. Besides spoken_in, FB15K-237 contains other relations that can be utilized to infer language such as the official_language spoken in a country. Rule 3 uses this relation to infer the language spoken in a TV program by first exploiting knowledge of its country of origin. Rules 5, 6 and 7 are longer rules containing 3 relations each in their body. Rule 5 infers a TV program’s country by first exploiting knowledge of one of its actor’s birth place and then determining which country the birth place belongs to. In terms of relation path: $P(\text{program}) \xrightarrow{\text{tv_program_actor}} A(\text{actor}) \xrightarrow{\text{born_in}} L(\text{birth place}) \xrightarrow{\text{located_in}} N(\text{nation})$. Rule 6 uses a film crew member’s marriage location instead to infer the region where the film was released. Rule 7 infers the marriage location of a celebrity by exploiting knowledge of where their friend got married.

Recursive Rules for WN18RR: Unlike FB15K-237, WN18RR is a hierarchical KG and thus sparser which calls for longer rules to be learned. A majority of WN18RR’s edges are concentrated in a few relations, e.g. hypernym, and most learned rules rely on such relations. Hypernym(X, Y) is true if Y , e.g. dog, is a kind of X , e.g. animal. Table 5 presents 3

	1) $\text{person_language}(P, L) \leftarrow \text{nationality}(P, N) \wedge \text{spoken_in}(L, N)$
	2) $\text{film_language}(F, L) \leftarrow \text{film_country}(F, C) \wedge \text{spoken_in}(L, C)$
	3) $\text{tv_program_language}(P, L) \leftarrow \text{country_of_tv_program}(P, N) \wedge \text{official_language}(N, L)$
FB15K-237	4) $\text{burial_place}(P, L) \leftarrow \text{nationality}(P, N) \wedge \text{located_in}(L, N)$
	5) $\text{country_of_tv_program}(P, N) \leftarrow \text{tv_program_actor}(P, A) \wedge \text{born_in}(A, L) \wedge \text{located_in}(L, N)$
	6) $\text{film_release_region}(F, R) \leftarrow \text{film_crew}(F, P) \wedge \text{marriage_location}(P, L) \wedge \text{located_in}(L, R)$
	7) $\text{marriage_location}(P, L) \leftarrow \text{celebrity_friends}(P, F) \wedge \text{marriage_location}(F, L') \wedge \text{location_adjoins}(L', L)$
	8) $\text{domain_topic}(X, Y) \leftarrow \text{hypernym}(X, U) \wedge \text{hypernym}(U, V) \wedge \text{hypernym}(W, V) \wedge \text{hypernym}(Z, W) \wedge \text{domain_topic}(Z, Y)$
WN18RR	9) $\text{derivation}(X, Y) \leftarrow \text{hypernym}(X, U) \wedge \text{hypernym}(V, U) \wedge \text{derivation}(W, V) \wedge \text{hypernym}(W, Z) \wedge \text{hypernym}(Y, Z)$
	10) $\text{hypernym}(X, Y) \leftarrow \text{member_meronym}(U, X) \wedge \text{hypernym}(U, V) \wedge \text{hypernym}(W, V) \wedge \text{member_meronym}(W, Z) \wedge \text{hypernym}(Z, Y)$

Table 5: Examples of LNN-MP’s learned rules.

learned rules for WN18RR. $\text{Domain_topic}(X, Y)$ is true if X is the category, e.g., computer science, of scientific concept Y , e.g., CPU. Rule 8 infers the domain topic of Y by first climbing down the hypernym is-a hierarchy 2 levels and then climbing up: $X \xrightarrow{\text{is-a}} U \xrightarrow{\text{is-a}} V \xleftarrow{\text{is-a}} W \xleftarrow{\text{is-a}} Z \xrightarrow{\text{domain}} Y$. Rules 9 and 10 are recursive, where the relation in the head of the rule also appears in the body. Here, derivation indicates related form of concept, e.g., yearly is the derivation of year, and member_meronym denotes a concept which is a member of another concept such as player and team.

6 Related Work

Knowledge Base Completion (KB) approaches have gained a lot of interest recently, due to their ability to handle the incompleteness of knowledge bases. Embedding-based techniques maps entities and relations to low-dimensional vector space to infer new facts. These techniques use neighborhood structure of an entity or relation to learn their corresponding embedding. Starting off with translational embeddings (Bordes et al. 2013b), embedding-based techniques have evolved to use complex vector spaces such as ComplEx (Trouillon et al. 2016), RotatE (Sun et al. 2019a), and QuatE (Zhang et al. 2019). While the performance of embedding-based techniques has improved over time, rule-learning techniques have gained attention due to its inherent ability for learning interpretable rules (Yang, Yang, and Cohen 2017b; Sadeghian et al. 2019b; Rocktäschel and Riedel 2017; Qu et al. 2020).

The core ideas in rule learning can be categorized into two groups based on their mechanism to select relations for rules. While *Chain of Mixtures (CM)* represents each relation in the body as a mixture, e.g. NeuralLP (Yang, Yang, and Cohen 2017a), DRUM (Sadeghian et al. 2019a), *Mixture of Paths (MP)* learns relation paths; e.g., MINERVA (Das et al. 2018), RNNLogic (Qu et al. 2021). Recent trends in both these types of rule learning approaches has shown significant increase in complexity for performance gains over their simpler precursors. Among the first to learn rules for KBC, NeuralLP (Yang, Yang, and Cohen 2017a) uses a long short-term memory (LSTM) (Hochreiter and Schmidhuber 1997) as its rule generator. DRUM (Sadeghian et al. 2019a) improves upon NeuralLP by learning multiple such rules obtained by running a bi-directional LSTM for more steps. MINERVA (Das et al. 2018) steps away from representing each body relation as a mixture, proposing instead to learn

the relation sequences appearing in paths connecting source to destination vertices using neural reinforcement learning.

RNNLogic (Qu et al. 2021) is the latest to adopt a path-based approach for KBC that consists of two modules, a rule-generator for suggesting high quality paths and a reasoning predictor that uses said paths to predict missing information. RNNLogic employs expectation-maximization for training where the E-step identifies useful paths per data instance (edge in the KG) by sampling from an intractable posterior while the M-step uses the per-instance useful paths to update the overall set of paths. Both DRUM and RNNLogic represent a significant increase in complexity of their respective approaches compared to NeuralLP and MINERVA.

Unlike these approaches, we propose to utilize *Logical Neural Networks (LNN)* (Riegel et al. 2020); a simple yet powerful neuro-symbolic approach which extends Boolean logic to the real-valued domain. On top of LNN, we propose two simple approaches for rule learning based on CP and MP. Furthermore, our approach allows for combining rule-based KBC with any KGE, which results in state-of-the-art performance across multiple datasets.

7 Conclusion

In this paper, we proposed an approach for rule-based KBC based on Logic Neural Networks (LNN); a neuro-symbolic differentiable framework for real-valued logic. In particular, we present two approaches for rule learning, one that represent rules using chains of predicate mixtures (LNN-CM), and another that uses mixtures of paths (LNN-MP). We show that both approaches can be implemented using LNN and neuro-symbolic AI, and result in better or comparable performance to state-of-the-art. Furthermore, our framework facilitates combining rule learning with knowledge graph embedding techniques to harness the best of both worlds.

Our experimental results across four benchmarks show that such a combination provides better results and establishes new state-of-the-art performance across multiple datasets. We also showed how easy it is to extract and interpret learned rules. With both LNN-CM and LNN-MP implemented within the same LNN framework, one avenue of future work would be to explore combinations of the two approaches given their good performance on a number of KBC benchmarks.

References

- Bollacker, K.; Evans, C.; Paritosh, P.; Sturge, T.; and Taylor, J. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, 1247–1250.
- Bordes, A.; Usunier, N.; Garcia-Duran, A.; Weston, J.; and Yakhnenko, O. 2013a. Translating embeddings for modeling multi-relational data. In *NeurIPS*.
- Bordes, A.; Usunier, N.; Garcia-Duran, A.; Weston, J.; and Yakhnenko, O. 2013b. Translating embeddings for modeling multi-relational data. *Advances in neural information processing systems*, 26.
- Das, R.; Dhuliawala, S.; Zaheer, M.; Vilnis, L.; Durugkar, I.; Krishnamurthy, A.; Smola, A.; and McCallum, A. 2018. Go for a Walk and Arrive at the Answer: Reasoning Over Paths in Knowledge Bases using Reinforcement Learning. In *ICLR*.
- Dettmers, T.; Minervini, P.; Stenetorp, P.; and Riedel, S. 2018. Convolutional 2d knowledge graph embeddings. In *Thirty-second AAAI conference on artificial intelligence*.
- Dong, H.; Mao, J.; Lin, T.; Wang, C.; Li, L.; and Zhou, D. 2019. Neural Logic Machines. In *ICLR*.
- Duchi, J.; Hazan, E.; and Singer, Y. 2011. Adaptive Subgradient Methods for Online Learning and Stochastic Optimization. *JMLR*.
- Frerix, T.; Nießner, M.; and Cremers, D. 2020. Homogeneous Linear Inequality Constraints for Neural Network Activations. In *CVPR Workshops*.
- Hochreiter, S.; and Schmidhuber, J. 1997. Long short-term memory. *Neural computation*.
- Kok, S.; and Domingos, P. 2007. Statistical predicate invention. In *Proceedings of the 24th international conference on Machine learning*, 433–440.
- Lacroix, T.; Usunier, N.; and Obozinski, G. 2018. Canonical tensor decomposition for knowledge base completion. In *International Conference on Machine Learning*, 2863–2872. PMLR.
- Lao, N.; Mitchell, T.; and Cohen, W. W. 2011. Random Walk Inference and Learning in A Large Scale Knowledge Base. In *EMNLP*.
- Lin, X. V.; Socher, R.; and Xiong, C. 2018. Multi-hop knowledge graph reasoning with reward shaping. In *EMNLP*.
- Minervini, P.; Riedel, S.; Stenetorp, P.; Grefenstette, E.; and Rocktäschel, T. 2020. Learning reasoning strategies in end-to-end differentiable proving. In *International Conference on Machine Learning*, 6938–6949. PMLR.
- Pascanu, R.; Mikolov, T.; and Bengio, Y. 2013. On the Difficulty of Training Recurrent Neural Networks. In *ICML*.
- Qu, M.; Chen, J.; Xhonneux, L.-P.; Bengio, Y.; and Tang, J. 2020. Rnnlogic: Learning logic rules for reasoning on knowledge graphs. *arXiv preprint arXiv:2010.04029*.
- Qu, M.; Chen, J.; Xhonneux, L.-P.; Bengio, Y.; and Tang, J. 2021. {RNNL}ogic: Learning Logic Rules for Reasoning on Knowledge Graphs. In *ICLR*.
- Rebele, T.; Suchanek, F.; Hoffart, J.; Biega, J.; Kuzey, E.; and Weikum, G. 2016. YAGO: A multilingual knowledge base from wikipedia, wordnet, and geonames. In *International semantic web conference*, 177–185. Springer.
- Riegel, R.; Gray, A.; Luus, F.; Khan, N.; Makondo, N.; Akhalwaya, I. Y.; Qian, H.; Fagin, R.; Barahona, F.; Sharma, U.; Ikbal, S.; Karanam, H.; Neelam, S.; Likhyan, A.; and Srivastava, S. 2020. Logical Neural Networks. *CoRR*.
- Rocktäschel, T.; and Riedel, S. 2017. End-to-end differentiable proving. *arXiv preprint arXiv:1705.11040*.
- Sadeghian, A.; Armandpour, M.; Ding, P.; and Wang, D. Z. 2019a. DRUM: End-To-End Differentiable Rule Mining On Knowledge Graphs. In *NeurIPS*.
- Sadeghian, A.; Armandpour, M.; Ding, P.; and Wang, D. Z. 2019b. Drum: End-to-end differentiable rule mining on knowledge graphs. *arXiv preprint arXiv:1911.00055*.
- Sun, Z.; Deng, Z.-H.; Nie, J.-Y.; and Tang, J. 2019a. Rotate: Knowledge graph embedding by relational rotation in complex space. *arXiv preprint arXiv:1902.10197*.
- Sun, Z.; Deng, Z.-H.; Nie, J.-Y.; and Tang, J. 2019b. RotateE: Knowledge Graph Embedding by Relational Rotation in Complex Space. In *ICLR*.
- Sun, Z.; Vashishth, S.; Sanyal, S.; Talukdar, P.; and Yang, Y. 2020. A Re-evaluation of Knowledge Graph Completion Methods. In *ACL*.
- Toutanova, K.; and Chen, D. 2015. Observed versus latent features for knowledge base and text inference. In *Proceedings of the 3rd workshop on continuous vector space models and their compositionality*, 57–66.
- Trinh, T. H.; Dai, A. M.; Luong, M.-T.; and Le, Q. V. 2018. Learning Longer-term Dependencies in RNNs with Auxiliary Losses. In *ICLR Workshops*.
- Trouillon, T.; Welbl, J.; Riedel, S.; Gaussier, É.; and Bouchard, G. 2016. Complex embeddings for simple link prediction. In *International conference on machine learning*, 2071–2080. PMLR.
- Yang, F.; Yang, Z.; and Cohen, W. W. 2017a. Differentiable Learning of Logical Rules for Knowledge Base Reasoning. In *NeurIPS*.
- Yang, F.; Yang, Z.; and Cohen, W. W. 2017b. Differentiable learning of logical rules for knowledge base reasoning. *arXiv preprint arXiv:1702.08367*.
- Zhang, S.; Tay, Y.; Yao, L.; and Liu, Q. 2019. Quaternion knowledge graph embeddings. *arXiv preprint arXiv:1904.10281*.