

Neural Bellman-Ford Networks: A General Graph Neural Network Framework for Link Prediction

Zhaocheng Zhu^{1,2}, Zuobai Zhang^{1,2}, Louis-Pascal Xhonneux^{1,2}, Jian Tang^{1,3,4}

Mila - Québec AI Institute¹, Université de Montréal²

HEC Montréal³, CIFAR AI Chair⁴

{zhaocheng.zhu, zuobai.zhang, louis-pascal.xhonneux}@mila.quebec
jian.tang@hec.ca

Abstract

Link prediction is a very fundamental task on graphs. Inspired by traditional path-based methods, in this paper we propose a general and flexible representation learning framework based on paths for link prediction. Specifically, we define the representation of a pair of nodes as the *generalized sum* of all path representations between the nodes, with each path representation as the *generalized product* of the edge representations in the path. Motivated by the Bellman-Ford algorithm for solving the shortest path problem, we show that the proposed path formulation can be efficiently solved by the generalized Bellman-Ford algorithm. To further improve the capacity of the path formulation, we propose the Neural Bellman-Ford Network (NBFNet), a general graph neural network framework that solves the path formulation with learned operators in the generalized Bellman-Ford algorithm. The NBFNet parameterizes the generalized Bellman-Ford algorithm with 3 neural components, namely INDICATOR, MESSAGE and AGGREGATE functions, which corresponds to the boundary condition, *multiplication* operator, and *summation* operator respectively¹. The NBFNet covers many traditional path-based methods, and can be applied to both homogeneous graphs and multi-relational graphs (e.g., knowledge graphs) in both transductive and inductive settings. Experiments on both homogeneous graphs and knowledge graphs show that the proposed NBFNet outperforms existing methods by a large margin in both transductive and inductive settings, achieving new state-of-the-art results².

1 Introduction

Predicting the interactions between nodes (a.k.a. link prediction) is a fundamental task in the field of graph machine learning. Given the ubiquitous existence of graphs, such a task has many applications, such as recommender system [34], knowledge graph completion [41] and drug repurposing [27].

Traditional methods of link prediction usually define different heuristic metrics over the paths between a pair of nodes. For example, Katz index [30] is defined as a weighted count of paths between two nodes. Personalized PageRank [42] measures the similarity of two nodes as the random walk probability from one to the other. Graph distance [37] uses the length of the shortest path between two nodes to predict their association. These methods can be directly applied to new graphs, i.e., inductive setting, enjoy good interpretability and scale up to large graphs. However, they are designed based on handcrafted metrics and may not be optimal for link prediction on real-world graphs.

¹Unless stated otherwise, we use *summation* and *multiplication* to refer the generalized operators in the path formulation, rather than the basic operations of arithmetic.

²Code is available at <https://github.com/DeepGraphLearning/NBFNet>

To address these limitations, some link prediction methods adopt graph neural networks (GNNs) [32, 48, 59] to automatically extract important features from local neighborhoods for link prediction. Thanks to the high expressiveness of GNNs, these methods have shown state-of-the-art performance. However, these methods can only be applied to predict new links on the training graph, i.e. transductive setting, and lack interpretability. While some recent methods [73, 55] extract features from local subgraphs with GNNs and support inductive setting, the scalability of these methods is compromised.

Therefore, we wonder if there exists an approach that enjoys the advantages of both traditional path-based methods and recent approaches based on graph neural networks, i.e., **generalization in the inductive setting, interpretability, high model capacity and scalability**.

In this paper, we propose such a solution. Inspired by traditional path-based methods, our goal is to develop a general and flexible representation learning framework for link prediction based on the paths between two nodes. Specifically, we define the representation of a pair of nodes as the *generalized sum* of all the path representations between them, where each path representation is defined as the *generalized product* of the edge representations in the path. Many link prediction methods, such as Katz index [30], personalized PageRank [42], graph distance [37], as well as graph theory algorithms like widest path [4] and most reliable path [4], are special instances of this path formulation with different *summation* and *multiplication* operators. Motivated by the polynomial-time algorithm for the shortest path problem [5], we show that such a formulation can be efficiently solved via the generalized Bellman-Ford algorithm [4] under mild conditions and scale up to large graphs.

The operators in the generalized Bellman-Ford algorithm—*summation* and *multiplication*—are handcrafted, which have limited flexibility. Therefore, we further propose the Neural Bellman-Ford Networks (NBFNet), a graph neural network framework that solves the above path formulation with learned operators in the generalized Bellman-Ford algorithm. Specifically, NBFNet parameterizes the generalized Bellman-Ford algorithm with three neural components, namely **INDICATOR, MESSAGE and AGGREGATE functions**. The INDICATOR function initializes a representation on each node, which is taken as the boundary condition of the generalized Bellman-Ford algorithm. The MESSAGE and the AGGREGATE functions learn the *multiplication* and *summation* operators respectively.

We show that the MESSAGE function can be defined according to the relational operators in knowledge graph embeddings [6, 68, 58, 31, 52], e.g., as a translation in Euclidean space induced by the relational operators of TransE [6]. The AGGREGATE function can be defined as learnable set aggregation functions [71, 65, 9]. With such parameterization, NBFNet can generalize to the inductive setting, meanwhile achieve one of the lowest time complexity among inductive GNN methods. A comparison of NBFNet and other GNN frameworks for link prediction is showed in Table 1. With other instantiations of MESSAGE and AGGREGATE functions, our framework can also recover some existing works on learning logic rules [69, 46] for link prediction on knowledge graphs (Table 2).

Our NBFNet framework can be applied to several link prediction variants, covering not only single-relational graphs (e.g., homogeneous graphs) but also multi-relational graphs (e.g., knowledge graphs). We empirically evaluate the proposed NBFNet for link prediction on homogeneous graphs and knowledge graphs in both transductive and inductive settings. Experimental results show that the proposed NBFNet outperforms existing state-of-the-art methods by a large margin in all settings, with an average relative performance gain of 18% on knowledge graph completion (HITS@1) and 22% on inductive relation prediction (HITS@10). We also show that the proposed NBFNet is indeed interpretable by visualizing the top-k relevant paths for link prediction on knowledge graphs.

Table 1: Comparison of GNN frameworks for link prediction. The time complexity refers to the *amortized time* for predicting a single edge or triplet. $|\mathcal{V}|$ is the number of nodes, $|\mathcal{E}|$ is the number of edges, and d is the dimension of representations. The wall time is measured on FB15k-237 test set with 40 CPU cores and 4 GPUs. We estimate the wall time of GraIL based on a downsampled test set.

Method	Inductive ³	Interpretable	Learned Representation	Time Complexity	Wall Time
VGAE [32] / RGCN [48]			✓	$O(d)$	18 secs
NeuralLP [69] / DRUM [46]	✓	✓		$O\left(\frac{ \mathcal{E} d}{ \mathcal{V} } + d^2\right)$	2.1 mins
SEAL [73] / GraIL [55]	✓		✓	$O(\mathcal{E} d^2)$	≈1 month
NBFNet	✓	✓	✓	$O\left(\frac{ \mathcal{E} d}{ \mathcal{V} } + d^2\right)$	4.0 mins

2 Related Work

Existing work on link prediction can be generally classified into 3 main paradigms: path-based methods, embedding methods, and graph neural networks.

Path-based Methods. Early methods on homogeneous graphs compute the similarity between two nodes based on the weighted count of paths (Katz index [30]), random walk probability (personalized PageRank [42]) or the length of the shortest path (graph distance [37]). SimRank [28] uses advanced metrics such as the expected meeting distance on homogeneous graphs, which is extended by PathSim [51] to heterogeneous graphs. On knowledge graphs, Path Ranking [35, 15] directly uses relational paths as symbolic features for prediction. Rule mining methods, such as NeuralLP [69] and DRUM [46], learn probabilistic logical rules to weight different paths. Path representation methods, such as Path-RNN [40] and its successors [11, 62], encode each path with recurrent neural networks (RNNs), and aggregate paths for prediction. However, these methods need to traverse an exponential number of paths and are limited to very short paths, e.g., ≤ 3 edges. To scale up path-based methods, All-Paths [57] proposes to efficiently aggregate all paths with dynamic programming. However, All-Paths is restricted to bilinear models and has limited model capacity. Another stream of works [64, 10, 22] learns an agent to collect useful paths for link prediction. While these methods can produce interpretable paths, they suffer from extremely sparse rewards and require careful engineering of the reward function [38] or the search strategy [50]. Some other works [8, 44] adopt variational inference to learn a path finder and a path reasoner for link prediction.

Embedding Methods. Embedding methods learn a distributed representation for each node and edge by preserving the edge structure of the graph. Representative methods include DeepWalk [43] and LINE [53] on homogeneous graphs, and TransE [6], DistMult [68] and RotatE [52] on knowledge graphs. Later works improve embedding methods with new score functions [58, 13, 31, 52, 54, 76] that capture common semantic patterns of the relations, or search the score function in a general design space [75]. Embedding methods achieve promising results on link prediction, and can be scaled to very large graphs using multiple GPUs [78]. However, embedding methods do not explicitly encode local subgraphs between node pairs and cannot be applied to the inductive setting.

Graph Neural Networks. Graph neural networks (GNNs) [47, 33, 60, 65] are a family of representation learning models that encode topological structures of graphs. For link prediction, the prevalent frameworks [32, 48, 12, 59] adopt an auto-encoder formulation, which uses GNNs to encode node representations, and decodes edges as a function over node pairs. Such frameworks are potentially inductive if the dataset provides node features, but are transductive only when node features are unavailable. Another stream of frameworks, such as SEAL [73] and GraIL [55], explicitly encodes the subgraph around each node pair for link prediction. While these frameworks are proved to be more powerful than the auto-encoder formulation [74] and can solve the inductive setting, they require to materialize a subgraph for each link, which is not scalable to large graphs. By contrast, our NBFNet explicitly captures the paths between two nodes for link prediction, meanwhile achieves a relatively low time complexity (Table 1). ID-GNN [70] formalizes link prediction as a conditional node classification task, and augments GNNs with the identity of the source node. While the architecture of NBFNet shares some spirits with ID-GNN, our model is motivated by the generalized Bellman-Ford algorithm and has theoretical connections with traditional path-based methods. There are also some works trying to scale up GNNs for link prediction by dynamically pruning the set of nodes in message passing [66, 20]. These methods are complementary to NBFNet, and may be incorporated into our method to further improve scalability.

3 Methodology

In this section, we first define a path formulation for link prediction. Our path formulation generalizes several traditional methods, and can be efficiently solved by the generalized Bellman-Ford algorithm. Then we propose Neural Bellman-Ford Networks to learn the path formulation with neural functions.

3.1 Path Formulation for Link Prediction

We consider the link prediction problem on both knowledge graphs and homogeneous graphs. A

³We consider the inductive setting where a model can generalize to entirely new graphs without node features.

knowledge graph is denoted by $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{R})$, where \mathcal{V} and \mathcal{E} represent the set of entities (nodes) and relations (edges) respectively, and \mathcal{R} is the set of relation types. We use $\mathcal{N}(u)$ to denote the set of nodes connected to u , and $\mathcal{E}(u)$ to denote the set of edges ending with node u . A homogeneous graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ can be viewed as a special case of knowledge graphs, with only one relation type for all edges. Throughout this paper, we use **bold** terms, $\mathbf{w}_q(e)$ or $\mathbf{h}_q(u, v)$, to denote vector representations, and *italic* terms, w_e or w_{uv} , to denote scalars like the weight of edge (u, v) in homogeneous graphs or triplet (u, r, v) in knowledge graphs. Without loss of generality, we derive our method based on knowledge graphs, while our method can also be applied to homogeneous graphs.

Path Formulation. Link prediction is aimed at predicting the existence of a query relation q between a head entity u and a tail entity v . From a representation learning perspective, this requires to learn a pair representation $\mathbf{h}_q(u, v)$, which captures the local subgraph structure between u and v w.r.t. the query relation q . In traditional methods, such a local structure is encoded by counting different types of random walks from u to v [35, 15]. Inspired by this construction, we formulate the pair representation as a *generalized sum* of path representations between u and v with a commutative *summation* operator \oplus . Each path representation $\mathbf{h}_q(P)$ is defined as a *generalized product* of the edge representations in the path with the *multiplication* operator \otimes .

$$\mathbf{h}_q(u, v) = \mathbf{h}_q(P_1) \oplus \mathbf{h}_q(P_2) \oplus \dots \oplus \mathbf{h}_q(P_{|\mathcal{P}_{uv}|})|_{P_i \in \mathcal{P}_{uv}} \triangleq \bigoplus_{P \in \mathcal{P}_{uv}} \mathbf{h}_q(P) \quad (1)$$

$$\mathbf{h}_q(P = (e_1, e_2, \dots, e_{|P|})) = \mathbf{w}_q(e_1) \otimes \mathbf{w}_q(e_2) \otimes \dots \otimes \mathbf{w}_q(e_{|P|}) \triangleq \bigotimes_{i=1}^{|P|} \mathbf{w}_q(e_i) \quad (2)$$

where \mathcal{P}_{uv} denotes the set of paths from u to v and $\mathbf{w}_q(e_i)$ is the representation of edge e_i . Note the *multiplication* operator \otimes is not required to be commutative (e.g., matrix multiplication), therefore we define \bigotimes to compute the product following the exact order. Intuitively, the path formulation can be interpreted as a depth-first-search (DFS) algorithm, where one searches all possible paths from u to v , computes their representations (Equation 2) and aggregates the results (Equation 1). Such a formulation is capable of modeling several traditional link prediction methods, as well as graph theory algorithms. Formally, Theorem 1-5 state the corresponding path formulations for 3 link prediction methods and 2 graph theory algorithms respectively. See Appendix A for proofs.

Theorem 1 *Katz index is a path formulation with $\oplus = +$, $\otimes = \times$ and $\mathbf{w}_q(e) = \beta w_e$.*

Theorem 2 *Personalized PageRank is a path formulation with $\oplus = +$, $\otimes = \times$ and $\mathbf{w}_q(e) = \alpha w_{uv} / \sum_{v' \in \mathcal{N}(u)} w_{uv'}$.*

Theorem 3 *Graph distance is a path formulation with $\oplus = \min$, $\otimes = +$ and $\mathbf{w}_q(e) = w_e$.*

Theorem 4 *Widest path is a path formulation with $\oplus = \max$, $\otimes = \min$ and $\mathbf{w}_q(e) = w_e$.*

Theorem 5 *Most reliable path is a path formulation with $\oplus = \max$, $\otimes = \times$ and $\mathbf{w}_q(e) = w_e$.*

Generalized Bellman-Ford Algorithm. While the above formulation is able to model important heuristics for link prediction, it is computationally expensive since the number of paths grows exponentially with the path length. Previous works [40, 11, 62] that directly computes the exponential number of paths can only afford a maximal path length of 3. A more scalable solution is to use the generalized Bellman-Ford algorithm [4]. Specifically, assuming the operators $\langle \oplus, \otimes \rangle$ satisfy a semiring system [21] with *summation identity* $\mathbb{0}_q$ and *multiplication identity* $\mathbb{1}_q$, we have the following algorithm.

$$\mathbf{h}_q^{(0)}(u, v) \leftarrow \mathbb{1}_q(u = v) \quad (3)$$

$$\mathbf{h}_q^{(t)}(u, v) \leftarrow \left(\bigoplus_{(x, r, v) \in \mathcal{E}(v)} \mathbf{h}_q^{(t-1)}(u, x) \otimes \mathbf{w}_q(x, r, v) \right) \oplus \mathbf{h}_q^{(0)}(u, v) \quad (4)$$

where $\mathbb{1}_q(u = v)$ is the *indicator* function that outputs $\mathbb{1}_q$ if $u = v$ and $\mathbb{0}_q$ otherwise. $\mathbf{w}_q(x, r, v)$ is the representation for edge $e = (x, r, v)$ and r is the relation type of the edge. Equation 3 is known as the boundary condition, while Equation 4 is known as the Bellman-Ford iteration. The high-level idea of the generalized Bellman-Ford algorithm is to **compute the pair representation $\mathbf{h}_q(u, v)$ for a given entity u , a given query relation q and all $v \in \mathcal{V}$ in parallel**, and reduce the

total computation by the distributive property of *multiplication* over *summation*. Since u and q are fixed in the generalized Bellman-Ford algorithm, we may abbreviate $\mathbf{h}_q^{(t)}(u, v)$ as $\mathbf{h}_v^{(t)}$ when the context is clear. When $\oplus = \min$ and $\otimes = +$, it recovers the original Bellman-Ford algorithm for the shortest path problem [5]. See Appendix B for preliminaries and the proof of the above algorithm.

Theorem 6 *Katz index, personalized PageRank, graph distance, widest path and most reliable path can be solved via the generalized Bellman-Ford algorithm.*

Table 2: Comparison of operators in NBFNet and other methods from the view of path formulation.

Class	Method	MESSAGE $\mathbf{w}_q(e_i) \otimes \mathbf{w}_q(e_j)$	AGGREGATE $\mathbf{h}_q(P_i) \oplus \mathbf{h}_q(P_j)$	INDICATOR $\mathbb{I}_q, \mathbb{I}_q$	Edge Representation $\mathbf{w}_q(e)$
Traditional Link Prediction	Katz Index [30]	$\mathbf{w}_q(e_i) \times \mathbf{w}_q(e_j)$	$\mathbf{h}_q(P_i) + \mathbf{h}_q(P_j)$	0, 1	βw_e
	Personalized PageRank [42]	$\mathbf{w}_q(e_i) \times \mathbf{w}_q(e_j)$	$\mathbf{h}_q(P_i) + \mathbf{h}_q(P_j)$	0, 1	$\alpha w_{uv} / \sum_{v' \in \mathcal{N}(u)} w_{uv'}$
	Graph Distance [37]	$\mathbf{w}_q(e_i) + \mathbf{w}_q(e_j)$	$\min(\mathbf{h}_q(P_i), \mathbf{h}_q(P_j))$	$+\infty, 0$	w_e
Graph Theory Algorithms	Widest Path [4]	$\min(\mathbf{w}_q(e_i), \mathbf{w}_q(e_j))$	$\max(\mathbf{h}_q(P_i), \mathbf{h}_q(P_j))$	$-\infty, +\infty$	w_e
	Most Reliable Path [4]	$\mathbf{w}_q(e_i) \times \mathbf{w}_q(e_j)$	$\max(\mathbf{h}_q(P_i), \mathbf{h}_q(P_j))$	0, 1	w_e
Logic Rules	NeurallP [69] / DRUM [46]	$\mathbf{w}_q(e_i) \times \mathbf{w}_q(e_j)$	$\mathbf{h}_q(P_i) + \mathbf{h}_q(P_j)$	0, 1	Weights learned by LSTM [23]
	NBFNet	Relational operators of knowledge graph embeddings [6, 68, 52]	Learned set aggregators [9]	Learned indicator functions	Learned relation embeddings

3.2 Neural Bellman-Ford Networks

While the generalized Bellman-Ford algorithm can solve many classical methods (Theorem 6), these methods instantiate the path formulation with handcrafted operators (Table 2), and may not be optimal for link prediction. To improve the capacity of path formulation, we propose a general framework, Neural Bellman-Ford Networks (NBFNet), to learn the operators in the pair representations.

Neural Parameterization. We relax the semiring assumption and parameterize the generalized Bellman-Ford algorithm (Equation 3 and 4) with 3 neural functions, namely INDICATOR, MESSAGE and AGGREGATE functions. The INDICATOR function replaces the *indicator* function $\mathbb{I}_q(u = v)$. The MESSAGE function replaces the binary *multiplication* operator \otimes . The AGGREGATE function is a permutation invariant function over sets that replaces the *n-ary summation* operator \oplus . Note that one may alternatively define AGGREGATE as the commutative binary operator \oplus and apply it to a sequence of messages. However, this will make the parameterization more complicated.

Algorithm 1 Neural Bellman-Ford Networks

Input: source node u , query relation q , #layers T

Output: pair representations $\mathbf{h}_q(u, v)$ for all $v \in \mathcal{V}$

```

1: for  $v \in \mathcal{V}$  do ▷ Boundary condition
2:    $\mathbf{h}_v^{(0)} \leftarrow \text{INDICATOR}(u, v, q)$ 
3: end for
4: for  $t \leftarrow 1$  to  $T$  do ▷ Bellman-Ford iteration
5:   for  $v \in \mathcal{V}$  do
6:      $\mathcal{M}_v^{(t)} \leftarrow \{\mathbf{h}_v^{(0)}\}$  ▷ Message augmentation
7:     for  $(x, r, v) \in \mathcal{E}(v)$  do
8:        $\mathbf{m}_{(x, r, v)}^{(t)} \leftarrow \text{MESSAGE}^{(t)}(\mathbf{h}_x^{(t-1)}, \mathbf{w}_q(x, r, v))$ 
9:        $\mathcal{M}_v^{(t)} \leftarrow \mathcal{M}_v^{(t)} \cup \{\mathbf{m}_{(x, r, v)}^{(t)}\}$ 
10:    end for
11:     $\mathbf{h}_v^{(t)} \leftarrow \text{AGGREGATE}^{(t)}(\mathcal{M}_v^{(t)})$ 
12:  end for
13: end for
14: return  $\mathbf{h}_v^{(T)}$  as  $\mathbf{h}_q(u, v)$  for all  $v \in \mathcal{V}$ 

```

Now consider the generalized Bellman-Ford algorithm for a given entity u and relation q . In this context, we abbreviate $\mathbf{h}_q^{(t)}(u, v)$ as $\mathbf{h}_v^{(t)}$, i.e., a representation on entity v in the t -th iteration. It should be stressed that $\mathbf{h}_v^{(t)}$ is still a pair representation, rather than a node representation. By substituting the neural functions into Equation 3 and 4, we get our Neural Bellman-Ford Networks.

$$\mathbf{h}_v^{(0)} \leftarrow \text{INDICATOR}(u, v, q) \quad (5)$$

$$\mathbf{h}_v^{(t)} \leftarrow \text{AGGREGATE} \left(\left\{ \text{MESSAGE} \left(\mathbf{h}_x^{(t-1)}, \mathbf{w}_q(x, r, v) \right) \mid (x, r, v) \in \mathcal{E}(v) \right\} \cup \{\mathbf{h}_v^{(0)}\} \right) \quad (6)$$

NBFNet can be interpreted as a novel GNN framework for learning pair representations. Compared to common GNN frameworks [32, 48] that compute the pair representation as two independent node representations $\mathbf{h}_q(u)$ and $\mathbf{h}_q(v)$, NBFNet initializes a representation on the source node u , and readouts the pair representation on the target node v . Intuitively, our framework can be viewed as a

source-specific message passing process, where every node learns a representation conditioned on the source node. The pseudo code of NBFNet is outlined in Algorithm 1.

Design Space. Now we discuss some principled designs for MESSAGE, AGGREGATE and INDICATOR functions by drawing insights from traditional methods. Note the potential design space for NBFNet is way larger than what is presented here, as one can always borrow MESSAGE and AGGREGATE from the arsenal of message-passing GNNs [19, 16, 60, 65].

For the MESSAGE function, traditional methods instantiate it as natural summation, natural multiplication or min over scalars. Therefore, we may use the vectorized version of summation or multiplication. Intuitively, summation of $\mathbf{h}_x^{(t-1)}$ and $\mathbf{w}_q(x, r, v)$ can be interpreted as a translation of $\mathbf{h}_x^{(t-1)}$ by $\mathbf{w}_q(x, r, v)$ in the pair representation space, while multiplication corresponds to scaling. Such transformations correspond to the relational operators [18, 45] in knowledge graph embeddings [6, 68, 58, 31, 52]. For example, translation and scaling are the relational operators used in TransE [6] and DistMult [68] respectively. We also consider the rotation operator in RotatE [52].

The AGGREGATE function is instantiated as natural summation, max or min in traditional methods, which are reminiscent of set aggregation functions [71, 65, 9] used in GNNs. Therefore, we specify the AGGREGATE function to be sum, mean, or max, followed by a linear transformation and a non-linear activation. We also consider the principal neighborhood aggregation (PNA) proposed in a recent work [9], which jointly learns the types and scales of the aggregation function.

The INDICATOR function is aimed at providing a non-trivial representation for the source node u as the boundary condition. Therefore, we learn a query embedding \mathbf{q} for $\textcircled{1}_q$ and define INDICATOR function as $\mathbb{1}(u = v) * \mathbf{q}$. Note it is also possible to additionally learn an embedding for $\textcircled{0}_q$. However, we find a single query embedding works better in practice.

The edge representations are instantiated as transition probabilities or length in traditional methods. We notice that an edge may have different contribution in answering different query relations. Therefore, we parameterize the edge representations as a linear function over the query relation, i.e., $\mathbf{w}_q(x, r, v) = \mathbf{W}_r \mathbf{q} + \mathbf{b}_r$. For homogeneous graphs or knowledge graphs with very few relations, we simplify the parameterization to $\mathbf{w}_q(x, r, v) = \mathbf{b}_r$ to prevent overfitting. Note that one may also parameterize $\mathbf{w}_q(x, r, v)$ with learnable entity embeddings \mathbf{x} and \mathbf{v} , but such a parameterization cannot solve the inductive setting. Similar to NeuralLP [69] & DRUM [46], we use different edge representations for different iterations, which is able to distinguish noncommutative edges in paths, e.g., father’s mother v.s. mother’s father.

Link Prediction. We now show how to apply the learned pair representations $\mathbf{h}_q(u, v)$ to the link prediction problem. We predict the conditional likelihood of the tail entity v as $p(v|u, q) = \sigma(f(\mathbf{h}_q(u, v)))$, where $\sigma(\cdot)$ is the sigmoid function and $f(\cdot)$ is a feed-forward neural network. The conditional likelihood of the head entity u can be predicted by $p(u|v, q^{-1}) = \sigma(f(\mathbf{h}_{q^{-1}}(v, u)))$ with the same model. Following previous works [6, 52], we minimize the negative log-likelihood of positive and negative triplets (Equation 7). The negative samples are generated according to Partial Completeness Assumption (PCA) [14], which corrupts one of the entities in a positive triplet to create a negative sample. For undirected graphs, we symmetrize the representations and define $p_q(u, v) = \sigma(f(\mathbf{h}_q(u, v) + \mathbf{h}_q(v, u)))$. Equation 8 shows the loss for homogeneous graphs.

$$\mathcal{L}_{KG} = -\log p(u, q, v) - \sum_{i=1}^n \frac{1}{n} \log(1 - p(u'_i, q, v'_i)) \quad (7)$$

$$\mathcal{L}_{homo} = -\log p(u, v) - \sum_{i=1}^n \frac{1}{n} \log(1 - p(u'_i, v'_i)), \quad (8)$$

where n is the number of negative samples per positive sample and (u'_i, q, v'_i) and (u'_i, v'_i) are the i -th negative samples for knowledge graphs and homogeneous graphs, respectively.

Time Complexity. One advantage of NBFNet is that it has a relatively low time complexity during inference⁴. Consider a scenario where a model is required to infer the conditional likelihood of all possible triplets $p(v|u, q)$. We group triplets with the same condition u, q together, where each group contains $|\mathcal{V}|$ triplets. For each group, we only need to execute Algorithm 1 once to get their

⁴Although the same analysis can be applied to training on a fixed number of samples, we note it is less instructive since one can trade-off samples for performance, and the trade-off varies from method to method.

predictions. Since a small constant number of iterations T is enough for NBFNet to converge (Table 6b), Algorithm 1 has a time complexity of $O(|\mathcal{E}|d + |\mathcal{V}|d^2)$, where d is the dimension of representations. Therefore, the amortized time complexity for a single triplet is $O\left(\frac{|\mathcal{E}|d}{|\mathcal{V}|} + d^2\right)$. For a detailed derivation of time complexity of other GNN frameworks, please refer to Appendix C.

4 Experiment

4.1 Experiment Setup

We evaluate NBFNet in three settings, knowledge graph completion, homogeneous graph link prediction and inductive relation prediction. The former two are transductive settings, while the last is an inductive setting. For knowledge graphs, we use FB15k-237 [56] and WN18RR [13]. We use the standard transductive splits [56, 13] and inductive splits [55] of these datasets. For homogeneous graphs, we use Cora, Citeseer and PubMed [49]. Following previous works [32, 12], we split the edges into train/valid/test with a ratio of 85:5:10. Statistics of datasets can be found in Appendix E. Additional experiments of NBFNet on OGB [25] datasets can be found in Appendix G.

Implementation Details. Our implementation generally follows the open source codebases of knowledge graph completion⁵ and homogeneous graph link prediction⁶. For knowledge graphs, we follow [69, 46] and augment each triplet $\langle u, q, v \rangle$ with a flipped triplet $\langle v, q^{-1}, u \rangle$. For homogeneous graphs, we follow [33, 32] and augment each node u with a self loop $\langle u, u \rangle$. We instantiate NBFNet with 6 layers, each with 32 hidden units. The feed-forward network $f(\cdot)$ is set to a 2-layer MLP with 64 hidden units. ReLU is used as the activation function for all hidden layers. We drop out edges that directly connect query node pairs during training to encourage the model to capture longer paths and prevent overfitting. Our model is trained on 4 Tesla V100 GPUs for 20 epochs. We select the models based on their performance on the validation set. See Appendix F for more details.

Evaluation. We follow the filtered ranking protocol [6] for knowledge graph completion. For a test triplet $\langle u, q, v \rangle$, we rank it against all negative triplets $\langle u, q, v' \rangle$ or $\langle u', q, v \rangle$ that do not appear in the knowledge graph. We report mean rank (MR), mean reciprocal rank (MRR) and HITS at N (H@N) for knowledge graph completion. For inductive relation prediction, we follow [55] and draw 50 negative triplets for each positive triplet and use the above filtered ranking. We report HITS@10 for inductive relation prediction. For homogeneous graph link prediction, we follow [32] and compare the positive edges against the same number of negative edges. We report area under the receiver operating characteristic curve (AUROC) and average precision (AP) for homogeneous graphs.

Baselines. We compare NBFNet against path-based methods, embedding methods, and GNNs. These include 11 baselines for knowledge graph completion, 10 baselines for homogeneous graph link prediction and 4 baselines for inductive relation prediction. Note the inductive setting only includes path-based methods and GNNs, since existing embedding methods cannot handle this setting.

4.2 Main Results

Table 3 summarizes the results on knowledge graph completion. NBFNet significantly outperforms existing methods on all metrics and both datasets. NBFNet achieves an average relative gain of 21% in HITS@1 compared to the best path-based method, DRUM [46], on two datasets. Since DRUM is a special instance of NBFNet with natural summation and multiplication operators, this indicates the importance of learning MESSAGE and AGGREGATE functions in NBFNet. NBFNet also outperforms the best embedding method, LowFER [1], with an average relative performance gain of 18% in HITS@1 on two datasets. Meanwhile, NBFNet requires much less parameters than embedding methods. NBFNet only uses 3M parameters on FB15k-237, while TransE needs 30M parameters. See Appendix D for details on the number of parameters.

Table 4 shows the results on homogeneous graph link prediction. NBFNet gets the best results on Cora and PubMed, meanwhile achieves competitive results on CiteSeer. Note CiteSeer is extremely sparse (Appendix E), which makes it hard to learn good representations with NBFNet. One thing to note here is that unlike other GNN methods, NBFNet does not use the node features provided by

⁵<https://github.com/DeepGraphLearning/KnowledgeGraphEmbedding>. MIT license.

⁶<https://github.com/tkipf/gae>. MIT license.

Table 3: Knowledge graph completion results. Results of NeuraLP and DRUM are taken from [46]. Results of RotatE, HAKE and LowFER are taken from their original papers [52, 76, 1]. Results of the other embedding methods are taken from [52]. Since GraIL has scalability issues in this setting, we evaluate it with 50 and 100 negative triplets for FB15k-237 and WN18RR respectively and report MR based on an unbiased estimation.

Class	Method	FB15k-237					WN18RR				
		MR	MRR	H@1	H@3	H@10	MR	MRR	H@1	H@3	H@10
Path-based	Path Ranking [35]	3521	0.174	0.119	0.186	0.285	22438	0.324	0.276	0.360	0.406
	NeuralLP [69]	-	0.240	-	-	0.362	-	0.435	0.371	0.434	0.566
	DRUM [46]	-	0.343	0.255	0.378	0.516	-	0.486	0.425	0.513	0.586
Embeddings	TransE [6]	357	0.294	-	-	0.465	3384	0.226	-	-	0.501
	DistMult [68]	254	0.241	0.155	0.263	0.419	5110	0.43	0.39	0.44	0.49
	ComplEx [58]	339	0.247	0.158	0.275	0.428	5261	0.44	0.41	0.46	0.51
	RotatE [52]	177	0.338	0.241	0.375	0.553	3340	0.476	0.428	0.492	0.571
	HAKE [76]	-	0.346	0.250	0.381	0.542	-	0.497	0.452	0.516	0.582
	LowFER [1]	-	0.359	0.266	0.396	0.544	-	0.465	0.434	0.479	0.526
GNNs	RGCN [48]	221	0.273	0.182	0.303	0.456	2719	0.402	0.345	0.437	0.494
	GraIL [55]	2053	-	-	-	-	2539	-	-	-	-
	NBFNet	114	0.415	0.321	0.454	0.599	636	0.551	0.497	0.573	0.666

Table 4: Homogeneous graph link prediction results. Results of VGAE and S-VGAE are taken from their original papers [32, 12].

Class	Method	Cora		Citeseer		PubMed	
		AUROC	AP	AUROC	AP	AUROC	AP
Path-based	Katz Index [30]	0.834	0.889	0.768	0.810	0.757	0.856
	Personalized PageRank [42]	0.845	0.899	0.762	0.814	0.763	0.860
	SimRank [28]	0.838	0.888	0.755	0.805	0.743	0.829
Embeddings	DeepWalk [43]	0.831	0.850	0.805	0.836	0.844	0.841
	LINE [53]	0.844	0.876	0.791	0.826	0.849	0.888
	node2vec [17]	0.872	0.879	0.838	0.868	0.891	0.914
GNNs	VGAE [32]	0.914	0.926	0.908	0.920	0.944	0.947
	S-VGAE [12]	0.941	0.941	0.947	0.952	0.960	0.960
	SEAL [73]	0.933	0.942	0.905	0.924	0.978	0.979
	TLC-GNN [67]	0.934	0.931	0.909	0.916	0.970	0.968
	NBFNet	0.956	0.962	0.923	0.936	0.983	0.982

Table 5: Inductive relation prediction results (HITS@10). V1-v4 corresponds to the 4 standard versions of inductive splits. Results of compared methods are taken from [55].

Class	Method	FB15k-237				WN18RR			
		v1	v2	v3	v4	v1	v2	v3	v4
Path-based	NeuralLP [16]	0.529	0.589	0.529	0.559	0.744	0.689	0.462	0.671
	DRUM [46]	0.529	0.587	0.529	0.559	0.744	0.689	0.462	0.671
	RuleN [39]	0.498	0.778	0.877	0.856	0.809	0.782	0.534	0.716
GNNs	GraIL [55]	0.642	0.818	0.828	0.893	0.825	0.787	0.584	0.734
	NBFNet	0.834	0.949	0.951	0.960	0.948	0.905	0.893	0.890

the datasets but is still able to outperform most other methods. We leave how to effectively combine node features and structural representations for link prediction as our future work.

Table 5 summarizes the results on inductive relation prediction. On all inductive splits of two datasets, NBFNet achieves the best result. NBFNet outperforms the previous best method, GraIL [55], with an average relative performance gain of 22% in HITS@10. Note that GraIL explicitly encodes the local subgraph surrounding each node pair and has a high time complexity (Appendix C). Usually, GraIL can at most encode a 2-hop subgraph, while our NBFNet can efficiently explore longer paths.

4.3 Ablation Study

MESSAGE & AGGREGATE Functions. Table 6a shows the results of different MESSAGE and AGGREGATE functions. Generally, NBFNet benefits from advanced embedding methods (DistMult,

RotatE > TransE) and aggregation functions (PNA > sum, mean, max). Among simple AGGREGATE functions (sum, mean, max), combinations of MESSAGE and AGGREGATE functions (TransE & max, DistMult & sum) that satisfy the semiring assumption⁷ of the generalized Bellman-Ford algorithm, achieve locally optimal performance. PNA significantly improves over simple counterparts, which highlights the importance of learning more powerful AGGREGATE functions.

Number of GNN Layers. Table 6b compares the results of NBFNet with different number of layers. Although it has been reported that GNNs with deep layers often result in significant performance drop [36, 77], we observe NBFNet does not have this issue. The performance increases monotonically with more layers, hitting a saturation after 6 layers. We conjecture the reason is that longer paths have negligible contribution, and paths not longer than 6 are enough for link prediction.

Performance by Relation Category. We break down the performance of NBFNet by the categories of query relations: one-to-one, one-to-many, many-to-one and many-to-many⁸. Table 6c shows the prediction results for each category. It is observed that NBFNet not only improves on easy one-to-one cases, but also on hard cases where there are multiple true answers for the query.

Table 6: Ablation studies of NBFNet on FB15k-237. Due to space constraints, we only report MRR here. For full results on all metrics, please refer to Appendix H.

(a) Different MESSAGE and AGGREGATE functions.					(b) Different number of layers.				
MESSAGE	AGGREGATE				Method	#Layers (T)			
	Sum	Mean	Max	PNA [9]		2	4	6	8
TransE [6]	0.297	0.310	0.377	0.383	NBFNet	0.345	0.409	0.415	0.416
DistMult [69]	0.388	0.384	0.374	0.415					
RotatE [52]	0.392	0.376	0.385	0.414					

(c) Performance w.r.t. relation category. The two scores are the rankings over heads and tails respectively.				
Method	Relation Category			
	1-to-1	1-to-N	N-to-1	N-to-N
TransE [6]	0.498/0.488	0.455/0.071	0.079/0.744	0.224/0.330
RotatE [51]	0.487/0.484	0.467/0.070	0.081/0.747	0.234/0.338
NBFNet	0.578/0.600	0.499/0.122	0.165/0.790	0.348/0.456

4.4 Path Interpretations of Predictions

One advantage of NBFNet is that we can interpret its predictions through paths, which may be important for users to understand and debug the model. Intuitively, the interpretations should contain paths that contribute most to the prediction $p(u, q, v)$. Following local interpretation methods [3, 72], we approximate the local landscape of NBFNet with a linear model over the set of all paths, i.e., 1st-order Taylor polynomial. We define the importance of a path as its weight in the linear model, which can be computed by the partial derivative of the prediction w.r.t. the path. Formally, the top-k path interpretations for $p(u, q, v)$ are defined as

$$P_1, P_2, \dots, P_k = \text{top-k}_{P \in \mathcal{P}_{uv}} \frac{\partial p(u, q, v)}{\partial P} \quad (9)$$

Note this formulation generalizes the definition of logical rules [69, 46] to non-linear models. While directly computing the importance of all paths is intractable, we approximate them with edge importance. Specifically, the importance of each path is approximated by the sum of the importance of edges in that path, where edge importance is obtained via auto differentiation. Then the top-k path interpretations are equivalent to the top-k longest paths on the edge importance graph, which can be solved by a Bellman-Ford-style beam search. Better approximation is left as a future work.

Table 7 visualizes path interpretations from FB15k-237 test set. While users may have different insights towards the visualization, here is our understanding. 1) In the first example, NBFNet learns

⁷Here semiring is discussed under the assumption of linear activation functions. Rigorously, no combination satisfies a semiring if we consider non-linearity in the model.

⁸The categories are defined same as [63]. We compute the average number of tails per head and the average number of heads per tail. The category is *one* if the average number is smaller than 1.5 and *many* otherwise.

soft logical entailment, such as $\text{impersonate}^{-1} \wedge \text{nationality} \implies \text{nationality}$ and $\text{ethnicity}^{-1} \wedge \text{distribution} \implies \text{nationality}$. 2) In second example, NBFNet performs analogical reasoning by leveraging the fact that *Florence* is similar to *Rome*. 3) In the last example, NBFNet extracts longer paths, since there is no obvious connection between *Pearl Harbor (film)* and *Japanese language*.

Table 7: Path interpretations of predictions on FB15k-237 test set. For each query triplet, we visualize the top-2 path interpretations and their weights. Inverse relations are denoted with a superscript $^{-1}$.

Query	$\langle u, q, v \rangle$: $\langle O. Hardy, \text{nationality}, U.S. \rangle$
0.243	$\langle O. Hardy, \text{impersonate}^{-1}, R. Little \rangle \wedge \langle R. Little, \text{nationality}, U.S. \rangle$
0.224	$\langle O. Hardy, \text{ethnicity}^{-1}, \text{Scottish American} \rangle \wedge \langle \text{Scottish American}, \text{distribution}, U.S. \rangle$
Query	$\langle u, q, v \rangle$: $\langle \text{Florence}, \text{vacationer}, D.C. Henrie \rangle$
0.251	$\langle \text{Florence}, \text{contain}^{-1}, \text{Italy} \rangle \wedge \langle \text{Italy}, \text{capital}, \text{Rome} \rangle \wedge \langle \text{Rome}, \text{vacationer}, D.C. Henrie \rangle$
0.183	$\langle \text{Florence}, \text{place live}^{-1}, G.F. Handel \rangle \wedge \langle G.F. Handel, \text{place live}, \text{Rome} \rangle \wedge \langle \text{Rome}, \text{vacationer}, D.C. Henrie \rangle$
Query	$\langle u, q, v \rangle$: $\langle \text{Pearl Harbor (film)}, \text{language}, \text{Japanese} \rangle$
0.211	$\langle \text{Pearl Harbor (film)}, \text{film actor}, C.-H. Tagawa \rangle \wedge \langle C.-H. Tagawa, \text{nationality}, \text{Japan} \rangle$ $\wedge \langle \text{Japan}, \text{country of origin}, \text{Yu-Gi-Oh!} \rangle \wedge \langle \text{Yu-Gi-Oh!}, \text{language}, \text{Japanese} \rangle$
0.208	$\langle \text{Pearl Harbor (film)}, \text{film actor}, C.-H. Tagawa \rangle \wedge \langle C.-H. Tagawa, \text{nationality}, \text{Japan} \rangle$ $\wedge \langle \text{Japan}, \text{official language}, \text{Japanese} \rangle$

5 Discussion and Conclusion

Limitations. There are a few limitations for NBFNet. First, the assumption of the generalized Bellman-Ford algorithm requires the operators $\langle \oplus, \otimes \rangle$ to satisfy a semiring. Due to the non-linear activation functions in neural networks, this assumption does not hold for NBFNet, and we do not have a theoretical guarantee on the loss incurred by this relaxation. Second, NBFNet is only verified on simple edge prediction, while there are other link prediction variants, e.g., complex logical queries with conjunctions (\wedge) and disjunctions (\vee) [18, 45]. In the future, we would like to how NBFNet approximates the path formulation, as well as apply NBFNet to other link prediction settings.

Social Impacts. Link prediction has a wide range of beneficial applications, including recommender systems, knowledge graph completion and drug repurposing. However, there are also some potentially negative impacts. First, NBFNet may encode the bias present in the training data, which leads to stereotyped predictions when the prediction is applied to a user on a social or e-commerce platform. Second, some harmful network activities could be augmented by powerful link prediction models, e.g., spamming, phishing, and social engineering. We expect future studies will mitigate these issues.

Conclusion. We present a representation learning framework based on paths for link prediction. Our path formulation generalizes several traditional methods, and can be efficiently solved via the generalized Bellman-Ford algorithm. To improve the capacity of the path formulation, we propose NBFNet, which parameterizes the generalized Bellman-Ford algorithm with learned INDICATOR, MESSAGE, AGGREGATE functions. Experiments on knowledge graphs and homogeneous graphs show that NBFNet outperforms a wide range of methods in both transductive and inductive settings.

Acknowledgements

We would like to thank Komal Teru for discussion on inductive relation prediction, Guyue Huang for discussion on fused message passing implementation, and Yao Lu for assistance on large-scale GPU training. We thank Meng Qu, Chence Shi and Minghao Xu for providing feedback on our manuscript.

This project is supported by the Natural Sciences and Engineering Research Council (NSERC) Discovery Grant, the Canada CIFAR AI Chair Program, collaboration grants between Microsoft Research and Mila, Samsung Electronics Co., Ltd., Amazon Faculty Research Award, Tencent AI Lab Rhino-Bird Gift Fund and a NRC Collaborative R&D Project (AI4D-CORE-06). This project was also partially funded by IVADO Fundamental Research Project grant PRF-2019-3583139727. The computation resource of this project is supported by Calcul Québec⁹ and Compute Canada¹⁰.

⁹<https://www.calculquebec.ca/>

¹⁰<https://www.computecanada.ca/>

References

- [1] Saadullah Amin, Stalin Varanasi, Katherine Ann Dunfield, and Günter Neumann. Lower: Low-rank bilinear pooling for link prediction. In *International Conference on Machine Learning*, pages 257–268. PMLR, 2020.
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [3] David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. How to explain individual classification decisions. *The Journal of Machine Learning Research*, 11:1803–1831, 2010.
- [4] John S Baras and George Theodorakopoulos. Path problems in networks. *Synthesis Lectures on Communication Networks*, 3(1):1–77, 2010.
- [5] Richard Bellman. On a routing problem. *Quarterly of applied mathematics*, 16(1):87–90, 1958.
- [6] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In *Advances in Neural Information Processing Systems*, pages 1–9, 2013.
- [7] Linlin Chao, Jianshan He, Taifeng Wang, and Wei Chu. PairRE: Knowledge graph embeddings via paired relation vectors. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4360–4369, 2021.
- [8] Wenhu Chen, Wenhan Xiong, Xifeng Yan, and William Yang Wang. Variational knowledge graph reasoning. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1823–1832, 2018.
- [9] Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal neighbourhood aggregation for graph nets. volume 33, 2020.
- [10] Rajarshi Das, Shehzaad Dhuliawala, Manzil Zaheer, Luke Vilnis, Ishan Durugkar, Akshay Krishnamurthy, Alex Smola, and Andrew McCallum. Go for a walk and arrive at the answer: Reasoning over paths in knowledge bases using reinforcement learning. In *International Conference on Learning Representations*, 2018.
- [11] Rajarshi Das, Arvind Neelakantan, David Belanger, and Andrew McCallum. Chains of reasoning over entities, relations, and text using recurrent neural networks. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 132–141, Valencia, Spain, April 2017. Association for Computational Linguistics.
- [12] Tim R Davidson, Luca Falorsi, Nicola De Cao, Thomas Kipf, and Jakub M Tomczak. Hyper-spherical variational auto-encoders. 2018.
- [13] Tim Dettmers, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. Convolutional 2d knowledge graph embeddings. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [14] Luis Antonio Galárraga, Christina Teflioudi, Katja Hose, and Fabian Suchanek. Amie: association rule mining under incomplete evidence in ontological knowledge bases. In *Proceedings of the 22nd international conference on World Wide Web*, pages 413–422, 2013.
- [15] Matt Gardner and Tom Mitchell. Efficient and expressive knowledge base completion using subgraph feature extraction. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1488–1498, 2015.
- [16] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International Conference on Machine Learning*, pages 1263–1272. PMLR, 2017.
- [17] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 855–864, 2016.

- [18] William L Hamilton, Payal Bajaj, Marinka Zitnik, Dan Jurafsky, and Jure Leskovec. Embedding logical queries on knowledge graphs. In *Advances in Neural Information Processing Systems*, pages 2030–2041, 2018.
- [19] William L Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, pages 1025–1035, 2017.
- [20] Zhen Han, Peng Chen, Yunpu Ma, and Volker Tresp. xerte: Explainable reasoning on temporal knowledge graphs for forecasting future links. 2021.
- [21] Udo Hebisch and Hanns Joachim Weinert. *Semirings: algebraic theory and applications in computer science*, volume 5. World Scientific, 1998.
- [22] Marcel Hildebrandt, Jorge Andres Quintero Serna, Yunpu Ma, Martin Ringsquandl, Mitchell Joblin, and Volker Tresp. Reasoning on knowledge graphs with debate dynamics. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 4123–4131, 2020.
- [23] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [24] Weihua Hu, Matthias Fey, Hongyu Ren, Maho Nakata, Yuxiao Dong, and Jure Leskovec. Ogb-lsc: A large-scale challenge for machine learning on graphs. *arXiv preprint arXiv:2103.09430*, 2021.
- [25] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *arXiv preprint arXiv:2005.00687*, 2020.
- [26] Guyue Huang, Guohao Dai, Yu Wang, and Huazhong Yang. Ge-spm: General-purpose sparse matrix-matrix multiplication on gpus for graph neural networks. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–12. IEEE, 2020.
- [27] Vassilis N Ioannidis, Da Zheng, and George Karypis. Few-shot link prediction via graph neural networks for covid-19 drug-repurposing. *arXiv preprint arXiv:2007.10261*, 2020.
- [28] Glen Jeh and Jennifer Widom. Simrank: a measure of structural-context similarity. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 538–543, 2002.
- [29] Glen Jeh and Jennifer Widom. Scaling personalized web search. In *Proceedings of the 12th international conference on World Wide Web*, pages 271–279, 2003.
- [30] Leo Katz. A new status index derived from sociometric analysis. *Psychometrika*, 18(1):39–43, 1953.
- [31] Seyed Mehran Kazemi and David Poole. Simple embedding for link prediction in knowledge graphs. In *Advances in Neural Information Processing Systems*, pages 4289–4300, 2018.
- [32] Thomas N Kipf and Max Welling. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308*, 2016.
- [33] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017.
- [34] Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, 2009.
- [35] Ni Lao and William W Cohen. Relational retrieval using a combination of path-constrained random walks. *Machine learning*, 81(1):53–67, 2010.
- [36] Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks for semi-supervised learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [37] David Liben-Nowell and Jon Kleinberg. The link-prediction problem for social networks. *Journal of the American society for information science and technology*, 58(7):1019–1031, 2007.
- [38] Xi Victoria Lin, Richard Socher, and Caiming Xiong. Multi-hop knowledge graph reasoning with reward shaping. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, EMNLP 2018, Brussels, Belgium, October 31-November 4, 2018*, 2018.

- [39] Christian Meilicke, Manuel Fink, Yanjie Wang, Daniel Ruffinelli, Rainer Gemulla, and Heiner Stuckenschmidt. Fine-grained evaluation of rule-and embedding-based systems for knowledge graph completion. In *International Semantic Web Conference*, pages 3–20. Springer, 2018.
- [40] Arvind Neelakantan, Benjamin Roth, and Andrew McCallum. Compositional vector space models for knowledge base completion. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 156–166, Beijing, China, July 2015. Association for Computational Linguistics.
- [41] Maximilian Nickel, Kevin Murphy, Volker Tresp, and Evgeniy Gabrilovich. A review of relational machine learning for knowledge graphs. *Proceedings of the IEEE*, 104(1):11–33, 2015.
- [42] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.
- [43] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 701–710, 2014.
- [44] Meng Qu, Junkun Chen, Louis-Pascal Xhonneux, Yoshua Bengio, and Jian Tang. Rnnlogic: Learning logic rules for reasoning on knowledge graphs. In *International Conference on Learning Representations*, 2021.
- [45] H Ren, W Hu, and J Leskovec. Query2box: Reasoning over knowledge graphs in vector space using box embeddings. In *International Conference on Learning Representations*, 2020.
- [46] Ali Sadeghian, Mohammadreza Armandpour, Patrick Ding, and Daisy Zhe Wang. Drum: End-to-end differentiable rule mining on knowledge graphs. volume 32, pages 15347–15357, 2019.
- [47] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008.
- [48] Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. Modeling relational data with graph convolutional networks. In *European semantic web conference*, pages 593–607. Springer, 2018.
- [49] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.
- [50] Yelong Shen, Jianshu Chen, Po-Sen Huang, Yuqing Guo, and Jianfeng Gao. M-walk: learning to walk over graphs using monte carlo tree search. In *Advances in Neural Information Processing Systems*, pages 6787–6798, 2018.
- [51] Yizhou Sun, Jiawei Han, Xifeng Yan, Philip S Yu, and Tianyi Wu. Pathsim: Meta path-based top-k similarity search in heterogeneous information networks. volume 4, pages 992–1003. VLDB Endowment, 2011.
- [52] Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. Rotate: Knowledge graph embedding by relational rotation in complex space. In *International Conference on Learning Representations*, 2019.
- [53] Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. Line: Large-scale information network embedding. In *Proceedings of the 24th international conference on World Wide Web*, pages 1067–1077, 2015.
- [54] Yun Tang, Jing Huang, Guangtao Wang, Xiaodong He, and Bowen Zhou. Orthogonal relation transforms with graph context modeling for knowledge graph embedding. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2713–2722, 2020.
- [55] Komal Teru, Etienne Denis, and Will Hamilton. Inductive relation prediction by subgraph reasoning. In *International Conference on Machine Learning*, pages 9448–9457. PMLR, 2020.
- [56] Kristina Toutanova and Danqi Chen. Observed versus latent features for knowledge base and text inference. In *Proceedings of the 3rd workshop on continuous vector space models and their compositionality*, pages 57–66, 2015.

- [57] Kristina Toutanova, Xi Victoria Lin, Wen-tau Yih, Hoifung Poon, and Chris Quirk. Compositional learning of embeddings for relation paths in knowledge base and text. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1434–1444, 2016.
- [58] Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. Complex embeddings for simple link prediction. In *International Conference on Machine Learning*, pages 2071–2080. PMLR, 2016.
- [59] Shikhar Vashishth, Soumya Sanyal, Vikram Nitin, and Partha Talukdar. Composition-based multi-relational graph convolutional networks. In *International Conference on Learning Representations*, 2020.
- [60] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018.
- [61] Andrew Viterbi. Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE transactions on Information Theory*, 13(2):260–269, 1967.
- [62] Hongwei Wang, Hongyu Ren, and Jure Leskovec. Relational message passing for knowledge graph completion. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 1697–1707, 2021.
- [63] Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. Knowledge graph embedding by translating on hyperplanes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 28, 2014.
- [64] Wenhan Xiong, Thien Hoang, and William Yang Wang. Deeppath: A reinforcement learning method for knowledge graph reasoning. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP 2017)*, Copenhagen, Denmark, September 2017. ACL.
- [65] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
- [66] Xiaoran Xu, Wei Feng, Yunsheng Jiang, Xiaohui Xie, Zhiqing Sun, and Zhi-Hong Deng. Dynamically pruned message passing networks for large-scale knowledge graph reasoning. In *International Conference on Learning Representations*, 2019.
- [67] Zuoyu Yan, Tengfei Ma, Liangcai Gao, Zhi Tang, and Chao Chen. Link prediction with persistent homology: An interactive view. In *International Conference on Machine Learning*, pages 11659–11669. PMLR, 2021.
- [68] Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding entities and relations for learning and inference in knowledge bases. In *International Conference on Learning Representations*, 2015.
- [69] Fan Yang, Zhilin Yang, and William W Cohen. Differentiable learning of logical rules for knowledge base reasoning. In *Advances in Neural Information Processing Systems*, pages 2316–2325, 2017.
- [70] Jiaxuan You, Jonathan M Gomes-Selman, Rex Ying, and Jure Leskovec. Identity-aware graph neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 10737–10745, 2021.
- [71] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. volume 30, 2017.
- [72] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014.
- [73] Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. volume 31, pages 5165–5175, 2018.
- [74] Muhan Zhang, Pan Li, Yinglong Xia, Kai Wang, and Long Jin. Revisiting graph neural networks for link prediction. *arXiv preprint arXiv:2010.16103*, 2020.
- [75] Yongqi Zhang, Quanming Yao, Wenyuan Dai, and Lei Chen. Autosf: Searching scoring functions for knowledge graph embedding. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*, pages 433–444. IEEE, 2020.

- [76] Zhanqiu Zhang, Jianyu Cai, Yongdong Zhang, and Jie Wang. Learning hierarchy-aware knowledge graph embeddings for link prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3065–3072, 2020.
- [77] Lingxiao Zhao and Leman Akoglu. Pairnorm: Tackling oversmoothing in gnns. In *International Conference on Learning Representations*, 2019.
- [78] Zhaocheng Zhu, Shizhen Xu, Meng Qu, and Jian Tang. Graphvite: A high-performance cpu-gpu hybrid system for node embedding. In *The World Wide Web Conference*, pages 2494–2504, 2019.

A Path Formulations for Traditional Methods

Here we demonstrate our path formulation is capable of modeling traditional link prediction methods like Katz index [30], personalized PageRank [42] and graph distance [37], as well as graph theory algorithms like widest path [4] and most reliable path [4].

Recall the path formulation is defined as

$$\mathbf{h}_q(u, v) = \mathbf{h}_q(P_1) \oplus \mathbf{h}_q(P_2) \oplus \dots \oplus \mathbf{h}_q(P_{|\mathcal{P}_{uv}|})|_{P_i \in \mathcal{P}_{uv}} \triangleq \bigoplus_{P \in \mathcal{P}_{uv}} \mathbf{h}_q(P) \quad (1)$$

$$\mathbf{h}_q(P = (e_1, e_2, \dots, e_{|P|})) = \mathbf{w}_q(e_1) \otimes \mathbf{w}_q(e_2) \otimes \dots \otimes \mathbf{w}_q(e_{|P|}) \triangleq \bigotimes_{i=1}^{|P|} \mathbf{w}_q(e_i) \quad (2)$$

which can be written in the following compact form

$$\mathbf{h}_q(u, v) = \bigoplus_{P \in \mathcal{P}_{uv}} \bigotimes_{i=1}^{|P|} \mathbf{w}_q(e_i) \quad (10)$$

A.1 Katz Index

The Katz index for a pair of nodes u, v is defined as a weighted count of paths between u and v , penalized by an attenuation factor $\beta \in (0, 1)$. Formally, it can be written as

$$\text{Katz}(u, v) = \sum_{t=1}^{\infty} \beta^t \mathbf{e}_u^\top \mathbf{A}^t \mathbf{e}_v \quad (11)$$

where \mathbf{A} denotes the adjacency matrix and $\mathbf{e}_u, \mathbf{e}_v$ denote the one-hot vector for nodes u, v respectively. The term $\mathbf{e}_u^\top \mathbf{A}^t \mathbf{e}_v$ counts all paths of length t between u , and v and shorter paths are assigned with larger weights.

Theorem 1 *Katz index is a path formulation with $\oplus = +$, $\otimes = \times$ and $\mathbf{w}_q(e) = \beta w_e$.*

Proof. We show that $\text{Katz}(u, v)$ can be transformed into a summation over all paths between u and v , where each path is represented by a product of damped edge weights in the path. Mathematically, it can be derived as

$$\text{Katz}(u, v) = \sum_{t=1}^{\infty} \beta^t \sum_{P \in \mathcal{P}_{uv}: |P|=t} \prod_{e \in P} w_e \quad (12)$$

$$= \sum_{P \in \mathcal{P}_{uv}} \prod_{e \in P} \beta w_e \quad (13)$$

Therefore, the Katz index can be viewed as a path formulation with the *summation* operator $+$, the *multiplication* operator \times and the edge representations βw_e . \square

A.2 Personalized PageRank

The personalized PageRank (PPR) for u computes the stationary distribution over nodes generated by an infinite random walker, where the walker moves to a neighbor node with probability α and returns to the source node u with probability $1 - \alpha$ at each step. The probability of a node v from a source node u has the following closed-form solution [29]

$$\text{PPR}(u, v) = (1 - \alpha) \sum_{t=1}^{\infty} \alpha^t \mathbf{e}_u^\top (\mathbf{D}^{-1} \mathbf{A})^t \mathbf{e}_v \quad (14)$$

where \mathbf{D} is the degree matrix and $\mathbf{D}^{-1} \mathbf{A}$ is the (random walk) normalized adjacency matrix. Note that $\mathbf{e}_u^\top (\mathbf{D}^{-1} \mathbf{A})^t \mathbf{e}_v$ computes the probability of t -step random walks from u to v .

Theorem 2 *Personalized PageRank is a path formulation with $\oplus = +$, $\otimes = \times$ and $\mathbf{w}_q(e) = \alpha w_{uv} / \sum_{v' \in \mathcal{N}(u)} w_{uv'}$.*

Proof. We omit the coefficient $1 - \alpha$, since it is always positive and has no effect on the ranking of different node pairs. Then we have

$$\text{PPR}(u, v) \propto \sum_{t=1}^{\infty} \alpha^t \sum_{P \in \mathcal{P}_{uv}: |P|=t} \prod_{(a,b) \in P} \frac{w_{ab}}{\sum_{b' \in \mathcal{N}(a)} w_{ab'}} \quad (15)$$

$$= \sum_{P \in \mathcal{P}_{uv}} \prod_{(a,b) \in P} \frac{\alpha w_{ab}}{\sum_{b' \in \mathcal{N}(a)} w_{ab'}} \quad (16)$$

where the *summation* operator is $+$, the *multiplication* operator is \times and edge representations are random walk probabilities scaled by α . \square

A.3 Graph Distance

Graph distance (GD) is defined as the minimum length of all paths between u and v .

Theorem 3 *Graph distance is a path formulation with $\oplus = \min$, $\otimes = +$ and $w_q(e) = w_e$.*

Proof. Since the length of a path is the sum of edge lengths in the path, we have

$$\text{GD}(u, v) = \min_{P \in \mathcal{P}_{uv}} \sum_{e \in P} w_e \quad (17)$$

Here the *summation* operator is \min , the *multiplication* operator is $+$ and the edge representations are the lengths of edges. \square

A.4 Widest Path

The widest path (WP), also known as the maximum capacity path, is aimed at finding a path between two given nodes, such that the path maximizes the minimum edge weight in the path.

Theorem 4 *Widest path is a path formulation with $\oplus = \max$, $\otimes = \min$ and $w_q(e) = w_e$.*

Proof. Given two nodes u and v , we can write the widest path as

$$\text{WP}(u, v) = \max_{P \in \mathcal{P}_{uv}} \min_{e \in P} w_e \quad (18)$$

Here the *summation* operator is \max , the *multiplication* operator is \min and the edge representations are plain edge weights. \square

A.5 Most Reliable Path

For a graph with non-negative edge probabilities, the most reliable path (MRP) is the path with maximal probability from a start node to an end node. This is also known as Viterbi algorithm [61] used in the maximum a posterior (MAP) inference of hidden Markov models (HMM).

Theorem 5 *Most reliable path is a path formulation with $\oplus = \max$, $\otimes = \times$ and $w_q(e) = w_e$.*

Proof. For a start node u and an end node v , the probability of their most reliable path is

$$\text{MRP}(u, v) = \max_{P \in \mathcal{P}_{uv}} \prod_{e \in P} w_e \quad (19)$$

Here the *summation* operator is \max , the *multiplication* operator is \times and the edge representations are edge probabilities. \square

B Generalized Bellman-Ford Algorithm

First, we prove that the path formulation can be efficiently solved by the generalized Bellman-Ford algorithm when the operators $\langle \oplus, \otimes \rangle$ satisfy a semiring. Then, we show that traditional methods satisfy the semiring assumption and therefore can be solved by the generalized Bellman-Ford algorithm.

B.1 Preliminaries on Semirings

Semirings are algebraic structures with two operators, *summation* \oplus and *multiplication* \otimes , that share similar properties with the natural summation and the natural multiplication defined on integers. Specifically, \oplus should be commutative, associative and have an identity element $\mathbb{0}$. \otimes should be associative and have an identity element $\mathbb{1}$. Mathematically, the *summation* \oplus satisfies

- **Commutative Property.** $a \oplus b = b \oplus a$
- **Associative Property.** $(a \oplus b) \oplus c = a \oplus (b \oplus c)$
- **Identity Element.** $a \oplus \mathbb{0} = a$

The *multiplication* \otimes satisfies

- **Associative Property.** $(a \otimes b) \otimes c = a \otimes (b \otimes c)$
- **Absorption Property.** $a \otimes \mathbb{0} = \mathbb{0} \otimes a = \mathbb{0}$
- **Identity Element.** $a \otimes \mathbb{1} = \mathbb{1} \otimes a = a$

Additionally, \otimes should be distributive over \oplus .

- **Distributive Property (Left).** $a \otimes (b \oplus c) = (a \otimes b) \oplus (a \otimes c)$
- **Distributive Property (Right).** $(b \oplus c) \otimes a = (b \otimes a) \oplus (c \otimes a)$

Note semirings differ from natural arithmetic operators in two aspects. First, the *summation* operator \oplus does not need to be invertible, e.g., min or max. Second, the *multiplication* operator \otimes does not need to be commutative nor invertible, e.g., matrix multiplication.

B.2 Generalized Bellman-Ford Algorithm for Path Formulation

Now we prove that the generalized Bellman-Ford algorithm computes the path formulation when the operators $\langle \oplus, \otimes \rangle$ satisfy a semiring. It should be stressed that the generalized Bellman-Ford algorithm for path problems has been proved in [4], and not a contribution of this paper. Here we apply the proof to our proposed path formulation.

The generalized Bellman-Ford algorithm computes the following iterations for all $v \in \mathcal{V}$

$$\mathbf{h}_q^{(0)}(u, v) \leftarrow \mathbb{1}_q(u = v) \quad (3)$$

$$\mathbf{h}_q^{(t)}(u, v) \leftarrow \left(\bigoplus_{(x,r,v) \in \mathcal{E}(v)} \mathbf{h}_q^{(t-1)}(u, x) \otimes \mathbf{w}_q(x, r, v) \right) \oplus \mathbf{h}_q^{(0)}(u, v) \quad (4)$$

Lemma 1 *After t Bellman-Ford iterations, the intermediate representation $\mathbf{h}_q^{(t)}(u, v)$ aggregates all path representations within a length of t edges for all v . That is*

$$\mathbf{h}_q^{(t)}(u, v) = \bigoplus_{P \in \mathcal{P}_{uv}: |P| \leq t} \bigotimes_{i=1}^{|P|} \mathbf{w}_q(e_i) \quad (20)$$

Proof. We prove Lemma 1 by induction. For the base case $t = 0$, there is a single path of length 0 from u to itself and no path to other nodes. Due to the product definition of path representations, a path of length 0 is equal to the *multiplication* identity $\mathbb{1}_q$. Similarly, a summation of no path is equal to the *summation* identity $\mathbb{0}_q$. Therefore, we have $\mathbf{h}_q^{(0)}(u, v) = \mathbb{1}_q(u = v) = \bigoplus_{P \in \mathcal{P}_{uv}: |P|=0} \bigotimes_{i=1}^{|P|} \mathbf{w}_q(e_i)$.

For the inductive case $t > 0$, we consider the second-to-last node x in each path if the path has a length larger than 0. To avoid overuse of brackets, we use the convention that \otimes and \bigotimes have a higher

priority than \oplus and \otimes .

$$\mathbf{h}_q^{(t)}(u, v) = \left(\bigoplus_{(x, r, v) \in \mathcal{E}(v)} \mathbf{h}_q^{(t-1)}(u, x) \otimes \mathbf{w}_q(x, r, v) \right) \oplus \mathbf{h}_q^{(0)}(u, v) \quad (21)$$

$$= \left[\bigoplus_{(x, r, v) \in \mathcal{E}(v)} \left(\bigoplus_{P \in \mathcal{P}_{ux}: |P| \leq t-1} \bigotimes_{i=1}^{|P|} \mathbf{w}_q(e_i) \right) \otimes \mathbf{w}_q(x, r, v) \right] \oplus \mathbf{h}_q^{(0)}(u, v) \quad (22)$$

$$= \left\{ \bigoplus_{(x, r, v) \in \mathcal{E}(v)} \left[\bigoplus_{P \in \mathcal{P}_{ux}: |P| \leq t-1} \left(\bigotimes_{i=1}^{|P|} \mathbf{w}_q(e_i) \right) \otimes \mathbf{w}_q(x, r, v) \right] \right\} \oplus \mathbf{h}_q^{(0)}(u, v) \quad (23)$$

$$= \left(\bigoplus_{P \in \mathcal{P}_{uv}: 1 \leq |P| \leq t} \bigotimes_{i=1}^{|P|} \mathbf{w}_q(e_i) \right) \oplus \left(\bigoplus_{P \in \mathcal{P}_{uv}: |P|=0} \bigotimes_{i=1}^{|P|} \mathbf{w}_q(e_i) \right) \quad (24)$$

$$= \bigoplus_{P \in \mathcal{P}_{uv}: |P| \leq t} \bigotimes_{i=1}^{|P|} \mathbf{w}_q(e_i), \quad (25)$$

where Equation 22 substitutes the inductive assumption for $\mathbf{h}_q^{(t-1)}(u, x)$, Equation 23 uses the distributive property of \otimes over \oplus . \square

By comparing Lemma 1 and Equation 10, we can see the intermediate representation converges to our path formulation $\lim_{t \rightarrow \infty} \mathbf{h}_q^{(t)}(u, v) = \mathbf{h}_q(u, v)$. More specifically, at most $|\mathcal{V}|$ iterations are required if we only consider simple paths, i.e., paths without repeating nodes. In practice, for link prediction we find it only takes a very small number of iterations (e.g., $T = 6$) to converge, since long paths make negligible contribution to the task.

B.3 Traditional Methods

Theorem 6 *Katz index, personalized PageRank, graph distance, widest path and most reliable path can be solved via the generalized Bellman-Ford algorithm.*

Proof. Given that the generalized Bellman-Ford algorithm solves the path formulation when $\langle \oplus, \otimes \rangle$ satisfy a semiring, we only need to show that the operators of the path formulations for traditional methods satisfy semiring structures.

Katz index (Theorem 1) and personalized PageRank (Theorem 2) use the natural summation $+$ and the natural multiplication \times , which obviously satisfy a semiring.

Graph distance (Theorem 3) uses \min for *summation* and $+$ for *multiplication*. The corresponding identities are $\textcircled{0} = +\infty$ and $\textcircled{1} = 0$. It is obvious that $+$ satisfies the associative property and has identity element 0.

- **Commutative Property.** $\min(a, b) = \min(b, a)$
- **Associative Property.** $\min(\min(a, b), c) = \min(a, \min(b, c))$
- **Identity Element.** $\min(a, +\infty) = a$
- **Absorption Property.** $a + \infty = \infty + a = +\infty$
- **Distributive Property (Left).** $a + \min(b, c) = \min(a + b, a + c)$
- **Distributive Property (Right).** $\min(b, c) + a = \min(b + a, c + a)$

Widest path (Theorem 4) uses \max for *summation* and \min for *multiplication*. The corresponding identities are $\textcircled{0} = -\infty$ and $\textcircled{1} = +\infty$. We have

- **Commutative Property.** $\max(a, b) = \max(b, a)$
- **Associative Property.** $\max(\max(a, b), c) = \max(a, \max(b, c))$
- **Identity Element.** $\max(a, -\infty) = a$
- **Associative Property.** $\min(\min(a, b), c) = \min(a, \min(b, c))$
- **Absorption Property.** $\min(a, -\infty) = \min(-\infty, a) = -\infty$
- **Identity Element.** $\min(a, +\infty) = \min(+\infty, a) = a$

- **Distributive Property (Left).** $\min(a, \max(b, c)) = \max(\min(a, b), \min(a, c))$
- **Distributive Property (Right).** $\min(\max(b, c), a) = \max(\min(b, a), \min(c, a))$

where the distributive property can be proved by enumerating all possible orders of a , b and c .

Most reliable path (Theorem 5) uses \max for *summation* and \times for *multiplication*. The corresponding identities are $\textcircled{0} = 0$ and $\textcircled{1} = 1$, since all path representations are probabilities in $[0, 1]$. It is obvious that \times satisfies the associative property, the absorption property and has identity element 0.

- **Commutative Property.** $\max(a, b) = \max(b, a)$
- **Associative Property.** $\max(\max(a, b), c) = \max(a, \max(b, c))$
- **Identity Element.** $\max(a, 0) = a$
- **Distributive Property (Left).** $a \times \max(b, c) = \max(a \times b, a \times c)$
- **Distributive Property (Right).** $\max(b, c) \times a = \max(b \times a, c \times a)$

where the identity element and the distributive property hold for non-negative elements. \square

C Time Complexity of GNN Frameworks

Here we prove the time complexity for NBFNet and other GNN frameworks.

C.1 NBFNet

Lemma 2 *The time complexity of one NBFNet run (Algorithm 1) is $O(T(|\mathcal{E}|d + |\mathcal{V}|d^2))$.*

Proof. We break the time complexity by INDICATOR, MESSAGE and AGGREGATE functions.

INDICATOR is called $|\mathcal{V}|$ times, and a single call to INDICATOR takes $O(d)$ time. MESSAGE is called $T(|\mathcal{E}| + |\mathcal{V}|)$ times, and a single call to MESSAGE, i.e., a relation operator, takes $O(d)$ time. AGGREGATE is called $T|\mathcal{V}|$ times over a total of $T|\mathcal{E}|$ messages with d dimensions. Each call to AGGREGATE additionally takes $O(d^2)$ time due to the linear transformations in the function.

Therefore, the total complexity is summed to $O(T(|\mathcal{E}|d + |\mathcal{V}|d^2))$. \square

In practice, we find a small constant T works well for link prediction, and Lemma 2 can be reduced to $O(|\mathcal{E}|d + |\mathcal{V}|d^2)$ time.

Now consider applying NBFNet to infer the likelihood of all possible triplets. Without loss of generality, assume we want to predict the tail entity for each head entity and relation $p(v|u, q)$. We group triplets with the same condition u, q together, where each group contains $|\mathcal{V}|$ triplets. For triplets in a group, we only need to execute Algorithm 1 once to get their predictions. Therefore, the amortized time for a single triplet is $O\left(\frac{|\mathcal{E}|d}{|\mathcal{V}|} + d^2\right)$.

C.2 VGAE / RGCN

RGCN is a message-passing GNN applied to multi-relational graphs, with the message function being a per-relation linear transformation. VGAE can be viewed as a special case of RGCN applied to single-relational graphs. The time complexity of RGCN is similar to Lemma 2, except that each call to the message function takes $O(d^2)$ time due to the linear transformation. Therefore, the total complexity is $O(T(|\mathcal{E}|d^2 + |\mathcal{V}|d^2))$, where T refers to the number of layers in RGCN. Since $|\mathcal{V}| \leq |\mathcal{E}|$, the complexity is reduced to $O(T|\mathcal{E}|d^2)$ ¹¹. In practice, T is a small constant and we get $O(|\mathcal{E}|d^2)$ complexity.

While directly executing RGCN once for each triplet is costly, a smart way to apply RGCN for inference is to first compute all node representations, and then perform link prediction with the node representations. The first step runs RGCN once for $|\mathcal{V}|^2|\mathcal{R}|$ triplets, while the second step takes $O(d)$ time. Therefore, the amortized time for a single triplet is $O\left(\frac{|\mathcal{E}|d^2}{|\mathcal{V}|^2|\mathcal{R}|} + d\right)$. For large graphs and reasonable choices of d , we have $|\mathcal{E}|d \leq |\mathcal{V}|^2|\mathcal{R}|$ and the amortized time can be reduced to $O(d)$.

¹¹By moving the linear transformations from the message function to the aggregation function, one can also get an implementation of RGCN with $O(T|\mathcal{V}||\mathcal{R}|d^2)$ time, which is better for dense graphs but worse for sparse graphs. For knowledge graph datasets used in this paper, the above $O(T|\mathcal{E}|d^2)$ implementation is better.

C.3 NeuralLP / DRUM

DRUM can be viewed as a special case of NBFNet with MESSAGE being Hadamard product and AGGREGATE being natural summation. NeuralLP is a special case of DRUM where the dimension d equals to 1. Since there is no linear transformation in their AGGREGATE functions, the amortized time complexity for the message passing part is $O\left(\frac{T|\mathcal{E}|d}{|\mathcal{V}|}\right)$. Both DRUM and NeuralLP additionally use an LSTM to learn the edge weights for each layer, which additionally costs $O(Td^2)$ time for T layers. T is small and can be ignored like in other methods. Therefore, the amortized time of two parts is summed to $O\left(\frac{|\mathcal{E}|d}{|\mathcal{V}|} + d^2\right)$.

C.4 SEAL / GraIL

GraIL first extracts a local subgraph surrounding the link, and then applies RGCN to the local subgraph. SEAL can be viewed as a special case of GraIL applied to single-relational graphs. Therefore, their amortized time is the same as that of one RGCN run, which is $O(|\mathcal{E}|d^2)$.

Note that one may still run GraIL on large graphs by restricting the local subgraphs to be very small, e.g., within 1-hop neighborhood of the query entities. However, this will severely harm the performance of link prediction. Moreover, most real-world graphs are small-world networks, and a moderate radius can easily cover a non-trivial number of nodes and edges, which costs a lot of time for GraIL.

D Number of Parameters

Table 8: Number of parameters in NBFNet. The number of parameters only grows with the number of relations $|\mathcal{R}|$, rather than the number of nodes $|\mathcal{V}|$ or edges $|\mathcal{E}|$. For FB15k-237 augmented with flipped triplets, $|\mathcal{R}|$ is 474. Our best configuration uses $T = 6$, $d = 32$ and hidden dimension $m = 64$.

	#Parameter	
	Analytic Formula	FB15k-237
INDICATOR	$ \mathcal{R} d$	15,168
MESSAGE	$T \mathcal{R} d(d+1)$	3,003,264
AGGREGATE	$Td(13d+3)$	80,448
$f(\cdot)$	$m(2d+1) + m + 1$	4,225
Total		3,103,105

One advantage of NBFNet is that it requires much less parameters than embedding methods. For example, on FB15k-237, NBFNet requires 3M parameters while TransE requires 30M parameters. Table 8 shows a break down of number of parameters in NBFNet. Generally, the number of parameters in NBFNet scales linearly w.r.t. the number of relations, regardless the number of entities in the graph, which makes NBFNet more parameter-efficient for large graphs.

E Statistics of Datasets

Dataset statistics of two transductive settings, i.e., knowledge graph completion and homogeneous graph link prediction, are summarized in Table 9 and 10. Dataset statistics of inductive relation prediction is summarized in Table 11.

We use the standard transductive splits [56, 13] and inductive splits [55] for knowledge graphs. For homogeneous graphs, we follow previous works [32, 12] and randomly split the edges into train/validation/test sets with a ratio of 85:5:10. All the homogeneous graphs used in this paper are undirected. Note that for inductive relation prediction, the original paper [55] actually uses a *transductive valid set* that shares the same set of fact triplets as the training set for hyperparameter tuning. The *inductive test set* contains entities, query triplets and fact triplets that never appear in the training set. The same data split is adopted in this paper for a fair comparison.

Table 9: Dataset statistics for knowledge graph completion.

Dataset	#Entity	#Relation	#Train	#Triplet #Validation	#Test
FB15k-237 [56]	14,541	237	272,115	17,535	20,466
WN18RR [13]	40,943	11	86,835	3,034	3,134

Table 10: Dataset statistics for homogeneous link prediction.

Dataset	#Node	#Train	#Edge #Validation	#Test
Cora [49]	2,708	4,614	271	544
CiteSeer [49]	3,327	4,022	236	474
PubMed [49]	19,717	37,687	2,216	4,435

Table 11: Dataset statistics for inductive relation prediction. Queries refer to the triplets that are used as training or test labels, while facts are the triplets used as training or test inputs. In the training sets, all queries are also provided as facts.

Dataset		#Relation	#Entity	Train #Query	#Fact	#Entity	Validation #Query	#Fact	#Entity	Test #Query	#Fact
FB15k-237 [55]	v1	180	1,594	4,245	4,245	1,594	489	4,245	1,093	205	1,993
	v2	200	2,608	9,739	9,739	2,608	1,166	9,739	1,660	478	4,145
	v3	215	3,668	17,986	17,986	3,668	2,194	17,986	2,501	865	7,406
	v4	219	4,707	27,203	27,203	4,707	3,352	27,203	3,051	1,424	11,714
WN18RR [55]	v1	9	2,746	5,410	5,410	2,746	630	5,410	922	188	1,618
	v2	10	6,954	15,262	15,262	6,954	1,838	15,262	2,757	441	4,011
	v3	11	12,078	25,901	25,901	12,078	3,097	25,901	5,084	605	6,327
	v4	9	3,861	7,940	7,940	3,861	934	7,940	7,084	1,429	12,334

F Implementation Details

Table 12: Hyperparameter configurations of NBFNet on different datasets. Adv. temperature corresponds to the temperature in self-adversarial negative sampling [52]. Note for FB15k-237 and WN18RR, we use the same hyperparameters for their transductive and inductive settings. We find our model configuration is robust across all datasets, therefore we only tune the learning hyperparameters for each dataset. All the hyperparameters are chosen by the performance on the validation set.

Hyperparameter		FB15k-237	WN18RR	Cora	CiteSeer	PubMed
GNN	#layer(T)	6	6	6	6	6
	hidden dim.	32	32	32	32	32
MLP	#layer	2	2	2	2	2
	hidden dim.	64	64	64	64	64
Batch	#positive	256	128	256	256	64
	#negative/#positive(n)	32	32	1	1	1
Learning	optimizer	Adam	Adam	Adam	Adam	Adam
	learning rate	5e-3	5e-3	5e-3	5e-3	5e-3
	#epoch	20	20	20	20	20
	adv. temperature	0.5	1	-	-	-

Our implementation generally follows the open source codebases of knowledge graph completion¹² and homogeneous graph link prediction¹³. Table 12 lists the hyperparameter configurations for different datasets. Table 13 shows the wall time of training and inference on different datasets.

Data Augmentation. For knowledge graphs, we follow previous works [69, 46] and augment each triplet $\langle u, q, v \rangle$ with a flipped triplet $\langle v, q^{-1}, u \rangle$. For homogeneous graphs, we follow previous works [33, 32] and augment each node u with a self loop $\langle u, u \rangle$.

Architecture Details. We apply Layer Normalization [2] and short cut connection to accelerate the training of NBFNet. Layer Normalization is applied after each AGGREGATE function. The feed-forward network $f(\cdot)$ is instantiated as a MLP. ReLU is used as the activation function for all hidden layers. For undirected graphs, we symmetrize the pair representation by taking the sum of $\mathbf{h}_q(u, v)$ and $\mathbf{h}_q(v, u)$.

¹²<https://github.com/DeepGraphLearning/KnowledgeGraphEmbedding>. MIT license.

¹³<https://github.com/tkipf/gae>. MIT license.

Training Details. We train NBFNet on 4 Tesla V100 GPUs with standard data parallelism. During training, we drop out edges that directly connect query node pairs to encourage the model to capture longer paths and prevent overfitting. We select the best checkpoint for each model based on its performance on the validation set. The selection criteria is MRR for knowledge graphs and AUROC for homogeneous graphs.

Fused Message Passing. To reduce memory footprint and better utilize GPU hardware, we follow the efficient implementation of GNNs [26] and implement customized PyTorch operators that combines MESSAGE and AGGREGATE functions into a single operation, without creating all messages explicitly. This reduces the memory complexity of NBFNet from $O(|\mathcal{E}|d)$ to $O(|\mathcal{V}|d)$.

Table 13: Wall time of NBFNet on different datasets and in different settings (Table 3, 4 and 5). For inductive setting, the total time over 4 split versions is reported.

Wall Time	Transductive					Inductive	
	FB15k-237	WN18RR	Cora	CiteSeer	PubMed	FB15k-237	WN18RR
Training	9.7 hrs	4.4 hrs	5.5 mins	5.3 mins	5.6 hrs	23 mins	41 mins
Inference	4.0 mins	2.4 mins	< 1 sec	< 1 sec	25 secs	6 secs	20 secs

G Experimental Results on OGB Datasets

To demonstrate the effectiveness of NBFNet on large-scale graphs, we additionally evaluate our method on two knowledge graph datasets from OGB [25], ogbl-biokg and WikiKG90M. We follow the standard evaluation protocol of OGB link property prediction, and compute the mean reciprocal rank (MRR) of the true entity against 1,000 negative entities.

G.1 Results on ogbl-biokg

Ogbl-biokg is a large biomedical knowledge graph that contains 93,773 entities, 51 relations and 5,088,434 triplets. We compare NBFNet with 6 embedding methods on this dataset. Note by the time of this work, only embedding methods are available for such large-scale datasets. Table 14 shows the results on ogbl-biokg. NBFNet achieves the best result compared to all methods reported on the official leaderboard¹⁴ with much fewer parameters. Note the previous best model AutoSF is based on architecture search and requires more computation resource than NBFNet for training.

Table 14: Knowledge graph completion results on ogbl-biokg. Results of compared methods are taken from the OGB leaderboard.

Class	Method	Test MRR	Validation MRR	#Params
Embeddings	TransE [6]	0.7452	0.7456	187,648,000
	DistMult [68]	0.8043	0.8055	187,648,000
	ComplEx [58]	0.8095	0.8105	187,648,000
	RotatE [52]	0.7989	0.7997	187,597,000
	AutoSF [75]	0.8309	0.8317	93,824,000
	PairRE [7]	0.8164	0.8172	187,750,000
GNNs	NBFNet	0.8317	0.8318	734,209

G.2 Results on WikiKG90M

WikiKG90M is an extremely large dataset used in OGB large-scale challenge [24], hold at KDD Cup 2021. It is a general-purpose knowledge graph containing 87,143,637 entities, 1,315 relations and 504,220,369 triplets.

To apply NBFNet to such a large scale, we use a bidirectional breath-first-search (BFS) algorithm to sample a local subgraph for each query. Given a query, we generate a k -hop neighborhood for each

¹⁴https://ogb.stanford.edu/docs/leader_linkprop/#ogbl-biokg

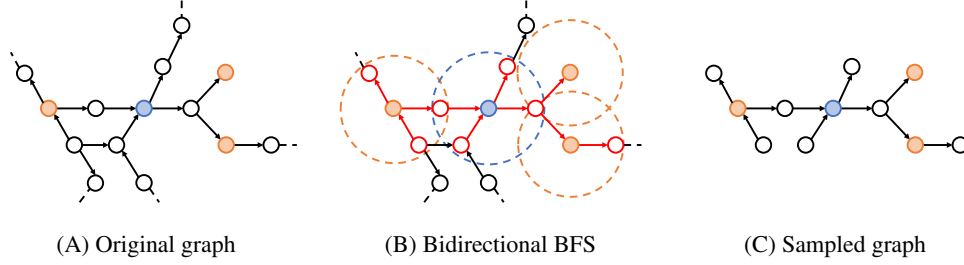


Figure 1: Illustration of bidirectional BFS sampling. For a **head entity** and multiple **tail candidates**, we use BFS to sample a k -hop neighborhood around each entity, regardless of the direction of edges. The neighborhood is denoted by dashed circles. The nodes and edges visited by the BFS algorithm are extracted to generate the sampled graph. Best viewed in color.

of the head entity and the candidate tail entities, based on a BFS search. The union of all generated neighborhoods is then collected as the sampled graph. With this sampling algorithm, any path within a length of $2k$ between the head entity and any tail candidate is guaranteed to be present in the sampled graph. See Figure 1 for illustration. While a standard single BFS algorithm computing the $2k$ -hop neighborhood of the head entity has the same guarantee, a bidirectional BFS algorithm significantly reduces the number of nodes and edges in the sampled graph.

We additionally downsample the neighbors when expanding the neighbors of an entity, to tackle entities with large degrees. For each entity visited during the BFS algorithm, we downsample its outgoing neighbors and incoming neighbors to m entities respectively.

Table 16 shows the results of NBFNet on WikiKG90M validation set. Our best single model uses $k = 2$ and $m = 100$. While the validation set requires to rank the true entity against 1,000 negative entities, in practice it is not mandatory to draw 1,000 negative samples for each positive sample during training. We find that reducing the negative samples from 1,000 to 20 and increasing the batch size from 4 to 64 provides a better result, although it creates a distribution shift between sampled graphs in training and validation. We leave further research of such distribution shift as a future work.

Table 15: Results of different MESSAGE and AGGREGATE functions on FB15k-237.

MESSAGE	AGGREGATE	MR	MRR	H@1	H@3	H@10
TransE [6]	Sum	191	0.297	0.217	0.321	0.453
	Mean	161	0.310	0.218	0.339	0.496
	Max	135	0.377	0.282	0.415	0.565
	PNA [9]	129	0.383	0.288	0.420	0.568
DistMult [68]	Sum	136	0.388	0.294	0.427	0.574
	Mean	132	0.384	0.287	0.425	0.577
	Max	136	0.374	0.279	0.412	0.563
	PNA [9]	114	0.415	0.321	0.454	0.599
RotatE [52]	Sum	129	0.392	0.298	0.429	0.580
	Mean	138	0.376	0.278	0.416	0.571
	Max	139	0.385	0.290	0.423	0.572
	PNA [9]	117	0.414	0.323	0.454	0.593

Table 16: Knowledge graph completion results on WikiKG90M validation set.

Model	Single Model	6 Model Ensemble
MRR	0.924	0.930

Table 17: Results of different number of layers on FB15k-237.

#Layers (T)	MR	MRR	H@1	H@3	H@10
2	191	0.345	0.261	0.377	0.510
4	119	0.409	0.315	0.450	0.592
6	114	0.415	0.321	0.454	0.599
8	115	0.416	0.322	0.457	0.599

H Ablation Study

Table 15 shows the full results of different MESSAGE and AGGREGATE functions. Table 17 shows the full results of NBFNet with different number of layers.