# Regular Path Query Evaluation Sharing a Reduced Transitive Closure Based on Graph Reduction

Inju Na
*School of Computing*
*KAIST*
Daejeon, Republic of Korea
ijna@kaist.ac.kr

Yang-Sae Moon[*]
*Department of Computer Science*
*Kangwon National University*
Chuncheon, Republic of Korea
ysmoon@kangwon.ac.kr

Ilyeop Yi, Kyu-Young Whang, Soon J. Hyun
*School of Computing*
*KAIST*
Daejeon, Republic of Korea
{iyyi, kywhang, sjhyun}@kaist.ac.kr

*Abstract*—**Regular path queries (RPQs) find pairs of vertices of paths satisfying given regular expressions on an edge-labeled, directed multigraph. When evaluating an RPQ, the evaluation of a Kleene closure (i.e., Kleene plus or Kleene star) is very expensive. Furthermore, when multiple RPQs include a Kleene closure as a common sub-query, repeated evaluations of the common sub-query cause serious performance degradation. In this paper, we present a novel concept of *RPQ-based graph reduction*, which significantly simplifies the original graph through edge-level and vertex-level reductions. Interestingly, RPQ-based graph reduction can replace the evaluation of the Kleene closure on the large original graph to that of the transitive closure to the small reduced graph. We then propose a reduced transitive closure (RTC) as a lightweight structure for efficiently sharing the result of a Kleene closure. We also present an RPQ evaluation algorithm, *RTCSharing*, which treats each clause in the disjunctive normal form of the given RPQ as a batch unit. If the batch units include a Kleene closure as a common sub-query, we share the lightweight RTC instead of the heavyweight result of the Kleene closure. RPQ-based graph reduction further enables us to formally represent the result of an RPQ including a Kleene closure as a relational algebra expression including the RTC. Through the formal expression, we optimize the evaluation of the batch unit by eliminating *useless* and *redundant operations* of the previous method. Experiments show that RTCSharing improves the performance significantly by up to 73.86 times compared with existing methods in terms of query response time.**

*Index Terms*—**regular path queries (RPQ), query optimization, graph reduction, reduced transitive closure (RTC), RTCSharing**

## I. INTRODUCTION

A regular path query (RPQ) is a regular expression that finds ordered pairs of vertices of paths satisfying the RPQ on an edge-labeled, directed multigraph [1]. Here, paths satisfy an RPQ when a sequence of labels matches the RPQ. The RPQ is attracting attention as an essential and important operation for graph data such as social networks and the Semantic Web and is supported by typical graph query languages such as SPARQL 1.1 [2] and Cypher [3]. RPQs can be used in

a variety of applications such as signal path detection in protein networks, recommending friends in social networks, and extracting information from linked open data [4], [5].

To evaluate an RPQ, it is necessary to traverse the graph from each vertex and perform pattern matching for labels of edges accessed during the traversal, finding the paths that satisfy the query. As such, evaluating an RPQ is a complex operation that combines graph traversal and pattern matching [6]. In particular, when an RPQ includes a Kleene closure (i.e., Kleene plus or Kleene star), the evaluation is more expensive because there can be long paths of varying lengths satisfying the RPQ [7]. Therefore, optimization of evaluating RPQs is actively being studied as an important issue [4], [5], [8]–[11]. Optimization of evaluating RPQs needs more attention when evaluating multiple RPQs. When there is a common sub-query among multiple RPQs, evaluating these RPQs individually leads to repeated evaluations of the common sub-query. This problem of repeated evaluations can be solved by evaluating the common sub-query once and sharing the result among the RPQs. However, when a common sub-query includes a Kleene closure, evaluating the common sub-query itself is expensive.

In this paper, we present a novel concept of *RPQ-based graph reduction*, which converts an original labeled multigraph $G$ to an unlabeled simple graph $\overline{G_R}$. RPQ-based graph reduction consists of two levels: 1) edge-level reduction producing an unlabeled graph $G_R$ from $G$ by mapping paths satisfying $R$ on $G$ to edges of $G_R$ and 2) vertex-level reduction producing a much reduced graph $\overline{G_R}$ from $G_R$ by representing each strongly connected component (SCC) as one vertex, which we define in Section 3. Interestingly, RPQ-based graph reduction can replace the evaluation of the Kleene closure on the large original graph $G$ to that of the transitive closure to the small reduced graph $\overline{G_R}$. Even though the graph reduction incurs a little overhead, the performance gain of using the reduced graph is much larger than the overhead. Based on the RPQ-based graph reduction, we then propose a reduced transitive closure (RTC) as a lightweight structure for efficiently sharing the result of a Kleene closure. Let $R_G^+$ (or $R_G^*$) be the evaluation result of a Kleene closure $R^+$ (or $R^*$) for any regular expression $R$ on the graph $G$. Since $R_G^*$ can be easily derived from $R_G^+$, we hereafter use $R^+$ and $R_G^+$ only unless confusion occurs. In the final algorithms (Algorithms 1 and 2), we also deal with $R^*$ by

simply deriving it from $R^+$. In Theorem 1, we formally show that $R_G^+$, the evaluation result of $R^+$ on the original graph $G$, can be easily calculated from the transitive closure (i.e., RTC) of the reduced graph $\overline{G_R}$.

We also propose an RPQ evaluation algorithm, Reduced Transitive Closure Sharing (*RTCSharing*), which converts the given RPQ to a logically equivalent disjunctive normal form (DNF) and evaluates it treating each clause as a batch unit. As in [15], we can convert all RPQs to a logically equivalent DNF treating each outermost Kleene closure as a literal. If multiple batch units include a Kleene closure as a common sub-query, we share the lightweight RTC instead of the heavyweight result of the Kleene closure. RPQ-based graph reduction further enables us to represent the result of an RPQ including a Kleene closure as a relational algebra expression in the form of a join sequence including the RTC. By representing the batch unit as a relational algebra expression, RTCSharing optimizes the evaluation of the batch unit by eliminating *redundant* and *useless operations* (see Section 4.2 for formal definitions) as follows: (1) it eliminates redundant operations (in effect, mostly duplicate check operations) that might occur in the next step removing redundant elements by unioning the intermediate results of each join step; and (2) it also eliminates useless operations caused by selection through *Prefix* and by the property of the reduced graph $\overline{G_R}$.

The contributions of this paper are summarized as follows:

- We present a novel notion of *RPQ-based graph reduction* that converts a labeled multigraph to an unlabeled simple graph for the efficient evaluation of a Kleene closure.
- We propose an RTC as a lightweight structure to efficiently share the result of a common sub-query among RPQs whose common sub-query is a Kleene closure $R^+$ for any regular expression $R$.
- We show that $R^+$ can be evaluated as the transitive closure of the edge-level reduced graph $G_R$ (Lemma 1), and further, can be calculated from the transitive closure of the two-level reduced graph $\overline{G_R}$ (i.e., the RTC) (Lemma 3 and Theorem 1).
- We propose an RPQ evaluation algorithm, *RTCSharing*, that uses each clause in the DNF of the given RPQ as a batch unit. RTCSharing shares the RTC instead of $R_G^+$ among batch units as the result of the common sub-query $R^+$.
- We show that the result of an RPQ including a Kleene closure can be represented as a relational algebra expression including the RTC and exploit it for optimizing the evaluation of the RPQ by eliminating *useless* and *redundant operations*.
- Through experiments with synthetic and real datasets, we show that *RTCSharing* significantly outperforms existing algorithms: 1) the recent algorithm that shares $R_G^+$ among RPQs [8] and 2) the naive algorithm that shares nothing among RPQs [5].

The paper is organized as follows. Section II describes the preliminary background. Section III proposes RPQ-based graph reduction and the RTC. Section IV describes the *RTC-Sharing* algorithm with the optimization. Section V shows the results of performance evaluation for the proposed methods in comparison with existing RPQ evaluation algorithms. Section VI summarizes the related work. Section VII concludes the paper.

## II. PRELIMINARIES

In this section, we introduce the background of regular path query research. Section II-A describes the data model used for regular path queries. Section II-B describes the definition of regular path queries and methods of evaluating them. Section II-C describes methods of evaluating multiple regular path queries.

### A. Data Model

The data used for regular path queries is an edge-labeled, directed multigraph $G$, which is defined as a 5-tuple ($V$, $E$, $f$, $\Sigma$, $l$) [1]. $V$ is a set of vertices. $E$ is a set of directional edges. $f: E \rightarrow V \times V$ is a function that maps each edge to an ordered pair of two vertices connected by the edge. $\Sigma$ is a set of labels. $l: E \rightarrow \Sigma$ is a function that maps each edge to its label. In $G$, multiple edges between any two vertices are allowed, but the labels of these edges must be distinct. Fig. 1 shows a graph, which will be used as an example graph throughout the paper.

TABLE I summarizes the notation related to $G$ to be used in this paper. Each vertex has a unique ID (VID). A vertex with VID $i$ is denoted by $v_i$. A label has a unique ID (LID). A label with LID $i$ is denoted by $l_i$. An edge from $v_s$ to $v_d$ with label $l_i$ is denoted by $e(v_s, l_i, v_d)$. A path from $v_s$ to $v_d$, which is a sequence of $e(v_s, l_{i_1}, v_{j_1})$, $e(v_{j_1}, l_{i_2}, v_{j_2})$, $\cdots$, $e(v_{j_{n-1}}, l_{i_n}, v_d)$ is denoted by $p(v_s, l_{i_1}, v_{j_1}, l_{i_2}, v_{j_2}, \cdots, l_{i_n}, v_d)$. A sequence of labels $l_{i_1} l_{i_2} \cdots l_{i_n}$ is called a path label [1], [5]. If it is not ambiguous, $p(v_s, l_{i_1}, v_{j_1}, l_{i_2}, v_{j_2}, \cdots, l_{i_n}, v_d)$ is simply denoted by $p(v_s, v_d)$.
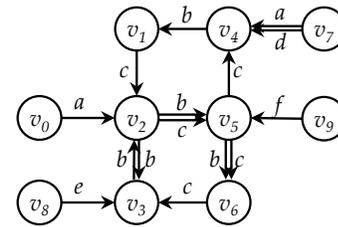


Fig. 1: An edge-labeled, directed multigraph.

TABLE I: The notation related to $G$.

| Notation | Description |
|---|---|
| $v_i$ | A vertex with VID i |
| $l_i$ | A label with LID i |
| $e(v_s, l_i, v_d)$ | An edge from $v_s$ to $v_d$ with label $l_i$ |
| $p(v_s, l_{i_1}, \cdots, l_{i_n}, v_d)$ or $p(v_s, v_d)$ | A path from $v_s$ to $v_d$, which is a sequence of $e(v_s, l_{i_1}, v_{j_1})$, $\cdots$, $e(v_{j_{n-1}}, l_{i_n}, v_d)$ |
| $l_{i_1} l_{i_2} \cdots l_{i_n}$ | The path label of a path $p(v_s, l_{i_1}, v_{j_1}, l_{i_2}, v_{j_2}, \cdots, l_{i_n}, v_d)$ |

## B. Regular Path Query(RPQ)

A regular path query (RPQ) $R$ on a graph $G$ is a regular expression over $\Sigma$ and finds a set of ordered vertex pairs (start vertex, end vertex) of the paths satisfying $R$ in $G$ [1] (Definition 1). That is, the evaluation result $R_G$ of RPQ $R$ on $G$ is a set of ordered vertex pairs as is defined in Definition 2.

**Definition 1.** *A path satisfies an RPQ $R$ when a path label of the path matches $R$ [1].*

**Definition 2.** *Given a graph $G$ and an RPQ $R$,*
$R_G = \{(v_i, v_j) \mid a\ p(v_i, v_j)\ that\ satisfies\ R\ exists\ in\ G\}.$

**Example 1.** Fig. 2 shows an example of $R_G$. $(d \cdot (b \cdot c)^+ \cdot c)_G$ is the result of RPQ $d \cdot (b \cdot c)^+ \cdot c$ on $G$ where $G$ is an example graph shown in Fig. 1. The left side of the figure shows the paths that satisfy $d \cdot (b \cdot c)^+ \cdot c$. The path labels of these paths are *dbcc*, *dbcbcc*, etc., all of which match $d \cdot (b \cdot c)^+ \cdot c$. The right side shows $(d \cdot (b \cdot c)^+ \cdot c)_G$. $(v_7, v_5)$ is included in $(d \cdot (b \cdot c)^+ \cdot c)_G$ because paths $p_1$, $p_3$, etc., satisfying $d \cdot (b \cdot c)^+ \cdot c$ exist between vertices $v_7$ and $v_5$. $(v_7, v_3)$ is also included in $(d \cdot (b \cdot c)^+ \cdot c)_G$ because paths $p_2$, $p_4$, etc., satisfying $d \cdot (b \cdot c)^+ \cdot c$, exist between vertices $v_7$ and $v_3$. □

To evaluate an RPQ $R$ on $G$, we traverse $G$ from each vertex $v_i$ and find the set of end vertices of the paths starting from $v_i$ that satisfy $R$. To find the paths satisfying $R$, pattern matching is done on the labels of the edges that are accessed during traversal. Finite automata are usually used for pattern matching [1], [4], [5], [10], [11]. Each traversal has a state of the finite automaton as the current state. If the state can be transited to the next state using the label of an edge originating from the current state, the traversal continues through that edge. If the current state is an accept state, (start vertex, end vertex) of the traversal is included in $R_G$. The traversal continues until there are no more traversable edges or the state is not able to be transited. Among algorithms based on finite automata, we use recent used one by Yakovets [5] in experiments in Section V.

**Example 2.** Fig. 3 shows an example of RPQ evaluation — evaluation of the RPQ $d \cdot (b \cdot c)^+ \cdot c$ using the finite automata (i.e.,



Fig. 2: An example of $R_G$.



Fig. 3: An example of RPQ evaluation.

NFA in Fig. 3) on the graph of Fig. 1. Through the traversal, $(v_7, v_5)$ and $(v_7, v_3)$ are included in $(d \cdot (b \cdot c)^+ \cdot c)_G$ since $p(v_7, d, v_4, b, v_1, c, v_2, c, v_5)$ and $p(v_7, d, v_4, b, v_1, c, v_2, b, v_5, c, v_6, c, v_3)$, which satisfies $d \cdot (b \cdot c)^+ \cdot c$ starting from $v_7$, are found. In the case of $p(v_7, d, v_4, b, v_1, c, v_2, b, v_3)$, there is an edge $e(v_3, b, v_2)$ accessible from $v_3$, but state transition cannot be done through the edge. Therefore, the traversal of the path is terminated. $p(v_7, d, v_4, b, v_1, c, v_2, b, v_5, c, v_4, b, v_1)$ is a special case. Since the end vertex $v_1$ of the path has already been visited by $p(v_7, d, v_4, b, v_1)$ in the same state (i.e., $q_2$ in Fig. 3), the subsequent traversals overlap with an earlier one. That is, the vertices that are included in $(d \cdot (b \cdot c)^+ \cdot c)_G$ by paths from $v_7$ found through subsequent traversals are duplicated. Therefore, if the end vertex of the path has already been visited in the same state from the start vertex of the path, the traversal is terminated to avoid duplication. □

As we can see in Example 2, evaluation of RPQ $R$ on $G$ is a complex operation since it requires both graph traversal and pattern matching [6]. In particular, evaluation of an RPQ with Kleene closures is more expensive and the result is larger than those of RPQs without them because the former can have long paths of varied lengths that satisfy the RPQ [7]. Time and space complexities of the naive RPQ evaluation depends on those of Kleene closure evaluation, which are $O(|V| \times |E|)$ and $O(|V|^2)$, respectively. Therefore, optimization of the evaluation of an RPQ with Kleene closures is a very important issue.

## C. Evaluation of Multiple RPQs

A naive method of evaluating multiple RPQs is to evaluate them individually. However, if there is a common sub-query among RPQs, the sub-query is evaluated repeatedly. To avoid the problem, Abul-Basher's method [8] shares the evaluation result of the common sub-query among the RPQs. We call the method *FullSharing* and use it as a baseline solution in the experiment. However, if a common sub-query is a Kleene closure, the one-time evaluation of the common sub-query itself is expensive and the size of the result could be large as mentioned in Section II-B. In addition, *useless* and *redundant operations* can occur when each RPQ is evaluated using the result of the common sub-query. We discuss the details of these problems and propose solutions in Section IV.

## III. RPQ-BASED GRAPH REDUCTION

In this section, we present a novel concept of *RPQ-based graph reduction*. Fig. 4 shows an overview of RPQ-based graph reduction. The first-level reduction ($G \rightarrow G_R$ of Section III-A) reduces a graph $G$ at the edge level for RPQ $R$, which maps the paths satisfying $R$ on $G$ to edges on $G_R$. The second-level reduction ($G_R \rightarrow \overline{G_R}$ of Section III-B) reduces $G_R$ at the vertex level, which maps each SCC of $G_R$ to a vertex of $\overline{G_R}$. Based on these graph reductions, we show that $R^+$ can be evaluated as the transitive closure of the edge-level reduced graph $G_R$ (Lemma 1), and further, can be calculated as the transitive closure of the two-level reduced graph $\overline{G_R}$ and the Cartesian product of vertices of SCCs mapped to vertices of $\overline{G_R}$ (Lemma 3 and Theorem 1). We also propose the transitive closure of $\overline{G_R}$ as a lightweight structure for efficiently sharing the result of $R^+$, which we call a reduced transitive closure (RTC of Section III-C).

TABLE II summarizes the notation related to $G_R$ and $\overline{G_R}$. Vertex, edge, and path notations of $G_R$ are the same as those of $G$ in TABLE I except that we use $e_R$ and $p_R$ instead of $e$ and $p$ to distinguish $G_R$ and $G$. Each SCC of $G_R$ has a unique ID (SID). An SCC with SID $i$ is denoted by $s_i$. If it is not ambiguous, the set of vertices in the SCC with SID $i$ is also denoted by $s_i$. Each vertex of $\overline{G_R}$ has the same unique ID as that of the SCC of $G_R$ mapped to it. A vertex with VID $i$ is denoted by $\overline{v_i}$. An edge from $\overline{v_s}$ to $\overline{v_d}$ is denoted by $\overline{e_R}(\overline{v_s}, \overline{v_d})$.

### A. Edge-Level Graph Reduction ($G \rightarrow G_R$)

The edge-level reduction maps all the paths satisfying $R$ between each pair of vertices of $G$ to one edge of $G_R$. Graph $G_R$ is an unlabeled, directed graph that reduces $G$ at the edge level for $R$. $G_R$ is defined as a 3-tuple ($V_R$, $E_R$, $f_R$). $V_R$ is a set of vertices: $\{v_i \mid (\exists v_j)((v_i, v_j) \in E_R \lor (v_j, v_i) \in E_R)\}$. $E_R$ is a set of edges: $\{e_R(v_i, v_j) \mid$ there exists $p(v_i, v_j)$ satisfying $R$ on $G\}$. $f_R: E_R \rightarrow V_R \times V_R$ is a function that maps each edge to an ordered pair of vertices connected by the edge. Reduction of $G$ to $G_R$ through the edge-level reduction has the following aspects.

- Vertices and edges that do not belong to a path satisfying $R$ are excluded. To obtain $R^+_G$ we only need the paths satisfying $R$. Thus, we don't need to consider vertices and edges not belonging to the paths satisfying $R$ on $G$.



Fig. 4: An overview of RPQ-based graph reduction.

TABLE II: The notation related to $G_R$ and $\overline{G_R}$.

| | Notation | Description |
|---|---|---|
| $G_R$ | $v_i$ | A vertex of $G_R$ with VID $i$ |
| | $e_R(v_s, v_d)$ | An edge of $G_R$ from $v_s$ to $v_d$ |
| | $p_R(v_s, v_d)$ | A path of $G_R$ from $v_s$ to $v_d$, which is a sequence of $e_R(v_s, v_{j_1}), \cdots, e_R(v_{j_{n-1}}, v_d)$ |
| | $s_i$ | An SCC of $G_R$ with SID $i$ or the set of vertices in the SCC |
| $\overline{G_R}$ | $\overline{v_i}$ | The vertex of $\overline{G_R}$ to which $s_i$ of $G_R$ is mapped (VID is $i$) |
| | $\overline{e_R}(\overline{v_s}, \overline{v_d})$ | An edge of $\overline{G_R}$ from $\overline{v_s}$ to $\overline{v_d}$ |

- Labeled graph $\rightarrow$ unlabeled graph. Since only the paths satisfying $R$ are mapped to edges, labels of edges are no longer needed. (i.e., every label is $R$.)
- Multigraph $\rightarrow$ simple graph. Since labels of edges are excluded, paths in the same direction between each vertex pair are mapped to one edge.

**Example 3.** Fig. 5 shows an example of edge-level graph reduction. Here, the graph $G$ in Fig. 1 is reduced at the edge level for $b \cdot c$. In $G$, the paths satisfying $b \cdot c$ are $p(v_2, v_4)$, $p(v_2, v_6)$, $p(v_3, v_5)$, $p(v_4, v_2)$, and $p(v_5, v_3)$. Since these paths are mapped to edges of $G_{b \cdot c}$, $E_{b \cdot c}$ becomes $\{e_{b \cdot c}(v_2, v_4), e_{b \cdot c}(v_2, v_6), e_{b \cdot c}(v_3, v_5), e_{b \cdot c}(v_4, v_2), e_{b \cdot c}(v_5, v_3)\}$. $\square$

We show that the problem of evaluating RPQ $R^+$ can be reduced to a problem of computing the transitive closure of the edge-level reduced graph $G_R$ (denoted by TC($G_R$)) in Lemma 1.

**Lemma 1 .** $R^+_G$ *is equivalent to TC($G_R$).*

*Proof.* We prove $R^+_G$ = TC($G_R$) by showing $R^+_G \subseteq$ TC($G_R$) and $R^+_G \supseteq$ TC($G_R$), respectively.

- $R^+_G \subseteq$ TC($G_R$)
  Suppose $(v_i, v_j) \in R^+_G$. Then, there exists a path $p(v_i, v_j)$ satisfying $R^n$ (concatenation of $R$ $n$ times, $n \geq 1$) on $G$. $p(v_i, v_j)$ can be represented as a sequence of $n$ paths $p(v_i, v_{k_1}), p(v_{k_1}, v_{k_2}), \cdots, p(v_{k_{n-1}}, v_j)$ each satisfying $R$. By edge-level reduction, all the paths satisfying $R$ between each pair of vertices maps to one edge in $G_R$. That is, because there exist $e_R(v_i, v_{k_1}), e_R(v_{k_1}, v_{k_2}), \cdots, e_R(v_{k_{n-1}}, v_j)$ in $G_R$, there exist $p_R(v_i, v_j)$ in $G_R$. Hence, $(v_i, v_j) \in$ TC($G_R$) and $R^+_G \subseteq$ TC($G_R$).
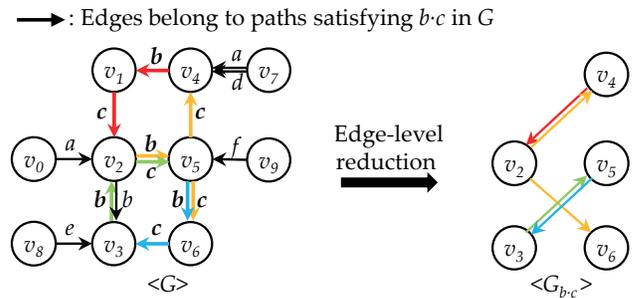


Fig. 5: An example of edge-level reduction.

- $R_G^+ \supseteq \text{TC}(G_R)$

  Suppose $(v_i, v_j) \in \text{TC}(G_R)$. Then, there exists a path $p_R(v_i, v_j)$ of length $n(n \geq 1)$ on $G_R$. $p_R(v_i, v_j)$ can be represented as a sequence of $n$ edges $e_R(v_i, v_{k_1})$, $e_R(v_{k_1}, v_{k_2})$, $\cdots$, $e_R(v_{k_{n-1}}, v_j)$. By the definition of $G_R$, there exist paths $p(v_i, v_{k_1})$, $p(v_{k_1}, v_{k_2})$, $\cdots$, $p(v_{k_{n-1}}, v_j)$ each satisfying $R$ in $G$. That is, there exists a path $p(v_i, v_j)$ satisfying $R^n$. Hence, $(v_i, v_j) \in R_G^+$ and $R_G^+ \supseteq \text{TC}(G_R)$. $\square$

**Example 4.** In Fig. 5, $(b \cdot c)_G^+$ is equivalent to $\text{TC}(G_{b \cdot c})$: $\{(v_2, v_2), (v_2, v_4), (v_2, v_6), (v_3, v_3), (v_3, v_5), (v_4, v_2), (v_4, v_4), (v_4, v_6), (v_5, v_3), (v_5, v_5)\}$. $\square$

### B. Vertex-Level Graph Reduction ($G_R \to \overline{G_R}$)

The vertex-level reduction maps each SCC of $G_R$ to one vertex of $\overline{G_R}$. The graph $\overline{G_R}$ is an unlabeled, directed graph. $\overline{G_R}$ is defined as a 3-tuple $(\overline{V_R}, \overline{E_R}, \overline{f_R})$. $\overline{V_R}$ is a set of vertices: $\{\overline{v_i} \mid (\exists s_j)(s_j \text{ is an SCC of } G_R \wedge i = j)\}$. $\overline{E_R}$ is a set of edges: $\{\overline{e_R}(\overline{v_i}, \overline{v_j}) \mid (\exists s_k)(\exists s_l)(\exists v_m)(\exists v_n)(s_k \text{ and } s_l \text{ are SCCs of } G_R \wedge v_m \in s_k \wedge v_n \in s_l \wedge e_R(v_m, v_n) \in E_R \wedge i = k \wedge j = l)\}$. $\overline{f_R}: \overline{E_R} \to \overline{V_R} \times \overline{V_R}$ is a function that maps each edge to an ordered pair of vertices connected by the edge. The vertex-level reduction has the following characteristics:

- Edges between any pair of vertices in the same SCC of $G_R$ are mapped to one self-loop edge in $\overline{G_R}$.
- Edges with the same direction between any pair of vertices in two different SCCs of $G_R$ are mapped to one edge in $\overline{G_R}$.

**Example 5.** Fig. 6 shows an example of vertex-level graph reduction. Here, $G_{b \cdot c}$ in Fig. 5 is reduced to $\overline{G_{b \cdot c}}$ at the vertex level. There are three SCCs in $G_{b \cdot c}$: $s_0$, $s_1$, $s_2$. Since each of these SCCs is mapped to one vertex of $\overline{G_{b \cdot c}}$, $\overline{V_{b \cdot c}}$ becomes $\{\overline{v_0}, \overline{v_1}, \overline{v_2}\}$. There exists an edge from $v_2$ in $s_0$ to $v_6$ in $s_1$; these vertices belong to different SCCs. This edge of $G_R$ is mapped to the edge $\overline{e_{b \cdot c}}(\overline{v_0}, \overline{v_1})$ of $\overline{G_R}$. The SCC $s_0$ $(s_2)$ has edges between the vertices that belong to $s_0$ $(s_2)$ itself. These edges are mapped to the edge $\overline{e_{b \cdot c}}(\overline{v_0}, \overline{v_0})$ $(\overline{e_{b \cdot c}}(\overline{v_2}, \overline{v_2}))$ in $\overline{G_R}$ constituting an self-loop edge. Thus, $\overline{E_{b \cdot c}}$ is $\{\overline{e_{b \cdot c}}(\overline{v_0}, \overline{v_0}), \overline{e_{b \cdot c}}(\overline{v_0}, \overline{v_1}), \overline{e_{b \cdot c}}(\overline{v_2}, \overline{v_2})\}$. $\square$

Using Lemma 3 below, we can efficiently compute the $\text{TC}(G_R)$ (and accordingly, $R_G^+$) by computing the transitive closure of $\overline{G_R}$ (denoted by $\text{TC}(\overline{G_R})$). Purdom [12] provided the following Lemma 2.



Fig. 6: An example of vertex-level reduction.

**Lemma 2** [12]. *For any pair of vertices $v_i$ and $v_j$, if there is a path from $v_i$ to $v_j$, there are paths from each vertex belonging to the same SCC as that of $v_i$ to all vertices belonging to the same SCC as that of $v_j$.*

We now introduce the following Lemma 3.

**Lemma 3 .** $\text{TC}(G_R) = \{(v_i, v_j) \mid (\overline{v_k}, \overline{v_l}) \in \text{TC}(\overline{G_R}) \wedge (v_i, v_j) \in s_k \times s_l\}$, *where $\overline{v_k}$, $\overline{v_l}$ are vertices of $\overline{G_R}$ to which $s_k$ and $s_l$ of $G_R$ are mapped, respectively.*

*Proof.*

- $\text{TC}(G_R) \subseteq \{(v_i, v_j) \mid (\overline{v_k}, \overline{v_l}) \in \text{TC}(\overline{G_R}) \wedge (v_i, v_j) \in s_k \times s_l\}$. Suppose $(v_i, v_j) \in \text{TC}(G_R)$. Then, there is a path from a vertex of $\overline{G_R}$ which is mapped to the SCC containing $v_i$ to a vertex of $\overline{G_R}$ which is mapped to the SCC containing $v_j$ by the definition of $\overline{G_R}$. Hence, $(v_i, v_j) \in \{(v_i, v_j) \mid (\overline{v_k}, \overline{v_l}) \in \text{TC}(\overline{G_R}) \wedge (v_i, v_j) \in s_k \times s_l\}$ holds.
- $\{(v_i, v_j) \mid (\overline{v_k}, \overline{v_l}) \in \text{TC}(\overline{G_R}) \wedge (v_i, v_j) \in s_k \times s_l\} \subseteq \text{TC}(G_R)$. Lemma 2 means that if there is a path from $v_i$ in $s_k$ to $v_j$ in $s_l$, then there exist a path from each vertex in $s_k$ to each vertex in $s_l$. Hence, $\{(v_i, v_j) \mid (\overline{v_k}, \overline{v_l}) \in \text{TC}(\overline{G_R}) \wedge (v_i, v_j) \in s_k \times s_l\} \subseteq \text{TC}(G_R)$ holds by Lemma 2. $\square$

The transitive closure algorithms in [12] and [13] are instances of implementation of Lemma 3 (without a formal correctness proof). Purdom [12] proposed an algorithm that essentially computes the transitive closure of $\overline{G_R}$, and then, the Cartesian product $s_k \times s_l$ without formally introducing the concept of the vertex-reduced graph, but instead, treating the nodes in an SCC of the original graph as an equivalent class. Nuutila [13] improved Purdom's algorithm by obtaining the transitive closure of $\overline{G_R}$ and the Cartesian product $s_k \times s_l$ in one step in an interleaved way. In this paper, we formalize the concept implied by the algorithm by Prudom [12] with a formal definition of the vertex reduced graph and the correctness proof.

**Theorem 1.** $R_G^+ = \{(v_i, v_j) \mid (\overline{v_k}, \overline{v_l}) \in \text{TC}(\overline{G_R}) \wedge (v_i, v_j) \in s_k \times s_l\}$ *where $\overline{v_k}$, $\overline{v_l}$ are vertices of $\overline{G_R}$ to which $s_k$ and $s_l$ of $G_R$ are mapped, respectively.*

*Proof.* Derived from Lemmas 1 and 3. $\square$

**Example 6.** In Fig. 6, $\text{TC}(\overline{G_{b \cdot c}})$ is $\{(\overline{v_0}, \overline{v_0}), (\overline{v_0}, \overline{v_1}), (\overline{v_2}, \overline{v_2})\}$. For each vertex pair $(\overline{v_i}, \overline{v_j})$, the union of the Cartesian product of $s_i$ and $s_j$ is $\{(v_2, v_2), (v_2, v_4), (v_4, v_4), (v_4, v_2), (v_2, v_6), (v_4, v_6), (v_3, v_3), (v_3, v_5), (v_5, v_3), (v_5, v_5)\}$, which is the same as $\text{TC}(G_{b \cdot c})$. $\square$

### C. Reduced Transitive Closure (RTC)

We use $\text{TC}(\overline{G_R})$ (denoted by $\overline{R_G^+}$) as a reduced transitive closure to share the result of a Kleene plus $R^+$. As illustrated in Section III-B, we can efficiently enumerate $R_G^+$ by using $\overline{R_G^+}$. Moreover, $\overline{R_G^+}$ is computationally simpler and smaller than $R_G^+$ as shown in TABLE III. Although both computing $\overline{R_G^+}$ on $\overline{G_R}$ and computing $R_G^+$ on $G$ find pairs of vertices by traversing a graph, the former is simpler than the latter because of the following differences:

- The size of the target graph to be traversed is smaller ($|\overline{V_R}| << |V_R|$ in general). $\overline{G_R}$ is a two-level reduced graph of $G$, which is reduced in size. Thus, the former target graph $\overline{G_R}$ is generally smaller than the latter target graph $G$.
- The operations performed when traversing the graph are simpler. The former traverses the graph using only operations that identify reachability. The latter, on the other hand, performs additional pattern matching operations for the labels of the edges as well. Therefore, the former operations are simpler.

TABLE III summarizes the computational and space complexity of $\overline{R_G^+}$ and $R_G^+$. $\overline{R_G^+}$ is computed on $\overline{G_R}$ after two-level graph reduction: $G \rightarrow G_R \rightarrow \overline{G_R}$. $R_G^+$ is computed using $R_G$ after evaluating $R$ on $G$. The main operation required when reducing the graph $G \rightarrow G_R$ is to evaluate $R$ on $G$. As explained in Section II-B, evaluating $R$ on $G$ is computationally simpler than and relatively negligible with evaluating $R^+$ on $G$. Therefore, we exclude the computation for reducing the graph $G \rightarrow G_R$ from the computational complexities of both $\overline{R_G^+}$ and $R_G^+$. The computational complexity of evaluating $R_G^+$ is $O(|V_R| \times |E_R|)$. The main operation required when reducing the graph $G_R \rightarrow \overline{G_R}$ is to find all SCCs of $G_R$. The most efficient method for this operation is known as the Tarjan's algorithm [14], whose computational complexity is $O(|V_R|+|E_R|) << O(|V_R| \times |E_R|)$. Since the overhead of reducing the graph $G_R \rightarrow \overline{G_R}$ is negligible compared with the computational complexity of evaluating $R^+$ on $G_R$, we exclude it from the comparison in TABLE III. The computational complexity of computing $\overline{R_G^+}$ on $\overline{G_R}$ is $O(|\overline{V_R}| \times |\overline{E_R}|)$ which is generally smaller than $O(|V_R| \times |E_R|)$. Therefore, $\overline{R_G^+}$ has smaller computational complexity than $R_G^+$. This observation is demonstrated in performance evaluation of Section V. In the worst case, $\overline{R_G^+}$ and $R_G^+$ are all vertex pairs of the target graph, so that the space complexity is $O(|\overline{V_R}|^2)$ and $O(|V_R|^2)$, respectively. By the vertex-level reduction $O(|\overline{V_R}|) << O(|V_R|)$ in general. That is, $\overline{R_G^+}$ has a generally smaller space complexity than $R_G^+$.

## IV. REDUCED TRANSITIVE CLOSURE SHARING (RTCSHARING)

In this section, we propose an RPQ evaluation algorithm, which we call *RTCSharing*. *RTCSharing* recursively evaluates RPQs calling EvalBatchUnit, which evaluates batch units of the RPQs. *RTCSharing* also share the RTC among the batch units. Section IV-A proposes *RTCSharing* and Section IV-B optimizes the evaluation of batch units.

### A. RTCSharing: an RPQ Evaluation Algorithm

In *RTCSharing*, we first convert the given RPQ to a logically equivalent disjunctive normal form (DNF). Since all logical formulas can be converted to a logically equivalent DNF [15],

we can convert all RPQs to a logically equivalent DNF treating each outermost Kleene closure as a literal. Then, we evaluate each clause in the DNF treating it as a batch unit. The batch unit is in the form of *Prefix·R⁺·Postfix* or *Prefix·R\*·Postfix* where *Prefix* and $R$ are any regular expressions that can include multiple or nested Kleene closures and *Postfix* is a regular expression without a Kleene closure, i.e., the $R^+$ or $R^*$ is the rightmost Kleene closure in the clause. For concise notation, we denote *Prefix* and *Postfix* as *Pre* and *Post*, respectively. If the batch unit includes multiple or nested Kleene closures, we recursively evaluate it using the result of the previous recursive step until the batch unit does not include any Kleene closures. It means that an escape hatch of the recursion is the batch unit that does not include any Kleene closures. When evaluating batch units including a Kleene closure as a common sub-query, we share the RTC among batch units using Theorem 1. That is, since we evaluate it only once, the performance degradation from the DNF transformation does not occur. We can further improve the performance by optimizing the evaluation order of the batch units, and we leave the optimization issue as a future work. We explain the details of the evaluation of the batch unit in Section IV-B.

Algorithm 1 shows the details of *RTCSharing*. First, in line 2, we convert the given RPQ $Q$ to a logically equivalent DNF *Q_DNF*. Here, we treat each outermost Kleene closure as a literal. Then, we evaluate each clause in *Q_DNF* treating it as the batch unit in lines 4 to 12 and union the result in line 13. In line 4, we decompose the batch unit *clause* into *Pre*, $R$, *Type* (+, \*, or NULL), and *Post*. If *clause* does not include a Kleene closure, both *Pre* and $R$ are $\epsilon$, *Type* is NULL, and *Post* is the entire clause. In line 6, we evaluate the entire clause *Post* calling EvalRPQwithoutKC, which uses any existing method [5]. Otherwise, $R^+$ or $R^*$ is the rightmost Kleene closure, *Type* is + or \*, and *Post* does not include any Kleene closures. In line 8, we evaluate *Pre* recursively calling RTCSharing with *Pre* as a query. If the RTC for $R$ exists, we

TABLE III: Complexity of $R_G^+$ and $\overline{R_G^+}$.

| Complexity | $R_G^+$ | $\overline{R_G^+}$ |
|---|---|---|
| Computational | $O(|V_R| \times |E_R|)$ | $O(|\overline{V_R}| \times |\overline{E_R}|)$ |
| Space | $O(|V_R|^2)$ | $O(|\overline{V_R}|^2)$ |

---

**Algorithm 1:** RTCSharing.

**Input:** Query $Q$
**Output:** The set of results $Q_G$

1   $Q_G \leftarrow \emptyset$
2   $Q\_DNF \leftarrow$ ConvertRPQtoDNF($Q$)
3   **foreach** *clause* $\in Q\_DNF$ **do**
4     *Pre*, $R$, *Type*, *Post* $\leftarrow$ DecomposeCL(*clause*)
5     **if** *Type* == NULL **then**
      /\* *clause* has no Kleene closure      \*/
6       *clause_G* $\leftarrow$ EvalRPQwithoutKC(*Post*)
7     **else**
      /\* *clause* has a Kleene closure      \*/
8       *Pre_G* $\leftarrow$ RTCSharing(*Pre*)
9       **if** the RTC for $R$ does not exist **then**
        /\* Compute RTC: $\overline{R_G^+}$ and *SCC*   \*/
10        $R_G \leftarrow$ RTCSharing($R$)
11        $\overline{R_G^+}$, *SCC* $\leftarrow$ Compute_RTC($R_G$)
12       *clause_G* $\leftarrow$ EvalBatchUnit(*Pre_G*, $\overline{R_G^+}$, *SCC*, *Type*, *Post*)
13     $Q_G \leftarrow Q_G \cup$ *clause_G*

reuse them. Otherwise, we compute and store them to share in lines 10 and 11. In line 10, we evaluate $R$ recursively calling RTCSharing with $R$ as a query. Since $R_G$ is the same with the edge set of $G_R$, we can compute the RTC using $R_G$. In line 11, we compute the RTC $\overline{R_G^+}$ calling Compute_RTC that uses Tarjan's algorithm [14]. In line 12, we evaluate the batch unit $Pre \cdot R^+ \cdot Post$ calling EvalBatchUnit. We note that, unlike $Pre$, $Post$ is not pre-evaluated. $Post$ is directly handled by EvalBatchUnit. EvalBatchUnit will be given by Algorithm 2.

**Example 7.** We show three examples of RPQ evaluation using *RTCSharing*. To focus on the method we propose, we only consider the sequential evaluation of the given RPQs in the DNF. Fig. 7 shows recursion trees of the RPQs.

- RPQ = $a$.
  (The query does not include a Kleene closure.)
  We first decompose the query so that both $Pre$ and $R$ is $\epsilon$, $Type$ is NULL, and $Post$ is $a$. We then evaluate the query calling EvalRPQwithoutKC in line 6.
- RPQ = $a \cdot (a \cdot b)^+ \cdot b$.
  (The query includes a Kleene closure.)
  We first decompose the query so that $Pre$ is $a$, $R$ is $a \cdot b$, $Type$ is +, and $Post$ is $b$. We evaluate $a$ in line 8 and compute the RTC for $a \cdot b$ in lines 10 and 11. Then, we evaluate the query calling EvalBatchUnit with $a_G$ and the RTC for $a \cdot b$ in line 12.
- RPQ = $(a \cdot b)^* \cdot b^+ \cdot (a \cdot b^+ \cdot c)^+$.
  (The query includes multiple and nested Kleene closures.)
  We first decompose the query so that $Pre$ is $(a \cdot b)^* \cdot b^+$, $R$ is $a \cdot b^+ \cdot c$, $Type$ is +, and $Post$ is $\epsilon$. We evaluate $(a \cdot b)^* \cdot b^+$ calling *RTCSharing* recursively in line 8. In this recursive step, $Pre$ is $(a \cdot b)^*$ and $R$ is $b$, $Type$ is + and $Post$ is $\epsilon$. When evaluating $(a \cdot b)^*$, we reuse the RTC for $a \cdot b$, which was already computed when evaluating the above query $a \cdot (a \cdot b)^+ \cdot b$. When evaluating $(a \cdot b^+ \cdot c)^+$, we also reuse the RTC for $b$, which was already computed when evaluating $(a \cdot b)^* \cdot b^+$. Finally, we evaluate the query calling EvalBatchUnit with $((a \cdot b)^* \cdot b^+)_G$ and the RTC for $a \cdot b^+ \cdot c$ in line 12. □

### B. Evaluation of the Batch Unit

When evaluating the batch unit, especially $Pre \cdot R^+$, *useless* and *redundant operations* can occur because of $Pre_G$. These are sources of performance degradation. Among vertex pairs in $\overline{R_G^+}$, we need only those that are connected from a vertex pair in $Pre_G$. That is, operations that evaluate $R^+$ starting from vertices that are not connected from any vertex pair in $Pre_G$ are useless. We call them *useless-1 operations*.

Vertex pairs in $Pre_G$ whose start vertices are the same can cause duplicate results in $(Pre \cdot R^+)_G$. The concatenations of those vertex pairs and vertex pairs in $R_G^+$ whose end vertices are the same produce the same result. Therefore, when evaluating $R^+$ starting from vertex pairs in $Pre_G$ whose start vertices are the same, operations that find vertex pairs in $R_G^+$ whose end vertices are the same are *redundant operations*. We classify the *redundant operations* as *redundant-1 operations* as in Definition 3 and *redundant-2 operations* as in Definition 4.

**Definition 3.** *Redundant-1 operations are those that find vertex pairs in $R_G^+$ whose end vertices are the same when evaluating $R^+$ starting from vertex pairs in $Pre_G$ whose start vertices are the same and end vertices belong to the same SCC.*

**Definition 4.** *Redundant-2 operations are those that find vertex pairs in $R_G^+$ whose end vertices are the same when evaluating $R^+$ starting from vertex pairs in $Pre_G$ whose start vertices are the same but end vertices belong to different SCCs.*

**Example 8.** Fig. 8 shows an example evaluation of $Pre \cdot R^+$ involving *useless-1*, *redundant-1*, and *redundant-2 operations*. The paths satisfying $R^+$ from the vertices belonging to $s_h$ to the vertices belonging to $s_g$ and $s_f$ are not connected from vertex pairs in $Pre_G$. Therefore, operations that find them are *useless-1 operations*.
Operations that evaluate $Pre \cdot R^+$ starting from $v_0$ via vertices in $s_0$ or those in $s_i$ produce the same results (i.e., $(v_0, v_k)$, $\cdots$, $(v_0, v_j)$). That is, these operations are *redundant operations*. Among them, operations via vertices in the same SCC are *redundant-1 operations* (e.g., the operations from $v_0$
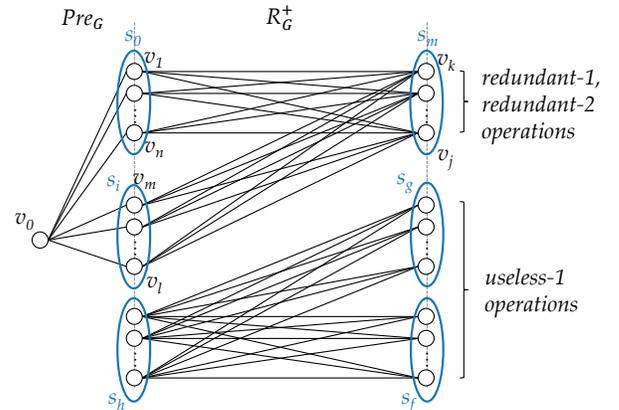


Fig. 7: Recursion trees of the example queries.



Fig. 8: An example of evaluating $Pre \cdot R^+$ involving *redundant* and *useless operations*.

via vertices in $s_0$). Otherwise, the operations via vertices in different SCCs are *redundant-2 operations* (e.g., the operations from $v_0$ via a vertex in $s_0$ and one in $s_i$). □

These *useless* and *redundant operations* are the major cause of performance degradation, and we need to optimize the evaluation of batch units. To efficiently evaluate the batch units, we first formally represent the result as a relational algebra expression including $\overline{R_G^+}$. We first define the notation:

- $R_G(START\_V, END\_V) = \{(v_i, v_j)|(v_i, v_j) \in R_G\}$: a relation corresponding to $R_G$ for any regular expression $R$
- $SCC(V, S) = \{(v_i, s_j)|v_i \in s_j\}$: a relation that represents the relationship between each vertex of $G_R$ and the SCC containing the vertex
- $\overline{R_G^+}(START\_S, END\_S) = \{(s_i, s_j)|(\overline{v_i}, \overline{v_j}) \in \overline{R_G^+}\}$: a relation corresponding to $\overline{R_G^+}$, the transitive closure of $\overline{G_R}$

Using the notation above, we represent an equivalent (1) as Lemma 4.

**Lemma 4 .**

$$(A \cdot B)_G(START\_V, END\_V)$$
$$=\pi_{A_G.START\_V, B_G.END\_V}\Big[A_G(START\_V, END\_V) \quad (1)$$
$$\bowtie_{A_G.END\_V=B_G.START\_V} B_G(START\_V, END\_V)\Big]$$

*Proof.* Because of space limitations in proof, we omit the attribute names of the relations.

- $(A \cdot B)_G \subseteq \pi_{A_G.START\_V, B_G.END\_V}(A_G \bowtie_{A_G.END\_V=B_G.START\_V} B_G)$. Suppose $(A \cdot B)_G \ni (v_i, v_j)$. Then, there exists a path $p(v_i, v_j)$ satisfying $A \cdot B$ on $G$. $p(v_i, v_j)$ can be decomposed into a path $p(v_i, v_k)$ satisfying $A$ and a path $p(v_k, v_j)$ satisfying $B$. Because $(v_i, v_k) \in A_G$ and $(v_k, v_j) \in B_G$, $(v_i, v_j) \in \pi_{A_G.START\_V, B_G.END\_V}(A_G \bowtie_{A_G.END\_V=B_G.START\_V} B_G)$.
- $(A \cdot B)_G \supseteq \pi_{A_G.START\_V, B_G.END\_V}(A_G \bowtie_{A_G.END\_V=B_G.START\_V} B_G)$. Suppose $\pi_{A_G.START\_V, B_G.END\_V}(A_G \bowtie_{A_G.END\_V=B_G.START\_V} B_G) \ni (v_i, v_j)$. Then, there must exist $(v_i, v_k)$ and $(v_k, v_j)$ where $(v_i, v_k) \in A_G$, $(v_k, v_j) \in B_G$. Because there exists a path $p(v_i, v_k)$ satisfying $A$ and a path $p(v_k, v_j)$ satisfying $B$ on $G$, there must exist a $p(v_i, v_j)$ satisfying $A \cdot B$. Therefore, $(v_i, v_j) \in (A \cdot B)_G$ holds. □

Theorem 2 shows a relational algebra expression for evaluation of RPQ $R^+$ on $G$ using the RTC. The result of RPQ $R^+$ on $G$ is represented as a relational algebra expression in Theorem 2. Equation (2) represents the union of Cartesian products of $s_i$ and $s_j$ for every vertex pair$(\overline{v_i}, \overline{v_j})$(i.e., an element in the RTC). That is, we can efficiently evaluate $R^+$ by using the RTC.

**Theorem 2.** *The result of RPQ $R^+$ on $G$ is represented as a relational algebra expression in (2) for any RPQ $R$. Here, $\rho_{SSCC}$ and $\rho_{ESCC}$ are renaming operations.*

$$R_G^+(START\_V, END\_V)$$
$$=\pi_{SSCC.V, ESCC.V}\Big[\rho_{SSCC}(SCC(V, S))$$
$$\bowtie_{S=START\_S} \overline{R_G^+}(START\_S, END\_S) \quad (2)$$
$$\bowtie_{END\_S=S} \rho_{ESCC}(SCC(V, S))\Big]$$

*Proof.* Theorem 1 can be represented as the following tuple relational calculus expression.

$$R_G^+ = \{res \mid (\exists rtc)(\exists sscc)(\exists escc)\ \overline{R_G^+}(rtc) \text{ AND}$$
$$SCC(sscc) \text{ AND } SCC(escc) \text{ AND}$$
$$sscc[S] = rtc[START\_S] \text{ AND } rtc[END\_S] = escc[S] \text{ AND}$$
$$/* (v_i, v_j) \in s_k \times s_l */$$
$$res[START\_V] = sscc[V] \text{ AND } res[END\_V] = escc[V]\}$$
$$/* \text{ projection } */$$

Equation (2) is derived from the above relational calculus expression in a straightforward manner [16]. □

We now expand the result of the batch unit to eliminate *useless* and *redundant operations*. Using Lemma 4 and (1), the result of the batch unit can be represented as in (3) to (5). Using (2) we can expand (4) as in (7) to (9). When evaluating (4), we can eliminate *useless-1 operations* by using $\overline{R_G^+}$ instead of the entire $R_G^+$ and evaluating $R^+$ starting only from tuples (i.e., vertex pairs) in $Pre_G(START\_V, END\_V)$.

$$(Pre \cdot R^+ \cdot Post)_G(START\_V, END\_V)$$
$$=\pi_{Pre_G.START\_V, Post_G.END\_V}\Big[Pre_G(START\_V, END\_V) \quad (3)$$
$$\bowtie_{Pre_G.END\_V=R_G^+.START\_V} R_G^+(START\_V, END\_V) \quad (4)$$
$$\bowtie_{R_G^+.END\_V=Post_G.START\_V} Post_G(START\_V, END\_V)\Big] \quad (5)$$

$$(Pre \cdot R^+ \cdot Post)_G(START\_V, END\_V)$$
$$=\pi_{Pre_G.START\_V, Post_G.END\_V}\Big[Pre_G(START\_V, END\_V) \quad (6)$$
$$\bowtie_{END\_V=V} SCC(V, S) \quad (7)$$
$$\bowtie_{S=START\_S} \overline{R_G^+}(START\_S, END\_S) \quad (8)$$
$$\bowtie_{END\_S=S} SCC(V, S) \quad (9)$$
$$\bowtie_{V=START\_V} Post_G(START\_V, END\_V)\Big] \quad (10)$$

To eliminate *redundant-1* and *redundant-2 operations*, we union the intermediate results, i.e., the RTCs, at each join step. For each vertex pair $(v_i, v_j)$ in $Pre_G$, (7) finds the SCC $s_k$ containing $v_j$ and returns $(v_i, s_k)$ as the result. The intermediate results of (7) are duplicated for vertex pairs whose start vertices are the same and end vertices belong to the same SCC. Thus, we can eliminate *redundant-1 operations* in the next join step by unioning (i.e., eliminating duplicates in) the intermediate results of (7). For each pair of vertex and SCC $(v_i, s_k)$ in $Pre_G \bowtie SCC$, (8) finds the SCC $s_l$ that is reachable from $s_k$ and returns $(v_i, s_l)$ as the result. If, in the results of (7), there exist paths from vertex pairs whose start vertices are the same but end vertices belong to different SCCs to vertices belonging to the same SCC, the intermediate results of (8) for those paths are duplicated. Therefore, we find unique SCCs containing end vertices of the paths satisfying $R^+$ by unioning (i.e., eliminating duplicates in) the intermediate results of (8). Thus, we can eliminate subsequent *redundant-2 operations*. Because sets of vertices belonging to different SCCs are

mutually disjoint, if the results of (8) are distinct, there are no duplicates in the result of (9). That is, an operation that checks duplicates for union when performing (9) is useless. We call them *useless-2 operations*. We can eliminate *useless-2 operations* by not performing duplicate checks.

**Example 9.** Fig. 9 shows an example of evaluating $Pre \cdot R^+$ where *useless-1*, *redundant-1*, and *redundant-2 operations* are eliminated. We evaluate $R^+$ starting from vertex pairs in $Pre_G$ using $\overline{R^+_G}$ instead of $R^+_G$. Therefore, the paths satisfying $R^+$ from the vertices belonging to $s_{n+1}$ to the vertices belonging to $s_k$ and $s_{k+1}$ are not found. Thus, *useless-1 operations* are eliminated.

We next eliminate duplicates of the intermediate results of (7). Therefore, vertex pairs in (6) whose start vertices are the same and end vertices belong to the same SCC become one pair of vertex and SCC in (7). That is, $(v_0, v_1), \cdots, (v_0, v_n)$ in (6) become $(v_0, s_0)$ in (7). Therefore, when evaluating $R^+$ starting from $(v_0, v_1), \cdots, (v_0, v_n)$ in $Pre_G$, we find $v_k, \cdots, v_j$ only once by going through $s_0$. Thus, *redundant-1 operations* are eliminated.

We also eliminate duplicates of the intermediate results of (8). Therefore, if there exist paths from vertex pairs whose start vertices are the same but end vertices belong to different SCCs to vertices belonging to a same SCC, those vertex pairs in (6) become one pair of vertex and SCC in (8). That is, $(v_0, v_1)$ and $(v_0, v_m)$ in (6) become $(v_0, s_m)$ in (8). Therefore, when evaluating $Pre \cdot R^+$ starting from $v_0$ via each vertex belonging to $s_0$ and $s_n$, respectively, we find $v_k, \cdots, v_j$ only once by going through the unique $s_m$. Thus, subsequent *redundant-2 operations* are eliminated. $\square$

Algorithm 2 shows EvalBatchUnit that evaluates a batch unit in an optimized way of eliminating *useless* and *redundant operations* explained so far. First, in line 1 to 3, we initialize variables storing the results of (7) to (10). Here, we deal with $Pre \cdot R^* \cdot Post$ by initializing $ResEq9$ with $Pre_G$. Then, in line 4, we eliminate *useless-1 operations*. By line 4, lines 5 to 12 that evaluate $R^+$ are executed only for vertex pairs existing in $Pre_G$. Otherwise, these operations are *useless-1 operations*. In lines 5 to 7, we eliminate *redundant-1 operations*. In line 5, for each
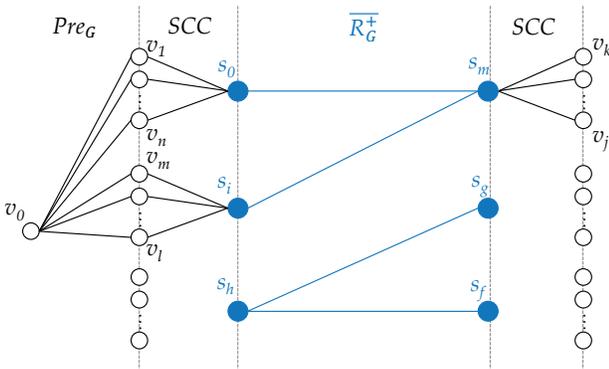


Fig. 9: An example of evaluating $Pre \cdot R^+$ eliminating *redundant* and *useless operations*.

---

**Algorithm 2:** EvalBatchUnit.

**Input:** $Pre_G$, $\overline{R^+_G}$, $SCC$, $Type$, $Post$
**Output:** $ResEq10$

/* Initialize variables storing the results of (7) to (10)    */
1  $ResEq7$, $ResEq8$, $ResEq9$, $ResEq10 \leftarrow \emptyset$
2  **if** $Type$ *is* * **then**
3  |  $ResEq9 \leftarrow Pre_G$          // Initialization for $Pre \cdot R^* \cdot Post$

/* Compute $Pre_G \bowtie R^+_G$: (7) to (9)    */
4  **foreach** $(v_i, v_j) \in Pre_G$ **do**
   |  /* Eliminate *useless-1 ops.*    */
5  |  $s_j \leftarrow \pi_S(\sigma_{V = v_j} SCC)$
6  |  **if** $(v_i, s_j) \notin ResEq7$ **then**    // duplicate check for (7)
7  |  |  Insert $(v_i, s_j)$ into $ResEq7$    // union the result of (7)
   |  |  /* Eliminate *redundant-1 ops.*    */
8  |  |  **foreach** $(s_j, s_k) \in \sigma_{START\_S = s_j} \overline{R^+_G}$ **do**
9  |  |  |  **if** $(v_i, s_k) \notin ResEq8$ **then**    // duplicate check for (8)
10 |  |  |  |  Insert $(v_i, s_k)$ into $ResEq8$    // union the result of (8)
   |  |  |  |  /* Eliminate *redundant-2 ops.*    */
11 |  |  |  |  **foreach** $(s_k, v_k) \in \sigma_{S = s_k} SCC(S, V)$ **do**
   |  |  |  |  |  /* Eliminate *useless-2 ops.*    */
12 |  |  |  |  |  Insert $(v_i, v_k)$ into $ResEq9$    // add the result of (9)

/* Compute $(Pre \cdot R^+)_G \bowtie Post_G$: (10)    */
13 **foreach** $(v_i, v_k) \in ResEq9$ **do**
14 |  **foreach** $(v_k, v_l) \in$ EvalRestrictedRPQ($Post$, $v_k$) **do**
15 |  |  **if** $(v_i, v_l) \notin ResEq10$ **then**    // duplicate check for (10)
16 |  |  |  Insert $(v_i, v_l)$ into $ResEq10$    // union the result of (10)

---

vertex pair $(v_i, v_j)$ we find the SCC $s_k$ of $G_R$ containing $v_j$. In lines 6 and 7, we union $(v_i, s_j)$ into $ResEq7$. Then, lines 8 to 12 are executed only when $(v_i, s_j)$ does not exist in $ResEq7$. Otherwise, these operations are *redundant-1 operations*. In lines 8 to 10, we eliminate *redundant-2 operations*. In line 8, we find $s_k$, which is reachable from $s_j$. In lines 9 and 10, we union $(v_i, s_k)$ into $ResEq8$. Lines 11 and 12 are executed only when $(v_i, s_k)$ does not exist in $ResEq8$. Otherwise, these operations are *redundant-2 operations*. In lines 11 and 12, we eliminate *useless-2 operations*. Here, we add $(v_i, v_k)$ into $ResEq9$ for each vertex $v_k$ contained in $s_k$ without duplicate checks, which are *useless-2 operations*. Lines 13 to 16 that evaluate $Post$ are executed only for vertex pairs existing in $ResEq9$ (i.e., $(Pre \cdot R^+)_G$). EvalRestrictedRPQ($Post$, $v_k$) finds paths satisfying $Post$ from the vertex $v_k$ on $G$ and returns $(v_k, v_l)$ for each path $p(v_k, v_l)$.

## V. PERFORMANCE EVALUATION

In this section, we focus on evaluating multiple RPQs in the form of the batch unit. With this, we can easily control parameters synthetically to generate complete test cases.

### A. Experimental Environment

We use synthetic datasets generated by the RMAT model [17] using TrillionG [18] and four real datasets: Yago2s [19], Robots [20], Advogato [21], and Youtube_Sampled [22]. Using TrillionG we generate synthetic graphs of various average vertex degrees per label (i.e., $\frac{|E|}{|V||\Sigma|}$) where the other characteristics except the degree

are kept the same. Since TrillionG generates edge-unlabeled, directed multigraphs, we randomly added a label to each edge to make edge-labeled graphs. We denote the RMAT graph with $2^{13}$ vertices and $2^{N+13}$ edges by RMAT_$N$. Yago2s, Robots, and Advogato are edge-labeled, directed multigraphs. Youtube_Sampled is a subset of Youtube, which is an edge-labeled, undirected multigraph. Since the data subject to RPQs is a directed graph, we randomly added a direction to each edge of Youtube_Sampled. We construct Youtube_Sampled from Youtube using random vertex sampling. $V$ is the set of randomly sampled vertices, and $E$ is the set of edges between sampled vertices. We simply denote Youtube_Sampled by Youtube. TABLE IV summarizes the statistics of the datasets.

We use synthetic multiple RPQ sets where each RPQ is in the form of the batch unit. To create a controlled environment, we simulate the effects of *Pre* and *Post* using single labels and model $R$ as a concatenation of labels in $\Sigma$ (i.e., a clause without Kleene closure) whose length varies from 1 to 3. First, we randomly select a total of 90 $R$s (one for each multiple RPQ set), 10 for each length, and then, randomly select from $\Sigma$ pairs of *Pre* and *Post* for each $R$. The number of RPQs in each multiple RPQ set is 1, 2, 4, 6, 8, and 10, and a larger multiple RPQ set contains smaller multiple RPQ sets.

The multiple RPQ evaluation methods to be tested for comparison are as follows. Since the source codes for *NoSharing* and *FullSharing* have not been released, we implement them ourselves based on the literature [5], [8].

- *RTCSharing*(*RTC*): A method sharing $\overline{R_G^+}$ among RPQs.
- *NoSharing*(*No*): A method individually evaluating RPQs using the single RPQ evaluation method proposed by Yakovets et al. [5].
- *FullSharing*(*Full*): A method sharing $R_G^+$ among RPQs proposed by Abul-Basher [8]

All experiments have been conducted on a Linux (kernel version: 2.6.32) machine equipped with an Intel Core i7-7700 CPU and 64GB main memory. All the multiple RPQ evaluation methods used in the experiment including *RTCSharing* have been implemented in C++.

### B. Performance Evaluation

In the experiment, we compare the performance of evaluating multiple RPQs whose common sub-query is a Kleene plus $R^+$. Section V-B1 compares the performance as we vary the average vertex degrees per label (i.e., $\frac{|E|}{|V||\Sigma|}$). In Section V-B2, we vary the number of RPQs constituting a multiple RPQ set to see the amortization effect of sharing the data, $R_G^+$ or $\overline{R_G^+}$, among RPQs.

TABLE IV: Statistics of datasets used in the experiments.

| Dataset | | $|V|$ | $|E|$ | $|\Sigma|$ | $\frac{|E|}{|V||\Sigma|}$ |
|---|---|---|---|---|---|
| Real graph datasets | Yago2s | 108,048,761 | 244,796,155 | 104 | 0.02 |
| | Robots | 1,725 | 3,596 | 4 | 0.52 |
| | Advogato | 6,541 | 51,127 | 3 | 2.61 |
| | Youtube | 1,600 | 91,343 | 5 | 11.42 |
| Synthetic graph datasets | RMAT_$N$ ($N = 0..6$) | $2^{13}$ | $2^{N+13}$ | 4 | $2^{N-2}$ |

The performance metrics are multiple RPQ sets' average query response time (*query response time* in short) and average size of data shared among RPQs (*shared data size* in short). The query response time includes the time taken 1) to construct the two-level reduced graph in *RTCSharing*, 2) to compute the shared data ($\overline{R_G^+}$ in *RTCSharing* or $R_G^+$ in *FullSharing*), and 3) to complete the evaluation of all RPQs. We also divide the query evaluation into three parts and compare each part between *RTCSharing* and *FullSharing* to show the effects of different aspects of *RTCSharing* on performance. First, to show that *RTCSharing* is simpler than *FullSharing* for computing the data shared among RPQs, we compare the computation time of $\overline{R_G^+}$ in *RTCSharing* with that of $R_G^+$ in *FullSharing*. This computation time is denoted by *Shared_Data*. Since the two methods compute $R_G$ identically, we exclude the computation time of $R_G$ from Shared_Data of both methods. Second, to show the effects of avoiding *redundant* and *useless operations* by representing each RPQ as a relational algebra expression and optimizing its evaluation, we compare the computation time of $Pre_G(START\_V, END\_V) \bowtie R_G^+(START\_V, END\_V)$ (denoted by $Pre_G \bowtie R_G^+$) of each method. The only difference between two methods in $Pre_G \bowtie R_G^+$ is the optimization related to *useless* and *redundant operations*. Finally, we compare the computation time of the remainder evaluation (i.e., computing $Pre_G$, $R_G$, and computing $(Pre \cdot R^+ \cdot Post)_G$ from $(Pre \cdot R^+)_G$) for which the comparison methods operate identically. Since *NoSharing* does not share data among RPQs, for *NoSharing*, we present only the overall query response time. The shared data size is the number of pairs in $\overline{R_G^+}$ for *RTCSharing* or that for *FullSharing*. Since the size of data shared among RPQs is not dependent on the number of RPQs, we present it only in Section V-B1.

*1) Experiment 1: Average vertex degree per label is varied.:* We experiment the graphs having different average vertex degree per label (*vertex degree* in short). We use synthetic and real datasets described in TABLE IV. We use multiple RPQ sets consisting of 4(median) RPQs among the multiple RPQ sets generated as described in Section V-A.

**Query response time** Fig. 10 shows (a) the query response time on the synthetic datasets and (b) the normalized query response time on the real datasets. Not only the vertex degree but also the size of the dataset (the number of vertices and the number of edges) affect the performance. When experimenting
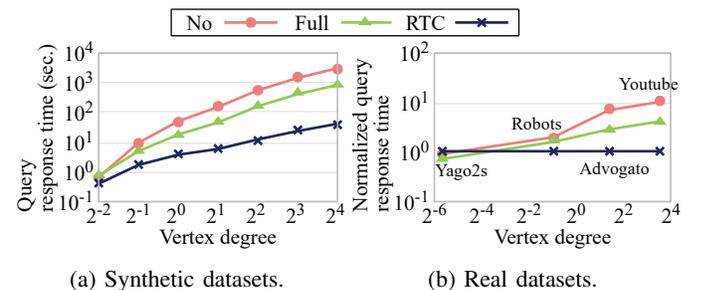


(a) Synthetic datasets.   (b) Real datasets.

Fig. 10: Computational performances of *No*, *Full*, and *RTC* as the vertex degree is varied ($\sharp$ RPQs = 4).

(a) Synthetic datasets.
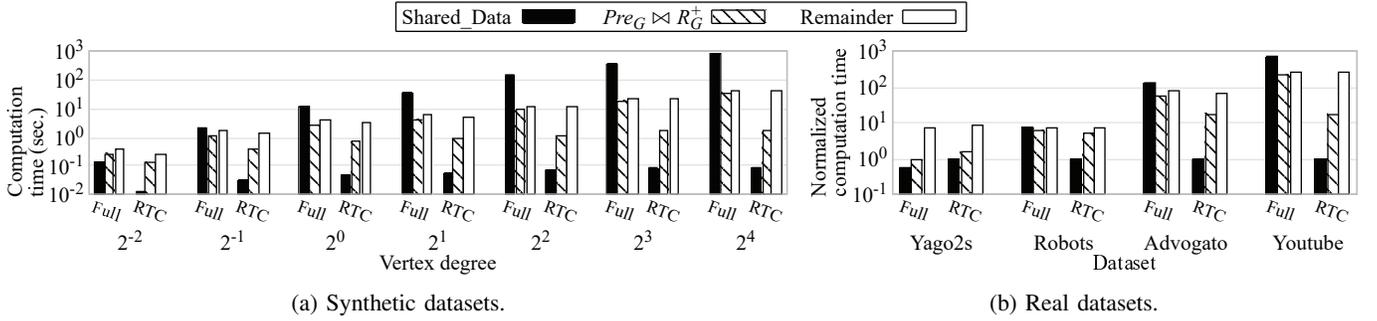
(b) Real datasets.

Fig. 11: Computation time of three parts of *Full* and *RTC* as the vertex degree is varied ($\sharp$ RPQs = 4).

on different real datasets, it is not easy to show the effect of the vertex degree only because the sizes of the datasets are different. To remedy this problem, we compare the relative performance by normalizing the query response time of the comparison methods by that of *RTCSharing* on the real datasets. As shown in Fig. 10, *RTCSharing* improves the query response time by 1.63 to 20.20 over *FullSharing* on all graphs except Yago2s.

Fig. 11 shows the computation time of three detailed parts. The figure confirms that *RTCSharing* beats *FullSharing* through the improvement of Shared_Data and $Pre_G \bowtie R_G^+$. *RTCSharing* improves Shared_Data by 7.78 to 9013.64 times and $Pre_G \bowtie R_G^+$ by 1.30 to 19.11 times over *FullSharing* on all graphs except Yago2s. Shared_Data is improved because the computation of $\overline{R_G^+}$ is simpler than that of $R_G^+$. $Pre_G \bowtie R_G^+$ is improved because *RTCSharing* eliminates *redundant-1*, *redundant-2* and *useless-1 operations* by processing and optimizing each RPQ as a relational algebra expression. When we compare *RTCSharing* with *NoSharing*, *RTCSharing* significantly improves the performance by up to 73.86 times over *NoSharing* on all graphs except Yago2s. The reason is that *NoSharing* computes $R_G^+$ repeatedly for each RPQ while *RTCSharing* avoids the overhead by sharing $\overline{R_G^+}$ among RPQs.

For Yago2s, however, Fig. 10(b) shows that *RTCSharing* is up to 1.36 times slower than *FullSharing* and 1.07 times slower than *NoSharing* in the query response time. Fig. 11(b) shows that it is slower 1.82 times in Shared_Data; 1.88 times in $Pre_G \bowtie R_G^+$. Yago2s is an exceptional case, in which vertex degree is extremely small, 0.02, and the average number of vertices in an SCC of $G_R$ is 1.00, which means that the vertex-level reduction is not very effective. Thus, $G_R$ and $\overline{G_R}$ are similar (almost the same) in size, so that the computation times of $R_G^+$ and $\overline{R_G^+}$ also are similar. However, *RTCSharing* has the overhead of reducing the graph, and thus, *RTCSharing* is slower than *FullSharing* in Shared_Data. *RTCSharing* is slower in $Pre_G \bowtie R_G^+$ because the average number of vertices in an SCC of $G_R$ is 1.00, and there are few *redundant-1* and *redundant-2 operations* that *RTCSharing* can eliminate; on the other hand, *RTCSharing* has additional join overheads.

In Figs. 10 and 11, we note that, as the vertex degree increases, *RTCSharing* improves the performance gradually. In Figs. 10(a) and 11(a) on the synthetic datasets, as the vertex degree increases from $2^{-2}$ to $2^4$, the ratio of the query response time of *FullSharing* over *RTCSharing* increases from

1.88 to 20.20; that of Shared_Data from 10.40 to 9013.64; that of $Pre_G \bowtie R_G^+$ from 2.03 to 19.11. Also in Figs. 10(b) and 11(b) on the real datasets, as the vertex degree increases from 0.02 to 0.52, 2.61, and 11.42, the ratio of the query response time of *FullSharing* over *RTCSharing* increases from 0.74 to 1.63, 2.92, and 4.20; that of Shared_Data from 0.55 to 7.78, 129.40, and 671.86; that of $Pre_G \bowtie R_G^+$ from 0.53 to 1.30, 2.95, and 11.97. The reasons are as follows. As the vertex degree increases, the average number of vertices of $G_R$ reduced to one vertex of $\overline{G_R}$ increases, so that the size of reduced graph $\overline{G_R}$ decreases. Thus, the computation time for $\overline{R_G^+}$ becomes smaller, and the ratio of *FullSharing* over *RTCSharing* also becomes large in Shared_Data. In addition, since the number of vertices in each SCC of $G_R$ increases, the number of *redundant-1* and *redundant-2 operations* that are eliminated in *RTCSharing* increases, and the ratio becomes larger in $Pre_G \bowtie R_G^+$. As a result, the ratio becomes larger in the query response time as well. We note, however, that the ratio is smaller in the query response time than in Shared_Data or $Pre_G \bowtie R_G^+$. The reason is as follows. As the vertex degree increases, the length of the path satisfying the RPQ becomes longer so that the number of final results increases, making the time required to find $(Pre \cdot R^+ \cdot Post)_G$ from $(Pre \cdot R^+)_G$ increase. Therefore, the portion of Remainder, which is largely the same in both methods, in the query response time increases, lowering the ratio of the query response time. When we compare *RTCSharing* with *NoSharing*, as the vertex degree increases, the ratio of *NoSharing* over *RTCSharing* in the query response time increases from 1.68 to 73.86 on the synthetic datasets and from 0.93 to 2.03, 7.13, and 10.45 on the real datasets. The reason is that the overhead of repeated computation of $R_G^+$ increases in *NoSharing* as the vertex degree increases.

**Shared data size** Fig. 12 shows (a) the shared data size on the synthetic datasets and (b) the normalized shared data size on the real datasets. As the vertex degree increases, the size of $R_G^+$ in *FullSharing* increases. On the other hand, even as the vertex degree increases, the size of $\overline{R_G^+}$ in *RTCSharing* does not increase much. The reason is that as the vertex degree increases, the average number of vertices of $G_R$ reduced to one vertex of $\overline{G_R}$ increases, so that the size of reduced graph $\overline{G_R}$ decreases (see Fig. 13). As a result, as the vertex degree increases from $2^{-2}$ to $2^4$ on the synthetic datasets, the shared
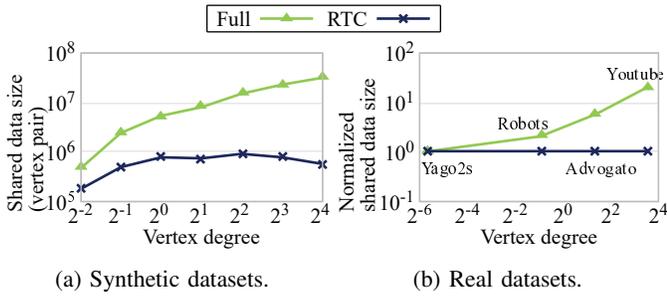
(a) Synthetic datasets.

(b) Real datasets.

Fig. 12: Space performances of *Full* ($R_G^+$) and *RTC* ($\overline{R_G^+}$) for datasets as the vertex degree is varied (♯ RPQs = 4).
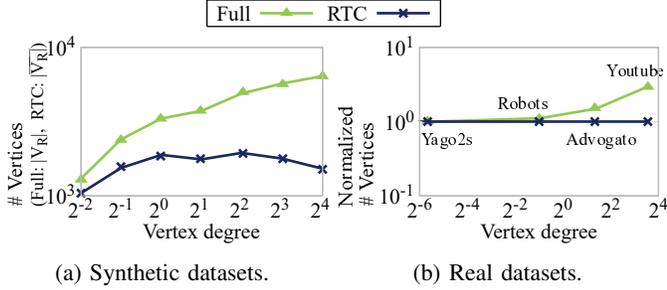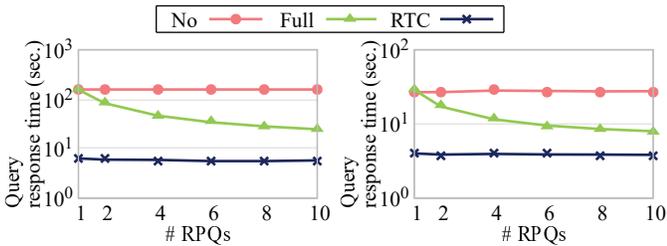


(a) Synthetic datasets.

(b) Real datasets.

Fig. 13: The number of vertices of *Full* ($|V_R|$) and *RTC* ($|\overline{V_R}|$) for datasets as the vertex degree is varied (♯ RPQs = 4).

data size of *FullSharing* over *RTCSharing* increases from 2.61 to 54.94. On real datasets, as the vertex degree increases from 0.02 to 0.52, 2.61, and 11.42 on the real datasets, that increases from 1.05 to 2.09, 5.87, and 20.23. The reason that the Yago2s' ratio is greater than 1.00 (the number of vertices in an SCC of $G_R$) is that Yago2s has two exceptional *R*s, which have high vertex degrees. for them, the size of $R_G^+$ over $\overline{R_G^+}$ is only 1.05.

*2) Experiment 2: The number of RPQs is varied.:* In this experiment, we vary the number of RPQs constituting each multiple RPQ set. We use RMAT_3 and Advogato, which have a median vertex degree among synthetic and real datasets, respectively. We use multiple RPQ sets consisting of 1, 2, 4, 6, 8 and 10 RPQs generated as described in Section V-A.

Fig. 14 shows the query response time on (a) a synthetic dataset (RMAT_3) and (b) a real dataset (Advogato) as the number of RPQs is varied. Fig. 15 shows the computation time of each part of *RTCSharing* and *FullSharing*. As the
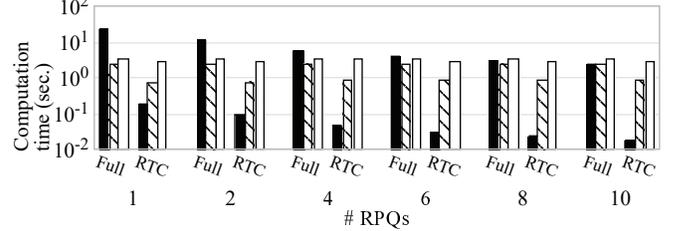


(a) Synthetic dataset (RMAT_3, vertex vertex degree = 2).

(b) Real dataset (Advogato, vertex degree = 2.61).

Fig. 14: Computational performances of *No*, *Full*, and *RTC* as the number of RPQs is varied.



(a) Synthetic dataset (RMAT_3, vertex degree = 2).



(b) Real dataset (Advogato, vertex degree = 2.61).

Fig. 15: Computation time of three parts of *Full* and *RTC* as the number of RPQs is varied.

number of RPQs increases, Shared_Data in both *RTCSharing* and *FullSharing* decreases. This is because the computation time of $R_G^+$ or $\overline{R_G^+}$ is amortized by the number of RPQs. In the case of *RTCSharing*, the amortization effect is not large since the portion of Shared_Data in the total time is much smaller than that of *FullSharing* due to efficiency of *RTCSharing* (see Fig. 15). As a result, as the number of RPQs increases from 1 to 10, the ratio of the query response time of *FullSharing* over *RTCSharing* decreases from 24.35 to 4.25 on the synthetic dataset and from 7.17 to 2.08 on the real dataset. The ratio of the query response time of *NoSharing* over *RTCSharing* slightly increases from 23.11 to 25.38 on the synthetic dataset and from 6.76 to 7.17 on the real dataset. This gradual improvement also comes from the amortization effect of Shared_Data in *RTCSharing* despite that the portion of Shared_Data in the query response time is small.

## VI. RELATED WORK

**Transitive closure** Purdom [12] and Nuutila [13] proposed algorithms for computing transitive closure by using SCCs, which essentially implemented Lemma 3, but without formalization. Since they are only for an unlabeled graph, they alone cannot be used to evaluate an RPQ. If combined with edge-level reduction, they can be used for evaluating the entire $R^+$. However, they cannot efficiently process the RPQ *Pre·$R^+$* due to incurring *useless and redundant operations*. Moreover, for multiple RPQs evaluation, they can only facilitate *FullSharing* and cannot directly support *RTCSharing*.

**Reachability query** The result of the RPQ $R^+$ on $G$ is the same with the result of the reachability query on $G_R$. The main costs of this query are the index construction time, the index size, and the query response time. Some of recent approaches [23], [24] answer the query using the index only. Some other recent approaches [25], [26] traverse a graph at

run-time if needed. If combined with edge-level reduction, those methods can also be used for evaluating the entire $R^+$. However, since all pairs of vertices $O(|V|^2)$ should be checked, they cannot be efficient methods. In addition, they cannot avoid *redundant* and *useless operations*.

**Evaluation methods of single RPQ queries** A few methods [4], [10] tried to reduce the number of edges unnecessarily accessed when traversing the graph. The method proposed by Koschmieder and Leser [10] compares the number of edges of each label that is present in the given RPQ and traverses the edge having the label that has the smallest number of edges first. Then, it continues to traverse other edges only when they are connected to those already traversed. Nguyen and Kim [4] considered not only the number of edges of each label but also the number of unique start and end vertices of the edges. Yannakakis [27] proposed a method that reduced the problem of evaluating an RPQ to a problem of finding the paths between given vertices by constructing a graph. This method is entirely different from the edge-level reduced graph used in this paper. Yakovets et al. [5] proposed a cost model for finding the optimal RPQ evaluation order and a method to efficiently evaluate the Kleene star $R^*$. The method is similar to our edge-level graph reduction in that it traverses pre-evaluated paths rather than edges, but does not consider vertex-level graph reduction nor evaluation of multiple RPQs. Fletcher et al. [9] and Tetzel et al. [28] proposed methods to evaluate RPQs using indexes. The former finds all the paths in the graph of the lengths less than a certain threshold in advance, creates an index using the path labels of these paths as the keys, and uses the index to evaluate an RPQ. The latter compares the performance of the method using the compressed index and the one using the uncompressed index. Pacaci et al. [29] focused on the evaluation over streaming graphs. All of these existing methods did not consider *redundant* and *useless operations* and focused mainly on single RPQ evaluation.

**Evaluation methods of multiple RPQs** Abul-Basher [8] proposed an optimization technique for evaluating multiple RPQs. This method finds a common sub-query of given multiple RPQs, evaluates it first, and shares the results among the RPQs to avoid repeated computations. This is the *FullSharing* method we used for comparison in Section V. However, as explained in Section V, when the common sub-query is $R^+$, its evaluation is costly. There are also *redundant* and *useless operations* when evaluating each RPQ using $R_G^+$.

## VII. Conclusions

In this paper, we propose a notion of *RPQ-based graph reduction* that replaces the evaluation of the Kleene closure on the large original graph $G$ to that of the transitive closure to the small graph $\overline{G_R}$. We showed that $R_G^+$ can be easily calculated from the transitive closure of $\overline{G_R}$ (i.e., RTC), which is computationally simpler and smaller than $R_G^+$ and $R_G^*$, in Theorem 1. We also proposed an RPQ evaluation algorithm, *RTCSharing*, that takes advantage of this RTC. *RTCSharing* treats each clause in the DNF of the given RPQ as a batch unit that is in the form of *Prefix·$R^+$·Postfix* or *Prefix·$R^*$·Postfix*. In

*RTCSharing*, we represent the batch unit as a relational algebra expression (join sequence) including the RTC and efficiently evaluate it sharing the RTC among batch units. We also eliminate *useless-1 operations* by evaluating $R^+$ only starting from vertex pairs in *Prefix$_G$*, *redundant-1* and *redundant-2 operations* by unioning on the intermediate result of each join step, and *useless-2 operations* by using the mutually disjoint property of SCCs. We formally explain that *useless-1*, *redundant-1*, and *redundant-2 operations* are caused by having *Prefix*, and *useless-2 operation* is caused by structural property of $\overline{R_G^+}$. Experiments using synthetic and real datasets show that *RTCSharing* significantly improves the performance by up to 73.86 times over *NoSharing* and 20.20 times over *FullSharing* in terms of average query response time.

## References

[1] Mendelzon, O. and Wood, P., "Finding regular simple paths in graph databases," *SIAM Journal on Computing*, Vol. 24, No. 6, pp. 1235-1258, 1995.

[2] Harris, S. and Seaborne, A., SPARQL 1.1 query language, W3C, Mar. 2013.

[3] Cypher, http://neo4j.com/ accessed on Dec. 20th, 2018.

[4] Nguyen, V. and Kim, K., "Efficient regular path query evaluation by splitting with unit-subquery cost matrix," *IEICE Transactions on Information and Systems*, Vol. 100, No. 10, pp. 2648-2652, 2017.

[5] Yakovets, N., Godfrey, P., and Gryz, J., "Query planning for evaluating SPARQL property paths," In *Proc. ACM Int'l Conf. on Management of Data (SIGMOD)*, pp. 1875-1889, San Francisco, California, June-July 2016.

[6] Grahne, G. and Thomo, A. "Query containment and rewriting using views for regular path queries under constraints," In *Proc. ACM Symposium on Principles of Database Systems (PODS)*, pp. 111-122, San Diego, California, June 2003.

[7] Grahne, G. and Thomo, A. "New rewritings and optimizations for regular path queries," In *Proc. IEEE Int'l Conf. on Database Theory (ICDT)*, pp. 242-258, Siena, Italy, Jan. 2003.

[8] Abul-Basher, Z., "Multiple-query optimization of regular path queries," In *Proc. IEEE Int'l Conf. on Data Engineering(ICDE)*, pp. 1426-1430, San Diego, California, Apr. 2017.

[9] Fletcher, G., Peters, J., and Poulovassilis, A., "Efficient regular path query evaluation using path indexes," In *Int'l Conf. on Extending Database Techonology(EDBT)*, pp. 636-639, Bordeaux, France, Mar. 2016.

[10] Koschmieder, A. and Leser, U., "Regular path queries on large graphs," In *Proc. Int'l Conf. on Scientific and Statistical Database Management(SSDBM)*, pp. 177-194, Chania, Greece, June 2012.

[11] Yakovets, N., Godfrey, P., and Gryz, J., "Waveguide: Evaluating SPARQL Property Path Queries," In *Proc. Int'l Conf. on Extending Database Techonology(EDBT)*, pp. 525-528, Brussels, Belgium, Mar. 2015.

[12] Purdom, P., "A transitive closure algorithm," *BIT Numerical Mathematics*, Vol. 10, No. 1, pp. 76-94, 1970.

[13] Nuutila, E., "An efficient transitive closure algorithm for cyclic digraphs," *Information Processing Letters*, Vol. 52, No. 4, pp. 207-213, Nov. 1994.

[14] Tarjan, R., "Depth first search and linear graph algorithms," *SIAM Journal on Computing*, Vol. 1, No. 2, pp. 146-160, 1972.

[15] Davey, A. B. and Priestley, A. H., *Introduction to Lattices and Order*, Cambridge University Press, 1990.

[16] Ullman, J., *Principles of database and knowledge-base systems*, Vol. 1, Computer Science Press, 1988.

[17] Chakrabarti, D., Zhan, Y., and Faloutsos, C., "R-MAT: A recursive model for graph mining," In *Proc. SIAM SDM*, pp.442-446, 2004.

[18] Park, H and Kim, M., "TrillionG: A trillion-scale synthetic graph generator using a recursive vector model," In *Proc. ACM SIGMOD*, pp. 913-928, 2017.

[19] Yago2s dataset, https://www.mpi-inf.mpg.de/departments/databases-and-information-systems/research/yago-naga/yago/downloads/ accessed on Dec. 2018.

[20] Robots dataset, http://www.trustlet.org/datasets/robots_net/ accessed on Dec. 2018.

[21] Advogato network dataset, http://konect.uni-koblenz.de/networks/advogato accessed on Dec. 2018.

[22] Youtube dataset, http://socialcomputing.asu.edu/pages/datasets accessed on Dec. 2018.

[23] Jin, R. and Wang, G., "Simple, fast, and scalable reachability oracle," *Proceedings of the VLDB Endowment*, Vol. 6, No. 14, pp. 1978-1989, 2013.

[24] Valstar, L., Fletcher, G., and Yoshida, Y., "Landmark indexing for evaluation of label-constrained reachability queries," In *Proc. ACM Int'l Conf. on Management of Data (SIGMOD)*, pp. 345-358, Chicago, Illinois, May 2017.

[25] Wei, H., Yu, J. X., Lu, C., and Jin, R., "Reachability Querying: An Independent Permutation Labeling Approach," *The VLDB Journal*, Vol. 27, No. 1, pp. 1-26, 2018.

[26] Su, J., Zhu, Q., Wei, H., and Yu, J. X. "Reachability Querying: Can it be even faster?," *IEEE Transactions on Kwowledge and Data Engineering*, Vol. 29, No. 3, pp. 683-697, 2017.

[27] Yannakakis, M., "Graph-theoretic methods in database theory," In *Proc. ACM Symposium on Principles of Database Systems (PODS)*, pp. 230-242, Nashville, Tennessee, Apr. 1990.

[28] Tetzel, F., Hannes, Vo., Paradies, M., and Lehner, W., "An analysis of the feasibility of graph compression techniques for indexing regular path queries," In *Proc. ACM Int'l Workshop on Graph Data-management Experiences and Systems (GRADES)*, pp.11-16, Chicago, IL., 2017.

[29] Pacaci, A., Bonifati, A., and Ozsu, T. M. "Regular Path Query Evaluation on Streaming Graphs," In *Proc. ACM Int'l Conf. on Management of Data*, pp. 1415-1430, Portland, OR, June 2020.