



Distance labeling: on parallelism, compression, and ordering

Wentao Li¹ · Miao Qiao² · Lu Qin¹ · Ying Zhang¹ · Lijun Chang³ · Xuemin Lin⁴

Received: 3 December 2020 / Revised: 14 May 2021 / Accepted: 1 August 2021
© The Author(s), under exclusive licence to Springer-Verlag GmbH Germany, part of Springer Nature 2021

Abstract

Distance labeling approaches are widely adopted to speed up the online performance of shortest-distance queries. The construction of the distance labeling, however, can be exhaustive, especially on big graphs. For a major category of large graphs, small-world networks, the state-of-the-art approach is pruned landmark labeling (PLL). PLL prunes distance labels based on a node order and directly constructs the pruned labels by performing breadth-first searches in the node order. The pruning technique, as well as the index construction, has a strong sequential nature which hinders PLL from being parallelized. It becomes an urgent issue on massive small-world networks whose index can hardly be constructed by a single thread within a reasonable time. This paper first scales distance labeling on small-world networks by proposing a parallel shortest-distance labeling (PSL) scheme. PSL insightfully converts the PLL's node-order dependency to a shortest-distance dependence, which leads to a propagation-based parallel labeling in D rounds where D denotes the diameter of the graph. To further scale up PSL, it is critical to reduce the index size. This paper proposes effective index compression techniques based on graph properties as well as label properties; it also explores best practices in using betweenness-based node order to reduce the index size. The efficient betweenness estimation of the graph nodes proposed may be of independent interest to graph practitioners. Extensive experimental results verify our efficiency on billion-scale graphs, near-linear speedup in a multi-core environment, and up to 94% reduction in the index size.

Keywords Shortest distance · 2-Hop labeling · Betweenness · Parallelism · Compression · Ordering

1 Introduction

Given a graph G , a shortest-distance query $q(s, t)$ reports a minimized length of an $s-t$ path on G . It is a fundamental primitive as either a main function or a building block of applications such as geographic navigation, Internet routing, socially tenuous group finding [41], influential community searching [29] and event detection [40]. Many of these applications cannot afford frequent online distance computations, and therefore, 2-hop labeling [17] and its variations prevail as indexing techniques.

The index size of 2-hop labeling, however, can be quadratic to the number n of the nodes in the graph. For each node v , 2-hop labeling selects a set of nodes as hubs and tags v with its distances to its hubs as labels. A query $q(s, t)$ minimizes, over all hubs r shared by s and t , the 2-hop distances from s to t via r , i.e., $\text{dist}(s, r) + \text{dist}(r, t)$. To report a precise distance, the shared hubs of s and t must hit—have a common node with—a shortest path between s and t . Such a requirement over all pairs, s and t , of nodes is called the 2-hop cover constraint. A label set that satisfies

✉ Miao Qiao
miao.qiao@auckland.ac.nz

Wentao Li
wentao.li@uts.edu.au

Lu Qin
lu.qin@uts.edu.au

Ying Zhang
ying.zhang@uts.edu.au

Lijun Chang
lijun.chang@sydney.edu.au

Xuemin Lin
lxue@cse.unsw.edu.au

¹ AAIL, FEIT, University of Technology Sydney, Sydney, Australia

² University of Auckland, Auckland, New Zealand

³ The University of Sydney, Sydney, Australia

⁴ The University of New South Wales, Sydney, Australia

the 2-hop cover constraint can have a cardinality quadratic to n , especially on dense graphs. For example, a clique necessitates $\Omega(n^2)$ labels.

Finding a global minimum index size of 2-hop labeling, unfortunately, is NP-hard [17]. A local minimum, instead, can be reached by iteratively pruning redundant labels.¹ A label of a node v is redundant if the remaining labels in the label set still satisfy the 2-hop cover constraint. The pruning technique, however, has a strong sequential nature—pruning one label will affect the redundancy of another label. Consider two nodes u and v on the same shortest path between two nodes s and t . The moment when both s and t have the hub set of $\{u, v\}$, all labels on s and t are redundant. After pruning the label on s with the hub u , however, both labels on s and t with the hub v become critical. Due to such a dependency, the order of the pruning has a great influence on the pruning outcome and effectiveness.

The optimization of the pruning order is based on graph properties. For example, the planarity and hierarchical structure of road networks have been well explored to reach a scalable solution (see [33] as an entrance). For a major category of real graphs, small-world [43,45] networks, the state-of-the-art approach is pruned landmark labeling (PLL) [4].

The key to PLL's success on small-world networks is to encode the highly clustered topology into a node order and construct/prune labels strictly following the node order.

1. PLL prunes labels based on a node order that prioritizes the high-centrality² nodes. The label on a node s to its hub t is pruned if their distance can be answered by labels from s and t to a higher ranked hub. Therefore, a high-centrality hub r is able to prune labels along a large number (due to the clustered topology of the graph) of shortest paths hit by r .
2. PLL prunes a majority of labels in an implicit way. In other words, PLL constructs *pruned* labels directly as opposed to following a construct-and-then-prune paradigm. This is done by performing a *pruned* breadth-first-search (BFS) sourced from a hub r with the assignment of r *sequentially* following the node order.

It is worth noting that the index construction of PLL is highly node-order dependent: the pruning procedure in the BFS of hub r is dependent on the pruned labels constructed for the predecessor, in the node order, of r . Such a strong sequential nature of PLL hinders its parallelization.

¹ In many labeling approaches, the labels are pruned in an implicit way—a label will not be generated if pruning it is guaranteed to be safe.

² The centrality can be defined with degree, closeness, and betweenness [31].

On the other hand, the index time becomes an urgent issue for massive small-world networks whose index can hardly be constructed by a single thread within a reasonable time. For example, for the graph SINA³ with 58 million nodes and 261 million edges, PLL cannot finish the indexing within 3 days. The same situation applies to UK⁴ which has 77 million nodes and 2.9 billion edges.

This paper focuses on the scalability issue of the 2-hop distance labeling of small-world networks. We propose non-trivial algorithms to parallelize the indexing process of PLL and further reduce the index size. The scalability of our proposed approach is confirmed by extensive experiments. Our contributions are summarized as follows.

- We propose a parallel shortest-distance labeling approach PSL upon a novel and insightful conversion from the node-order label dependency of PLL to a shortest-distance label dependency. This conversion leads to a non-trivial propagation based labeling process. The algorithm completes in D rounds where D denotes the diameter of the graph—small for small-world networks. The resulting labels are identical to those constructed in the sequential algorithm of PLL.
- We provide two compression techniques to reduce the index size. The first one is based on graph properties and is thus applicable to all 2-hop labeling approaches; the second one explores the property of PSL, which leads to significant index reduction.
- We further explore best practices in using betweenness-based node order to reduce the index size. Given the quadratic time (infeasible for big graphs) in computing exact betweenness, we introduce k -betweenness—betweenness on paths with no more than k hops—to allow (i) an efficient sampling-based approximation and (ii) a holistic optimization of the node order for index reduction. The novel and efficient sampling-based approximate computation of node betweenness is the key to this reduction and may be of independent interest.
- We conduct extensive experiments to verify the performance of the proposed techniques. In a single-core environment, our index reduction technique dramatically shrinks the index size and improves the index time. In a multi-core environment, our PSL approach shows near-linear speed-up in parallelism. The proposed techniques jointly enable the index construction on networks with billion scale offline, which verifies the efficiency of the proposed approach.

³ <http://networkrepository.com/index.php>.

⁴ <http://law.di.unimi.it>.

The rest of the paper is organized as follows. Section 2 introduces the state-of-the-art 2-hop labeling approach on small-world networks. Section 3 devises a distance labeling algorithm. Section 4 introduces two index reduction techniques. Section 5 computes the betweenness-based node order by proposing novel approximation algorithms, which further reduced the index size. Section 6 summarizes related works. Section 7 experimentally evaluates our proposed approaches on real small-world networks, and Sect. 8 concludes the paper.

2 Preliminary

2.1 Shortest-distance problem

Let G be an unweighted graph with vertex set V_G and edge set E_G . Denote by n and m the number $|V_G|$ of nodes and $|E_G|$ of edges in the graph, respectively. For each node $v \in V_G$, denote by $N(v) = \{u | (u, v) \in E_G\}$ the **neighbors** of v and $\deg(v) = |N(v)|$ the degree of node v in G . We mainly use undirected graphs in the paper; Appendix B extends our techniques to directed graphs. Without loss of generality, we assume a connected graph G . Our techniques can be extended to disconnected graphs easily.

Let $p(s, t) = \{v_1, v_2, \dots, v_k\}$ with $s = v_1$ to $t = v_k$. p is a path on G if, for $\forall 1 \leq i \leq k$, $\text{edge}(v_i, v_{i+1}) \in E_G$. For an $i \in [1, k]$, denote by $p(s, t) = p(s, v_i) + p(v_i, t)$ the concatenation of two paths. The length of a path $p(s, t)$ is the number of edges on the path, i.e., $|p(s, t)| = k - 1$. The shortest path between s and t is the path with shortest length. The shortest length is the length of the shortest path, denoted as $\text{dist}_G(s, t)$. Given a graph G , a **point-to-point distance query** $q(s, t)$ with $s, t \in V$ returns the shortest distance $\text{dist}_G(s, t)$ between s and t . When the context is clear, we use $V, E, N(v), \deg(v), \text{dist}(s, t)$ to represent the node set, edge set, neighbor set of v , the degree of v and the shortest distance from s to t , respectively, for simplicity.

Example 1 Figure 1 shows a network $G = (V, E)$ with 12 nodes and 23 edges. The neighbors of v_6 are $N(v_6) = \{v_2, v_3, v_7\}$. Two paths between v_4 and v_6 are $p_1(v_4, v_6) = \{v_4, v_3, v_6\}$, $p_2(v_4, v_6) = \{v_4, v_1, v_2, v_6\}$. The shortest path $p_1(v_4, v_6)$ has the shortest length 2.

2.2 2-Hop labeling for distance queries

To efficiently process point-to-point distance queries, 2-hop labeling approach [17] precomputes the distances from each node to preselected hub nodes and uses the 2-hop distances via hubs to answer a query. Below we introduce the 2-hop labeling approach that has been slightly updated [4,20,26] to enable label reduction.

A **labeling function** L maps each node $v \in V$ to a **label set** $L(v)$. $L(v)$ consists of a set of **label entries** where each entry is a key/value pair $(u, \text{dist}(v, u))$ with a node $u \in V$ and the distance between v and u . The **hub nodes** of v are the projections of $L(v)$ on the key, i.e., $C(v) = \{u | (u, \text{dist}(v, u)) \in L(v)\}$. C is called the **hub function** of L . $\{L(v) | v \in V\}$ is a 2-hop labeling if L satisfies the **2-hop cover constraint** below.

Definition 1 (2-hop Cover Constraint [17]) A labeling function L satisfies the 2-hop cover constraint if for any node pair s and t , $C(s) \cap C(t)$ shares a node with a shortest path from s to t .

For a 2-hop labeling L , the **label size** $|L(v)|$ of a node v is the number of entries in $L(v)$. Denote by δ the largest label size of G , i.e., $\delta = \max_{v \in V} (|L(v)|)$.

Given a 2-hop labeling L , a distance query $q(s, t)$ is answered with $\text{Query}(s, t, L)$ defined below.

$$\text{Query}(s, t, L) = \min_{u \in C(s) \cap C(t)} \text{dist}(s, u) + \text{dist}(u, t)$$

Lemma 1 For a 2-hop labeling L that satisfies the 2-hop cover constraint, $\text{Query}(s, t, L) = \text{dist}(s, t)$.

Proof See Appendix A. \square

Assume that the label entries in each label set are stored in the ascending order of the key. The online cost of answering $q(s, t)$ is on retrieving and merging the entries in $L(s)$ and $L(t)$. Thus, the query time complexity is $O(|L(s)| + |L(t)|)$.

2.3 Pruned landmark labeling approach

Pruned landmark labeling approach (PLL) is the state-of-the-art 2-hop labeling approach on small-world networks.

Node Order. The effectiveness of PLL heavily relies on a *total order* r on V , called the node order. For two nodes u and v , if $r(u) > r(v)$, we say u has a higher rank than v . With the node order defined, we can safely rename the nodes such that

$$r(v_1) > r(v_2) > \dots > r(v_n).$$

A highly prevalent node order is degree-based: the order r prioritizes nodes with higher degrees and breaks ties based on original node ID. Specifically, for any two nodes v and v' , $r(v) > r(v')$ if

- $\deg(v) > \deg(v')$ or
- $\deg(v) = \deg(v')$ and $\text{ID}(v) > \text{ID}(v')$.

This paper uses degree-based node order by default unless another node order is specified in the context.

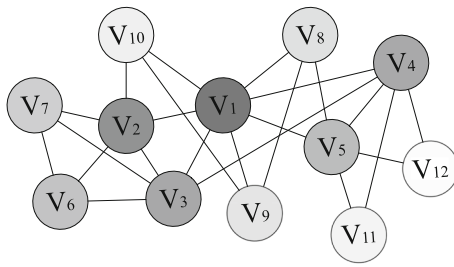


Fig. 1 Graph G

Example 2 We rank the nodes in Fig. 1 according to their degrees. When two nodes have the same degree, the tie is broken by the original ID of the node. We re-order the nodes such that $r(v_1) > r(v_2) > \dots > r(v_{12})$.

PLL with Pruned BFS. Algorithm 1 shows the process of PLL. Given a graph G and a node order v_1, v_2, \dots, v_n , PLL constructs a pruned 2-hop labeling L^{PLL} in n rounds (Line 1). In the i -th round, $i \in [1, n]$, PLL performs a pruned BFS search (a standard BFS search apart from Lines 6–8) sourced from v_i . To prune the BFS, PLL tests if the existing index can already report the distance between v_i and a node u (Line 6). If yes, u will neither be labeled nor expanded in this round (Line 7); otherwise, a label with hub v_i will be added to u (Line 8) and u will be expanded right away (Lines 9–12). Obviously, on the nodes that are either unexpanded or unreached, the labels with hub v_i are conceptually pruned.

Lemma 2 [4] *The index of PLL satisfies the 2-hop cover constraint.*

Algorithm 1: PLL

```

Input: Graph  $G(V, E)$ 
Output: The index  $L^{PLL}$ 
1 for  $i = 1, 2, \dots, n$  do
2    $Q \leftarrow$  a queue with only one element  $v_i$ ;
3    $\text{dist}(v_i) \leftarrow 0$  and  $\text{dist}(v) \leftarrow \infty, \forall v \in V \setminus v_i$ ;
4   while  $Q \neq \emptyset$  do
5      $u \leftarrow Q.\text{pop}()$ ;
6     if  $\text{Query}(v_i, u, L^{PLL}) \leq \text{dist}(u)$  then
7        $\text{continue}$ ;
8      $L^{PLL}(u) \leftarrow L^{PLL}(u) \cup \{(v_i, \text{dist}(u))\}$ ;
9     for  $w \in N(u)$  do
10      if  $\text{dist}(w) = \infty$  then
11         $\text{dist}(w) \leftarrow \text{dist}(u) + 1$ ;
12         $Q.\text{push}(w)$ ;
13 return  $L^{PLL}$ ;

```

The runtime of PLL for labeling large graphs can be very long. As shown in Line 6 of Algorithm 1, the query function to calculate the distance between v_i and u (i.e.,

$Q(v_i, u, L^{PLL})$) takes $O(\delta)$ time. The number of function calls is $\sum_{u \in V} \sum_{r \in C(u)} \text{deg}(u)$, which may reach δm in the worst case. This leads to a rather high time cost in terms of function calls for PLL, which is confirmed by our extensive empirical studies: PLL takes more than 3 days for labeling the graph SINA with 58 million nodes and 261 million edges.

Remarks. Note that PLL can work with any total order on V . Since there are $|V|!$ different total orders on V (the number of permutations of nodes in V), the selection of the node order in optimizing the space and/or temporal efficiency of PLL remains an open problem. It has been suggested by existing literature [31] that betweenness-centrality-based node order may be better than degree-based node order; however, improving PLL based on betweenness centrality faces the two challenges listed below.

- The computation of the exact betweenness centrality is as expensive as computing pairwise shortest distances, which is unaffordable on large graphs.
- The best practice of optimizing PLL based on approximate betweenness, that is, cost-effectively estimating the betweenness to reduce the index size of PLL is yet to be explored.

Part of this paper will dedicate to exploring betweenness centrality in forming a better distance index. Specifically, Sect. 2.4 will introduce the definition of betweenness centrality; Sect. 5 will propose a better algorithm for distance index construction based on a novel sampling-based betweenness centrality computation.

2.4 Node order: betweenness centrality

This paper actively explores the computation and application of betweenness-centrality-based node order in improving the efficiency of distance indexing. This section introduces betweenness centrality related concepts.

Given a graph $G(V, E)$ and two nodes $s, t \in V$, denote by $\sigma_{s,t}$ the number of different shortest paths between s and t (note that two different shortest paths between s and t can overlap on some nodes). Denoted by $\sigma_{s,t}(v)$, for $\forall v \in V$, the number of, among all the s - t shortest paths, shortest paths through v . If $s = t$, then $\sigma_{s,t} = 1$; if $v = s$ or $v = t$, then $\sigma_{s,t}(v) = 0$. We now define, for each node $v \in V$, its betweenness $\text{bc}(v)$.

Definition 2 (*Betweenness*) Given graph $G(V, E)$, for $\forall v \in V$, betweenness centrality

$$\text{bc}(v) = \sum_{s,t \in V} \frac{\sigma_{s,t}(v)}{\sigma_{s,t}}$$

Example 3 For node pair v_3, v_{10} in Fig. 1, there are two shortest paths between them: $p_1(v_3, v_{10}) = \{v_3, v_2, v_{10}\}$,

and $p_2(v_3, v_{10}) = \{v_3, v_1, v_{10}\}$. Then, $\sigma_{v_3, v_{10}} = 2$. Since only $p_1(v_3, v_{10})$ passes through v_2 , then $\sigma_{v_3, v_{10}}(v_2) = 1$. $\sigma_{v_3, v_{10}}(v_9) = 0$ because there is no v_3 - v_{10} shortest path through v_9 .

The betweenness centrality costs quadratic [13] time to compute, which is expensive for big graphs. For efficiently estimating betweenness centrality, we also resort to k -betweenness [14], a variation of betweenness centrality. Given a positive integer k , k -betweenness is defined for each node $v \in V$ by considering only shortest paths whose lengths are no more than k .

Definition 3 (*k-Betweenness* [14]) Given a graph $G(V, E)$ and an integer $k \geq 0$, the k -betweenness

$$\text{kbc}(v) = \sum_{s,t \in V, \text{dist}(s,t) \leq k} \frac{\sigma_{s,t}(v)}{\sigma_{s,t}}, \text{ for each } v \in V.$$

k -betweenness is a meaningful approximation of the betweenness centrality since k -betweenness is exactly the betweenness centrality when k approaches the diameter (the longest shortest path) of the graph. It was proposed since paths of long distances are less likely to form new edges, e.g., friendships in a social network [12] or a joint work in a collaboration network.

Example 4 If k is set to 2, then node pair v_{10}, v_{12} contributes nothing to k -betweenness of other nodes since their shortest distance $\text{dist}(v_{10}, v_{12}) = 3$.

$\text{kbc}(v)$ is an aggregation over all the shortest paths of lengths no larger than k . To simplify the discussions on the computation of $\text{kbc}(v)$, we introduce the concept of partial k -betweenness $\text{kbc}_s(v)$ which is the portion of $\text{kbc}(v)$ contributed by paths starting from node s .

Definition 4 (*Partial k-Betweenness*) Given graph $G(V, E)$, an integer $k \geq 0$ and a node $s \in V$,

$$\text{for } \forall v \in V, \text{kbc}_s(v) = \sum_{t \in V, \text{dist}(s,t) \leq k} \frac{\sigma_{s,t}(v)}{\sigma_{s,t}}.$$

Note that k -betweenness can be easily derived from partial k -betweenness

$$\text{kbc}(v) = \sum_{s \in V} \text{kbc}_s(v).$$

Therefore, the computation of $\text{kbc}(v)$ boils down to computing $\text{kbc}_s(v)$ for each node s of G . According to the definition of $\text{kbc}_s(v)$, it can be observed that $\text{kbc}_s(v)$ becomes zero if v is not included in any shortest path sourced at s with $\leq k$ length; thus, we have Lemma 3.

Lemma 3 If $\text{dist}(s, v) \geq k$ or $\text{dist}(s, v) = 0$ then $\text{kbc}_s(v) = 0$.

Proof If $\text{dist}(s, v) = 0$, $\text{kbc}_s(v) = 0$ since $\sigma_{v,t}(v) = 0$ for any $t \in V$. If $\text{dist}(s, v) = k$, v can only be the end point of any s - t shortest path $p(s, t)$ with $|p(s, t)| \leq k$. Then, $\sigma_{s,v}(v) = 0$. If $\text{dist}(s, v) > k$, there is no s - t path via v with $|p(s, t)| \leq k$, and $\sigma_{s,t}(v) = 0$. Thus, $\text{kbc}_s(v) = 0$. \square

An exact algorithm to compute k -betweenness is presented in [14]. The basic idea is to perform a graph traversal sourced from each $s \in V$ to compute $\text{kbc}_s(v)$, for $\forall v \in V$. Compared to betweenness computation, calculating k -betweenness only needs to visit nodes within a distance k to each source, thus improving the efficiency. However, for small-world graphs, the number of nodes within distance k (k greater than 2) to each source may still be large [42], which makes the exact k -betweenness computation for large graphs undesirable.

3 Parallelized distance labeling

Sections 3.1 and 3.2 revisit PLL to identify the label properties and order dependency. Section 3.3 transforms the order dependency in PLL to distance dependency. By utilizing the distance dependency, Sect. 3.4 proposes a practical approach in constructing the index in parallel.

3.1 Label property

The labels of PLL show an important node-order property.

Theorem 1 For any two nodes $\forall u, v \in V$, v is a hub of u under PLL, i.e., $(v, \text{dist}(v, u)) \in L^{\text{PLL}}(u)$, if and only if v is the highest ranked node on all the shortest paths from u to v .

Proof Let S be the set of nodes on all the shortest paths from u to w . Let w be the highest ranked node in S .

We prove that all nodes in S have w as their hubs in L^{PLL} by contradiction. Assume that there is a node z in S such that z does not have a hub of w in L^{PLL} . Consider the round of Algorithm 1 where the pruned BFS sourced w is performing. Let L' be the snapshot of the PLL label set right before the round begins. Given that z has no hub of w , then either

- z is explicitly pruned: $\text{Query}(z, w, L') = \text{dist}(w, z)$, or
- z is implicitly pruned: z is not reached since there is at least a node z' on the shortest path from w to z explicitly pruned with $\text{Query}(z', w, L') = \text{dist}(w, z')$.

In either case, it requires a common hub between w and z (or z') to produce the query result, which is impossible since (i) $z, z' \in S$ and (ii) w has the highest rank in S and (iii) L' does not include any hub ranked higher than w . Contradiction.

Since all nodes in S have w as their hubs in L^{PLL} , we prove the two directions of the theorem in two cases: (1) if $w = v$,

that is, v is the highest ranked node in S , then v is a hub of $u \in S$ and (2) if $r(w) > r(v)$, when before the pruned BFS sourced from v is performed, w is already a common hub of u and v . As w is on the shortest path between u and v , the label with hub v on u is pruned and not in PLL. \square

Lemmas 4–6 are derived from Theorem 1.

Lemma 4 *If v is a hub of u , $r(v) > r(u)$.*

Proof Since v has the highest rank on a shortest path from v to u (Theorem 1), $r(v) > r(u)$. \square

Lemma 5 *For $\forall u \in V$, $(u, 0) \in L^{\text{PLL}}(u)$.*

Proof We make the path as $p(u, u)$ and according to Theorem 1, $(u, 0)$ will be always inserted to $L^{\text{PLL}}(u)$. \square

Lemma 6 *For $\forall (u, v) \in E$, $(u, 1) \in L^{\text{PLL}}(v)$, if $r(u) > r(v)$; otherwise, $(v, 1) \in L^{\text{PLL}}(u)$.*

Proof Let $p(u, v)$ be the path with an edge. According to Theorem 1, the higher ranked node is the hub node. \square

3.2 Order dependency

To see the dependency among the labels, we partition the labels in L^{PLL} according to their hub nodes. Let v_1, v_2, \dots, v_n be the node order under which label set L^{PLL} was constructed.

We define two sets with particular meanings. Recall that PLL has n rounds where the i -th round performs a pruned BFS sourced from v_i . We denote by $L_{<i}^{\text{PLL}}(u)$ the snapshot of $L^{\text{PLL}}(u)$ at the beginning of the i -th round and by $L_i^{\text{PLL}}(u)$ the incremental label of u built in the i -th round.

Definition 5 (*Order Specific Label Set*)

$$L_i^{\text{PLL}}(u) = \{(v_i, \text{dist}(v_i, u)) \in L^{\text{PLL}}\},$$

for $\forall i \in [1, n], u \in V$. Let $L_i^{\text{PLL}} = \bigcup_{u \in V} L_i^{\text{PLL}}(u)$.

Definition 6 (*Order Partial Label Set*)

$$L_{<i}^{\text{PLL}}(u) = \{(v_j, \text{dist}(v_j, u)) \in L^{\text{PLL}} \mid j < i\},$$

for $\forall i \in [1, n + 1], u \in V$. Let $L_{<i}^{\text{PLL}} = \bigcup_{u \in V} L_{<i}^{\text{PLL}}(u)$. $L_{<n+1}^{\text{PLL}} = L^{\text{PLL}}$.

The following lemma shows that the pruning condition in Algorithm 1 leads to an order dependency among labels.

Lemma 7 (*Order Dependency*) $L_i^{\text{PLL}}(u)$ depends on $L_{<i}^{\text{PLL}}(u)$. Specifically, $L_i^{\text{PLL}}(u) =$

$$\begin{cases} \{(v_i, \text{dist}(v_i, u))\} & \text{if } \text{Query}(v_i, u, L_{<i}^{\text{PLL}}) > \text{dist}(v_i, u); \\ \emptyset & \text{otherwise.} \end{cases}$$

Proof Let S be the set of nodes on the shortest path from v_i to u (including v_i and u). Let w be the node with the highest rank in S . If $v_i = w$, according to Theorem 1, (i) v_i is a hub of u and (ii) for $\forall v \in S \setminus v_i$, v is not a hub of v_i , and thus $\text{Query}(v_i, u, L_{<i}^{\text{PLL}}) > \text{dist}(v_i, u)$. If $r(v_i) < r(w)$, then v_i is not a hub of u and label $(w, \text{dist}(w, v_i)), (w, \text{dist}(w, u)) \in L_{<i}^{\text{PLL}}$ and thus $\text{Query}(v_i, u, L_{<i}^{\text{PLL}}) = \text{dist}(v_i, u)$. \square

Lemma 7 shows that $L_i^{\text{PLL}}(u)$ depends on $L_{<i}^{\text{PLL}}(u)$ while $L_{<i}^{\text{PLL}}(u)$ depends on $L_{i-1}^{\text{PLL}}(u)$. Such a convoluted dependency can hardly be removed as long as the labels are built in the node order.

Example 5 Table 1 illustrates how PLL constructs the index. A cell at the row of v_i and the column of v_j records the order specific label of v_i at the j -th round. In column v_1 , pruned BFS inserts v_1 into $L_1^{\text{PLL}}(u), \forall u \in V$. In column v_2 , PLL performs pruned BFS and v_2 becomes the hub of $\{v_2, v_3, v_6, v_7, v_{10}\}$ due to the pruning condition of $L_1^{\text{PLL}} = \{L_1^{\text{PLL}}(u) \mid u \in V\}$. The order dependency in the column v_7 : partial set $L_{<7}^{\text{PLL}} = \bigcup_{i < 7, u \in V} L_i^{\text{PLL}}(u)$ prunes the labels on all nodes except on v_7 .

3.3 Distance dependency

To break the order dependency in the label construction, consider the pruning condition of Line 6, Algorithm 1. When $\text{Query}(v_i, u, L_{<i}^{\text{PLL}}) = \text{dist}(u, v_i)$ prunes a node label on u , there must be two labels on u and v_i , respectively, to a common hub w such that $\text{dist}(u, w) + \text{dist}(w, v_i) = \text{dist}(u, v_i)$. Therefore, $\text{dist}(u, w)$ and $\text{dist}(w, v_i)$ must be no greater than $\text{dist}(u, v_i)$. In other words, all the labels with distances greater than $\text{dist}(u, v_i)$ have no effect on the query result of $\text{Query}(v_i, u, L_{<i}^{\text{PLL}})$ and the corresponding pruning outcomes.

From the above intuition, we group the label entries in L^{PLL} based on their label distances. The rearranged label sets will pave the way to our parallel shortest-distance labeling (PSL) approach (Sect. 3.4) and are thus called PSL label sets. Let D be the diameter of the graph G .

Definition 7 (*Distance Specific Label Set*)

$$L_d^{\text{PSL}}(u) = \{(v, \text{dist}(v, u)) \in L^{\text{PLL}}(u) \mid \text{dist}(v, u) = d\},$$

for $\forall u \in V, d \in [1, D]$. Let $L_d^{\text{PSL}} = \{L_d^{\text{PSL}}(u) \mid u \in V\}$.

Similarly, the partial label of a node then becomes the set of label entries with distance less than a certain distance and is defined in Definition 8.

Definition 8 (*Distance Partial Label Set*)

$$L_{<d}^{\text{PSL}}(u) = \{(v, \text{dist}(v, u)) \in L^{\text{PLL}}(u) \mid \text{dist}(v, u) < d\},$$

Table 1 The index of PLL and PSL

ID	PSL														
	PLL						PSL								
	v_1	v_2	v_3	v_4	v_5	v_6	v_7	v_8	v_9	v_{10}	v_{11}	v_{12}	0	1	2
v_1	$(v_1, 0)$	-	-	-	-	-	-	-	-	-	-	-	$(v_1, 0)$	-	-
v_2	$(v_1, 1)$	$(v_2, 0)$	-	-	-	-	-	-	-	-	-	-	$(v_2, 0)$	$(v_1, 1)$	-
v_3	$(v_1, 1)$	$(v_2, 1)$	$(v_3, 0)$	-	-	-	-	-	-	-	-	-	$(v_3, 0)$	$(v_1, 1)$	$(v_2, 1)$
v_4	$(v_1, 1)$	-	$(v_3, 1)$	$(v_4, 0)$	-	-	-	-	-	-	-	-	$(v_4, 0)$	$(v_1, 1)$	$(v_3, 1)$
v_5	$(v_1, 1)$	-	-	$(v_4, 1)$	$(v_5, 0)$	-	-	-	-	-	-	-	$(v_5, 0)$	$(v_1, 1)$	$(v_4, 1)$
v_6	$(v_1, 2)$	$(v_2, 1)$	$(v_3, 1)$	-	-	$(v_6, 0)$	-	-	-	-	-	-	$(v_6, 0)$	$(v_2, 1)$	$(v_1, 2)$
v_7	$(v_1, 2)$	$(v_2, 1)$	$(v_3, 1)$	-	-	$(v_6, 1)$	$(v_7, 0)$	-	-	-	-	-	$(v_7, 0)$	$(v_2, 1)$	$(v_6, 1)$
v_8	$(v_1, 1)$	-	-	-	$(v_5, 1)$	-	-	$(v_8, 0)$	-	-	-	-	$(v_8, 0)$	$(v_1, 1)$	$(v_5, 1)$
v_9	$(v_1, 1)$	-	-	-	-	-	-	$(v_8, 1)$	$(v_9, 0)$	-	-	-	$(v_9, 0)$	$(v_1, 1)$	$(v_8, 1)$
v_{10}	$(v_1, 1)$	$(v_2, 1)$	-	-	-	-	-	-	$(v_9, 1)$	$(v_{10}, 0)$	-	-	$(v_{10}, 0)$	$(v_1, 1)$	$(v_2, 1)$
v_{11}	$(v_1, 2)$	-	$(v_3, 2)$	$(v_4, 1)$	$(v_5, 1)$	-	-	-	-	-	$(v_{11}, 0)$	-	$(v_{11}, 0)$	$(v_4, 1)$	$(v_5, 1)$
v_{12}	$(v_1, 2)$	-	$(v_3, 2)$	$(v_4, 1)$	$(v_5, 1)$	-	-	-	-	-	-	$(v_{12}, 0)$	$(v_{12}, 0)$	$(v_4, 1)$	$(v_5, 1)$

for $\forall u \in V, d \in [1, D + 1]$. Let $L^{\text{PSL}} = \bigcup_{u \in V} L^{\text{PSL}}(u)$. In particular, $L^{\text{PSL}}(u) = L^{\text{PSL}}_{<D+1}(u)$.

The equivalence of the index L^{PLL} and the novel index L^{PSL} is given in Theorem 2.

Theorem 2 $L^{\text{PSL}} = L^{\text{PLL}}$.

Proof Since all the label $(v, \text{dist}(v, u))$ in L^{PLL} has $\text{dist}(v, u) \leq D$, L^{PSL} includes all labels in L^{PLL} and has no other labels according to the definition. \square

Example 6 Table 1 shows a rearrangement of labels in PLL. A cell with row v_i and column j denotes label set of $L^{\text{PSL}}_j(v_i)$ —the PLL labels of v_i whose distances are j .

Distance Dependency. Definitions 7 and 8 provide us an opportunity in removing the order dependency in the label construction process.

Theorem 3 (Distance Dependency) $L^{\text{PSL}}_d(u)$ depends on $L^{\text{PSL}}_{<d}$. Specifically, given a node u , for a node $v \in V$ with $r(v) > r(u)$ and $\text{dist}(u, v) = d$, $(v, \text{dist}(v, u)) \in L^{\text{PSL}}_d(u)$ if and only if $\text{Query}(u, v, L^{\text{PSL}}_{<d}) > d$.

Proof Consider a node v with $\text{dist}(u, v) = d$. Denote by S the set of nodes on the shortest paths from u to v and let w be the highest ranked node in S . According to Theorem 1, we have two exclusive cases:

- $w = v$ if and only if v is the hub of u ;
- $w \neq v$ means that

- w is the hub of both u and v , and
- $\text{dist}(u, w), \text{dist}(w, v) < d$,

and therefore, $\text{Query}(u, v, L^{\text{PSL}}_{<d}) = d$.

Therefore, if $(v, \text{dist}(v, u)) \notin L^{\text{PSL}}_d(u)$, namely, v is not a hub of u , then $w \neq v$, and then $\text{Query}(u, v, L^{\text{PSL}}_{<d}) = d$. Besides, if $(v, \text{dist}(v, u)) \in L^{\text{PSL}}_d(u)$, namely, v is a hub of u , v is the highest ranked node in S and therefore, no other node in S can be a hub of v , that is, $\text{Query}(u, v, L^{\text{PSL}}_{<d}) > d$. \square

By transforming the order dependency to distance dependency, it is possible to complete the index construction in D rounds where D denotes the diameter of the graph.

Example 7 In Table 1, each row corresponds to the partial label of a node, while each column corresponds to the incremental labels regarding each distance value. When $d = 0$, each node add to itself since the distance between itself is zero. When $d = 1$, we either add nodes that are 1-hop away to a node u or prune the 1-hop away nodes. Note that according to Lemma 4, only higher ranking nodes can be hubs of lower ranking nodes. When $d = 2$, we either add nodes that are 2-hop away to a node u or prune the 2-hop away nodes.

For instance, if $u = v_{11}$, the node v_1 that is 2-hop away is added into $L_2^{\text{PSL}}(v_{11})$. But node v_8 is pruned since we can make use of v_5 , which is less than two hops away with v_8 , to prune it.

3.4 The parallelized labeling method

To apply Theorem 3 to generate $L_d^{\text{PSL}}(u)$, all the node pairs with distance equal to d are to be examined which is also expensive. This section provides a practical algorithm, parallel shortest-distance labeling (PSL), to construct the index L^{PSL} in label propagations.

Propagation-Based Label Construction. This section provides a positive answer to the following question: can we build the distance specific label $L_d^{\text{PSL}}(u)$ by gathering the labels of its neighbors, namely, $L_{d-1}^{\text{PSL}}(v)$, for $v \in N(u)$? We formally show that $\bigcup_{v \in N(u)} L_{d-1}^{\text{PSL}}(v)$ is sufficient to create $L_d^{\text{PSL}}(u)$ in Lemma 8.

Lemma 8 All the hub nodes of labels in $L_d^{\text{PSL}}(u)$ appear in labels $\bigcup_{v \in N(u)} L_{d-1}^{\text{PSL}}(v)$ as hub nodes.

Proof We show that if a node is not a hub of any node $v \in N(u)$ in $L_{d-1}^{\text{PSL}}(v)$, then it is not a hub of u in $L_d^{\text{PSL}}(u)$. Let $w \neq u$ be a hub of u in $L_d^{\text{PSL}}(u)$ but is not a hub of any node $v \in N(u)$ in $L_{d-1}^{\text{PSL}}(v)$. Note that the PLL was built in a BFS search. Consider the round when the pruned BFS search is sourced from w . Since $w \neq u$ and w is a hub of u , there is a shortest path from w to u such that w is a hub of all nodes on the path. Let s be the predecessor of u on the shortest path. $s \in N(v)$ and $(w, \text{dist}(w, s)) \in L^{\text{PLL}}$. Since $\text{dist}(w, s) = d - 1$, w is a hub of $L_{d-1}^{\text{PSL}}(s)$, contradiction. \square

Pruning Conditions. From Lemma 8, we can construct $L^{\text{PSL}}(u)$ in an iterative way and the initial condition is given in Lemma 5 by inserting u to the label $L_0^{\text{PSL}}(u)$ as its own hub. However, pouring all nodes in $\bigcup_{v \in N(u)} L_{d-1}^{\text{PSL}}(v)$ directly into $L_d^{\text{PSL}}(u)$ produces a large set of candidate labels. Therefore, we propose two rules to prune unnecessary label entries.

Lemma 9 A hub w in the label set $\bigcup_{v \in N(u)} L_{d-1}^{\text{PSL}}(v)$ is not a hub of u if $r(w) < r(u)$.

Proof Lemma 4. \square

Lemma 10 A hub w in the label set $\bigcup_{v \in N(u)} L_{d-1}^{\text{PSL}}(v)$ is not a hub of u in $L_d^{\text{PSL}}(u)$ if $\text{Query}(w, u, L_{<d}^{\text{PSL}}) \leq d$.

Proof If $\text{Query}(w, u, L_{<d}^{\text{PSL}}) < d$, then $\text{dist}(w, u) < d$, w is not a hub of u with distance $\text{dist}(w, u) = d$. If $\text{Query}(w, u, L_{<d}^{\text{PSL}}) = d$, we discuss in two cases.

- $\text{dist}(w, u) < d$, w is not a hub of u with distance d .
- $\text{dist}(w, u) = d$, there is a node z on the shortest path between w and u with $r(z) > r(w)$. According to Theorem 1, w is not be a hub of u in L^{PLL} .

Therefore, w is not a hub of u if $\text{Query}(w, u, L_{<d}^{\text{PSL}}) \leq d$. \square

Based on the above pruning rules, we propose our label propagation function to find the exact $L_d^{\text{PSL}}(u)$, $\forall u \in V$.

Denote by $C_d(v)$ the set of hub nodes in label set $L_d^{\text{PSL}}(v)$, for $\forall v \in V$ and $d \in [1, D + 1]$.

Theorem 4 (Label Propagation Function)

$$L_d^{\text{PSL}}(u) = \bigcup_{w \in C_{d-1}(v), \text{ for } \forall v \in N(u)} L_d^{\text{PSL}}(u, w) \tag{1}$$

where $L_d^{\text{PSL}}(u, w) =$

$$\begin{cases} \emptyset, & \text{if } r(w) < r(u) \text{ or } \text{Query}(w, u, L_{<d}^{\text{PSL}}) \leq d; \\ \{(w, \text{dist}(w, u))\}, & \text{otherwise.} \end{cases} \tag{2}$$

Proof Denote by L' the label set computed from Equation (1). We show that $L' = L_d^{\text{PSL}}(u)$ in two directions. Due to the correctness of Lemma 8 and the pruning conditions, the label set $L_d^{\text{PSL}}(u) \subseteq L'$. The follow parts prove $L' \subseteq L_d^{\text{PSL}}(u)$. Let $(w, \text{dist}(w, u))$ be a label in L' . Equation (2) shows that $r(w) > r(u)$ and $\text{Query}(w, u, L_{<d}^{\text{PSL}}) > d$.

If in L^{PLL} , w is not a hub of u , then according to Theorem 1, there is a node s that in S —the set of all nodes in the shortest path between w and u —with $r(s) > r(w) > r(u)$. Therefore, $\text{dist}(w, s), \text{dist}(s, u) < d$ and $\text{dist}(w, u) \leq d$, and thus, $\text{Query}(w, u, L_{<d}^{\text{PSL}}) \leq d$, contradiction.

Therefore, w is a hub of u in L^{PLL} . Besides, if $\text{dist}(w, u) < d$, $\text{Query}(w, u, L_{<d}^{\text{PSL}}) = \text{dist}(w, u) < d$, contradiction. Thus, $\text{dist}(w, u) = d$. Now we have proved that w is a hub of u in L^{PLL} with $\text{dist}(w, u) = d$, i.e., w is a hub of u in $L_d^{\text{PSL}}(u)$ which completes the proof. \square

The PSL Algorithm. Algorithm 2 puts all parts of PSL together. $L_0^{\text{PSL}}(u)$ is obtained by add u to itself (Line 1). Then, for each edge, the higher ranked node v is added into lower ranked node u to form $L_1^{\text{PSL}}(u)$ according to Lemma 6 (Lines 2–4). If L_{d-1}^{PSL} is empty—the path with length $d - 1$ is covered by $L_{<d-1}^{\text{PSL}}$ —we find the final index (Line 6). Otherwise, nodes are parallelly processed to build L_d^{PSL} for $d > 1$ (Lines 7–12): each node u first finds its superset $\text{cand}(u)$ (Lemma 8) (Line 8) and then, pruning conditions 9–10 apply (Lines 10–11). Entry $(w, \text{dist}(w, u))$ is then added to $L_d^{\text{PSL}}(u)$ (Lines 11–12).

Example 8 In Fig. 2a, each node $u \in V$ is added to $L_0^{\text{PSL}}(u)$ for $d = 0$. In Fig. 2b, for each edge (u, v) , v is added to $L_1^{\text{PSL}}(u)$ if $r(v) > r(u)$. For instance, $L_1^{\text{PSL}}(v_3) = \{(v_1, 1), (v_2, 1)\}$, $L_1^{\text{PSL}}(v_2) = \{(v_1, 1)\}$, $L_1^{\text{PSL}}(v_7) = \{(v_2, 1), (v_3, 1), (v_6, 1)\}$. In Fig. 2c, for each node u , hubs in $\{L_1^{\text{PSL}}(w) | w \in N(u)\}$ are candidate hubs and then added to $L_2^{\text{PSL}}(u)$ if the pass pruning conditions. v_6 has three neighbors v_2, v_3, v_7 . Then, candidate nodes are $\{v_1, v_2, v_3, v_6, v_7\}$.

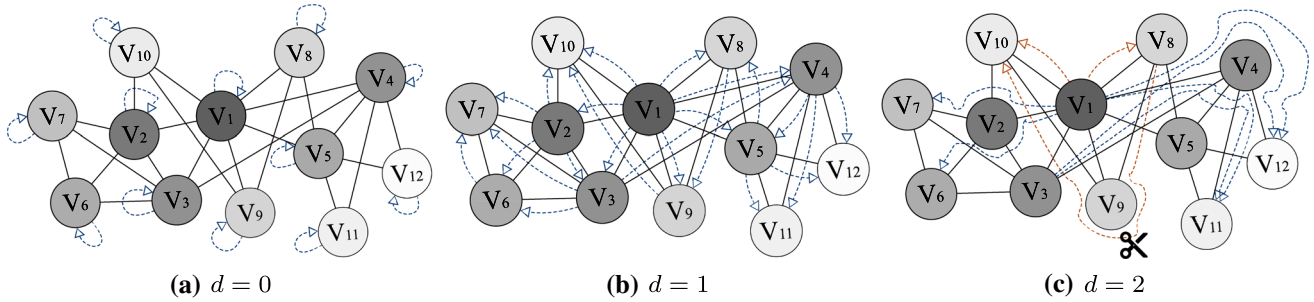


Fig. 2 The execution of PSL from $d = 0$, $d = 1$ to $d = 2$

Algorithm 2: PSL

Input: Graph $G(V, E)$
Output: The index L^{PSL}

```

1 Insert  $(u, 0)$  into  $L_0^{PSL}(u), \forall u \in V;$ 
2 for  $(u, v) \in E$  do
3   if  $r(u) > r(v)$  then Insert  $(u, 1)$  into  $L_1^{PSL}(v);$ 
4   else Insert  $(v, 1)$  into  $L_1^{PSL}(u);$ 
5  $d \leftarrow 2;$ 
6 while  $L_{d-1}^{PSL}$  is not empty do
7   for  $u \in V$  in parallel do
8     cand( $u$ )  $\leftarrow$  hubs in  $L_{d-1}^{PSL}(v), \forall v \in N(u);$ 
9     for each node  $w \in$  cand( $u$ ) do
10      // Pruning Condition Lemma 9
11      if  $r(w) < r(u)$  then continue;
12      // Pruning Condition Lemma 10
13      if Query( $w, u, L_{<d}^{PSL}$ )  $\leq d$  then continue;
14      Insert  $(w, d)$  into  $L_d^{PSL}(u);$ 
15    $d \leftarrow d + 1;$ 
16 return  $L^{PSL};$ 

```

$(v_1, 2)$ will be put into $L_2^{PSL}(v_6)$ since the current index gives the answer ∞ and $r(v_1) > r(v_6)$. $\{v_2, v_3, v_6\}$ will be pruned by the current index while v_7 will be pruned since $r(v_7) < r(v_6)$. Therefore, $L_2^{PSL}(v_6) = \{(v_1, 2)\}$.

Theorem 5 *The time complexity of PSL under one core environment is $O(\delta^2 \cdot m)$.*

Proof Let $L^{PSL} = L_{<D+1}^{PSL} = L^{PLL}$. For each label in $L^{PSL}(v)$, it has been collected by each of v 's neighbors once as candidates (Line 11). For each candidate, a query (Line 15) is called in $O(\delta)$ time. The total cost is $\sum_{v \in V} \delta d(v) \times \delta = O(\delta^2 m)$. \square

4 Index size reduction

Parallel index construction reduces the index time while leaving the index size $L^{PSL} = L^{PLL}$ unchanged. This section improves the scalability of the PSL by reducing the index size. Section 4.1 reduces the graph size using the equivalence

relationships among nodes. Section 4.2 optimizes the index size of PSL based on an observation on the label distribution.

4.1 Equivalence relation reduction

We consider the equivalence of two nodes u and v based on their neighbors. Depending on whether u and v have an edge, we define two types of equivalence relations.

Definition 9 (Node Equivalence Relations) For $u, v \in V$,

- $u \simeq_1 v$ if $N(u) = N(v)$;
- $u \simeq_2 v$ if $N(u) \cup \{u\} = N(v) \cup \{v\}$.

It can be verified that \simeq_1 and \simeq_2 are equivalence relations. Their reflexive, symmetric and transitive properties are inherited from the equality operator over node sets.

Since $u \notin N(u)$, $u \simeq_1 v$ requires that u and v have no edge while $u \simeq_2 v$ requires that u and v must have an edge.

Each equivalence relation partitions V into disjoint equivalent classes: the equivalent class of a node v includes all the nodes that are equivalent to v . We say an equivalent class is non-trivial if it includes at least two nodes. Definition 10 obtains nodes in non-trivial equivalent classes under the two equivalence relations and Lemma 11 shows that these non-trivial equivalent classes are disjoint.

Definition 10 Define three vertex sets V_1, V_2 and V_3 with

- $V_1 = \{u \in V \mid \text{there exists } v \neq u \text{ such that } u \simeq_1 v\}$
- $V_2 = \{u \in V \mid \text{there exists } v \neq u \text{ such that } u \simeq_2 v\}$
- $V_3 = V \setminus V_1 \setminus V_2$.

Example 9 In Fig. 3, $V_1 = \{v_{11}, v_{12}\}$ since $N(v_{11}) = N(v_{12}) = \{v_4, v_5\}$; $V_2 = \{v_6, v_7\}$ since $\{N(v_6) \cup v_6\} = \{N(v_7) \cup v_7\} = \{v_2, v_3, v_6, v_7\}$.

Lemma 11 V_1, V_2 and V_3 is a partition of the graph G .

Proof Since V_3 is the complement of $V_1 \cup V_2$, the three vertex sets jointly cover V . It remains to prove that $V_1 \cap V_2 = \emptyset$.

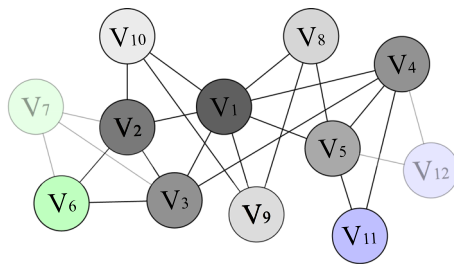


Fig. 3 Equivalence relation reduction

Let u be a node $u \in V_1 \cap V_2$. According to the definition, there exist two nodes $v \neq u$ and $w \neq u$ such that $u \simeq_1 v$ and $u \simeq_2 w$. In other words, $N(u) = N(v)$ and $N(u) \cup \{u\} = N(w) \cup \{w\}$. Since v has no edge to u while w has an edge to u , $v \neq w$. Thus, $w \in N(u) = N(v)$, namely, there is an edge between w and v . Since $v \in N(w) \setminus \{u\} \subseteq N(u)$, u and v have an edge, contradiction. Therefore, $V_1 \cap V_2 = \emptyset$. \square

According to Lemma 11, each node belongs to at most one non-trivial equivalence class constructed under the two equivalence relations. Therefore, we define the mapping function f that maps a node to the node with the smallest node ID in the corresponding non-trivial equivalent class.

Definition 11

$$f(u) = \begin{cases} \min\{v|v \simeq_1 u\}, & \text{if } u \in V_1; \\ \min\{v|v \simeq_2 u\}, & \text{if } u \in V_2; \\ u, & \text{if } u \in V_3; \end{cases} \quad (3)$$

Example 10 In Fig. 3, $f(v_{11}) = f(v_{12}) = v_{11}$; $f(v_6) = f(v_7) = v_6$; $f(u) = u$, for $u \in V_3$.

Graph Reduction. We compress the graph by eliminating all the nodes u in V_1 and V_2 and their incident edges unless $f(u) = u$. This operation transforms G to its subgraph G^s .

Example 11 In Fig. 3, $f(v_7) \neq v_7$, we delete v_7 . Similarly, $f(v_{12}) \neq v_{12}$, we delete v_{12} . Nodes u with $f(u) = u$ are kept.

Lemma 12 For any two nodes s, t with $f(s) \neq f(t)$, $\text{dist}_G(s, t) = \text{dist}_{G^s}(f(s), f(t))$.

Proof Let $p(s, t) = \{v_1, v_2, \dots, v_k\}$ be a shortest path on G from s to t and let $p^s(s, t) = \{f(v_1), f(v_2), \dots, f(v_k)\}$.

This paragraph proves that for any nodes x and y on p with $x \neq y$, $f(x) \neq f(y)$. We first show that for all $v \neq t$

on p , $f(v) \neq f(t)$: if otherwise the predecessor $pre(v)$ of v on the path $p—pre(v)$ exists since $f(s) \neq f(t)$ —can link to t directly and then reduces the path length, contradiction. Therefore, any node v with $f(v) \neq f(t)$ has a successor on p . Secondly, let $u \neq t$ be a node on p ; denote by S the equivalent class of u ; let z be the last node in S on the path. $suc(z)$, the successor of z on the path exists since $f(u) = f(z) \neq f(t)$ (from the first point). There is an edge from u to $suc(z)$ since (1) z has an edge to $suc(z)$, (2) $u, z \in S$ and (3) $suc(z) \notin S$. Thus, if $suc(z)$ is not the successor of u , then p is not a shortest path. Therefore, all nodes on p have different $f(\cdot)$ values.

It is easy to verify that if $f(u) \neq f(v)$ and there is an edge between u and v , then there is an edge between $f(u)$ and $f(v)$. Thus, $p^s(s, t)$ is a path on G^s . Since G^s is a subgraph of G , $\text{dist}_G(s, t) \leq \text{dist}_{G^s}(s, t) \leq \text{dist}_G(s, t)$. \square

Example 12 Denote by $F(V') = \{v \in V' | v = f(v)\}$ the remained nodes in a set under equivalence reduction. Table 2 shows the effectiveness of the equivalence relations on index reduction. For YOUT (TPD), around 33.15% (17.67%) and 0.45% (0.67%) of nodes are eliminated by the first and the second equivalence relation, respectively, and the index size is reduced by 31.13% (16.16%).

Query Processing. With the compressed graph, the query processing has to be adapted. We answer query $q(s, t)$ in the following four cases. (1) If $s = t$, return 0. (2) If $f(s) = f(t)$ under $s \simeq_1 t$ then return 2. (3) If $f(s) = f(t)$ under $s \simeq_2 t$, return 1. (4) Otherwise, return $q(f(s), f(t))$ in G^s .

4.2 Local minimum set elimination

The index reducing technique in this section is motivated by an observation on the PLL label distribution.

For PLL with nodes ordered in node degrees, Fig. 4 shows the label size distribution of two representative small-world networks: YouTube (denoted by YOUT) is a social network and UK-Tpd (denoted by TPD) is a web graph. The maximum degrees of YOUT and TPD are 91,751 and 63,731, respectively. It can be observed that low-degree nodes have significantly larger label sizes than the high degree nodes. This observation motivates the elimination of node labels on the nodes ranked lowest among its neighbors.

Definition 12 (Local Minimum Set) A node is local minimum node if it has the lowest rank among its neighbors. Local

Table 2 Reduced index size with equivalence relations

Dataset	Number of reduced nodes			Index space (MB)	
	$ V $	$ V_1 \setminus F(V_1) $	$ V_2 \setminus F(V_2) $	Before	After
YOUT	3,223,590	1,068,666	14,405	2141.512	1474.86
TPD	1,766,010	312,166	11,912	1783.192	1495.05

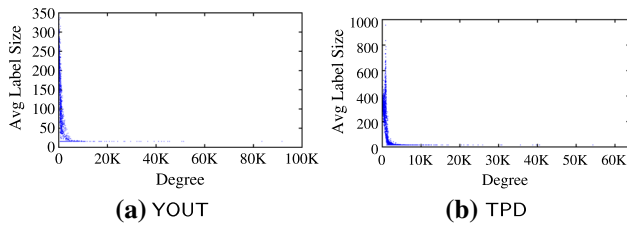


Fig. 4 PLL: degree and label size

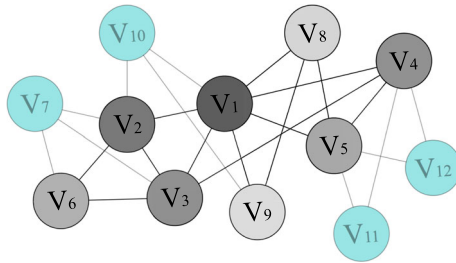


Fig. 5 Local minimum set

minimum set constitutes of local minimum nodes:

$$M(G) = \{u \in V | \text{for } \forall v \in N(u), r(u) < r(v)\}.$$

Example 13 In Fig. 5, $M(G) = \{v_7, v_{10}, v_{11}, v_{12}\}$. For example, node v_7 has the lowest rank among its neighbors.

An ideal property of a local minimum node v is that v is referred to by no node other than v itself as a hub.

Lemma 13 For any node $\forall v \in M(G)$ and any node $\forall u \in V$, v is a hub of u in L^{PSL} if and only if $v = u$.

Proof According to Theorem 1, v is a hub of u if v is the highest ranked node in S —the set of all nodes on the shortest path from u to v . Unless $u = v$, for any shortest path from u to v , there is a node $w \in N(v)$ on the path. If v is a local minimum node, $r(v) < r(w)$ and v cannot be a hub of u . \square

Construct Labels for $V \setminus M(G)$. Lemma 13 shows that removing nodes in $M(G)$ does not affect the label set of any node in $V \setminus M(G)$. However, in our propagation-based label construction, $L_d^{\text{PSL}}(v)$ is built from $L_{d-1}^{\text{PSL}}(u)$, $\forall u \in N(v)$. In other words, for a node $u \in N(v) \cap M(G)$, without $L_{d-1}^{\text{PSL}}(u)$ we cannot construct $L_d^{\text{PSL}}(u)$ using Theorem 4.

To tackle the above problem, the key finding is that nodes in $M(G)$ are independent. That is, there is no edge between nodes in $M(G)$. Thus, a node u with some of its neighbor from $M(G)$ can be saved by resorting to u 's two-hop neighbors via nodes in $M(G)$. These 2-hop neighbors will certainly fall in $V \setminus M(G)$, and their labels are ready for use.

Definition 13 (Generalized Neighbors) Given a node $v \in V \setminus M(G)$, we define two neighbor sets. $N^1(v) = N(v) \setminus M(G)$

Table 3 Reduced index size with local minimum set

Dataset	Node number		Index space (MB)	
	$ V $	$ M(G) $	Before	After
YOUT	3,223,590	2,287,357	2141.512	1234.377
TPD	1,766,010	1,151,224	1783.192	989.567

includes the neighbors of v that fall in $V \setminus M(G)$ and $N^2(v) = \{w | w \in (N(u) \setminus \{v\}), \forall u \in (N(v) \cap M(G))\}$ includes the two-hop neighbors of v connected via nodes in $M(G)$.

Example 14 In Fig. 5, since $v_9 \in V \setminus M(G)$, $N^1(v_9) = \{v_1, v_8\}$, $N^2(v_9) = \{v_1, v_2\}$.

We show that the generalized neighbors are not in $M(G)$.

Lemma 14 Given a node $v \in V \setminus M(G)$, $N^1(v) \cap M(G) = \emptyset$ and $N^2(v) \cap M(G) = \emptyset$.

Proof $N^1(v) \cap M(G) = \emptyset$ by Definition 13. Let $x \in N^2(v)$ be a node expanded from $y \in N(v) \cap M(G)$. If $x \in M(G)$, then $r(y) < r(x)$ and $r(x) < r(y)$, contradiction. \square

Example 15 In Fig. 5, $N^2(v_9) = \{v_1, v_2\}$, which are all in the set $V \setminus M(G)$.

We show a label propagation function on $V \setminus M(G)$ below.

For $\forall v \in V$ and $d \in [1, D + 1]$, denote, by $C_d(v)$, the set of hub nodes in label set $L_d^{\text{PSL}}(v)$.

Theorem 6 For each node $u \in V \setminus M(G)$

$$L_d^{\text{PSL}}(u) = \bigcup_{\substack{w \in C_{d-1}(v), \text{ for } \forall v \in N^1(u) \\ w \in C_{d-2}(v'), \text{ for } \forall v' \in N^2(u)}} L_d^{\text{PSL}}(u, w), \quad (4)$$

where $L_d^{\text{PSL}}(u, w) =$

$$\begin{cases} \emptyset, & \text{if } r(w) < r(u) \text{ or } \text{Query}(w, u, L_{<d}^{\text{PSL}}) \leq \text{dist}(w, u); \\ \{w, \text{dist}(w, u)\}, & \text{otherwise.} \end{cases} \quad (5)$$

Proof Let L'' be the labels drawn from Eq. (4). We reuse the proof of Theorem 4 by showing that the hubs L' constructed in Eq. (1) is a subset of the hubs in L'' . According to Lemma 8, $\bigcup_{v' \in N^2(u)} C_{d-2}(v')$ is a super set of $\bigcup_{v \in N(v) \cap M(G)} C_{d-1}(v)$, besides, $N(u) = N^1(u) \cup (N(u) \cap M(G))$, thus $\bigcup_{v \in N(u)} C_{d-1}(v) \subseteq \bigcup_{v \in N^1(u)} C_{d-1}(v) \cup \bigcup_{v' \in N^2(u)} C_{d-2}(v')$ which completes the proof. \square

Example 16 Table 3 shows the effectiveness on reducing the index size using local minimum set. For YOUT (TPD), the local minimum set eliminates about 70.95% (65.18%) nodes and reduces the index size by 42.4% (44.5%).

Table 4 Local minimum set: index and query time

Dataset	Index time (s)		Query time (s)	
	Before	After	Before	After
YOUT	23.805	15.786	1.13E-06	1.71E-06
TPD	18.997	13.71	1.80E-06	3.71E-06

Query Processing. The reduced index provides the labels for nodes in $V \setminus M(G)$. When it comes to query processing, we can recover the labels of nodes in $M(G)$ with the union of the labels of neighbors. For a query $q(s, t)$, without loss of generality, if $s \in M(G)$ and $t \in V \setminus M(G)$, we swap s and t . To reduce the online cost, we use a hash join to produce the 2-hop distances. Let H be a table of size $|V \setminus M(G)|$ where $H(w)$ records the labeled distance in $L^{\text{PSL}}(s)$ with hub w . $H(w) = \infty$ if w is not a hub of s . Since the label set $L^{\text{PSL}}(s)$ may not be available, we construct H in two cases.

- If $s \in V \setminus M(G)$, we hash the labels in $L^{\text{PSL}}(s)$ by letting $H(v) = \text{dist}(s, v)$ for each hub v of s .
- Otherwise, we construct labels of s by visiting neighbors $w \in N(s)$ of s and update $H(v)$ with $\text{dist}(v, w) + 1$ for each hub v of w — $H(v)$ only keeps the minimum value along the updates.

After H being constructed, we generate labels of t in a similar way and instead of updating the table H , we fetch the value stored in the table H under the same hub node and then compose a 2-hop distance.

Note that the hash table H can be maintained across different queries without initialization: we keep a dirty log and recover H after processing each query.

Lemma 15 When $s, t \in M(G)$, the time cost of distance query is $O(\sum_{a \in N(s)} |L^{\text{PSL}}(a)| + \sum_{b \in N(t)} |L^{\text{PSL}}(b)|)$.

Proof For s , we store nodes in $\{L^{\text{PSL}}(a) | a \in N(s)\}$ in H . For t , we scan the nodes in $\{L^{\text{PSL}}(b) | b \in N(t)\}$ to gain the distance. The linear scan takes $O(|\{L^{\text{PSL}}(a) | a \in N(s)\}| + |\{L^{\text{PSL}}(b) | b \in N(t)\}|)$ time in total. \square

Example 17 Table 4 shows the index time and query time in a 45-core environment. Local minimum set technique reduces, for YOUT (TPD), the index time by 33.69% (27.83%) at a cost of $1.5 \times$ ($2.06 \times$) query time. The trade-off is worthwhile since the query time is still in micro-seconds.

5 Index optimization with betweenness-based node order

The PSL proposed in Sect. 3 parallelizes PLL in a multi-core environment, and the main bottleneck of this labeling method

is the unaffordable index size. The two index reduction techniques proposed in Sect. 4 are built upon a node order which is, by default, degree based. To further reduce the index size, this section investigates the application of betweenness-based node order in PSL. As suggested by [31] and verified by our preliminary experimentation (Exp-9, Sect. 7), betweenness-based node order leads to a smaller index size. The difficulty in applying the betweenness-centrality to PSL is twofold. 1) The computation of the betweenness centrality for all the nodes is computationally expensive ($O(mn)$ [13]) for big graphs. 2) The best practice of cost-effectively optimizing PSL with approximate betweenness is unknown. A better estimation of k -betweenness leads to a smaller index size; however, improving estimation precision can be exhaustive. Section 5.1 first proposes a sampling-based approach for estimating k -betweenness; to further improve the estimation efficiency, Sect. 5.2 presents a pool-based sampling algorithm. Section 5.3 introduces an algorithm in engaging the betweenness estimation in PSL for index reduction.

5.1 Basic sampling

Exact k -betweenness of a node $v \in V$ summarizes the partial betweenness $\text{kbc}_s(v)$ over all source nodes s in V . However, only a small number of sources s contribute to the computation of $\text{kbc}(v)$: Lemma 3 shows that nodes s with $\text{dist}(s, v) = 0$ or $\text{dist}(s, v) \geq k$ has $\text{kbc}_s(v) = 0$. These useless sources can be safely removed for v .

Definition 14 (*k-Reachable Set*) The k -reachable set of a vertex $v \in V$ is defined as $R(v) = \{s | 0 < \text{dist}(s, v) < k\}$.

Example 18 To estimate $\text{kbc}(v_9)$ (with $k = 2$) in Fig. 1, we only consider source nodes in $R(v_9) = \{v_1, v_2, v_3, v_4, v_5, v_8, v_{10}\}$ since nodes w outside $R(v_9)$ make the partial betweenness $\text{kbc}_w(v_9)$ zero.

Under the framework of betweenness approximation [7], random samples need to be selected from $R(v)$. Suppose we randomly select some nodes S from $R(v)$ for a node v . For each sample $s \in S$, $\text{kbc}_s(v)$ can be computed by undertaking a graph traversal from s [14]. We estimate $\text{kbc}(v)$ with

$$\widetilde{\text{kbc}}(v) = \left(\sum_{s \in S} \text{kbc}_s(v) \right) \cdot \frac{|R(v)|}{|S|}.$$

Lemma 16 shows that $\widetilde{\text{kbc}}(v)$ is an unbiased estimator of $\text{kbc}(v)$.

Lemma 16 $E(\widetilde{\text{kbc}}(v)) = \text{kbc}(v)$, for $\forall v \in V$.

Proof For $\forall s \in R(v)$, we define a random variable $X_s = \text{kbc}_s(v) \cdot |R(v)|$. We have $E(X_s) = \sum_{s \in R(v)} \frac{1}{|R(v)|} X_s =$

$\sum_{s \in R(v)} \frac{1}{|R(v)|} \text{kbc}_s(v) \cdot |R(v)| = \text{kbc}(v)$. When aggregating X_s over samples s in S , we have $E(\widetilde{\text{kbc}}(v)) = E(\sum_{s \in S} \text{kbc}_s(v) \cdot \frac{|R(v)|}{|S|}) = E(\sum_{s \in S} X_s \cdot \frac{1}{|S|}) = E(X_s) = \text{kbc}(v)$. \square

Lemma 17 Suppose $K = \max_{s \in S} (\text{kbc}_s(v))$,

$$P(|\widetilde{\text{kbc}}(v) - \text{kbc}(v)| > \epsilon) \leq 2 \exp\left(-2|S| \cdot \left(\frac{\epsilon}{K \cdot |R(v)|}\right)^2\right).$$

Proof Let X_1, X_2, \dots, X_q be independent random variables with values in $[a, b]$, and $\bar{X} = \frac{X_1 + X_2 + \dots + X_q}{q}$, then $P(|E(\bar{X}) - \bar{X}| > \epsilon) \leq 2 \exp(-2q \cdot (\frac{\epsilon}{b-a})^2)$ by Hoeffding's inequality [23]. For any $s \in S$, we define $X_s = \text{kbc}_s(v) \cdot |R(v)|$, then $E(X_s) = \text{kbc}(v)$, $\bar{X}_s = \text{kbc}(v)$, $q = |S|$, $a = 0$, $b = K \cdot |R(v)|$. Plugging these terms into Hoeffding's inequality proves the lemma. \square

Discussion. Given a node $v \in V$, by Lemma 17, for a given $\epsilon \in R^+$ and $\delta \in (0, 1)$, if $|S| \geq \frac{\ln(\frac{2}{\delta}) \cdot K^2 \cdot |R(v)|^2}{2\epsilon^2}$, we obtain an estimation of $\text{kbc}(v)$ within an additive error ϵ with a probability at least δ [23]. The required sample size, unfortunately, is very large. Although there are techniques to reduce the sample size [11,37,38], there are two drawbacks of this basic sampling approach: (i) each node v needs to compute $|R(v)|$ to reach an unbiased estimator $\widetilde{\text{kbc}}(v)$ of $\text{kbc}(v)$; (ii) node v precomputes $R(v)$ to select samples. The cost of obtaining $R(v)$ (and $|R(v)|$) for $\forall v \in V$ by performing $n = |V|$ BFS (with a length limited to k) is no better than the exact k -betweenness computation.

5.2 Pool-based sampling

Size Estimation. To solve the first drawback of the above sampling method, we select a pool S_{size} of nodes to approximate $|R(v)|$ for all nodes $v \in V$. Specifically, for each node $s \in S_{\text{size}}$, we conduct a k -bounded BFS from s which only visits nodes that are $< k$ hops away from s . Suppose v has been visited $n_{\text{size}}(v)$ times by the k -bounded BFS from s (that is, there are $n_{\text{size}}(v)$ samples in S_{size} that belong to $R(v)$), we estimate $|R(v)|$ with

$$\widetilde{R}(v) = n_{\text{size}}(v) \cdot \frac{n}{|S_{\text{size}}|}.$$

Lemma 18 shows that $\widetilde{R}(v)$ is an unbiased estimator of $|R(v)|$.

Lemma 18 $E(\widetilde{R}(v)) = |R(v)|$, for $\forall v \in V$.

Proof For each sample $s \in S_{\text{size}}$, we define a random variable X_s to indicate whether s is in $R(v)$, that is, $X_s = \begin{cases} 1, & \text{if } s \in R(v) \\ 0, & \text{otherwise} \end{cases}$. Then, $P(X_s = 1) = \frac{|R(v)|}{n}$ and $E(X_s) = \frac{|R(v)|}{n}$. Thus, $E(n_{\text{size}}(v)) = \sum_{s \in S_{\text{size}}} E(X_s) = \frac{|R(v)|}{n} \cdot |S_{\text{size}}|$, and $E(\widetilde{R}(v)) = |R(v)|$. \square

Algorithm 3: Size Estimation

Input: Graph $G(V, E)$, hop k , time budget T_s
Output: $S_{\text{size}}, n_{\text{size}}(v)$ for $\forall v \in V$

```

1  $S_{\text{size}} \leftarrow \emptyset$ ;
2  $n_{\text{size}}(v) \leftarrow 0$ , for  $\forall v \in V$ ;
3 for sampling time  $\leq T_s$  do
4    $s \leftarrow$  a node chosen uniformly at random from  $V$ ;
5    $S_{\text{size}} \leftarrow S_{\text{size}} \cup \{s\}$ ;
6   Let  $\sigma_{s,s} \leftarrow 1$  and  $\text{dist}(s) \leftarrow 0$ ;
7   For  $\forall v \in V \setminus \{s\}$ , let  $\sigma_{s,v} \leftarrow 0$  and  $\text{dist}(v) \leftarrow -1$ ;
8    $\text{curr} \leftarrow \{s\}$ ;  $\text{next} \leftarrow \emptyset$ ;
9   for  $i = 0, 1, \dots, k - 2$  do
10    for  $\forall v \in \text{curr}$  and  $\forall w \in N(v)$  do
11      if  $\text{dist}(w) = -1$  then
12         $\text{dist}(w) \leftarrow \text{dist}(v) + 1$ ;
13         $n_{\text{size}}(w) \leftarrow n_{\text{size}}(w) + 1$ ;
14         $\text{next} \leftarrow \text{next} \cup \{w\}$ ;
15     $\text{curr} \leftarrow \text{next}$ ,  $\text{next} \leftarrow \emptyset$ ;
16 return  $S_{\text{size}}, n_{\text{size}}(v)$  for  $\forall v \in V$ ;
```

Algorithm 3 estimates, for $\forall v \in V$, the size $n_{\text{size}}(v)$ of $R(v)$, within the sampling time budget T_s . For a random sample s (Line 4), we append s in S_{size} (Line 5) and perform a k -bounded BFS from s (Lines 6–15). For each newly visited node v (i.e., $s \in R(v)$), $n_{\text{size}}(v)$ is increased by 1 (Line 13). The process continues until the sampling time goes beyond the budget T_s (Line 3).

Partial Betweenness Estimation. To solve the second drawback of the basic sampling approach, we select a pool of nodes S_{bc} to compute k -partial betweenness for all nodes $v \in V$. Among the samples in S_{bc} , suppose $n_{\text{bc}}(v)$ nodes (denoted as S_v) are included in $R(v)$ for a certain v , we summarize the partial k -betweenness of v over S_v to obtain $\kappa(v)$, for each individual node $v \in V$; $\kappa(v)$ shall be used to estimate $\text{kbc}(v)$.

$$\kappa(v) = \sum_{s \in S_v} \text{kbc}_s(v).$$

In this way, we avoid sampling from $R(v)$ for each node in V .

Algorithm 4 estimate $\kappa(v)$ and $n_{\text{bc}}(v)$, for $\forall v \in V$ (Lines 1–27) within time budget T_s . We repeatedly select a sample s (Line 3) uniformly at random, until the time budget T_s is consumed (Line 2). Given s , we follow the method introduced in [14] to compute $\text{kbc}_s(v)$, for $\forall v \in V$. Note that in this process, $\text{kbc}_s(v)$ is added to $\kappa(v)$, and $n_{\text{bc}}(v)$ is increased by 1.

We first conduct a k -bounded BFS from s (Lines 1–18), aiming at computing $\sigma_{s,v}$, the number of shortest paths from s to v , for $\forall v \in R(s)$ (equivalently $s \in R(v)$). Specifically, we use curr and next to store nodes expanded in the current round and the nodes to expand in the next round. $\sigma_{s,v}$ is initialized with zero for all $v \in V$, except for s , whose $\sigma_{s,s}$

is set to 1; $\text{dist}(v)$ is initialized to -1 for all v , except for s , which is set to 0 (Lines 6–7). We first insert s into curr to start the BFS (Line 8), and then we explore nodes within distance k to s (Line 9): for each node $v \in \text{curr}$ (Line 10), we check v 's neighbor w (Line 11). If w is not visited before (Line 12), $\text{dist}(w)$ is updated to $\text{dist}(v) + 1$ (Line 13), and w is appended to next and S (Lines 14–15). If w is one hop farther than v regarding s , we increase $\sigma_{s,w}$ by adding $\sigma_{s,v}$ to it (Lines 16–17). Then, next is assigned to curr for the next round (Line 18).

When all nodes within distance k to s have been stored in stack S , we perform a backward BFS to compute $\text{kbc}_s(v)$, for $\forall v \in V$ (Lines 19–27). Specifically, $\text{kbc}_s(v)$ is initialized as zero (Line 19), and we visit nodes w in S reversely—in the order of non-increasing distance to s (Line 21). For each neighbor v of w , if v is one hop closer than w regarding s , $\text{kbc}_s(w)$ is used to update $\text{kbc}_s(v)$ (Lines 22–24). For each $v \neq s$, $\text{kbc}_s(v)$ is added to $\kappa(v)$, and $n_{\text{bc}}(v)$ is increased by 1 (Lines 26–27).

Algorithm 4: Partial Betweenness Estimation

Input: Graph $G(V, E)$, hop k , time budget T_s
Output: $n_{\text{bc}}(v), \kappa(v)$ for $\forall v \in V$

```

1  $n_{\text{bc}}(v) \leftarrow 0, \kappa(v) \leftarrow 0$ , for  $\forall v \in V$ ;
2 while sampling time  $\leq T_s$  do
3    $s \leftarrow$  a random node in  $V$ ;
4   // Forward BFS
5    $\text{curr} \leftarrow \emptyset, \text{next} \leftarrow \emptyset$ ;
6    $S \leftarrow$  an empty stack;
7   For  $\forall v \in V \setminus \{s\}, \sigma_{s,v} \leftarrow 0, \text{dist}(v) \leftarrow -1$ ;
8    $\sigma_{s,s} \leftarrow 1, \text{dist}(s) \leftarrow 0$ ;
9    $\text{curr} \leftarrow \text{curr} \cup \{s\}$ ;
10  for  $i = 0, 1, \dots, k - 1$  do
11    for  $\forall v \in \text{curr}$  do
12      for  $\forall w \in N(v)$  do
13        if  $\text{dist}(w) = -1$  then
14           $\text{dist}(w) \leftarrow \text{dist}(v) + 1$ ;
15           $\text{next} \leftarrow \text{next} \cup \{w\}$ ;
16           $S \leftarrow S \cup \{w\}$ ;
17        if  $\text{dist}(w) = \text{dist}(v) + 1$  then
18           $\sigma_{s,w} \leftarrow \sigma_{s,w} + \sigma_{s,v}$ ;
19     $\text{curr} \leftarrow \text{next}, \text{next} \leftarrow \emptyset$ ;
20  // Backward BFS
21   $\text{kbc}_s(v) \leftarrow 0, \forall v \in V$ ;
22  while  $S \neq \emptyset$  do
23     $w \leftarrow$  pop from  $S$ ;
24    for  $v \in N(w)$  do
25      if  $\text{dist}(w) \neq \text{dist}(v) + 1$  then continue;
26       $\text{kbc}_s(v) \leftarrow \text{kbc}_s(v) + \frac{\sigma_{s,v}}{\sigma_{s,w}} \cdot (1 + \text{kbc}_s(w))$ ;
27      if  $v \neq s$  then
28         $\kappa(v) \leftarrow \kappa(v) + \text{kbc}_s(v)$ ;
29         $n_{\text{bc}}(v) \leftarrow n_{\text{bc}}(v) + 1$ ;
30  return  $n_{\text{bc}}(v), \kappa(v)$  for  $\forall v \in V$ ;

```

Algorithm 5: Order Generation

Input: Graph $G(V, E)$, hop k , time budget T, θ
Output: $r(v)$ for $\forall v \in V$

```

1  $S_{\text{size}}, n_{\text{size}}(v) \leftarrow$  Algorithm 3( $G, k, \theta T$ ), for  $\forall v \in V$ ;
2  $n_{\text{bc}}(v), \kappa(v) \leftarrow$  Algorithm 4( $G, k, (1 - \theta)T$ ), for  $\forall v \in V$ ;
3  $\widetilde{\text{kbc}}(v) \leftarrow \frac{\kappa(v)}{n_{\text{bc}}(v)} \cdot (n_{\text{size}}(v) \cdot \frac{n}{|S_{\text{size}}|})$ , for  $\forall v \in V$ ;
4 Generate  $r(v)$  in non-increasing order of  $\widetilde{\text{kbc}}(v)$ ;
5 return  $r(v)$  for  $\forall v \in V$ ;

```

To analyze the estimation accuracy of the pool-based sampling, we focus on $S_v = \{v \in S_{\text{bc}} | v \in R(v)\}$ and its size $n_{\text{bc}}(v) = |S_v|$, for each $v \in V$. We show that for each v , the size $n_{\text{bc}}(v)$ is proportional to $|R(v)|$.

Lemma 19 $E(n_{\text{bc}}(v)) = |S_{\text{bc}}| \times \frac{|R(v)|}{n}$.

Proof For a node that is chosen uniformly at random from V , it falls in $R(v)$ with probability $\frac{|R(v)|}{n}$. Aggregating this probability over all nodes in S_{bc} derives the expectation $E(n_{\text{bc}}(v)) = |S_{\text{bc}}| \times \frac{|R(v)|}{n}$. \square

Order Generation. With the outputs of size estimation and partial betweenness estimation, we show that

$$\widetilde{\text{kbc}}(v) = \frac{\kappa(v)}{n_{\text{bc}}(v)} \cdot (n_{\text{size}}(v) \cdot \frac{n}{|S_{\text{size}}|}) \tag{6}$$

is an unbiased estimator of $\text{kbc}(v)$.

Lemma 20 $E(\widetilde{\text{kbc}}(v)) = \text{kbc}(v)$, for $\forall v \in V$.

Proof Given a node $v \in V$, we define a random variable $X_s = \text{kbc}_s(v) \cdot n_{\text{size}}(v) \cdot \frac{n}{|S_{\text{size}}|}$, for $\forall s \in R(v)$. Then, $E(X_s) = \frac{1}{|R(v)|} \cdot \sum_{s \in R(v)} \text{kbc}_s(v) \cdot E(n_{\text{size}}(v) \cdot \frac{n}{|S_{\text{size}}|}) = \sum_{s \in R(v)} \text{kbc}_s(v) = \text{kbc}(v)$ (size estimation and betweenness estimation are independent). Suppose samples S_{bc} are used to estimate the betweenness, among which nodes $S_v \subseteq S_{\text{bc}}$ are included in $R(v)$. The size of S_v is $n_{\text{bc}}(v)$. Then, $E(\widetilde{\text{kbc}}(v)) = E(\sum_{s \in S_v} \text{kbc}_s(v) \cdot \frac{n_{\text{size}}(v)}{n_{\text{bc}}(v)} \cdot \frac{n}{|S_{\text{size}}|}) = E(\sum_{s \in S_v} X_s \cdot \frac{1}{n_{\text{bc}}(v)}) = E(X_s) = \text{kbc}(v)$. \square

By applying Lemma 17, the accuracy is given below.

Lemma 21 Suppose $K = \max_{s \in S_v} (\text{kbc}_s(v))$,

$$P(|\widetilde{\text{kbc}}(v) - \text{kbc}(v)| > \epsilon) \leq 2 \exp(-2n_{\text{bc}}(v) \cdot (\frac{\epsilon}{K \cdot |R(v)|})^2).$$

With the estimation $\widetilde{\text{kbc}}(v)$ of $\text{kbc}(v)$ computed for each node $v \in V$, the betweenness-based node order r is set such that for any $u, v \in V, r(u) > r(v)$ if

- $\widetilde{\text{kbc}}(u) > \widetilde{\text{kbc}}(v)$;
- $\widetilde{\text{kbc}}(u) = \widetilde{\text{kbc}}(v)$, $\text{ID}(u) > \text{ID}(v)$.

Algorithm 5 shows the order generation algorithm. First, Algorithm 3 (with sampling time budget θT) and Algorithm 4 (with sampling time budget $(1 - \theta)T$) are called to estimate size and partial betweenness of each node (Lines 1–2). Then, $\widetilde{\text{kbc}}(v)$ is computed based on Eq. (6) (Line 3). Finally, $r(v)$ is determined by the above rule (Line 4). The parameter θ controls the time used in Algorithm 3 and Algorithm 4. In practice, we set the parameter θ as 0.2 since it leads to a good effect when k -betweenness is used for ordering nodes.

Lemma 22 *The time cost of Algorithm 5 is $O(|S|(n + m))$ where $|S|$ is the number of samples used in stage 1 and stage 2.*

Remarks. In Algorithm 5, instead of giving a pre-defined sample size, the sample size is controlled adaptively by the sampling time—the estimation accuracy will improve if more time is given.

5.3 Improved betweenness-based node order

Recall that the index reduction techniques proposed in Sect. 4 remove the local minimum set $M(G)$ which, in the current node order, is determined by the k -betweenness of nodes in V . The remaining nodes $V \setminus M(G)$, however, may have different k -betweenness in the updated graph structure. Therefore, it is desirable to recompute k -betweenness for $V \setminus M(G)$ once $M(G)$ is eliminated for a better approximation.

Virtual Graph. To recompute k -betweenness, one challenge is that if two nodes $u, v \in V \setminus M(G)$ are connected only by nodes in $M(G)$, then u and v are disconnected in the updated graph. To this end, given a graph $G(V, E)$, we define a virtual graph $\overline{G}(\overline{V}, \overline{E})$ with $\overline{V} = V \setminus M(G)$ as its node set. Recall Definition 13, we add two types of edges to \overline{E} for each node $u \in \overline{V}$, (u, v) with $l(u, v) = 1$ for every $v \in N^1(v)$ and (u, v) with $l(u, v) = 2$ for every $v \in N^2(v) \setminus N^1(v)$. The edge set is sufficient to keep the connectivity: due to the property of the local minimum reduction, if $u \in M(G)$ for some node $u \in V$, then $N(v) \cap M(G) = \emptyset$; therefore, to retain the connectivity, we only need to consider $u, v \in \overline{V}$ that are connected only by one node in $M(G)$.

Example 19 In Fig. 5, since $v_8 \in N^1(v_9)$, we have the edge (v_9, v_8) with weight 1 in \overline{G} ; since $v_2 \in N^2(v_9) \setminus N^1(v_9)$, we have the edge (v_9, v_2) with weight 2 in \overline{G} .

⁵ For the convenience of presentation, we replace an edge (u, v) of length 2 with two unit-weighted edges $(u, w), (w, v)$ with a new node w interpolated in between.

Algorithm 6: Improved Order Generation

Input: Graph $G(V, E)$, hop k , sample time T, θ
Output: $r(v)$, for $\forall v \in V$

- 1 $\text{kbc}(v), \forall v \in V \leftarrow \text{Algorithm 5}(G, k, (1 - \theta)T)$;
- 2 $M(G) \leftarrow \emptyset$;
- 3 **for** $v \in V$ **do**
- 4 **for** $w \in N(v)$ **do**
- 5 **if** $\text{kbc}(v) > \text{kbc}(w)$ **then** continue;
- 6 $M(G) \leftarrow M(G) \cup \{v\}$;
- 7 $\overline{V} \leftarrow V \setminus M(G)$;
- 8 $\overline{E} \leftarrow \{(u, v) | u, v \in V \setminus M(G)\}$;
- 9 **for** $v \in \overline{V}$ **do**
- 10 $N^1(v) \leftarrow \emptyset, N^2(v) \leftarrow \emptyset$;
- 11 **for** $w \in N(v) \cap \overline{V}$ **do** insert $w \rightarrow N^1(v)$;
- 12 **for** $w \in N(v) \cap M(G)$ **do**
- 13 **for** $u \in N(w)$ **do**
- 14 **if** $u \neq v$ **then**
- 15 Insert $u \rightarrow N^2(v)$;
- 16 **for** $w \in N^1(v)$ **do**
- 17 Insert edge (v, w) with weight 1 in \overline{E} ;
- 18 **for** $w \in N^2(v) \setminus N^1(v)$ **do**
- 19 Insert edge (v, w) with weight 2 in \overline{E} ;
- 20 **For** $\forall v \in \overline{V}, S_{\text{size}}, n_{\text{size}}(v) \leftarrow \text{Algorithm 3}(\overline{G}(\overline{V}, \overline{E}), k, \theta T)$;
- 21 **For** $\forall v \in \overline{V}, \kappa(v), n_{\text{bc}}(v) \leftarrow \text{Algorithm 4}(\overline{G}(\overline{V}, \overline{E}), k, \theta T)$;
- 22 $\widetilde{\text{kbc}}(v) \leftarrow \frac{\kappa(v)}{n_{\text{bc}}(v)} \cdot (n_{\text{size}}(v) \cdot \frac{n}{|S_{\text{size}}|})$, for $\forall v \in \overline{V}$;
- 23 Sort $r(v)$ in non-increasing order of $\widetilde{\text{kbc}}(v)$, for $\forall v \in \overline{V}$;
- 24 Set $r(v)$ as the minimum among its neighbors, for $\forall v \in \{V \setminus \overline{V}\}$;
- 25 **return** $r(v)$, for $\forall v \in V$;

Improved Order Generation. The sampling algorithm considering local minimum set $M(G)$ elimination is in Algorithm 6. Similar to Algorithm 5, we set θ to 0.2 in practice. First, Algorithm 5 is invoked to compute $\widetilde{\text{kbc}}(v), \forall v \in V$ (Line 1), and nodes v with the minimum $\widetilde{\text{kbc}}(v)$ among $N(v)$ constitute $M(G)$ (Lines 2–6). Then, for each node $v \in \overline{V} = V \setminus M(G)$, $N^1(v)$ and $N^2(v)$ are formed according to Definition 13 (Lines 8–15), where edges \overline{E} are formed in Lines 8 and 16–19.

Afterward, Algorithm 3 approximates the size of $R(v)$ in \overline{G} , for $\forall v \in \overline{V}$; Algorithm 4 obtains k -partial betweenness $\widetilde{\text{kbc}}(v)$ in \overline{G} , for $\forall v \in \overline{V}$. The k -betweenness of $v \in \overline{V}$ in \overline{G} is then computed by the outputs of the above two algorithms (Line 22). For nodes in \overline{V} , their orders are defined by $\widetilde{\text{kbc}}(v)$ in \overline{G} (Line 24), while we enforce nodes in $V \setminus \overline{V}$ to have the minimum orders—these nodes remain to be local minimum set after the re-computing (Line 25).

6 Related work

Indexing shortest distances for fast online query processing has been extensively studied. A recent experimental comparison on distance labeling algorithms can be found in [31].

Distance Labeling on Small-world Networks. To index shortest distances for small-world networks, existing solutions either build a partial index to assist the online search algorithms [5,20,22] or build a complete index to fully support the distance query [4,26]. The solutions in the latter category require a larger index but will obtain much faster query processing time.

In the first category, Is-label approach first determines the vertex hierarchy through the independent set and then creates the label for each node by this hierarchy structure [20]. Tree decomposition is used in [5] to discover the core-fringe structure of social networks, and then index is created on these two separate parts. Shortest path trees of high-degree nodes are used [22] as index to guide the online searching to process the distance query.

In the second category, PLL [4] constructs the index by performing pruned BFS whose detail is given in Sect. 2.3. The hop doubling approach in [26] applies generation rules to join the short paths to long paths, until the whole paths are covered. Compared to PLL, the algorithm proposed in [26] uses less memory but will spend much more index time.

Distance Labeling on Road Networks. For distance indexing approaches on road networks, the approach in [2] constructs the index by eliminating the high ranking nodes and add it to the labels of its neighbors. The approach proposed by Wei [46] first decomposes the graph into a tree as the index, and then the distance of two nodes is answered through this index using dynamic programming. The pruned highway labeling approach proposed by Akiba et al. [3] decomposes the road network into disjoint paths and the label of a node include the distance to some nodes of the paths. A hierarchical hop-based index is proposed in [33] to answer shortest-distance queries in a road network with bounded query processing time and index size. More details about the distance query on road networks can be found in [31,47].

Approximate Distance Labeling. For approximate distance labeling algorithms, the basic idea is to select nodes as landmarks and then precompute the distances from the landmarks to all the other nodes. The distance between any node pair can be estimated using triangle inequality [15,35]. Online processing on landmarks is used to improve the precision [36,44]. However, on small-world networks, the relative error becomes significant since the distances are bounded by the small diameter.

Betweenness Computation. Betweenness was proposed by Freeman [19]. The best exact computation algorithm incurs $O(nm)$ [13], which has been confirmed to be almost optimal for both sparse [10] and dense graphs [1]. Due to the complexity, numerous approximation algorithms are given to trade accuracy for speed. Pioneering work was done in [24] by using a sampling-based approach, and subsequent studies aimed at reducing the sampling costs [11,37,38]. For example, Matteo et al. [37] applied VC-dimension theory

to calculate the sample size required to achieve the desired approximation. The computed sample size is independent of the number of vertices but depends only on the graph diameter (i.e., the longest shortest path in the graph). To eliminate the dependence on the graph diameter and to further reduce the required sample size, Matteo et al. [38] used the concepts of Rademacher averages and pseudodimension to accelerate the betweenness approximation. An experimental comparison of approximate algorithms is presented in the literature [6] to validate the efficiency and accuracy of various methods. Another line of direction investigates variations of betweenness to reduce the computation costs [14,18,34]. k -betweenness used in this paper belongs to this category [14], and we devise approximation algorithms to compute it fast.

k -betweenness is an approximate notion of betweenness: when k reaches the graph diameter (the length of the longest shortest path in the graph), k -betweenness becomes betweenness. We use k -betweenness to holistically optimize the node order given the time resource in computing the node order. An adequate k strikes a balance between (i) the gap between the k -betweenness and betweenness and (ii) the gap introduced by the sampling-based estimation of the k -betweenness—a larger k reduces the first gap while increases the second gap. In this paper, k is carefully chosen to holistically optimize the node order. As a type of centrality measures, k -betweenness can be used to identify important nodes in networks, such as biological networks [25], virus propagation networks [32], terrorist networks [16], and transportation networks [21]. The approximation algorithms proposed in the paper can produce elegant estimation in a given sampling time budget and thus be beneficial for the tasks on the above networks.

Extensions from [30]. This work is an extension of the conference version [30]. Compared to [30], we make the following novel contributions. (1) In Sect. 2.3, the benefits (index reduction) and challenges (quadratic computation) of using betweenness-based node order in distance labeling are discussed. (2) In Sect. 2.4, betweenness-related concepts are introduced, including betweenness and its variation k -betweenness. (3) In Sect. 5, approximation algorithms for k -betweenness computation and how to use betweenness estimation in the index construction process are presented. (4) In Sect. 7, corresponding experiments are conducted to verify the effectiveness of distance labeling using betweenness-based node order.

7 Experimental results

In this section, we first validate the effects of parallelism and compression techniques in Sect. 7.1, followed by the evaluation of k -betweenness as a node order in Sect. 7.2.

All algorithms used in the experiments were implemented in C++ and compiled with GNU GCC 4.8.5 and -O3 level

Table 5 The description of the datasets

Name	Dataset	n	m	Type
DELI	Delicious ⁶	536,109	1,365,961	Social Network
GP	GPlus ⁶	211,188	1,506,896	Social Network
LAST	Lastfm ⁶	1,191,806	4,519,330	Social Network
GOOG	Google ⁷	875,713	5,105,039	Web Graph
AMAZ	Amazon ⁸	735,323	5,158,388	Social Network
DIGG	Digg ⁶	770,800	5,907,132	Social Network
FLIX	Flixster ⁹	2,523,386	7,918,801	Social Network
TREC	Trec ⁹	1,601,787	8,063,026	Web Graph
YOUT	YouTube ⁹	3,223,589	9,375,374	Social Network
SKIT	Skitter ⁹	1,696,415	11,095,298	Internet Topology
TWIT	Twitter ⁷	456,631	14,855,875	Social Network
HUDO	Hudong ⁶	1,984,485	14,869,484	Web Graph
PET	Petster ⁹	623,766	15,699,276	Social Network
BAID	Baidu ⁶	2,141,301	17,794,839	Web Graph
TPD	UK-Tpd ⁸	1,766,010	18,244,650	Web Graph
DBLP	DBLP ⁹	1,314,050	18,986,618	Coauthorship
TOPC	Topcats ⁷	1,791,489	28,511,807	Web Graph
POK	Pokec ⁷	1,632,803	30,622,564	Social Network
FLIC	Flickr ⁹	2,302,925	33,140,017	Social Network
HOST	UK-Host ⁸	4,769,354	50,829,923	Web Graph
STAC	Stack ⁷	6,024,271	63,497,050	Interaction
LJ	Ljournal ⁸	5,363,260	79,023,142	Social Network
FB	Facebook ⁶	58,790,783	92,208,195	Social Network
INDO	Indochina ⁸	7,414,866	194,109,311	Web Graph
SINA	Sina ⁶	58,655,850	261,321,071	Social Network
WIKI	Wiki ⁹	12,150,976	378,142,420	Web Graph
ARAB	Arabic ⁸	22,744,080	639,999,458	Web Graph
IT	IT-2004 ⁸	41,291,594	1,150,725,436	Web Graph
SK	SK-2005 ⁸	50,636,154	1,949,412,601	Web Graph
UK	UK-2006 ⁸	77,741,046	2,965,197,340	Web Graph

optimization. All experiments were conducted on a machine with 48 CPU cores and 384 GB main memory running Linux (Red Hat Linux 4.8.5, 64 bit). Each CPU core is Intel Xeon 2.1GHz. The parallelized programs are supported by the OpenMP framework. We set the cut-off time as 24 hours.

7.1 Test of parallelism and compression

Algorithms. We compare our proposed algorithms against the state-of-the-art algorithm PLL [4]. Our techniques include the following three methods:

- PSL: the parallelized distance labeling technique introduced in Sect. 3.
- PSL⁺: PSL with the equivalence relation elimination technique as introduced in Sect. 4.1.

- PSL*: PSL with the equivalence relation elimination technique plus the local minimal set elimination technique as introduced in Sect. 4.2.

Datasets. We conducted experiments on 30 real-world graphs whose properties are shown in Table 5. The largest graph has more than 2.9 billion edges. The datasets are from various types of small-world networks including social networks, web graphs, Internet topology graphs, coauthorship graphs, and interaction networks. All graphs were downloaded from Network Repository⁶ [39], Stanford Large

⁶ <http://networkrepository.com/index.php>.

Network Dataset Collection⁷ [28], Laboratory for Web Algorithms⁸ [8,9], and the Koblenz Network Collection⁹ [27].

Exp 1: Index Time on a Single Core. We compare the index time of PLL with PSL, PSL⁺ and PSL* on a single core. Note that the bit-parallel technique introduced in [4] is used for all methods since it is a separate optimization which can be plugged into all distance labeling methods.

Figure 6 shows that PSL has an index time comparable to PLL, while PSL⁺ and PSL* reduce the index time of PLL—a by-product of the index reduction. For example, on the dataset ARAB, PSL⁺ and PSL* successfully constructed the index while PLL and PSL failed.

Exp 2: Index Time on Multiple Cores. Fig. 7 shows the index time of PSL, PSL⁺ and PSL* on 45 cores. Compared to the single-core results shown in Fig. 6, all the three methods have a significant speedup. This speedup allows PSL to index multiple massive graphs, e.g., LJ, ARAB and SK, that cannot be indexed on a single core. PSL* succeeded in indexing all the graphs while both PSL and PSL⁺ failed on FB and UK—thanks to the index reduction. The results show that the parallelism together with the index reduction techniques scales up the distance labeling to handle larger graphs.

Exp 3: Index Size. Figure 8 shows the index size of PLL, PSL, PSL⁺, and PSL*. The label size of PLL and PSL is the same, which conforms to the analysis in Sect. 3.3. Both index reduction techniques are effective. PSL⁺ reduces the index size of PSL on SK by more than 50%. Moreover, only PSL* can index massive graphs such as UK while the other approaches ran out of memory. This verified the effectiveness of our index reduction approaches.

Exp 4: Query Time. We compare the average query time of PSL, PSL⁺ and PSL* on 10⁶ random queries. Figure 9 shows that PSL⁺ and PSL* have a query time comparable to PSL. For PSL⁺, the additional query cost on checking equivalence relations is negligible. Since G^s is smaller than G , the query time of PSL⁺ is sometimes smaller than PSL. For example, the query time of PSL⁺ on DELI is 1.17E−6 s while the query time of PSL is 1.31E−6 s. For PSL*, although the labels of nodes in $M(G)$ need to be constructed on-the-fly, the query time of PSL* is within twice the query time of PSL on average, remaining in micro-second level.

Exp 5: Indexing Speedup on Multi-Cores. The speedup of the index time of an approach on x cores is calculated by

$$\text{speedup} = \frac{\text{The index time of the approach with 1 core}}{\text{The index time of the approach with } x \text{ cores}}.$$

According to the above equation, when the core number is 1, the speedup is constantly 1; when an approach fails in

indexing on 1 core within the time limit, its speedup cannot be derived. Figure 10 shows the index time speedup of PSL, PSL⁺ and PSL* with the core number varying from 1, 12, 23, 34, to 45 on six networks, DBLP, POK, LJ, FB, WIKI, and SK, respectively. A near linear speedup has been observed for all the three approaches along with the increasing number of cores. The speedup of each approach is relatively stable over different graphs. On 45 cores, PSL shows, over all datasets, an average speedup of 30 and a maximum speedup of 32, PSL⁺ shows average 28 and maximum 31 while PSL* shows average 27 and maximum 31. The index reduction techniques have little influence on the speedup: the lines of the three approaches clutter, especially on DBLP. A mild slowdown in the speedup when the core number gets close to 45 can be explained by the imbalance resource allocation introduced by more cores. The index size reduction techniques can be critical: PSL failed on FB even when 45 cores were engaged, while PSL* removed redundant nodes to achieve a completion.

Exp 6: Scalability on Index Time. We randomly divided the nodes of a graph into 5 groups; each group consisted of 1/5 of the nodes. We created 5 graphs, while the i -th test case is the induced subgraph on the first i node groups. The experiments were performed on the 5 graphs, respectively.

Figure 11 shows that the index time of PSL* increases almost linearly with the number of nodes of the graph. For example, the index time is about 48 times on 100% nodes than on 20% nodes of DBLP and is about 8 times for FB. For PSL and PSL⁺, although there is a situation where these two methods fail to create the index, the index time increases smoothly when the number of nodes increases. Therefore, the above results justify the scalability of PSL for index time.

Exp 7: Scalability on Index Size. The setting is the same as the former experiment. Figure 12 shows that the space consumption grows smoothly with the graph size for all three methods. For example, the index space on 100% nodes of DBLP is about 184.6, 251.2, 182.5 times larger than that on 20% nodes for PSL, PSL⁺, and PSL*, respectively. Therefore, the smooth increase of the index space shows the scalability of PSL for the index size.

Exp 8: Scalability on Query Time. Figure 13 shows that the query time of the proposed approaches grows smoothly with the graph size. For example, on LJ, the query time on 100% nodes is about 368.34, 372.92, and 546.26 times larger than that on 20% nodes for PSL, PSL⁺, and PSL*, respectively. Other graphs show a similar trend. Combining the above experiments on the scalability test, we draw the conclusion that the proposed methods all show excellent scalability.

7.2 Test on node ordering

Algorithms. For PSL, we use PSL_D to denote PSL whose order is determined by degrees and PSL_B to denote PSL whose

⁷ <http://snap.stanford.edu/data/>.

⁸ <http://law.di.unimi.it>.

⁹ <http://konect.uni-koblenz.de/>.

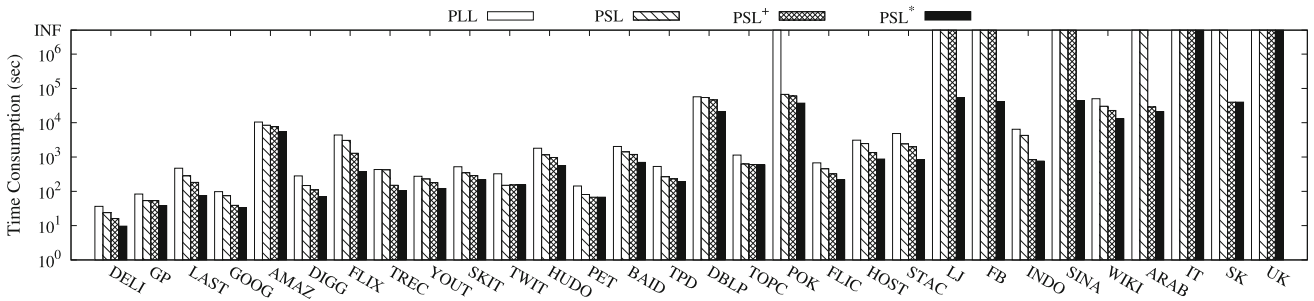


Fig. 6 The comparison of the index time on one core

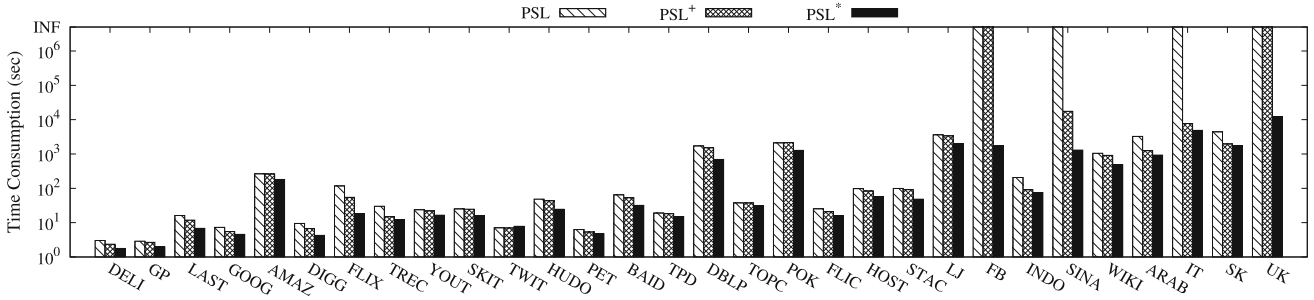


Fig. 7 The comparison of the index time on 45 cores

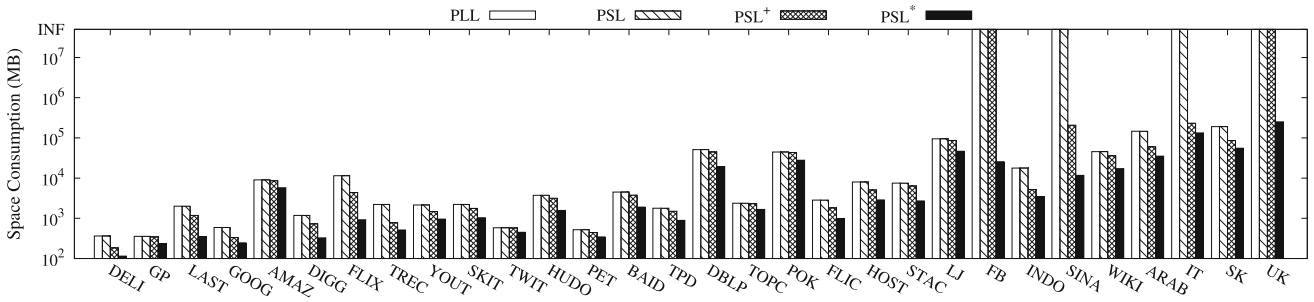


Fig. 8 The comparison of the index size

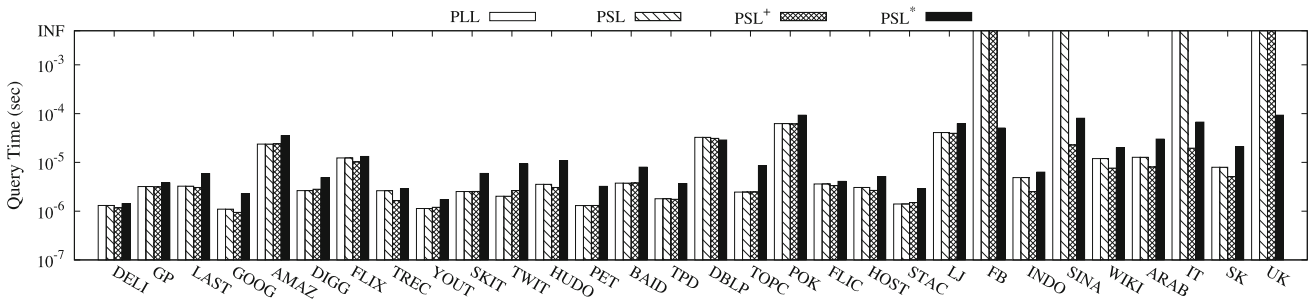


Fig. 9 The comparison of the query time

Table 6 The Description of Added Datasets

Name	Dataset	<i>n</i>	<i>m</i>	Type
UK75	UK-2007-05 ⁸	105,896,555	3,738,733,648	Web Graph
UKOQ	UK-2007 ⁸	133,633,040	5,507,679,822	Web Graph

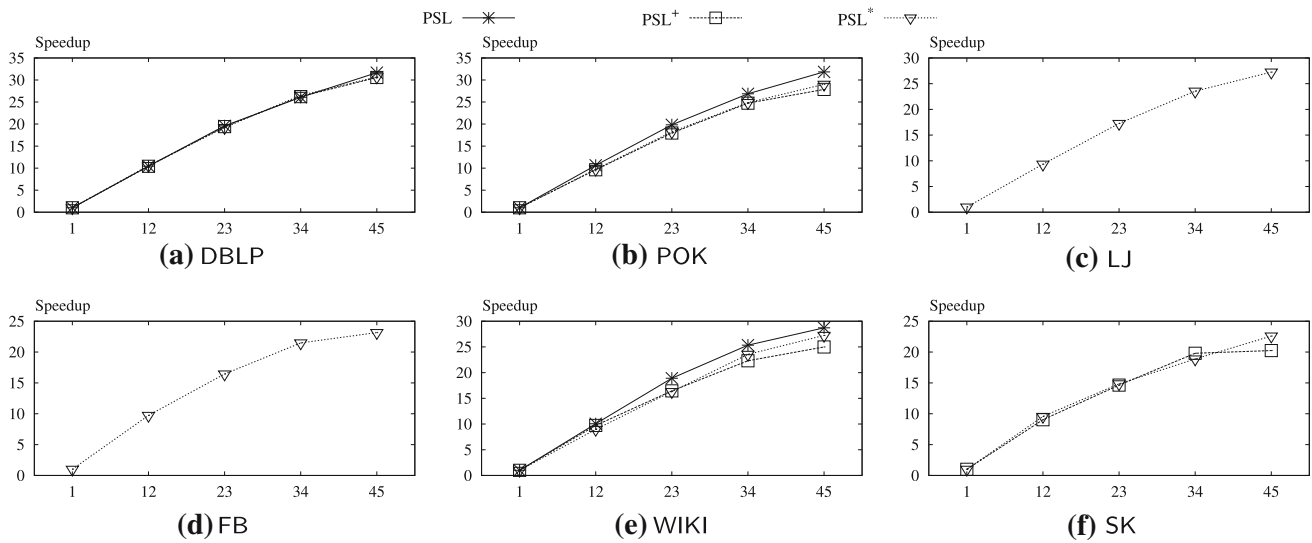


Fig. 10 The effect of core number on the index time

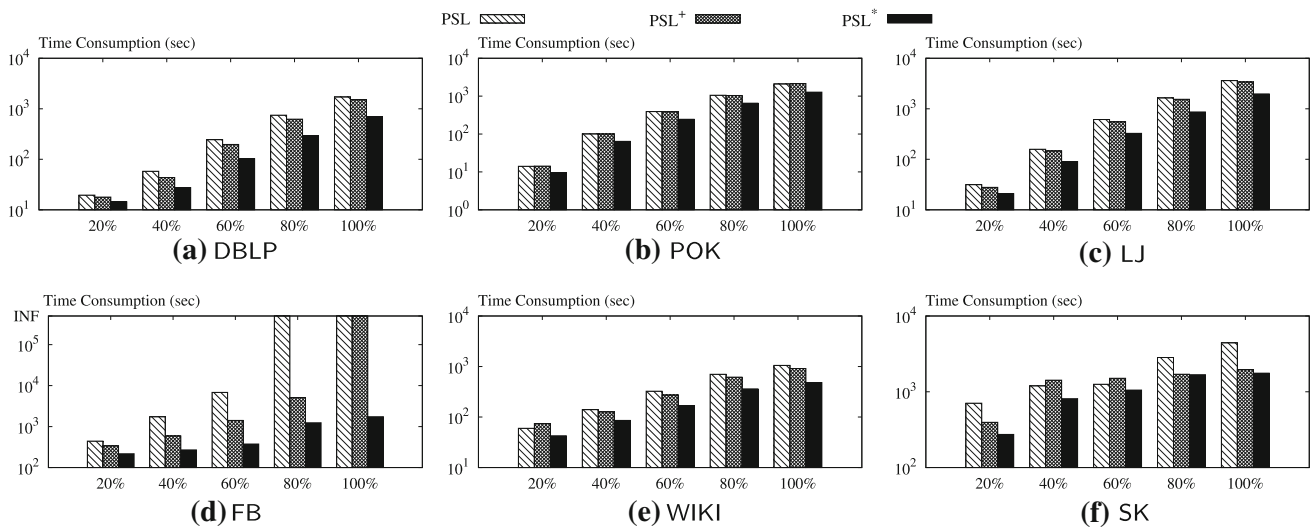


Fig. 11 The test of scalability for the index time

order is determined by k -betweenness. Furthermore, to test the effect of removing local minimum set on computing k -betweenness, we impose different node orders on PSL*, which includes the following three methods:

- PSL*_D: PSL* using degrees to determine node order.
- PSL*_B: PSL* using k -betweenness computed by the pool-based sampling method (Algorithm 5) for ordering.
- PSL*_I: PSL* using k -betweenness computed by the improved sampling method (Algorithm 6) for ordering.

Datasets. Experiments were performed on 30 real-world graphs in Table 5. Moreover, to further test the effect of different ordering methods, we provide two additional datasets,

as shown in Table 6. The largest added graph has more than 5.5 billion edges.

Exp 9: Degree-based and Betweenness-based Node Orders on PSL. We study the effect of node orders (using degree and betweenness) on PSL index sizes. Among them, we obtain the node orders determined by betweenness in two ways: PSL_B whose order is determined by our proposed k -betweenness algorithm, and we set the parameter k to 4; and PSL_C whose order is determined by classical betweenness. We use the method ABRA¹⁰ in [38] to estimate the classical between-

¹⁰ We chose ABRA for two reasons. First, as pointed out in [38], ABRA outperforms the method of [37]. Second, ABRA can be terminated at any time during execution, which leads to a fair comparison with our method. The source code of ABRA is also the code used in the literature [6] and has been implemented in parallel with OpenMP.

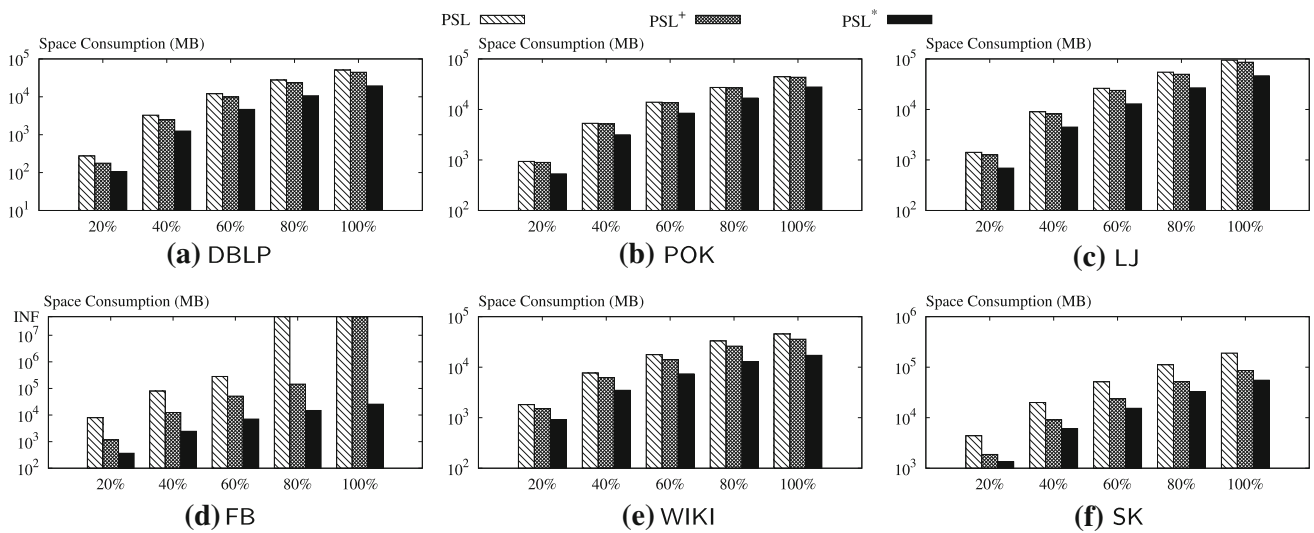


Fig. 12 The test of scalability for the index size

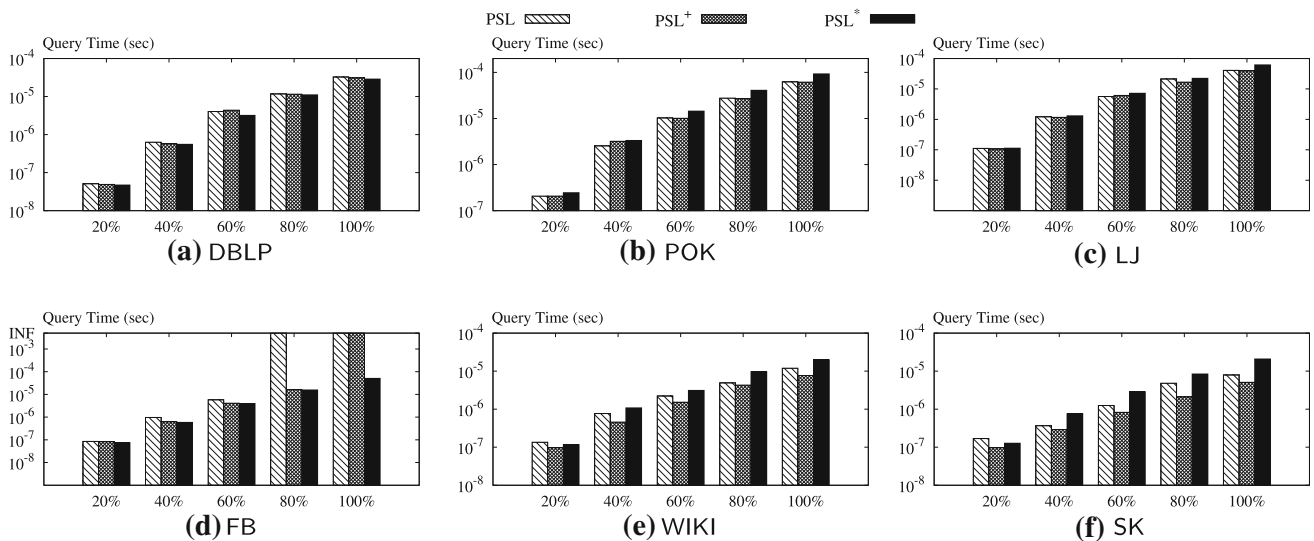


Fig. 13 The test of scalability for the query time

ness values, and its parameters are set according to those in [38]. For PSL_B and PSL_C , we stop sampling when the sampling time exceeds the same time threshold. We compared the index sizes of PSL_D , PSL_C and PSL_B on all graphs where PSL_D can create indexes. The results are shown in Fig. 14.

First, we compare PSL_B with PSL_D to show that using betweenness is superior to using degree as the node order. As can be seen in Fig. 14, the index size of PSL_B is always smaller than that of PSL_D , and the index size of PSL_B can be more than five times smaller than that of PSL_D on ARAB. These results show that setting betweenness to node order is useful for reducing the index size.

Then, we compare PSL_B with PSL_C to illustrate the necessity of the proposed k -betweenness approximation algorithm. In 19 out of 24 graphs, the index size of PSL_B is smaller than

that of PSL_C (by a factor of up to 2.45 on WIKI); on other graphs, the index size of PSL_B is comparable to that of PSL_C . This result illustrates why new betweenness approximation methods need to be designed for distance labeling: replacing classical betweenness with k -betweenness leads to a considerable reduction in the index size of PSL_B compared to PSL_C , especially for large graphs.

Exp 10: Effect of Node Order on the Index Size of PSL^* .

The primary goal of determining the node order using k -betweenness is to reduce the index size—on a multiple-core environment, the failure of labeling methods mainly results from the unaffordable index size. We compared PSL^* using different node ordering methods, where the hop number k is set to 4, and the sampling time T is set to 3600 s. The effect

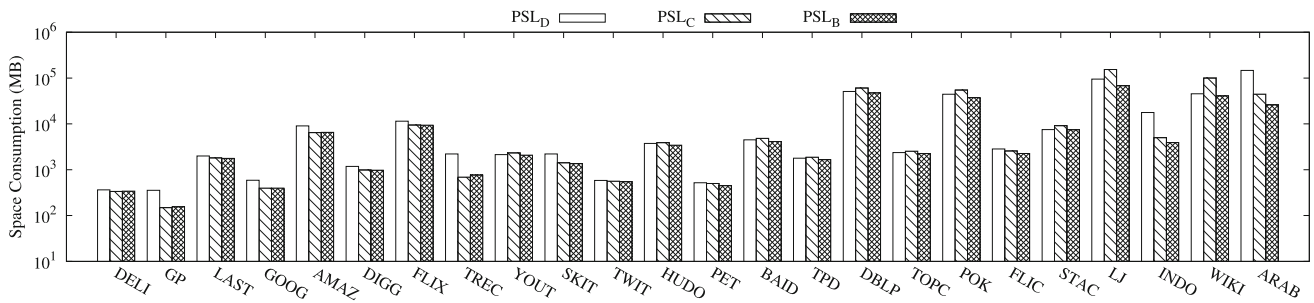


Fig. 14 The comparison of node order degree and betweenness

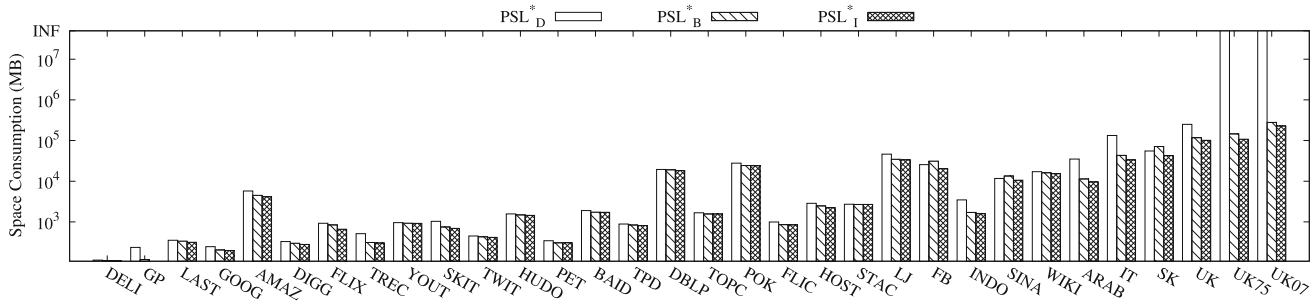


Fig. 15 The effect of the node order on the index size

of parameters T and k on the index size will be discussed later, and the results on all 32 graphs are given in Fig. 15.

Figure 15 indicates that replacing degrees (PSL^*_D) with k -betweenness (PSL^*_B and PSL^*_I) enables the indexing on large graphs UK75 and UK07. This demonstrates the meaning of adopting k -betweenness as a node order. Moreover, on the 30 graphs where PSL^*_D finished labeling, the index size of PSL^*_I is, on average, 1.48 times smaller on average than that of PSL^*_D , and the size of PSL^*_D is reduced by about 4 times at most.

We then verify that it is useful to consider the local minimum set ($M(G)$) elimination in the computation of k -betweenness. As shown in Fig. 15, the index size of PSL^*_B can be sometimes larger than that of PSL^*_D : PSL^*_B 's index size is 1.22 times and 1.3 times that of PSL^*_D on FB and SK, respectively. In contrast, the index of PSL^*_I is always smaller than that of PSL^*_D . Furthermore, the index size of PSL^*_I is, on average, 1.12 times smaller than that of PSL^*_B , and the size of PSL^*_B is reduced by more than 1.67 times at most. These results are encouraging because it shows that taking $M(G)$ into account can effectively reduce the index size under the same sampling time.

Exp 11: Effect of Node Order on Query Time of PSL^* . Figure 16 compares PSL^*_D , PSL^*_B , and PSL^*_I in query time. On average, PSL^*_B takes 0.91 times as long as PSL^*_D , while PSL^*_I takes only 1.04 times as long as PSL^*_D . This means that reducing the size does not affect the query time— PSL^*_B shortens the query time of PSL^*_D , and PSL^*_I 's query time is close to that of PSL^*_D .

Exp 12: Effect of Node Order on Index Time of PSL^* .

We show the index time (including one-hour sampling time for PSL^*_B and PSL^*_I) for different ordering methods, and the results are given in Fig. 17. On all the graphs, the index time of PSL^*_B does not exceed the index time of PSL^*_D by more than 2 hours, while the index time of PSL^*_I does not exceed the index time of PSL^*_D by more than 1.5 hours. Note that the additional overhead in index time is acceptable: on the one hand, we need time to estimate k -betweenness, and on the other hand, adopting k -betweenness as the node order does not significantly improve the index time.

It is also interesting to observe that on some graphs the index time is reduced when we replace PSL^*_D by PSL^*_I : on UK, the index time of PSL^*_D is 11986.17 s while the time is 9599.761 s for PSL^*_I . Furthermore, note that PSL^*_D cannot index on large graphs such as UK75 and UK07 due to the exhaustive index size. These results support the idea of adopting k -betweenness as the node order, provided that index time can be dramatically reduced in a multi-core environment.

Exp 13: Effect of k on Index Size. We examine the effect of hop number k on the index size, where k is the parameter that defines k -betweenness. Since PSL^*_I performs better than PSL^*_B in reducing the index size, we only present the results of PSL^*_I . We varied the number k from 2, 3, 4, 5, to 6, and the results are shown in Fig. 18. Note that the red line in the figures is the index size when k is set as the diameter of the graph.

Figure 18 shows that different graphs have different trends: as k increases, the index size first decreases and then

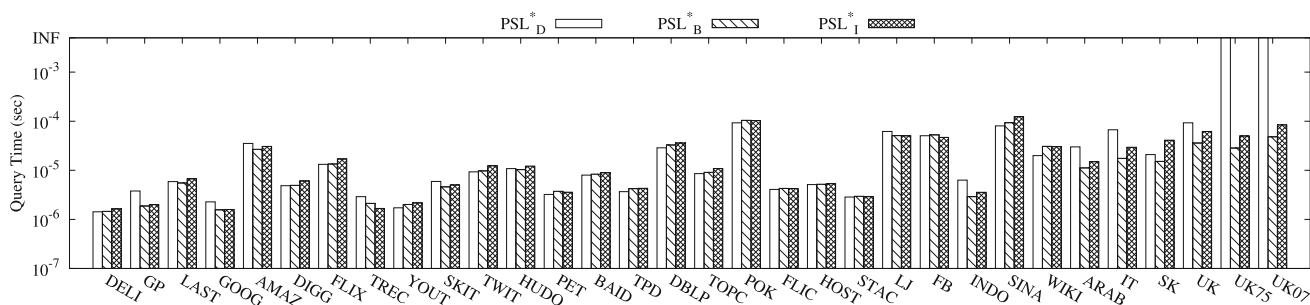


Fig. 16 The effect of the node order on the query time

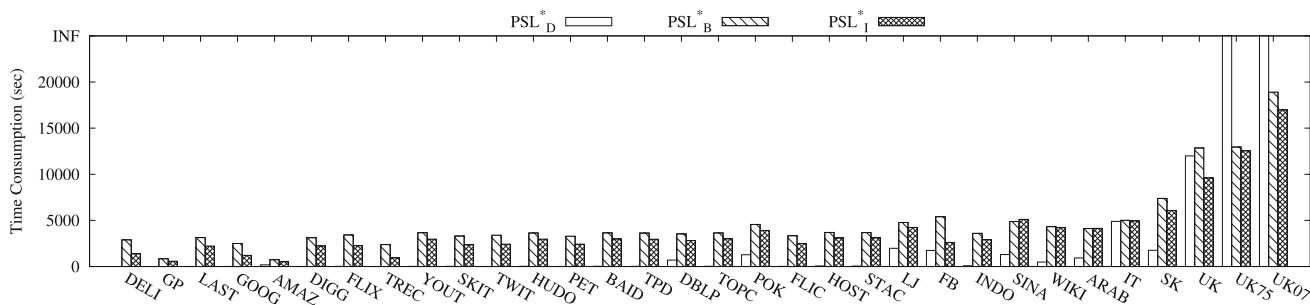


Fig. 17 The effect of the node order on the index time

increases on DBLP, POK, LJ, and WIKI; on FB, the index size decreases continuously; on SK, the index size first increases then decreases. The different trends suggest that we adopt k -betweenness rather than betweenness (when k is infinite) is desirable: given a limited sampling time, a larger k does not imply a smaller index. Furthermore, setting k to 4 allows a reasonably small index size on all graphs, and 4 is the default hop number for k -betweenness.

Exp 14: Effect of k on Query Time. We examine the effect of hop number k on the query time, where all experimental settings are the same as Exp 13. Figure 19 shows that different graphs have different trends: as k increases, the query time first decreases and then increases on POK, LJ, and FB; the query time first increases then decreases on WIKI; the query time fluctuates on DBLP and SK. Also, by comparing with the query time obtained using betweenness (when k is set to infinity), we find that the query time obtained using k -betweenness is comparable. This shows that using k -betweenness as the node order can reduce the index size without sacrificing the query time.

Exp 15: Effect of Sampling Time T . PSL_I adopts a sampling-based algorithm to approximate k -betweenness. Instead of giving the total sample size, PSL_I provides the time limit T for the sampling process. To evaluate the effect of sampling time T on the index size of PSL_I, we changed T from 900, 1800, 3600, 5400, to 7200 s, and the results are given in Fig. 20.

On all the graphs, the index size does not increase as more sampling time is given. This is reasonable, as an increasing

T leads to a more accurate estimation of k -betweenness. Furthermore, for some graphs, such as DBLP and FB, the index size reduce smoothly after one hour, which explains why 3600 s is the default sampling time for PSL_I. However, on large graphs such as LJ and SK, the index size keeps decreasing. This verifies the benefits of our method in handling large graphs when more sampling time is given.

Exp 16: Effect of Sampling Time T on Query Time. We examine the effect of sampling time T on the query time, where all experimental settings are the same as the Exp 15. Figure 21 shows that the difference in query time across all graphs is insignificant when T is changed: despite the different trends in query time on various graphs, the ratio between the maximum and minimum query time on all graphs does not exceed 2.83. This result further highlights that k -betweenness as node order can guarantee good query time while reducing index size.

Exp 17: Overall Index Size Reduction Ratio.

After the reduction of index time using multi-core parallelization, the study of index size reduction becomes important. This paper proposes two ways in reducing the index size of PSL_D: (i) index compression by removing redundant information (e.g., equivalent relationship reduction and local minimum set elimination); and (ii) setting the node order using k -betweenness. The final PSL_I combines the above two reduction techniques. To further highlight the significance of size reduction, we compared the index size between PSL_D and PSL_I. We use the metric *ratio* to show the percentage

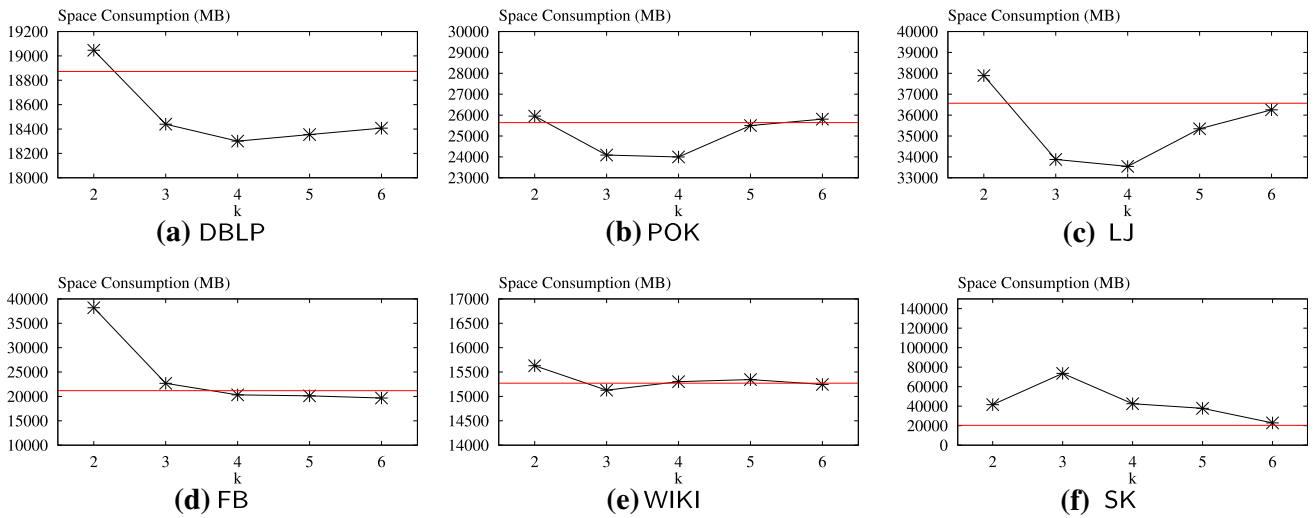


Fig. 18 The effect of the hop number k on the index size (PSL^*_1)

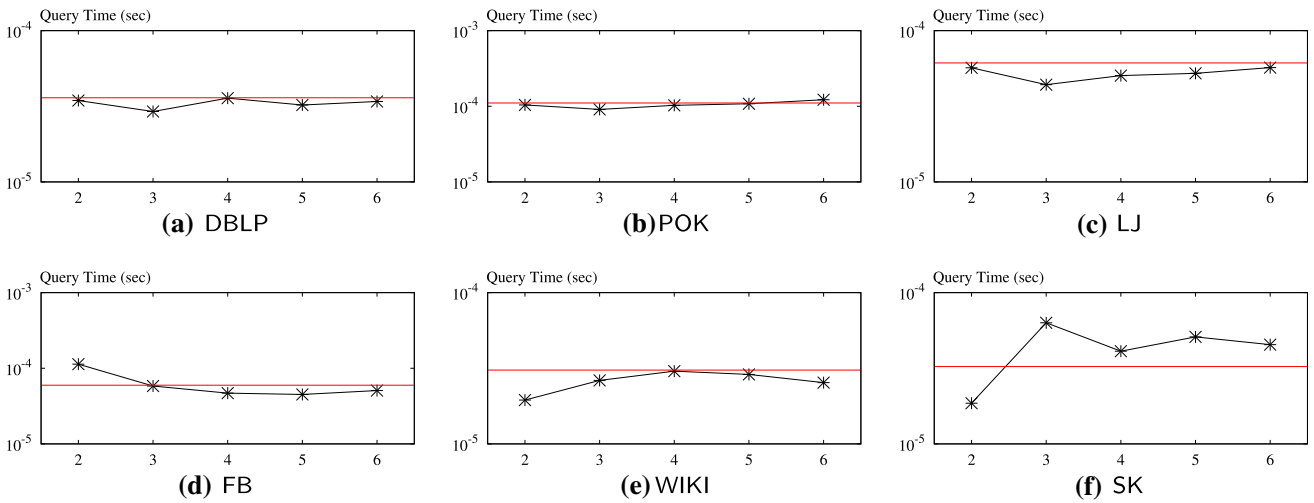


Fig. 19 The effect of the hop number k on the query time

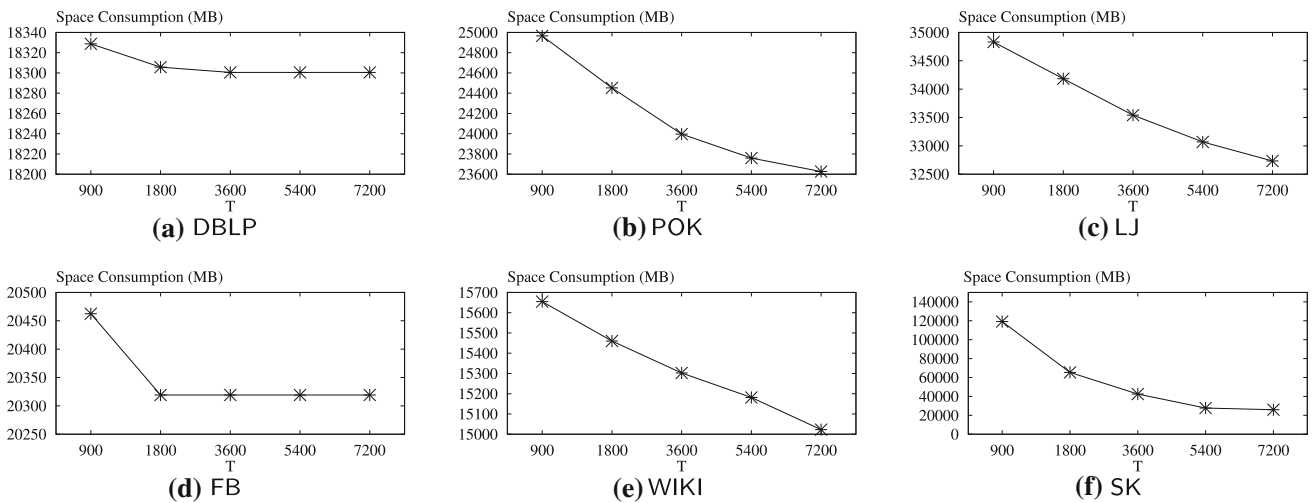


Fig. 20 The effect of the sampling time T on the index size (PSL^*_1)

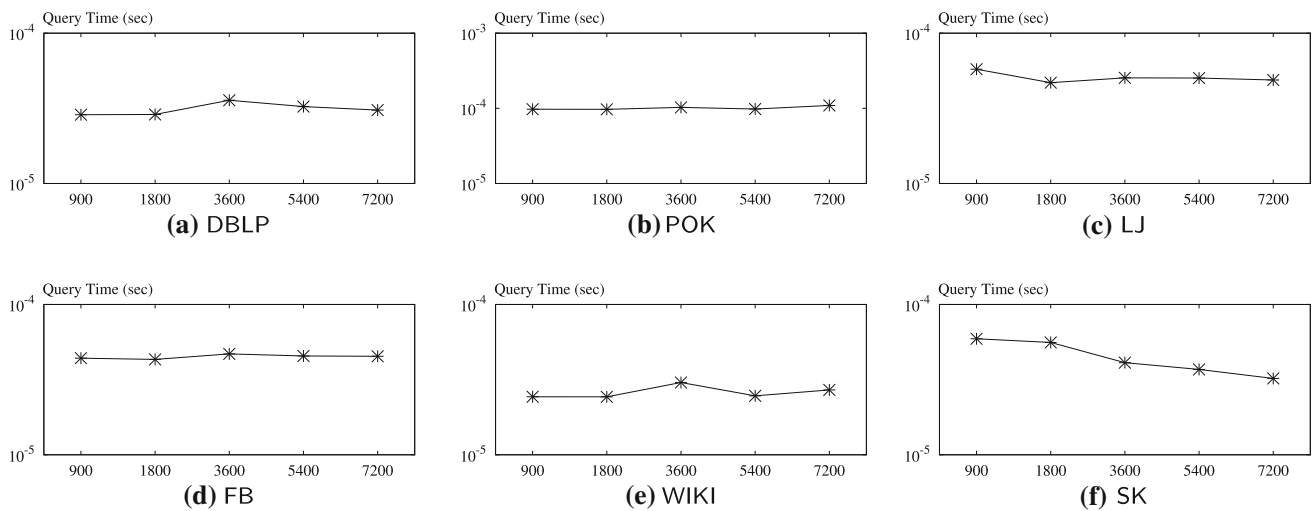


Fig. 21 The effect of sampling time T on the query time

Table 7 Overall index size reduction ratio

Name	PSL _D (MB)	PSL* ₁ (MB)	Ratio %	Name	PSL _D (MB)	PSL* ₁ (MB)	Ratio %
DELI	364.05	110.474	69.65	GP	355.24	107.495	69.74
LAST	1997.52	314.168	84.27	GOOG	589.11	196.834	66.59
AMAZ	9025.18	4150.954	54.01	DIGG	1178.41	278.791	76.34
FLIX	11444.23	657.548	94.25	TREC	2208.96	300.507	86.40
YOUT	2141.51	919.85	57.05	SKIT	2209.91	691.007	68.73
TWIT	582.58	414.14	28.91	HUDD	3738.98	1442.207	61.43
PET	519.32	305.785	41.12	BAID	4493.61	1717.496	61.78
TPD	1783.19	809.296	54.62	DBLP	50996.04	18300.504	64.11
TOPC	2365.64	1557.537	34.16	POK	44414.19	23996.3	45.97
FLIC	2839.96	845.43	70.23	HOST	8005.39	2230.982	72.13
STAC	7495.68	2686.053	64.17	LJ	94950.66	33542.286	64.67
INDO	17731.95	1581.035	91.08	WIKI	45447.10	15302.388	66.33
ARAB	146394.20	9587.213	93.45	SK	190216.16	42533.59	77.64

of index size reduction that PSL*₁ achieves compared to PSL_D,

where $\text{ratio} = 100\% - \frac{\text{Index size of PSL}^*_1}{\text{Index size of PSL}_D} \times 100\%$. Table 7 lists the ratio on all graphs that PSL_D can complete the labeling process.

Table 7 shows that PSL*₁ can compress PSL_D's index size by 94.25% on FLIX—PSL_D's index size is decreased by more than an order of magnitude. Furthermore, PSL*₁ can build the index on large graphs where PSL_D fails, demonstrating the necessity in using index reduction techniques for distance labeling even in a multi-core environment.

8 Conclusions

In this paper, we propose a novel parallelized labeling scheme for distance queries on small-world networks. Our method

accelerates the index construction by concurrently creating labels with the same label distances. Moreover, the index size is reduced by removing redundant nodes from the graph and removing labels of local minimum sets from the index. Scalable approximation algorithms for k -betweenness computation is proposed, so that k -betweenness can be used as a node order to further reduce the index size. Extensive experimental results illustrate the superior efficiency of our approach. In particular, our approach enables the building of the index for networks at billion scales. Experimental results verify the near-linear speedup of our algorithms in a multi-core environment.

Acknowledgements Miao Qiao is supported by Marsden Fund UOA 1732, Royal Society of New Zealand and Catalyst: Strategic Fund 3721519 from Government Funding, Ministry of Business Innovation and Employment. Lu Qin is supported by ARC FT200100787 and DP210101347. Ying Zhang is supported by ARC DP180103096

and FT170100128. Lijun Chang is supported by ARC DP160101513 and FT180100256. Xuemin Lin is supported by NSFC61232006, 2018YFB1003504, ARC DP200101338, DP180103096, and DP170101628.

A Proof of Lemma 1

According to triangle inequality, for any node $u \in V$, $\text{dist}(s, u) + \text{dist}(u, t) \geq \text{dist}(s, t)$. For a node u' on a shortest path from s to t , $\text{dist}(s, t) = \text{dist}(s, u') + \text{dist}(u', t)$. Since $C(s) \cap C(t)$ shares a node with a shortest path from s to t , $\min_{v \in C(s) \cap C(t)} \text{dist}(s, v) + \text{dist}(v, t) = \text{dist}(s, t)$.

B Extend PSL to directed graphs

For directed graphs, each node $v \in V$ is associated with a set of hub nodes $C_{\text{IN}}(v)$, where $w \in C_{\text{IN}}(v)$ can reach v and another set of hub nodes $C_{\text{OUT}}(v)$, where v can reach $w \in C_{\text{OUT}}(v)$. Combined with the distance, we obtain two labels $L_{\text{IN}}(v) = \{(u, \text{dist}(u, v)) | u \in C_{\text{IN}}(v)\}$ and $L_{\text{OUT}}(v) = \{(u, \text{dist}(v, u)) | u \in C_{\text{OUT}}(v)\}$ for the node v . To compute the labels $L_{\text{OUT}}(v)$, we run PSL on G ; to compute $L_{\text{IN}}(v)$, we reverse the edge direction of graph and run PSL on the reversed graph. To process the distance query $q(s, t)$, we make use of $\text{Query}(s, t, L)$ defined in the following equation.

$$\text{Query}(s, t, L) = \min_{u \in C_{\text{OUT}}(s) \cap C_{\text{IN}}(t)} (\text{dist}(s, u) + \text{dist}(u, t)).$$

References

1. Abboud, A., Grandoni, F., Williams, V.V.: Subcubic equivalences between graph centrality problems, APSP and diameter. In: Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms, pp. 1681–1697. SIAM (2014)
2. Abraham, I., Delling, D., Goldberg, A.V., Werneck, R.F.: Hierarchical hub labelings for shortest paths. In: European Symposium on Algorithms, pp. 24–35. Springer (2012)
3. Akiba, T., Iwata, Y., Kawarabayashi, K., Kawata, Y.: Fast shortest-path distance queries on road networks by pruned highway labeling. In: 2014 Proceedings of the Sixteenth Workshop on Algorithm Engineering and Experiments (ALENEX), pp. 147–154. SIAM (2014)
4. Akiba, T., Iwata, Y., Yoshida, Y.: Fast exact shortest-path distance queries on large networks by pruned landmark labeling. In: Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data, pp. 349–360. ACM (2013)
5. Akiba, T., Sommer, C., Kawarabayashi, K.: Shortest-path queries for complex networks: exploiting low tree-width outside the core. In: Proceedings of the 15th International Conference on Extending Database Technology, pp. 144–155. ACM (2012)
6. AlGhamdi, Z., Jamour, F., Skiadopoulos, S., Kalnis, P.: A benchmark for betweenness centrality approximation algorithms on large graphs. In: Proceedings of the 29th International Conference on Scientific and Statistical Database Management, pp. 1–12 (2017)
7. Bader, D.A., Kintali, S., Madduri, K., Mihail, M.: Approximating betweenness centrality. In: International Workshop on Algorithms and Models for the Web-Graph, pp. 124–137. Springer (2007)
8. Boldi, P., Rosa, M., Santini, M., Vigna, S.: Layered label propagation: a multiresolution coordinate-free ordering for compressing social networks. In: Srinivasan, S., Ramamritham, K., Kumar, A., Ravindra, M.P., Bertino, E., Kumar, R. (eds.) Proceedings of the 20th International Conference on World Wide Web, pp. 587–596. ACM Press (2011)
9. Boldi, P., Vigna, S.: The webgraph framework I: compression techniques. In: Proceedings of the Thirteenth International World Wide Web Conference (WWW 2004), pp. 595–601. ACM Press, Manhattan, USA (2004)
10. Borassi, M., Crescenzi, P., Habib, M.: Into the square: on the complexity of some quadratic-time solvable problems. Electron. Notes Theor. Comput. Sci. **322**, 51–67 (2016)
11. Borassi, M., Natale, E.: KADABRA is an adaptive algorithm for betweenness via random approximation. J. Exp. Algorithmics (JEA) **24**(1), 1–35 (2019)
12. Borgatti, S.P., Everett, M.G.: A graph-theoretic perspective on centrality. Soc. Netw. **28**(4), 466–484 (2006)
13. Brandes, U.: A faster algorithm for betweenness centrality. J. Math. Sociol. **25**(2), 163–177 (2001)
14. Brandes, U.: On variants of shortest-path betweenness centrality and their generic computation. Soc. Netw. **30**(2), 136–145 (2008)
15. Chen, W., Sommer, C., Teng, S.-H., Wang, Y.: A compact routing scheme and approximate distance oracle for power-law graphs. ACM Trans. Algorithms (TALG) **9**(1), 4 (2012)
16. Coffman, T., Greenblatt, S., Marcus, S.: Graph-based technologies for intelligence analysis. Commun. ACM **47**(3), 45–47 (2004)
17. Cohen, E., Halperin, E., Kaplan, H., Zwick, U.: Reachability and distance queries via 2-hop labels. In: Proceedings of the Thirteenth Annual ACM-SIAM Symposium on Discrete Algorithms, pp. 937–946. Society for Industrial and Applied Mathematics (2002)
18. Dolev, S., Elovici, Y., Puzis, R.: Routing betweenness centrality. J. ACM (JACM) **57**(4), 1–27 (2010)
19. Freeman, L.C.: A set of measures of centrality based on betweenness. Sociometry **40**, 35–41 (1977)
20. Fu, A.W.C., Wu, H., Cheng, J., Wong, R.C.W.: Is-label: an independent-set based labeling scheme for point-to-point distance querying. Proc. VLDB Endow. **6**(6), 457–468 (2013)
21. Guimera, R., Mossa, S., Turtschi, A., Amaral, L.A.N.: The worldwide air transportation network: anomalous centrality, community structure, and cities global roles. Proc. Natl. Acad. Sci. **102**(22), 7794–7799 (2005)
22. Hayashi, T., Akiba, T., Kawarabayashi, K.: Fully dynamic shortest-path distance query acceleration on massive networks. In: Proceedings of the 25th ACM International Conference on Information and Knowledge Management, pp. 1533–1542. ACM (2016)
23. Hoeffding, W.: Probability inequalities for sums of bounded random variables. In: The Collected Works of Wassily Hoeffding, pp. 409–426. Springer (1994)
24. Jacob, R., Koschützki, D., Lehmann, K.A., Peeters, L., Tenfelde-Podehl, D.: Algorithms for centrality indices. In: Network Analysis, pp. 62–82. Springer (2005)
25. Jeong, H., Mason, S.P., Barabási, A.-L., Oltvai, Z.N.: Lethality and centrality in protein networks. Nature **411**(6833), 41–42 (2001)
26. Jiang, M., Fu, A.W.C., Wong, R.C.W., Xu, Y.: Hop doubling label indexing for point-to-point distance querying on scale-free networks. Proc. VLDB Endow. **7**(12), 1203–1214 (2014)
27. Kunegis, J.: Konect: the koblenz network collection. In: Proceedings of the 22nd International Conference on World Wide Web, pp. 1343–1350. ACM (2013)
28. Jure, L., Andrej, K.: SNAP datasets: Stanford large network dataset collection. <http://snap.stanford.edu/data>, June (2014)

29. Li, J., Wang, X., Deng, K., Yang, X., Sellis, T., Yu, J.X.: Most influential community search over large social networks. In: 2017 IEEE 33rd International Conference on Data Engineering (ICDE), pp. 871–882. IEEE (2017)
30. Li, W., Qiao, M., Qin, L., Zhang, Y., Chang, L., Lin, X.: Scaling distance labeling on small-world networks. In: Proceedings of the 2019 International Conference on Management of Data, pp. 1060–1077 (2019)
31. Li, Y., Leong Hou, U., Yiu, M.L., Kou, N.M., et al.: An experimental study on hub labeling based shortest path algorithms. *Proc. VLDB Endow.* **11**(4), 445–457 (2017)
32. Liljeros, F., Edling, C.R., Amaral, L.A., Stanley, H.E., Åberg, Y.: The web of human sexual contacts. *Nature* **411**(6840), 907–908 (2001)
33. Ouyang, D., Qin, L., Chang, L., Lin, X., Zhang, Y., Zhu, Q.: When hierarchy meets 2-hop-labeling: efficient shortest distance queries on road networks. In: Proceedings of the 2018 International Conference on Management of Data, pp. 709–724. ACM (2018)
34. Pfeffer, J., Carley, K.M.: k-centralities: local approximations of global measures based on shortest paths. In: Proceedings of the 21st International Conference on World Wide Web, pp. 1043–1050 (2012)
35. Potamias, M., Bonchi, F., Castillo, C., Gionis, A.: Fast shortest path distance estimation in large networks. In: Proceedings of the 18th ACM Conference on Information and Knowledge Management, pp. 867–876. ACM (2009)
36. Qiao, M., Cheng, H., Chang, L., Yu, J.X.: Approximate shortest distance computing: a query-dependent local landmark scheme. *IEEE Trans. Knowl. Data Eng.* **26**(1), 55–68 (2014)
37. Riondato, M., Kornaropoulos, E.M.: Fast approximation of betweenness centrality through sampling. *Data Min. Knowl. Disc.* **30**(2), 438–475 (2016)
38. Riondato, M., Upfal, E.: Abra: approximating betweenness centrality in static and dynamic graphs with rademacher averages. *ACM Trans. Knowl. Disc. Data (TKDD)* **12**(5), 1–38 (2018)
39. Rossi, R., Ahmed, N.: The network data repository with interactive graph analytics and visualization. In: Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence (2015)
40. Polina, R., Aris, A., Aristides, G., Nikolaj, T.: Event detection in activity networks. In: Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 1176–1185. ACM (2014)
41. Shen, C.Y., Huang, L.H., Yang, D.N., Shuai, H.H., Lee, W.C., Chen, M.S.: On finding socially tenuous groups for online social networks. In: Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 415–424. ACM (2017)
42. Then, M., Kaufmann, M., Chirigati, F., Hoang-Vu, T.-A., Pham, K., Kemper, A., Neumann, T., Vo, H.T.: The more the merrier: efficient multi-source graph traversal. *Proc. VLDB Endow.* **8**(4), 449–460 (2014)
43. Travers, J., Milgram, S.: The small world problem. *Psychol. Today* **1**(1), 61–67 (1967)
44. Tretyakov, K., Armas-Cervantes, A., García-Bañuelos, L., Vilo, J., Dumas, M.: Fast fully dynamic landmark-based estimation of shortest path distances in very large graphs. In: Proceedings of the 20th ACM International Conference on Information and Knowledge Management, pp. 1785–1794. ACM (2011)
45. Watts, D.J., Strogatz, S.H.: Collective dynamics of ‘small-world’ networks. *Nature* **393**(6684), 440 (1998)
46. Wei, F.: TEDI: efficient shortest path query answering on graphs. In: Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data, pp. 99–110. ACM (2010)
47. Wu, L., Xiao, X., Deng, D., Cong, G., Zhu, A.D., Zhou, S.: Shortest path and distance queries on road networks: an experimental evaluation. *Proc. VLDB Endow.* **5**(5), 406–417 (2012)

Publisher’s Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.