

Simple, Fast, and Scalable Reachability Oracle

Ruoming Jin
Department of Computer Science
Kent State University
jin@cs.kent.edu

Guan Wang
Department of Computer Science
Kent State University
gwang@cs.kent.edu

ABSTRACT

A reachability oracle (or hop labeling) assigns each vertex v two sets of vertices: $L_{out}(v)$ and $L_{in}(v)$, such that u reaches v iff $L_{out}(u) \cap L_{in}(v) \neq \emptyset$. Despite their simplicity and elegance, reachability oracles have failed to achieve efficiency in more than ten years since their introduction: The main problem is high construction cost, which stems from a set-cover framework and the need to materialize transitive closure. In this paper, we present two simple and efficient labeling algorithms, *Hierarchical-Labeling* and *Distribution-Labeling*, which can work on massive real-world graphs: Their construction time is an order of magnitude faster than the set-cover based labeling approach, and transitive closure materialization is not needed. On large graphs, their index sizes and their query performance can now beat the state-of-the-art transitive closure compression and online search approaches.

1. INTRODUCTION

As one of the most fundamental graph operators, reachability has drawn much research interest in recent years [7, 37, 11, 35, 23, 12, 8, 22, 40, 38, 6, 36, 21, 9] and seems to continue fascinating researchers with new focuses [20, 39, 24] and new variants [16, 13, 32]. The basic reachability query answers whether a vertex u can reach another vertex v using a simple path ($u \rightarrow v$) in a directed graph. The majority of the existing reachability computation approaches belong to either transitive closure materialization (compression) [2, 27, 37, 23, 36] or online search [7, 35, 38]. The transitive closure compression approaches tend to be faster but generally have difficulty scaling to massive graphs due to the precomputation and/or memory cost. Online search is slower (often by one or two orders of magnitude) but can work on large graphs [38, 20]. The latest research [20] introduces a unified SCARAB method based on a “reachability backbone” (similar to the highway in the transportation network) to deal with their limitations: it can both help scale the transitive closure approaches and speed up online search. However, the query performance of transitive closure approaches tends to be slowed down and they may still not work if the size of the reachability backbone remains too large [20].

The reachability oracle, more commonly known as hop labeling,

[14, 34] is an interesting third category of approaches which lie between transitive closure materialization and online search. Each vertex v is labeled with two sets: $L_{out}(v)$, which contains hops (vertices) v can reach, and $L_{in}(v)$, which contains hops that can reach v . Given $L_{out}(u)$ and $L_{in}(v)$, but nothing else, we can compute if u reaches v by determining whether there is at least a common hop, $L_{out}(u) \cap L_{in}(v) \neq \emptyset$. The idea is simple, elegant, and seems very promising: Hop labeling can be considered as a factorization of the binary matrix of transitive closure; thus, it should be able to deliver more compact indices than the transitive closure and also offer fast query performance.

Unfortunately, after more than ten years since its first proposal [14] and a list of worthy attempts [30, 11, 12, 22, 6], hop labeling (or reachability oracle), still eludes us and fails to meet its expectations. Despite its appealing theoretical nature, recent studies [20, 36, 13, 38] all seem to confirm its inability to handle real-world large graphs: hop labeling is expensive to construct, taking much longer time than other approaches, and can barely work on large graphs, due to prohibitive memory cost of the construction algorithm. Many studies [20, 36, 13, 38] also demonstrate up to an order of magnitude slower query performance compared with the fastest transitive closure compression approaches (though we discover the underlying reason is mainly due to the implementation of hop labeling L_{out} and L_{in} ; employing a sorted vector/array instead of a set can significantly eliminate the query performance gap).

The high construction cost of the reachability oracle is inherent to the existing labeling algorithms and directly results in the scalability bottleneck. In order to minimize the labeling size, many algorithms [14, 30, 11, 22, 6] rely on a greedy set-cover procedure, which involves two costly operators: 1) repetitively finding densest subgraphs from a large number of bipartite graphs; and 2) materialization of the entire transitive closure. The latter is needed since each reachability pair needs to be explicitly covered by a selected hop. Even with concise transitive closure representation, such as using geometric format [11], or reducing the covered pairs using 3-hop [22, 6], the overall construction complexity is still close to or greater than $O(n^3)$, which is still too expensive for large graphs. Alternative labeling algorithms [34, 12] try to use graph separators, but only special graph classes consisting of small graph separators, such as planar graphs, can adopt such techniques well [34].

Can the reachability oracle be practical? Is it a purely theoretical concept which can only work on small toy graphs, or it is a powerful tool which can shape reality and work on real-world large graphs with millions of vertices and edges? Arguably, this is one of the most important unsolved puzzles in reachability computation. This work resolves these questions by presenting two simple and efficient labeling algorithms, *Hierarchical-Labeling* and *Distribution-Labeling*, which can work on massive real-world graphs.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Articles from this volume were invited to present their results at The 39th International Conference on Very Large Data Bases, August 26th - 30th 2013, Riva del Garda, Trento, Italy.
Proceedings of the VLDB Endowment, Vol. 6, No. 14
Copyright 2013 VLDB Endowment 2150-8097/13/14... \$ 10.00.

Their construction speeds are as fast as the state-of-the-art transitive closure compression approaches, there is no expensive transitive closure materialization, dense subgraph detection, or greedy set-cover procedure, there is no need for graph separators, and on large graphs, their index sizes and their query performance beat the state-of-the-art transitive closure compression and online search approaches [23, 36, 20, 36, 13, 38].

2. RELATED WORK

To compute the reachability, the directed graph is typically transformed into a DAG (directed acyclic graph) by coalescing strongly connected components into vertices, avoiding the trivial case where vertices reach each other in a strongly connected component. The size of the DAG is often much smaller than that of the original graph and is more convenient for reachability indexing. Let $G = (V, E)$ be the DAG for a reachability query, with number of vertices $n = |V|$ and number of edges $m = |E|$.

2.1 Transitive Closure and Online Search

There are two extremes in computing reachability. At one end, the entire transitive closure (TC) is precomputed and fully materialized (often in a binary matrix). Since the reachability between any pair is recorded, reachability can be answered in constant time, though $O(n^2)$ storage is prohibitive for large graphs. At the other end, DFS/BFS can be employed. Though it does not need an additional index, its query answering time is too slow for large graphs. As mentioned before, the majority of the reachability computation approaches aim to either compress the transitive closure [2, 19, 27, 37, 23, 9, 36, 21] or speed up the online search [7, 35, 38].

Transitive Closure Compression: This family of approaches aims to compress the transitive closure – each vertex u records a compact representation of $TC(u)$, i.e., all the vertices it reaches. The reachability from vertex u to v is computed by checking vertex v against $TC(u)$. Representative approaches include chain compression [19, 8], interval or tree compression [2, 27], dual-labeling [37], path-tree [23], and bit-vector compression [36]. Using interval-compress as an example, any contiguous vertex segment in the original $TC(u)$ is represented by an interval. For instance, if $TC(u)$ is $\{1, 2, 3, 4, 8, 9, 10\}$, it can be represented as two intervals: $[1, 4]$ and $[8, 10]$. Existing studies [36, 38, 20] have shown these approaches are the fastest in terms of query answering since checking against transitive closure $TC(u)$ is typically quite simple (linear scan or binary search suffices); in particular, the interval and path-tree approaches seem to be the best in terms of query answering performance. However, the transitive closure materialization, despite compression, is still costly. The index size is often the reason these approaches are not scalable on large graphs [38, 20].

Fast Online Search: Instead of materializing the transitive closure, this set of approaches [7, 35, 38] aims to speed up the online search. To achieve this, auxiliary labeling information per vertex is precomputed and utilized for pruning the search space. Using the state-of-the-art GRAIL [38] as an example, each vertex is assigned multiple interval labels where each interval is computed by a random depth-first traversal. The interval can help determine whether a vertex in the search space can be immediately pruned because it never reaches the destination vertex v .

The pre-computation of the auxiliary labeling information in these approaches is generally quite light; the index size is also small. Thus, these approaches can be applicable to very large graphs. However, the query performance is not appealing; even the state-of-the-art GRAIL can be easily one or two orders of magnitude slower than the transitive closure compression approaches [38, 20].

2.2 Reachability Oracle

The reachability oracle [14, 34], also referred to as hop labeling, was pioneered by Cohen *et al.* [14]. Though it also encodes transitive closure, it does not explicitly compress the transitive closure of each individual vertex independently (unlike the transitive closure compression approaches). Here, each vertex v is labeled with two sets: $L_{out}(v)$, which contains hops (vertices) v can reach; and $L_{in}(v)$, which contain hops that can reach v . Given only $L_{out}(u)$ and $L_{in}(v)$, but nothing else, we can compute if u reaches v by determining whether there is a common hop, $L_{out}(u) \cap L_{in}(v)$. In fact, a reachability oracle can be considered as a factorization of the binary matrix of transitive closure [22]; thus more compact indices are expected from such a scheme.

The seminal 2-hop labeling [14] aims to minimize the reachability oracle size, which is the total label size $\sum(|L_{out}(u)| + |L_{in}(u)|)$. It employs an approximate (greedy) algorithm based on set-covering which can produce a reachability oracle with size no larger than the optimal one by a logarithmic factor. The optimal 2-hop index size is conjectured to be $\tilde{O}(nm^{1/2})$. The major problem of the 2-hop indexing approach is its high construction costs: It needs to iteratively find dense subgraphs from a large number of bipartite graphs (representing the covering of transitive closure). Its computational cost is $O(n^3|TC|)$, where $|TC|$ is the total size of transitive closure. A number of approaches have sought to reduce construction costs through speeding up the set cover procedure [30], using concise transitive closure representation [11], or reducing the covered pairs using 3-hop [22, 6]. However, they still need to repetitively find the densest subgraphs and materialize the transitive closure.

2.3 Reachability Backbone and SCARAB

In the latest study [20], the authors introduce a general framework, referred to as SCARAB (SCALing ReachABILITY), for scaling the existing reachability indices (including both transitive closure compression and hop labeling approaches) and for speeding up the online search approaches. The central idea is to leverage a “reachability backbone”, which carries the major “reachability flow” information. The reachability backbone is similar in spirit to the highway structure used in several state-of-the-art shortest path distance computation methods on road networks [4, 28, 29]. However, the SCARAB work [20] is one of the first studies to construct and utilize such structure in the reachability computation. The reachability backbone definition and its discovery are quite different from the highway structures in [4, 28, 29].

Formally, the reachability backbone $G^* = (V^*, E^*)$ of graph G is defined as a subgraph of the transitive closure of G ($E^* \subseteq TC(G)$), such that for any reachable (u, v) pair, there must exist local neighbors $u^* \in V^*, v^* \in V^*$ with respect to locality threshold ϵ , i.e., $d(u, u^*) \leq \epsilon$ and $d(v^*, v) \leq \epsilon$, and $u^* \rightarrow v^*$. Here $d(u, u^*)$ is the shortest path distance from u to u^* where the weight of each edge is unit. Two algorithms are developed to approximate the minimal backbone, one based on set-cover and the other based on BFS. The latter, referred to as *FastCover*, is particularly efficient and effective, with time complexity $O(\sum_{v \in V} |N_\epsilon(v)| \log |N_\epsilon(v)| + |E_\epsilon(v)|)$, where $N_\epsilon(v)$ ($E_\epsilon(v)$) is the set of vertices (edges) v can reach in ϵ steps. Experiments show that even with $\epsilon = 2$, the size of the reachability backbone is significantly smaller than the original graph (about 1/10 the number of vertices of the original graph).

Though the scaling approach is quite effective for helping deal with large graphs, it is still constrained by the power of the original index approaches. For many large graphs, the reachability backbone can still be too large for them to process as shown in the experiment study in [20].

We also note that in [13], a new variant of reachability queries, k -hop reachability, is introduced and studied. It asks whether vertex u can reach v within k steps. This problem can be considered a generalization of the basic reachability, where $k = \infty$. A k -reach indexing approach is developed and the study shows that approach can handle basic reachability quite effectively (with comparable query performance to the fastest transitive closure compression approaches on small graphs). The k -reach indexing approach is based on vertex cover (a set of vertices covers all the edges in the graph), and it actually produces a reachability backbone with $\epsilon = 1$ as defined in [20]. But this study directly materializes the transitive closure between any pair of vertices in the vertex cover, where in [20], the existing reachability indices are used. Thus, for very large graphs where the vertex cover is often large, the pair-wise reachability materialization is not feasible.

2.4 Other Related Works

Distance 2-HOP Labeling: The 2-hop labeling method proposed by Cohen *et al.* [14] can also handle the exact distance labeling. Here, each vertex u records a list of intermediate vertices $OUT(u)$, which it can reach along with their (shortest) distances, and a list of intermediate vertices $IN(u)$, which can reach it along with their distances. To answer the point-to-point shortest distance query from u to v , we simply need to check all the common intermediate vertices between $OUT(u)$ and $IN(v)$ and choose the vertex p , such that $dist(u, p) + dist(p, v)$ is minimized for all $p \in OUT(u) \cap IN(v)$. However, its computational cost (similar to the reachability 2-hop labeling) is too expensive even for graphs with hundreds of thousands of vertices.

Recently, Abraham *et al.* [1] have developed a fast and practical algorithm to heuristically construct the distance labeling on large road networks. In particular, they utilize *contraction hierarchies* (CH) [18], which transform the original graph into a level-wise structure, and then assign the maximum-rank vertex on the shortest path between s and t as the hop for s and t . However, the core of CH needs to iteratively remove vertices and then add *shortcuts* for fast shortest path computation. Due to the power-law property, such operation easily becomes very expensive for general graphs. For example, to remove a vertex with thousands of neighbors may require checking millions of potential shortcuts.

Relationship to the latest reachability labeling [10] and distance labeling [3] papers:

We have recently become aware that Cheng *et al.* [10] have developed a reachability labeling approach, referred to as *TF-label*. Their approach is similar to the *Hierarchical Labeling* (HL) approach being introduced in this work. In particular, it can be considered a special case of HL where $\epsilon = 1$ (Section 4). The hierarchy being constructed in [10] is based on iteratively extracting a reachability backbone with $\epsilon = 1$, inspired by independent sets. A similar approach has been used in their earlier work on distance labeling, referred to as *IS-Label* [17]. In this paper, the hierarchy structure is extracted based on the reachability backbone approach [20], which has been shown to be effective and efficient for scaling reachability computation. In another recent work [3], Akiba *et al.* have proposed a distance labeling approach, referred to as the *Pruned Landmark*. This approach is similar in spirit to the *Distribution Labeling* (DL) approach. However, DL performs BFS in both directions (forward and reverse) to handle reachability labeling. The condition for assigning labels is also different. Finally, we would like to point out that both Hierarchical Labeling (HL) and Distribution Labeling (DL) are proposed independently of [10] and [3].

3. APPROACH OVERVIEW

In a reachability oracle of graph G , each vertex v is labeled with two sets: $L_{out}(v)$, which contains hops (vertices) v can reach; and $L_{in}(v)$, which contain hops that can reach v . A labeling is *complete* if and only if for any vertex pair where $u \rightarrow v$, $L_{out}(u) \cap L_{in}(v) \neq \emptyset$. The goal is to minimize the total label size, i.e., $\sum(|L_{out}(u)| + |L_{in}(u)|)$. A smaller reachability oracle not only helps to fit the index in main memory, but also speeds up the query processing (with $O(|L_{out}(u)| + |L_{in}(v)|)$ time complexity).

As we mentioned before, though the existing set-cover based approaches [14, 30, 11, 22, 6] can achieve approximate optimal labeling size within a logarithmic factor, their computational and memory costs are prohibitively expensive for large graphs. The labeling process not only needs to materialize the transitive closure, but it also uses an iterative set-cover procedure which repetitively invokes dense subgraph detection. The reason for such a complicated algorithm is that the following two criteria need to be met: 1) a labeling must be complete, and 2) we wish the labeling to be minimal. The existing approach [14, 22] essentially transforms the labeling problem into a set cover problem with the cost of constructing the ground set (which is the entire transitive closure) and dynamic generation and selection of good candidate sets (through dense subgraph detection).

To achieve efficient labeling which can work on massive graphs, the following issues require appropriate handling:

1. (Completeness without Transitive Closure): Can we guarantee labeling completeness without materialization of the transitive closure? Even compact [11] or reduced [22] materialization can be expensive for large graphs. Thus, the key is whether a labeling process can avoid the need to explicitly check whether a reachable pair (against some form of transitive closure) is covered by the existing labeling.

2. (Compactness without Optimization): Without the set-cover, it seems difficult to produce bounded approximate optimal labeling. But this does not mean that a compact reachability oracle cannot be produced. Clearly, each vertex should not record every valid hop in the labeling. In the set-cover framework, a price is computed to determine whether a vertex should be added to certain vertex labels. What other criteria can help determine the importance of hops so that each vertex can be more selective in what it records?

In this paper, we investigate how the hierarchical structure of a DAG can help produce a complete and compact reachability oracle. The basic idea is as follows: assuming a DAG can be represented in a hierarchical (multi-level) structure, such that the lower-level reachability needs to go through upper-level (but not vice versa), then we can somehow recursively broadcast the upper-level labels to lower-level labels. In other words, the labels of lower-level vertices (L_{in} and L_{out}) can directly utilize the already computed labels in the upper-level. Thus, on one side, by using the hierarchical structure, the completeness of labeling can be automatically guaranteed. On the other side, it provides an importance score (the level) of every hop, and each vertex only records those hops whose levels are higher than or equal to its own level. We note that there have been several studies [31, 28, 15, 26, 5, 1] using the hierarchical structure for shortest path distance computation on road networks; however, how to construct and utilize the hierarchical structure for reachability computation has not been fully addressed. To the best of our knowledge, this is the first study to construct a fast and scalable reachability oracle based on hierarchical DAG decomposition.

Now, to turn such an idea into a fast labeling algorithm for reachability oracle, the following two research questions need to be answered: 1) What hierarchical structure representation of a DAG can

be used? 2) How should L_{out} and L_{in} be computed efficiently using a given hierarchical structure? In this paper, we introduce two algorithms based on different hierarchical structures of a DAG:

Hierarchical-Labeling (Section 4): In this approach, the hierarchical structure is produced by a recursive reachability backbone approach, i.e., finding a reachability backbone G^* from the original graph G and then applying the backbone extraction algorithm on G^* . Recall that the reachability backbone is introduced by the latest SCARAB framework [20] which aims to scale the existing reachability computation approaches. Here we apply it recursively to provide a hierarchical DAG decomposition. Given this, a fast labeling algorithm is designed to quickly compute L_{in} and L_{out} one vertex by one vertex in a level-wise fashion.

Distribution-Labeling (Section 5): In this approach, the sophisticated reachability backbone hierarchy is replaced with the simplest hierarchy – a total order, i.e., each vertex is assigned a unique level in the hierarchy structure. Given this, instead of computing L_{in} and L_{out} one vertex at a time, the labeling algorithm will *distribute* the hop one by one (from higher order to lower order) to L_{in} and L_{out} of other vertices. The worst case computation complexity of this labeling algorithm is $O(n(n+m))$ (of the same order as transitive closure computation), though in practice it is much faster than the transitive closure computation.

4. HIERARCHICAL LABELING

Before we proceed to discuss the *Hierarchical Labeling* approach, let us formally introduce the one-side reachability backbone (first defined in [20] for scaling the existing reachability computation), which serves as the basis for hierarchical DAG decomposition and the labeling algorithm.

DEFINITION 1. (One-Side Reachability Backbone [20]) Given DAG G , and local threshold ϵ , the one-side reachability backbone $G^* = (V^*, E^*)$ is defined as follows: 1) $V^* \subseteq V$, such that for any vertex pair (u, v) in G with $d(u, v) = \epsilon$, there is a vertex v^* with $d(u, v^*) \leq \epsilon$ and $d(v^*, v) \leq \epsilon$; 2) E^* includes the edges which link vertex pair (u^*, v^*) in V^* with $d(u^*, v^*) \leq \epsilon + 1$.

Note that E^* can be simplified as a transitive reduction [20] (the minimal edge set preserving the reachability). Since computing transitive reduction is as expensive as transitive closure, rules like the following can be applied: $(u^*, v^*) \in E^*$ can be removed if there is another intermediate vertex $x \in V^*$ (not u^* and v^*) with $d(u^*, x) \leq \epsilon$ and $d(x, v^*) \leq \epsilon$. Also, for any two vertices u and v , if their distance is no higher than ϵ (local threshold), we refer to them as being a local pair (or being local to one another).

EXAMPLE 4.1. As a simple example, let V^* be a vertex cover of G , i.e., at least one end of an edge in E is in V^* ; and let E^* contain all edges $(u^*, v^*) \in V^* \times V^*$, such that $d(u^*, v^*) \leq 2$. Then, $G^* = (V^*, E^*)$ is one-side reachability backbone with $\epsilon = 1$. In Figure 1(b), G_1 is the reachability backbone of graph G_0 (Figure 1(a)) for $\epsilon = 2$.

The important property of the one-side reachability backbone is that for any non-local pair (u, v) : $u \rightarrow v$ and $d(u, v) > \epsilon$, there always exists $u^* \in V^*$ and $v^* \in V^*$, such that $d(u, u^*) \leq \epsilon$, $d(v^*, v) \leq \epsilon$, and $u^* \rightarrow v^*$. This property will serve as the key tool for recursively computing L_{out} and L_{in} . In [20], the authors develop the *FastCover* algorithm employing ϵ -step BFS for each vertex for discovering the one-side reachability backbone. They also show that when $\epsilon = 2$, the backbone can already be significantly reduced. To simplify our discussion, in this paper, we will focus on using the reachability backbone with $\epsilon = 2$ though the approach can be applied to other locality threshold values.

4.1 Hierarchical DAG Decomposition and Labeling Algorithm

Let us start with the hierarchical DAG decomposition which is based on the reachability backbone.

DEFINITION 2. (Hierarchical DAG Decomposition) Given DAG $G = (V, E)$, a vertex hierarchy is defined as $V_0 = V \supset V_1 \supset V_2 \supset \dots \supset V_h$, with corresponding edge sets $E_0, E_1, E_2 \dots E_h$, such that $G_i = (V_i, E_i)$ is the (one-side) reachability backbone of $G_{i-1} = (V_{i-1}, E_{i-1})$, where $0 < i \leq h$. The final graph $G_h = (V_h, E_h)$ is referred to as the core graph.

Intuitively, the vertex hierarchy shows the relative importance of vertices in terms of reachability computation. The lower level reachability computation can be resolved using the higher level vertices, but not the other way around. In other words, the reachability (backbone) property is preserved through the vertex hierarchy.

LEMMA 1. Assuming $u \in V_i, v \in V_i$, u reaches v in G ($u \xrightarrow{G} v$) iff u reaches v in G_i ($u \xrightarrow{G_i} v$). Furthermore, for any non-local vertex pairs $(u_i, v_i) \in V_i$, $d(u_i, v_i | G_i) > \epsilon$ (the distance in G_i), there always exists $u_{i+1} \in V_{i+1}$ and $v_{i+1} \in V_{i+1}$, such that $d(u_i, u_{i+1} | G_i) \leq \epsilon$, $d(v_{i+1}, v_i | G_i) \leq \epsilon$, and $u_{i+1} \xrightarrow{G_{i+1}} v_{i+1}$.

Proof Sketch: The first claim: assuming $u \in V_i, v \in V_i$, u reaches v in G ($u \xrightarrow{G} v$) iff u reaches v in G_i ($u \xrightarrow{G_i} v$), can be proved by induction. The base case where $i = 1$ is clearly true based on the reachability backbone definition (the reachability backbone will preserve the reachability between vertices in the backbone as they appear in the original graph). Assuming this is true for all $i < k$, then it also holds to be true for $i = k$. This is because for any $u \in V_i, v \in V_i$, we must have $u \in V_{i-1}$ and $v \in V_{i-1}$. Based on the reachability backbone definition, we have $u \xrightarrow{G_{i-1}} v$ iff $u \xrightarrow{G_{i-1}} v$. Then based on the induction, we have G ($u \xrightarrow{G} v$) iff u reaches v in G_i ($u \xrightarrow{G_i} v$). The second claim directly follows the reachability definition. \square

EXAMPLE 4.2. Figure 1 shows a vertex hierarchy for DAG G_0 (a), where $V_1 = \{5, 7, 9, \dots, 40\}$ (b) and $V_2 = \{7, 25, 35, 40\}$ (c). G_1 is the (one-side) reachability backbone of G_0 and G_2 is the corresponding (one-side) reachability backbone of G_1 .

To utilize the hierarchical decomposition for labeling, let us further introduce a few notations related to the vertex hierarchy. Each vertex v is assigned to a unique level: $level(v) = i$ iff $v \in V_i \setminus V_{i+1}$, where $0 \leq i \leq h$ and $V_{h+1} = \emptyset$. (Later, we will show that each vertex is labeled at its corresponding level using G_i and labels of vertices from higher levels). Assuming v is at level i , i.e., $level(v) = i$, let $N_{out}^k(v | G_i)$ ($N_{in}^k(v | G_i)$) be the v 's k -degree outgoing (incoming) neighborhood, which includes all the vertices v can reach (reaching v) within k steps in G_i . Finally, for any vertex v at level $i < h$, its corresponding outgoing (incoming) backbone vertex set $\mathcal{B}_{out}^\epsilon(v)$ ($\mathcal{B}_{in}^\epsilon(v)$) is defined as:

$$\mathcal{B}_{out}^\epsilon(v) = \{u \in V_{i+1} | d(v, u | G_i) \leq \epsilon \text{ and there is no other vertex } x \in V_{i+1}, \text{ where } d(v, x | G_i) \leq \epsilon \wedge d(x, u | G_i) \leq \epsilon (v \rightarrow x \rightarrow u)\} \quad (1)$$

$$\mathcal{B}_{in}^\epsilon(v) = \{u \in V_{i+1} | d(u, v | G_i) \leq \epsilon \text{ and there is no other vertex } y \in V_{i+1}, \text{ where } d(u, y | G_i) \leq \epsilon \wedge d(y, v | G_i) \leq \epsilon (u \rightarrow y \rightarrow v)\} \quad (2)$$

Now, let us see how the labeling algorithm works given the hierarchical decomposition. Contrary to the decomposition process which proceeds from the lower level to higher level (like peeling), the labeling performs from the higher level to the lower level.

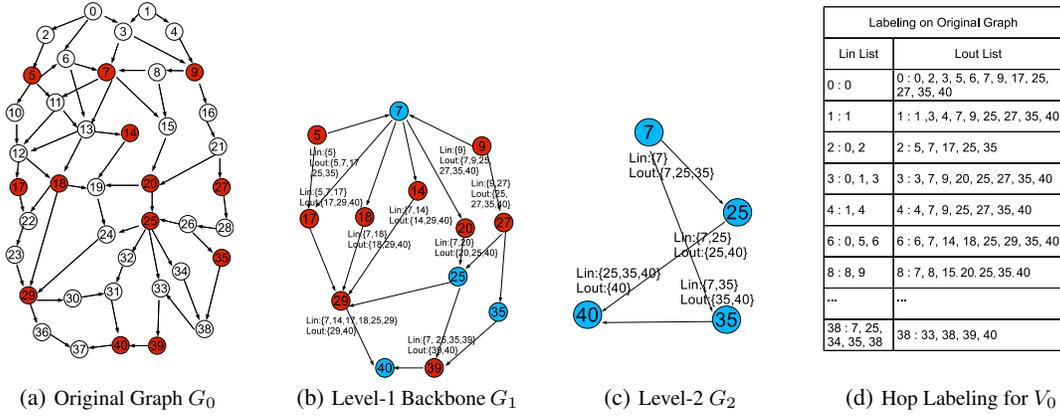


Figure 1: Running Examples of Hierarchical-Labeling

Specifically, it first labels the core graph G_h and then iteratively labels the vertex at level $h - 1$ to level 0.

Labeling Core Graph G_h : Theoretically, the diameter of the core graph G_h is no more than ϵ (the pairwise distance between any vertex pair in G_h is no more than ϵ), and thus no more reachability backbone is needed ($V_{h+1} = \emptyset$). In this case, for a vertex $v \in V_h$ ($level(v) = h$), the basic labeling can be as simple as follows:

$$L_{out}(v) = N_{out}^{\lceil \epsilon/2 \rceil}(v|G_h); \quad L_{in}(v) = N_{in}^{\lceil \epsilon/2 \rceil}(v|G_h) \quad (3)$$

The labeling is clearly complete for G_h as any reachable pair is within distance ϵ . Alternatively, since the core graph is typically rather small, we can also employ the existing 2-hop labeling algorithm [14, 30] to perform the labeling for core graphs. Given this, practically, the decomposition can be stopped when the vertex set V_h is small enough (typically less than $10K$) instead of making its diameter less than or equal to ϵ .

Labeling Vertices with Lower Level i ($0 \leq i < h$): After the core graph is labeled, the remaining vertices will be labeled in a level-wise fashion from higher level $h - 1$ to lower level (until level 0). For each vertex v at level $0 \leq i < h$, assuming all vertices in the higher level ($> i$) have been labeled (L_{out} and L_{in}), then the following simple rule can be utilized for labeling v :

$$L_{out}(v) = N_{out}^{\lceil \epsilon/2 \rceil}(v|G_i) \cup \left(\bigcup_{u \in \mathcal{B}_{out}^\epsilon(v|G_i)} L_{out}(u) \right) \quad (4)$$

$$L_{in}(v) = N_{in}^{\lceil \epsilon/2 \rceil}(v|G_i) \cup \left(\bigcup_{u \in \mathcal{B}_{in}^\epsilon(v|G_i)} L_{in}(u) \right) \quad (5)$$

Basically, the label of $L_{out}(v)$ ($L_{in}(v)$) at level i consists of two parts: the outgoing (incoming) $\lceil \epsilon/2 \rceil$ -degree neighbors of v in G_i , and the labels from its corresponding outgoing (incoming) backbone vertex set $\mathcal{B}_{out}^\epsilon(v|G_i)$ ($\mathcal{B}_{in}^\epsilon(v|G_i)$). In particular, if $\epsilon = 2$ (the typical locality threshold), then each vertex v basically records its direct outgoing (incoming) neighbors in G_i and the labels from its backbone vertex set.

Overall Algorithm: Algorithm 1 sketches the complete Hierarchical-Labeling approach. Basically, we first perform the recursive hierarchical DAG decomposition (Line 1). Then, the vertices at the core graph G_h will be labeled either by Formula 3 or using the existing 2-hop labeling approach (Line 2). Finally, the while-loop performs the labeling from higher level $h - 1$ to lower level 0 iteratively (Lines 4-10), where each vertex v in the level i (Lines 5 - 9) will be labelled based on Formulas 4 and 5.

EXAMPLE 4.3. Figure 1 illustrates the Hierarchical-Labeling process, where Figure 1(c) shows the labeling of core graphs. Note that for simplicity, each vertex by default records itself in both L_{in} and L_{out} , and $\epsilon = 2$. Figure 1(b) shows the labeling for vertices in V_1 ; and Table 1(c) illustrates the labeling of a few vertices in V_0 .

Algorithm 1 Hierarchical-Labeling($G = (V, E)$)

- 1: Perform Hierarchical Decomposition of G based on Definition 2;
- 2: Labeling core graph G_h ;
- 3: $i \leftarrow h - 1$;
- 4: **while** $i \geq 0$ {Labeling V_i from higher level to lower} **do**
- 5: **for each** $v \in V_i \setminus V_{i+1}$ {labeling each vertex specific for V_i } **do**
- 6: $L_{out}(v) \leftarrow N_{out}^{\lceil \epsilon/2 \rceil}(v|G_i) \cup (\bigcup_{u \in \mathcal{B}_{out}^\epsilon(v|G_i)} L_{out}(u))$
- 7: $L_{in}(v) \leftarrow N_{in}^{\lceil \epsilon/2 \rceil}(v|G_i) \cup (\bigcup_{u \in \mathcal{B}_{in}^\epsilon(v|G_i)} L_{in}(u))$
- 8: **end for**
- 9: $i \leftarrow i - 1$;
- 10: **end while**

4.2 Algorithm Correctness and Complexity

In the following, we first prove the correctness of the Hierarchical-Labeling algorithm, that is, that it produces a complete labeling: for any vertex pair (u, v) , $u \rightarrow v$ iff $L_{out}(u) \cap L_{in}(v) \neq \emptyset$. We then discuss its time complexity.

THEOREM 1. The Hierarchical-Labeling approach (Algorithm 1) produces a complete labeling for each vertex v in graph G , such that for any vertex pair (u, v) : $u \rightarrow v$ iff $L_{out}(u) \cap L_{in}(v) \neq \emptyset$.

Proof Sketch: We prove the correctness through induction: assuming Algorithm 1 produces the correct labeling for V_{i+1} , then it produces the correct labeling for V_i . Basically if for any vertex pair u^* and v^* in V_{i+1} , $u^* \rightarrow v^*$ iff $L_{out}(u^*) \cap L_{in}(v^*) \neq \emptyset$, then we would like to show that for any vertex pair u and v in V_i , this also holds. To prove this, we consider four different cases for any u and v in V_{i+1} : 1) $u \in V_i \setminus V_{i+1}$ and $v \in V_i \setminus V_{i+1}$; 2) $u \in V_i \setminus V_{i+1}$ and $v \in V_{i+1}$; 3) $u \in V_{i+1}$ and $v \in V_i \setminus V_{i+1}$; and 4) $u \in V_{i+1}$ and $v \in V_{i+1}$. Since case 4 trivially holds based on the reduction and cases 2 and 3 are symmetric, we will focus on proving cases 1 and 2.

Case 1 ($u \in V_i \setminus V_{i+1}$ and $v \in V_i \setminus V_{i+1}$): We observe: 1) $u \rightarrow v$ with $d(u, v) \leq \epsilon$ (local pair) iff there is $x \in V_i$, such that $d(u, x) \leq \lceil \frac{\epsilon}{2} \rceil$ and $d(x, v) \leq \lceil \frac{\epsilon}{2} \rceil$, i.e., $N_{out}^{\lceil \epsilon/2 \rceil}(u|G_i) \cap N_{in}^{\lceil \epsilon/2 \rceil}(v|G_i) \neq \emptyset$; and 2) $u \rightarrow v$ with $d(u, v) > \epsilon$ (non-local pair)

iff there are backbone vertices $u^*, v^* \in V_{i+1}$, such that $d(u, u^*) \leq \epsilon$, $d(v^*, v) \leq \epsilon$ and $u^* \rightarrow v^*$. That is, $L_{out}(u^*) \cap L_{in}(v^*) \neq \emptyset$ iff there are $x \in \mathcal{B}_{out}^\epsilon(u|G_i)$ and $y \in \mathcal{B}_{in}^\epsilon(v|G_i)$, such that $x \rightarrow y$, i.e., $L_{out}(x) \cap L_{in}(y) \neq \emptyset$ (if there is $x \in V_{i+1}$, such that $d(u, x) \leq \epsilon$ and $d(x, u^*) \leq \epsilon$, then we can always use x to replace u^* for the above claim; ($u^* \rightarrow v^*$ then $x \rightarrow v^*$))

iff $(\bigcup_{u \in \mathcal{B}_{out}^\epsilon(v|G_i)} L_{out}(u)) \cap (\bigcup_{u \in \mathcal{B}_{in}^\epsilon(v|G_i)} L_{out}(u)) \neq \emptyset$.

Case 2 ($u \in V_i \setminus V_{i+1}$ and $v \in V_{i+1}$): We observe 1) $u \rightarrow v$ with $d(u, v) \leq \epsilon$ (local pair) iff either $v \in \mathcal{B}_{out}^\epsilon(u|G_i)$ ($v \in L_{out}(u)$ and $v \in L_{in}(v)$), or there is $x \in \mathcal{B}_{out}^\epsilon(v|G_i)$, such that $x \rightarrow v$, i.e. $L_{out}(x) \cap L_{in}(v) \neq \emptyset$

iff $(\bigcup_{u \in \mathcal{B}_{out}^\epsilon(v|G_i)} L_{out}(u)) \cap (\bigcup_{u \in \mathcal{B}_{in}^\epsilon(v|G_i)} L_{out}(u)) \neq \emptyset$; and 2) $u \rightarrow v$ with $d(u, v) > \epsilon$ (non-local pair) iff there exists x such that $x \in \mathcal{B}_{out}^\epsilon(v|G_i)$ and $x \rightarrow v$, i.e. $L_{out}(x) \cap L_{in}(v) \neq \emptyset$ iff $(\bigcup_{u \in \mathcal{B}_{out}^\epsilon(v|G_i)} L_{out}(u)) \cap (\bigcup_{u \in \mathcal{B}_{in}^\epsilon(v|G_i)} L_{out}(u)) \neq \emptyset$.

Thus, in all cases, we have the correct labeling for any vertex pair u and v in V_{i+1} . Now, the core labeling is correct either based on the basic case where the graph diameter is no more than ϵ or based on the existing 2-hop labeling approaches [14, 30]. Together with the above induction rule, we have for any vertex pair in $V = V_0$, the label is complete and we thus prove the claim. \square

Complexity Analysis: The computational complexity of Algorithm 1 comes from three components: 1) the hierarchical DAG decomposition, 2) the core graph labeling, and 3) the remaining vertex labeling for levels from $h - 1$ to 0. For the first component, as we mentioned earlier, we can employ the *FastCover* algorithm [20] iteratively to extract the reachability backbone vertices V_i and their corresponding graph G_i . The *FastCover* algorithm is very efficient and to extract G_{i+1} from G_i , it just needs to traverse the ϵ neighbors of each vertex in G_{i+1} . Its complexity is $O(\sum_{v \in V} |N_{out}^\epsilon(v|G_i)| \log |N_{out}^\epsilon(v|G_i)| + |E_{out}^\epsilon(v|G_i)|)$, where $E_{out}^\epsilon(v|G_i)$ is the set of edges v can reach in ϵ steps. Also, we note that in practice, the vertex set V_i shrinks very quickly and after a few iterations (5 or 6 typically for $\epsilon = 2$), the number of backbone vertices is on the order of thousands (Section 6). We can also limit the total number of iterations, such as bounding h to be 10 and/or stop the decomposition when the V_i is smaller than some limit such as $10K$. For the second component, if the diameter is smaller than ϵ and Formula 3 is employed, it also has a linear cost: $O(\sum_{v \in V} (|N_{out}^\epsilon(v|G_h)| + |E_{out}^\epsilon(v|G_h)| + |N_{in}^\epsilon(v|G_h)| + |E_{in}^\epsilon(v|G_h)|))$. If we employ the existing 2-hop labeling approach [14, 30], the cost can be $O(|V_h|^4)$. However, since $|V_h|$ is rather small, the cost can be acceptable and in practice (Section 6), it is also quite efficient. Finally, the cost to assign labels for all the remaining vertices is linear to their neighborhood cardinality and the labeling size of each vertex. It can be written as $O(\sum_{v \in V_i \setminus V_{i+1}} (|N_{out}^\epsilon(v|G_h)| + |E_{out}^\epsilon(v|G_h)| + |N_{in}^\epsilon(v|G_h)| + |E_{in}^\epsilon(v|G_h)|) + ML)$, where M is the maximal number of vertices in the backbone vertex set and L is the maximal number of vertices in any L_{in} or L_{out} .

We note that for large graphs, the last component typically dominates the total computational cost as we need to perform list merge (set-union) operations to generate L_{out} and L_{in} for each vertex. However, compared with the existing hop labeling approach, Hierarchical-Labeling is significantly cheaper as there is no need for materializing transitive closure and the set-cover algorithm. The experimental study (Section 6) finds that the labeling size produced by the Hierarchical-Labeling approach is comparable to that produced by the expensive set-cover based optimization.

5. DISTRIBUTION LABELING

The Hierarchical-Labeling approach provides a fast alternative to produce a *complete* reachability oracle. Its labeling is dependent on a reachability-based hierarchical decomposition and follows a process similar to the classical transitive closure computation [33], where the transitive closure of all incoming neighbors are merged to produce the new transitive closure. However, the potential issue is that when merging L_{out} and L_{in} of higher level vertices for the lower level vertices, this approach does not (and cannot) check

whether any hop is redundant, i.e., their removal can still produce a complete labeling. Given the current framework, it is hard to evaluate the importance of each individual hop as they being cascaded into lower level vertices. Recall that for a vertex v , when computing its $L_{out}(v)$ and $L_{in}(v)$, its corresponding backbone vertex sets ($\mathcal{B}_{out}^\epsilon(v)$ and $\mathcal{B}_{in}^\epsilon(v)$) only eliminate those redundant backbones if they can be linked through a local vertex (Formulas 1 and 2). Thus even if $u \in \mathcal{B}_{out}^\epsilon(v)$, it may still be redundant as there is another vertex $u' \in \mathcal{B}_{out}^\epsilon(v)$ such that $u' \rightarrow u$ (but $d(u', u)$ is large). However, this issue is related to the difficulty of computing transitive reduction as mentioned earlier.

In light of these issues, we ponder the following: Can we perform labeling without the recursive hierarchical decomposition? Can we explicitly confirm the ‘‘power’’ or ‘‘importance’’ of an individual hop as it is being added into L_{out} and L_{in} ? In this work, we provide positive answers to these questions and along the way, we discover a simple, fast, and elegant labeling algorithm, referred to as *Distribution-Labeling*: 1) the recursive hierarchical decomposition is replaced with a simple total order of vertices (the order criterion can be as simple as a basic function of vertex degree); 2) each hop is explicitly verified to be added into L_{out} and L_{in} only when it can cover some additional reachable pairs, i.e., it is non-redundant. Surprisingly, the labeling size produced by this approach is even smaller than the set-cover approach on the benchmarking graphs used in the recent reachability studies (Section 6).

5.1 Hop Coverage and Labeling Basis

We first formally define the ‘‘covering power’’ of a hop and then study the relationship of two vertices in terms of their ‘‘covering power’’.

DEFINITION 3. (Hop Coverage) For vertex v , its **coverage** $Cov(v)$ is defined as $TC^{-1}(v) \times TC(v) = \{(u, w) : u \rightarrow v \text{ and } v \rightarrow w\}$. Note that $TC^{-1}(v)$ is the reverse transitive closure of v which includes all the vertices reaching v . If for any pair in $(u, w) \in Cov(v)$, $L_{out}(u) \cap L_{in}(w) \neq \emptyset$, then we say $Cov(v)$ is covered by the labeling. We also say $Cov(v)$ can be covered by v if each vertex u reaching v ($u \in TC^{-1}(v)$) has $v \in L_{out}(u)$ and each vertex w being reached by v has $v \in L_{in}(w)$ ($w \in TC(v)$).

Given this, the labeling L_{out} and L_{in} is complete if it covers $Cov(V) = \cup_{v \in V} Cov(v)$, i.e., for any $(u, w) \in Cov(V)$, To achieve a complete labeling, let us start with $Cov(v, v') = Cov(v) \cup Cov(v')$. We study how to use only v and v' to cover $Cov(v, v')$. Specifically, we consider the following question: *assuming v has been recorded by $L_{out}(u)$ for every $u \in TC^{-1}(v)$ and by $L_{in}(w)$ for every $w \in TC(v)$, then in order to cover the reachability pairs in $Cov(v, v')$ and only v' can serve as the hop, what vertices should record v' in their L_{out} and L_{in} ?*

To answer this question, we consider three cases: 1) v and v' are incomparable, i.e., $v \not\rightarrow v'$ and $v' \not\rightarrow v$; 2) $v' \rightarrow v$; and 3) $v \rightarrow v'$. For the first case, the labeling is straightforward: each $u \in TC^{-1}(v')$ needs to record $v' \in L_{out}(u)$ and each $w \in TC(v')$ needs to record $v' \in L_{in}(w)$. Note that in the worst case, this is needed in order to recover pairs as $TC^{-1}(v') \times \{v'\}$ and $\{v'\} \times TC(v')$. For Cases 2 and 3, Lemma 2 provides the answer.

LEMMA 2. Let $L_{out}(u) = \{v\}$ for every $u \in TC^{-1}(v)$ and $L_{in}(w) = \{v\}$ for every $w \in TC(v)$. If $v' \rightarrow v$, then with $L_{out}(u) = \{v, v'\}$ for $u \in TC^{-1}(v')$ and $L_{in}(w) = \{v'\}$ for $w \in TC(v') \setminus TC(v)$ (other labels remain the same), $Cov(\{v, v'\})$ is covered (using only hops v and v'). If $v \rightarrow v'$, then with $L_{out}(u) = \{v'\}$ for $u \in TC^{-1}(v') \setminus TC^{-1}(v)$ and $L_{in}(w) = \{v, v'\}$ for $w \in TC(v')$ (other labels remain the same), $Cov(\{v, v'\})$ is covered (using only hops v and v').

Vertex Order: The vertex order can be considered an extreme hierarchical decomposition, where each level contains only one vertex. Furthermore, the higher level the vertex, then the more important it is, the earlier it will be selected for covering, and the more vertices that are likely to record it in their L_{out} and L_{in} lists. There are many approaches for determining the vertex order. For instance, if following the set-cover framework, the vertex can be dynamically selected to be the *cheapest* in covering new pairs, i.e., $\frac{|TC^{-1}(v_i) \setminus TC^{-1}(X)| + |TC(v_i) \setminus TC(Y)|}{|Cov(V_s \cup \{v_i\}) \setminus Cov(V_s)|}$. However, this is computationally expensive. We may also use $|Cov(v_i)|$ which measures the covering power of vertex v , but this still needs to compute transitive closure. In this study, we found the following rank function, $(|N_{out}(v)| + 1) \times (|N_{in}(v)| + 1)$, which measures the vertex pairs with distance no more than 2 being covered by v , is a good candidate and can provides compact labeling. Indeed, we have used a similar criterion in [20] for selecting reachability backbone. In the experimental evaluation (Section 6), we will also use this rank function for computing the distribution labeling.

Labeling L_{out} and L_{in} : Given vertex v_i , we need to find (1) $u \in TC^{-1}(v_i) \setminus TC^{-1}(X)$, i.e., the vertices reaching v_i but not reaching by v such that $v \rightarrow v_i$ and it has a higher order (already being processed); and (2) $w \in TC(v_i) \setminus TC(Y)$, i.e., the vertices which can be reached by v_i but cannot be reached by v such that $v_i \rightarrow v$ and it has a higher order. The straightforward way for solving (1) is to perform a reversed traversal and visit (expand) the vertices based on the reversed topological order; then once the visited vertex has a higher order than v_i , all its descendents (including itself) will be colored (flagged) to be excluded from adding v_i to L_{out} ; thus v_i will be added to L_{out} for all uncolored vertices during the reverse traversal process. A similar ordered traversal process can be used for solving (2). However, the (reverse) ordered traversal needs a priority queue which results in $O(|V| \log |V| + |E|)$ complexity at each iteration. In this work, we utilize a more efficient approach that can effectively prune the traversal space and avoid the priority queue, which is illustrated in Algorithm 2.

Algorithm 2 Distribution-Labeling($G=(V,E)$)

```

1: Rank vertices in  $G$  in certain order;
2: for each  $v_i \in V$  {from higher order to lower} do
3:   Perform Reverse BFS starting from  $v_i$ , and for each vertex  $u$  being
   visited:
4:   if  $L_{out}(u) \cap L_{in}(v_i) \neq \emptyset$  then
5:     Do not add  $v_i$  to  $L_{out}(u)$  nor expand  $u$ ;
6:   else
7:     Add  $v_i$  into  $L_{out}(u)$  and expand  $u$  in the reverse BFS;
8:   end if
9:   Perform BFS starting from  $v_i$ , and for each vertex  $w$  being visited:
10:  if  $L_{in}(w) \cap L_{out}(v_i) \neq \emptyset$  then
11:    Do not add  $v_i$  to  $L_{in}(w)$  nor expand  $w$ ;
12:  else
13:    Add  $v_i$  into  $L_{in}(w)$  and expand  $w$  in the BFS;
14:  end if
15: end for

```

In Algorithm 2, the iteration labeling process is sketched in the foreach loop (Lines 2 to 15). The main procedure in computing $u \in TC^{-1}(v_i) \setminus TC^{-1}(X)$ for labeling L_{out} is outlined in Lines 3 – 8. The main idea is that when visiting a vertex u , once $L_{out}(u) \cap L_{in}(v_i)$ is no longer empty, we can simply exclude u and its descendents from consideration, i.e., $u \in TC^{-1}(X)$ (Lines 4 – 6). Intuitively, this is because there exists a vertex v , such that $u \rightarrow v \rightarrow v_i$ and has order higher than v_i . Similarly, the procedure that computes $w \in TC(v_i) \setminus TC(Y)$ for labeling L_{in} is outlined in Lines 9 – 14. Here, the condition $L_{in}(w) \cap L_{out}(v_i) \neq \emptyset$ is utilized to prune w and its descendents to determine L_{in} labeling.

Figure 2 illustrates the labeling process based on Algorithm 2 for the first three vertices 13, 7, and 25.

5.3 Completeness and Compactness

In the following, we discuss the labeling completeness (correctness), compactness (non-redundancy), and time complexity.

THEOREM 3. (Completeness) *The Distribution-Labeling algorithm (Algorithm 2) produces a complete L_{out} and L_{in} labeling, i.e., for any vertex pair (u, v) , $u \rightarrow v$ iff $L_{out}(u) \cap L_{in}(v) \neq \emptyset$.*

Proof Sketch: 1) $u \in TC^{-1}(v_i) \setminus TC^{-1}(X)$ and 2) $w \in TC(v_i) \setminus TC(Y)$. They are symmetric and we will focus on 1). Note for $u \in TC^{-1}(v_i) \setminus TC^{-1}(X)$, we need to exclude vertex u' such that $u' \rightarrow v \rightarrow v_i$, where v is already processed (has higher order than v_i). Assuming the labeling is complete for $Cov(V_s)$, where $V_s = \{v_1, \dots, v_{i-1}\}$, then $L_{out}(u') \cap L_{in}(v_i) \neq \emptyset$ (Line 4). If u' should be excluded, then its descendents from the BFS traversal will also be true and should also be excluded. Furthermore, the reverse BFS can visit all vertices where this condition does not hold, i.e., $L_{out}(u) \cap L_{in}(v_i) = \emptyset$, and thus $u \in TC^{-1}(v_i) \setminus TC^{-1}(X)$. \square

Theorem 3 shows that the Distribution-Labeling algorithm is correct; but how compact is the labeling? The following theorem shows an interesting *non-redundant* property of the produced labeling, i.e., no hop can be removed from L_{in} or L_{out} while preserving completeness. We note that this property has not been investigated before in the existing studies on reachability oracle and hop labeling [14, 30, 11, 12, 22, 6].

THEOREM 4. (Non-Redundancy) *The Distribution-Labeling algorithm (Algorithm 2) produces a non-redundant L_{out} and L_{in} labeling, i.e., if any hop h is removed from a L_{out} or L_{in} label set, then the labeling becomes incomplete.*

Proof Sketch: We will show that 1) for any $u \in TC^{-1}(v_i) \setminus TC^{-1}(X)$, v_i cannot be removed from L_{out} ; and 2) for any $w \in TC(v_i) \setminus TC(Y)$, v_i cannot be removed from L_{in} . Note that when v_i is being added to $L_{out}(u)$ and $L_{in}(w)$, it is non-redundant as the new labeling at least covers $(TC^{-1}(v_i) \setminus TC^{-1}(X)) \times \{v_i\}$ and $\{v_i\} \times TC(v_i) \setminus TC(Y)$.

However, will any later processed vertex v_j , such that $i < j$, make v_i redundant? The answer is no because in this case (still focusing on the above covered pairs by v_i), $u \rightarrow v_j \rightarrow v_i$ (or $w \leftarrow v_j \leftarrow v_i$), but the order of v_i is higher than v_j and v_j will not be added v_i into its L_{out} or L_{in} . In other words, for any vertex pair in $(TC^{-1}(v_i) \setminus TC^{-1}(X)) \times \{v_i\}$ or $\{v_i\} \times TC(v_i) \setminus TC(Y)$, v_i is the only hop linking these pairs, i.e., $L_{out}(u) \cap L_{in}(v_i) = \{v_i\}$ and $L_{out}(v_i) \cap L_{out}(u) = \{v_i\}$. Thus, v_i is non-redundant for all the vertices recording it as label, i.e., $L_{out}(u), u \in TC^{-1}(v_i) \setminus TC^{-1}(X)$ and $L_{in}(w), w \in TC(v_i) \setminus TC(Y)$. \square

As we discussed earlier, Hierarchical-Labeling does not have this property; we can see this through counter-examples. For instance, in Figure 1(b), 17 is redundant for $L_{out}(5)$. However, to remove these cases, the transitive reduction would have to be performed, which is expensive. Furthermore, whether the labels produced by the existing set-cover based approach [14] are redundant or not remains an open question.

Time Complexity: The worst case computational complexity of Algorithm 2 can be written as $O(|V|(|V| + |E|)L)$, where L is the maximal labeling size. However, the conditions in Line 4 and 10 can significantly prune the search space, and L is typically rather small, the *Distribution-Labeling* can perform labeling very efficiently. In the experimental study (Section 6), we will show Algorithm 2 is on average more than an order of magnitude faster than

Small Real Graph			Large Real Graph		
Dataset	V	E	Dataset	V	E
agrocyc	12684	13408	citeseer	693,947	312,282
amaz	3710	3600	citeseerx	6,540,399	15,011,259
antra	12499	13104	cit-Patents	3,774,768	16,518,947
arxiv	21608	116805	email	231,000	223,004
ecoo	12620	13350	go.uniprot	6,967,956	34,770,235
hpycyc	4771	5859	lj	971,232	1,024,140
human	38811	39576	mapped_100K	2,658,702	2,660,628
kegg	3617	3908	mapped_1M	9,387,448	9,440,404
mtbrv	9602	10245	uniprotenc_100m	16,087,295	16,087,293
nasa	5605	7735	uniprotenc_150m	25,037,600	25,037,598
p2p	48438	55349	uniprotenc_22m	1,595,444	1,595,442
reactome	901	846	web	371,764	517,805
vchocyc	9491	10143	wiki	2,281,879	2,311,570
xmark	6080	7028			

Table 1: Real datasets

the existing hop labeling and has comparable or faster labeling time than the state-of-the-art reachability indexing approaches on large graphs. Its labeling size is also small and surprisingly, even smaller than the greedy set-cover based labeling approaches in most of the cases. This may be an evidence that the labeling of the existing set-cover based approach [14] is redundant.

6. EXPERIMENTAL EVALUATION

In this section, we empirically evaluate the *Hierarchical-Labeling* and *Distribution-Labeling* labeling algorithms against the state-of-the-art reachability computation approaches on a range of real graphs which have been widely used for studying reachability [38, 36, 21, 20, 13]. Particularly, we are interested in the following questions (in terms of the query efficiency, construction cost, and index size): 1) How do the reachability oracle approaches perform compared with the transitive closure compression and online search approaches? 2) How do these two approaches perform compared with the existing 2-hop approaches assuming the latter one can complete the labeling? 3) How do these two methods (Hierarchical-Labeling and Distribution-Labeling) compare with one another?

6.1 Experimental Setup

To answer these questions, we evaluate the *Hierarchical-Labeling* (HL) and *Distribution-Labeling* (DL) labeling algorithms against the state-of-the-art reachability computation approaches:

1) *PathTree* (PT) [23], an improved version of Agrawal’s tree-interval method [2]; 2) *Nuutila’s Interval* (INT) [27], a transitive closure compression method, recently demonstrated to be one of the fastest reachability computation methods [36]; 3) *PAWH-8* (PW8) [36], the latest bit-vector compression method for transitive closure compression [36] and *PAWH-8* is its best variant [36]. 4) *K-Reach* (KR) [13], a latest vertex-cover based approach for general reachability computation, i.e., determine whether two vertices are within distance k . Here k is set to be the total number of vertices in the graph for the basic reachability. 5) *GRAIL* (GL) [38], a scalable reachability indexing approach using random DFS labeling (the number of intervals is set at 5, as suggested by authors). 6) *2HOP* (2HOP) [14], Cohen *et al.*’s 2-hop labeling approach;

Here, Path-Tree (1), Interval (2), and PAWH-8(3) are the state-of-the-art transitive closure compression approaches; K-Reach (4) is the latest general reachability approach and has been shown to be very capable in dealing with basic reachability [13] (it can also be considered as transitive closure compression as it materializes the transitive closure for the vertex-cover, a subset of vertices); GRAIL (5) is the state-of-the-art online search approach; and 2HOP (6) is the existing set-cover based hop labeling approach. In addition, we also include the latest SCARAB method [20] for scaling PathTree and speeding up GRAIL, referred to as *PATH-TREE** (PT*) and *GRAIL** (GL*), respectively. The locality parameter ϵ is set at 2 for

SCARAB. We also add the comparison with the latest reachability labeling *TF-label* (TF) [10], and the latest distance labeling method *Pruned Landmark* (PL) [3].

All the methods (including source code) except 2HOP and PL are either downloaded from authors’ websites or provided by the authors directly. We have implemented 2HOP, Hierarchical-Labeling (HL), Distribution-Labeling (DL), and PL; and 2HOP has been improved with several fast heuristics [30, 22] to speed up its construction time. All these algorithms are implemented in C++ based on the Standard Template Library (STL). We also downloaded and tested *IS-Label* (for distance computation) [17]. However, its query performance is on average more than 3 orders of magnitude slower than that of the reachability labeling and transitive closure compression methods; we omit reporting its results here.

In the experiments, we focus on reporting the three key measures for reachability computation: query time, construction time, and index size. For the query time, both *equal* and *random* reachability query workload are used. The equal query workload has about 50% positive (reachable pairs) and about 50% negative (unreachable pairs) queries. Positive queries are generated by sampling the transitive closure. Also the query time is the running time of a total of 100,000 reachability queries.

All experiments are performed on multi-core machines with Intel Xeon x5650 CPUs and 48GB RAM.

6.2 Experimental Results

In the following, we report the experimental results on small graphs first and then on large graphs. These graphs have been widely used for studying reachability computation [37, 11, 23, 22, 40, 38, 6, 36, 21, 13, 20]. In Table 1, the first three columns give the names, number of vertices and number of edges for the coalesced DAGs derived from each original graph. The last three columns give similar information for large real graphs.

Small Graphs: For small graphs, due to space limitation, we only report the query times of equal query load and construction times. The complete experimental results consisting of random query load, and index size can be found in the complete technical report [25]. (Their relative performance in terms of query times on the random query load and index size are comparable to those of large graphs.)

Table 2 reports the query times of the reachability oracle approaches (2HOP, Hierarchical-Labeling, and Distribution-Labeling) against the state-of-the-art transitive closure compression approaches (PAWH-8, INTERVAL, PATH-TREE, K-REACH), and online search (GRAIL), as well as some of their SCARAB counterparts, including *GRAIL** (GL*) and *PATH-TREE** (PT*) using the *equal* query load. We also compare them with the latest reachability labeling *TF-label* (TF) [10], and the latest distance labeling method *Pruned Landmark* (PL) [3].

We make the following important observations on the query time: 1) On small graphs, PATH-TREE outperforms other methods, though K-REACH is fairly close (as it is quite similar to the transitive closure materialization). Interestingly, the reachability oracle methods turn out to be quite comparable. In particular, the Distribution-Labeling (DL) is consistently about 2 times slower than PATH-TREE, and even faster than the other transitive closure compression approaches, INTERVAL and PAWH-8, on equal query load. 2) Compared to the existing set-cover based labeling approach 2HOP, Hierarchical-Labeling (HL) is quite comparable (slightly slower), but the query time of Distribution-Labeling (DL) is only 2/3 of that of 2HOP. 3) The reachability oracle approaches are slightly slower on the random query load than on the equal query load. This is because to determine vertex u cannot reach vertex v , the query processing has to completely scan $L_{out}(u)$ and $L_{in}(v)$. 4)

Dataset	GL	GL*	PT	PT*	KR	PW8	INT	2HOP	PL	TF	HL	DL
agrocyc	189.8	115.5	1.1	13.2	1.5	7.9	2.6	3.8	77.9	19.4	4.3	2.6
amaze	343.5	28.8	1.2	8.7	1.4	3.5	3.1	3.0	155.1	2.6	2.9	2.5
anthra	124.2	92.6	1.3	13.2	1.5	7.7	2.6	3.8	160.1	16.5	3.9	2.6
arxiv	282.5	214.2	2.8	15.2	—	11.4	3.7	7.0	194.2	17.3	11.8	3.4
ecoo	122.1	125.7	1.1	13.5	1.5	7.7	2.7	3.9	100.3	18.6	4.4	1.4
hpycyc	87.8	25.0	1.1	13.8	1.5	8.6	1.5	3.8	139.4	14.2	4.0	1.3
human	185.4	89.5	1.2	16.5	1.8	4.4	3.2	3.6	88.5	19.3	2.5	1.7
kegg	272.1	44.9	1.2	16.5	1.5	4.8	2.6	3.2	94.2	2.7	3.4	2.0
mtbrv	115.4	118.0	1.1	13.6	1.5	7.2	2.6	3.9	213.6	4.5	5.1	2.2
nasa	135.8	126.8	1.4	21.5	2.2	18.4	4.9	4.4	135.5	3.4	4.1	3.6
p2p	1117.8	308.4	1.2	6.5	—	3.3	1.7	2.0	118.1	3.2	3.6	2.1
reactome	111.8	22.8	1.1	11.6	1.8	12.1	3.0	3.1	125.5	2.0	2.9	2.1
vchocyc	107.6	97.0	1.0	13.6	1.5	7.9	2.6	3.7	102.8	16.9	3.8	2.5
xmark	134.7	255.9	1.4	21.6	1.9	35.8	4.9	5.9	152.8	5.8	6.5	4.4

Table 2: Query Time (ms) Based on Equal Query of Small Real Datasets

Dataset	GL	GL*	PT	PT*	KR	PW8	INT	2HOP	PL	TF	HL	DL
agrocyc	22.6	2.3	128.1	44.0	284.5	5.0	3.7	245.6	38.3	28.5	120.8	12.6
amaze	7.4	0.9	357.4	18.4	330.4	4.5	3.2	2672.2	14.4	3.9	43.2	4.1
anthra	14.1	2.3	88.2	41.4	246.3	4.1	2.9	241.0	40.8	44.9	89.4	12.4
arxiv	45.0	22.9	131873.0	8615.7	—	100.6	90.8	145332.0	814.5	759.5	618.6	38.2
ecoo	12.8	2.2	94.5	43.1	282.1	5.0	3.7	254.6	48.3	30.9	92.2	12.5
hpycyc	4.7	0.9	39.0	11.2	223.6	2.7	1.8	199.2	23.6	16.5	41.5	5.2
human	71.2	4.7	298.2	75.5	296.5	5.3	4.1	417.5	113.2	96.6	155.2	37.4
kegg	4.1	1.0	436.0	19.4	411.8	5.6	2.2	2878.0	16.3	4.3	48.3	2.4
mtbrv	9.4	1.7	71.7	30.9	249.3	2.2	3.0	208.3	41.5	14.6	115.3	9.8
nasa	10.2	2.2	49.4	21.1	1637.5	9.8	6.1	835.9	29.9	10.5	143.5	8.9
p2p	69.8	10.8	2942.6	432.6	—	14.8	25.6	21618.9	175.3	45.9	564.5	51.1
reactome	1.2	0.7	8.3	6.9	35.1	1.2	0.8	161.4	1.6	1.2	25.4	1.0
vchocyc	9.3	1.7	70.3	30.9	260.1	4.4	3.2	224.3	2.6	45.4	65.3	9.5
xmark	11.1	1.8	109.3	22.0	806.2	10.4	5.9	1557.0	41.6	19.5	53.2	8.7

Table 3: Construction Time (ms) of Small Real Datasets

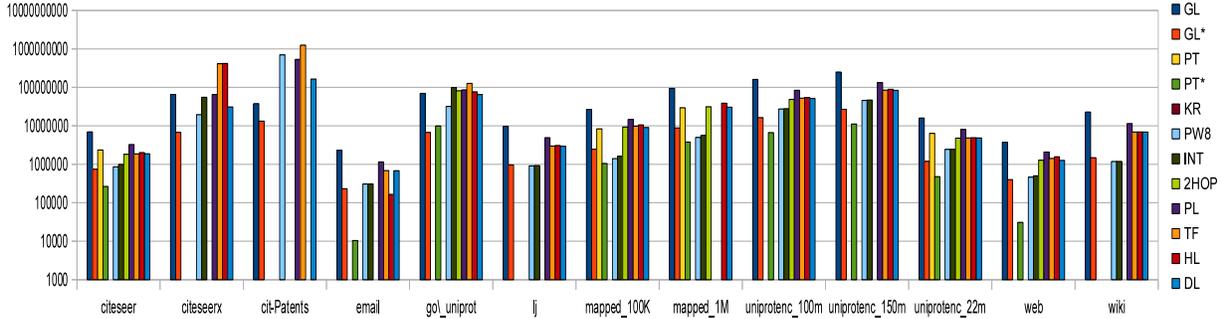


Figure 3: Index Size on Large Real Graphs (in terms of the number of integers used in the indices)

Dataset	GL	GL*	PT	PT*	KR	PW8	INT	2HOP	PL	TF	HL	DL
citeseer	63.4	42.4	4.9	26.9	—	20.6	12.3	4.5	294.5	6.1	7.7	5.3
citeseerx	2012.3	20230.9	—	—	—	76.3	8.8	—	235.8	37.2	210.2	7.7
cit-Patents	403.9	711.1	—	—	—	2538.9	—	—	148.2	56.7	—	53.2
email	575.9	30.2	—	10.2	—	13.8	5.7	—	99.5	3.0	2.2	3.0
go_uniprot	77.6	80.4	—	29.3	—	41.9	17.0	16.0	113.2	1752.5	6.2	12.7
lj	11972.2	2137.6	—	—	—	10.4	4.9	—	150.2	4.2	4.7	4.5
mapped_100K	253.8	92.3	6.7	25.1	—	90.6	6.0	5.1	7.1	15.0	5.1	7.2
mapped_1M	762.2	99.0	8.4	26.4	—	46.6	6.3	5.6	—	—	6.1	12.4
uniprotenc_100m	82.7	37.3	—	31.9	—	29.8	19.8	—	212.8	9.3	5.7	9.0
uniprotenc_150m	79.1	45.6	—	35.8	—	31.1	20.3	—	183.8	16.0	6.5	7.8
uniprotenc_22m	53.0	23.6	6.2	23.7	—	23.1	15.5	47598.4	104.6	6.2	4.4	5.9
web	2369.6	1041.9	—	11.6	—	13.1	5.5	3.0	87.1	4.2	3.5	3.9
wiki	100313.0	25655.5	—	—	—	8.6	6.0	—	396.5	4.5	3.6	4.9

Table 4: Query Time (ms) Based on Equal Query of Large Real Datasets

Due to additional distance comparison cost, the Pruned Landmark (PL) is fairly slow; its query performance is close to the GRAIL. 5) Both Hierarchical Labeling (HL) and Distribution Labeling (DL) are faster than the TF-labeling (TF) on the equal query load; though their performance become closer on the random workload [25].

Further analysis shows that the TF-labeling (TF) tends to scan much more labels than HL and DL to answer the positive queries. This is reflected in the equal query load as the random query load is mainly dominated by the negative queries. The number of labeling-scan operations for negative queries are determined by the labeling size,

Dataset	GL	GL*	PT	PT*	KR	PW8	INT	2HOP	PL	TF	HL	DL
citeseer	40.2	21.4	4.4	22.6	—	12.4	9.6	7.0	142.2	1.2	4.7	7.1
citeseerx	2585.6	719.1	—	—	—	39.8	13.4	—	109.6	55.4	23.7	11.9
cit-Patents	501.5	517.2	—	—	—	1766.3	—	—	101.5	66.5	—	48.1
email	754.4	8.3	—	11.3	—	14.0	10.1	—	72.7	5.8	3.5	5.0
go_uniprot	47.6	29.8	—	26.2	—	52.5	20.8	13.0	156.4	18.7	12.0	23.5
lj	829613.0	448.9	—	—	—	21.2	11.7	—	98.7	8.2	5.7	7.6
mapped_100K	52.4	23.6	5.9	20.8	—	4.9	5.0	6.5	123.0	7.2	6.7	9.4
mapped_1M	55.0	24.7	8.7	23.9	—	5.6	6.7	7.1	-	—	9.8	9.9
uniprotenc_100m	53.0	33.6	—	29.7	—	28.3	20.1	—	167.0	14.2	7.5	10.8
uniprotenc_150m	56.6	33.0	—	31.9	—	29.1	23.1	—	332.2	16.0	10.7	11.3
uniprotenc_22m	40.5	25.6	9.1	24.8	—	21.9	15.2	4.4	112.6	7.7	5.9	8.7
web	61295.3	386.8	—	18.5	—	22.3	9.2	4.4	70.4	8.3	4.7	6.3
wiki	76336.7	28.2	—	—	—	6.2	8.8	—	373.2	8.1	5.9	9.3

Table 5: Query Time (ms) Based on Random Query of Large Real Datasets

Dataset	GL	GL*	PT	PT*	KR	PW8	INT	2HOP	PL	TF	HL	DL
citeseer	2.0	0.2	18.0	1.1	—	0.5	0.3	14.1	1.7	0.7	2.2	0.5
citeseerx	17.6	3.7	—	—	—	17.0	7.0	—	26.6	207.0	182.1	9.9
cit-Patents	15.7	7.7	—	—	—	935.5	—	—	59.5	79.4	—	114.6
email	0.6	0.1	—	1.2	—	0.2	0.1	—	0.5	0.1	0.2	0.1
go_uniprot	32.4	3.2	—	5038.4	—	34.4	20.7	252.5	65.6	57.8	279.1	16.7
lj	2.6	0.2	—	—	—	0.7	0.8	—	2.0	0.6	1.2	0.3
mapped_100K	6.2	0.7	26.7	5.1	—	0.4	0.4	9.8	7.0	2.0	10.1	1.9
mapped_1M	28.3	2.8	103.3	23.1	—	2.4	3.8	52.2	-	—	45.5	6.9
uniprotenc_100m	66.3	5.1	—	8168.2	—	16.3	11.6	1029.0	33.9	49.4	67.3	13.9
uniprotenc_150m	101.6	8.7	—	18840.1	—	27.2	18.9	—	38.7	34.2	119.6	21.0
uniprotenc_22m	5.0	0.3	9801.7	41.2	—	1.4	1.1	102.7	3.2	2.0	5.2	1.0
web	1.0	0.1	—	46.2	—	0.6	1.8	32017.7	1.1	0.4	1.4	0.6
wiki	7.1	0.3	—	—	—	0.4	0.7	—	5.2	0.9	5.4	1.5

Table 6: Construction Time (Second) of Large Real Datasets

and DL typically has the smaller index size than TF while HL is close to TF (see later discussion on index sizes). This suggests the vertex order (hierarchy) defined in DL (HL) is more effective for reachability answering than TF (especially for positive queries).

Table 3 shows the construction time of different reachability indices on small graphs. We observe K-REACH and 2HOP are the slowest. This is understandable as K-REACH needs to perform vertex-cover discovery and materialize the transitive closure for the vertex-cover; and 2HOP needs to perform the expensive greedy set-cover and completely materialize the transitive closure. INTERVAL and PAWH-8 turn out to be the fastest and even faster than the online search GRAIL approach as the later still needs to perform random DFS a few times (in this study, we choose the number to be 5 as being used in [38]). Both Hierarchical-Labeling (HL) and Distribution-Labeling (DL) are much more efficient in labeling: The Hierarchical-Labeling is on average 5 times faster than 2HOP whereas the Distribution-Labeling is consistently 20 times faster (and in some case more than two order of magnitude faster) than 2HOP. In fact, it has even faster construction time than GRAIL and quite comparable to the INTERVAL and PAWH-8. The TF-labeling (TF) and the Pruned Landmark (PL) are typically faster than Hierarchical-Labeling (HL) but slower than Distribution-Labeling (DL). This is understandable as TF is simpler than Hierarchical-Labeling (HL) and PL needs additional computation cost with respect to Distribution-Labeling (DL).

Large Graphs: Large graphs provide the real challenge for the reachability computation. We observe that only three methods, GRAIL (GL), PAWH-8 (PW8), and Distribution-Labeling (DL) are able to handle all these graphs (GRAIL* is the SCARAB variant for speeding up query performance). Hierarchical-Labeling (HL), Pruned Landmark (PL), INTERVAL (INT), and TF-Labeling (TF) can work on 12, and PATH-TREE* (PT*) can work on 9, out of 13 large graphs. K-REACH (KR), PATH-TREE (PT) and 2-HOP (2HOP) fails on most of the large graphs. For K-REACH and PATH-TREE, their labeling size are too large to be materialized

in the main memory; for 2-HOP, its running time for these graphs often exceeds the 24-hour time limit.

Tables 4 and 5 report the query time using the *equal* and *random* query load, respectively. We make the following observations: 1) On large graphs, the transitive closure compression approaches, even on the graphs they can work, become significant slower. This is expected as the compressed transitive closure $TC(v)$ becomes larger, its search (linear or binary) becomes more expensive. Now, the advantage of the reachability oracle becomes clear as they become the fastest in terms of query time (even faster than PATH-TREE and INTERVAL, and consistently more than 5 times faster than PAWH-8). 2) Compared with the original 2HOP labeling, both Hierarchical-Labeling and Distribution-Labeling have comparable query performance on the graphs which they all can run. 3) The latest TF-labeling (TF) is slower than both Hierarchical-Labeling (HL) and Distribution-Labeling (DL) on most of the large graphs and for both equal and random query load.

Tables 6 shows the construction time on large graphs for all methods. We observe that PAWH-8 and INTERVAL are very fast though as the graph becomes larger, they become slower or cannot finish. Distribution-Labeling turns out to be quite comparable (fastest on several graphs). Hierarchical-Labeling can work on 8 out of 9 graphs and it shows significant improvement on 2 out of 5 graphs which 2HOP can also process. Distribution-Labeling is on average of one order of magnitude performance faster than 2HOP on these five graphs. Also, on average, the construction time of TF-labeling (TF) is comparable to that of Hierarchical-Labeling (HL), and is slower than that of Distribution-Labeling (DL).

Figure 3 shows the index size of different approaches. The results are quite consistent with the results on the small graphs on those graphs they can work. For most cases, PAWH-8 and INTERVAL have the smallest index size. 2HOP, Hierarchical-Labeling and Distribution-Labeling also perform well (better than GRAIL and K-Reach). The labeling sizes of 2HOP, Hierarchical-Labeling and Distribution-Labeling are quite comparable; Distribution-Labeling

has smaller labeling size than Hierarchical Labeling and very close to (or better than) 2HOP on the graphs it can run. Finally, on average, the index size of TF-labeling is quite close to that of Hierarchical Labeling (HL), but slightly higher than Distribution-Labeling (DL).

7. CONCLUSION

In this paper, by introducing two simple, elegant, and effective labeling approaches, *Hierarchical Labeling* and *Distribution Labeling*, we are able to resolve an important open question in reachability computation: the reachability oracle can be a powerful tool (or even the most useful one) to handle real, very large graphs. Our experimental results demonstrate that they can perform on graphs with millions of vertices/edges (scalable), are quickest in answering reachability queries on large graphs (fast), and have comparable or better labeling size as the set-cover based optimization approaches (compact). In the future, we will investigate the labeling on dynamic graphs and how to apply them on more general reachability computation, such as the k -reach problem.

Acknowledgment: The work was mainly done while the authors worked and/or visited at GraphSQL Inc. The work was also partially supported by the National Science Foundation under CAREER grant no. IIS-0953950 and by the Ohio Supercomputer Center under Grant no. PGS0218.

8. REFERENCES

- [1] I. Abraham, D. Delling, A. V. Goldberg, and R. F. Werneck. A hub-based labeling algorithm for shortest paths in road networks. In *SEA '11*, 2011.
- [2] R. Agrawal, A. Borgida, and H. V. Jagadish. Efficient mgmt. transitive relationships in large data and knowledge bases. In *SIGMOD*, 1989.
- [3] T. Akiba, Y. Iwata, and Y. Yoshida. Fast exact shortest-path distance queries on large networks by pruned landmark labeling. In *SIGMOD*, 2013.
- [4] H. Bast, S. Funke, P. Sanders, and D. Schultes. Fast Routing in Road Networks with Transit Nodes. *Science*, 316:566–, April 2007.
- [5] R. Bauer, D. Delling, P. Sanders, D. Schieferdecker, D. Schultes, and D. Wagner. Combining hierarchical and goal-directed speed-up techniques for dijkstra's algorithm. *J. Exp. Algorithmics*, 15, March 2010.
- [6] J. Cai and C. K. Poon. Path-hop: efficiently indexing large graphs for reachability queries. In *CIKM '10*, 2010.
- [7] L. Chen, A. Gupta, and M. E. Kurul. Stack-based algorithms for pattern matching on dags. In *VLDB '05*, pages 493–504, 2005.
- [8] Y. Chen and Y. Chen. An efficient algorithm for answering graph reachability queries. In *ICDE*, 2008.
- [9] Y. Chen and Y. Chen. Decomposing dags into spanning trees: A new way to compress transitive closures. In *ICDE '11*, 2011.
- [10] J. Cheng, S. Huang, H. Wu, and A. Wai-Chee Fu. Tf-label: a topological-folding labeling scheme for reachability querying in a large graph. In *SIGMOD*, 2013.
- [11] J. Cheng, J. X. Yu, X. Lin, H. Wang, and P. S. Yu. Fast computation of reachability labeling for large graphs. In *EDBT*, 2006.
- [12] J. Cheng, J. X. Yu, X. Lin, H. Wang, and P. S. Yu. Fast computing reachability labelings for large graphs with high compression rate. In *EDBT*, 2008.
- [13] James Cheng, Zechao Shang, Hong Cheng, Haixun Wang, and Jeffrey Xu Yu. K-reach: who is in your small world. *Proc. VLDB Endow.*, 5(11), July 2012.
- [14] Edith Cohen, Eran Halperin, Haim Kaplan, and Uri Zwick. Reachability and distance queries via 2-hop labels. *SIAM J. Comput.*, 32(5):1338–1355, 2003.
- [15] Daniel Delling, Martin Holzer, Kirill Mller, Frank Schulz, and Dorothea Wagner. High-performance multi-level graphs. In *9th DIMACS Implementation Challenge*, pages 52–65, 2006.
- [16] W. Fan, X. Wang, and Y. Wu. Performance guarantees for distributed reachability queries. *Proc. VLDB Endow.*, 5(11), July 2012.
- [17] A. Fu, H. Wu, J. Cheng, S. Chu, and R. Wong. Is-label: an independent-set based labeling scheme for point-to-point distance querying on large graphs. *PVLDB*, 2013.
- [18] R. Geisberger, P. Sanders, D. Schultes, and D. Delling. Contraction hierarchies: faster and simpler hierarchical routing in road networks. In *Proceedings of the 7th international conference on Experimental algorithms*, 2008.
- [19] H. V. Jagadish. A compression technique to materialize transitive closure. *ACM Trans. Database Syst.*, 15(4):558–598, 1990.
- [20] R. Jin, N. Ruan, S. Dey, and J. Y. Xu. Scarab: scaling reachability computation on large graphs. In *SIGMOD '12*, 2012.
- [21] R. Jin, N. Ruan, Y. Xiang, and H. Wang. Path-tree: An efficient reachability indexing scheme for large directed graphs. *TODS*, 36(1), 2011.
- [22] R. Jin, Y. Xiang, N. Ruan, and D. Fuhry. 3-hop: a high-compression indexing scheme for reachability query. In *SIGMOD '09*, 2009.
- [23] R. Jin, Y. Xiang, N. Ruan, and H. Wang. Efficiently answering reachability queries on very large directed graphs. In *SIGMOD '08*, 2008.
- [24] Ruoming Jin, Lin Liu, Bolin Ding, and Haixun Wang. Distance-constraint reachability computation in uncertain graphs. *Proc. VLDB Endow.*, 4(9), June 2011.
- [25] Ruoming Jin and Guan Wang. Simple, fast, and scalable reachability oracle. *CoRR*, abs/1305.0502, 2013.
- [26] H. Kriegel, P. Kröger, M. Renz, and T. Schmidt. Hierarchical graph embedding for efficient query processing in very large traffic networks. In *SSDBM '08*, 2008.
- [27] E. Nuutila. *Efficient Transitive Closure Computation in Large Digraphs*. PhD thesis, Finnish Academy of Technology, 1995.
- [28] P. Sanders and D. Schultes. Highway hierarchies hasten exact shortest path queries. In *17th Eur. Symp. Algorithms (ESA)*, 2005.
- [29] Peter Sanders and Dominik Schultes. Engineering highway hierarchies. *J. Exp. Algorithmics*, 17:1.6:1.1–1.6:1.40, September 2012.
- [30] R. Schenkel, A. Theobald, and G. Weikum. HOPI: An efficient connection index for complex XML document collections. In *EDBT*, 2004.
- [31] S. Shekhar, A. Fetterer, and B. Goyal. Materialization trade-offs in hierarchical shortest path algorithms. In *SSD '97*, 1997.
- [32] H. Shirani-Mehr, F. Banaei-Kashani, and C. Shahabi. Efficient reachability query evaluation in large spatiotemporal contact datasets. *Proc. VLDB Endow.*, 5(9), May 2012.
- [33] K. Simon. An improved algorithm for transitive closure on acyclic digraphs. *Theor. Comput. Sci.*, 58(1-3):325–346, 1988.
- [34] Mikkel Thorup. Compact oracles for reachability and approximate distances in planar digraphs. *J. ACM*, 51(6):993–1024, November 2004.
- [35] S. TriBl and U. Leser. Fast and practical indexing and querying of very large graphs. In *SIGMOD '07*, 2007.
- [36] S. J. van Schaik and O. de Moor. A memory efficient reachability data structure through bit vector compression. In *SIGMOD '11*, pages 913–924, 2011.
- [37] H. Wang, H. He, J. Yang, P. S. Yu, and J. X. Yu. Dual labeling: Answering graph reachability queries in constant time. In *ICDE '06*, page 75, 2006.
- [38] H. Yildirim, V. Chaoji, and M. J. Zaki. Grail: Scalable reachability index for large graphs. *PVLDB*, pages 276–284, 2010.
- [39] Zhiwei Zhang, Jeffrey Xu Yu, Lu Qin, Qing Zhu, and Xiaofang Zhou. I/o cost minimization: reachability queries processing over massive graphs. In *EDBT '12*, 2012.
- [40] L. Zhu, B. Choi, B. He, J. X. Yu, and W. K. Ng. A uniform framework for ad-hoc indexes to answer reachability queries on large graphs. In *DASFAA '09*, 2009.