

A Reachability Index for Recursive Label-Concatenated Graph Queries

Chao Zhang*, Angela Bonifati*, Hugo Kapp†, Vlad Ioan Haprian† and Jean-Pierre Lozi†

*Lyon 1 University, Lyon, France

†Oracle Labs, Zürich, Switzerland

{chao.zhang, angela.bonifati}@univ-lyon1.fr, {hugo.kapp, vlad.haprian, jean-pierre.lozi}@oracle.com

Abstract—Reachability queries checking the existence of a path from a source node to a target node are fundamental operators for querying and processing graph data. Current approaches for index-based evaluation of reachability queries either focus on plain reachability or constraint-based reachability with only alternation of labels. In this paper, for the first time we study the problem of index-based processing for recursive label-concatenated reachability queries, referred to as RLC queries. These queries check the existence of a path that can satisfy the constraint defined by a concatenation of at most k edge labels under the Kleene plus. Many practical graph database and network analysis applications exhibit RLC queries. However, their evaluation remains prohibitive in current graph database engines.

We introduce the RLC index, the first reachability index to efficiently process RLC queries. The RLC index checks whether the source vertex can reach an intermediate vertex that can also reach the target vertex under a recursive label-concatenated constraint. We propose an indexing algorithm to build the RLC index, which guarantees the soundness and the completeness of query execution and avoids recording redundant index entries. Comprehensive experiments on real-world graphs show that the RLC index can significantly reduce both the offline processing cost and the memory overhead of transitive closure, while improving query processing up to six orders of magnitude over online traversals. Finally, our open-source implementation of the RLC index significantly outperforms current mainstream graph engines for evaluating RLC queries.

Index Terms—reachability index, graph query, graph databases, RLC queries

I. INTRODUCTION

Graphs have been the natural choice of data representation in various domains [1], *e.g.*, social, biochemical, fraud detection and transportation networks, and reachability queries are fundamental graph operators [2]. Plain reachability queries check whether there exists a path from a source vertex to a target vertex, for which various indexing techniques have been proposed [3]–[21]. To facilitate the representation of different types of relationships in real-world applications, *edge-labeled graphs* and *property graphs*, where labels can be assigned to edges, are more widely adopted nowadays than unlabeled graphs. Such advanced graph models allow users to add path constraints when defining reachability queries, which play a key role in graph analytics. However, current index-based approaches focus on constraint-based reachability with only alternation [22]–[26]. In this paper, we consider for the first time reachability queries with a complex path constraint corresponding to a *concatenation* of edge labels under

the Kleene plus, referred to as *recursive label-concatenated queries* (RLC queries). RLC queries are a prominent subset of regular path queries, for which indexing techniques have not been understood yet. As such, they are significant counterparts of queries with alternation of labels as path constraints [24] and queries without path constraints, *i.e.*, plain reachability queries [3]. We further motivate RLC queries by means of a running example.

Example 1: Figure 1 shows a property graph inspired by a real-world use case encoding an interleaved social and professional network along with information of bank accounts of persons. RLC queries can be used to detect fraud and money laundering patterns among financial transactions. For instance, the query $Q1(A_{14}, A_{19}, (\text{debits}, \text{credits})^+)$ checks whether there is a path from account A_{14} to A_{19} such that the label sequence of the path is a concatenation of an arbitrary number of occurrences of $(\text{debits}, \text{credits})$, which can lead to detect suspicious patterns of money transfers between these accounts. The RLC query $Q1((A_{14}, A_{19}, (\text{debits}, \text{credits})^+)$ evaluates to *true* because of the existence of the path $(A_{14}, \text{debits}, E_{15}, \text{credits}, A_{17}, \text{debits}, E_{18}, \text{credits}, A_{19})$. Another example is $Q2(P_{10}, P_{13}, (\text{knows}, \text{knows}, \text{worksFor})^+)$ that evaluates to *false* because there is no path from P_{10} to P_{13} satisfying the constraint.

RLC queries are also frequently occurring in real-world query logs, *e.g.*, Wikidata Query Logs [27], which is the largest repository of open-source graph queries (of the order of 500M queries). In particular, RLC queries often timed out in these logs [27] thus showing the limitations of graph query engines to efficiently evaluate them. Moreover, Neo4j (v4.3) [28] and TigerGraph (v3.3) [29], two of the mainstream graph data processing engines, do not yet support RLC queries in their current version. On the other hand, these systems have already identified the need to support these queries in the near future by following the developments of the Standard Graph Query Language (GQL) [30]. RLC queries can be expressed in Gremlin supported by TinkerPop-Enabled Graph Systems [31], *e.g.*, Amazon Neptune [32], in PGQL [33] supported by Oracle PGX [34], [35], and in SPARQL 1.1 (ASK query) supported by Virtuoso [36] and Apache Jena [37], among the others. However, many of these systems cannot efficiently evaluate RLC queries yet as shown in our experimental study.

To the best of our knowledge, little research has been

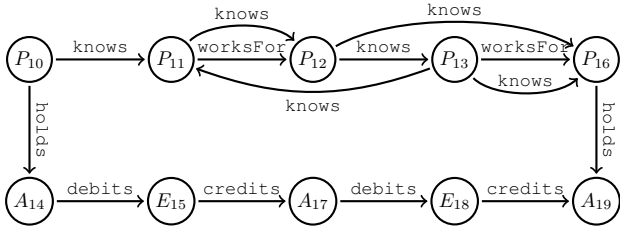


Fig. 1. A social and professional network for illustrating RLC queries.

carried out on an index-based solution to evaluate RLC queries, whereas indexing is a desirable asset of future graph processing systems allowing to improve and predict the performances of graph queries [38]. Previous plain reachability indexes have been well studied [3]–[18], but are not usable to evaluate RLC queries because of the lack of edge label information. More recent indexing methods have been focused on alternation of edge labels [22]–[26], instead of concatenation in RLC queries.

In this work, we propose the RLC index, the first reachability index tailored for RLC queries. In the RLC index, we assign pairs of vertices and succinct label sequences to each vertex. Then, the RLC index processes an RLC query by checking whether the source vertex can reach an intermediate vertex that can also reach the target vertex under the path constraint.

Challenge C1. One of the challenges for building an index to process RLC queries is that there can be infinite sequences of edge labels from vertex s to vertex t due to the presence of cycles on paths from s to t . Thus, building an index for arbitrary RLC queries can be hard since the number of label sequences for cyclic graphs is exponentially growing and potentially infinite, which leads to unrealistic indexing time and index size. We overcame this issue by leveraging a practical observation: the number of edge labels concatenated under the Kleene plus (or star) is typically bounded in real-world query logs, *e.g.*, [27]. We refer to the maximum number of concatenated labels in a workload of RLC queries as recursive k . Given an arbitrary recursive k , we show that the RLC index can be correctly built to evaluate any RLC query with a recursive concatenation of *at most* k arbitrary edge labels. Note that recursive k only depends on the number of concatenated labels and not on the actual length of any path selected by RLC queries.

Challenge C2. Although abundant indexing algorithms [3]–[21] have been proposed for plain reachability, which are then extended to alternation-based reachability queries [22]–[26], none of them is able to build an index for RLC queries because of the completely different path constraints. More precisely, indexing reachability with recursive concatenation-based path-constraints requires recording *sequences* of edge labels. However, indexing reachability with alternation-based path-constraints requires *sets* of edge labels, or plain reachability without edge labels. Such a fundamental difference illustrates the hardness of the problem and shows the inadequacy of existing indexing algorithms to build the RLC index. We design a novel indexing algorithm that builds the RLC index efficiently over large and highly cyclic graphs with millions

of vertices and edges as shown in our experiments. The indexing algorithm searches and records paths with recursive concatenation-based constraints, and during each searching step, identifies the condition that allows to prune search space as well as avoid recording redundant index entries.

Contribution. Our contributions are summarized as follows:

- We study the problem of indexing for RLC queries, *i.e.*, reachability queries with a complex path constraint consisting of the Kleene plus over a concatenation of edge labels. We propose a novel instantiation of the problem to address challenge C1.
- We propose the RLC index, the first reachability index for processing RLC queries, along with its indexing algorithm to address challenge C2. We prove that the latter builds a sound and complete RLC index for an arbitrary recursive k , where redundant index entries can be greedily removed.
- Our comprehensive experiments using highly cyclic real-world graphs show that the RLC index can be efficiently built and can significantly reduce the memory overhead of the (extended) transitive closure. Moreover, the RLC index is capable of answering RLC queries efficiently, up to six orders of magnitude over online traversals. We also demonstrate the speed-up and the generality of using the RLC index to accelerate query processing on mainstream graph engines. Our code is available as open source¹.

The rest of the paper is organized as follows. Section II presents the related work while Section III introduces the RLC queries. Section IV presents the theoretical foundation and Section V describes the RLC index and its indexing algorithm. Section VI presents the experimental assessment of RLC queries using the RLC index. Section VII concludes our work. All related proofs are included in our online technical report [39] due to the space constraint.

II. RELATED WORK

The two most fundamental graph queries [40] are graph pattern matching queries and regular path queries (aka navigational queries). The former matches a query graph against a graph database while the latter recursively navigates a graph according to a regular expression. The two fundamental queries are completely orthogonal, and RLC queries belong to the class of the latter. Thus, we discuss regular path queries and related indexes in this section. We note approaches [41]–[57] for efficient evaluating graph pattern matching queries. We refer readers to surveys [58], [59]. We do not elaborate on them further and do not consider them in our experimental evaluation due to the completely different query types.

Plain Reachability Index. Given an unlabeled graph $G = (V, E)$ and a pair of vertices (s, t) , a plain reachability query asks whether there exists a path from s to t . The existing approaches lie between two extremes, *i.e.*, online traversals and the transitive closure. Various indexes have been proposed. Comprehensive surveys can be found in [18], [60], [61] Plain

¹Open Source Link: <https://github.com/g-rpqs/rlc-index>

reachability indexes mainly fall into two categories [17], [18]: (1) index-only approaches, *e.g.*, *Chain Cover* [4], [7], *Tree Cover* [3], *Dual Labeling* [6], *Path-Tree Cover* [8], 2-Hop labeling [5], TFL [10], TOL [19], PLL [21], 3-Hop labeling [9], among others; (2) index-with-graph-traversal approaches, such as *Tree+SSPI* [12], *GRIPP* [13], *GRAIL* [14], *Ferrari* [15], *IP* [17], and *BFL* [18].

RLC queries are different from plain reachability queries because they are evaluated on labeled graphs to find the existence of a path satisfying additional recursive label-concatenated constraints. Thus, the indexes used to evaluate plain reachability queries, *e.g.*, 2-hop labeling [5], are not suitable for RLC queries. More precisely, indexing techniques for plain reachability queries only record information about graph structure but ignore information of edge labels.

Alternation-Based Reachability Index. Reachability queries with a path constraint that is based on alternation of edge labels (instead of concatenation as in our work), are known as LCR queries in the literature. Index-based solutions for LCR queries have been extensively studied in the last decade.

Jin *et. al* [22] presented the first result on LCR queries. To compress the generalized transitive closure that records reachable pairs and sets of path-labels, the authors proposed sufficient label sets and a spanning tree with partial transitive closure. The Zou *et. al* [23] method finds all strongly connected components (SCCs) in an input graph and replaces each SCC with a bipartite graph to obtain an edge labeled DAG. Valstar *et. al* [24] proposed a landmark-based index, where the generalized transitive closure for a set of high-degree vertices called landmarks are built and an online traversal is applied to answer LCR queries, which is accelerated by hitting landmarks. The state-of-the-art indexing techniques for LCR queries are the Peng *et. al* [25] method and the Chen *et. al* [26] method. Peng *et. al* [25] proposed the LC 2-hop labeling, which extends the 2-hop labeling framework through adding minimal sets of path-labels for each entry in the 2-hop labeling. Chen *et. al* [26] proposed a recursive method to handle LCR queries, where an input graph is recursively decomposed into spanning trees and graph summaries.

LCR queries and RLC queries have completely different regular expressions as path constraints, which makes the corresponding indexing problem inherently different. The expression in LCR queries is an alternation of edge labels, while the one in RLC queries is a concatenation of edge labels. The completely different path constraint makes LCR indexes inapplicable for processing RLC queries. The difference on path constraints also makes the indexing algorithms for LCR queries and RLC queries fundamentally different. Specifically, for an LCR indexing algorithm, it is sufficient to traverse any cycle in a graph only once. Conversely, in the case of RLC queries, a cycle, especially a self loop, might need to be traversed multiple times depending on label sequences along paths. Therefore, indexing approaches for LCR queries are not applicable to indexing RLC queries.

Regular Path Queries. Regular path queries correspond to queries generating node pairs according to path constraints

specified using regular expressions. Under the *simple* path semantics (non-repeated vertices or edges in a path), it is NP-complete to check the existence of a path satisfying a regular expression [62]. Thus, a recent work [63] focused on an approximate solution for evaluating regular simple path queries. By restricting regular expressions or graph instances, there exist tractable cases [62], [64], including LCR queries but not RLC queries. When it comes to the *arbitrary* path semantics (allowing repeated vertices or edges in a path), regular path queries for generating node pairs can be processed by using automata-based techniques [65], [66] or bi-directional BFSs from rare labels [67]. Optimal solutions can be used for sub-classes of regular expressions, *e.g.*, a matrix-based method [68] for the case without recursive concatenation, and a B+tree index for the case with a bounded path length (non-recursive path constraints) [69] and a fixed path pattern [70]. There also exist in the literature partial evaluation for distributed graphs [71], and incremental approaches for streaming graphs [72], [73]. Compared to all these works, we focus on designing an index-based solution to handle reachability queries with path constraints of recursive concatenation over a static and centralized graph under the arbitrary path semantics, which is an open challenge in the design of reachability indexes. To the best of our knowledge, our work is the first of its kind focusing on the design of a reachability index for such queries.

III. PROBLEM STATEMENT

An edge-labeled graph is $G = (V, E, \mathbb{L})$, where \mathbb{L} is a finite set of labels, V a finite set of vertices, and $E \subseteq V \times \mathbb{L} \times V$ a finite set of labeled edges. For the graph in Fig. 1, we have $\mathbb{L} = \{\text{knows}, \text{worksFor}, \text{debts}, \text{credits}, \text{holds}\}$, and that $e_1 = (P_{10}, \text{knows}, P_{11})$ is a labeled edge. We use $\lambda : E \rightarrow \mathbb{L}$ to denote the mapping from an edge to its label, *e.g.*, $\lambda(e_1) = \text{knows}$. The frequently used symbols are summarized in Table I.

In this paper, we consider the *arbitrary path* semantics [40] for evaluating RLC queries, *i.e.*, vertices or edges can appear more than once along the path. The arbitrary path semantics is widely adopted in practical graph query languages, *e.g.*, SPARQL 1.1 and PGQL, because evaluating a reachability query with an arbitrary regular expression as the path constraint is tractable under the arbitrary path semantics but is NP-complete under the simple path semantics, *e.g.*, RLC queries. It should be noteworthy that LCR queries with alternation-based path constraints belong to a special class (*trC* class [64]), which remain tractable under the simple path semantics.

An arbitrary path p in G is a vertex-edge alternating sequence $p(v_0, v_n) = (v_0, e_1, \dots, e_n, v_n)$, where $n \geq 1$, and $v_0, \dots, v_n \in V$, $e_1, \dots, e_n \in E$, and $|p(v_0, v_n)| = n$ that is the length of the path. For $p(v_0, v_n)$, v_0 is the source vertex and v_n is the target vertex. If there exists a path from v_0 to v_n , then v_0 reaches v_n , denoted as $v_0 \rightsquigarrow v_n$. The label sequence of the path $p(v_0, v_n)$ is $\Lambda(p(v_0, v_n)) = (\lambda(e_1), \dots, \lambda(e_n))$. When the context is clear, we also use $\Lambda(u, v)$ to denote the sequence of edge labels of a path from u to v .

A. Minimum Repeats

We use l_i to denote an edge label, $L = (l_1, \dots, l_n)$ a label sequence, $|L| = n$ the length of L , and ϵ the empty label sequence, *i.e.*, $|\epsilon| = 0$. We use \circ to denote the concatenation of label sequences (or labels), *i.e.*, $(l_1, \dots, l_i) \circ (l_{i+1}, \dots, l_n) = (l_1, \dots, l_n)$, or $L \circ L = L^2$. For the empty label sequence ϵ , we define $L \circ \epsilon = \epsilon \circ L = L$.

The sequence $L' = (l'_1, \dots, l'_{n'})$, $|L'| = n'$ is a *repeat* of $L = (l_1, \dots, l_n)$, if there exists an integer z , such that $\frac{n}{n'} = z \geq 1$, and $l'_j = l_{j+i \times n'}$ for every $j \in (1, \dots, n')$ and $i \in (0, \dots, z-1)$. A repeat L' of L is minimum if L' has the shortest length of all the repeats of L . The *minimum repeat* (MR) of L is denoted as $MR(L)$ that is also a sequence of edge labels. For example, given the path $p = (P_{10}, \text{knows}, P_{11}, \text{worksFor}, P_{12}, \text{knows}, P_{13}, \text{worksFor}, P_{16})$ in Fig. 1, we have $MR(p(P_{10}, P_{16})) = (\text{knows}, \text{worksFor})$. If $L = MR(L)$, we also say L itself is a minimum repeat. Given a positive integer k and a label sequence L , if $|MR(L)| \leq k$, then we say L has a non-empty k-MR that is $MR(L)$.

Lemma 1: For a label sequence L , $MR(L)$ is unique.

Suppose L has two minimum repeats, then only the shorter is $MR(L)$. Therefore, we have Lemma 1.

B. RLC Query

We consider the path constraint $L^+ = (l_1, \dots, l_k)^+$, where '+' is the Kleene plus, *i.e.*, one-or-more concatenations of the sequence $L = (l_1, \dots, l_k)$. A label sequence $\Lambda(u, v)$ of a path $p(u, v)$ satisfies a label-constraint L^+ , if and only if $MR(\Lambda(u, v)) = L$. If such a path $p(u, v)$ exists, then u can reach v with the constraint L^+ , denoted as $u \xrightarrow{L^+} v$, otherwise $u \not\xrightarrow{L^+} v$. Notice that reachability queries with path-constraints based on recursive concatenation could ask for additional constraints related to path length. For instance, if the path constraint is $L^+ = (\text{knows}, \text{knows})^+$, then the query would need to additionally check whether the path length is even. In general, such fragment of queries, *i.e.*, queries with L^+ s.t. $L \neq MR(L)$, impose an additional constraint on path length leading to the complicated even-path problem (NP-complete [62], [64], [74] for simple paths), which is beyond the scope of our paper.

Definition 1 (RLC Query): Given an edge-labeled directed graph $G = (V, E, \mathbb{L})$ and a recursive k , an RLC query is a triple (s, t, L^+) , $L = (l_1, \dots, l_j)$, where $s, t \in V$, $L = MR(L)$, $j \leq k$, and $l_i \in \mathbb{L}$ for $i \in (1, \dots, j)$. If $s \xrightarrow{L^+} t$, then the answer to the query is *true*. Otherwise, the answer is *false*.

For the sake of simplicity in this paper we focus on the RLC queries with the Kleene plus. Queries (s, t, L^*) with the Kleene star can be trivially reduced to queries with the Kleene plus (s, t, L^+) through checking whether s is equal to t . Our method is thus applicable to RLC queries with the Kleene star.

Given an RLC query $Q(s, t, L^+)$, under the arbitrary path semantics, two naive approaches can be used to evaluate Q . As RLC queries are path queries with regular expressions, in the first approach RLC queries can be evaluated by online graph traversals, *e.g.*, BFS, guided by a minimized NFA

TABLE I
FREQUENTLY USED SYMBOLS.

Notation	Description
p , or $p(u, v)$	a path, or the path from u to v
\circ	concatenation of labels or label sequences
$\Lambda(u, v)$, or $\Lambda(p(u, v))$	the label sequence of a path from u to v
L	a label sequence
L^+	a label constraint
$MR(L)$	the minimum repeat of a label sequence L
k	the upper bound of the number of labels in L^+
$S^k(u, v)$	the concise set of minimum repeats from u to v
$u \xrightarrow{L^+} v$, or $u \not\xrightarrow{L^+} v$	u reaches v through an L^+ -path, or otherwise
$u \rightsquigarrow v$, or $u \not\rightsquigarrow v$	u reaches v , or otherwise
$in(v)$, or $out(v)$	the set of vertices that can reach v , or v can reach
$aid(v)$	the access id of vertex v by the indexing algorithm
$v_i^{(j)}$	a vertex with vertex id i and access id j

(nondeterministic finite automaton) [61] that is constructed according to the regular expression in an RLC query. The second approach leads to pre-computing the transitive closure, that for each pair of vertices (s, t) records whether $s \rightsquigarrow t$ and all the label sequences from s to t . Notice that the naive approach to build the transitive closure is not usable in our case because cycles may exist on the path from s to t such that the BFS from s can generate infinite label sequences for paths to t . To address this issue, we adopt an *extended transitive closure*, which is presented in Section VI. As demonstrated in our experiments, these two solutions require either too much query time or storage space and are impractical for large graphs.

C. Indexing Problem

Our goal is to build an index to efficiently process RLC queries. The indexing problem is summarized as follows.

Problem 3.1: Given an edge-labeled graph G , the indexing problem is to build a reachability index for processing RLC queries on G , such that the size of the index is minimal and the correctness of query processing is preserved.

We observe that recording MRs, instead of raw label sequences of paths in G , can reduce the storage space, and such a strategy does not violate the correctness of query processing. The main benefits are twofold: (1) MRs are not longer than raw label sequences; (2) different raw label sequences may have the same MR. For example, in Fig. 1, there exist two paths from P_{10} to P_{16} having the label sequence $(\text{knows}, \text{knows}, \text{knows}, \text{knows})$ and $(\text{knows}, \text{knows}, \text{knows})$, which have the same MR knows .

Definition 2 (Concise Label Sequences): Let $\mathbb{P}(s, t)$ be the set of all paths from s to t . The concise set of label sequences from vertex s to t , denoted as $S^k(s, t)$, is the set of k-MRs of all label sequences from s to t , *i.e.*, $S^k(s, t) = \{L | p \in \mathbb{P}(s, t), L = MR(\Lambda(p)), |L| \leq k\}$.

To process RLC queries, we need to compute and record the concise label sequences. We have Proposition 1 by definition.

Proposition 1: $s \xrightarrow{L^+} t, |L| \leq k$ in G iff $L \in S^k(s, t)$.

For example, in Fig. 1, we have $S^2(P_{12}, P_{16}) = \{(\text{knows}), (\text{knows}, \text{worksFor})\}$. With $S^2(P_{12}, P_{16})$, RLC queries from P_{12} to P_{16} can be processed correctly.

IV. KERNEL-BASED SEARCH

In this section, we deal with the following question: *how to compute concise label sequences?* The problem for computing a concise label sequence is that if a cycle exists on a path from s to t , there exist infinite paths from s to t , which makes the computation of $S^k(s, t)$ infeasible, e.g., $|\mathbb{P}(P_{11}, P_{13})|$ in Fig. 1 is infinite. We overcome this issue by leveraging the upper bound of recursively concatenated labels in a constraint, i.e., recursive k . We observed that we don't have to compute all possible label sequences for paths going from P_{11} to P_{13} as the set of label sequences L such that $|MR(L)| \leq k$ is actually finite. In general, let v be an intermediate vertex that a forward breadth-first search from s is visiting. The main idea is that when the path from s to v reaches a specific length, we can decide whether we need to further explore the outgoing neighbours of v . Moreover, if the outgoing neighbours of v are worth exploring, the following search can be guided by a specific label constraint. In the following, we first provide an illustrating example, and then formally define the specific constraint that is used to guide the subsequent search.

Example 2 (Illustrating Example): Consider the graph in Fig. 1. Assume we need to compute $S^2(P_{11}, P_{13})$, i.e., $k = 2$, and we perform a breadth-first search from P_{11} . When P_{13} is visited for the first time, we add (knows) and (worksFor, knows) into $S^2(P_{11}, P_{13})$. After that, when the search depth reaches $2k = 4$, i.e., P_{12} is visited for the second time, we can have 4 different label sequences, which are $L_1 = (\text{knows}, \text{knows}, \text{knows}, \text{knows})$, $L_2 = (\text{knows}, \text{knows}, \text{knows}, \text{worksFor})$, $L_3 = (\text{worksFor}, \text{knows}, \text{knows}, \text{knows})$, and $L_4 = (\text{worksFor}, \text{knows}, \text{knows}, \text{worksFor})$. Given this, all the 4 label sequences except L_1 do not need to be expanded anymore, because their expansions cannot produce a minimum repeat whose length is not larger than 2. After this, the following search continued with L_1 is guided by (knows)⁺ that is computed from L_1 . However, because there already exists (knows) in $S^2(P_{11}, P_{13})$, the search terminates.

Definition 3: If a label sequence L can be represented as $L = (L')^h \circ L''$, where $h \geq 2$, $L' \neq \epsilon$ and $MR(L') = L'$, and L'' is ϵ or a proper prefix of L' , then L has the *kernel* L' and the *tail* L'' .

For example, the label sequence (knows, knows, knows, knows) from P_{11} has a kernel knows and a tail ϵ .

Kernel-based search. When a kernel has been determined at a vertex that is being visited, the subsequent search to compute $S^k(s, t)$ can be guided by the Kleene plus of the kernel, e.g., (knows)⁺ is used to guide the search in Example 2. We call this approach KBS (*kernel-based search*) in the remainder of this paper. In a nutshell, KBS consists of two phases: (1) *kernel-search* and (2) *kernel-BFS*, where the first phase is to compute kernels, and the second to perform kernel-guided BFS. We show in Theorem 1 that KBS can compute a sound and complete $S^k(s, t)$. The proof of Theorem 1 is included in our technical report [39] due to the space limit.

Theorem 1: Given a path p from u to v and a positive integer

k , p has a non-empty k -MR if and only if one of the following conditions is satisfied,

- Case 1: $|p| \leq k$. $MR(\Lambda(p))$ is the k -MR of p ;
- Case 2: $k < |p| \leq 2k$. If $|MR(\Lambda(p))| \leq k$, $MR(\Lambda(p))$ is the k -MR of p ;
- Case 3: $|p| > 2k$. Let x be the intermediate vertex on p , s.t. $|p(u, x)| = 2k$. If $\Lambda(p(u, x))$ has a kernel L' and a tail L'' , and $MR(L'' \circ \Lambda(p(x, v))) = L'$, then L' is the k -MR of p .

We discuss below two strategies to compute kernels based on Theorem 1, namely *lazy* KBS and *eager* KBS, and explain why eager KBS is better than lazy KBS, which is used in our indexing algorithm presented in Section V-B.

Lazy KBS. Theorem 1 can be transformed into an algorithm to find kernels, i.e., for a source vertex we generate all paths of length $2k$, and then compute all the kernels of these paths. This strategy is referred to as lazy KBS, which means kernels are correctly determined when the length of paths reaches $2k$, e.g., lazy KBS is used in Example 2.

Eager KBS. In contrast to the lazy strategy, we can determine kernel candidates earlier, instead of valid kernels that require the length of paths to be $2k$. The main idea is to treat any k -MR that is computed using any path p , $|p| \leq k$ as a kernel candidate which is then used to guide KBS. Although an invalid kernel may be included, the search guided by the invalid kernel will not reach a target vertex through a path of which the k -MR is the invalid kernel. Thus, the result computed by the eager strategy is still sound and complete.

Example 3: Consider the example of computing $S^2(P_{10}, P_{13})$ in Fig. 1. Using the eager strategy, when P_{12} is visited for the first time, two kernel candidates can be determined, i.e., (knows) and (knows, worksFor). Although (knows, worksFor)⁺ is an invalid kernel, the search guided by it cannot reach P_{13} .

The key advantage of the eager strategy over the lazy strategy is that it allows us to advance KBS from the kernel-search phase to the kernel-BFS phase. This can make KBS more efficient because generating all label sequences of length $2k$ from a source vertex is more expensive than the case of paths of length k , especially on a dense graph.

V. RLC INDEX

In this section, we present the RLC index, and the corresponding query and indexing algorithm.

A. Overarching Idea

Given an RLC query (s, t, L^+) , $|L| \leq k$, the idea is to check whether there exists a path (s, \dots, u, \dots, t) whose label sequence satisfies the label constraint L^+ , where u is an intermediate vertex in p . In other words, the query is answered by concatenating two MRs of the sub-paths of p , i.e., $MR(\Lambda(s, u))$ and $MR(\Lambda(u, t))$.

Definition 4 (RLC Index): Let $G = (V, E, \mathbb{L})$ be an edge-labeled graph and recursive k be a positive integer. The RLC index of G assigns to each vertex $v \in V$ two sets: $\mathcal{L}_{in}(v) = \{(u, L') \mid u \rightsquigarrow v, L' \in S^k(u, v)\}$, and $\mathcal{L}_{out}(v) = \{(w, L'') \mid v \rightsquigarrow w, L'' \in S^k(v, w)\}$. Therefore, there is a path

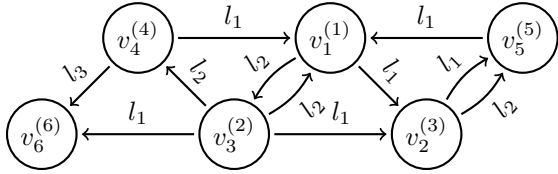


Fig. 2. A graph instance G for illustrating the RLC index.

$p(s, t)$ satisfying an arbitrary constraint L^+ , $|L| \leq k$, if and only if one of the following cases is satisfied,

- Case 1: $\exists(x, L') \in \mathcal{L}_{out}(s)$ and $\exists(x, L'') \in \mathcal{L}_{in}(t)$, such that $L' = L'' = L$;
- Case 2: $\exists(t, L''') \in \mathcal{L}_{out}(s)$ or $\exists(s, L''') \in \mathcal{L}_{in}(t)$, such that $L''' = L$.

Example 4 (Running Example of the RLC Index): Consider the graph G shown in Fig. 2. The RLC index with recursive $k = 2$ for G is presented in Table II. We have $Q_1(v_3, v_6, (l_2, l_1)^+) = true$ because $\exists(v_1, (l_2, l_1)) \in \mathcal{L}_{out}(v_3)$ and $\exists(v_1, (l_2, l_1)) \in \mathcal{L}_{in}(v_6)$. Indeed, there exists the path $(v_3, l_2, v_4, l_1, v_1, l_2, v_3, l_1, v_6)$ from v_3 to v_6 in the graph in Fig. 2. For $Q_2(v_1, v_2, (l_2, l_1)^+)$, the answer is *true* because $\exists(v_1, (l_2, l_1)) \in \mathcal{L}_{in}(v_2)$. Given $Q_3(v_1, v_3, (l_1)^+)$, the answer is *false*. Although v_1 can reach v_3 , e.g., $\exists(v_1, l_2) \in \mathcal{L}_{in}(v_3)$, the constraint $(l_1)^+$ of Q_3 cannot be satisfied.

Whereas our indexing framework leverages the canonical 2-hop labeling framework for plain reachability queries [5], indexing RLC queries is inherently more challenging due to the presence of complex recursive label concatenations, which calls for the design of a novel indexing algorithm.

In order to save storage space, when building an RLC index, redundant index entries have to be removed as many as possible, i.e., building a minimal (or condensed) RLC index. The underlying idea is that if there exists a path p such that $u \xrightarrow{L^+} v$, then the RLC index only records the reachability information of the path p once, i.e., either through Case 1 or Case 2 in Definition 4. This leads to the following definition.

Definition 5 (Condensed RLC Index): The RLC index is condensed, if for every index entry $(s, L) \in \mathcal{L}_{in}(t)$ (or $(t, L) \in \mathcal{L}_{out}(s)$), there do not exist index entries $(u, L') \in \mathcal{L}_{out}(s)$ and $(u, L'') \in \mathcal{L}_{in}(t)$ such that $L = L' = L''$.

We focus on designing an indexing algorithm that can build a correct (sound and complete) and condensed RLC index.

B. Query and Indexing Algorithm

The query algorithm is presented in Algorithm 1, where we use I to denote an index entry. Each index entry I has the schema (vid, mr) , where vid represents vertex id and mr recorded minimal repeat. Given an RLC query (s, t, L^+) , to efficiently find $(u, L') \in \mathcal{L}_{out}(s)$ and $(u, L'') \in \mathcal{L}_{in}(t)$, we execute a merge join over $\mathcal{L}_{out}(s)$ and $\mathcal{L}_{in}(t)$, shown at line 4 in Algorithm 1. The output of the merge join is a set of index entry pairs (I', I'') , s.t. $I'.vid = I''.vid$. Case 1 of the RLC index (see Definition 4) is checked at line 5. Case 2 is checked at line 3. If one of these cases can be satisfied, the answer *true* will be returned immediately. Otherwise, index

TABLE II
THE RLC INDEX FOR THE GRAPH IN FIG. 2.

v	$\mathcal{L}_{in}(v)$	$\mathcal{L}_{out}(v)$
v_1	-	$(v_1, l_2), (v_1, l_1), (v_1, (l_2, l_1))$
v_2	$(v_1, l_1), (v_1, (l_2, l_1))$	$(v_1, (l_2, l_1)), (v_1, l_1)$
v_3	$(v_1, l_2), (v_1, (l_1, l_2))$	$(v_1, l_2), (v_1, (l_2, l_1)), (v_1, l_1), (v_3, (l_1, l_2))$
v_4	(v_1, l_2)	$(v_1, l_1), (v_3, (l_1, l_2))$
v_5	$(v_1, (l_1, l_2)), (v_1, l_1), (v_3, (l_1, l_2)), (v_2, l_2)$	$(v_1, l_1), (v_3, (l_1, l_2))$
v_6	$(v_1, (l_2, l_1)), (v_3, l_1), (v_3, (l_2, l_3)), (v_4, l_3)$	-

Algorithm 1: Query Algorithm.

```

1 procedure Query( $s, t, L^+$ )
2   if  $\exists(t, L) \in \mathcal{L}_{out}(s)$  or  $\exists(s, L) \in \mathcal{L}_{in}(t)$  then
3     return true;
4   for  $(I', I'') \in mergeJoin(\mathcal{L}_{out}(s), \mathcal{L}_{in}(t))$  do
5     if  $I'.mr = L$  and  $I''.mr = L$  then
6       return true;
7   return false;
```

entries in $\mathcal{L}_{out}(s)$ and $\mathcal{L}_{in}(t)$ are exhaustively merged, and the answer *false* will be returned at last.

In the remainder of this subsection, we present an indexing algorithm (Algorithm 2) to build the RLC index that is sound, complete and condensed. We use v_i to denote a vertex with id i . Given a graph $G(V, E, \mathbb{L})$, the indexing algorithm mainly performs backward and forward KBS from each vertex in V to create index entries, and *pruning rules* are applied to accelerate index building as well as remove redundant index entries.

Indexing using KBS. We explain below how the backward KBS creates \mathcal{L}_{out} -entries. The forward KBS follows the same procedure, except that \mathcal{L}_{in} -entries will be created. Suppose that the backward KBS from vertex v_i is visiting v . If $|MR(\Lambda(v, v_i))| \leq k$, then we add $(v_i, MR(\Lambda(v, v_i)))$ into $\mathcal{L}_{out}(v)$. Although there may be cycles in a graph, the KBS will not go on forever, because when the depth of the search reaches k , the KBS will be transformed from its kernel-search phase into its kernel-BFS phase that is then guided by the Kleene plus of a kernel candidate, such that the KBS can terminate if any invalid label (or state) transition is met, or a vertex being visited with a label l_i of the kernel has been already visited with a label l_j of the kernel, s.t. $i = j$.

KBSs are executed from each vertex in V and the execution follows a specific order. The idea is to start with vertices that have more connections to others, allowing such vertices to be intermediate hops to remove redundancy. In the RLC index, we leverage the IN-OUT strategy, i.e., sorting vertices according to $(|out(v)| + 1) \times (|in(v)| + 1)$ in descending order, which is known as an efficient and effective strategy for various reachability indexes based on the 2-hop labeling framework. The id of vertex v in the sorted list is referred to as the access id of v , denoted as $aid(v)$ starting from 1, e.g., for the graph in Fig. 2, the sorted list is $(v_1, v_3, v_2, v_4, v_5, v_6)$, where $aid(v_3) = 2$, or simply $v_3^{(2)}$ in Fig. 2.

Example 5 (Running Example of Indexing): Consider the

Algorithm 2: Indexing Algorithm.

```
1 procedure kernelBasedSearch( $v, k$ )
2   for  $(L, vSet) \in$  backwardKernelSearch( $v, k$ ) do
3      $\backslash$  backwardKernelBFS( $v, vSet, L$ );
4   for  $(L, vSet) \in$  forwardKernelSearch( $v, k$ ) do
5      $\backslash$  forwardKernelBFS( $v, vSet, L$ );
6 procedure backwardKernelSearch( $v, k$ )
7    $q \leftarrow$  an empty queue of (vertex, label sequence);
8    $q.enqueue(v, \epsilon)$ ;
9    $map \leftarrow$  a map of (kernel candidates, vertex set);
10  while  $q$  is not empty do
11     $(x, seq) \leftarrow q.dequeue()$ ;
12    for in-coming edge  $e(y, x)$  to  $x$  do
13       $seq' \leftarrow \lambda(e(y, x)) \circ seq$ ;  $L \leftarrow MR(seq')$ ;
14      insert( $y, v, L$ );
15       $map.get(L).add(x)$ ;
16      if  $|seq'| < k$  then
17         $\backslash$   $q.enqueue(y, seq')$ ;
18  return  $map$ ;
19 procedure insert( $s, t, L$ )
20  if  $aid(t) > aid(s)$  or  $Query(s, t, L^+)$  then
21     $\backslash$  return false;
22  else
23    add( $t, L$ ) into  $\mathcal{L}_{out}(s)$ ;
24    return true;
25 procedure backwardKernelBFS( $v, vSet, L$ )
26   $q \leftarrow$  an empty queue of (vertex, integer);
27  for  $x \in vSet$  do
28     $\backslash$  mark  $x$  as visited by state 1,  $q.enqueue(x, |L|)$ ;
29  while  $q$  is not empty do
30     $(x, i) \leftarrow q.dequeue()$ ,  $i \leftarrow i - 1$ ;
31    if  $i = 0$  then  $i = |L|$ ;
32    label  $l \leftarrow L.get(i)$ ;
33    for in-coming edge  $e(y, x)$  to  $x$  do
34      if  $l \neq \lambda(e(y, x))$  or  $y$  was visited by state  $i$  then
35         $\backslash$  continue;
36      if  $i = 1$  and insert( $y, v, L$ ) then
37         $\backslash$  continue;
38     $q.enqueue(y, i)$ ; mark  $y$  visited by state  $i$ ;
```

graph in Fig. 2 and the RLC index in Table II with recursive $k = 2$. The KBSs are executed from each vertex in the order of $(v_1, v_3, v_2, v_4, v_5, v_6)$. We explain below the backward KBS from v_1 , which is the first search of the indexing algorithm. The traversal of depth 1 of this backward KBS visits v_4 and creates (v_1, l_1) in $\mathcal{L}_{out}(v_4)$, visits v_3 and creates (v_1, l_2) in $\mathcal{L}_{out}(v_3)$, and visits v_5 and creates $(v_1, l_1) \in \mathcal{L}_{out}(v_5)$. The traversal of depth 2 creates $(v_1, (l_2, l_1))$ in $\mathcal{L}_{out}(v_3)$, (v_1, l_2) in $\mathcal{L}_{out}(v_1)$, $(v_1, l_1) \in \mathcal{L}_{out}(v_2)$, and $(v_1, (l_2, l_1)) \in \mathcal{L}_{out}(v_2)$. Then, the kernel-search phase of this KBS terminates as the depth of the search reaches 2, which generates kernel candidate l_1 with a set of frontier vertices $\{v_4, v_5, v_2\}$, kernel candidate l_2 with a set of frontier vertices $\{v_3, v_1\}$, and kernel candidate (l_2, l_1) with a set of frontier vertices (v_3, v_2) . After this, this KBS is turned into three kernel-BFSs guided by $(l_1)^+$, $(l_2)^+$, and $(l_2, l_1)^+$ with the corresponding frontier vertices. The kernel-BFS terminates under the case of an invalid label transition or a repeated visiting. For example, the label of the

incoming edge of v_3 is l_2 , which is an invalid state transition of $(l_2, l_1)^+$ in the backward kernel-BFS from v_1 , such that the kernel-BFS guided by $(l_2, l_1)^+$ terminates at v_3 . For another example, index entry $(v_1, l_1) \in \mathcal{L}_{out}(v_1)$ is created when v_1 is visited for the first time by the kernel-BFS from v_1 guided by $(l_1)^+$, but this kernel-BFS will not continue when it visits v_5 that has already been visited with the label l_1 .

Pruning Rules. To speed up index construction and remove redundant index entries, we apply pruning rules during KBSs. For ease of presentation, we present the rules for backward KBSs, and the same rules apply for forward ones.

- **PR1:** If the k -MR of an index entry that needs to be recorded can be acquired from the current snapshot of the RLC index, then the index entry can be skipped.
- **PR2:** If vertex v_i is visited by the backward KBS performed from vertex $v_{i'}$ s.t. $aid(v_{i'}) > aid(v_i)$, then the corresponding index entry can be skipped.
- **PR3:** If vertex v_i is visited by the kernel-BFS phase of a backward KBS performed from vertex $v_{i'}$, and PR1 (or PR2) is triggered, then vertex v_i and all the vertices in $in(v_i)$ are skipped.

The correctness of Algorithm 2 with pruning rules is guaranteed by Theorem 3 presented in Section V-C.

Example 6 (Running Example of Pruning Rules): Consider the forward KBS from v_3 for the graph in Fig. 2. It can visit v_2 through label sequence (l_2, l_1) , such that it tries to create $(v_3, (l_2, l_1))$ in $\mathcal{L}_{in}(v_2)$. However, there already exist $(v_1, (l_2, l_1)) \in \mathcal{L}_{out}(v_3)$ and $(v_1, (l_2, l_1)) \in \mathcal{L}_{in}(v_2)$, such that $Q(v_3, v_2, (l_2, l_1)^+) = true$ with the current snapshot of the RLC index, i.e., the reachability information has already been recorded. Thus, the index entry $(v_3, (l_2, l_1))$ in $\mathcal{L}_{in}(v_2)$ is pruned according to PR1. As an example of PR2, consider the backward KBS from v_2 . It can visit v_1 through path $(v_1, l_2, v_3, l_1, v_2)$, such that it tries to create $(v_2, (l_2, l_1))$ in $\mathcal{L}_{out}(v_1)$. Given $aid(v_2) > aid(v_1)$, such that the index entry can be pruned by PR2. Consider the forward KBS from v_2 for an example of PR3. It visits v_2 through path $(v_2, l_2, v_5, l_1, v_1, l_2, v_3, l_1, v_2)$, where at the edge (v_5, l_1, v_1) the KBS is transformed from the kernel-search phase to the kernel-BFS phase that is then guided by $(l_2, l_1)^+$. When v_2 visits itself for the first time, the KBS tries to create index entry $(v_2, (l_2, l_1)) \in \mathcal{L}_{in}(v_2)$. However, this index entry can be pruned by PR1 because of $(v_1, (l_2, l_1)) \in \mathcal{L}_{out}(v_2)$ and $(v_1, (l_2, l_1)) \in (v_2)$. Then, PR3 is triggered, which means the kernel-BFS from v_2 with the kernel $(l_2, l_1)^+$ can terminate.

The indexing algorithm is presented in Algorithm 2. For ease of presentation, each procedure focuses on the backward case, and the forward case can be obtained by trivial modifications, e.g., replacing in-coming edges with out-going edges. We use the KMP algorithm [75] to compute the minimum repeat of a label sequence, i.e., $MR()$ at line 13 in Algorithm 2. The indexing algorithm performs backward and forward KBS from each vertex. The KBS from a vertex v consists of two phases: kernel-search (line 6 to line 18) and kernel-BFS (line 25 to line 38). The kernel-search returns for each vertex v all kernel candidates and a set of frontier vertices

vSet. The kernel-BFS is performed for each kernel candidate, *i.e.*, a BFS with vertices in *vSet* as frontier vertices guided by a kernel candidate. PR1 and PR2 are included at line 20, which can be triggered by both kernel-search and kernel-BFS. However, RP3 can only be triggered by kernel-BFS, which is implemented at line 36, *i.e.*, if the `insert` function returns true, then vertices can be pruned during the search. However, PR3 is not related to kernel-search that focuses on generating kernels, such that the result returned by the `insert` function at line 14 is not taken into account.

C. Analysis of RLC Index

We present in Theorem 2 that pruning rules can guarantee the condensed property of the RLC index, and in Theorem 3 that the RLC index constructed by Algorithm 2 is correct, *i.e.*, sound and complete. The proofs of Theorem 2 and 3 are included in our technical report [39] due to the space limit.

Theorem 2: With pruning rules, the RLC index is condensed.

Theorem 3: Given an edge-labeled graph G and the RLC index of G with a positive integer k built by Algorithm 2, there exists a path from vertex s to vertex t in G which satisfies a label constraint L^+ , $|L| \leq k$, if and only if one of the following condition is satisfied

- (1) $\exists(x, L) \in \mathcal{L}_{out}(s)$ and $\exists(x, L) \in \mathcal{L}_{in}(t)$;
- (2) $\exists(t, L) \in \mathcal{L}_{out}(s)$, or $\exists(s, L) \in \mathcal{L}_{in}(t)$.

We analyze the complexity of the RLC index below.

Query time. Given a query $Q(s, t, L^+)$, the time complexity of Algorithm 1 is $O(|\mathcal{L}_{out}(s)| + |\mathcal{L}_{in}(t)|)$, because we only need to take $O(|\mathcal{L}_{out}(s)| + |\mathcal{L}_{in}(t)|)$ time to apply the merge join to find $(x, L) \in \mathcal{L}_{out}(s)$ and $(x, L) \in \mathcal{L}_{in}(t)$. Note that index entries in $\mathcal{L}_{out}(s)$ and $\mathcal{L}_{in}(t)$ have already been sorted according to the access id of vertices, such that we do not need to sort index entries when applying the merge join.

Index size. The index size can be $O(|V|^2|\mathbb{L}|^k)$ in the worst case, since each $\mathcal{L}_{in}(v)$ or $\mathcal{L}_{out}(v)$ can contain $O(|V|C)$ index entries, where $C = O(|\mathbb{L}|^k)$ is the number of distinct minimum repeats for all label sequences derived from $|\mathbb{L}|$ of length up to k . C can be computed as follows, $C = \sum_{i=1}^k F(i)$, where $F(i) = |\mathbb{L}|^i - (\sum_{j \in f(i), j \neq i} F(j))$ with $F(1) = |\mathbb{L}|$ and $f(i)$ the set of factors of integer i .

Indexing time. In Algorithm 2, we perform a KBS from each vertex, and each KBS consists of two phases: the kernel-search phase and the kernel-BFS phase. Performing a kernel-search of depth k from a vertex requires $O(|\mathbb{L}|^k|V|^k)$ time, and generates $O(|L|^k)$ kernel candidates as discussed in the index size analysis. Each kernel candidate requires a kernel-BFS taking $O(|E|k)$ time. Hence the time complexity for performing a KBS is $O(|\mathbb{L}|^k|V|^k + |L|^k|E|k)$. The total index time is $O(|V|^{k+1}|\mathbb{L}|^k + |L|^k|V||E|k)$ in the worst case.

Notice that these complexities resemble the complexity classes of previous label-constraint reachability indexes [24], [25] for LCR queries.

VI. EXPERIMENTAL EVALUATION

In this section, we study the performance of the RLC index. We first used real-world graphs to evaluate the indexing

TABLE III
OVERVIEW OF REAL-WORLD GRAPHS.

Dataset	$ V $	$ E $	$ \mathbb{L} $	Synthetic Labels	Loop Count	Triangle Count
Advogato (AD)	6K	51K	3		4K	98K
Soc-Epinions (EP)	75K	508K	8	✓	0	1.6M
Twitter-ICWSM (TW)	465K	834K	8	✓	0	38K
Web-NotreDame (WN)	325K	1.4M	8	✓	27K	8.9M
Web-Stanford (WS)	281K	2M	8	✓	0	11M
Web-Google (WG)	875K	5M	8	✓	0	13M
Wiki-Talk (WT)	2.3M	5M	8	✓	0	9M
Web-BerkStan (WB)	685K	7M	8	✓	0	64M
Wiki-hyperlink (WH)	1.7M	28.5M	8	✓	4K	52M
Pokec (PR)	1.6M	30.6M	8	✓	0	32M
StackOverflow (SO)	2.6M	63.4M	3		15M	114M
LiveJournal (LJ)	4.8M	68.9M	50	✓	0	285M
Wiki-link-fr (WF)	3.3M	123.7M	25	✓	19K	30B

time, index size, and query time of the RLC index, where we focus on practical graphs and workloads. Then, we used synthetic graphs to conduct a comprehensive study of the impact of different characteristics on the RLC index, including label set size, average degree, and the number of vertices. Finally, we compared the query time of the RLC index with existing systems, where we additionally consider more types of reachability queries from real-world query logs in order to demonstrate the generality of our approach.

a) Baselines: To the best of our knowledge, the RLC index is the first indexing technique designed for processing recursive label-concatenated reachability queries, and indices for other types of reachability queries are not usable in our context because of specific path constraints defined in the RLC queries. Thus, the chosen baselines for the RLC index are online graph traversals guided by NFAs (see Section III-B). We consider both BFS (breadth-first search) and BiBFS (bidirectional BFS) as the underlying online traversal methods, and simply refer to the baselines as BFS and BiBFS in this section. We note that DFS (depth-first search) is an alternative to BFS with the same time complexity but is not as efficient as BiBFS. In addition, we also include an extended transitive closure as a baseline, referred to as ETC. The indexing algorithm of ETC performs a forward KBS from each vertex without pruning rules, and records for every reachable pair of vertices (u, v) any k-MR of any path $p(u, v)$. In ETC, we use a hashmap to store reachable pairs of vertices and the corresponding set of k-MRs. There are two differences between the indexing algorithm of ETC and the one of the RLC index: (1) only forward KBS is used for building ETC, instead of forward and backward KBS for the RLC index, and (2) none of the pruning rules is applied for building ETC.

b) Datasets: We use real-world datasets and synthetic graphs in our experiments. We present the statistics of real-world datasets in Table III (sorted according to $|E|$), which are from either the SNAP [76] or the KONECT [77] project. We also include for each graph the loop count (cycles of length 1) and the triangle count (cycles of length 3) shown in the last two columns in Table III. We generate synthetic labels for graphs that do not have labels on their edges, indicated by the column of synthetic labels in Table III. The edge labels have been generated according to the Zipfian distribution [78] with exponent 2. The synthetic graphs used in our experiments

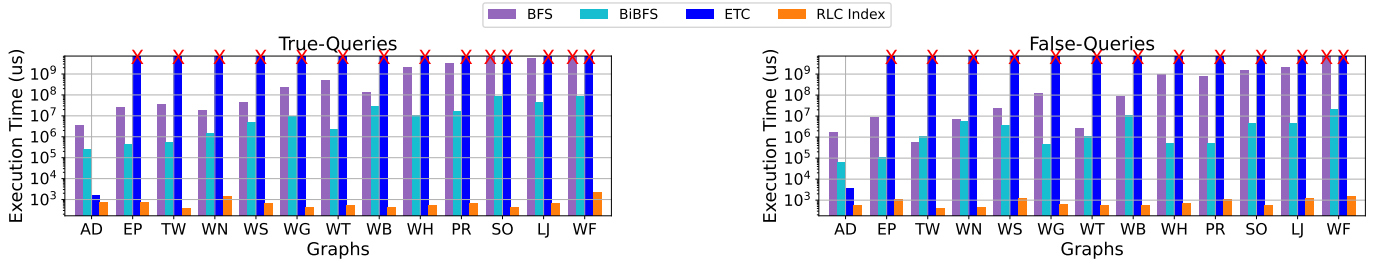


Fig. 3. Execution time of 1000 true-queries or 1000 false-queries on real-world graphs.

follows two different modes, namely the *Erdős-Rényi* (ER) model and the *Barabási-Albert* (BA) model. We use JGraphT [79] to generate the ER- and BA-graphs. The main difference between ER-graphs and BA-graphs is that ER-graphs have an almost uniform degree-distribution while BA-graphs have a skew in it because BA-graphs contain a complete sub-graphs. The method to assign labels to edges in synthetic graphs is the same as the one used for real-world graphs.

c) Query generation: As a common practice to evaluate a reachability index in the literature, *e.g.*, [24]–[26], we generate for each real-world graph comprehensive query sets, each of which contains 1000 true-queries and 1000 false-queries, respectively. We explain the method for query generation as follows. We uniformly select a source vertex s and a target vertex t , and also uniformly choose a label constraint L^+ . Then, a bidirectional breadth-first search is conducted to test whether s reaches t under the constraint of L^+ . If the test returns *true*, we add $(s, t, L^+, true)$ to the true-query set, otherwise we add it into the false-query set. After that, we generate another (s, t, L^+) , and repeat the above procedure until the completion of the two query sets.

d) Implementation and Setting: Our open-source implementation has been done in Java 11, spanning baseline solutions and the RLC index. We run experiments on a machine with 8 virtual CPUs of 2.40GHz, and 128GB main memory, where the heap size of JVM is configured to be 120GB.

A. Performance on Real-World Graphs

In this section, we analyze the performance of the RLC index on real-world graphs. We first compare the RLC index with ETC in terms of indexing time and index size, and with BFS and BiBFS in terms of query time. The recursive k value is set to 2, and each query in the respective workload has a recursive concatenation of 2 labels. The goal of the first experiment is to understand the performance of the RLC index for practical RLC query workloads whose length of recursive concatenation is at most 2, as in practical property paths [27]. After that, we conduct experiments on real-world graphs with different recursive k values to understand the impact for general queries.

Indexing time. Table IV shows the indexing time of the RLC index and ETC on real-world graphs, where “-” indicates that the method timed out on the graph. Building ETC cannot be completed in 24 hours for real-world graphs (or it runs out of memory) except for the AD graph with the least number of edges. The RLC index for the AD graph can be built in 0.7s,

TABLE IV
INDEXING TIME (IT) AND INDEX SIZE (IS).

Dataset	RLC Index		ETC	
	IT (s)	IS (MB)	IT (s)	IS (MB)
AD	0.7	1.9	2216.1	2798.7
EP	22.6	29.3	-	-
TW	8.1	93.5	-	-
WN	33.1	122.6	-	-
WS	53.5	173.9	-	-
WG	101.3	403.6	-	-
WT	812.9	607.1	-	-
WB	167.1	474.2	-	-
WH	3707.2	1319.1	-	-
PR	3104.1	1212.6	-	-
SO	57072.5	844.2	-	-
LJ	18240.9	6248.1	-	-
WF	51338.7	6467.9	-	-

leading to a four-orders-of-magnitude improvement over ETC. The indexing time improvement of the RLC index over ETC is mainly due to the pruning rules that skip vertices in graph traversals when building the RLC index. Thus, this experiment also shows the significant impact of pruning rules in terms of indexing time. The indexing time of the RLC index for the first 10 graphs is at most 1 hour. The last three graphs, *i.e.*, the SO graph, the LJ graph, and the WF graph are more challenging than the others, not only because they have more vertices and edges, but also because they have a larger number of loops and triangles, as shown in Table III. The SO graph has the longest indexing time due to its highly dense and cyclic character, *i.e.*, it has 15M loops and 114M triangles. Although the WF graph has much fewer loops than the SO graph, it contains 30B triangles. Consequently, the indexing time of the WF graph is at the same order of magnitude as the one of the SO graph. While it has more vertices, triangles, and edge labels than the SO graph, the LJ graph requires a lower indexing time.

Index size. Table IV shows the size of the RLC index for real-world graphs. The size of the RLC index is much smaller than the size of ETC for the AD graph (that is the only graph ETC can be built within 24 hours). The index size improvement is because of not only the indexing schema of the RLC index but also the application of pruning rules that can avoid recording redundant index entries. The effectiveness of pruning rules can also be observed between the PR graph and the SO graph. Although the SO graph is larger than the PR graph in terms of the number of vertices and the number of edges, the index size of the SO graph is smaller than the index size of the PR graph. Thus, this experiment also demonstrates the significant impact of pruning rules in terms of index size.

Query time. Fig. 3 shows the execution time of the true-

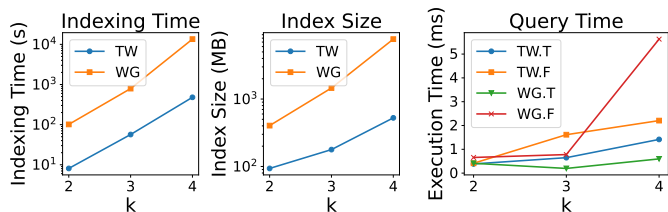


Fig. 4. RLC index performance with different recursive k values.

query set and the false-query set. In general, the execution time of a query set of 1000 queries using the RLC index is around 1 millisecond for all the graphs in Table III, except for the WF graph (that has the largest number of edges) for which around 2 milliseconds. Query execution using BFS times out for the true-queries on both the SO graph and the WF graph, and for the false queries on the WF graph. In addition, we only report the query time of ETC for the AD graph as ETC cannot be built for all the other graphs. As shown in Fig. 3, the RLC index shows an up to six-orders-of-magnitude improvement over BFS and four-orders-of-magnitude improvement over BiBFS. The query time using the RLC index is slightly faster than ETC, because the larger number of reachable pairs of vertices recorded in ETC leads to a slight overhead of checking for the reachability between a source s and a target t and also the existence of a minimum repeat of a path from s to t .

Different recursive k values. We analyze the impact of the recursive k values on real-world graphs. For space reasons, we have focused on the TW and WG datasets, *i.e.*, graphs from Twitter and Google. The indexing results are shown in Figure 4, where we also report the query time of 1000 true-queries and 1000 false-queries with a recursive concatenation of 2, 3, and 4 labels. As expected, the indexing time and index size of the RLC index increase when the recursive k value increases, and the larger index size with more path constraints can also lead to the increase of query time. The main reason is due to the fact that the number of path-constraints (kernels) will exponentially grow as the increasing of k , and the kernel-search phase of the indexing algorithm needs to take into consideration the potential path-constraints that exist in the graph. Notice that the increasing rate of index size is much slower than the increasing rate of indexing time, which means that the number of paths that can satisfy path constraints of recursive concatenation is not increasing as the increase of the length of the constraints, *i.e.*, the recursive k value. The main reason is that in real-world graphs only a few labels have a large number of occurrence (the Zipfian distribution) as shown in existing benchmarks *e.g.*, gMark [78]. Consequently, a long concatenation of edge labels cannot repeat frequently due to the lack of desired labels. Thus, RLC queries with a large recursive k value may not need indexing.

Summary and outlook. Our indexing algorithm designed for an arbitrary recursive k value can efficiently build the RLC index for processing practical RLC queries that are difficult to evaluate in modern graph query engines. Specifically, the

recursive concatenation length of RLC queries in recent open-source query logs [27] is not larger than 2 and such practical RLC queries appear quite often in time-out logs. As shown in Table IV and Fig. 3, the RLC index with a recursive k of 2 can be efficiently built even for large and highly dense real-world graphs, and can significantly improve the processing time of such timed out queries in practice.

B. Impact of Graph Characteristics

In this section, we focus on analyzing the performance of the RLC index on synthetic graphs (ER-graphs and BA-graphs) with different characteristics, namely average degree, label set size, and the number of vertices. The recursive k value is set to 2 in this section, and more experiments about different k values on synthetic graphs are presented in our technical report [39]. We generate for each graph a query set of 1000 true-queries and a query set of 1000 false-queries, which are referred to as $ER.T$ and $ER.F$ for an ER-graph, and $BA.T$ and $BA.F$ for a BA-graph.

1) Impact of label set size and average degree: In this experiment, we use BA-graphs and ER-graphs with 1M vertices, and we vary the average degree d in (2, 3, 4, 5), and label set size $|\mathbb{L}|$ in (8, 12, 16, 20, 24, 28, 32, 36), *e.g.*, a graph with $d = 5$ and $|\mathbb{L}| = 16$ has 1M vertices, 5M edges, and 16 distinct edge labels. We aim at analyzing indexing time, index size, and query time of the RLC index as the increase of d and $|\mathbb{L}|$. The experimental results for ER-graphs and BA-graphs are reported in Fig. 5. We discuss the results below.

Indexing time. We observe in Fig. 5 that the indexing time for the used ER-graphs and BA-graphs with a fixed d shows a linear increase (for most cases) as $|\mathbb{L}|$ increases. This can be understood as follows. When $|\mathbb{L}|$ increases, the number of possible minimum repeats increases, requiring more time for the kernel-search phase of a KBS in the indexing algorithm to traverse the graph and generate potential kernel candidates, resulting in more kernel-BFS executions. The total number of possible minimum repeats can be quadratic in $|\mathbb{L}|$ in the worst case when the input parameter k of the RLC index is two. Furthermore, because there are more edges to traverse, the indexing time for both ER-graphs and BA-graphs with a fixed $|\mathbb{L}|$ increases linearly as d increases.

Index size. As illustrated in Fig. 5, an increase in average degree d can result in a larger index size for both ER-graphs and BA-graphs. The reason is that a vertex s can reach a vertex t through more paths, leading to a higher number of minimum repeats being recorded. As the size of the label set grows, the corresponding impact on ER-graphs is different from the one on BA-graphs. Specifically, the increase is negligible for ER-graphs with a small d , *e.g.*, 2, and becomes more noticeable for ER-graphs with a large d , *e.g.*, 5. For any d , however, we see a clear linear increase in index size with the growth in $|\mathbb{L}|$ for BA-graphs. This is because a BA-graph comprises a complete sub-graph, and vertices inside the complete sub-graph have higher degrees. Therefore, the KBSs executed from such vertices can create more index entries as $|\mathbb{L}|$ grows, as it can reach other vertices through paths with more distinct

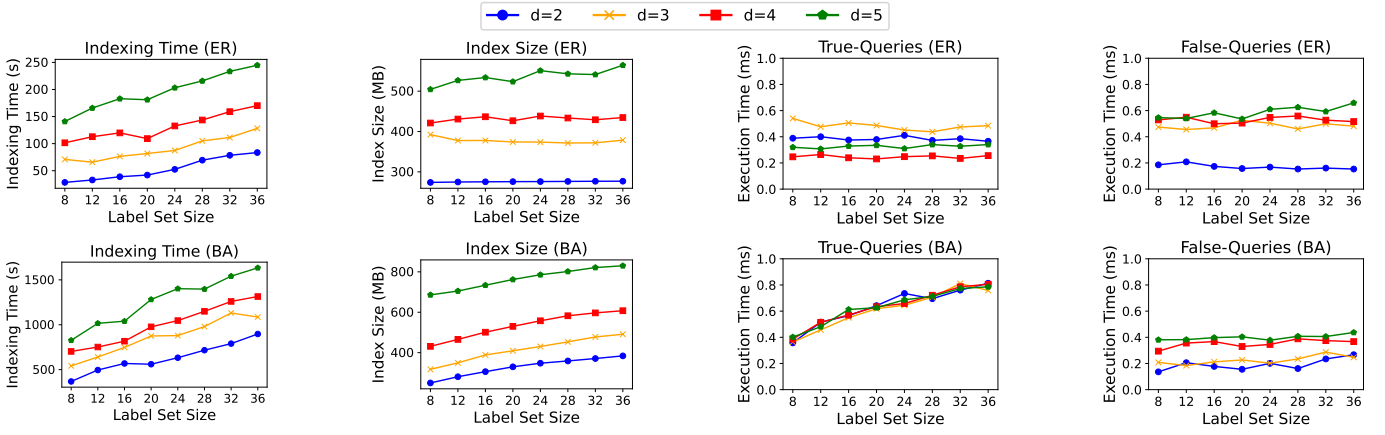


Fig. 5. Indexing time, index size, and execution time for graphs with $|V| = 1M$, varying d , and varying $|\mathbb{L}|$.

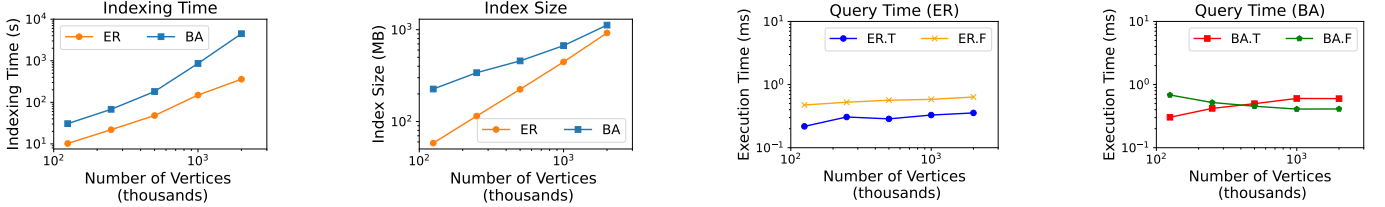


Fig. 6. Indexing time, index size, and query execution time for graphs with $d = 5$, $|\mathbb{L}| = 16$, and varying $|V|$.

minimum repeats. However, because of the uniform degree distribution, the increase in the number of minimum repeats of paths from a vertex s to a vertex t due to an increase in $|\mathbb{L}|$ is not significant for ER-graphs when d is small, but the corresponding impact becomes stronger when d is larger.

Query time. As shown in Fig. 5 the growth of $|\mathbb{L}|$ has a different impact on query time. More precisely, when $|\mathbb{L}|$ rises, the execution time of both true- and false-queries for ER-graphs remains steady. When it comes to BA-graphs, increasing $|\mathbb{L}|$ can lead to a minor boost in true-query execution time but has almost no impact on false query execution time. This is due to the fact that the vertices in the complete sub-graph of a BA-graph can reach much more vertices than the vertices outside the complete sub-graph in the BA-graph, leading to a skew in the distribution of vertices in index entries, *i.e.*, many index entries have the same vertex. Furthermore, when $|\mathbb{L}|$ grows, the number of minimum repeats also increases, which makes the skew higher. As a result, processing true-queries will encounter situations where the query algorithm searches for a particular minimum repeat in a significant number of index entries with the same vertex. However, for false-queries, the query result can be returned instantly if there are no index entries with the same vertex. Fig. 5 also shows that d has a negligible impact on the execution time of the true-queries on the BA graph. The main reason is that the number of index entries for some vertices in the BA graph does not increase significantly w.r.t the increase of d because of the skew in the degree distribution and a fixed $|\mathbb{L}|$.

2) *Scalability:* In this experiment, we use BA-graphs and ER-graphs with average degree 5, 16 edge labels, and vary the number of vertices in (125K, 250K, 500K, 1M, 2M). The goal is to analyze the scalability of the RLC index in terms of

$|V|$. The results of indexing time, index size, and query time for both ER-graphs and BA-graphs are reported in Fig. 6.

Indexing time and index size. Fig. 6 shows that indexing time and index size grows with the increase of the number of vertices. The main reason is that graphs with more vertices require more KBS iterations, which increases indexing time and also the number of index entries. We observe that indexing BA-graphs is more expensive than indexing ER-graphs because of the presence of a complete sub-graph in the former. We also observe that the different topological structures between BA-graphs and ER-graphs have different impacts on the increasing rate of index size w.r.t $|V|$. The uniform degree distribution in ER-graphs makes the index size increase with a sharper rate than the one of BA-graphs, because BA-graphs contain a significant number of vertices of high degree, which is also growing as $|V|$ increases, and the indexing algorithm starts building the index from these vertices such that index entries containing these vertices can be leveraged to prune redundant index entries that need to be inserted later on.

Query time. Fig. 6 shows that query time on the ER-graphs and true-query time on the BA-graphs slightly increase as the number of vertices grows, as expected on larger graphs. In addition, the false-query time is higher than the true-query time on ER-graphs, and the true-query time is higher than the false-query time on BA-graphs. We interpret the results as follows. Given a query (s, t, L^+) on a graph which has a uniform distribution in vertex degree, the index entries (v, mr) in both $\mathcal{L}_{out}(s)$ and $\mathcal{L}_{in}(t)$ also have a uniform distribution in terms of v . Thus, the query algorithm (based on merge join) need to perform an exhaustive search in $\mathcal{L}_{out}(s)$ and $\mathcal{L}_{in}(t)$ to find index entries with the same vertex v , which

results in false-queries taking longer time to execute than true-queries. However, when the distribution of vertex degree is skewed in BA-graphs, the index entries in $\mathcal{L}_{out}(s)$ and $\mathcal{L}_{in}(t)$ are dominated by vertices of high degree. In addition, as there are many paths passing through high-degree vertices with distinct minimum repeats, the number of index entries with such vertices is also large. Thus, the query algorithm can perform a faster search for false-queries than true-queries, because the number of distinct vertices in both $\mathcal{L}_{out}(s)$ and $\mathcal{L}_{in}(t)$ is not large. Notice that the number of vertices of high degree also increases as $|V|$ increases. For false-queries in BA-graphs, the dominance of index entries with high-degree vertices becomes stronger as $|V|$ increases. Consequently, the number of index entries in $\mathcal{L}_{out}(s)$ and $\mathcal{L}_{in}(t)$ with distinct vertex id can decrease, which makes the merge-join algorithm execute faster. The query algorithm, on the other hand, needs to select a specific mr among index entries with vertex u of a high degree to process true-queries, which takes more time.

C. Comparison with Existing Systems

In this section, we focus on evaluating how much improvement the RLC index can provide over mainstream graph processing systems. We recall that many current graph query engines fail to evaluate RLC queries, thus we focused on three systems that are able to evaluate these queries on property graphs and RDF graphs. In order to show how the index performs with varying types of queries, such as longer concatenations or path queries frequently occurring in real-world query logs, we consider the following queries: **Q1** being a single label under the Kleene plus a^+ ; **Q2** consisting of a concatenation of length 2 $(a \circ b)^+$; **Q3** having the expression $(a \circ b \circ c)^+$ thus a concatenation of length 3. In this case, we build the index with $k = 3$ to support all the three RLC queries, especially Q3 having the longest concatenation. For the sake of completeness, an extended reachability query with the constraint $a^+ \circ b^+$ is also included in this experiment, referred to as **Q4**. We have employed Q4 to study the generality and applicability of the RLC index to a wide range of regular path queries in real-world graph query logs [27]. To deal with this query, we use the RLC index in combination with an online traversal to continuously check whether intermediately visited vertices can satisfy the path constraint.

Three graph engines, including commercial and open-sourced ones, used in the experiments have been selected as representatives of the few available graph engines capable of evaluating RLC queries. We do not reveal the identity of all the systems as some are proprietary and we cannot release performance data. Anonymized engines are referred to as Sys1 and Sys2 in the results, and the third one is Virtuoso Open-Source Edition (v7.2.6.3233). For the systems that support multi-threaded query evaluation, we set the system to single-threaded to ensure a fair comparison with our approach, which is single-threaded only. For Virtuoso, we disabled the transaction logging to avoid additional overheads, and configured it to work entirely in memory by setting the maximum amount of memory for transitive queries to the available memory of

TABLE V
SPEED-UPS (SU) AND WORKLOAD SIZE BREAK-EVEN POINTS (BEP) OF THE RLC INDEX OVER GRAPH ENGINES.

Sys.	RLC Query						Extended Query	
	Q1		Q2		Q3		Q4	
	SU	BEP	SU	BEP	SU	BEP	SU	BEP
Sys1	1200x	84100	10400x	34000	18400x	9400	34000x	300
Sys2	3000x	34900	202000x	1700	1300000x	130	104000x	98
Virtuoso	597x	180000	4900x	71700	38100000x	5	-	-

the server. Note that these systems have their own indexes by default, which we leave the configuration unchanged, e.g., Virtuoso 7 has column-wise indexes by default. These indexes are not suitable for RLC queries, as confirmed by our comparative analysis.

We use the WN graph as a representative of real-world graphs, which has a moderate number of vertices and edges, along with a sufficient number of cycles to evaluate. We build one RLC index with the parameter $k = 3$ for the WN graph and use it to process all four queries. In this case, the RLC index of the WN graph can be built in 5.9 minutes and takes up 821 megabytes. We run each query using each system within the 10-minutes time limit, and we repeat the execution 20 times and report the median of the query execution time. We use two metrics to evaluate the improvement of the RLC index, namely speed-ups (SU) and the workload size break-even points (BEP). SU shows the query time improvement of the RLC index over an included graph system. BEP indicates the number of queries that make the indexing time of the RLC index pay off. Table V shows the results, where “-” indicates that the query execution of this system timed out. The RLC index shows that using a single index lookup can gain significant speed-ups over included systems for processing Q3, which has the longest concatenation under the Kleene plus. We use the BEP value to understand the amortized cost of using the RLC index to accelerate Q3 processing. In particular, executing Q3 130 times on Sys2 is equivalent to the time it takes to build and query the RLC index the same number of times. The RLC index can also significantly improve the execution time for Q1, Q2, and Q4.

VII. CONCLUSION

In this paper, we introduced RLC queries, reachability queries with a path constraint based on the Kleene plus over a concatenation of edge labels. RLC queries are becoming relevant in real-world applications, such as social networks, financial transactions and SPARQL endpoints. In order to efficiently process RLC queries online, we propose the RLC index, the first reachability index suitable for these queries. We design an indexing algorithm with pruning rules to not only accelerate the index construction but also remove redundant index entries. Our comprehensive experimental study demonstrates that the RLC index allows to strike a balance between online processing (full online traversals) and offline computation (a materialized transitive closure). Last but not least, our open-source implementation of the RLC index can significantly speed up the processing of RLC queries in current mainstream graph engines.

REFERENCES

- [1] M. E. J. Newman, *Networks: An Introduction*. Oxford University Press, 2010.
- [2] S. Sahu, A. Mhedhbi, S. Salihoglu, J. Lin, and M. T. Özsu, “The ubiquity of large graphs and surprising challenges of graph processing,” *Proc. VLDB Endow.*, vol. 11, no. 4, p. 420–431, dec 2017. [Online]. Available: <https://doi.org/10.1145/3186728.3164139>
- [3] R. Agrawal, A. Borgida, and H. V. Jagadish, “Efficient management of transitive relationships in large data and knowledge bases,” in *Proceedings of the 1989 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD ’89. New York, NY, USA: Association for Computing Machinery, 1989, p. 253–262. [Online]. Available: <https://doi.org/10.1145/67544.66950>
- [4] H. V. Jagadish, “A compression technique to materialize transitive closure,” *ACM Trans. Database Syst.*, vol. 15, no. 4, p. 558–598, Dec. 1990. [Online]. Available: <https://doi.org/10.1145/99935.99944>
- [5] E. Cohen, E. Halperin, H. Kaplan, and U. Zwick, “Reachability and distance queries via 2-hop labels,” in *Proceedings of the Thirteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, ser. SODA ’02. USA: Society for Industrial and Applied Mathematics, 2002, p. 937–946.
- [6] Haixun Wang, Hao He, Jun Yang, P. S. Yu, and J. X. Yu, “Dual labeling: Answering graph reachability queries in constant time,” in *22nd International Conference on Data Engineering (ICDE’06)*, 2006, pp. 75–75.
- [7] Y. Chen and Y. Chen, “An efficient algorithm for answering graph reachability queries,” in *2008 IEEE 24th International Conference on Data Engineering*, 2008, pp. 893–902.
- [8] R. Jin, Y. Xiang, N. Ruan, and H. Wang, “Efficiently answering reachability queries on very large directed graphs,” in *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD ’08. New York, NY, USA: Association for Computing Machinery, 2008, p. 595–608. [Online]. Available: <https://doi.org/10.1145/1376616.1376677>
- [9] R. Jin, Y. Xiang, N. Ruan, and D. Fuhry, “3-hop: A high-compression indexing scheme for reachability query,” in *Proceedings of the 2009 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD ’09. New York, NY, USA: Association for Computing Machinery, 2009, p. 813–826. [Online]. Available: <https://doi.org/10.1145/1559845.1559930>
- [10] J. Cheng, S. Huang, H. Wu, and A. W.-C. Fu, “Tf-label: A topological-folding labeling scheme for reachability querying in a large graph,” in *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD ’13. New York, NY, USA: Association for Computing Machinery, 2013, p. 193–204. [Online]. Available: <https://doi.org/10.1145/2463676.2465286>
- [11] R. Jin and G. Wang, “Simple, fast, and scalable reachability oracle,” *Proc. VLDB Endow.*, vol. 6, no. 14, p. 1978–1989, Sep. 2013. [Online]. Available: <https://doi.org/10.14778/2556549.2556578>
- [12] L. Chen, A. Gupta, and M. E. Kurul, “Stack-based algorithms for pattern matching on dags,” in *Proceedings of the 31st International Conference on Very Large Data Bases*, ser. VLDB ’05. VLDB Endowment, 2005, p. 493–504.
- [13] S. Trißl and U. Leser, “Fast and practical indexing and querying of very large graphs,” in *Proceedings of the 2007 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD ’07. New York, NY, USA: Association for Computing Machinery, 2007, p. 845–856. [Online]. Available: <https://doi.org/10.1145/1247480.1247573>
- [14] H. Yildirim, V. Chaoji, and M. J. Zaki, “Grail: Scalable reachability index for large graphs,” *Proc. VLDB Endow.*, vol. 3, no. 1–2, p. 276–284, Sep. 2010. [Online]. Available: <https://doi.org/10.14778/1920841.1920879>
- [15] S. Seufert, A. Anand, S. Bedathur, and G. Weikum, “Ferrari: Flexible and efficient reachability range assignment for graph indexing,” in *2013 IEEE 29th International Conference on Data Engineering (ICDE)*, 2013, pp. 1009–1020.
- [16] R. R. Veloso, L. Cerf, W. Meira, and M. J. Zaki, “Reachability queries in very large graphs: A fast refined online search approach,” in *EDBT*, 2014.
- [17] H. Wei, J. X. Yu, C. Lu, and R. Jin, “Reachability querying: An independent permutation labeling approach,” *Proc. VLDB Endow.*, vol. 7, no. 12, p. 1191–1202, Aug. 2014. [Online]. Available: <https://doi.org/10.14778/2732977.2732992>
- [18] J. Su, Q. Zhu, H. Wei, and J. X. Yu, “Reachability querying: Can it be even faster?” *IEEE Transactions on Knowledge and Data Engineering*, vol. 29, no. 3, pp. 683–697, 2017.
- [19] A. D. Zhu, W. Lin, S. Wang, and X. Xiao, “Reachability queries on large dynamic graphs: A total order approach,” in *Proceedings of the 2014 ACM International Conference on Management of Data*, ser. SIGMOD ’14. New York, NY, USA: Association for Computing Machinery, 2014, p. 1323–1334. [Online]. Available: <https://doi.org/10.1145/2588555.2612181>
- [20] J. Cai and C. K. Poon, “Path-hop: Efficiently indexing large graphs for reachability queries,” in *Proceedings of the 19th ACM International Conference on Information and Knowledge Management*, ser. CIKM ’10. New York, NY, USA: Association for Computing Machinery, 2010, p. 119–128. [Online]. Available: <https://doi.org/10.1145/1871437.1871457>
- [21] Y. Yano, T. Akiba, Y. Iwata, and Y. Yoshida, “Fast and scalable reachability queries on graphs by pruned labeling with landmarks and paths,” in *Proceedings of the 22nd ACM International Conference on Information and Knowledge Management*, ser. CIKM ’13. New York, NY, USA: Association for Computing Machinery, 2013, p. 1601–1606. [Online]. Available: <https://doi.org/10.1145/2505515.2505724>
- [22] R. Jin, H. Hong, H. Wang, N. Ruan, and Y. Xiang, “Computing label-constraint reachability in graph databases,” in *Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD ’10. New York, NY, USA: Association for Computing Machinery, 2010, p. 123–134. [Online]. Available: <https://doi.org/10.1145/1807167.1807183>
- [23] L. Zou, K. Xu, J. X. Yu, L. Chen, Y. Xiao, and D. Zhao, “Efficient processing of label-constraint reachability queries in large graphs,” *Information Systems*, vol. 40, pp. 47–66, 2014. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0306437913001257>
- [24] L. D. Valstar, G. H. Fletcher, and Y. Yoshida, “Landmark indexing for evaluation of label-constrained reachability queries,” in *Proceedings of the 2017 ACM International Conference on Management of Data*, ser. SIGMOD ’17. New York, NY, USA: Association for Computing Machinery, 2017, p. 345–358. [Online]. Available: <https://doi.org/10.1145/3035918.3035955>
- [25] Y. Peng, Y. Zhang, X. Lin, L. Qin, and W. Zhang, “Answering billion-scale label-constrained reachability queries within microsecond,” *Proc. VLDB Endow.*, vol. 13, no. 6, p. 812–825, Feb. 2020. [Online]. Available: <https://doi.org/10.14778/3380750.3380753>
- [26] Y. Chen and G. Singh, “Graph indexing for efficient evaluation of label-constrained reachability queries,” *ACM Trans. Database Syst.*, 2021.
- [27] A. Bonifati, W. Martens, and T. Timm, “Navigating the maze of wikidata query logs,” in *The World Wide Web Conference*, ser. WWW ’19. New York, NY, USA: Association for Computing Machinery, 2019, p. 127–138. [Online]. Available: <https://doi.org/10.1145/3308558.3313472>
- [28] “Neo4j,” <http://www.opencypher.org>.
- [29] A. Deutsch, Y. Xu, M. Wu, and V. E. Lee, “Tigergraph: A native mpp graph database,” *ArXiv*, vol. abs/1901.08248, 2019.
- [30] “Graph query language gql standard,” <https://www.gqlstandards.org/>.
- [31] “Apache tinkerepop data system providers,” <https://tinkerpop.apache.org/providers.html>.
- [32] B. R. Bebee, D. Choi, A. Gupta, A. Gutmans, A. Khandelwal, Y. Kiran, S. Mallidi, B. McGaughy, M. Personick, K. J. Rajan, S. Rondelli, A. Ryazanov, M. Schmidt, K. Sengupta, B. B. Thompson, D. Vaidya, and S. X. Wang, “Amazon neptune: Graph data management in the cloud,” in *SEMWEB*, 2018.
- [33] O. van Rest, S. Hong, J. Kim, X. Meng, and H. Chafi, “Pgql: A property graph query language,” in *Proceedings of the Fourth International Workshop on Graph Data Management Experiences and Systems*, ser. GRADES ’16. New York, NY, USA: Association for Computing Machinery, 2016. [Online]. Available: <https://doi.org/10.1145/2960414.2960421>
- [34] S. Hong, S. Depner, T. Manhardt, J. Van Der Lugt, M. Verstraaten, and H. Chafi, “Pgx.d: a fast distributed graph processing engine,” in *SC ’15: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, 2015, pp. 1–12.
- [35] V. Trigonakis, J.-P. Lozi, N. P. Roth, I. Psaroudakis, A. Delamare, V. Haprian, C. Iorgulescu, P. Koupy, J. Lee, S. Hong, and H. Chafi, “adfs: An almost depth-first-search distributed graph-querying system,” in *2021 USENIX Annual Technical Conference (USENIX ATC 21)*. USENIX Association, Jul. 2021, pp. 209–224. [Online]. Available: <https://www.usenix.org/conference/atc21/presentation/trigonakis>
- [36] “Virtuoso,” <http://vos.openlinksw.com/owiki/wiki/VOS>.

- [37] "Apache jena," <https://jena.apache.org>.
- [38] S. Sakr, A. Bonifati, H. Voigt, A. Iosup, K. Ammar, R. Angles, W. G. Aref, M. Arenas, M. Besta, P. A. Boncz, K. Daudjee, E. D. Valle, S. Dumbrava, O. Hartig, B. Haslhofer, T. Hegeman, J. Hidders, K. Hose, A. Iammitchi, V. Kalavri, H. Kapp, W. Martens, M. T. Özsu, E. Peukert, S. Plantikow, M. Ragab, M. Ripeanu, S. Salihoglu, C. Schulz, P. Selmer, J. F. Sequeda, J. Shinavier, G. Szárnyas, R. Tommasini, A. Tumeo, A. Uta, A. L. Varbanescu, H. Wu, N. Yakovets, D. Yan, and E. Yoneki, "The future is big graphs: a community view on graph processing systems," *Commun. ACM*, vol. 64, no. 9, pp. 62–71, 2021.
- [39] C. Zhang, A. Bonifati, H. Kapp, V. I. Haprian, and J.-P. Lozi, "A reachability index for recursive label-concatenated graph queries," 2022. [Online]. Available: <https://arxiv.org/abs/2203.08606>
- [40] R. Angles, M. Arenas, P. Barceló, A. Hogan, J. Reutter, and D. Vrgoč, "Foundations of modern query languages for graph databases," *ACM Comput. Surv.*, vol. 50, no. 5, Sep. 2017. [Online]. Available: <https://doi.org/10.1145/3104031>
- [41] J. R. Ullmann, "An algorithm for subgraph isomorphism," *Journal of the ACM (JACM)*, vol. 23, pp. 31–42, 1976.
- [42] L. Cordella, P. Foggia, C. Sansone, and M. Vento, "A (sub)graph isomorphism algorithm for matching large graphs," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 26, no. 10, pp. 1367–1372, 2004.
- [43] A. Jüttner and P. Madarasi, "Vf2++—an improved subgraph isomorphism algorithm," *Discrete Applied Mathematics*, vol. 242, pp. 69–81, 2018, computational Advances in Combinatorial Optimization. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0166218X18300829>
- [44] V. Carletti, P. Foggia, A. Saggese, and M. Vento, "Challenging the time complexity of exact subgraph isomorphism for huge and dense graphs with vf3," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 4, pp. 804–818, 2018.
- [45] H. Shang, Y. Zhang, X. Lin, and J. X. Yu, "Taming verification hardness: An efficient algorithm for testing subgraph isomorphism," *Proc. VLDB Endow.*, vol. 1, no. 1, p. 364–375, aug 2008. [Online]. Available: <https://doi.org/10.14778/1453856.1453899>
- [46] S. Zhang, S. Li, and J. Yang, "Gaddi: Distance index based subgraph matching in biological networks," in *Proceedings of the 12th International Conference on Extending Database Technology: Advances in Database Technology*, ser. EDBT '09. New York, NY, USA: Association for Computing Machinery, 2009, p. 192–203. [Online]. Available: <https://doi.org/10.1145/1516360.1516384>
- [47] P. Zhao and J. Han, "On graph query optimization in large networks," *Proc. VLDB Endow.*, vol. 3, no. 1–2, p. 340–351, sep 2010. [Online]. Available: <https://doi.org/10.14778/1920841.1920887>
- [48] H. He and A. K. Singh, "Graphs-at-a-time: Query language and access methods for graph databases," in *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD '08. New York, NY, USA: Association for Computing Machinery, 2008, p. 405–418. [Online]. Available: <https://doi.org/10.1145/1376616.1376660>
- [49] W.-S. Han, J. Lee, and J.-H. Lee, "Turbo_{iso}: Towards ultrafast and robust subgraph isomorphism search in large graph databases," in *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD '13. New York, NY, USA: Association for Computing Machinery, 2013, p. 337–348. [Online]. Available: <https://doi.org/10.1145/2463676.2465300>
- [50] X. Ren and J. Wang, "Exploiting vertex relationships in speeding up subgraph isomorphism over large graphs," *Proc. VLDB Endow.*, vol. 8, no. 5, p. 617–628, jan 2015. [Online]. Available: <https://doi.org/10.14778/2735479.2735493>
- [51] F. Bi, L. Chang, X. Lin, L. Qin, and W. Zhang, "Efficient subgraph matching by postponing cartesian products," in *Proceedings of the 2016 International Conference on Management of Data*, ser. SIGMOD '16. New York, NY, USA: Association for Computing Machinery, 2016, p. 1199–1214. [Online]. Available: <https://doi.org/10.1145/2882903.2915236>
- [52] C. R. Rivero and H. M. Jamil, "Efficient and scalable labeled subgraph matching using sgmatch," *Knowl. Inf. Syst.*, vol. 51, no. 1, p. 61–87, apr 2017. [Online]. Available: <https://doi.org/10.1007/s10115-016-0968-2>
- [53] B. Bhattarai, H. Liu, and H. H. Huang, "Ceci: Compact embedding cluster index for scalable subgraph matching," in *Proceedings of the 2019 International Conference on Management of Data*, ser. SIGMOD '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 1447–1462. [Online]. Available: <https://doi.org/10.1145/3299869.3300086>
- [54] M. Han, H. Kim, G. Gu, K. Park, and W.-S. Han, "Efficient subgraph matching: Harmonizing dynamic programming, adaptive matching order, and failing set together," in *Proceedings of the 2019 International Conference on Management of Data*, ser. SIGMOD '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 1429–1446. [Online]. Available: <https://doi.org/10.1145/3299869.3319880>
- [55] S. Sun and Q. Luo, "Subgraph matching with effective matching order and indexing," *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 1, pp. 491–505, 2022.
- [56] S. Sun, X. Sun, Y. Che, Q. Luo, and B. He, "Rapidmatch: A holistic approach to subgraph query processing," *Proc. VLDB Endow.*, vol. 14, no. 2, p. 176–188, oct 2020. [Online]. Available: <https://doi.org/10.14778/3425879.3425888>
- [57] H. Kim, S. Min, K. Park, X. Lin, S.-H. Hong, and W.-S. Han, "Idar: Fast supergraph search using dag integration," *Proc. VLDB Endow.*, vol. 13, no. 9, p. 1456–1468, may 2020. [Online]. Available: <https://doi.org/10.14778/3397230.3397241>
- [58] J. Lee, W.-S. Han, R. Kasperovics, and J.-H. Lee, "An in-depth comparison of subgraph isomorphism algorithms in graph databases," *Proc. VLDB Endow.*, vol. 6, no. 2, p. 133–144, dec 2012. [Online]. Available: <https://doi.org/10.14778/2535568.2448946>
- [59] S. Sun and Q. Luo, "In-memory subgraph matching: An in-depth study," in *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 1083–1098. [Online]. Available: <https://doi.org/10.1145/3318464.3380581>
- [60] J. X. Yu and J. Cheng, "Graph reachability queries: A survey," in *Managing and Mining Graph Data*. Springer, 2010, pp. 181–215.
- [61] A. Bonifati, G. Fletcher, H. Voigt, N. Yakovets, and H. V. Jagadish, *Querying Graphs*. Morgan & Claypool Publishers, 2018.
- [62] A. O. Mendelzon and P. T. Wood, "Finding regular simple paths in graph databases," in *Proceedings of the 15th International Conference on Very Large Data Bases*, ser. VLDB '89. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1989, p. 185–193.
- [63] S. Wadhwa, A. Prasad, S. Ranu, A. Bagchi, and S. Bedathur, "Efficiently answering regular simple path queries on large labeled networks," in *Proceedings of the 2019 International Conference on Management of Data*, ser. SIGMOD '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 1463–1480. [Online]. Available: <https://doi.org/10.1145/3299869.3319882>
- [64] G. Bagan, A. Bonifati, and B. Groz, "A trichotomy for regular simple path queries on graphs," in *Proceedings of the 32nd ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, ser. PODS '13. New York, NY, USA: Association for Computing Machinery, 2013, p. 261–272. [Online]. Available: <https://doi.org/10.1145/2463664.2467795>
- [65] P. T. Wood, "Query languages for graph databases," *SIGMOD Rec.*, vol. 41, no. 1, p. 50–60, Apr. 2012. [Online]. Available: <https://doi.org/10.1145/2206869.2206879>
- [66] P. Barceló Baeza, "Querying graph databases," in *Proceedings of the 32nd ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, ser. PODS '13. New York, NY, USA: Association for Computing Machinery, 2013, p. 175–188. [Online]. Available: <https://doi.org/10.1145/2463664.2465216>
- [67] A. Koschmieder and U. Leser, "Regular path queries on large graphs," in *Proceedings of the 24th International Conference on Scientific and Statistical Database Management*, ser. SSDBM'12. Berlin, Heidelberg: Springer-Verlag, 2012, p. 177–194. [Online]. Available: https://doi.org/10.1007/978-3-642-31235-9_12
- [68] W. Fan, J. Li, S. Ma, N. Tang, and Y. Wu, "Adding regular expressions to graph reachability and pattern queries," in *2011 IEEE 27th International Conference on Data Engineering*, 2011, pp. 39–50.
- [69] G. Fletcher, J. J. Peters, and A. Poullovassilis, "Efficient regular path query evaluation using path indexes," in *EDBT*, 2016.
- [70] J. Kuijpers, G. Fletcher, T. Lindaaker, and N. Yakovets, "Path indexing in the cypher query pipeline," in *EDBT*, 2021.
- [71] W. Fan, X. Wang, and Y. Wu, "Performance guarantees for distributed reachability queries," *Proc. VLDB Endow.*, vol. 5, no. 11, p. 1304–1316, Jul. 2012. [Online]. Available: <https://doi.org/10.14778/2350229.2350248>
- [72] A. Pacaci, A. Bonifati, and M. T. Özsu, "Regular path query evaluation on streaming graphs," in *Proceedings of the 2020 ACM SIGMOD*

International Conference on Management of Data, ser. SIGMOD '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 1415–1430. [Online]. Available: <https://doi.org/10.1145/3318464.3389733>

- [73] A. Pacaci, A. Bonifati, and M. T. Özsu, “Evaluating complex queries on streaming graphs,” 2021.
- [74] A. S. LaPaugh and C. H. Papadimitriou, “The even-path problem for graphs and digraphs,” *Networks*, vol. 14, pp. 507–513, 1984.
- [75] D. Knuth, J. H. Morris, and V. Pratt, “Fast pattern matching in strings,” *SIAM J. Comput.*, vol. 6, pp. 323–350, 1977.
- [76] J. Leskovec and A. Krevl, “SNAP Datasets: Stanford large network dataset collection,” <http://snap.stanford.edu/data>, Jun. 2014.
- [77] J. Kunegis, “Konect: The koblenz network collection,” in *Proceedings of the 22nd International Conference on World Wide Web*, ser. WWW '13 Companion. New York, NY, USA: Association for Computing Machinery, 2013, p. 1343–1350. [Online]. Available: <https://doi.org/10.1145/2487788.2488173>
- [78] G. Bagan, A. Bonifati, R. Ciucanu, G. H. L. Fletcher, A. Lemay, and N. Advokaat, “gMark: Schema-driven generation of graphs and queries,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 29, no. 4, pp. 856–869, 2017.
- [79] D. Michail, J. Kinable, B. Naveh, and J. V. Sichi, “JgraphT—a java library for graph data structures and algorithms,” *ACM Trans. Math. Softw.*, vol. 46, no. 2, may 2020. [Online]. Available: <https://doi.org/10.1145/3381449>

APPENDIX A PROOF OF THEOREM 1

A proof of Theorem 1 is provided in this section. Before that, we first show the kernel of a label sequence is unique in Lemma 2 that will be used in the proof of Theorem 1.

Lemma 2: If L has a kernel, then the kernel is unique.

Proof: The proof is based on induction. The statement is if a label sequence L of length n has a kernel, then the kernel is unique. It is trivial to prove the initial case $|L| = 2$. Assuming the case n is true, then we show the case $n+1$ is also true. Let $|L| = n+1$. We use \bar{L} to denote the label sequence obtained by removing the last label of L , i.e., $|\bar{L}| = n$. We prove the case $n+1$ below.

Assuming L can have two different kernels L_1 and L_2 , and $L_1 \neq L_2$, then $L = (L_1)^{h_1} \circ L'_1$, $h_1 \geq 2$ and $L = (L_2)^{h_2} \circ L'_2$, $h_2 \geq 2$. If $|L'_1| = 0$ and $|L'_2| = 0$, then L has two MRs, which is contradictory (Lemma 1). If $|L'_1| \neq 0$ and $|L'_2| \neq 0$, then \bar{L} has kernels L_1 and L_2 , which contradicts to the case n . The remaining cases are that only one of L'_1 and L'_2 has a length of 0. W.l.o.g. consider $|L'_1| = 0$ and $|L'_2| \neq 0$. Given this, if $h_1 > 2$, then \bar{L} also has kernels L_1 and L_2 , which is contradictory. Then we have $h_1 = 2$ and $|L'_1| = 0$, i.e.,

$$L = L_1 \circ L_1. \quad (1)$$

In addition, we have

$$L = (L_2)^{h_2} \circ L'_2, h_2 \geq 2, |L'_2| \neq 0. \quad (2)$$

Let $L = (l_1, \dots, l_{2|L_1|})$ and $|L_1| = a|L_2| + b$, $1 \leq a, b < |L_2|$. According to Equ. (1), we have $l_i = l_{i+|L_1|}$, $1 \leq i \leq |L_1|$ and $l_{i'} = l_{i'-|L_1|}$, $|L_1| < i' \leq 2|L_1|$, which means $l_i = l_{i+a|L_2|+b}$ and $l_{i'} = l_{i'-a|L_2|-b}$. Based on Equ. (2), we have $l_i = l_{i+b}$ and $l_{i'} = l_{i'-b}$. Given this, consider the following two cases: case (i) if $2|L_1| \bmod b = 0$, then $|MR(L_1 \circ L_1)| = b \neq |L_1|$ that contradicts to the fact that L_1 is the unique MR of L ; case (ii) if $2|L_1| \bmod b \neq 0$, then $L = (L_3)^{h_3} \circ L'_3$, where either $|L_3| = b$, $h_3 \geq 2$, and $|L'_3| \neq 0$, or $|L_3| < b$ and $h_3 > 2$.

Note that, in the two sub-cases of case (ii), L'_3 is ϵ , or a proper prefix of L_3 . Therefore, \bar{L} has a kernel L_3 , $|L_3| \leq b$. However, \bar{L} also has a kernel L_2 and $|L_2| > b \geq |L_3|$, which is also a contradiction. ■

Based on Lemma 2, we prove Theorem 1 below.

Proof of Theorem 1: It is not difficult to prove Case 1 and Case 2. We focus on Case 3 below. For ease of presentation, let $\Lambda(u, x) = \Lambda(p(u, x))$ and $\Lambda(x, v) = \Lambda(p(x, v))$.

(Sufficiency) Because $\Lambda(u, x)$ has the kernel L' and the tail L'' , such that $\Lambda(u, x) = (L')^h \circ L''$, $h \geq 2$. Thus, we have $MR(\Lambda(u, x) \circ \Lambda(x, v)) = MR((L')^h \circ L'' \circ \Lambda(x, v)) = L'$, otherwise $MR(L'' \circ \Lambda(x, v)) \neq L'$. In addition, we have $|L'| \leq k$ because $|\Lambda(u, x)| = 2k$. Thus, L' is the k-MR of p .

(Necessity) We show that p does not have a non-empty k-MR in the following two cases.

- Case (i): $\Lambda(u, x)$ does not have a kernel and a tail. Assume that p can have a non-empty k-MR L''' in this case. Because $|L'''| \leq k$ and $|\Lambda(u, x)| = 2k$, then $\Lambda(u, x)$ has a kernel and a tail, which is contradiction to the case definition.
- Case (ii): $\Lambda(u, x)$ has a kernel L' and a tail L'' , but $MR(L'' \circ \Lambda(x, v)) \neq L'$. Assume that p has a non-empty k-MR L''' in this case. Knowing that $|L'''| \leq k$ and $|\Lambda(u, x)| = 2k$, we have $L''' = L'$ because the kernel of $\Lambda(u, x)$ is unique (Lemma 2). Therefore, $MR((L')^h \circ L'' \circ \Lambda(x, v)) = L'$, $h \geq 2$, which means $MR(L'' \circ \Lambda(x, v)) = L'$, that is also a contradiction. ■

APPENDIX B PROOF OF THEOREM 2 AND 3

We prove Theorem 2 and Theorem 3 in this section. Before proceeding further, we first present the following lemmas that will be used to prove the two theorems.

Lemma 3: Given a path $p(s, t)$ having a k-MR L . If the KBS from s can visit t (or the KBS from t can visit s), then the k-MR L of $p(s, t)$ must be recorded in the index.

Proof: If the KBS from s can visit t , then regardless of whether PR1 or PR2 is applied, the k-MR L of $p(s, t)$ must be recorded. ■

Lemma 4: Given two paths $p(s, u)$ and $p(u, t)$ with k-MR L in a graph, where $aid(u) < aid(s)$ and $aid(u) < aid(t)$. We have: if $aid(u) \leq i$, then the k-MR of the path from s to t through vertex u is recorded by Algorithm 2 in the i -th snapshot of the RLC index that is computed after performing KBS from a vertex with access id i .

Proof: The proof is based on induction. It is trivial to prove the initial case $i = 1$. We assume the case $i = n$ is true and prove the case $i = n+1$ below. We only need to show the case $aid(u) = n+1$. Let $p(s, t) = (s, \dots, u, \dots, t)$.

Assuming the backward KBS from u does not visit s . Then PR3 is triggered, such that there exists vertex w , $aid(w) < aid(u)$, and $p(s, w)$ and $p(w, u)$ have the k-MR L . This case can be reduced to the case $i = n$, because the k-MR of path (w, \dots, u, \dots, t) is L and $aid(w) < aid(u) < n+1$. In the same way, we can also prove the case if the forward KBS from u does not visit t .

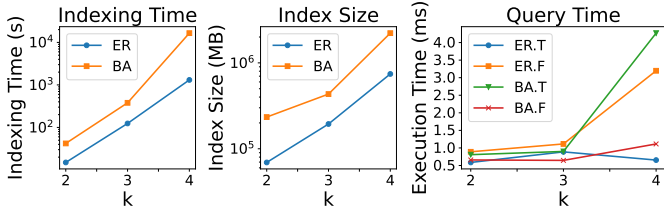


Fig. 7. Evaluation of the RLC index with varying k .

We consider the case that both the backward KBS and the forward KBS from u can visit s and t . For $p(s, u)$, we have the following two cases: Case (1) $\exists(u, L) \in \mathcal{L}_{out}(s)$; Case (2) $\exists(v, L) \in \mathcal{L}_{out}(s)$ and $\exists(v, L) \in \mathcal{L}_{in}(u)$, $aid(v) < aid(u)$. For Case (2), both $p(s, v)$ and $p(v, t)$ have the k-MR L , such that this case can be reduced to the case $i = n$ as $aid(v) < aid(u) = n + 1$. Then we only need to consider Case (1), i.e., $\exists(u, L) \in \mathcal{L}_{out}(s)$. In the same way, for $p(u, t)$, we only need to consider the case $\exists(u, L) \in \mathcal{L}_{in}(s)$. Given this, the k-MR of the path from s to t is recorded by having $(u, L) \in \mathcal{L}_{out}(s)$ and $(u, L) \in \mathcal{L}_{in}(t)$. ■

Lemma 5: Given a path p from s to t with a k-MR L . If the index entry $(t, L) \in \mathcal{L}_{out}(s)$ (or $(s, L) \in \mathcal{L}_{in}(t)$) is pruned because of PR3, then we have one of the following two cases:

- $\exists(s, L) \in \mathcal{L}_{in}(t)$ (or $\exists(t, L) \in \mathcal{L}_{out}(s)$);
- $\exists(v, L) \in \mathcal{L}_{out}(s)$ and $\exists(v, L) \in \mathcal{L}_{in}(t)$, such that $aid(v) < aid(t)$ (or $aid(v) < aid(s)$).

Proof: There exists two cases: $aid(t) \leq aid(s)$ or $aid(t) > aid(s)$. We prove the case of $aid(t) \leq aid(s)$. The proof for the other case follows the same sketch. Let $p(s, t) = (s, \dots, u, \dots, t)$, such that PR3 can be triggered. W.l.o.g. let $aid(u) < aid(t)$ (if $aid(u) \geq aid(t)$ and PR3 is triggered, then there exists vertex w , $aid(w) < aid(u)$, which can be reduced to the case of $aid(u) < aid(t)$). Given this, we have path $p(s, u)$ and $p(u, t)$ have the same k-MR L according to the definition of PR3. Then we have three cases: Case (1) $aid(s) > aid(u)$; Case (2) $aid(s) = aid(u)$; Case (3) $aid(s) < aid(u)$. Case (1) can be proved by Lemma 4, because $aid(s) > aid(u)$, $aid(t) > aid(u)$, and both $p(s, u)$ and $p(u, t)$ have k-MR L . Case (2) can be proved by Lemma 3 because the backward KBS from t can visit u , i.e., $aid(s) = aid(u)$. Case (3) can also be proved by 3 if the forward KBS from s can visit t . The only case left is that $aid(s) < aid(u)$ and the forward KBS from s cannot visit t because of PR3. In this case, there must exist vertex v , $aid(v) < aid(s) < aid(u) = n + 1$, and the k-MR of $p(s, v)$ and $p(v, u)$ is L . Then we have $aid(v) < aid(s)$ and $aid(v) < aid(t)$, and both path $p(s, v)$ and (v, \dots, u, \dots, t) have the k-MR L , which can be proved by Lemma 4. ■

Proof of Theorem 2: Assuming there exists index entry $(t, L) \in \mathcal{L}_{out}(s)$ in the RLC index, and there also exist $(u, L) \in \mathcal{L}_{out}(s)$ and $(u, L) \in \mathcal{L}_{in}(t)$. Then we have $aid(u) \geq aid(t)$, otherwise $(t, L) \in \mathcal{L}_{out}(s)$ can be pruned. Given this, $(u, L) \in \mathcal{L}_{in}(t)$ cannot exist because the backward KBS from t is performed earlier than the forward KBS from u , which

means we have either $(t, L) \in \mathcal{L}_{out}(u)$, or $(v, L) \in \mathcal{L}_{out}(u)$ and $(v, L) \in \mathcal{L}_{in}(t)$, such that $(u, L) \in \mathcal{L}_{in}(t)$ is pruned. The proof follows the same sketch if $(s, L) \in \mathcal{L}_{in}(t)$ is considered. ■

Proof of Theorem 3: (Sufficiency) It is straightforward.

(Necessity) Let p be the path from s to t with the k-MR L . W.l.o.g. let the backward KBS from t be performed first. Then we have two cases: the backward KBS from t can visit or cannot visit s . In the first case, the k-MR L of path p must be recorded according to Lemma 3. In the second case, PR3 must be triggered. According to Lemma 5, we have the k-MR of p is also recorded. ■

APPENDIX C IMPACT OF k

In this experiment, we aim at analyzing the impact of k on the index. We use a BA-graph and an ER-graph with 125K vertices, average degree 5, 16 edge labels. We build RLC indexes for the BA-graph and the ER-graph with k in (2, 3, 4). For each graph, we evaluate a true-query set of 1000 queries and a false-query set of 1000 queries using the three different indexes built with the three different k values. The results of indexing time, index size, and query time are reported in Fig. 7.

Overall results. Fig. 7 shows that the indexing time and index size of both types of synthetic graphs rise exponentially as k grows. The fundamental reason is that as k increases, the number of possible minimum repeats that have to be considered in graph traversals during indexing increases exponentially, which is an inherent step in building a complete reachability index for RLC queries. The exponential increase of minimum repeats w.r.t k also has an impact on query time, particularly on the true-queries of BA-graphs and the false-queries on ER-graphs. This is mainly due to larger index size, and the true-queries on BA-graphs and the false-queries on ER-graphs are more expensive to process than their opposites, respectively.

APPENDIX D REMARKS

An alternative version of the RLC index allowed to concatenate different minimum repeats to answer an RLC query, i.e., in Case 1 of Definition 4, L' can be different from L'' in the alternative version. However, such a design will prevent the use of PR3 that can prune vertices and avoid redundant traversals. The reason is related to the technical proof of Theorem 3. In a nutshell, if concatenating different minimum repeats is allowed and PR3 is also allowed to apply, then the index might not be complete, i.e., missing some index entries. Therefore, we can only allow one of the two designs, i.e., either allowing concatenating different minimum repeats or allowing applying PR3. In the previous version, we allowed the former one. Consequently, the indexing time of the alternative version is much longer than the version introduced in this paper. For instance, for the smallest graph used in our experiments (the AD graph presented in Section VI), the indexing time of the

alternative version is 32x slower than the current one. Thus, we focus on concatenating the same minimum repeats.