

---

# Inductive Logical Query Answering in Knowledge Graphs

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Formulating and answering logical queries is a standard communication interface  
2 for knowledge graphs (KGs) and their representations. Alleviating the notorious  
3 incompleteness of real-world KGs, neural methods achieved impressive results  
4 in link prediction and complex query answering tasks by learning representations  
5 of entities, relations, and queries. Still, most existing query answering methods  
6 are inherently transductive and cannot be generalized to KGs containing new en-  
7 tities without retraining entity embeddings. In this work, we study the inductive  
8 query answering task where inference is performed on a graph containing new  
9 entities with queries over both seen and unseen entities. To this end, we devise  
10 two mechanisms leveraging inductive *node* and *relational structure* representations  
11 powered by graph neural networks (GNNs). Experimentally, we show that induc-  
12 tive models are able to perform logical reasoning at inference time over unseen  
13 nodes generalizing to graphs up to 500% larger than training ones. Exploring the  
14 efficiency–effectiveness trade-off, we find the inductive *relational structure* method  
15 generally achieves higher performance, while the inductive *node representation*  
16 method is able to answer complex queries in the *inference-only* regime without any  
17 training on queries and scale to graphs of millions of nodes.

## 18 1 Introduction

19 Traditionally, querying knowledge graphs (KGs) is performed via databases using structured query  
20 languages like SPARQL. Databases can answer complex queries relatively fast under the assumption  
21 of *completeness*, i.e., there is no missing information in the graph. In practice, however, KGs are  
22 notoriously incomplete [29]. Embedding-based methods that learn vector representations of entities  
23 and relations are known to be effective in *simple link prediction* predicting heads or tails of query  
24 patterns (*head, relation, ?*), e.g., (*Einstein, graduate, ?*), forming the field of *KG completion* [1, 14].

25 Complex queries are graph patterns expressed in a subset of first-order logic (FOL) with operators  
26 such as intersection ( $\wedge$ ), union ( $\vee$ ), negation ( $\neg$ ) and existentially quantified ( $\exists$ ) variables<sup>1</sup>, e.g.,  
27  $?U.\exists V : \text{Win}(\text{NobelPrize}, V) \wedge \text{Citizen}(\text{USA}, V) \wedge \text{Graduate}(V, U)$  (Fig. 1). Complex queries  
28 define a superset of link prediction on KGs. The conventional link prediction task can be viewed as a  
29 complex query with a single triplet pattern without logic operators, e.g.,  $\text{Citizen}(\text{USA}, V)$ , which  
30 we also denote as a *projection* query.

31 To tackle complex queries on incomplete knowledge graphs, *query embedding* methods are proposed  
32 to execute logic operations in the latent space, including variants that employ geometric [12, 20, 35],  
33 probabilistic [21, 7], neural-symbolic [23, 6, 4], neural [18, 3], and GNN [9, 2] approaches for  
34 learning entity, relation, and query representations.

---

<sup>1</sup>The universal quantifier ( $\forall$ ) is often discarded as in real-world KGs there is no node connected to all others.

Where did US citizens with Nobel Prize graduate?  $q = v. \exists u. \text{Win}(\text{Nobel Prize}, u) \wedge \text{Citizen}(\text{USA}, u) \wedge \text{Graduate}(u, v)$

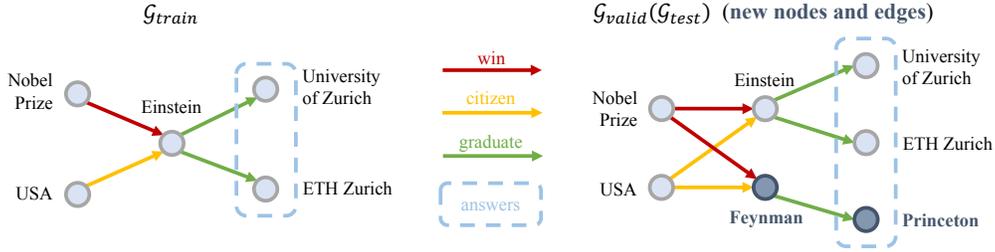


Figure 1: Inductive query answering problem: at inference time, the graph is updated with new nodes Feynman and Princeton and edges such that the same query now has more answers.

35 However, this very fact of learning a separate embedding for each entity makes those methods  
 36 inherently *transductive* i.e., they are bound to the space of learned entities and can not be generalized  
 37 to unseen entities without retraining the whole embedding matrix which can be prohibitively expensive  
 38 in large graphs. The problem is illustrated in Fig. 1: given a graph about Einstein and a logical  
 39 query *Where did US citizens with Nobel Prize graduate?*, transductive QE methods learn to execute  
 40 logical operators and return the answer set {University of Zurich, ETH Zurich}. Then, the  
 41 graph is updated with new nodes and edges about Feynman and Princeton, and the same query now  
 42 has more correct answers {University of Zurich, ETH Zurich, Princeton} as new unseen  
 43 entities satisfy the query as well.

44 Such *inductive inference* is not possible for transductive models as they do not have representations for  
 45 new Feynman and Princeton nodes. In the extreme case, inference graphs might be disconnected  
 46 from the training one and only share the set of relations. Therefore, inductive capabilities are a key  
 47 factor to enable transferring trained query answering models onto updated or entirely new KGs.

48 In this work, we study answering complex queries in the inductive setting, where the model has to deal  
 49 with unseen entities at inference time. Inspired by recent advancement in inductive link prediction on  
 50 KGs [36, 10], we devise two solutions for learning inductive representations for complex query: 1)  
 51 The first solution, NodePiece-QE, extends the inductive node representation model NodePiece [10]  
 52 for complex query answering. NodePiece-QE learns **inductive representations of each entity** as a  
 53 function of tokens from a fixed-size vocabulary, and answers complex query with a non-parametric  
 54 logical query executor [4]. The advantages of NodePiece-QE are that it only needs to be trained  
 55 on simple link prediction data, answers complex queries in the *inference-only* mode, and that it can  
 56 scale to large KGs. 2) The second solution, NBFNet-QE, extends the inductive link prediction model  
 57 NBFNet [36] for complex query answering. NBFNet-QE learns **inductive representations of the**  
 58 **relational structure** without entity embeddings, and uses the relational structure between the query  
 59 constants and the answers to make the prediction. NBFNet-QE can be trained end-to-end on complex  
 60 queries, achieves much better performance than NodePiece-QE, but struggles to scale to large KGs.

61 To the best of our knowledge, this is the first work to study complex logical query answering in the  
 62 inductive setting. Conducting experiments on a novel benchmarking suite of 10 datasets, we find that  
 63 (i) both inductive solutions exhibit non-trivial performance answering logical queries over unseen  
 64 entities and query patterns; (ii) inductive models demonstrate out-of-distribution generalization  
 65 capabilities to graphs up to 500% larger than training ones; (iii) akin to updatable databases, inductive  
 66 methods can successfully find new correct answers to known training queries after adding new  
 67 nodes and edges; (iv) the inductive *node representation* method scales to answering logical queries  
 68 over a graph of 2M nodes with 500k new, unseen nodes; (v) GNN-based models still exhibit  
 69 difficulties [17, 32] generalizing to larger graphs than they were originally trained on.

## 70 2 Related Work

71 **Knowledge Graph Completion.** Knowledge graph completion, a.k.a. simple link prediction, has  
 72 been widely studied in the *transductive* paradigm [5, 30, 24, 34], i.e., when inference is performed  
 73 on the same training graph with a fixed set of entities. Generally, these methods learn a shallow  
 74 embedding vector for each entity. We refer the audience to respective surveys [1, 14] covering

75 dozens of transductive embedding methods. The emergence of message passing [11] and Graph  
 76 Neural Networks (GNNs) has led to more advanced, *inductive* representation learning approaches  
 77 that model entity or triplet representations as a function of the graph structure in its neighborhood.  
 78 GraIL [25] learns triplet representations based on the subgraph structure surrounding the two entities.  
 79 NeuralLP [31], DRUM [22] and NBFNet [36] learn the pairwise entity representations based on  
 80 the set of relation paths between two entities. NodePiece [10] learns entity representations from a  
 81 fixed-size vocabulary of tokens that can be anchor nodes in a graph or relation types.

82 **Complex Query Answering.** In the complex (multi-hop) query answering setup with logical  
 83 operators, existing models employ different approaches, e.g., geometric [12, 20, 35], probabilistic [21,  
 84 7], neural-symbolic [23, 6, 4], neural [18, 3], and GNN [9, 2]. Still, all the approaches are created and  
 85 evaluated exclusively in the transductive mode where the set of entities does not change at inference  
 86 time. To the best of our knowledge, there is no related work in inductive logical query answering  
 87 when an inference graph contains new entities. With our work, we aim to bridge this gap and extend  
 88 inductive representation learning algorithms to logical query answering. In particular, we focus on  
 89 the inductive setup where an inference graph is a superset of a training graph<sup>2</sup> such that (i) inference  
 90 queries require reasoning over both seen and new entities; (ii) original training queries might have  
 91 more correct answers at inference time with the addition of new entities.

### 92 3 Preliminaries and Problem Definition

93 **Knowledge Graph and Inductive Setup.** Given a finite set of entities  $\mathcal{E}$ , a finite set of relations  $\mathcal{R}$ ,  
 94 and a set of triples (edges)  $\mathcal{T} = (\mathcal{E} \times \mathcal{R} \times \mathcal{E})$ , a knowledge graph  $\mathcal{G}$  is defined as  $\mathcal{G} = (\mathcal{E}, \mathcal{R}, \mathcal{T})$ .  
 95 Accounting for the inductive setup, we define a *training* graph  $\mathcal{G}_{train} = (\mathcal{E}_{train}, \mathcal{R}, \mathcal{T}_{train})$  and *inference*  
 96 graph  $\mathcal{G}_{inf} = (\mathcal{E}_{inf}, \mathcal{R}, \mathcal{T}_{inf})$  such that  $\mathcal{E}_{train} \subset \mathcal{E}_{inf}$  and  $\mathcal{T}_{train} \subset \mathcal{T}_{inf}$ . That is, the *inference* graph extends  
 97 the training graph with new entities and edges<sup>3</sup>. The inference graph  $\mathcal{G}_{inf}$  is an incomplete part of the  
 98 not observable complete graph  $\hat{\mathcal{G}}_{inf} = (\mathcal{E}_{inf}, \mathcal{R}, \hat{\mathcal{T}}_{inf})$  with  $\hat{\mathcal{T}}_{inf} = \mathcal{T}_{inf} \cup \mathcal{T}_{pred}$  whose missing triples  
 99  $\mathcal{T}_{pred}$  have to be predicted at inference time.

100 **First-Order Logic Queries.** Applied to KGs, a first-order logic (FOL) query  $\mathcal{Q}$  is a formula that  
 101 consists of constants  $\mathcal{C}$  ( $\mathcal{C} \subseteq \mathcal{E}$ ), variables  $\mathcal{V}$  ( $\mathcal{V} \subseteq \mathcal{E}$ , existentially quantified), relation *projections*  
 102  $R(a, b)$  denoting a binary function over constants or variables, and logic symbols ( $\exists, \wedge, \vee, \neg$ ). The  
 103 answers  $A_{\mathcal{G}}(\mathcal{Q})$  to the query  $\mathcal{Q}$  are assignments of variables in a formula such that the instantiated  
 104 query formula is a subgraph of the complete graph  $\hat{\mathcal{G}}$ .

105 Fig. 1 illustrates the logical form of a query *Where did US citizens with Nobel Prize graduate?* as  
 106  $?U.\exists V : \text{Win}(\text{NobelPrize}, V) \wedge \text{Citizen}(\text{USA}, V) \wedge \text{Graduate}(V, U)$  where *NobelPrize* and *USA*  
 107 are *constants*; *Win*, *Citizen*, *Graduate* are *relation projections* (labeled edges);  $V, U$  - *variables*  
 108 such that  $V$  is an existentially quantified free variable and  $U$  is the projected bound *target* of the  
 109 query. Common for the literature, we aim at predicting assignments of the query *target* whereas  
 110 assignments of intermediate variables might not always be explicitly interpreted depending on the  
 111 model architecture. In the example, the answer set  $A_{\mathcal{G}}(\mathcal{Q})$  is a binding of a target variable  $U$  to  
 112 constants *University of Zurich* and *ETH Zurich*.

113 **Inductive FOL Queries.** In the standard transductive query answering setup, query constants and  
 114 variables at both training and inference time belong to the same set of entities, i.e.,  $\mathcal{C}_{train} = \mathcal{C}_{inf} \subseteq$   
 115  $\mathcal{E}$ ,  $\mathcal{V}_{train} = \mathcal{V}_{inf} \subseteq \mathcal{E}$ . In the inductive setup covered in this work, query constants and variables at  
 116 inference time belong to a different and larger set of entities  $\mathcal{E}_{inf}$  from the inference graph  $\mathcal{G}_{inf}$ , i.e.,  
 117  $\mathcal{C}_{train} \subseteq \mathcal{E}_{train}$ ,  $\mathcal{V}_{train} \subseteq \mathcal{E}_{train}$  but  $\mathcal{C}_{inf} \subseteq \mathcal{E}_{inf}$ ,  $\mathcal{V}_{inf} \subseteq \mathcal{E}_{inf}$ . This also leads to the fact that training queries  
 118 executed over the inference graph might have more correct answers, i.e.,  $A_{\mathcal{G}_{train}}(\mathcal{Q}) \subseteq A_{\mathcal{G}_{inf}}(\mathcal{Q})$ . For  
 119 example (cf. Fig. 1), the inference graph is updated with new nodes *Feynman*, *Princeton* and their  
 120 new respective edges. The same query now has a larger set of intermediate variables satisfying the  
 121 formula (*Feynman*) and an additional correct answer *Princeton*. Therefore, inductive generalization  
 122 is essential for obtaining representations of such new nodes and enabling logical reasoning over both  
 123 seen and new nodes, i.e., finding more answers to known queries in larger graphs or answering new  
 124 queries with new constants. In the following section, we describe two approaches for achieving  
 125 inductive generalization with different parameterization strategies.

<sup>2</sup>The set of relation types is fixed.

<sup>3</sup>Note that the set of relation types  $\mathcal{R}$  remains the same.

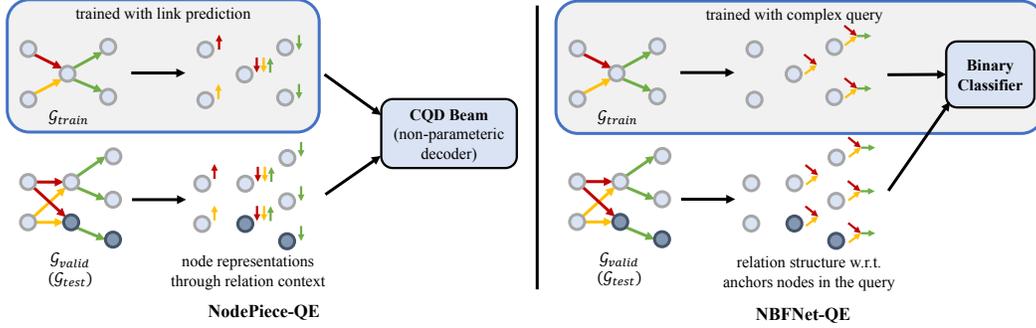


Figure 2: Inductive node representation (NodePiece-QE, left) and relational structure (NBFNet-QE, right) strategies for complex logical query answering. In NodePiece-QE, we obtain inductive node representations through the invariant set of tokens (here, through incident relation types). NodePiece-QE is the inference-only approach and is pre-trained with simple  $1p$  link prediction and can be directly applied to inductive complex queries with a non-parametric decoder (e.g., CQD Beam). In NBFNet-QE, we learn the the relative structure of each node w.r.t. the anchor nodes in the query. NBFNet-QE is trainable end-to-end with *complex queries*.

## 126 4 Method

127 **Inductive Representations of Complex Queries.** Given a complex query  $Q = (\mathcal{C}, \mathcal{R}_Q, \mathcal{G})$ , the  
 128 goal is to rank all possible entities according to the query. From a representation learning perspective,  
 129 this requires us to learn a conditional representation function  $f(e|\mathcal{C}, \mathcal{R}_Q, \mathcal{G})$  for each entity  $e \in \mathcal{E}$ .  
 130 Transductive methods learn a shallow embedding for each answer entity  $e \in \mathcal{E}$ , and, therefore, cannot  
 131 generalize to unseen entities. For inductive methods, the function  $f(e|\mathcal{C}, \mathcal{R}_Q, \mathcal{G})$  should generalize to  
 132 some unseen answer entity  $e'$  (or unseen constant entity  $c' \in \mathcal{C}'$ ) at inference time. Here, we discuss  
 133 two solutions for devising such an inductive function.

134 The first solution is to **parameterize the representation of each entity  $e$  as a function of an**  
 135 **invariant vocabulary** of *tokens* that does not change at training and inference. Particularly, the  
 136 vocabulary might consist of unique relation types  $\mathcal{R}$  that are always the same for  $\mathcal{G}_{train}$  and  $\mathcal{G}_{inf}$ , and  
 137 we are able to infer the representation of an unseen answer entity (or an unseen constant entity) as a  
 138 function of its incident relations (cf. Fig. 2 left). The idea has been studied in NodePiece [10] for  
 139 simple link prediction. Here, we adopt a similar idea to learn inductive entity representations for  
 140 complex query answering. Once we obtain the representations for unseen entities, we can use any  
 141 off-the-shelf decoding method (e.g., CQD-Beam [4]) for predicting the answer to the complex query.  
 142 We denote this strategy as NodePiece-QE.

143 The second solution is to **parameterize  $f(e|\mathcal{C}, \mathcal{R}_Q, \mathcal{G})$  as a function of the relational structure.**  
 144 Intuitively, an answer of a complex query can be decided solely based on the relational structure  
 145 between the query constants and the answer (Fig. 1). Even after anonymizing entity names (and,  
 146 hence, not learning any explicit entity embedding), we are still able to infer `Princeton` as an answer  
 147 since it forms a distinctive relational structure  $\rightarrow$  with the query constants and conforms to the query  
 148 structure. Similarly, intermediate nodes will be deemed correct if they follow a relational structure  
 149  $\rightarrow$ . In other words, we do not need to know the answer node is `Princeton`, but only need to know  
 150 the relative position of `Princeton` w.r.t. the constants like `Nobel Prize` and `USA`. Based on this  
 151 idea, we design  $f(e|\mathcal{C}, \mathcal{R}_Q, \mathcal{G})$  to be a relational structure search function. Such an idea has been  
 152 studied in Neural Bellman-Ford Networks (NBFNet) [36] to search for a single relation in simple  
 153 link prediction. Here, we chain several NBFNet instances with differentiable logic operations to learn  
 154 inductive complex query in an end-to-end fashion. We denote this strategy as NBFNet-QE.

### 155 4.1 NodePiece-QE: Inductive Node Representation

156 Here, we aim at reconstructing node representations for seen and new entities without learning shallow  
 157 node embedding vectors. To this end, we employ NodePiece [10], a compositional tokenization  
 158 approach that learns an invariant vocabulary of *tokens* shared between training and inference graphs.  
 159 Formally, given a vocabulary of tokens  $t_i \in T$ , each entity  $e_i$  is deterministically hashed into a set of

160 representative tokens  $e_i = [t_1, \dots, t_k]$ . An entity vector  $e_i$  is then obtained as a function of token  
 161 embeddings  $e_i = f_\theta([t_i, \dots, t_k])$ ,  $t_i \in \mathbf{T}^{|T| \times d}$  where the encoder function  $f_\theta : \mathbb{R}^{k \times d} \rightarrow \mathbb{R}^d$  is  
 162 parameterized with a neural network  $\theta$ .

163 Since the set of relation types  $\mathcal{R}$  is invariant for training and inference graphs, we can learn relation  
 164 embeddings  $\mathbf{R}^{|\mathcal{R}| \times d}$  and our vocabulary of learnable tokens  $T$  is comprised of distinct relation types  
 165 such that entities are hashed into a set of unique incident relation types. For example (cf. Fig. 2 left),  
 166 a middle node from a training graph  $\mathcal{G}_{train}$  is hashed with a set of relations  $e_i = [\uparrow\downarrow]$  that stands for  
 167 two unique incoming relations  $\downarrow$  and one unique outgoing relation  $\uparrow$ . Passing the hashes through  
 168  $f_\theta$ , we can reconstruct the whole entity embedding matrix  $\mathbf{E}^{|\mathcal{E}_{train}| \times d}$ . Additionally, it is possible to  
 169 enrich entity and relation embeddings by passing them through a relational GNN encoder [28] over a  
 170 target graph  $\mathcal{G}$ :  $\mathbf{E}', \mathbf{R}' = \text{GNN}(\mathbf{E}, \mathbf{R}, \mathcal{G})$ . In both ways, the entity embedding matrix  $\mathbf{E}$  encodes a  
 171 *joint* probability distribution  $p(h, r, t)$  for all triples in a graph.

172 Having a uniform featurization mechanism for both seen and unseen entities, it is now possible to  
 173 apply any previously-transductive complex query answering model with learnable entity embeddings  
 174 and logical operators [20, 9, 21, 6]. Moreover, it was recently shown [4] that a combination of  
 175 simple link prediction pre-training and a non-parametric logical executor allows to effectively answer  
 176 complex FOL queries in the *inference-only* regime without training on any complex query sample.  
 177 We adopt this Continuous Query Decomposition algorithm with beam search (CQD-Beam) as the  
 178 main query answering decoder. CQD-Beam relies only on entity and relation embeddings  $\mathbf{E}, \mathbf{R}$   
 179 pre-trained on a simple *1p* link prediction task. Then, given a complex query, CQD-Beam applies  
 180 *t-norms* and *t-conorms* [16] that execute conjunctions ( $\wedge$ ) and disjunctions ( $\vee$ ) as non-parametric  
 181 algebraic operations in the embedding space, respectively.

182 In our inductive setup (Fig. 2), we train a NodePiece encoder  $f_\theta$  and relation embeddings  $\mathbf{R}$  (and  
 183 optionally a GNN) on the *1p* link prediction task over the training graph  $\mathcal{G}_{train}$ . We then apply the  
 184 learned encoder to materialize entity representations of the inference graph  $\mathbf{E}^{|\mathcal{E}_{inf}| \times d}$  and send them  
 185 to CQD-Beam that performs a non-parametric decoding of complex FOL queries over new inference  
 186 entities. The inference-only nature of NodePiece-QE is designed to be challenging and probing the  
 187 abilities for zero-shot generalization in performing complex logical reasoning over larger graphs.

## 188 4.2 NBFNet-QE: Inductive Relational Structure Representation

189 The second strategy relies on learning inductive relational structure representations instead of explicit  
 190 node representations. Having the same set of relation types  $\mathcal{R}$  at training and inference time, we can  
 191 parameterize each entity based on the relative relational structure between it and the anchor nodes  
 192 in a given query. For instance (Fig. 2 right), given a query with a particular relational structure  $\rightarrow$   
 193 and a set of anchor nodes, the representation of each node captures its relational structure relative  
 194 to the anchor nodes. Each neighborhood expansion step is equivalent to the *projection* step. In our  
 195 example, immediate neighboring nodes will capture the intersection pattern  $\rightarrow$ , and further nodes, in  
 196 turn, capture the extended *intersection-projection* structure  $\rightarrow$ .

197 Therefore, a node is likely to be an answer if its captured (or predicted) relational structure conforms  
 198 with the query relational structure. As long as the set of relations is fixed, relation *projection* is  
 199 performed in the same way for training or new unseen nodes. The idea of a one-hop (*1p*) projection  
 200 for simple link prediction has been proposed by Neural Bellman-Ford Networks (NBFNet) [36].

201 In particular, given a relation *projection* query  $(h, r, ?)$ , NBFNet assigns unique initial states  $\mathbf{h}^{(0)}$  to  
 202 all nodes in a graph by applying an indicator function  $\mathbf{h}_e^{(0)} = \text{INDICATOR}(h, v, r)$ , i.e., a head node  $h$   
 203 is initialized with a learnable relation embedding  $\mathbf{r}$  and all other nodes are initialized with zeros. Then,  
 204 NBFNet applies  $L$  relational message passing GNN layers where each layer  $l$  has its own learnable  
 205 relation embedding matrix  $\mathbf{R}_l$  obtained as a projection of the initial relation  $\mathbf{R}_l = \mathbf{W}_l \mathbf{r} + \mathbf{b}_l$ . Final  
 206 layer representations  $\mathbf{h}^{(L)}$  are passed through an MLP and activation function  $\sigma$  to get a probability  
 207 distribution over all nodes in a graph  $p(t|h, r) = \sigma(\text{MLP}(\mathbf{h}^{(L)}))$ . As each query spawns a uniquely  
 208 initialized graph and message passing procedure, NBFNet is seen to be applying a *labeling trick* [33]  
 209 to model a conditional probability distribution  $p(t|h, r)$  which is provably more expressive than a  
 210 joint distribution  $p(h, r, t)$  produced by standard graph encoders.

211 Applied to complex queries, chaining  $k$  NBFNet instances allows to answer  $k$ -hop projection queries,  
 212 e.g., two instances for  $2p$  queries. NBFNet-QE employs NBFNet as a *trainable* projection operator

213 and endows it with differentiable, non-parametric *product* logic for modeling conjunction ( $\wedge$ ),  
 214 disjunction ( $\vee$ ), and negation ( $\neg$ ) over the *fuzzy sets* of all entities  $x \in [0, 1]^{\mathcal{E}}$ , i.e., after applying a  
 215 logical operator (discussed in Appendix A), each entity’s degree of truth is associated with a scalar  
 216 in range  $[0, 1]$ . For the next hop projection, the indicator function initializes a node state with a  
 217 relation vector  $r_i$  weighted by a scalar probability predicted in the previous hop  $x_e$ :  $h_e^{(0)} = x_e r_i$ .  
 218 Differentiable logical operators allow training NBFNet-QE end-to-end on complex queries.

## 219 5 Experiments

220 We designed the experimental agenda to demonstrate that inductive representation strategies are able  
 221 to: (1) answer complex logical queries over new, unseen entities at inference time, i.e., when query  
 222 anchors are new nodes (Section 5.2); (2) predict new correct answers for known *training* queries when  
 223 executed over larger inference graphs, i.e., when query anchors come from the training graph but  
 224 variables and answers belong to the larger inference graph (Section 5.3); (3) generalize to inference  
 225 graphs of up to 500% larger than training graphs; (4) scale to inductive query answering over graphs  
 226 of millions of nodes when updated with 500k new nodes and 5M new edges (Section 5.4).

### 227 5.1 Setup & Dataset

228 **Dataset.** Due to the absence of inductive logical query benchmarks, we create a novel suite of  
 229 datasets<sup>4</sup> based on FB15k-237 [26] (open license) and following the BetaE [21] query sampling  
 230 methodology. Given a source graph with  $\mathcal{E}$  entities, we sample  $|\mathcal{E}_{train}| = r \cdot |\mathcal{E}|$ ,  $r \in [0.1, 0.9]$  nodes  
 231 to induce a training graph  $\mathcal{G}_{train}$ . For validation and test graphs, we split the remaining set of entities  
 232 into two non-overlapping sets each with  $\frac{1-r}{2}|\mathcal{E}|$  nodes. We then merge training and unseen nodes  
 233 into the inference set of nodes  $\mathcal{E}_{inf}$  and induce inference graphs for validation and test from those sets,  
 234 respectively, i.e.,  $\mathcal{E}_{inf}^{val} = \mathcal{E}_{train} \cup \mathcal{E}_{val}$  and  $\mathcal{E}_{inf}^{test} = \mathcal{E}_{train} \cup \mathcal{E}_{test}$ . That is, validation and test inference  
 235 graphs both extend the training graph but their sets of new entities are disjoint. Finally, we sample and  
 236 remove 15% of edges  $\mathcal{T}_{pred}$  in the inference graphs as missing edges for link prediction pre-training  
 237 and query sampling. Overall, we sample 9 such datasets varying  $r$  to obtain ratios of inference graph  
 238 size to the training graph  $|\mathcal{E}_{inf}|/|\mathcal{E}_{train}|$  from 105% to 550%.

239 For each dataset, we employ the query sampler from BetaE [21] to extract 14 typical query types  
 240  $1p/2p/3p/2i/3i/ip/pi/2u/up/2in/3in/inp/pin/pni$ . Training queries are sampled from the training graph  
 241  $\mathcal{G}_{train}$ , validation and test queries are sampled from their respective inference graphs  $\mathcal{G}_{inf}$  where at  
 242 least one edge belongs to  $\mathcal{T}_{pred}$  and has to be predicted at inference time.

243 As inference graphs extend training graphs, training queries are very likely to have new answers being  
 244 executed over  $\mathcal{G}_{inf}$  with simple graph traversal and without any link prediction. We create an additional  
 245 set of true answers for all training queries executed over the test inference graph  $\mathcal{G}_{inf}^{test}$  to measure the  
 246 generalization capabilities of query answering models. This is designed to be an inference task and  
 247 extends the *faithfulness* experiment [23]. Dataset statistics can be found in Appendix B.

248 **Evaluation Protocol.** Following the literature [21], query answers are separated into two sets: *easy*  
 249 *answers* that only require graph traversal over existing edges, and *hard answers* that require inferring  
 250 missing links to achieve the answer node. For the main experiment, evaluation involves ranking of  
 251 *hard* answers against all entities having easy ones filtered out. For evaluating training queries on  
 252 inference graphs, we only have *easy* answers and rank them against all entities. We report Hits@10  
 253 as the main performance metric on different query types.

254 **Implementation Details.** All NodePiece [10]-based models were pre-trained until convergence on  
 255 a simple  $1p$  link prediction task with the relations-only vocabulary and entity tokenization, MLP  
 256 encoder, and ComplEx [27] scoring function. We used a 2-layer CompGCN [28] as an optional  
 257 message passing encoder on top of NodePiece features. The non-parametric CQD-Beam [4] de-  
 258 coder for answering complex queries is tuned for each query type based on the validation set of  
 259 queries, most of the setups employ a *product t-norm*, sigmoid entity score normalization, and beam  
 260 size of 32. Following the literature, the NBFNet-QE models were trained on 10 query patterns  
 261 ( $1p/2p/3p/2i/3i/2in/3in/inp/pin/pni$ ) where  $ip/pi/2u/up$  are only seen at inference time. Each model  
 262 employs a 4-layer NBFNet [36] as a trainable projection operator with DistMult [30] composi-

<sup>4</sup>Available to reviewers, will be published under the CC0 license

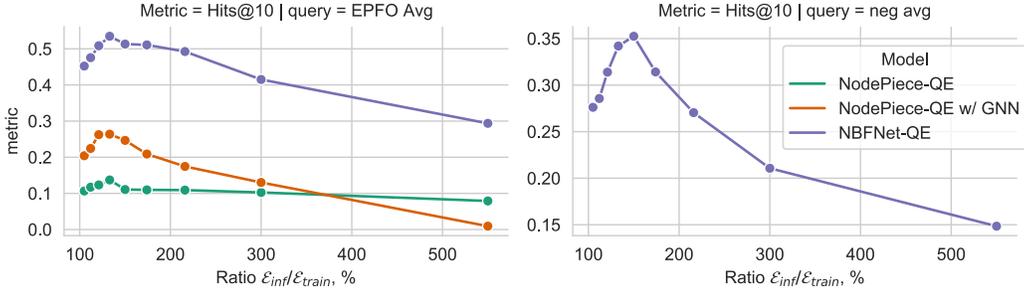


Figure 3: Aggregated Hits@10 performance of **test queries** (involving unseen entities) executed on inference graphs of different ratios compared to training graphs. NodePiece-based models are *inference-only* and support EPFO queries, NBFNet-QE is trainable and supports negation queries.

Table 1: Test Hits@10 results (%) on answering inductive FOL queries when  $\mathcal{E}_{inf}/\mathcal{E}_{train} = 175\%$ .  $avg_p$  is the average on EPFO queries ( $\wedge, \vee$ ).  $avg_n$  is the average on queries with negation.

Model	$avg_p$	$avg_n$	1p	2p	3p	2i	3i	pi	ip	2u	up	2in	3in	inp	pin	pni
<i>Inference-only</i>																
NodePiece-QE	11.0	-	21.3	8.9	5.1	13.0	14.7	9.8	8.7	9.7	7.5	-	-	-	-	-
NodePiece-QE w/ GNN	20.9	-	34.4	15.6	10.5	28.7	33.4	19.2	16.2	17.8	12.2	-	-	-	-	-
<i>Trainable</i>																
NBFNet-QE	51.1	31.4	66.1	40.9	31.2	73.0	83.3	58.3	41.3	37.8	27.8	31.1	44.3	28.4	25.2	28.0

263 tion function and PNA [8] aggregation. Other logical operators ( $\wedge, \vee, \neg$ ) are executed with the  
 264 non-parametric *product t-norm*. Both NodePiece-QE and NBFNet-QE are implemented<sup>5</sup> with Py-  
 265 Torch [19] and trained with the Adam [15] optimizer. NodePiece-QE models were pre-trained  
 266 and evaluated on a single Tesla V100 32 GB GPU whereas NBFNet-QE models were trained and  
 267 evaluated on 4 Tesla V100 16GB. All hyperparameters are listed in Appendix D.

## 268 5.2 Complex Query Answering over Unseen Entities on Differently Sized Inference Graphs

269 First, we probe *inference-only* NodePiece-based embedding models and *trainable* NBFNet-QE in  
 270 the inductive setup, i.e., query answering over unseen nodes requiring link prediction over unseen  
 271 nodes. Table 1 summarizes the results on a reference dataset with ratio  $\mathcal{E}_{inf}/\mathcal{E}_{train}$  of 175% while  
 272 Fig. 3 illustrates a bigger picture on all datasets (we provide a detailed breakdown by query type for  
 273 all splits in Appendix C). We observe that even inference-only models pre-trained solely on simple *1p*  
 274 link prediction exhibit non-trivial performance in answering queries with unseen entities. Paired with  
 275 an additional GNN encoder, the inference-only baseline exhibits significantly better performance  
 276 over all query types and inference graphs up to 300% larger than training graphs.

277 The trainable NBFNet-QE models expectedly outperform non-trainable baselines and can tackle  
 278 queries with negation ( $\neg$ ). Here, we confirm that the *labeling trick* [33] and conditional  $p(t|h, r)$   
 279 modeling better capture the relation projection problem than joint  $p(h, r, t)$  encoding approaches.

280 Still, all evaluated models with message passing, both inference-only NodePiece-QE with GNN  
 281 and trainable NBFNet-QE, suffer from increasing the size of the inference graph and having more  
 282 unseen entities. Reaching best results on  $\mathcal{E}_{inf}/\mathcal{E}_{train}$  ratios around 130%, both approaches steadily  
 283 deteriorate up until final 550% by 20 absolute Hits@10 points on EPFO queries and negation queries.  
 284 We attribute this deterioration to the known generalization issues [17, 32] of message passing GNNs  
 285 when performing inference over much larger graph than the network has seen during training. On  
 286 the other hand, a simple NodePiece-QE model without message passing retains similar performance  
 287 independently of the inference graph size.

288 Lastly, we observe that lower performance of inference-only NodePiece models can be also attributed  
 289 to underfitting (cf. train graph charts in Fig. 4). Although *1p* link predictors were trained until  
 290 convergence (on the inductive validation set of missing triples), the performance of training queries

<sup>5</sup>Source code is available to reviewers

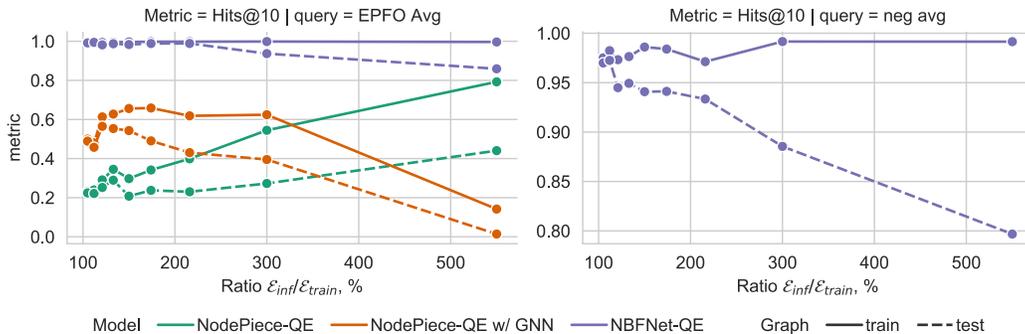


Figure 4: Aggregated Hits@10 performance of **training queries** on the original training and extended test inference graphs where queries have new correct answers. NodePiece-based models are *inference-only* and support EPFO queries, NBFNet-QE is trainable and supports negation queries.

291 on training graphs with *easy answers* that require only relation traversal without predicting missing  
 292 edges is not yet saturated. This fact suggests that better fitting entity featurization (obtained by  
 293 NodePiece or other strategies) could further improve the test performance in the inference-only  
 294 regime. We leave the search of such strategies for future work.

### 295 5.3 Predicting New Answers for Training Queries on Larger Inference Graphs

296 Simulating the incremental addition of new edges in graph databases, we evaluate the performance  
 297 of our inference-only and trainable QE models on *training* queries on the original training graph  
 298 and extended inference graph (with added test edges). As databases are able to immediately retrieve  
 299 new answers to known queries after updating the graph, we aim at exploring and quantifying this  
 300 behaviour of neural reasoning models. In this experiment, we probe training queries and their *easy*  
 301 *answers* that require performing only graph traversal without predicting missing links in the inference  
 302 graph. While execution of training queries over the *training* graph indicates how well the model  
 303 could fit training data, executing training queries over the bigger *inference* graph with new entities  
 304 aims to capture basic reasoning capabilities of QE models in the inductive regime.

305 Particular challenges arising when executing training queries over a bigger graph are: (1) the same  
 306 queries can have more correct answers as more new nodes and edges satisfying the query pattern  
 307 might have been added (as in Fig. 1); (2) more new entities create a “distractor” setting with more  
 308 false positives. Generally, evaluation of training queries on the inference graph can be considered  
 309 as an extended version of the *faithfulness* [23] evaluation that captures how well a trained model  
 310 can answer original training queries, i.e., memorization capacity. In all 9 datasets, most of training  
 311 queries have at least one new correct answer in the inference graph (more details in Appendix B).

312 Fig. 4 illustrates the performance of the inference-only NodePiece-QE (without and with GNN) and  
 313 trainable NBFNet-QE. Generally, NBFNet-QE fits the training query data almost perfectly confirming  
 314 the original finding [36] that NBFNet can perform graph traversal akin to symbolic rule-based models.  
 315 NBFNet-QE can also find new correct answers on graphs up to 300% larger than training ones. Then,  
 316 the performance quickly deteriorates which we attribute to the *distractor* factor with more unseen  
 317 entities and the already mentioned generalization issue on larger inference graphs.

318 The inference-only NodePiece-QE models, as expected, do not fully fit the training data as they were  
 319 never trained on complex queries. Still, the inference-only models exhibit non-trivial performance  
 320 in finding more answers on graphs up to 200% larger than training ones with relatively small  
 321 performance margins compared to training queries. The most surprising observation is that GNN-free  
 322 NodePiece-QE models improve the performance on both training and inference graphs as the graphs  
 323 (and the  $\mathcal{E}_{inf}/\mathcal{E}_{train}$  ratio) grow larger while GNN-enriched models steadily deteriorate. We attribute  
 324 this growth to the relation-based NodePiece tokenization and its learned features that tend to be more  
 325 discriminative in larger inference graphs where new nodes have smaller degree and thus can be better  
 326 identified by their incident relation types. We provide more experimental results for each dataset ratio  
 327 with breakdown by query type in Appendix C.

328 **5.4 Scaling to Millions of Nodes on WikiKG-QE**

329 Finally, we perform a scalability experiment evaluating complex query answering in the inductive  
 330 mode on a new large dataset *WikiKG-QE* constructed from OGB WikiKG 2 [13] (CC0 license).  
 331 While the original task is transductive link prediction, we split the graph into a training graph of  
 332 1.5M entities (5.8M edges, 512 unique relation types) and validation (test) graphs of 500k unseen  
 333 nodes (5M known and 600k missing edges) each. The resulting validation (test) inference graphs are  
 334 therefore of 2M entities and 11M edges with the  $\mathcal{E}_{inf}/\mathcal{E}_{train}$  ratio of 133% (details are in Appendix B).

335 None of GNN-enabled models can scale to such sizes, so we use a basic inference-only NodePiece-  
 336 QE. Due to the problem size, we only sample 10k EPFO queries of each type from the *test inference*  
 337 graph to run in the inference-only regime. Each query has at least one missing edge to be predicted at  
 338 inference. The answers are ranked against all 2M entities in the filtered setting (in contrast to the  
 339 OGB task that ranks against 1000 pre-computed negative samples) and Hits@100 as the target metric.

340 We pre-train a NodePiece encoder (in addition to relation types, we tokenize nodes with a vocabulary  
 341 of 20k anchor nodes, total 3M parameters in the encoder) with the ComplEx decoder on *Ip* link  
 342 prediction over the training graph for 1.5M steps (see Appendix D for hyperparameters). Then, the  
 343 graph is extended with 500k new nodes and 5M new edges forming the inference graph. Then, using  
 344 the pre-trained encoder, we materialize representations of entities (both seen and new) and relations  
 345 from this inference graph. Finally, CQD-Beam executes the queries against the bigger inference  
 346 graph extended with 500k new nodes and 5M new edges.

Table 2: Test Hits@100 of NodePiece-QE on WikiKG-QE (2M nodes, 11M edges including 500k new nodes and 5M new edges) in the inference-only regime.  $avg_p$  is the average on EPFO queries.

Model	$avg_p$	1p	2p	3p	2i	3i	pi	ip	2u	up
NodePiece-QE	6.6	19.0	2.1	1.8	11.0	15.8	4.4	2.8	1.3	1.5

347 As shown in Table 2, we find a non-trivial performance of the inference-only model on EPFO queries  
 348 demonstrating that inductive *node representation* QE models are able to scale to graphs with hundreds  
 349 of thousands of new nodes and millions of new edges in the zero-shot fashion. That is, answering  
 350 complex queries over unseen entities is available right upon updating the graph without the need  
 351 to retrain a model. This fact paves the way for the concept of *neural graph databases* capable of  
 352 performing zero-shot inference over updatable graphs without expensive retraining.

353 **6 Limitations and Future Work**

354 **Limitations.** With the two proposed inductive query answering strategies, we observe a common  
 355 trade-off between the performance and computational complexity. That is, inductive *node repre-*  
 356 *sentation* models like NodePiece-QE are fast, scalable, and can be executed in the inference-only  
 357 regime but underperform compared to the inductive *relational structure representation* models like  
 358 NBFNet-QE. On the other hand, NBFNet-QE incurs high computational costs due to executing each  
 359 query on a uniquely initialized graph instance. Alleviating this issue is a key to scalability.

360 **Societal Impact.** The inductive setup assumes running inference on (partly) unseen data, that is, the  
 361 nature of this unseen data might be out-of-distribution, unknown and potentially malicious. This fact  
 362 has to be taken into account when evaluating predictions and overall system trustworthiness.

363 **Conclusion and Future Work.** In this work, we defined the problem of inductive complex logical  
 364 query answering and proposed two possible parameterization strategies based on *node* and *relational*  
 365 *structure* representations to deal with new, unseen entities at inference time. Experiments demon-  
 366 strated that both strategies are able to answer complex logical queries over unseen entities as well as  
 367 identify new answers on larger inference graphs. In the future work, we plan to extend the inductive  
 368 setup to completely disjoint training and inference graphs, expand the set of supported logical query  
 369 patterns aligned with popular queries over real-world KGs, enable reasoning over continuous features  
 370 like texts and numbers, support more KG modalities like hypergraphs and hyper-relational graphs,  
 371 and further explore the concept of neural graph databases.

## References

- 372
- 373 [1] Mehdi Ali, Max Berrendorf, Charles Tapley Hoyt, Laurent Vermue, Mikhail Galkin, Sahand  
374 Sharifzadeh, Asja Fischer, Volker Tresp, and Jens Lehmann. Bringing light into the dark: A  
375 large-scale evaluation of knowledge graph embedding models under a unified framework. *IEEE*  
376 *Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- 377 [2] Dimitrios Alivanistos, Max Berrendorf, Michael Cochez, and Mikhail Galkin. Query embedding  
378 on hyper-relational knowledge graphs. In *International Conference on Learning Representations*,  
379 2022.
- 380 [3] Alfonso Amayuelas, Shuai Zhang, Xi Susie Rao, and Ce Zhang. Neural methods for logical  
381 reasoning over knowledge graphs. In *International Conference on Learning Representations*,  
382 2022.
- 383 [4] Erik Arakelyan, Daniel Daza, Pasquale Minervini, and Michael Cochez. Complex query  
384 answering with neural link predictors. In *International Conference on Learning Representations*,  
385 2021.
- 386 [5] Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, and Oksana Yakhnenko.  
387 Translating embeddings for modeling multi-relational data. In Christopher J. C. Burges, Léon  
388 Bottou, Zoubin Ghahramani, and Kilian Q. Weinberger, editors, *Advances in Neural Information*  
389 *Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems*  
390 *2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*,  
391 pages 2787–2795, 2013.
- 392 [6] Xuelu Chen, Ziniu Hu, and Yizhou Sun. Fuzzy logic based logical query answering on  
393 knowledge graph. In *International Conference on Machine Learning*. PMLR, 2021.
- 394 [7] Nurendra Choudhary, Nikhil Rao, Sumeet Katariya, Karthik Subbian, and Chandan K. Reddy.  
395 Probabilistic entity representation model for reasoning over knowledge graphs. In *Thirty-Fifth*  
396 *Conference on Neural Information Processing Systems*, 2021.
- 397 [8] Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal  
398 neighbourhood aggregation for graph nets. *Advances in Neural Information Processing Systems*,  
399 33:13260–13271, 2020.
- 400 [9] Daniel Daza and Michael Cochez. Message passing query embedding. *arXiv preprint*  
401 *arXiv:2002.02406*, 2020.
- 402 [10] Mikhail Galkin, Etienne Denis, Jiapeng Wu, and William L. Hamilton. Nodepiece: Composi-  
403 tional and parameter-efficient representations of large knowledge graphs. In *International*  
404 *Conference on Learning Representations*, 2022.
- 405 [11] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl.  
406 Neural message passing for quantum chemistry. In *Proceedings of the 34th International*  
407 *Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*,  
408 volume 70 of *Proceedings of Machine Learning Research*, pages 1263–1272. PMLR, 2017.
- 409 [12] Will Hamilton, Payal Bajaj, Marinka Zitnik, Dan Jurafsky, and Jure Leskovec. Embedding  
410 logical queries on knowledge graphs. *Advances in Neural Information Processing Systems*, 31,  
411 2018.
- 412 [13] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele  
413 Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs.  
414 In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-  
415 Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference*  
416 *on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*,  
417 2020.
- 418 [14] Shaoxiong Ji, Shirui Pan, Erik Cambria, Pekka Marttinen, and Philip S. Yu. A survey on  
419 knowledge graphs: Representation, acquisition and applications. *IEEE Transactions on Neural*  
420 *Networks and Learning Systems*, 33(2):494–514, 2022.

- 421 [15] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*  
422 (*Poster*), 2015.
- 423 [16] Erich-Peter Klement, Radko Mesiar, and Endre Pap. Triangular norms. position paper I: basic  
424 analytical and algebraic properties. *Fuzzy Sets Syst.*, 143(1):5–26, 2004.
- 425 [17] Boris Knyazev, Graham W. Taylor, and Mohamed R. Amer. Understanding attention and  
426 generalization in graph neural networks. In *Advances in Neural Information Processing*  
427 *Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS*  
428 *2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 4204–4214, 2019.
- 429 [18] Bhushan Kotnis, Carolin Lawrence, and Mathias Niepert. Answering complex queries in  
430 knowledge graphs with bidirectional sequence encoders. *CoRR*, abs/2004.02596, 2020.
- 431 [19] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,  
432 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas  
433 Köpf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,  
434 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-  
435 performance deep learning library. In *Advances in Neural Information Processing Systems 32:*  
436 *Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December*  
437 *8-14, 2019, Vancouver, BC, Canada*, pages 8024–8035, 2019.
- 438 [20] Hongyu Ren, Weihua Hu, and Jure Leskovec. Query2box: Reasoning over knowledge graphs in  
439 vector space using box embeddings. In *International Conference on Learning Representations*,  
440 2019.
- 441 [21] Hongyu Ren and Jure Leskovec. Beta embeddings for multi-hop logical reasoning in knowledge  
442 graphs. *Advances in Neural Information Processing Systems*, 33, 2020.
- 443 [22] Ali Sadeghian, Mohammadreza Armandpour, Patrick Ding, and Daisy Zhe Wang. Drum:  
444 End-to-end differentiable rule mining on knowledge graphs. *Advances in Neural Information*  
445 *Processing Systems*, 32, 2019.
- 446 [23] Haitian Sun, Andrew O. Arnold, Tania Bedrax-Weiss, Fernando Pereira, and William W.  
447 Cohen. Faithful embeddings for knowledge base queries. In Hugo Larochelle, Marc’Aurelio  
448 Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural*  
449 *Information Processing Systems 33: Annual Conference on Neural Information Processing*  
450 *Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- 451 [24] Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. Rotate: Knowledge graph em-  
452 bedding by relational rotation in complex space. In *7th International Conference on Learning*  
453 *Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
- 454 [25] Komal K. Teru, Etienne Denis, and Will Hamilton. Inductive relation prediction by subgraph  
455 reasoning. In *Proceedings of the 37th International Conference on Machine Learning, ICML*  
456 *2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*,  
457 pages 9448–9457. PMLR, 2020.
- 458 [26] Kristina Toutanova and Danqi Chen. Observed versus latent features for knowledge base and  
459 text inference. In *Proceedings of the 3rd Workshop on Continuous Vector Space Models and*  
460 *their Compositionality*, pages 57–66, Beijing, China, July 2015. Association for Computational  
461 Linguistics.
- 462 [27] Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard.  
463 Complex embeddings for simple link prediction. In *International conference on machine*  
464 *learning*, pages 2071–2080. PMLR, 2016.
- 465 [28] Shikhar Vashishth, Soumya Sanyal, Vikram Nitin, and Partha Talukdar. Composition-based  
466 multi-relational graph convolutional networks. In *International Conference on Learning Repr-*  
467 *esentations*, 2020.

- 468 [29] Robert West, Evgeniy Gabrilovich, Kevin Murphy, Shaohua Sun, Rahul Gupta, and Dekang  
469 Lin. Knowledge base completion via search-based question answering. In Chin-Wan Chung,  
470 Andrei Z. Broder, Kyuseok Shim, and Torsten Suel, editors, *23rd International World Wide  
471 Web Conference, WWW '14, Seoul, Republic of Korea, April 7-11, 2014*, pages 515–526. ACM,  
472 2014.
- 473 [30] Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding entities and  
474 relations for learning and inference in knowledge bases. In *3rd International Conference on  
475 Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track  
476 Proceedings*, 2015.
- 477 [31] Fan Yang, Zhilin Yang, and William W. Cohen. Differentiable learning of logical rules for  
478 knowledge base reasoning. In *Advances in Neural Information Processing Systems 30: Annual  
479 Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach,  
480 CA, USA*, pages 2319–2328, 2017.
- 481 [32] Gilad Yehudai, Ethan Fetaya, Eli A. Meirom, Gal Chechik, and Haggai Maron. From local  
482 structures to size generalization in graph neural networks. In *Proceedings of the 38th Interna-  
483 tional Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume  
484 139 of *Proceedings of Machine Learning Research*, pages 11975–11986. PMLR, 2021.
- 485 [33] Muhan Zhang, Pan Li, Yinglong Xia, Kai Wang, and Long Jin. Labeling trick: A theory of using  
486 graph neural networks for multi-node representation learning. In M. Ranzato, A. Beygelzimer,  
487 Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information  
488 Processing Systems*, volume 34, pages 9061–9073. Curran Associates, Inc., 2021.
- 489 [34] Shuai Zhang, Yi Tay, Lina Yao, and Qi Liu. Quaternion knowledge graph embeddings. In  
490 Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox,  
491 and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual  
492 Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14,  
493 2019, Vancouver, BC, Canada*, pages 2731–2741, 2019.
- 494 [35] Zhanqiu Zhang, Jie Wang, Jiajun Chen, Shuiwang Ji, and Feng Wu. Cone: Cone embeddings  
495 for multi-hop reasoning over knowledge graphs. *Advances in Neural Information Processing  
496 Systems*, 34, 2021.
- 497 [36] Zhaocheng Zhu, Zuobai Zhang, Louis-Pascal Xhonneux, and Jian Tang. Neural bellman-ford  
498 networks: A general graph neural network framework for link prediction. *Advances in Neural  
499 Information Processing Systems*, 34, 2021.

500 **Checklist**

- 501 1. For all authors...
- 502 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
503 contributions and scope? [Yes]
- 504 (b) Did you describe the limitations of your work? [Yes] See Section 6
- 505 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See  
506 Section 6
- 507 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
508 them? [Yes]
- 509 2. If you are including theoretical results...
- 510 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 511 (b) Did you include complete proofs of all theoretical results? [N/A]
- 512 3. If you ran experiments...
- 513 (a) Did you include the code, data, and instructions needed to reproduce the main exper-  
514 imental results (either in the supplemental material or as a URL)? [Yes] Code and  
515 sample data are included in the supplementary material
- 516 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
517 were chosen)? [Yes] Dataset creation process is described in Section 5.1 with more  
518 details in Appendix B. Hyperparameters are specified in Appendix D.
- 519 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
520 ments multiple times)? [No] We observe negligible variance w.r.t. random seeds
- 521 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
522 of GPUs, internal cluster, or cloud provider)? [Yes] Training details are specified in  
523 Section 5.1 and in Appendix D.
- 524 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 525 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 526 (b) Did you mention the license of the assets? [Yes]
- 527 (c) Did you include any new assets either in the supplemental material or as a URL?  
528 [Yes] Due to the overall size, we include a sample of the benchmarking suite in the  
529 supplemental material and will openly publish the whole dataset.
- 530 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
531 using/curating? [N/A] No personal data involved
- 532 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
533 information or offensive content? [Yes] The datasets are anonymized, we discuss it in  
534 Appendix B.
- 535 5. If you used crowdsourcing or conducted research with human subjects...
- 536 (a) Did you include the full text of instructions given to participants and screenshots, if  
537 applicable? [N/A]
- 538 (b) Did you describe any potential participant risks, with links to Institutional Review  
539 Board (IRB) approvals, if applicable? [N/A]
- 540 (c) Did you include the estimated hourly wage paid to participants and the total amount  
541 spent on participant compensation? [N/A]