
EFO_k-CQA: Towards Knowledge Graph Complex Query Answering beyond Set Operation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 To answer complex queries on knowledge graphs, logical reasoning over incomplete
2 knowledge is required due to the open-world assumption. Learning-based methods
3 are essential because they are capable of generalizing over unobserved knowledge.
4 Therefore, an appropriate dataset is fundamental to both obtaining and evaluating
5 such methods under this paradigm. In this paper, we propose a comprehensive
6 framework for data generation, model training, and method evaluation that covers
7 the combinatorial space of Existential First-order Queries with multiple variables
8 (EFO_k). The combinatorial query space in our framework significantly extends
9 those defined by set operations in the existing literature. Additionally, we construct
10 a dataset, EFO_k-CQA, with 741 query types for empirical evaluation, and our
11 benchmark results provide new insights into how query hardness affects the results.
12 Furthermore, we demonstrate that the existing dataset construction process is
13 systematically biased that hinders the appropriate development of query-answering
14 methods, highlighting the importance of our work. Our code and data are provided
15 in <https://anonymous.4open.science/r/EFOK-CQA/README.md>.

16 1 Introduction

17 The Knowledge Graph (KG) is a powerful database that encodes relational knowledge into a graph
18 representation [34, 31], supporting downstream tasks [41, 8] with essential factual knowledge. How-
19 ever, KGs suffer from incompleteness during its construction [34, 7], which is formally acknowledged
20 as Open World Assumption (OWA) [19]. The task of Complex Query Answering (CQA) proposed
21 recently has attracted much research interest [13, 28]. This task ambitiously aims to answer database-
22 level complex queries described by logical complex connectives (conjunction \wedge , disjunction \vee ,
23 and negation \neg) and quantifiers¹ (existential \exists) [37, 27, 18]. However, CQA on KGs differs from
24 query answering on databases in two aspects: (1) traditional query answering algorithms obtain
25 incomplete answers because of the incomplete KG [13]; (2) the huge size of the knowledge graph
26 limits the scalability of traditional algorithms [26]. Therefore, learning-based methods dominate
27 the CQA tasks because they can empirically generalize to unseen knowledge as well as prevent the
28 resource-demanding symbolic search.

29 The thriving of learning-based methods also puts an urgent request on high-quality datasets and
30 benchmarks. In the previous study, datasets are developed by progressively expanding the **syntactical**

¹The universal quantifier is usually not considered in query answering tasks, as a common practice from both CQA on KG [37, 27] and database query answering [25]

31 **expressiveness**, where conjunction [13], union [26], negation [28], and other operators [20] are taken
 32 into account sequentially. In particular, the dataset proposed in [28] contains all logical connectives
 33 and becomes the standard training set for model development. [36] proposed a large evaluation
 34 benchmark EFO-1-QA that systematically evaluates the combinatorial generalizability of CQA
 35 models on such queries. More related works are included in Appendix A.

36 However, the queries in aforementioned datasets [28, 36] are recently justified as “Tree-Form”
 37 queries [39] as they rely on the tree combinations of set operations. Compared to the well-established
 38 TPC-H decision support benchmark [25] for database query processing, queries in existing CQA
 39 benchmarks [28, 36] have two common shortcomings: (1) lack of **combinatorial answers**: only
 40 one variable is queried, and (2) lack of **structural hardness**: all existing queries subject to the
 41 structure-based tractability [29, 39]. It is rather questionable whether existing CQA data under such
 42 limited scope can support the future development of methodologies for general decision support with
 43 open-world knowledge.

44 The goal of this paper is to establish a new framework that addresses the aforementioned shortcomings
 45 to support further research in complex query answering on knowledge graphs. Our framework is
 46 formally motivated by the well-established investigation of constraint satisfaction problems, which
 47 all queries can be formulated as. In general, the contribution of our work is four folds.

48 **Complete coverage** We capture the complete Existential First Order (EFO) queries from their
 49 rigorous definitions, underscoring both **combinatorial hardness** and **structural hardness**
 50 and extending the existing coverage [36] which covers only a subset of EFO₁ query. The
 51 captured query family is denoted as EFO_k where *k* stands for multiple variables.

52 **Curated datasets** We derive EFO_k-CQA dataset, a non-exclusive extension of the previous EFO-1-
 53 QA benchmark [36] and contains 741 types of query. We design several rules to guarantee
 54 that our dataset includes high-quality nontrivial queries, particularly those that contain
 55 multiple query variables and are not structure-based tractable.

56 **Convenient implementation** We implement the entire pipeline for query generation, answer sam-
 57 pling, model training and inference, and evaluation for the undiscussed scenarios of **combi-
 58 natorial answers**. Our pipeline is backward compatible, which supports both set operation-
 59 based methods and more recent ones.

60 **Results and findings** We evaluate six representative CQA methods on our benchmark. Our results
 61 refresh the previous empirical findings and further reveal the structural bias of previous data.

62 2 Problem definition

63 2.1 Existential first order (EFO) queries on knowledge graphs

64 Given a set \mathcal{E} of entities and a set \mathcal{R} of relations, a knowledge graph \mathcal{KG} encodes knowledge as set
 65 of factual triple $\mathcal{KG} = \{(h, r, t)\} \subset \mathcal{E} \times \mathcal{R} \times \mathcal{E}$. According to the OWA, the knowledge graph that
 66 we have observed \mathcal{KG}_o is only part of the real knowledge graph, meaning that $\mathcal{KG}_o \subset \mathcal{KG}$.

67 The existing research only focuses on the logical formulas without universal quantifiers [27, 35]. We
 68 then offer the definition of it based on strict first order logic.

69 **Definition 1** (Term). *A term is either a variable x or an entity $a \in \mathcal{E}$.*

70 **Definition 2** (Atomic formula). *ϕ is an atomic formula if $\phi = r(h, t)$, where $r \in \mathcal{R}$ is a relation, h
 71 and t are two terms.*

72 **Definition 3** (Existential first order formula). *The set of the existential formulas is the smallest set Φ
 73 that satisfies the following:*

- 74 (i) For atomic formula $r(h, t)$, itself and its negation $r(h, t)$, $\neg r(h, t) \in \Phi$
- 75 (ii) If $\phi, \psi \in \Phi$, then $(\phi \wedge \psi), (\phi \vee \psi) \in \Phi$
- 76 (iii) If $\phi \in \Phi$ and x_i is any variable, then $\exists x_i \phi \in \Phi$.

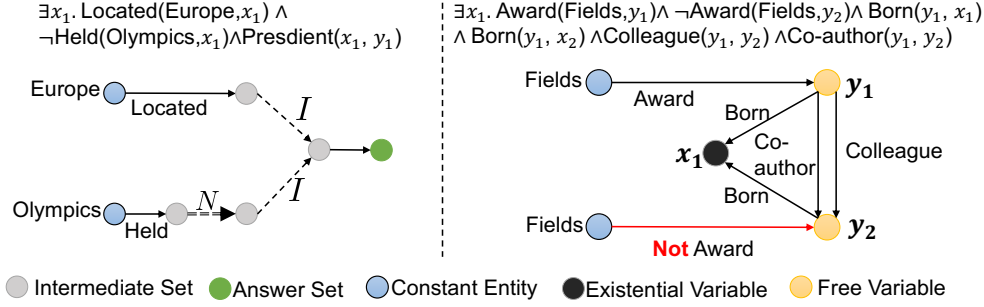


Figure 1: Operator Tree versus Query Graph. **Left:** An operator tree representing a given query “List the presidents of European countries that have never held the Olympics” [28]; **Right:** A query graph representing a given query “Find a pair of persons who are both colleagues and co-authors and were born in the same country, with one having awarded the fields medal while the another not”, which is both a multigraph and a cyclic graph, containing two free variables.

77 **Definition 4** (Free variable). *If a variable y is not associated with a quantifier, it is called a free*
 78 *variable, otherwise, it is called a bounded variable. We write $\phi(y_1, \dots, y_k)$ to indicate y_1, \dots, y_k*
 79 *are the free variables of ϕ .*

80 **Definition 5** (Sentence and query). *A formula ϕ is a sentence if it contains no free variable, otherwise,*
 81 *it is called a query. In this paper, we always consider formula with free variable, thus, we use formula*
 82 *and query interchangeably.*

83 **Definition 6** (Substitution). *For a_1, \dots, a_k , where $a_i \in \mathcal{E}$, we write $\phi(a_1/y_1, \dots, a_k/y_k)$ or simply*
 84 *$\phi(a_1, \dots, a_k)$ for the result of simultaneously replacing all free occurrence of y_i in ϕ by a_i , $i =$*
 85 *$1, \dots, k$.*

86 **Definition 7** (Answer of an EFO query). *For a given existential query $\phi(y_1, \dots, y_k)$, its answer is a*
 87 *set that defined by*

$$\mathcal{A}[\phi(y_1, \dots, y_k)] = \{(a_1, \dots, a_k) \mid a_i \in \mathcal{E}, i = 1, \dots, k, \phi(a_1, \dots, a_k) \text{ is True}\}$$

88 **Definition 8** (Disjunctive Normal Form (DNF)). *For any existential formula $\phi(y_1, \dots, y_k)$, it can*
 89 *be converted to the Disjunctive normal form as shown below:*

$$\phi(y_1, \dots, y_k) = \gamma_1(y_1, \dots, y_k) \vee \dots \vee \gamma_m(y_1, \dots, y_k) \quad (1)$$

$$\gamma_i(y_1, \dots, y_k) = \exists x_1, \dots, x_n. \rho_{i1} \wedge \dots \wedge \rho_{in} \quad (2)$$

90 *where ρ_{ij} is either an atomic formula or the negation of an atomic formula, x_i is called an existential*
 91 *variable.*

92 DNF form has a strong property that $\mathcal{A}[\phi(y_1, \dots, y_k)] = \cup_{i=1}^m \mathcal{A}[\gamma_i(y_1, \dots, y_k)]$, which allows
 93 us to only consider conjunctive formulas γ_i and then aggregate those answers to retrieve the final
 94 answers. This practical technique has been used in many previous research [22, 27]. Therefore, we
 95 only discuss conjunctive formulas in the rest of this paper.

96 2.2 Constraint satisfaction problem for EFO queries

97 Formally, a constraint satisfaction problem (CSP) \mathcal{P} can be represented by a triple $\mathcal{P} = (X, D, C)$
 98 where $X = (x_1, \dots, x_n)$ is an n -tuple of variables, $D = (D_1, \dots, D_n)$ is the corresponding n -tuple
 99 of domains, $C = (C_1, \dots, C_t)$ is t -tuple constraint, each constraint C_i is a pair of (S_i, R_{S_i}) where
 100 S_i is a set of variables $S_i = \{x_{i_j}\}$ and R_{S_i} is the constraint over those variables [29].

101 Historically, there are strong parallels between CSP and conjunctive queries in knowledge bases [10,
 102 17]. The terms correspond to the variable set X . The domain D_i of a constant entity contains only
 103 itself, while it is the whole entity set \mathcal{E} for other variables. Each constraint C_i is binary that is induced
 104 by an atomic formula or its negation, for example, for an atomic formula $r(h, t)$, we have $S_i = \{h, t\}$,
 105 $R_{S_i} = \{(h, t) \mid h, t \in \mathcal{E}, (h, r, t) \in \mathcal{KG}\}$. Finally, by the definition of existential quantifier, we only
 106 consider the answer of free variable, rather than tracking all terms within the existential formulas.

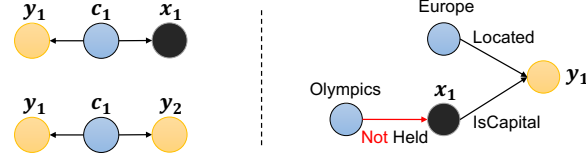


Figure 2: Left: Example of trivial abstract query graph, in the upper left graph, the x_1 is redundant violating Assumption 13, in the bottom left graph, answers for the whole query can be decomposed to answer two free variables y_1 and y_2 alone, violating Assumption 14. Right: Example of new query graph that is not included in previous benchmark [36] even though it can be represented by operator-tree. The representation of query graph follows Figure 1.

107 **Definition 9** (CSP answer of conjunctive formula). *For a conjunctive formula γ in Equation 2 with k*
 108 *free variables and n existential variables, the answer set of it formulated as CSP instance is:*

$$\overline{\mathcal{A}}[\gamma(y_1, \dots, y_k)] = \mathcal{A}[\gamma^*(y_1, \dots, y_{n+k})], \text{ where } \gamma^* = \rho_{i1} \wedge \dots \wedge \rho_{it}$$

109 This shows that the inference of existential formulas is easier than solving CSP instances since the
 110 existential variables do not need to be kept track of.

111 2.3 The representation of query

112 To give an explicit representation of existential formula, [13] firstly proposes to represent a formula
 113 by operator tree, where each node represents the answer set for a sub-query, and the logic operators in
 114 it naturally represent set operations. This method allows for the recursive computation from constant
 115 entity to the final answer set in a bottom-up manner [28]. However, this representation method is
 116 inherently directed, acyclic, and simple, therefore more recent research breaks these constraints by
 117 being bidirectional [21, 37] or being cyclic or multi [39]. To meet these new requirements, they
 118 propose to represent the formula by the query graph [39], which inherits the convention of constraint
 119 network in representing CSP instance. We utilize this design and further extend it to represent EFO_k
 120 formula that contains multiple free variables. We provide the illustration and comparison of the
 121 operator tree and the query graph in Figure 1, where we show the strong expressiveness of the query
 122 graph. We also provide the formal definition of query graph as follows:

123 **Definition 10** (Query graph). *Let γ be a conjunctive formula in equation 2, its query graph is defined*
 124 *by $G(\gamma) = \{(h, r, t, \{T, F\})\}$, where an atomic formula $\rho = r(h, t)$ in γ corresponds to (h, r, t, T)*
 125 *and $\rho = \neg r(h, t)$ corresponds to (h, r, t, F) .*

126 Therefore, any conjunctive formulas can be represented by a query graph, in the rest of the paper, we
 127 use query graphs and conjunctive formulas interchangeably.

128 3 The combinatorial space of EFO_k queries

129 Although previous research has given a systematic investigation in the combinatorial space of operator
 130 trees [36], the combinatorial space of the query graph is much more challenging due to the extremely
 131 large search space and the lack of explicit recursive formulation. To tackle this issue on a strong
 132 theoretical background, we put forward additional assumptions to exclude trivial query graphs. Such
 133 assumptions or restrictions also exist in the previous dataset and benchmark [28, 36]. Specifically,
 134 we propose to split the task of generating data into two levels, the abstract level, and the grounded
 135 level. At the abstract level, we create *abstract query graph*, at the grounded level, we provide the
 136 abstract query graph with the relation and constant and instantiate it as a query graph. In this section,
 137 we elaborate on how we investigate the scope of the nontrivial EFO_k query of interest step by step.

138 3.1 Nontrivial abstract query graph of EFO_k

139 The abstract query graph is the ungrounded query graph without information of certain knowledge
 140 graphs, and we give an example in Figure 3.

141 **Definition 11** (Abstract query graph). *The abstract query graph $\mathcal{G} = (V, E, f, g)$ is a directed*
 142 *graph with three node types, $\{\mathbf{Constant Entity}, \mathbf{Existential Variable}, \mathbf{Free variable}\}$, and two edge*
 143 *types, $\{\mathbf{positive}, \mathbf{negative}\}$. The V is the set of nodes, E is the set of directed edges, f is the function*
 144 *maps node to node type, g is the function maps edge to edge type.*

145 **Definition 12** (Grounding). *For an abstract query graph \mathcal{G} , a grounding is a function I that maps it*
 146 *into a query graph $I(\mathcal{G})$.*

147 We propose two assumptions of the abstract query graph as follows:

148 **Assumption 13** (No redundancy). *For a abstract query graph \mathcal{G} , there is not a subgraph $\mathcal{G}_s \subsetneq \mathcal{G}$*
 149 *such that for every grounding I , $\mathcal{A}[I(\mathcal{G})] = \mathcal{A}[I(\mathcal{G}_s)]$.*

150 **Assumption 14** (No decomposition). *For an abstract query graph \mathcal{G} , there are no such two*
 151 *subgraphs $\mathcal{G}_1, \mathcal{G}_2$, satisfying that $\mathcal{G}_1, \mathcal{G}_2 \subsetneq \mathcal{G}$, such that for every instantiation I , $\mathcal{A}[I(\mathcal{G})] =$*
 152 *$\mathcal{A}[I(\mathcal{G}_1)] \times \mathcal{A}[I(\mathcal{G}_2)]$, where the \times represents the Cartesian product.*

153 We note that the assumption 14 inherits the idea of the **structural** decomposition technique in
 154 CSP [11], which allows for solving a CSP instance by solving several sub-problems and combining
 155 the answer together based on topology property. Additionally, meeting these two assumptions in the
 156 grounded query graph is extremely computationally costly which we aim to avoid in practice.

157 We provide some easy examples to be excluded for violating the assumptions above in Figure 2.

158 3.2 Nontrivial query graph of \mathbf{EFO}_k

159 Similarly, we propose two assumptions on the query graph.

160 **Assumption 15** (Meaningful negation). *For any negative edge e in query graph G , we require*
 161 *removing it results in different CSP answers: $\overline{\mathcal{A}}[G - e] \neq \overline{\mathcal{A}}[G]$.²*

162 Assumption 15 treats negation separately because of the fact that for any \mathcal{KG} , any relation $r \in \mathcal{R}$,
 163 there is $|\{(h, t) | h, t \in \mathcal{E}, (h, r, t) \in \mathcal{KG}\}| \ll \mathcal{E}^2$, which means that the constraint induced by the
 164 negation of an atomic formula is much less “strict” than the one induced by a positive atomic formula.

165 **Assumption 16** (Appropriate answer size). *There is a constant $M \ll \mathcal{E}$ to bound the candidate set*
 166 *for each free variable f_i in G , such that for any i , $|\{(a_i \in \mathcal{E} | (a_1, \dots, a_i, \dots, a_k) \in \mathcal{A}[G])\}| \leq M$.*

167 We note the Assumption 16 **extends** the “bounded negation” assumption in the previous dataset [28,
 168 36]. We give an example “Find a city that is located in Europe and is the capital of a country that has
 169 not held the Olympics” in Figure 2, where the candidate set of x_1 is in fact bounded by its relation
 170 with the y_1 variable but not from the bottom “Olympics” constant, hence, this query is excluded in
 171 their dataset due to the directionality of operator tree.

172 Overall, the scope of the formula investigated in this paper surpasses the previous \mathbf{EFO} -1-QA
 173 benchmark because of: (1). We include the \mathbf{EFO}_k formula with multiple free variables for the first
 174 time; (2). We include the whole family of \mathbf{EFO}_1 query, many of them can not be represented by
 175 operator tree; (3) Our assumption is more systematic than previous ones as shown by the example in
 176 Figure 2. More details are offered in Appendix D.3.

177 4 Framework

178 We develop a versatile framework that supports five key functionalities fundamental to the whole
 179 CQA task: (1) Enumeration of nontrivial abstract query graphs as discussed in Section 3; (2) Sample
 180 grounding for the abstract query graph; (3) Compute answer for any query graph efficiently; (4)
 181 Support implementation of existing CQA models; (5) Conduct evaluation including newly introduced
 182 \mathbf{EFO}_k queries with multiple free variables. We explain each functionality in the following. An
 183 illustration of the first three functionalities is given in Figure 3.

²Ideally, we should expect them to have different answers as the existential formulas, however, this is computation costly and difficult to sample in practice, which is further discussed in Appendix D.

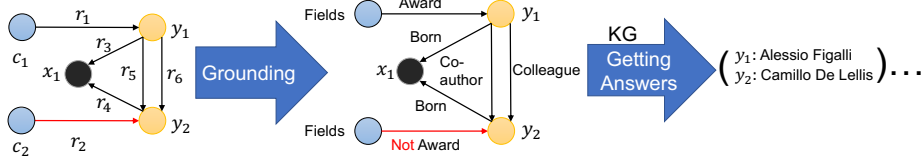


Figure 3: Illustration of the functionality of our framework. Left: abstract query graph, Middle: query graph, Right: answer of query.

184 4.1 Enumerate abstract query graph

185 As discussed in Section 3, we are able to abide by those assumptions as well as **enumerate** all
 186 possible query graphs within a given search space where certain parameters, including the number of
 187 constants, free variables, existential variables, and the number of edges are all given. Additionally,
 188 we apply the graph isomorphism algorithm to avoid duplicated query graphs being generated. More
 189 details for our generation method are provided in Appendix D.1.

190 4.2 Ground abstract query graph

191 To ground an abstract query graph \mathcal{G} and comply with the assumption 15, we split the abstract query
 192 graph into two parts, the positive part and the negative part, $\mathcal{G} = \mathcal{G}_p \cup \mathcal{G}_n$. Then the grounding
 193 process is also split into two steps: 1. Sample grounding for the positive subgraph \mathcal{G}_p and compute
 194 its answer 2. Ground the \mathcal{G}_n to decrease the answer got in the first step. Details in Appendix D.2.

195 Finally, to fulfill the assumption 16, we follow the previous practice of manually filtering out queries
 196 that have more than 100 answers [28, 36], as we have introduced the EFO_k queries, we slightly
 197 soften this constraint to be no more than $100 \times k$ answers.

198 4.3 Answer for existential formula

199 As illustrated in Section 2.2, the answer to an existential formula can be solved by a CSP solver,
 200 however, we also show in Definition 9 that CSP requires keeping track of the existential variables and
 201 it leads to huge computation costs. Thus, we develop our own algorithm following the standard solving
 202 technique of CSP, which ensures consistency conditions in the first step, and do the backtracking to get
 203 the final answers in the second step. Finally, we select part of our sampled queries and double-check
 204 it with the CSP solver <https://github.com/python-constraint/python-constraint>.

205 4.4 Learning-based methods

206 As the query graph is an extension to the operator tree regarding the express ability to existential
 207 formulas, we are able to reproduce CQA models that are initially implemented by the operator tree
 208 in our new framework. Specifically, since the operator tree is directed and acyclic, we compute its
 209 topology ordering that allows for step-by-step computation in the query graph. This algorithm is
 210 illustrated in detail in the Appendix F. We note our implementation coincides with the original one.

211 Conversely, for the newly proposed models that are based on query graphs, the original operator
 212 tree framework is not able to implement them, while our framework is powerful enough. We have
 213 therefore clearly shown that the query graph representation is more powerful than the previous
 214 operator tree and is able to support arbitrary existential formulas as explained in Section 2.3.

215 4.5 Evaluation protocol

216 As we have mentioned in Section 2.1, there is an observed knowledge graph \mathcal{KG}_o and a full knowledge
 217 graph \mathcal{KG} . Thus, there is a set of observed answers \mathcal{A}_o and a set of full answers \mathcal{A} correspondingly.
 218 Since the goal of CQA is to tackle the challenge of OWA, it has been a common practice to evaluate
 219 CQA models by the “hard” answers $\mathcal{A}_h = \mathcal{A} - \mathcal{A}_o$ [26, 27]. However, to the best of our knowledge,

220 there has not been a systematic evaluation protocol for EFO_k queries, thus we leverage this idea and
221 propose three types of different metrics to fill the research gap in the area of evaluation of queries
222 with multiple free variables, and thus have combinatorial answers.

223 **Marginal.** For any free variable f_i , its full answer is $\mathcal{A}^{f_i} = \{a_i \in \mathcal{E} | (a_1, \dots, a_i, \dots, a_k) \in \mathcal{A}\}$, the
224 observed answer of it $\mathcal{A}_o^{f_i}$ is defined similarly. This is termed “solution projection” in CSP theory [12]
225 to evaluate whether the locally retrieved answer can be extended to an answer for the whole problem.
226 Then, we rank the hard answer $\mathcal{A}_h^{f_i} = \mathcal{A}^{f_i} - \mathcal{A}_o^{f_i}$ ³, against those non-answers $\mathcal{E} - \mathcal{A}^{f_i} - \mathcal{A}_o^{f_i}$ and
227 use the ranking to compute standard metrics like MRR, HIT@K for every free variable. Finally, the
228 metric on the whole query graph is taken as the average of the metric on all free variables. We note
229 that this metric is an extension of the previous design proposed by [20]. However, this metric has the
230 inherent drawback that it fails to evaluate the combinatorial answer by the k -length tuple and thus
231 fails to find the correspondence among free variables.

232 **Multiply.** Because of the limitation of the marginal metric discussed above, we propose to evaluate
233 the combinatorial answer by each k -length tuple (a_1, \dots, a_k) in the hard answer set \mathcal{A}_h . Specifically,
234 we rank each a_i in the corresponding node f_i the same as the marginal metric. Then, we propose the
235 $HIT@n^k$ metric, it is 1 if all a_i is ranked in the top n in the corresponding node f_i , and 0 otherwise.

236 **Joint.** Finally, we note these metrics above are not the standard way of evaluation, which is based on
237 a joint ranking for all the \mathcal{E}^k combinations of the entire search space. We propose to estimate the
238 joint ranking in a closed form given certain assumptions, see Appendix E for the proof and details.

239 5 The EFO_k -CQA dataset and benchmark results

240 5.1 The EFO_k -CQA dataset

241 With the help of our framework developed in Section 4, we are able to develop a new dataset called
242 EFO_k -CQA, whose combinatorial space is parameterized by the number of constants, existential and
243 free variables, and the number of edges. EFO_k -CQA dataset includes 741 different abstract query
244 graphs in total. The parameters and the generation process, as well as its statistics, are detailed in
245 Appendix D.4.

246 Then, we conduct experiments on our new EFO_k -CQA dataset with six representative CQA models
247 including BetaE [28], LogicE [24], and ConE [40], which are built on the operator tree, CQD [2],
248 LMPNN [35], and FIT [39] which are built on query graph. The experiments are conducted in two
249 parts, (1). the queries with one free variable, specifically, including those that can not be represented
250 by operator tree; (2). the queries that contain multiple free variables.

251 We have made some adaptations to the implementation of CQA models, allowing them to infer EFO_k
252 queries, full detail is offered in Appendix F. The experiment is conducted on a standard knowledge
253 graph FB15k-237 [32] and additional experiments on other standard knowledge graphs FB15k and
254 NELL are presented in Appendix H.

255 5.2 Benchmark results for $k = 1$

256 Because of the great number of abstract query graphs, we follow [36] to group query graphs by three
257 factors: (1). the number of constant entities; (2). the number of existential variables, and (3). the
258 topology of the query graph⁴. The result is shown in Table 1.

259 **Structure analysis.** Firstly, we find a clear monotonic trend that adding constant entities makes a
260 query easier while adding existing variables makes a query harder, which the previous research [36]
261 fails to uncover. Besides, we are the first to consider the topology of query graphs: when the number

³We note $\mathcal{A}_h^{f_i}$ can be empty for some free variable or even for all free variables, making these marginal metrics not reliable, details in Appendix E.

⁴We make a further constraint in our EFO_k -CQA dataset that the total edge is at most as many as the number of nodes, thus, a graph can not be both a multigraph and a cyclic graph.

Table 1: HIT@10 scores(%) for inferring queries with one free variable on FB15k-237. We denote e as the number of existential variables and c as the number of constant entities. SDAG represents the Simple Directed Acyclic Graph, Multi for multigraph, and Cyclic for the cyclic graph. AVG.(c) and AVG.(e) is the average score of queries with the number of constant entities / existential variables fixed.

Model	$e \backslash c$		0			1			2			AVG.(c)	AVG.
			SDAG	SDAG	Multi	SDAG	Multi	Cyclic	SDAG	Multi	Cyclic		
BetaE	1		31.4	33.0	22.3	21.1	17.7	30.7	22.1			36.4	
	2		57.2	36.2	35.5	29.3	29.4	45.3	32.5				
	3		80.0	53.1	53.6	38.2	37.8	58.2	42.1				
	AVG.(e)		59.3	43.8	40.6	33.8	32.7	49.3					
LogicE	1		34.4	34.9	23.0	21.4	17.4	30.3	22.4		36.7		
	2		60.0	38.4	36.8	29.8	29.3	45.3	33.0				
	3		83.0	55.5	55.5	38.5	37.8	57.8	42.4				
	AVG.(e)		62.2	46.0	42.0	34.2	32.6	49.1					
ConE	1		34.9	35.4	23.6	21.8	18.4	34.2	23.5		39.0		
	2		61.0	39.1	38.4	32.0	31.5	50.2	35.2				
	3		84.8	56.7	57.1	41.1	40.0	63.4	44.9				
	AVG.(e)		63.4	47.0	43.5	36.5	34.7	54.1					
CQD	1		39.0	34.2	17.6	17.4	12.7	28.7	18.7		35.9		
	2		50.7	33.8	33.6	28.4	28.4	45.7	31.4				
	3		58.4	49.6	52.4	39.3	39.1	60.4	42.6				
	AVG.(e)		50.7	41.4	38.4	33.8	32.4	50.2					
LMPNN	1		38.6	37.8	21.8	22.9	17.8	31.7	23.2		35.8		
	2		62.2	40.2	35.0	30.8	28.1	44.4	32.5				
	3		86.6	56.9	51.9	38.3	35.3	55.8	40.8				
	AVG.(e)		65.4	47.8	39.6	34.5	30.8	48.0					
FIT	1		38.7	42.7	32.5	26.1	22.5	41.5	28.8		47.0		
	2		65.5	47.7	48.2	39.7	40.1	56.5	43.4				
	3		84.2	63.9	63.5	50.5	50.4	63.5	53.6				
	AVG.(e)		65.8	54.7	51.5	44.9	43.7	57.5					

of constants and existential variables is fixed, we have found the originally investigated queries that correspond to Simple Directed Acyclic Graphs (SDAG) are generally easier than the multigraphs ones but harder than the cyclic graph ones. This is an intriguing result that greatly deviates from traditional CSP theory in close world which finds that the cyclic graph is NP-complete, while the acyclic graph is tractable [6]. Our conjecture for this intriguing result in the open world is that the cyclic graph contains one more constraint than SDAG that serves as a source of information for CQA models, while the multigraph tightens an existing constraint and thus makes the query harder.

Model analysis. For models that are built on operator tree, including BetaE, LogicE, and ConE, their relative performance is steady among all breakdowns and is consistent with their reported score in the original dataset [28], showing similar generalizability. However, for models that are built on query graphs, including CQD, LMPNN, and FIT, we have found that LMPNN performs generally better than CQD in SDAG, but falls behind CQD in multigraphs and cyclic graphs. We assume the reason behind this is that LMPNN requires training while CQD does not, however, the original dataset are **biased** which only considers SDAG, leading to the result that LMPNN doesn't generalize well to the unseen tasks with different topology property. We expect future CQA models may use our framework to address this issue of biased data and generalize better to more complex queries.

We note FIT is designed to infer all EFO₁ queries and is indeed able to outperform other models in almost all breakdowns, however, its performance comes with the price of computational cost, and

Table 2: HIT@10 scores(%) of three different types for answering queries with two free variables on FB15k-237. The constant number is fixed to be two. e is the number of existential variables. The SDAG, Multi, and Cyclic are the same as Table 1.

Model	HIT@10 Type	$e = 0$		$e = 1$			$e = 2$			AVG.
		SDAG	Multi	SDAG	Multi	Cyclic	SDAG	Multi	Cyclic	
BetaE	Marginal	54.5	50.2	49.5	46.0	58.8	37.2	35.5	58.3	43.8
	Multiply	27.3	22.4	22.3	16.9	26.2	16.9	13.9	25.7	18.3
	Joint	6.3	5.4	5.2	4.2	10.8	2.2	2.3	9.5	4.5
LogicE	Marginal	58.2	50.9	52.2	47.4	60.4	37.7	35.8	59.2	44.6
	Multiply	32.1	23.1	24.9	18.1	28.3	18.1	14.8	26.6	19.5
	Joint	6.8	6.0	6.1	4.5	12.3	2.5	2.7	10.3	5.1
ConE	Marginal	60.3	53.8	54.2	50.3	66.2	40.1	38.5	63.7	47.7
	Multiply	33.7	25.2	26.1	19.8	32.1	19.5	16.3	30.3	21.5
	Joint	6.7	6.4	6.2	4.8	12.6	2.6	2.7	10.9	5.3
CQD	Marginal	50.4	46.5	49.1	45.6	59.7	33.5	33.1	61.5	42.8
	Multiply	28.9	23.4	25.4	19.5	31.3	17.8	16.0	30.5	21.0
	Joint	8.0	8.0	7.4	6.0	13.9	3.6	3.9	12.0	6.4
LMPNN	Marginal	58.4	51.1	54.9	49.2	64.7	39.6	36.1	58.7	45.4
	Multiply	35.0	26.7	29.2	21.7	33.4	21.4	17.0	28.4	22.2
	Joint	7.6	7.5	7.1	5.3	12.9	2.8	2.9	9.5	5.2
FIT	Marginal	64.3	61.0	63.1	60.7	58.5	49.0	49.1	60.2	54.3
	Multiply	39.7	32.2	35.9	27.8	27.4	29.5	26.8	32.4	29.2
	Joint	7.4	9.0	7.8	6.5	10.1	3.7	4.6	10.6	6.4

280 face challenges in cyclic graph where it degenerates to enumeration: which we further explain in
 281 Appendix F.

282 5.3 Benchmark results for $k = 2$

283 As we have explained in Section 4.5, we propose three kinds of metrics, marginal ones, multiply
 284 ones, and joint ones, from easy to hard, to evaluate the performance of a model in the scenario of
 285 multiple variables. The evaluation result is shown in Table 2. As the effect of the number of constant
 286 variables is quite clear, we remove it and add the metrics based on HIT@10 as the new factor.

287 For the impact regarding the number of existential variables and the topology property of the query
 288 graph, we find the result is similar to Table 1, which may be explained by the fact that those models
 289 are all initially designed to infer queries with one free variable. For the three metrics we have
 290 proposed, we have identified a clear difficulty difference among them though they generally show
 291 similar trends. The scores of joint HIT@10 are pretty low, indicating the great hardness of answering
 292 queries with multiple variables. Moreover, we have found that FIT falls behind other models in some
 293 breakdowns which are mostly cyclic graphs, corroborating our discussion in Section 5.2.

294 6 Conclusion

295 In this paper, we make a thorough investigation of the family of EFO_k formulas based on strong
 296 theoretical background. We then present a new powerful framework that supports several functionali-
 297 ties essential to CQA task, with this help, we build the EFO_k -CQA dataset that greatly extends the
 298 previous dataset and benchmark. Our evaluation result brings new empirical findings and reflects the
 299 biased selection in the previous dataset impairs the performance of CQA models, emphasizing the
 300 contribution of our work.

References

- [1] Dimitrios Alivanistos, Max Berrendorf, Michael Cochez, and Mikhail Galkin. Query Embedding on Hyper-relational Knowledge Graphs, September 2022. arXiv:2106.08166 [cs].
- [2] Erik Arakelyan, Daniel Daza, Pasquale Minervini, and Michael Cochez. Complex Query Answering with Neural Link Predictors. In *International Conference on Learning Representations*, 2020.
- [3] Jiaxin Bai, Zihao Wang, Hongming Zhang, and Yangqiu Song. Query2Particles: Knowledge Graph Reasoning with Particle Embeddings. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 2703–2714, 2022.
- [4] Yushi Bai, Xin Lv, Juanzi Li, and Lei Hou. Answering Complex Logical Queries on Knowledge Graphs via Query Computation Tree Optimization. In *Proceedings of the 40th International Conference on Machine Learning*, pages 1472–1491. PMLR, July 2023. ISSN: 2640-3498.
- [5] 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*, volume 26. Curran Associates, Inc., 2013.
- [6] Clément Carbonnel and Martin C Cooper. Tractability in constraint satisfaction problems: a survey. *Constraints*, 21(2):115–144, 2016. Publisher: Springer.
- [7] Andrew Carlson, Justin Betteridge, Bryan Kisiel, Burr Settles, Estevam Hruschka, and Tom Mitchell. Toward an architecture for never-ending language learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 24, pages 1306–1313, 2010. Issue: 1.
- [8] Lisa Ehrlinger and Wolfram Wöb. Towards a definition of knowledge graphs. *SEMANTiCS (Posters, Demos, SuCCESS)*, 48(1-4):2, 2016.
- [9] Michael Galkin, Zhaocheng Zhu, Hongyu Ren, and Jian Tang. Inductive logical query answering in knowledge graphs. *Advances in Neural Information Processing Systems*, 35:15230–15243, 2022.
- [10] Georg Gottlob, Nicola Leone, and Francesco Scarcello. Hypertree decompositions and tractable queries. In *Proceedings of the eighteenth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 21–32, 1999.
- [11] Georg Gottlob, Nicola Leone, and Francesco Scarcello. A comparison of structural CSP decomposition methods. *Artificial Intelligence*, 124(2):243–282, December 2000.
- [12] Gianluigi Greco and Francesco Scarcello. On The Power of Tree Projections: Structural Tractability of Enumerating CSP Solutions. *Constraints*, 18(1):38–74, January 2013. arXiv:1005.1567 [cs].
- [13] Will Hamilton, Payal Bajaj, Marinka Zitnik, Dan Jurafsky, and Jure Leskovec. Embedding logical queries on knowledge graphs. *Advances in neural information processing systems*, 31, 2018.
- [14] Zhiwei Hu, Víctor Gutiérrez-Basulto, Zhiliang Xiang, Xiaoli Li, and Jeff Pan. *Type-aware Embeddings for Multi-Hop Reasoning over Knowledge Graphs*. May 2022.
- [15] Qian Huang, Hongyu Ren, and Jure Leskovec. Few-shot relational reasoning via connection subgraph pretraining. *Advances in Neural Information Processing Systems*, 35:6397–6409, 2022.
- [16] Zhen Jia, Soumajit Pramanik, Rishiraj Saha Roy, and Gerhard Weikum. Complex Temporal Question Answering on Knowledge Graphs. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management, CIKM ’21*, pages 792–802, New York, NY, USA, 2021. Association for Computing Machinery.

- 346 [17] Phokion G Kolaitis and Moshe Y Vardi. Conjunctive-query containment and constraint satisfac-
347 tion. In *Proceedings of the seventeenth ACM SIGACT-SIGMOD-SIGART symposium on*
348 *Principles of database systems*, pages 205–213, 1998.
- 349 [18] Jure Leskovec. Databases as Graphs: Predictive Queries for Declarative Machine Learning. In
350 *Proceedings of the 42nd ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database*
351 *Systems*, PODS '23, page 1, New York, NY, USA, 2023. Association for Computing Machinery.
352 event-place: Seattle, WA, USA.
- 353 [19] Leonid Libkin and Cristina Sirangelo. Open and Closed World Assumptions in Data Exchange.
354 *Description Logics*, 477, 2009.
- 355 [20] Lihui Liu, Boxin Du, Heng Ji, ChengXiang Zhai, and Hanghang Tong. Neural-Answering
356 Logical Queries on Knowledge Graphs. In *Proceedings of the 27th ACM SIGKDD Conference*
357 *on Knowledge Discovery & Data Mining*, pages 1087–1097, 2021.
- 358 [21] Xiao Liu, Shiyu Zhao, Kai Su, Yukuo Cen, Jiezhong Qiu, Mengdi Zhang, Wei Wu, Yuxiao
359 Dong, and Jie Tang. Mask and Reason: Pre-Training Knowledge Graph Transformers for
360 Complex Logical Queries. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge*
361 *Discovery and Data Mining*, pages 1120–1130, August 2022. arXiv:2208.07638 [cs].
- 362 [22] Xiao Long, Liansheng Zhuang, Li Aodi, Shafei Wang, and Houqiang Li. Neural-based Mixture
363 Probabilistic Query Embedding for Answering FOL queries on Knowledge Graphs. 2022.
- 364 [23] Haoran Luo, Yuhao Yang, Gengxian Zhou, Yikai Guo, Tianyu Yao, Zichen Tang, Xueyuan Lin,
365 Kaiyang Wan, and others. NQE: N-ary Query Embedding for Complex Query Answering over
366 Hyper-relational Knowledge Graphs. *arXiv preprint arXiv:2211.13469*, 2022.
- 367 [24] Francois Luus, Prithviraj Sen, Pavan Kapanipathi, Ryan Riegel, Ndivhuwo Makondo, Thabang
368 Lebeso, and Alexander Gray. Logic embeddings for complex query answering. *arXiv preprint*
369 *arXiv:2103.00418*, 2021.
- 370 [25] Meikel Poess and Chris Floyd. New TPC benchmarks for decision support and web commerce.
371 *ACM Sigmod Record*, 29(4):64–71, 2000. Publisher: ACM New York, NY, USA.
- 372 [26] H Ren, W Hu, and J Leskovec. Query2box: Reasoning Over Knowledge Graphs In Vector Space
373 Using Box Embeddings. In *International Conference on Learning Representations (ICLR)*,
374 2020.
- 375 [27] Hongyu Ren, Mikhail Galkin, Michael Cochez, Zhaocheng Zhu, and Jure Leskovec. Neural
376 Graph Reasoning: Complex Logical Query Answering Meets Graph Databases, March 2023.
377 arXiv:2303.14617 [cs].
- 378 [28] Hongyu Ren and Jure Leskovec. Beta embeddings for multi-hop logical reasoning in knowledge
379 graphs. *Advances in Neural Information Processing Systems*, 33:19716–19726, 2020.
- 380 [29] Francesca Rossi, Peter van Beek, and Toby Walsh. *Handbook of Constraint Programming*.
381 Elsevier Science Inc., USA, 2006.
- 382 [30] Apoorv Saxena, Soumen Chakrabarti, and Partha Talukdar. Question Answering Over Temporal
383 Knowledge Graphs, June 2021. arXiv:2106.01515 [cs].
- 384 [31] Fabian M Suchanek, Gjergji Kasneci, and Gerhard Weikum. Yago: a core of semantic knowl-
385 edge. In *Proceedings of the 16th international conference on World Wide Web*, pages 697–706,
386 2007.
- 387 [32] Kristina Toutanova and Danqi Chen. Observed versus latent features for knowledge base and
388 text inference. In *Proceedings of the 3rd workshop on continuous vector space models and their*
389 *compositionality*, pages 57–66, 2015.

- 390 [33] Kush R. Varshney. Trustworthy machine learning and artificial intelligence. *XRDS: Crossroads*,
391 *The ACM Magazine for Students*, 25(3):26–29, 2019.
- 392 [34] Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Com-*
393 *munications of the ACM*, 57(10):78–85, 2014. Publisher: ACM New York, NY, USA.
- 394 [35] Zihao Wang, Yangqiu Song, Ginny Wong, and Simon See. Logical Message Passing Networks
395 with One-hop Inference on Atomic Formulas. In *The Eleventh International Conference on*
396 *Learning Representations*, 2023.
- 397 [36] Zihao Wang, Hang Yin, and Yangqiu Song. Benchmarking the Combinatorial Generalizability
398 of Complex Query Answering on Knowledge Graphs. *Proceedings of the Neural Information*
399 *Processing Systems Track on Datasets and Benchmarks*, 1, December 2021.
- 400 [37] Zihao Wang, Hang Yin, and Yangqiu Song. Logical Queries on Knowledge Graphs: Emerging
401 Interface of Incomplete Relational Data. *Data Engineering*, page 3, 2022.
- 402 [38] Zezhong Xu, Wen Zhang, Peng Ye, Hui Chen, and Huajun Chen. Neural-Symbolic Entangled
403 Framework for Complex Query Answering, September 2022. arXiv:2209.08779 [cs].
- 404 [39] Hang Yin, Zihao Wang, and Yangqiu Song. On Existential First Order Queries Inference on
405 Knowledge Graphs, April 2023. arXiv:2304.07063 [cs].
- 406 [40] Zhanqiu Zhang, Jie Wang, Jiajun Chen, Shuiwang Ji, and Feng Wu. Cone: Cone embeddings
407 for multi-hop reasoning over knowledge graphs. *Advances in Neural Information Processing*
408 *Systems*, 34:19172–19183, 2021.
- 409 [41] Tao Zhou, Jie Ren, Matúš Medo, and Yi-Cheng Zhang. Bipartite network projection and
410 personal recommendation. *Physical review E*, 76(4):046115, 2007. Publisher: APS.

411 Checklist

- 412 1. For all authors...
- 413 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
414 contributions and scope? [Yes]
- 415 (b) Did you describe the limitations of your work? [Yes] We can not handle queries with
416 the universal quantifier.
- 417 (c) Did you discuss any potential negative societal impacts of your work? [No] We believe
418 there is no negative social impact.
- 419 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
420 them? [Yes]
- 421 2. If you are including theoretical results...
- 422 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Clear assump-
423 tions are made in Section 3 to define the scope of the query we investigate.
- 424 (b) Did you include complete proofs of all theoretical results? [Yes]
- 425 3. If you ran experiments (e.g. for benchmarks)...
- 426 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
427 mental results (either in the supplemental material or as a URL)? [Yes] We have given
428 the link in the abstract.
- 429 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
430 were chosen)? [Yes] This is in Appendix F.
- 431 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
432 ments multiple times)? [No] However, we have evaluated CQA models in the previous
433 dataset and the result is similar to the scores in original paper.

- 434 (d) Did you include the total amount of compute and the type of resources used (e.g., type
435 of GPUs, internal cluster, or cloud provider)? [Yes]
- 436 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 437 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 438 (b) Did you mention the license of the assets? [No] They are all open datasets.
- 439 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 440 (d) Did you discuss whether and how consent was obtained from people whose data you're
441 using/curating? [N/A]
- 442 (e) Did you discuss whether the data you are using/curating contains personally identifiable
443 information or offensive content? [N/A]
- 444 5. If you used crowdsourcing or conducted research with human subjects...
- 445 (a) Did you include the full text of instructions given to participants and screenshots, if
446 applicable? [N/A] We have not used crowdsourcing.
- 447 (b) Did you describe any potential participant risks, with links to Institutional Review
448 Board (IRB) approvals, if applicable? [N/A]
- 449 (c) Did you include the estimated hourly wage paid to participants and the total amount
450 spent on participant compensation? [N/A]