

# Efficient and Exact Local Search for Random Walk Based Top-K Proximity Query in Large Graphs

Yubao Wu, Ruoming Jin, and Xiang Zhang

**Abstract**—Top- $k$  proximity query in large graphs is a fundamental problem with a wide range of applications. Various random walk based measures have been proposed to measure the proximity between different nodes. Although these measures are effective, efficiently computing them on large graphs is a challenging task. In this paper, we develop an efficient and exact local search method, FLoS (Fast Local Search), for top- $k$  proximity query in large graphs. FLoS guarantees the exactness of the solution. Moreover, it can be applied to a variety of commonly used proximity measures. FLoS is based on the *no local optimum* property of proximity measures. We show that many measures have no local optimum. Utilizing this property, we introduce several operations to manipulate transition probabilities and develop tight lower and upper bounds on the proximity values. The lower and upper bounds monotonically converge to the exact proximity value when more nodes are visited. We further extend FLoS to measures having local optimum by utilizing relationship among different measures. We perform comprehensive experiments on real and synthetic large graphs to evaluate the efficiency and effectiveness of the proposed method.

**Index Terms**—Local search, proximity search, random walk, top- $k$  search, nearest neighbors

## 1 INTRODUCTION

GIVEN a large graph and a query node, finding its  $k$ -nearest-neighbor ( $k$ NN) is a primitive operation that has recently attracted intensive research interests [1], [2], [3], [4]. In general, there are two challenges in top- $k$  proximity query. One is to design proximity measures that can effectively capture the similarity between nodes. Another challenge is to develop efficient algorithms to compute the top- $k$  nodes for a given measure.

Designing effective proximity (similarity) measures is a nontrivial task. Random walk based measures have been shown to be effective in many applications. Some examples include discounted/truncated hitting time [1], [5], penalized hitting probability [6], [7], random walk with restart [8], [9], RoundTripRank [10], and absorption probability [11].

Although various proximity measures have been developed, how to efficiently compute them remains a challenging problem. For most random walk based measures, a naive method requires matrix inversion, which is prohibitive for large graphs. Two *global* approaches have been developed. One applies the power iteration method over the entire graph [12], [4], [13]. Another approach precomputes and stores the inversion of a matrix [8], [14], [15]. The precomputing step is usually expensive and needs to be repeated whenever the graph changes.

To improve the efficiency, local methods have been developed [1], [5], [7], [9]. The idea is to visit the nodes near the query node and dynamically expand the search range. Node proximities are estimated based on local information

only. Without using the global information, however, most existing local search methods cannot guarantee to find the exact solution. Moreover, they are usually designed for specific measures and cannot be generalized to other measures.

In this paper, we propose FLoS (Fast Local Search), a simple and unified local search method for efficient and exact top- $k$  proximity query in large graphs. FLoS has the following properties.

- *Exact*: It guarantees to find the exact top- $k$  nodes.
- *Unified*: It is a general method that can be applied to a variety of random walk based proximity measures. Most existing methods are designed for specific measures.
- *Efficient*: It uses a simple local search strategy that needs neither preprocessing nor iterating over the entire graph. Experimental results show that it is orders of magnitude faster than alternatives.

The key idea behind FLoS is that we can develop upper and lower bounds on the proximity of the nodes near the query node. These bounds can be dynamically updated when a larger portion of the graph is explored and will finally converge to the exact proximity value. The top- $k$  nodes can be identified once the differences between their upper and lower bounds are small enough to distinguish them from the remaining nodes.

The theoretical basis of FLoS relies on the *no local optimum* property of proximity measures. That is, given a query node  $q$ , for any node  $i$  ( $i \neq q$ ) in the graph,  $i$  always has a neighbor that is closer to  $q$  than  $i$  is. We show that many measures have no local optimum. This property ensures that the proximity of unvisited nodes is bounded by the maximum proximity (or minimum proximity for

---

• Y. Wu and X. Zhang are with the Department of Electrical Engineering and Computer Science, Case Western Reserve University, Cleveland, OH, 44106. E-mail: {yubao.wu, xiang.zhang}@case.edu  
 • R. Jin is with the Computer Science Department, Kent State University, Kent, OH, 44240. E-mail: jin@cs.kent.edu

some measures) in the boundary of the visited nodes. It can be utilized to find the top- $k$  nodes without exploring the entire graph under the assumption that the exact proximity can be computed based on local information. However, for most measures, the exact proximity cannot be computed without searching the entire graph. To tackle this challenge, we introduce several simple operations to modify transition probabilities, which enable developing upper and lower bounds on the proximity of visited nodes. The developed upper (lower) bounds monotonically decrease (increase) when more nodes are visited. We further study the relationship among different measures and show that FLoS can also be applied to measures having local optimum. Extensive experimental results show that, for a variety of measures, FLoS can dramatically improve the efficiency compared to the state-of-the-art methods.

## 2 RELATED WORK

Various random walk based proximity measures have been proposed recently [16]. Examples include truncated or discounted hitting time [5], [17], [1], penalized hitting probability [6], [7], random walk with restart [8], [1], effective importance (degree normalized random walk with restart) [9], RoundTripRank [10], and absorption probability [11]. The Katz score [18] captures multiple paths between two nodes and is closely related to the random walk based proximity measures [13].

The basic approach for proximity query is to use the power iteration method [12]. An improved iteration method designed for random walk with restart decomposes the proximity into random walk probabilities of different length [4]. It tries to reduce the number of iterations by estimating the proximity based on the information collected so far. The iteration method can also be improved by the prioritized execution of the iterative computation, where the node with the largest residual proximity value is updated first [13]. Another approach precomputes the information needed for proximity estimation during the query process [8], [14], [15]. However, this step is time consuming and becomes infeasible when the graph is large or constantly changing. Graph embedding method embeds nodes into geometric space so that node proximities can be preserved as much as possible [19]. The embedding step is also time-consuming. Moreover, the proximities in the new space are not exactly the same as the ones in the original graph.

Based on the intuition that nodes near the query node tend to have high proximity, local search methods try to visit a small number of nodes to approximate the proximities. Best-first [7] and depth-first [20] search strategies simply extract a fixed number of nodes near the query node. An approximate local search algorithm is proposed for truncated hitting time [5]. The key idea is to develop upper and lower bounds that can be used to approximate the proximities of local nodes. A similar local search algorithm is developed for personalized PageRank and degree normalized personalized PageRank [1]. The push style method is first developed in [21] for random walk with restart, and later improved by [22], [2] for the top- $k$  query problem. Starting from the query node, the push style method propagates the proximity value to the nodes in the neighborhood of

TABLE 1  
Main symbols

Symbols	Definitions
$G(V, E)$	undirected graph $G$ with node set $V$ and edge set $E$
$N_i$	neighbors of node $i$
$w_{i,j}$	weight of edge $(i, j)$
$w_i$	degree of node $i$ , $w_i = \sum_{j \in N_i} w_{i,j}$
$q$	query node
$\mathbf{e}$	$n \times 1$ vector with $e_q = 1$ and $e_i = 0$ if $i \neq q$ , where $n =  V $
$k$	number of returned nodes
$S$	a set of nodes
$\bar{S}$	complement of $S$ : $\bar{S} = V \setminus S$
$\delta S$	boundary of $S$ : $\{i \in S \mid \exists j \in N_i \cap \bar{S}\}$
$\delta \bar{S}$	boundary of $\bar{S}$ : $\{i \in \bar{S} \mid \exists j \in N_i \cap S\}$
$\mathbf{r}$	$n \times 1$ vector, $r_i$ : proximity of node $i$ w.r.t. the query node $q$
$\bar{\mathbf{r}}$	upper bound of $\mathbf{r}$ : $\bar{r}_i \geq r_i, \forall i \in S$
$\underline{\mathbf{r}}$	lower bound of $\mathbf{r}$ : $\underline{r}_i \leq r_i, \forall i \in S$
$p_{i,j}$	transition probability from node $i$ to $j$
$\mathbf{P}$	transition probability matrix: $P_{q,j} = 0$ ; $P_{i,j} = p_{i,j}$ if $i \neq q$
$\mathbf{d}, r_d$	a dummy node $d$ with constant proximity value $r_d$
$c$	decay factor in PHP, DHT, RWR, EI, or RT
$u \rightsquigarrow v$	node $u$ can reach node $v$ in the transition graph
$u \not\rightsquigarrow v$	node $u$ cannot reach node $v$ in the transition graph
$u \rightsquigarrow S$	node $u$ can reach at least one node in $S$
$u \not\rightsquigarrow S$	node $u$ cannot reach any node in $S$

the query node, and obtains approximate proximity values for them. This basic idea has been adapted to compute the top- $k$  nodes for effective importance [9], RoundTripRank [10], and Katz score [23]. Most of the existing local search methods cannot guarantee the exactness. Moreover, all of them are designed for specific proximity measures. It is unclear whether the local search methods can be generalized to other measures.

## 3 NO LOCAL OPTIMUM PROPERTY

In this section, we first introduce the basic concept of no local optimum property of proximity measures and discuss how it can be used to bound the proximity of the unvisited nodes. Then we study whether commonly used measures have no local optimum and discuss the relationship between them. Table 1 lists the main symbols and their definitions. The top- $k$  query problem is defined as

**Definition 1. [Top- $k$  Query Problem]** Suppose that we have an undirected graph  $G(V, E)$ , a query node  $q$  and a number  $k$ . Let  $\mathbf{r}_i$  represent the proximity of node  $i$  with regard to the query node  $q$ . The top- $k$  query problem aims at finding a node set  $K \subseteq V \setminus \{q\}$  such that  $|K| = k$  and  $\mathbf{r}_i \geq \mathbf{r}_j$ , for any node  $i \in K$  and  $j \in V \setminus (K \cup \{q\})$ .

### 3.1 Theoretical Basis

Note that for some measures, such as penalized hitting probability, random walk with restart, effective importance, RoundTripRank, Katz score and absorption probability, the larger the proximity the closer the nodes. In this case, no local optimum means no local maximum. For other measures, such as discounted or truncated hitting time, the smaller the proximity the closer the nodes. In this case, no local optimum means no local minimum.



Fig. 1. An example graph and its transition graph

Given an undirected and edge weighted graph  $G = (V, E)$  and a query node  $q \in V$ , let  $\mathbf{r}$  be the proximity vector with  $\mathbf{r}_i$  representing the proximity of node  $i \in V$  with respect to query node  $q$ .

**Definition 2. [No Local Maximum]** A proximity measure has no local maximum if for any node  $i \neq q$ , there exists a neighbor node  $j$  of  $i$  (i.e.,  $j \in N_i$ ), such that  $\mathbf{r}_j > \mathbf{r}_i$ .

**Definition 3. [No Local Minimum]** A proximity measure has no local minimum if for any node  $i \neq q$ , there exists a neighbor node  $j$  of  $i$  (i.e.,  $j \in N_i$ ), such that  $\mathbf{r}_j < \mathbf{r}_i$ .

We say that a proximity measure has no local optimum if it has no local maximum or minimum. In Section 3.2, we will examine whether the commonly used proximity measures have no local optimum. Unless otherwise mentioned, in the next, we assume that the larger the proximity the closer the nodes, and focus on the no local maximum property. All conclusions can also be applied to the proximity measures with no local minimum.

Let  $S$  be a set of nodes, and  $\bar{S} = V \setminus S$  be the remaining nodes. We use  $\delta S = \{i \in S \mid \exists j \in N_i \cap \bar{S}\}$  to denote the boundary of  $S$ , and  $\delta \bar{S} = \{i \in \bar{S} \mid \exists j \in N_i \cap S\}$  to denote the boundary of  $\bar{S}$ .

Figure 1(a) shows an undirected graph with 8 nodes. Suppose that the node set  $S = \{1, 2, 3, 4\}$ , then we have  $\bar{S} = \{5, 6, 7, 8\}$ ,  $\delta S = \{3, 4\}$ , and  $\delta \bar{S} = \{5, 6, 7\}$ .

**Theorem 1.** Let  $S$  be a node set containing the query node, and  $u$  be the node with the largest proximity in  $\delta S$ . If a proximity has no local maximum, we have that  $\mathbf{r}_u > \mathbf{r}_j$  ( $\forall j \in \bar{S}$ ).

*Proof:* Suppose otherwise. We have that  $\exists j \in \bar{S}$ , such that  $\mathbf{r}_u \leq \mathbf{r}_j$ . Now suppose that node  $v$  is the node with the largest proximity in  $\bar{S}$ . We have that  $\forall i \in \bar{S} \cup \delta S$ ,  $\mathbf{r}_v \geq \mathbf{r}_i$ . The neighbors of node  $v$  must exist in  $\bar{S} \cup \delta S$ , i.e.,  $N_v \subseteq \bar{S} \cup \delta S$ . Therefore, we have  $\mathbf{r}_v \geq \mathbf{r}_i$  ( $\forall i \in N_v$ ), which means node  $v$  is a local maximum. This contradicts the assumption.  $\square$

Based on Theorem 1, assuming that we already have the exact proximity vector  $\mathbf{r}$ , we can design a simple local search strategy as shown in Algorithm 1 to find the top- $k$  nodes. It starts from the query node  $q$  and uses  $K$  and  $S$  to store the top- $k$  nodes and visited nodes respectively. In each iteration, the algorithm finds the node  $u$  that has the largest proximity in  $S \setminus K$  and expands  $S$  to the neighbors of node  $u$ . Since  $\delta S \subseteq S \setminus K$ , the maximum proximity value in  $S \setminus K$  must be no less than the maximum proximity value in  $\delta S$ , and greater than the maximum proximity value in the unvisited nodes  $\bar{S}$  based on Theorem 1. Thus,  $K$  contains the top- $k$  nodes. The algorithm continues until  $|K| = k + 1$ .

Let  $h$  be the average number of neighbors of a node. In each iteration  $t$ , on average  $h$  nodes are added to  $S$ . This takes  $O(h \log ht)$  time for a sorted list. The overall complexity of Algorithm 1 is thus  $O(\sum_{t=1}^k h \log ht) = O(hk \log hk)$ .

#### Algorithm 1 The basic top- $k$ local search algorithm

**Input:**  $G(V, E)$ , query node  $q$ , proximity vector  $\mathbf{r}$ , number  $k$   
**Output:** Top- $k$  node set  $K$

```

1:  $S \leftarrow \{q\}; K \leftarrow \{\}$ ;
2: while  $|K| < k + 1$  do
3:    $u \leftarrow \operatorname{argmax}_{i \in S \setminus K} \mathbf{r}_i$ ;
4:    $K \leftarrow K \cup \{u\}; S \leftarrow S \cup N_u$ ;
5: return  $K \setminus \{q\}$ ;

```

### 3.2 Measures With and Without Local Optimum

Table 2 summarizes whether the commonly used proximity measures have no local optimum property. Next, we use penalized hitting probability (PHP) [6], [7] as an example to illustrate that it has no local maximum. We use  $w_i$  to denote the degree of node  $i$ , and  $w_{i,j}$  to denote the edge weight between  $i$  and  $j$ . The transition probability from  $i$  to  $j$  is thus  $p_{i,j} = w_{i,j} / w_i$ .

Suppose the undirected graph in Figure 1(a) has unit edge weight. Node 3 has degree 3, thus its transition probability to node 4 is  $p_{3,4} = 1/3$ . Based on these transition probabilities, we can construct the corresponding transition graph as shown in Figure 1(b). In the transition graph, each directed edge and the number on the edge represent the transition probability from one node to the other.

PHP penalizes the random walk for each additional step. Given a query node  $q$ , let  $\mathbf{r}$  denote the PHP proximity vector, with  $\mathbf{r}_i$  representing the proximity value of node  $i$ . PHP can be defined recursively as

$$\mathbf{r}_i = \begin{cases} 1, & \text{if } i = q, \\ c \sum_{j \in N_i} p_{i,j} \mathbf{r}_j, & \text{if } i \neq q, \end{cases}$$

where  $c$  ( $0 < c < 1$ ) is the decay factor in the random walk process. In [6],  $c = e^{-1}$  is used as the decay factor. The query node  $q$  has constant proximity value 1, and there is no transition probability going out of the query node. For example, there is no outgoing edges from the query node 1 in the transition graph in Figure 1(b).

Let  $\mathbf{P}$  be the transition probability matrix with

$$\mathbf{P}_{i,j} = \begin{cases} 0, & \text{if } i = q, \\ p_{i,j}, & \text{if } i \neq q. \end{cases}$$

Then the above recursive definition can be expressed as the following matrix form

$$\mathbf{r} = c\mathbf{P}\mathbf{r} + \mathbf{e},$$

where  $\mathbf{e}_i = 1$  if  $i = q$ , and  $\mathbf{e}_i = 0$  if  $i \neq q$ .

**Lemma 1.** PHP has no local maximum.

*Proof:* Suppose that node  $i$  is a local maximum. We have  $\mathbf{r}_i = c \sum_{j \in N_i} p_{i,j} \mathbf{r}_j \leq c \sum_{j \in N_i} p_{i,j} \mathbf{r}_i = c\mathbf{r}_i < \mathbf{r}_i$ . We get a contradiction that  $\mathbf{r}_i < \mathbf{r}_i$ .  $\square$

The proofs for whether other proximity measures in Table 2 have local optimum can be found in the Appendices. In particular, in Appendix A, we show that EI, DHT, and THT have no local optimum. In Appendix G, we show that RWR, RT, KZ, and AP have local maximum.

Some proximity measures have inherent relationship. The following theorem says that PHP, EI, and DHT are equivalent in terms of ranking.

**Theorem 2.** PHP, EI, and DHT give the same ranking results.

Please see Appendix A for the proof.

TABLE 2  
No local optimum property of some measures

Proximity measures	Abbr.	Ref.	Property
Penalized hitting probability	PHP	[6]	No local maximum
Effective importance	EI	[9]	No local maximum
Discounted hitting time	DHT	[1]	No local minimum
Truncated hitting time	THT	[5]	No local minimum (within L hops)
Random walk with restart	RWR	[8]	Local maximum
RoundTripRank	RT	[10]	Local maximum
Katz score	KZ	[18]	Local maximum
Absorption probability	AP	[11]	Local maximum

From the discussion in this section, we know that if a proximity measure has no local optimum, we can apply local search described in Algorithm 1 to find the top- $k$  nodes under the assumption that the proximity values of all the nodes are given. However, without exploring the entire graph, the exact values of all the nodes are unknown. To tackle this challenge, we can develop lower and upper bounds on the proximity values of visited nodes. When more nodes are visited, the lower and upper bounds become tighter and eventually converge to the exact proximity value.

Next, we will use PHP, which has no local optimum, as a concrete example to explain how to derive the lower and upper bounds. The strategy can be applied to other proximity measures with no local optimum. How to apply the local search strategy to the measures having local optimum will be discussed in Section 6.

## 4 BOUNDING THE PROXIMITY

To develop the lower and upper bounds, we introduce three basic operations, i.e., deletion, restoration, and destination change of transition probability. We show how proximities change if we modify transition probabilities according to these operations. Then we discuss how to derive lower and upper bounds based on them.

### 4.1 Modifying Transition Probability

We first introduce the operation of deleting a transition probability. Figure 2(a) shows an example, in which the original graph is on the top, with node 1 being the query node and transition probabilities  $p_{2,1} = p_{2,3} = 0.5$  and  $p_{3,2} = 1$ . After deleting  $p_{2,3}$ , the resulting graph is shown at the bottom. Note that deleting a transition probability is different from deleting an edge. Deleting an edge will change the transition probabilities on the remaining graph, while deleting a transition probability will not.

**Theorem 3.** Deleting a transition probability will not increase the proximity of any node.

Please see Appendix B for the proof.

Continue with the example in Figure 2(a). Suppose that the decay factor  $c = 0.5$ . The original PHP proximity vector is  $\mathbf{r} = [1, 2/7, 1/7]$ . After deleting  $p_{2,3}$ , the new proximity vector is  $\mathbf{r}' = [1, 1/4, 1/8]$ .

It can be shown in a similar way that if we restore the transition probability as shown in Figure 2(a), the proximities will not decrease.

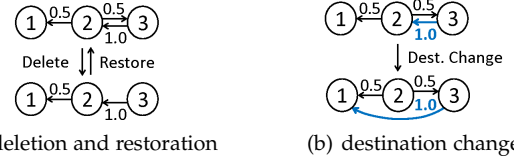


Fig. 2. Basic operations on transition probability

**Theorem 4.** Restoring a deleted transition probability will not decrease the proximity of any node.

*Proof:* Omitted.  $\square$

Figure 2(b) shows an example in which we change the destination of the original transition probability  $p_{3,2}$  from node 2 to 1. Thus,  $p_{3,1}$  is set to  $p_{3,2}$ , and  $p_{3,2}$  is set to 0.

**Theorem 5.** Changing the destination of transition probability  $p_{i,j}$  from node  $j$  to node  $u$  with  $\mathbf{r}_u \geq \mathbf{r}_j$  ( $\mathbf{r}_u \leq \mathbf{r}_j$ ) will not decrease (increase) the proximity of any node.

Please see Appendix B for the proof.

Let us continue with the example in Figure 2(b), where node 1 is the query node. After we change the destination of  $p_{3,2}$  from node 2 to 1, the proximity values should be non-decreasing. With a decay factor  $c = 0.5$ , the proximity vectors before and after the destination change are  $\mathbf{r} = [1, 2/7, 1/7]$  and  $\mathbf{r}' = [1, 3/8, 1/2]$  respectively.

The proofs for other measures having no local optimum are similar and omitted here.

### 4.2 Lower Bound

Let  $S$ ,  $\bar{S}$ ,  $\delta S$  and  $\delta \bar{S}$  represent the set of visited nodes, the set of unvisited nodes, the boundary of  $S$ , and the boundary of  $\bar{S}$ , respectively. From Theorem 3, if we delete all transition probabilities  $\{p_{i,j} : i \text{ or } j \in \bar{S}\}$  in the original graph, the proximity value of any node  $u$  computed using the resulting graph,  $\mathbf{r}_u$ , will be less than or equal to its original value  $\mathbf{r}_u$ , i.e.,  $\mathbf{r}_u \leq \mathbf{r}_u$ . Thus,  $\mathbf{r}$  can be used as the lower bound of  $\mathbf{r}$ .

Let us take the undirected graph in Figure 1(a) for example. Its transition graph is shown in Figure 3(a), with node 1 being the query node and transition probabilities shown on the edges. Suppose that the current set of visited nodes  $S = \{1, 2, 3, 4\}$ . Thus  $\bar{S} = \{5, 6, 7, 8\}$ ,  $\delta S = \{3, 4\}$ , and  $\delta \bar{S} = \{5, 6, 7\}$ . The nodes in  $S$  but not in  $\delta S$  are black. The nodes in  $\delta S$  are gray. The nodes in  $\bar{S}$  are white.

Figure 3(b) shows the resulting transition graph after deleting all transition probabilities  $\{p_{i,j} : i \text{ or } j \in \bar{S}\}$ . The proximity values  $\mathbf{r}$  computed based on Figure 3(b) will lower bound the original proximity values  $\mathbf{r}$  for the nodes in  $S$ .

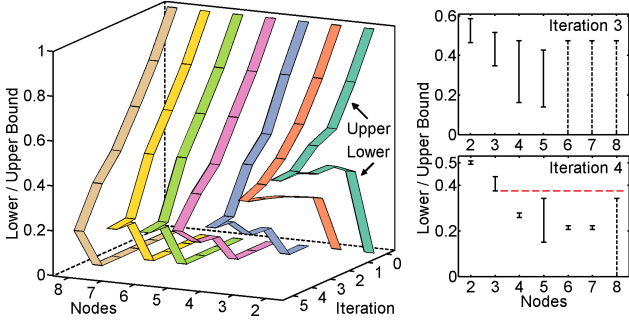
### 4.3 Upper Bound

From Theorem 5, if we change the destination of the transition probabilities  $\{p_{i,j} : i \in \delta S, j \in \delta \bar{S}\}$  to a newly added dummy node  $d$  with a constant proximity value  $\mathbf{r}_d$  and  $\mathbf{r}_d > \mathbf{r}_v$  ( $\forall v \in \bar{S}$ ), the proximity value of any node  $u$  computed using the resulting graph,  $\bar{\mathbf{r}}_u$ , will be greater than or equal to its original value  $\mathbf{r}_u$ , i.e.,  $\bar{\mathbf{r}}_u \geq \mathbf{r}_u$ . Thus,  $\bar{\mathbf{r}}$  can be used as the upper bound of  $\mathbf{r}$ .

Continue with the example in Figure 3(a). In the original graph,  $\delta S = \{3, 4\}$ ,  $\delta \bar{S} = \{5, 6, 7\}$ ,  $p_{3,5} = 1/3$ ,  $p_{4,6} = p_{4,7} = 1/4$ . The left figure in Figure 3(c) shows the resulting graph after we change all the transition probabilities going from  $\delta S$  to





Fig. 4. Lower and upper bounds in different iterations (PHP:  $q=1, c=0.8$ )TABLE 3  
Newly visited nodes in each iteration

Iteration	1	2	3	4	5
Newly visited nodes	{2, 3}	{4}	{5}	{6, 7}	{8}

is close to the exact solution, the algorithm will converge very fast. In our method, between two adjacent iterations, the proximity values of the visited nodes are very close. Therefore, updating the proximity is very efficient.

Algorithm 6 shows the termination criterion. We select the  $k$  nodes in  $S \setminus (\delta S \cup \{q\})$  with largest  $\underline{r}$  values. If the minimum lower bound of the selected nodes is greater than or equal to the maximum upper bound of the remaining visited nodes, the selected nodes will be the top- $k$  nodes in the entire graph. This is because the maximum proximity of unvisited nodes is bounded by the maximum proximity in  $\delta S$ , which is in turn bounded by the maximum upper bound in  $\delta S$ .

Figure 4 shows the lower and upper bounds at different iterations using the example graph in Figure 1(a). One iteration represents one local expansion process. The newly visited nodes in each iteration are listed in Table 3.

The left figure in Figure 4 shows how the lower and upper bounds change through local expansions. Query node 1 has constant proximity value 1.0 thus is not shown. It can be seen that the bounds monotonically change and eventually converge to the exact proximity value when all the nodes are visited. The monotonicity of the bounds is proved theoretically in next two sections.

The right figure in Figure 4 shows the lower and upper bounds in iteration 3 (at the top) and 4 (at the bottom). The interval from the lower to upper bounds is indicated by the vertical line segment. The interval of the bounds for the unvisited node is indicated by the dashed vertical line. In iteration 3, nodes {6, 7, 8} are unvisited, and their upper bound is the upper bound for node 4, which is the maximum upper bound for the boundary nodes {4, 5}. In iteration 4, the bounds become tighter, and the minimum lower bound of nodes {2, 3} is larger than the maximum upper bound of the remaining nodes {4, 5, 6, 7, 8}, which is indicated by the horizontal red dashed line. Therefore, nodes {2, 3} are guaranteed to be the top-2 nodes after iteration 4, even though node 8 is still unvisited.

## 5.2 Monotonicity of the Lower Bound

We first consider the monotonicity of the lower bound. Let  $S^{t-1}$  and  $S^t$  represent the set of visited nodes in iterations  $(t-1)$  and  $t$  respectively. In the next, we prove that the

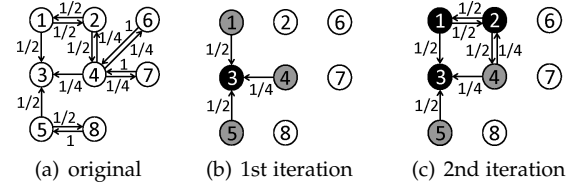


Fig. 5. Example transition graphs between two adjacent iterations for analyzing lower bound monotonicity with query node 3.

lower bound is monotonically non-decreasing when more nodes are visited, i.e.,  $\underline{r}_i^t \geq \underline{r}_i^{t-1}$  ( $i \in S^{t-1}$ ).

Given a directed transition graph, we say that a node  $u$  can reach a node  $v$  if there exists a sequence of adjacent nodes (i.e., a path) which starts from  $u$  and ends at  $v$ . For example, in the transition graph in Figure 5(a), node 1 can reach node 6, but node 5 cannot reach node 6. We use  $u \rightsquigarrow v$  to denote that node  $u$  can reach  $v$ , and  $u \nrightarrow v$  to denote that node  $u$  cannot reach  $v$ . We also use  $u \rightsquigarrow S$  to denote that node  $u$  can reach at least one node in  $S$ , and  $u \nrightarrow S$  to denote that node  $u$  cannot reach any node in  $S$ .

From iteration  $(t-1)$  to  $t$ , we only restore some transition probabilities in  $\{p_{i,j} : i \text{ or } j \in S^t \setminus S^{t-1}\}$ . The following Theorem 6 says that if node  $i$  can reach at least one of the newly added nodes, the lower bound of node  $i$  is strictly increasing. If node  $i$  cannot reach any of the newly added nodes, the lower bound value of node  $i$  will not change during the iteration.

Figure 5 shows an example. Figure 5(a) shows the full transition graph when the query is node 3. Figures 5(b) and 5(c) show the transition graphs constructed in the first and second iteration of Algorithm 2 respectively.  $S^1 = \{3, 1, 4, 5\}$ ,  $S^2 = \{3, 1, 4, 5, 2\}$ , and node 2 is the newly visited node in the second iteration. We can see that node 5 cannot reach node 2. Thus, the lower bound value  $\underline{r}_5$  is unchanged. The lower bound values of nodes {1, 4} are strictly increasing.

**Theorem 6. (Monotonicity of the lower bound)** For any node  $i \in S^{t-1}$ , we have that

$$\begin{cases} \underline{r}_i^t = \underline{r}_i^{t-1}, & \text{if } i \nrightarrow S^t \setminus S^{t-1}, \\ \underline{r}_i^t > \underline{r}_i^{t-1}, & \text{if } i \rightsquigarrow S^t \setminus S^{t-1}, \end{cases}$$

where  $\rightsquigarrow$  and  $\nrightarrow$  represent the reachability in the transition graph at the  $t$ -th iteration.

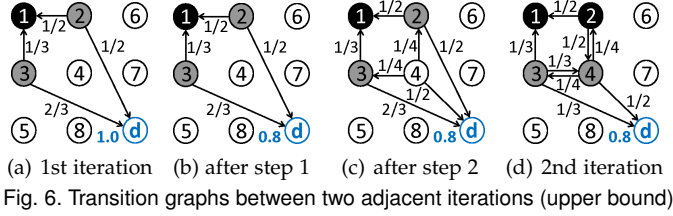
Please see Appendix C for the proof.

## 5.3 Monotonicity of the Upper Bound

In this subsection, we analyze the monotonicity of the upper bound. Specifically, we prove that the upper bound values are strictly increasing until they converge to the exact proximity values. That is, for any node  $i \in S^{t-1}$ ,  $\bar{r}_i^t < \bar{r}_i^{t-1}$  until  $\bar{r}_i^{t-1} = \underline{r}_i$ .

From iteration  $(t-1)$  to  $t$ , we decrease the proximity value of the dummy node and add new nodes in  $S^t \setminus S^{t-1}$ . After adding the new nodes, the transition probabilities need to be updated accordingly. Specifically, we need to

- 1) Decrease the proximity value of the dummy node from  $\bar{r}_d^{t-1}$  to  $\bar{r}_d^t$ ;
- 2) Add the transition probabilities  $\{p_{i,j}\}$  from the newly added nodes  $i (\in S^t \setminus S^{t-1})$  to nodes  $j (\in S^t)$ , and  $\{p_{i,d}\}$  from  $i$  to the dummy node  $d$ ;



- 3) Add the transition probabilities  $\{p_{j,i}\}$  from nodes  $j(\in \delta S^{t-1})$  to the newly added nodes  $i(\in S^t \setminus S^{t-1})$ , and remove their correspondences in  $\{p_{j,d}\}$ .

An example is shown in Figure 6. Figure 6(a) shows the transition graph for the first iteration. Figures 6(b), 6(c) and 6(d) show the resulting graphs after applying steps 1, 2, and 3 respectively. The graph in Figure 6(d) is the final transition graph for the next iteration.

The upper bound values will monotonically change at each step. In step 1, reducing the value of  $r_d$  will not increase the upper bound values. Applying step 2 will not change the upper bound values for the nodes in  $S^{t-1}$ , since all the newly added transition probabilities begin from nodes  $i(\in S^t \setminus S^{t-1})$ . In step 3, we resets the transition probabilities from nodes  $j(\in \delta S^{t-1})$  to the newly added nodes  $i(\in S^t \setminus S^{t-1})$ . This is equivalent to destination change, i.e., changing  $\{p_{j,d}\}$  to  $\{p_{j,i}\}$ . Moreover, we have that  $\bar{r}_d^t \geq \bar{r}_i^t$ . Thus, in step 3, the upper bound values will not increase. We provide rigorous analysis for the three steps in Appendix D.

**Theorem 7. (Monotonicity of the upper bound)** For any node  $i \in S^{t-1}$ , we have that

$$\begin{cases} \bar{r}_i^{t-1} = \bar{r}_i^t, & \text{if } i \rightsquigarrow d, \\ \bar{r}_i^{t-1} > \bar{r}_i^t, & \text{if } i \rightsquigarrow d, \end{cases}$$

where  $\rightsquigarrow$  and  $\rightsquigarrow$  represent the reachability in the transition graph at the  $t$ -th iteration.

Please see Appendix D for the proof.

If we have that  $i \rightsquigarrow d$  in the transition graph at the  $t$ -th iteration, we have that  $i \rightsquigarrow d$  in the transition graph at any future iteration. Thus,  $\bar{r}_i^t$  will not change during the future iterations. Since  $\bar{r}_i^t$  converges to the exact proximity value  $r_i$  when the entire graph is visited. We must have that  $\bar{r}_i^t = r_i$  when  $i \rightsquigarrow d$ . In conclusion, the upper bound strictly decreases until it converges to the exact proximity value.

The lower and upper bounds can be further tightened by adding self-loop transition probabilities to the nodes in  $\delta S$ . Please see Appendix E for more details.

#### 5.4 Complexity

Assume Algorithm 2 executes in  $\beta$  iterations. Let  $h$  be the average number of neighbors of a node. The LocalExpansion step takes  $O(ht)$  time to find the node to expand at the  $t$ -th iteration. To update the lower bound, updating  $\mathbf{P}$  needs  $O(h^2)$  operations, and updating  $\mathbf{r}$  and  $\mathbf{e}$  needs  $O(h)$  operations. Subgraph induced by  $S$  has  $O(h^2t)$  edges, so matrix  $\mathbf{P}$  has  $O(h^2t)$  non-zero entries. Therefore using the iterative method to solve linear equations takes  $O(\alpha h^2t)$  time, where  $\alpha$  is the number of iterations used in IterativeMethod. Thus the overall complexity of UpdateLowerBound in the  $t$ -th iteration is  $O(\alpha h^2t)$ . The complexity of UpdateUpperBound function is the same as that of UpdateLowerBound. In the CheckTermination-Criterion step, finding the nodes with largest lower bounds

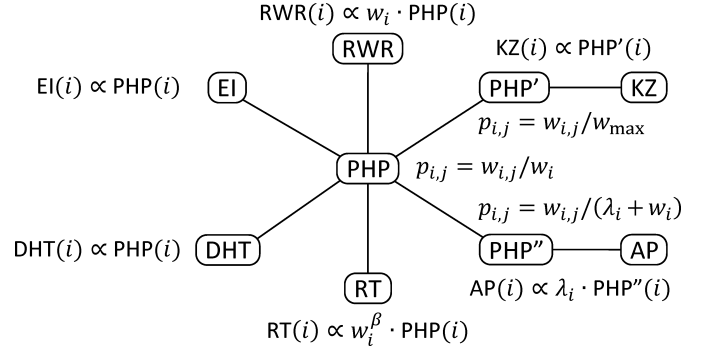


Fig. 7. Relationships between PHP and other proximity measures

takes  $O(ht)$  time. Therefore, the overall complexity of FLoS is  $O(\sum_{t=1}^{\beta} (\alpha h^2t + ht)) = O(\alpha h^2\beta^2)$ .

At each iteration, FLoS visits  $h$  new nodes on average. In the worst case, where the whole graph is visited, FLoS needs to run  $\beta = n/h$  iterations. Thus, the worst case complexity of FLoS is  $O(\alpha h^2\beta^2) = O(\alpha n^2)$ .

In the above complexity analysis, the number of iterations  $\beta$  is proportional to the number of visited nodes. Appendix F provides theoretical analysis of the number of visited nodes. In Section 8, we show experimental results on the number of visited nodes using real graphs.

Note that so far, we have used PHP to illustrate the key principles underlying the fast local search method. EI and DHT are equivalent with PHP thus there is no need to develop algorithm for them. For THT, deleting a transition probability will not increase the proximity of any node. Therefore, when we delete all the transition probabilities  $\{p_{i,j} : i \text{ or } j \in S\}$  in the original transition graph, the proximity value of any node computed based on the modified transition graph will be the lower bound. For the upper bound, we add a dummy node with value  $L$ , which is the largest possible proximity value of THT. All other processes are similar to those of PHP and omitted here.

## 6 EXTENSIONS OF FLoS TO THE PROXIMITY MEASURES HAVING LOCAL MAXIMUM

In this section, we study how to extend the FLoS method to random walk with restart, RoundTripRank, Katz score, and absorption probability.

The idea is to use the relationships between PHP and these proximity measures. Figure 7 summarizes the relationships between PHP and other proximity measures. For RWR, its proximity is proportional to the PHP proximity multiplied by the node degree. For RT, its proximity is proportional to the PHP proximity multiplied by the node degree to the power of  $\beta$ , where  $\beta$  is a constant in RT. For KZ, we define a new proximity measure  $\text{PHP}'$ , which is a variant of PHP. There is a simple relationship between KZ and  $\text{PHP}'$ . For AP, we define another new proximity measure  $\text{PHP}''$ , which is also a variant of PHP. There is a simple relationship between AP and  $\text{PHP}''$ . Compared with PHP, the transition probabilities in  $\text{PHP}'$  and  $\text{PHP}''$  are changed. In PHP, the transition probability is normalized by the node degree  $w_i$ . In  $\text{PHP}'$ , the transition probability is normalized by the maximum degree  $w_{\max}$ . In  $\text{PHP}''$ , the transition probability is normalized by the value  $(\lambda_i + w_i)$ ,

where  $\lambda_i$  is a constant in AP. Thus, the FLoS algorithm for PHP can be readily modified for PHP' and PHP''. Appendix G provides proofs for these relationships.

Next, we use RWR as an example to show how to extend FLoS to these proximity measures by using the relationship. Suppose that node  $v \in \delta S^t$  has the largest PHP proximity value. Based on Theorem 1, for any node  $i \in \bar{S}^t$ , we have that  $\text{PHP}(i) \leq \text{PHP}(v)$ . Let  $w(\bar{S}^t)$  denote the maximum degree of unvisited nodes in  $\bar{S}^t$ . We have that  $w_i \cdot \text{PHP}(i) \leq w(\bar{S}^t) \cdot \text{PHP}(i) \leq w(\bar{S}^t) \cdot \text{PHP}(v)$ . Therefore, if we maintain the maximum degree of unvisited nodes, we can develop the upper bound for the proximity values of unvisited nodes.

Specifically, we can apply FLoS to RWR as follows. In Algorithm 3, we can change line 1 to the following line.

1:  $u \leftarrow \arg\max_{i \in \delta S^{t-1}} w_i \cdot (\mathbf{r}_i^{t-1} + \bar{\mathbf{r}}_i^{t-1});$

In Algorithm 6, we can change line 2 and 3 to the following two lines.

2:  $K \leftarrow k$  nodes in  $S^t \setminus (\delta S^t \cup \{q\})$  with largest  $w_i \cdot \mathbf{r}_i^t$  values;  
 3: **if**  $\min_{i \in K} w_i \cdot \mathbf{r}_i^t \geq \max_{i \in S^t \setminus (K \cup \{q\})} w_i \cdot \mathbf{r}_i^t$  **and**  $\min_{i \in K} w_i \cdot \mathbf{r}_i^t \geq w(\bar{S}^t) \cdot \max_{i \in \delta S^t} \bar{\mathbf{r}}_i^t$  **then**  $\text{bStop} \leftarrow \text{true};$

All other processes remain the same.

For other proximity measures, we extend the FLoS algorithm in a similar way. Please see Appendix G for further details.

## 7 TOP-K REVERSE-PROXIMITY QUERY PROBLEM

In this section, we study the top- $k$  reverse-proximity query problem [24] and discuss how FLoS can be applied to solve it efficiently.

Given a query node  $q$ , we can compute the proximity values of all other nodes. We can also use each node  $i$  as the query, and compute the proximity value of  $q$ . We refer to this proximity of node  $q$  as the reverse proximity of node  $i$ . The top- $k$  reverse-proximity query problem aims at finding the top- $k$  nodes that are ranked by the reverse proximity. In Table 2, only EI and KZ are symmetric and all other proximity measures are not symmetric. The top- $k$  reverse-proximity query problem is different from the top- $k$  proximity query problem when the proximity measure is not symmetric.

Note that the top- $k$  reverse-proximity query problem is different from the reverse top- $k$  problem studied in [25]. Given a query node  $q$ , the reverse top- $k$  problem aims at finding all the nodes that have  $q$  in their top- $k$  proximity sets. In this paper, we study the top- $k$  reverse-proximity query problem, which aims at finding the top- $k$  nodes ranked by the reverse proximity [24].

The top- $k$  reverse-proximity query problem has been studied when RWR is used as the proximity measure [24]. In a recent paper [10], the original and reverse proximity values in RWR are interpreted as importance and specificity respectively. If node  $i$  has large RWR proximity value when the query node is  $q$ , node  $i$  is important for node  $q$ . On the other hand, if node  $q$  has large RWR proximity value when the query node is  $i$ , node  $i$  is specific for node  $q$ . The authors show that ranking by the combination of two directions performs better than ranking by one direction.

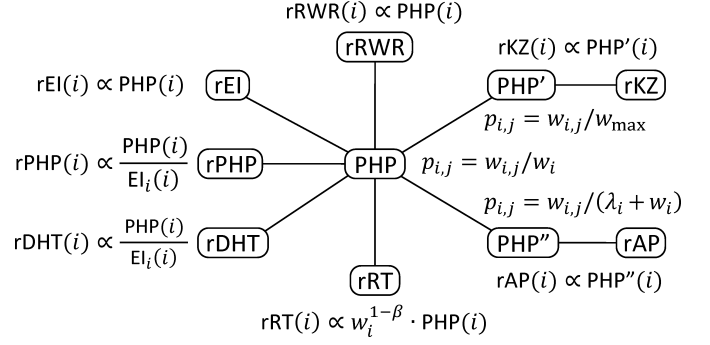


Fig. 8. Relationships between PHP and the reverse proximity measures

The naive method to solve the top- $k$  reverse-proximity query problem is as follows. First, each node is used as the query node, and the proximity value of node  $q$  is computed by the iterative method. Then the top- $k$  nodes with largest reverse proximity values are selected. Suppose that the iterative method takes  $O(\alpha m)$  for each query node. The naive method takes time  $O(\alpha mn)$ , where  $\alpha$  is the number of iterations in the iterative method,  $m$  is the number of edges, and  $n$  is the number of nodes. This is expensive and prohibitive for large graphs.

For the RWR proximity measure, it is shown that the reverse proximity vector can be computed using the iterative method, which has the same complexity  $O(\alpha m)$  as computing the original proximity vector [25]. However, the iterative method is still expensive since it needs to iterate over the entire graph. Moreover, it is unclear how to compute the reverse proximity vectors for other measures in a similar way.

To extend FLoS to the reverse proximity measures, we use the relationships between PHP and the reverse proximity measures. Figure 8 summarizes these relationships. rPHP, rRWR, rEI, rDHT, rRT, rKZ, and rAP represent the reversed version of their corresponding measures. Appendix H provides the proofs for these relationships. Based on these relationships, we can develop the bounds for the reverse proximity values based on the bounds for the PHP or its variant proximity values. Appendix H shows more details about how to extend the FLoS algorithm to the reverse proximity measures.

## 8 EXPERIMENTAL RESULTS

In this section, we present extensive experimental results on evaluating the performance of the FLoS algorithm. The datasets are shown in Table 4. The real datasets are publicly available from the website <http://snap.stanford.edu/data/>. The synthetic datasets are generated using the Erdős-Rényi random graph (RAND) model [26] and R-MAT model [27] with different parameters. All programs are written in C++. All experiments are performed on a server with 32G memory, Intel Xeon 3.2GHz CPU, and Redhat 4.1.2 OS.

### 8.1 State-of-the-Art Methods

The measures we use include PHP, EI, RWR, RT, KZ, THT, and AP. We compare FLoS with the state-of-the-art methods for each measure as summarized in Table 5. These methods are categorized into global and local methods.



TABLE 4  
Datasets used in the experiments

	Datasets	Abbr.	#Nodes	#Edges
Real	Amazon	AZ	334,863	925,872
	DBLP	DP	317,080	1,049,866
	Youtube	YT	1,134,890	2,987,624
	LiveJournal	LJ	3,997,962	34,681,189
Synthetic	In-memory	-	Varying size	
		-	Varying density	
	Disk-resident	-	Varying size	

TABLE 5  
State-of-the-art methods used for comparison

Our methods (Exact)	State-of-the-art methods			
	Abbr.	Key idea	Ref.	Exactness
FLoS_PHP	GI_PHP	Global iteration	[12]	Exact
	DNE	Local search	[7]	Approx.
	NN_EI	Local search	[9]	Exact
	LS_EI	Local search	[1]	Approx.
FLoS_RWR	GI_RWR	Global iteration	[12]	Exact
	GE_RWR	Graph embedding	[19]	Approx.
	Castanet	Improved GI	[4]	Exact
	K-dash	Matrix inversion	[14]	Exact
	LS_RWR	Local search	[1]	Approx.
FLoS_RT	GI_RT	Global iteration	[12]	Exact
	LS_RT	Local search	[10]	Approx.
FLoS_KZ	GI_KZ	Global iteration	[12]	Exact
	LS_KZ	Local search	[23]	Approx.
	AA_KZ	Improved GI	[13]	Exact
FLoS_THT	GI_THT	Global iteration	[12]	Exact
	LS_THT	Local search	[5]	Approx.
FLoS_AP	GI_AP	Global iteration	[12]	Exact
FLoS_rPHP	GI_rPHP	Global iteration	[12]	Exact

The global iteration (GI) method directly applies the iterative method on the entire graph [12]. It guarantees to find the exact top- $k$  nodes. The graph embedding (GE) method can answer the query in constant time after embedding [19]. It can only be applied to RWR. However, the embedding process is very time consuming. Moreover, it only returns approximate results. The Castanet algorithm is specifically designed for RWR. It improves the GI method and guarantees the exactness of the results [4]. AA\_KZ improves the global iteration method by prioritized execution of the iterative computation and also guarantees the exactness of the results [13]. K-dash is the state-of-the-art matrix-based method for RWR which guarantees result exactness [14]. Note that K-dash and GE can only be applied on two medium-sized real graphs because of the expensive preprocessing step.

Dynamic neighborhood expansion (DNE) method applies a best-first expansion strategy to find the top- $k$  nodes using PHP [7]. This strategy is heuristic and does not guarantee to find the exact solution. The number of visited nodes is fixed to 4,000 in the experiments. NN\_EI applies the push style method [21], [2] in local search, and guarantees the exactness of the top- $k$  results [9]. Since PHP and EI are equivalent in terms of ranking, we can compare the methods for PHP and EI directly. LS\_RWR applies the dynamic programming technique [28] to develop bounds in local search [1]. It returns approximate results. LS\_EI is based on LS\_RWR and has similar performance [1]. LS\_RT

leverages the push style method [21] developed for RWR to estimate the bounds and find the approximate top- $k$  nodes with largest RoundTripRank proximity values. LS\_KZ locally searches a small portion of the graph and adapts the push style method [21] to find the approximate top- $k$  results for the Katz score. LS\_THT is a local search method for THT [5].

The decay factors in PHP, RWR, EI, and RT are all set to 0.5. The decay factor in KZ is set to  $0.99/w_{\max}$ . In RT, we set the parameter  $\beta = 0.4$ . In AP, we set the parameter  $\lambda_i = 10$  for any node  $i$ . The truncated length in THT is set to 10.

We use FLoS\_rPHP to denote the FLoS method for reverse PHP. Reverse RWR gives the same ranking as PHP, so we only evaluate FLoS\_PHP. The FLoS method for reverse RT is quite similar to that for RT, thus we only evaluate FLoS\_RT. When we set the parameter  $\lambda_i = 10$  for any node  $i$ , AP becomes symmetric. Thus AP and reverse AP give the same ranking results, and we only evaluate FLoS\_AP.

## 8.2 Evaluation on Real Graphs

We study the efficiency of the selected methods on real graphs when varying the number of returned nodes  $k$ . For each  $k$ , we repeat the experiments  $10^3$  times, each with a randomly picked query node. The average running time is reported. For methods using the iteration procedure in Algorithm 7, the termination threshold is set to  $\tau = 10^{-5}$ . We also perform experiments using a fixed number of 10 iterations. The results are similar and omitted here.

### 8.2.1 Evaluation of FLoS\_PHP

Figure 9 shows the running time of different methods for PHP. The running time of DNE is almost a constant for different  $k$ , because it visits a fixed number of nodes. The running time of NN\_EI increases when  $k$  increases. FLoS\_PHP is more efficient than NN\_EI, which demonstrates that the bounds of FLoS are tighter. LS\_EI has a constant running time. This is because it extracts the cluster containing the query node. Note that LS\_EI takes tens of hours in the preprocessing step to cluster the graphs.

Figure 11(a) shows the ratio between the number of visited nodes using FLoS\_PHP and total number of nodes in the graph. The value indicated by the bar is the average ratio of  $10^3$  queries. The minimum and maximum ratios are also shown in the figure. As can be seen from the figure, only a very small part of the graph is needed for FLoS to find the exact solution. Moreover, the ratio decreases when the graph size increases. This indicates that FLoS is more effective for larger graphs.

### 8.2.2 Evaluation of FLoS\_RWR

Figure 10 shows the running time for RWR. K-dash has the best performance after precomputing the matrix inversion as shown in Figures 10(a) and 10(b). The precomputing step of K-dash takes tens of hours for the medium-sized AZ and DP graphs and cannot be applied to the other two larger graphs. GE\_RWR also has fast response time. However, as discussed before, its embedding step is time consuming and not applicable to larger graphs. Moreover, it does not find the exact solution. Castanet method cuts the running time from the GI method by 72% to 91%. LS\_RWR method has

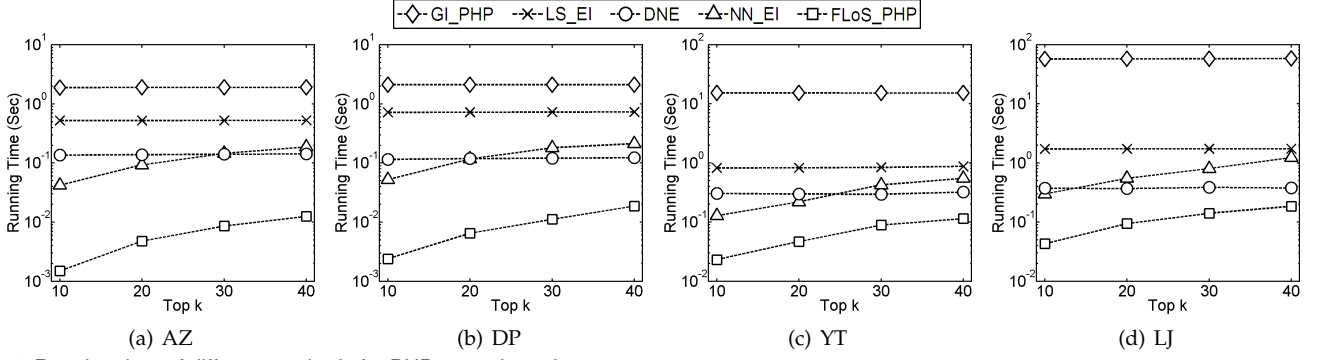


Fig. 9. Running time of different methods for PHP on real graphs

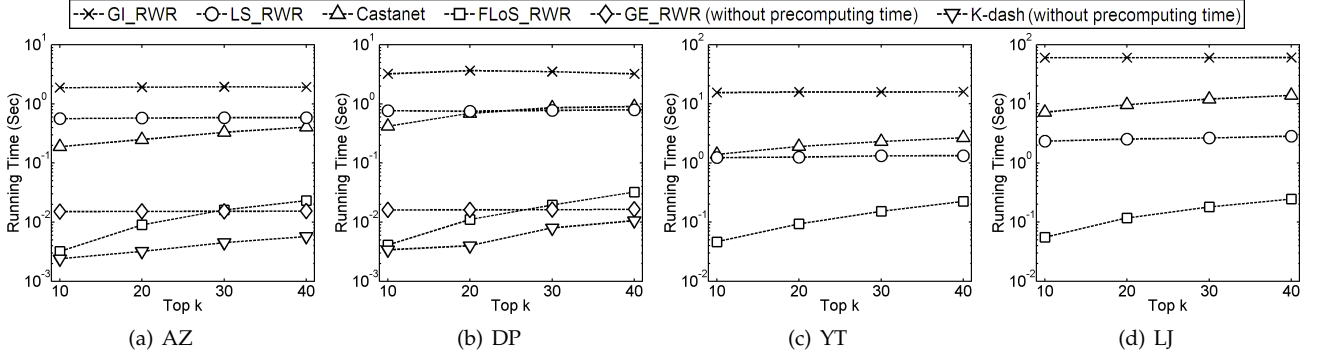


Fig. 10. Running time of different methods for RWR on real graphs

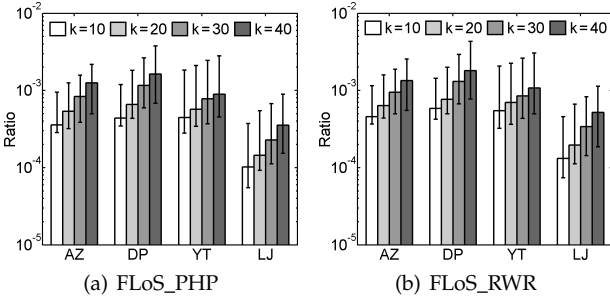


Fig. 11. Ratio between the number of visited nodes and the total number of nodes on real graphs

constant running time, and it needs tens of hours in the precomputing step to cluster the graphs.

Figure 11(b) shows the ratio of the number of visited nodes of the FLoS\_RWR method. The results are similar to that of Figure 11(a).

### 8.2.3 Evaluation of FLoS\_RT

Figure 12(a) shows the running time of different methods for RT. The number on the right side of the rectangle legend indicates the value of  $k$ . Since the GI\_RT method has almost constant running time for different  $k$ , we only show the result when  $k=10$ . The running time of FLoS\_RT and LS\_RT increases when  $k$  increases. FLoS\_RT is the most efficient method. FLoS\_RT is about 1 order of magnitude faster than LS\_RT, and 2 orders of magnitude faster than GI\_RT. LS\_RT uses the push style method to develop the bounds, which are looser than those of FLoS\_RT.

### 8.2.4 Evaluation of FLoS\_KZ

Figure 12(b) shows the running time of different methods for KZ. We also only show the running time when  $k=10$  for

the GI\_KZ method since it has almost constant running time for different  $k$ . FLoS\_KZ, LS\_KZ, and AA\_KZ methods all have increasing running time when increasing  $k$ . FLoS\_KZ is about 1-2 orders of magnitude faster than LS\_KZ and AA\_KZ. LS\_KZ uses the push style method to develop the bounds, which are not as tight as those of FLoS\_KZ. The results also demonstrate that the bounds in AA\_KZ are looser than those of FLoS\_KZ.

### 8.2.5 Evaluation of FLoS\_THT

Figure 12(c) shows the running time for THT. As we can see, FLoS\_THT runs faster than LS\_THT, which is specifically designed to speed up the computation for THT. This is because the lower and upper bounds of FLoS\_THT are tighter than those of LS\_THT. Both of the two local search methods are 2 to 3 orders of magnitude faster than GI\_THT.

### 8.2.6 Evaluation of FLoS\_AP

Figure 12(d) shows the running time for AP. Similar to the results of other proximity measures, FLoS\_AP runs 2-3 orders of magnitude faster than the GI\_AP method.

### 8.2.7 Evaluation of FLoS\_rPHP

In FLoS\_rPHP, we pre-compute the exact values  $E_i(i)$  for each node  $i$  by the K-dash method [14]. The precomputation step takes 28.5 and 34.6 hours for two medium-sized graphs, AZ and DP. Thus we did not apply FLoS\_rPHP on the large graphs.

Figure 12(e) shows the running time of our local search method and the global iteration method for reverse PHP. Similar to the results for other proximity measures, FLoS\_rPHP runs 2-3 orders of magnitude faster than the GI\_rPHP method.

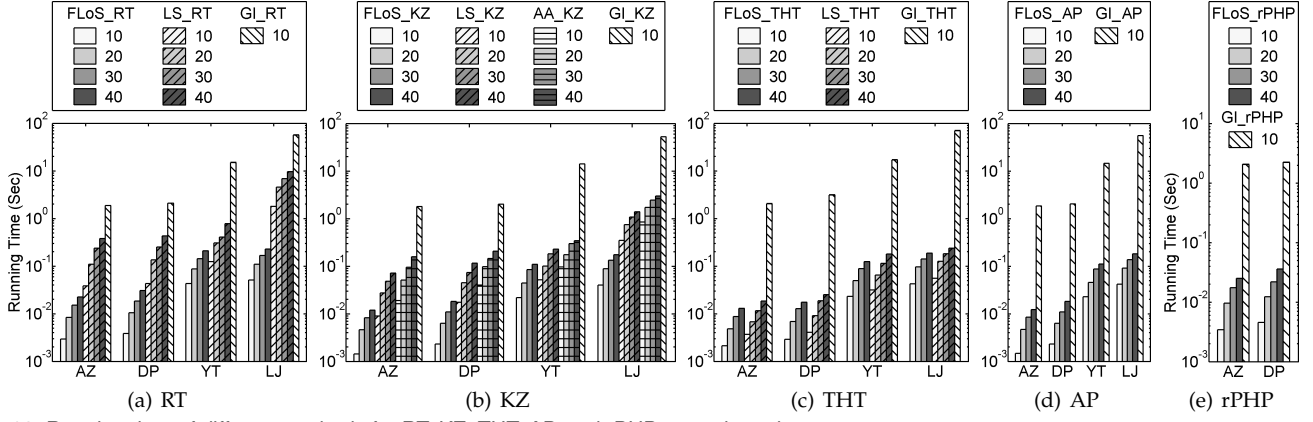


Fig. 12. Running time of different methods for RT, KZ, THT, AP, and rPHP on real graphs

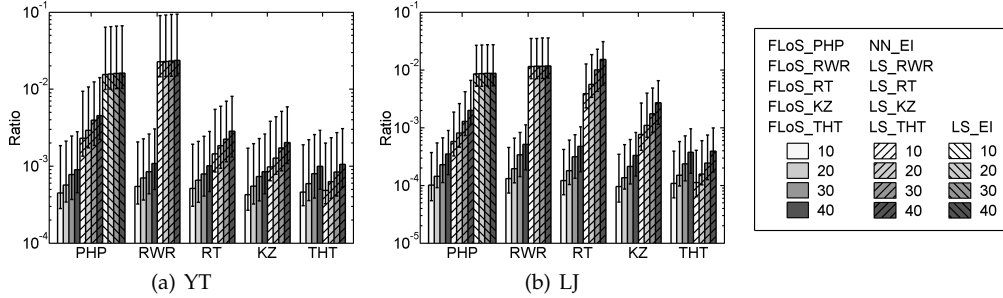


Fig. 13. Ratio between the number of visited nodes and the total number of nodes on real graphs for the local search methods

### 8.2.8 Number of Visited Nodes in Local Search Methods

In this subsection, we study the number of visited nodes of different local search methods on real graphs. The number of visited nodes in the DNE method is fixed, thus it is not included. Figure 13(a) shows the ratio between the number of visited nodes using different local search methods and total number of nodes in the YT graph. Figure 13(b) shows that in the LJ graph. The value indicated by the bar is the average ratio of  $10^3$  queries. The minimum and maximum ratios are also shown in the figure. As can be seen from the figure, other local search methods need to visit larger number of nodes than the FLoS methods do. This demonstrates the tightness of the bounds in the FLoS methods. We also can observe that the LS\_EI and LS\_RWR methods visit relatively large number of nodes and the ratio is stable when the number  $k$  changes. This is because in each expansion of the LS\_EI and LS\_RWR methods, all the nodes in one cluster will be visited. Thus, they need to visit larger number of nodes.

## 8.3 Evaluation on In-Memory Synthetic Graphs

We generate synthetic graphs with different parameters to evaluate the selected methods. More specifically, we study two types of graphs: Erdős-Rényi random graph (RAND) [26] and scale-free graph based on the R-MAT model [27]. There are two parameters, the size and density of the graphs. We study how these two parameters affect the running time of different methods for PHP, RWR, RT, and KZ.

We download the graph generator available from the website <https://github.com/dhruvbird/GTgraph> and use the default parameters to generate two series of graphs with varying size and varying density, using RAND and R-MAT respectively. The graphs with varying size have the

TABLE 6  
Statistics of in-memory synthetic graphs

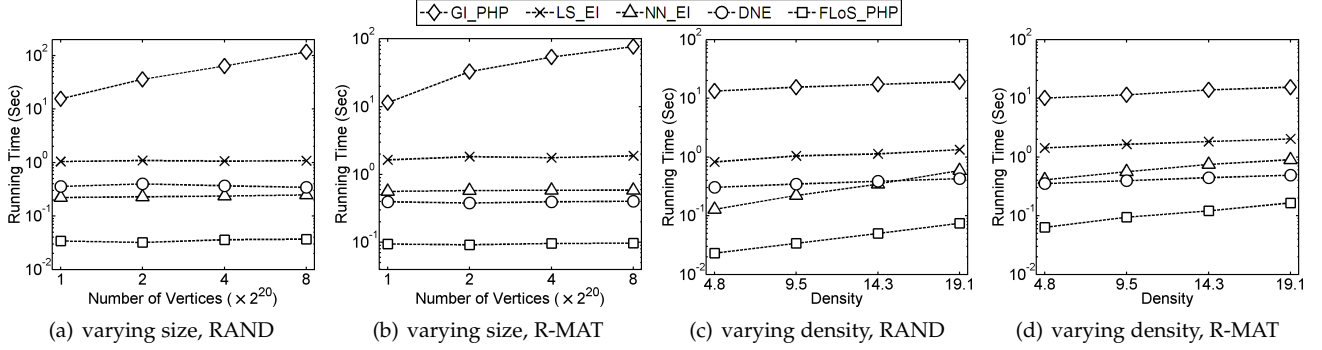
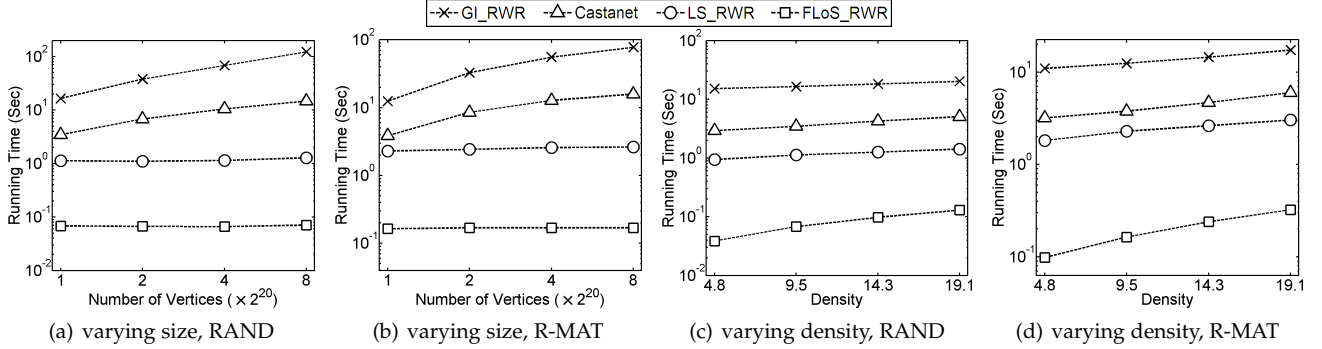
Varying size	#Nodes	$1 \times 2^{20}$	$2 \times 2^{20}$	$4 \times 2^{20}$	$8 \times 2^{20}$
	#Edges	$1 \times 10^7$	$2 \times 10^7$	$4 \times 10^7$	$8 \times 10^7$
	Density	9.5	9.5	9.5	9.5
Varying density	#Nodes	$1 \times 2^{20}$	$1 \times 2^{20}$	$1 \times 2^{20}$	$1 \times 2^{20}$
	#Edges	$5 \times 10^6$	$10 \times 10^6$	$15 \times 10^6$	$20 \times 10^6$
	Density	4.8	9.5	14.3	19.1

same density but different number of nodes. The graphs with varying density have the same number of nodes but different densities. The statistics are shown in Table 6.

We apply the selected methods for PHP, RWR, RT and KZ on these graphs with  $k=20$ . For each graph, we repeat the query  $10^3$  times with randomly picked query nodes, and report the average running time.

### 8.3.1 Evaluation of FLoS\_PHP

Figure 14(a) shows the running time of the selected methods for PHP on the series of RAND graphs with varying size. The running time of GI\_PHP increases as the number of nodes increases. FLoS\_PHP, DNE, NN\_EI and LS\_EI all have almost constant running time when the number of nodes increases. This is because these methods only search locally. When the density of the graph is fixed, adding more nodes to the graph will not change the size of the search space of these methods. Figure 14(b) shows the running time on the series of R-MAT graphs with varying size. Similar trends are observed. Comparing Figure 14(a) and 14(b), GI\_PHP has less running time on R-MAT than on RAND graphs, while other methods have more. The reason is that R-MAT graphs have the power-law distribution, thus it is easier for FLoS\_PHP, DNE, NN\_EI and LS\_EI to encounter hub nodes with larger degree when expanding subgraph. The faster

Fig. 14. Running time of different methods for PHP on in-memory synthetic graphs ( $k=20$ )Fig. 15. Running time of different methods for RWR on in-memory synthetic graphs ( $k=20$ )

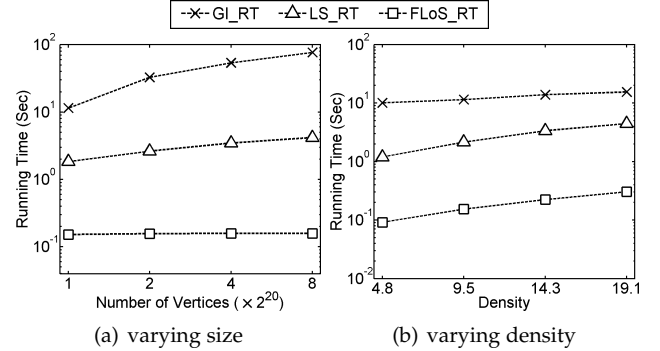
performance of GI\_PHP on R-MAT may be because of the greater data locality due to the hub node.

Figure 14(c) shows the running time of the selected methods for PHP on the series of RAND graphs with varying density. The running time of all the methods increases as the density increases. FLoS\_PHP and NN\_EI have increasing running time because the number of visited nodes in these two methods increases when the density becomes larger. LS\_EI has increasing running time because the number of nodes and edges increases in local clusters. Figure 14(d) shows the running time on the series of R-MAT graphs with varying density. Similar trends are observed.

### 8.3.2 Evaluation of FLoS\_RWR

Figure 15(a) shows the running time of the selected methods for RWR on the series of RAND graphs with varying size. The running time of GI\_RWR and Castanet increases as the number of nodes increases. Castanet method cuts the running time from the GI method by 69% to 88%. FLoS\_RWR and LS\_RWR both have almost constant running time when the number of nodes increases. This is because FLoS\_RWR and LS\_RWR only search locally. Figure 15(b) shows the running time on the series of R-MAT graphs with varying size. Similar trends are observed. Comparing Figure 15(a) and 15(b), GI\_RWR has less running time on the R-MAT graphs than on the RAND graphs, while other methods have more. The reason is similar as what discussed previously.

Figure 15(c) shows the running time on the series of RAND graphs with varying density. The running time of all the methods increases as the density increases. Figure 15(d) shows the running time on the series of R-MAT graphs with varying density. Similar trends are observed.

Fig. 16. Running time of different methods for RT on in-memory synthetic graphs (R-MAT,  $k=20$ )

### 8.3.3 Evaluation of FLoS\_RT

Figure 16(a) shows the running time of the selected methods for RT on the series of R-MAT graph with varying size. The running time of GI\_RT increases as the number of nodes increases. FLoS\_RT has almost constant running time when the number of nodes increases. Because it only searches locally. LS\_RT has increasing running time. LS\_RT needs to find the node with the largest residual proximity value in each iteration. When the number of nodes in the graph increases, the search space may also increase. This may be the reason why it has a slightly increasing running time.

Figure 16(b) shows the running time of the selected methods for RT on the series of R-MAT graph with varying density. The running time of all the methods increases as the density increases. Both FLoS\_RT and LS\_RT have increasing running time because they will visit more nodes in a graph with larger density.

### 8.3.4 Evaluation of FLoS\_KZ

Figure 17(a) shows the running time of the selected methods for KZ on the series of R-MAT graph with varying size. The

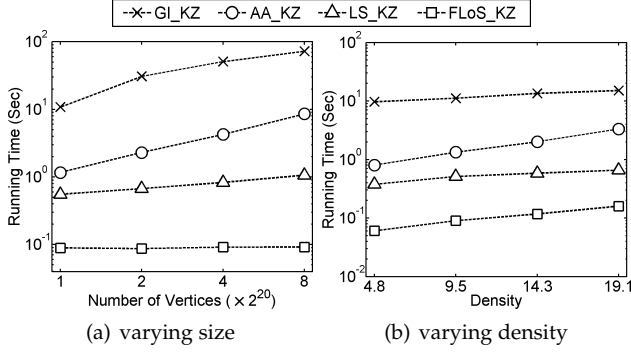


Fig. 17. Running time of different methods for KZ on in-memory synthetic graphs (R-MAT,  $k=20$ )

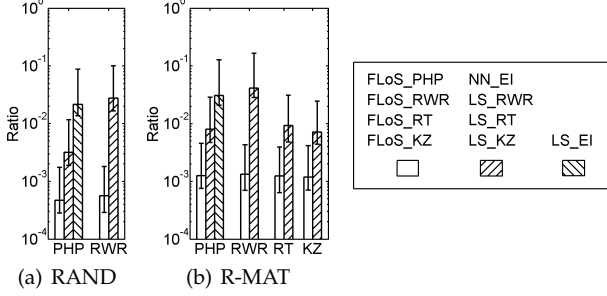


Fig. 18. Ratio between the number of visited nodes and the total number of nodes on synthetic graphs for the local search methods (2<sup>20</sup> nodes and 10<sup>7</sup> edges,  $k=20$ )

running time of GI\_KZ increases as the number of nodes increases. FLoS\_KZ has almost constant running time when the number of nodes increases. LS\_KZ and AA\_KZ both have increasing running time when increasing graph size. LS\_KZ needs to update the node with the largest residual proximity value in each iteration, thus it has a slightly increasing running time when the graph size increases. In AA\_KZ, computing the upper bound of each node requires linear time  $O(m)$ . Thus it has increasing running time.

Figure 17(b) shows the running time of the selected methods for KZ on the series of R-MAT graph with varying density. The running time of all the methods increases as the density increases. The reason why FLoS\_KZ, LS\_KZ and AA\_KZ have increasing running time is that they need to visit more nodes when increasing the graph density.

### 8.3.5 Number of Visited Nodes in Local Search Methods

In this subsection, we study the number of visited nodes using different local search methods on synthetic graphs. We use the synthetic graphs with 2<sup>20</sup> nodes and 10<sup>7</sup> edges. The number of query nodes is fixed to  $k=20$ . Figure 18(a) shows the ratio between the number of visited nodes using different local search methods for PHP and RWR and total number of nodes in the RAND graph. Figure 18(b) shows the ratio between the number of visited nodes using different local search methods for PHP, RWR, RT and KZ and total number of nodes in the R-MAT graph. The value indicated by the bar is the average ratio of 10<sup>3</sup> queries. The minimum and maximum ratios are also shown in the figure. As can be seen from the figure, other local search methods need to visit larger number of nodes than the FLoS methods do. This demonstrates the tightness of the bounds in the FLoS methods.

TABLE 7  
Statistics of disk-resident synthetic graphs

#Nodes	$16 \times 2^{20}$	$32 \times 2^{20}$	$48 \times 2^{20}$	$64 \times 2^{20}$
#Edges	$16 \times 10^7$	$32 \times 10^7$	$48 \times 10^7$	$64 \times 10^7$
Disk size	3.1 G	6.5 G	9.9 G	13.2 G

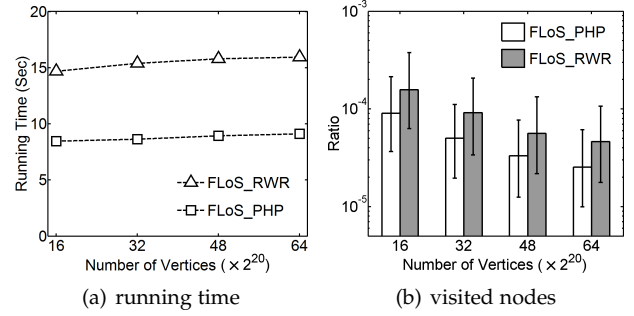


Fig. 19. Results of FLoS\_PHP and FLoS\_RWR on disk-resident synthetic graphs ( $k=20$ )

## 8.4 Evaluation on Disk-Resident Synthetic Graphs

What if the graphs are too large to fit into memory? To test the performance of FLoS on disk-resident graphs, we generate disk-resident R-MAT graphs, whose statistics are in Table 7. We use the open source Neo4j (available from <http://www.neo4j.org>) version 2.0 graph database. The FLoS method for disk-resident graphs only calls some basic query functions provided by Neo4j, such as, querying the neighbors of one node. And the remaining work is the same as that for in-memory graphs. We apply the FLoS\_PHP and FLoS\_RWR methods on the disk-resident graphs with  $k=20$ . We repeat the query 10<sup>3</sup> times with randomly picked query nodes and report the average running time. In the experiments, we restrict the memory usage to 2 GB.

Figure 19(a) shows the running time of FLoS\_PHP and FLoS\_RWR. From the figure, we can see that FLoS can process disk-resident graphs in tens of seconds. The reason is that FLoS only needs to find the neighbors of visited nodes and the transition probabilities on the edges. These results also verify that FLoS has almost constant running time when the number of nodes increases. Figure 19(b) shows the ratio of the number of visited nodes to the total number of nodes in the graph. FLoS only needs to explore a small portion of the whole graph to return the top- $k$  nodes. When the graph size becomes larger, the portion of visited nodes becomes smaller.

## 9 CONCLUSION

Top- $k$  nodes query in large graphs is a fundamental problem that has attracted intensive research interests. Existing methods need expensive preprocessing steps or are designed for specific proximity measures. In this paper, we propose a unified method, FLoS, which adopts a local search strategy to find the exact top- $k$  nodes efficiently. FLoS is based on the no local optimum property of proximity measures. By exploiting the relationship among different proximity measures, we can also extend FLoS to the proximity measures having local optimum. FLoS can be further extended to solve the top- $k$  reverse-proximity query problem. Extensive experimental results demonstrate that FLoS enables efficient and exact query for a variety of random walk based proximity measures.

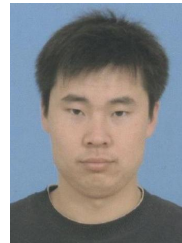


## ACKNOWLEDGMENTS

This work was partially supported by the National Science Foundation grants IIS-1162374, IIS-1218036, IIS-0953950, the NIH/NIGMS grant R01GM103309, and the OSC (Ohio Supercomputer Center) grant PGS0218.

## REFERENCES

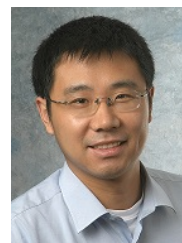
- [1] P. Sarkar and A. W. Moore, "Fast nearest-neighbor search in disk-resident graphs," in *KDD*, 2010, pp. 513–522.
- [2] S. Chakrabarti, A. Pathak, and M. Gupta, "Index design and query processing for graph conductance search," *VLDB J.*, vol. 20, no. 3, pp. 445–470, 2011.
- [3] P. Lee, L. V. Lakshmanan, and J. X. Yu, "On top-k structural similarity search," in *ICDE*, 2012, pp. 774–785.
- [4] Y. Fujiwara, M. Nakatsuji, H. Shiokawa, T. Mishima, and M. Onizuka, "Efficient ad-hoc search for personalized PageRank," in *SIGMOD*, 2013, pp. 445–456.
- [5] P. Sarkar and A. W. Moore, "A tractable approach to finding closest truncated-commute-time neighbors in large graphs," in *UAI*, 2007, pp. 335–343.
- [6] Z. Guan, J. Wu, Q. Zhang, A. Singh, and X. Yan, "Assessing and ranking structural correlations in graphs," in *SIGMOD*, 2011, pp. 937–948.
- [7] C. Zhang, L. Shou, K. Chen, G. Chen, and Y. Bei, "Evaluating geo-social influence in location-based social networks," in *CIKM*, 2012, pp. 1442–1451.
- [8] H. Tong, C. Faloutsos, and J.-Y. Pan, "Fast random walk with restart and its applications," in *ICDM*, 2006, pp. 613–622.
- [9] P. Bogdanov and A. Singh, "Accurate and scalable nearest neighbors in large networks based on effective importance," in *CIKM*, 2013, pp. 523–528.
- [10] Y. Fang, K.-C. Chang, and H. W. Lauw, "RoundTripRank: Graph-based proximity with importance and specificity," in *ICDE*, 2013, pp. 613–624.
- [11] X.-M. Wu, Z. Li, A. M. So, J. Wright, and S.-F. Chang, "Learning with partially absorbing random walks," in *NIPS*, 2012, pp. 3077–3085.
- [12] Y. Saad, *Iterative methods for sparse linear systems*. SIAM, 2003.
- [13] S. Khemmarat and L. Gao, "Fast top-k path-based relevance query on massive graphs," in *ICDE*, 2014, pp. 316–327.
- [14] Y. Fujiwara, M. Nakatsuji, M. Onizuka, and M. Kitsuregawa, "Fast and exact top-k search for random walk with restart," *PVLDB*, vol. 5, no. 5, pp. 442–453, 2012.
- [15] Y. Fujiwara, M. Nakatsuji, T. Yamamuro, H. Shiokawa, and M. Onizuka, "Efficient personalized PageRank with accuracy assurance," in *KDD*, 2012, pp. 15–23.
- [16] S. Cohen, B. Kimelfeld, and G. Koutrika, "A survey on proximity measures for social networks," in *Search Computing*, 2012, pp. 191–206.
- [17] P. Sarkar, A. W. Moore, and A. Prakash, "Fast incremental proximity search in large graphs," in *ICML*, 2008, pp. 896–903.
- [18] L. Katz, "A new status index derived from sociometric analysis," *Psychometrika*, vol. 18, no. 1, pp. 39–43, 1953.
- [19] X. Zhao, A. Chang, A. D. Sarma, H. Zheng, and B. Y. Zhao, "On the embeddability of random walk distances," *PVLDB*, vol. 6, no. 14, pp. 1690–1701, 2013.
- [20] Q. Mei, D. Zhou, and K. Church, "Query suggestion using hitting time," in *CIKM*, 2008, pp. 469–478.
- [21] P. Berkhin, "Bookmark-coloring algorithm for personalized PageRank computing," *Internet Mathematics*, vol. 3, no. 1, pp. 41–62, 2006.
- [22] M. Gupta, A. Pathak, and S. Chakrabarti, "Fast algorithms for top-k personalized PageRank queries," in *WWW*, 2008, pp. 1225–1226.
- [23] P. Esfandiar, F. Bonchi, D. F. Gleich *et al.*, "Fast Katz and commutators: Efficient estimation of social relatedness in large networks," in *Algorithms and Models for the Web-Graph*, 2010, pp. 132–145.
- [24] A. A. Benczur, K. Csalogany, T. Sarlos, and M. Uher, "SpamRank - Fully automatic link spam detection work in progress," in *AIRWeb*, 2005.
- [25] A. W. Yu, N. Mamoulis, and H. Su, "Reverse top-k search using random walk with restart," *PVLDB*, vol. 7, no. 5, pp. 401–412, 2014.
- [26] P. Erdős and A. Rényi, "On the evolution of random graphs," *Magyar Tud. Akad. Mat. Kutató Int. Közl.*, vol. 5, pp. 17–61, 1960.
- [27] D. Chakrabarti, Y. Zhan, and C. Faloutsos, "R-MAT: A recursive model for graph mining," in *SDM*, 2004, pp. 442–446.
- [28] G. Jeh and J. Widom, "Scaling personalized web search," in *WWW*, 2003, pp. 271–279.
- [29] C. Meyer, *Matrix analysis and applied linear algebra*. SIAM, 2000.
- [30] E. A. Guillemin, *Introductory circuit theory*. John Wiley & Sons, 1953.
- [31] G. Jeh and J. Widom, "SimRank: A measure of structural-context similarity," in *KDD*, 2002, pp. 538–543.



**Yubao Wu** received the Bachelor's and Master's degrees both in Dalian University of Technology, China. He is a fourth year Ph.D. student in the Department of Electrical Engineering and Computer Science, Case Western Reserve University. His research interests include big data analytics, data mining and bioinformatics.



**Ruoming Jin** received the doctor's degree in Computing Science from the Ohio State University in 2005. He is an associate professor in the Computer Science Department at Kent State University. His research interests are on Data Mining, Database, Biomedical Informatics and Cloud Computing.



**Xiang Zhang** received the doctor's degree in Computer Science from the University of North Carolina at Chapel Hill in 2011. He is the T&D Schroeder Assistant professor in the Department of Electrical Engineering and Computer Science at Case Western Reserve University. His research bridges the areas of data mining, database and bioinformatics.

## APPENDIX A

### EI, DHT, AND THT

Effective importance (EI) [9] is the degree normalized version of random walk with restart [8], which can be defined as

$$\mathbf{r}_i = \begin{cases} c \sum_{j \in N_i} p_{i,j} \mathbf{r}_j + \frac{1-c}{w_i}, & \text{if } i=q, \\ c \sum_{j \in N_i} p_{i,j} \mathbf{r}_j, & \text{if } i \neq q, \end{cases}$$

where  $c$  ( $0 < c < 1$ ) is a constant decay factor.

**Lemma 2.** EI has no local maximum.

*Proof:* Suppose that node  $i$  is a local maximum. We have  $\mathbf{r}_i = c \sum_{j \in N_i} p_{i,j} \mathbf{r}_j \leq c \sum_{j \in N_i} p_{i,j} \mathbf{r}_i = c \mathbf{r}_i < \mathbf{r}_i$ . Thus we get a contradiction that  $\mathbf{r}_i < \mathbf{r}_i$ .  $\square$

Discounted hitting time (DHT) [1] is a variant of hitting time [20], which can be defined as

$$\mathbf{r}_i = \begin{cases} 0, & \text{if } i=q, \\ 1 + c \sum_{j \in N_i} p_{i,j} \mathbf{r}_j, & \text{if } i \neq q, \end{cases}$$

where  $c$  ( $0 < c < 1$ ) is a constant decay factor.

**Lemma 3.** DHT has no local minimum.

*Proof:* The maximum discounted hitting time that a node could have is  $\frac{1}{1-c}$ , when it cannot reach  $q$ . For connected graph, we have  $\mathbf{r}_i < \frac{1}{1-c}$ . Now suppose that node  $i$  is a local minimum. We have that  $\mathbf{r}_i = 1 + c \sum_{j \in N_i} p_{i,j} \mathbf{r}_j \geq 1 + c \sum_{j \in N_i} p_{i,j} \mathbf{r}_i = 1 + c \mathbf{r}_i$ . Thus  $\mathbf{r}_i \geq \frac{1}{1-c}$ , which contradicts the fact that  $\mathbf{r}_i < \frac{1}{1-c}$ .  $\square$

Truncated hitting time (THT) [5] is another variant of hitting time [20]. The L-truncated hitting time only considers paths of length less than  $L$ . Let  $\mathbf{r}_i^L$  and  $\mathbf{r}_i^{L-1}$  be the  $L$ - and  $(L-1)$ -truncated hitting time of node  $i$  respectively when the query is  $q$ . The L-truncated hitting time can be defined as

$$\mathbf{r}_i^L = \begin{cases} 0, & \text{if } i=q, \\ 1 + \sum_{j \in N_i} p_{i,j} \mathbf{r}_j^{L-1}, & \text{if } i \neq q. \end{cases}$$

If a node is no less than  $L$  hops away from the query node, its proximity is set to  $L$ . If node  $i$  is within  $L$  hops away from the query node, we have that  $\mathbf{r}_i^L < L$ .

**Lemma 4.** THT has no local minimum for the nodes within  $L$  hops from the query node.

*Proof:* Suppose that node  $i$  is within  $L$  hops from  $q$ . We can prove that  $\forall j \in N_i, 1 + \mathbf{r}_j^{L-1} \geq \mathbf{r}_j^L$ , and there exists at least one node  $j \in N_i$  such that  $1 + \mathbf{r}_j^{L-1} > \mathbf{r}_j^L$  by mathematical induction on  $L$ . Now suppose that  $i$  is a local minimum, i.e.,  $\mathbf{r}_i^L \geq \mathbf{r}_i^L$  ( $\forall j \in N_i$ ). We have that  $\mathbf{r}_i^L = 1 + \sum_{j \in N_i} p_{i,j} \mathbf{r}_j^{L-1} = \sum_{j \in N_i} p_{i,j} (1 + \mathbf{r}_j^{L-1}) > \sum_{j \in N_i} p_{i,j} \mathbf{r}_j^L \geq \sum_{j \in N_i} p_{i,j} \mathbf{r}_i^L = \mathbf{r}_i^L$ . We get a contradiction that  $\mathbf{r}_i^L > \mathbf{r}_i^L$ .  $\square$

The proof of Theorem 2 is as follows.

*Proof:* First, we show that PHP and EI are equivalent. Suppose the decay factors in PHP and EI are both set to be  $c$ . Then PHP and EI have the same recursive definition for any node  $i \neq q$ , and we have  $\frac{\text{EI}(i)}{\text{PHP}(i)} = \frac{\text{EI}(q)}{\text{PHP}(q)}$ , which is a constant when  $q$  is fixed. Thus  $\text{PHP}(i)$  is a linear function of  $\text{EI}(i)$ .

Next, we show that PHP and DHT are equivalent. Suppose the decay factors in PHP and DHT are both set to be  $c$ . The relationship between PHP and DHT is  $\text{PHP}(i) = 1 - (1-c) \cdot \text{DHT}(i)$ , which can be proved by substituting it in the recursive definition of PHP. Thus  $\text{PHP}(i)$  is a linear function of  $\text{DHT}(i)$ .  $\square$

## APPENDIX B

### PROOFS OF THEOREMS 3 AND 5 IN SECTION 4.1

The proof of Theorem 3 is as follows.

*Proof:* Let  $\mathbf{P}$  be the original transition probability matrix of PHP. Deleting  $\mathbf{P}_{i,j}$  from  $\mathbf{P}$  is the same as setting  $\mathbf{P}_{i,j}$  to 0. Let  $\mathbf{P}'$  represent the resulting matrix. The PHP proximity vector  $\mathbf{r}$  is computed based on  $\mathbf{P}$ , and  $\mathbf{r}'$  is based on  $\mathbf{P}'$ .

Let  $\Delta \mathbf{r} = \mathbf{r} - \mathbf{r}'$ , and  $\Delta \mathbf{P} = \mathbf{P} - \mathbf{P}'$ . Note that  $\Delta \mathbf{P}$  has only one non-zero element  $\Delta \mathbf{P}_{i,j} = \mathbf{P}_{i,j}$ . We have  $\Delta \mathbf{r} = c(\mathbf{P}' + \Delta \mathbf{P})\mathbf{r} - c\mathbf{P}'\mathbf{r}' = c\mathbf{P}'\Delta \mathbf{r} + c\Delta \mathbf{P}\mathbf{r} = c\mathbf{P}'\Delta \mathbf{r} + \mathbf{e}'$ , where  $\mathbf{e}'$  is a vector with the only non-zero element  $\mathbf{e}'_i = c\mathbf{P}_{i,j}\mathbf{r}_j$ . The solution of the previous equation is  $\Delta \mathbf{r} = (\mathbf{I} - c\mathbf{P}')^{-1}\mathbf{e}'$ . The elements of  $c\mathbf{P}'$  are non-negative and  $\|c\mathbf{P}'\|_\infty < 1$ . We have that  $(\mathbf{I} - c\mathbf{P}')^{-1} = \sum_{l=0}^{\infty} (c\mathbf{P}')^l \geq 0$  [29]. Thus the solution  $\Delta \mathbf{r}$  must be non-negative. This completes the proof.  $\square$

The proof of Theorem 5 is as follows.

*Proof:* Let  $\mathbf{P}$  be the original transition probability matrix of PHP. Changing the destination of transition probability  $p_{i,j}$  from node  $j$  to node  $u$  is the same as adding  $\mathbf{P}_{i,j}$  to  $\mathbf{P}_{i,u}$  and setting  $\mathbf{P}_{i,j}$  to 0. Let  $\mathbf{P}'$  represent the resulting matrix. PHP proximity  $\mathbf{r}$  is computed based on  $\mathbf{P}$ , and  $\mathbf{r}'$  is based on  $\mathbf{P}'$ .

Let  $\Delta \mathbf{r} = \mathbf{r}' - \mathbf{r}$ , and  $\Delta \mathbf{P} = \mathbf{P} - \mathbf{P}'$ . Note that  $\Delta \mathbf{P}$  has only two non-zero elements  $\Delta \mathbf{P}_{i,j} = \mathbf{P}_{i,j}$  and  $\Delta \mathbf{P}_{i,u} = -\mathbf{P}_{i,j}$ . We have that  $\Delta \mathbf{r} = c\mathbf{P}'\mathbf{r}' - c(\mathbf{P}' + \Delta \mathbf{P})\mathbf{r} = c\mathbf{P}'\Delta \mathbf{r} - c\Delta \mathbf{P}\mathbf{r} = c\mathbf{P}'\Delta \mathbf{r} + \mathbf{e}'$ , where  $\mathbf{e}'$  is a vector with the only non-zero element  $\mathbf{e}'_i = c\mathbf{P}_{i,j}(\mathbf{r}_u - \mathbf{r}_j)$ . If  $\mathbf{r}_u \geq \mathbf{r}_j$ ,  $\Delta \mathbf{r}$  must be non-negative. Otherwise, it is non-positive. This completes the proof.  $\square$

## APPENDIX C

### PROOF OF THEOREM 6 IN SECTION 5.2

The proof of Theorem 6 is as follows.

*Proof:* Let  $\mathbf{P}^{t-1}$  and  $\mathbf{P}^t$  be the transition probability matrices at the  $(t-1)$ -th and  $t$ -th iterations respectively. The lower bound  $\mathbf{r}^{t-1}$  is computed based on  $\mathbf{P}^{t-1}$ , and  $\mathbf{r}^t$  is computed based on  $\mathbf{P}^t$ . Since the sizes of matrices  $\mathbf{P}^{t-1}$  and  $\mathbf{P}^t$  are  $|S^{t-1}| \times |S^{t-1}|$  and  $|S^t| \times |S^t|$  respectively, we extend  $\mathbf{P}^{t-1}$  with  $|S^t \setminus S^{t-1}|$  columns and  $|S^t \setminus S^{t-1}|$  rows. The newly added elements are set to 0. After extending,  $\mathbf{P}^{t-1}$  and  $\mathbf{P}^t$  are of the same size. Similarly, we extend  $\mathbf{r}^{t-1}$  with  $|S^t \setminus S^{t-1}|$  elements, which are also set to 0. Thus,  $\mathbf{r}^{t-1}$  and  $\mathbf{r}^t$  are of the same size  $|S^t| \times 1$  after extending. For any node  $i \in S^{t-1}$ , we have  $\mathbf{r}_i^{t-1} > 0$ . Figures 20(a) and 20(b) illustrate these matrices and vectors. We have the following two equations.

$$\begin{cases} \mathbf{r}^{t-1} = c\mathbf{P}^{t-1}\mathbf{r}^{t-1} + \mathbf{e}, \\ \mathbf{r}^t = c\mathbf{P}^t\mathbf{r}^t + \mathbf{e}, \end{cases}$$

where  $\mathbf{e}_q = 1$  and  $\mathbf{e}_i = 0$  if  $i \in S^t \setminus \{q\}$ .

Let  $\Delta \mathbf{r} = \mathbf{r}^t - \mathbf{r}^{t-1}$  and  $\Delta \mathbf{P} = \mathbf{P}^t - \mathbf{P}^{t-1}$ . Note that  $\Delta \mathbf{P}_{i,j} = 0$  if  $i, j \in S^{t-1}$  and  $\Delta \mathbf{P}_{i,j} = \mathbf{P}_{i,j}^{t-1}$  otherwise. Figure 20(c) illustrates the matrix  $\Delta \mathbf{P}$  and vector  $\Delta \mathbf{r}$ . We have that  $\Delta \mathbf{r} = c\mathbf{P}^t\mathbf{r}^t - c\mathbf{P}^{t-1}\mathbf{r}^{t-1} = c\mathbf{P}^t\Delta \mathbf{r} + c\Delta \mathbf{P}\mathbf{r}^{t-1} = c\mathbf{P}^t\Delta \mathbf{r} + \mathbf{e}'$ , where  $\mathbf{e}' = c\Delta \mathbf{P}\mathbf{r}^{t-1}$ .

Note that  $\Delta \mathbf{P}_{i,j} = 0$  if  $i, j \in S^{t-1}$ . We also have that  $\mathbf{r}_j^{t-1} = 0$  if  $j \in S^t \setminus S^{t-1}$ . Thus,  $\mathbf{e}'_i = 0$  if  $i \in S^{t-1}$ . For any node  $i \in S^t \setminus S^{t-1}$ , it must connect to at least one node in  $S^{t-1}$ . Thus, for each node  $i \in S^t \setminus S^{t-1}$ ,  $\Delta \mathbf{P}_{i,j} > 0$  for some node  $j \in S^{t-1}$ .

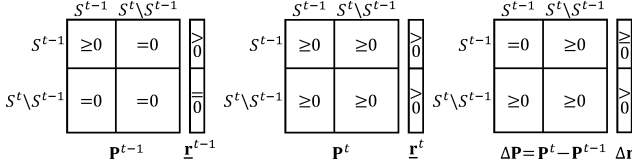
(a)  $(t-1)$ -th iteration (b)  $t$ -th iteration (c) difference

Fig. 20. Transition probability matrices between two adjacent iterations (lower bound)

We also have that  $\mathbf{r}_j^{t-1} > 0$  if  $j \in S^{t-1}$ . Therefore,  $\mathbf{e}'_i > 0$  if  $i \in S^t \setminus S^{t-1}$ . Thus, we have

$$\begin{cases} \mathbf{e}'_i = 0, & \text{if } i \in S^{t-1}, \\ \mathbf{e}'_i > 0, & \text{if } i \in S^t \setminus S^{t-1}. \end{cases}$$

We can get that  $\Delta \mathbf{r} = (\mathbf{I} - c\mathbf{P}^t)^{-1} \mathbf{e}'$ . The elements of  $c\mathbf{P}^t$  are non-negative and  $\|c\mathbf{P}^t\|_\infty < 1$ . We have that  $(\mathbf{I} - c\mathbf{P}^t)^{-1} = \sum_{l=0}^{\infty} c^l (\mathbf{P}^t)^l$ , where  $(\mathbf{P}^t)^l$  represents the matrix  $\mathbf{P}^t$  to the power of  $l$ . Each element  $[(\mathbf{P}^t)^l]_{i,j}$  represents the probability of transiting from node  $i$  to  $j$  in the  $l$ -th step.

Consider a node  $i \in S^{t-1}$ . If  $i \rightsquigarrow S^t \setminus S^{t-1}$ , we have  $[(\mathbf{P}^t)^l]_{i,j} > 0$  for some  $l$  and some node  $j \in S^t \setminus S^{t-1}$ . So,  $[(\mathbf{P}^t)^l \mathbf{e}']_i > 0$  and  $\Delta \mathbf{r}_i > 0$ . If  $i \rightsquigarrow S^t \setminus S^{t-1}$ , we have  $[(\mathbf{P}^t)^l]_{i,j} = 0$  for any  $l$  and any node  $j \in S^t \setminus S^{t-1}$ . So,  $[(\mathbf{P}^t)^l \mathbf{e}']_i = 0$  and  $\Delta \mathbf{r}_i = 0$ .  $\square$

## APPENDIX D

### PROOF OF THEOREM 7 IN SECTION 5.3

In this subsection, we will prove Theorem 7, which shows the monotonicity of the upper bounds. In Section 5.3, from iteration  $(t-1)$  to  $t$ , the transition graph is modified in three steps. We will analyze the three steps rigidly in Lemmas 5, 6, and 7 respectively. Only the process from iteration  $(t-1)$  to  $t$  is considered in the three lemmas. After the discussion of the three lemmas, we will use mathematical induction on each iteration to prove the monotonicity of the upper bound at all the iterations. Thus, in Lemmas 5, 6, and 7, we assume that  $t \geq 2$  and  $\bar{\mathbf{r}}_d^{t-1} > \bar{\mathbf{r}}_i^{t-1}$  for any node  $i \in \delta S^{t-1}$ . The idea is that the condition is first assumed to be satisfied at the  $(t-1)$ -th iteration, and we will prove that it is also satisfied at the  $t$ -th iteration.

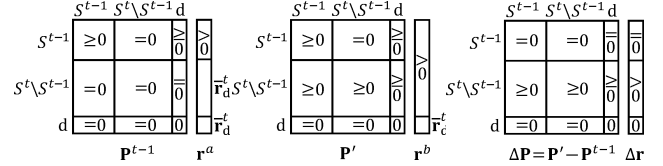
In Lemmas 5, 6, and 7, we use the following symbols. At the  $(t-1)$ -th iteration, we use  $\mathbf{P}^{t-1}$  to denote the transition probability matrix, and  $\bar{\mathbf{r}}^{t-1}$  to denote the upper bound vector. After applying step 1, the transition probability matrix does not change and we use  $\mathbf{r}^a$  to represent the new upper bound vector. After applying step 2, the new transition probability matrix is denoted as  $\mathbf{P}'$  and the new upper bound vector is denoted as  $\mathbf{r}^b$ . After applying step 3, we use  $\mathbf{P}^t$  to denote the transition probability matrix, and  $\bar{\mathbf{r}}^t$  to denote the upper bound vector.

**Lemma 5. (Step 1)** For any node  $i \in S^{t-1}$ , we have

$$\begin{cases} \bar{\mathbf{r}}_i^{t-1} = \mathbf{r}_i^a, & \text{if } i \rightsquigarrow d, \\ \bar{\mathbf{r}}_i^{t-1} > \mathbf{r}_i^a, & \text{if } i \rightsquigarrow d, \end{cases}$$

where  $\rightsquigarrow$  and  $\rightsquigarrow$  represent the reachability in the transition graph at the  $(t-1)$ -th iteration.

*Proof:* In step 1, we decrease the proximity value of the dummy node from  $\bar{\mathbf{r}}_d^{t-1}$  to  $\bar{\mathbf{r}}_d^t = \max_{i \in \delta S^{t-1}} \bar{\mathbf{r}}_i^{t-1}$ . By the assumption,  $\bar{\mathbf{r}}_d^{t-1} > \bar{\mathbf{r}}_i^{t-1}$  for any node  $i \in \delta S^{t-1}$ . Thus,  $\bar{\mathbf{r}}_d^{t-1} > \bar{\mathbf{r}}_d^t$ .



(a) after step 1 (b) after step 2 (c) difference

Fig. 21. Transition probability matrices between two adjacent iterations (upper bound, step 2)

The upper bound vectors  $\bar{\mathbf{r}}^{t-1}$  and  $\mathbf{r}^a$  both are computed based on  $\mathbf{P}^{t-1}$  but with different  $\mathbf{e}$  vectors. The size of matrix  $\mathbf{P}^{t-1}$  is  $(|S^{t-1}|+1) \times (|S^{t-1}|+1)$ . The sizes of vectors  $\bar{\mathbf{r}}^{t-1}$  and  $\mathbf{r}^a$  both are  $(|S^{t-1}|+1) \times 1$ . We have the following two equations.

$$\begin{cases} \bar{\mathbf{r}}^{t-1} = c\mathbf{P}^{t-1} \bar{\mathbf{r}}^{t-1} + \mathbf{e}^{t-1}, \\ \mathbf{r}^a = c\mathbf{P}^{t-1} \mathbf{r}^a + \mathbf{e}^a. \end{cases}$$

In the first equation,  $\mathbf{e}_q^{t-1} = 1$ ,  $\mathbf{e}_d^{t-1} = \bar{\mathbf{r}}_d^{t-1}$ , and  $\mathbf{e}_i^{t-1} = 0$  if  $i \in S^{t-1} \setminus \{q\}$ . In the second equation,  $\mathbf{e}_q^a = 1$ ,  $\mathbf{e}_d^a = \bar{\mathbf{r}}_d^t$ , and  $\mathbf{e}_i^a = 0$  if  $i \in S^{t-1} \setminus \{q\}$ .

Let  $\Delta \mathbf{r} = \bar{\mathbf{r}}^{t-1} - \mathbf{r}^a$ . We have that  $\Delta \mathbf{r} = c\mathbf{P}^{t-1} \Delta \mathbf{r} + \mathbf{e}'$ , where  $\mathbf{e}' = \mathbf{e}^{t-1} - \mathbf{e}^a$ . Since  $\bar{\mathbf{r}}_d^{t-1} > \bar{\mathbf{r}}_d^t$ , we have that

$$\begin{cases} \mathbf{e}'_i = 0, & \text{if } i \in S^{t-1}, \\ \mathbf{e}'_i > 0, & \text{if } i = d. \end{cases}$$

We can get that  $\Delta \mathbf{r} = (\mathbf{I} - c\mathbf{P}^{t-1})^{-1} \mathbf{e}'$ . The elements of  $c\mathbf{P}^{t-1}$  are non-negative and  $\|c\mathbf{P}^{t-1}\|_\infty < 1$ . We have that  $(\mathbf{I} - c\mathbf{P}^{t-1})^{-1} = \sum_{l=0}^{\infty} c^l (\mathbf{P}^{t-1})^l$ , where  $(\mathbf{P}^{t-1})^l$  denotes matrix  $\mathbf{P}^{t-1}$  to the power of  $l$ . Each element  $[(\mathbf{P}^{t-1})^l]_{i,d}$  represents the probability of transiting from node  $i$  to  $d$  in the  $l$ -th step.

Consider a node  $i \in S^{t-1}$ . If  $i \rightsquigarrow d$ , we have that  $[(\mathbf{P}^{t-1})^l]_{i,d} > 0$  for some  $l$ . Thus,  $\Delta \mathbf{r}_i > 0$ , i.e.,  $\bar{\mathbf{r}}_i^{t-1} > \mathbf{r}_i^a$ . If  $i \rightsquigarrow d$ , we have that  $[(\mathbf{P}^{t-1})^l]_{i,d} = 0$  for any  $l$ . Thus,  $\Delta \mathbf{r}_i = 0$ .  $\square$

**Lemma 6. (Step 2)**  $\begin{cases} \mathbf{r}_i^a = \mathbf{r}_i^b, & \text{if } i \in S^{t-1}, \\ \bar{\mathbf{r}}_d^t > \mathbf{r}_i^b, & \text{if } i \in S^t \setminus S^{t-1}. \end{cases}$

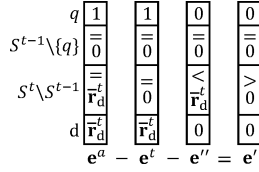
*Proof:* In step 2, we add transition probabilities  $\{p_{i,j}\}$  from the newly added nodes  $i \in (S^t \setminus S^{t-1})$  to nodes  $j \in (S^t)$ , and  $\{p_{i,d}\}$  from  $i$  to the dummy node  $d$ .

Since the sizes of matrices  $\mathbf{P}^{t-1}$  and  $\mathbf{P}'$  are  $(|S^{t-1}|+1) \times (|S^{t-1}|+1)$  and  $(|S^t|+1) \times (|S^t|+1)$  respectively, we extend  $\mathbf{P}^{t-1}$  with  $|S^t \setminus S^{t-1}|$  columns and  $|S^t \setminus S^{t-1}|$  rows as shown in Figure 21(a). The newly added elements are set to 0. Thus, matrices  $\mathbf{P}^{t-1}$  and  $\mathbf{P}'$  are of the same size after extension. Similarly, we extend  $\mathbf{r}^a$  with  $|S^t \setminus S^{t-1}|$  elements, which are set to  $\bar{\mathbf{r}}_d^t$ . Thus,  $\mathbf{r}^a$  and  $\mathbf{r}^b$  are of the same size  $(|S^t|+1) \times 1$  after extension. For any node  $i \in S^{t-1}$ , we must have that  $\mathbf{r}_i^a > 0$ . Figures 21(a) and 21(b) illustrate these matrices and vectors. We have the following two equations.

$$\begin{cases} \mathbf{r}^a = c\mathbf{P}^{t-1} \mathbf{r}^a + \mathbf{e}^a, \\ \mathbf{r}^b = c\mathbf{P}' \mathbf{r}^b + \mathbf{e}^b. \end{cases}$$

In the first equation,  $\mathbf{e}_q^a = 1$ ,  $\mathbf{e}_i^a = 0$  if  $i \in S^{t-1} \setminus \{q\}$  and  $\mathbf{e}_i^a = \bar{\mathbf{r}}_d^t$  if  $i \in (S^t \setminus S^{t-1}) \cup \{d\}$ . In the second equation,  $\mathbf{e}_q^b = 1$ ,  $\mathbf{e}_d^b = \bar{\mathbf{r}}_d^t$ , and  $\mathbf{e}_i^b = 0$  if  $i \in S^t \setminus \{q\}$ .

Let  $\Delta \mathbf{r} = \mathbf{r}^a - \mathbf{r}^b$  and  $\Delta \mathbf{P} = \mathbf{P}' - \mathbf{P}^{t-1}$ . Note that  $\Delta \mathbf{P}_{i,j} = 0$  if  $i \in S^{t-1} \cup \{d\}$  and  $\Delta \mathbf{P}_{i,j} = \mathbf{P}'_{i,j}$  otherwise. Figure 21(c) illustrates the matrix  $\Delta \mathbf{P}$  and vector  $\Delta \mathbf{r}$ . We have that  $\Delta \mathbf{r} = c\mathbf{P}^{t-1} \mathbf{r}^a - c\mathbf{P}' \mathbf{r}^b + \mathbf{e}^a - \mathbf{e}^b = c\mathbf{P}' \Delta \mathbf{r} - c\Delta \mathbf{P} \mathbf{r}^a + \mathbf{e}^a - \mathbf{e}^b = c\mathbf{P}' \Delta \mathbf{r} + \mathbf{e}'$ , where  $\mathbf{e}' = \mathbf{e}^a - \mathbf{e}^b - c\Delta \mathbf{P} \mathbf{r}^a$ .

Fig. 22. The vectors  $\mathbf{e}^a$ ,  $\mathbf{e}^t$ ,  $\mathbf{e}''$ , and  $\mathbf{e}'$ 

Let  $\mathbf{e}'' = c\Delta\mathbf{P}\mathbf{r}^a$ . Note that  $\Delta\mathbf{P}_{i,j} = 0$  if  $i \in S^{t-1} \cup \{d\}$ . Thus,  $\mathbf{e}'' = 0$  if  $i \in S^{t-1} \cup \{d\}$ . We have that  $\Delta\mathbf{P}_{i,j} = 0$  if  $i \in S^t \setminus S^{t-1}$  and  $j \in S^{t-1} \setminus \delta S^{t-1}$  and  $\Delta\mathbf{P}_{i,j} \geq 0$  if  $i \in S^t \setminus S^{t-1}$  and  $j \in \delta S^{t-1}$ . For any node  $j \in \delta S^{t-1}$ ,  $j \rightsquigarrow d$ . By Lemma 5, we have that  $\bar{\mathbf{r}}_j^{t-1} > \mathbf{r}_j^a$ . Since  $\bar{\mathbf{r}}_d^a = \max_{i \in \delta S^{t-1}} \bar{\mathbf{r}}_i^{t-1}$ ,  $\mathbf{r}_j^a < \bar{\mathbf{r}}_d^a$  for any node  $j \in \delta S^{t-1}$ . We also have that  $\mathbf{r}_j^a = \bar{\mathbf{r}}_d^a$  if  $j \in (S^t \setminus S^{t-1}) \cup \{d\}$ . The sum of the elements in the  $i$ -th ( $i \in S^t \setminus S^{t-1}$ ) row of  $\Delta\mathbf{P}$  equals 1. Thus,  $\mathbf{e}_i'' < \bar{\mathbf{r}}_d^a$  if  $i \in S^t \setminus S^{t-1}$ . Since  $\mathbf{e}' = \mathbf{e}^a - \mathbf{e}^t - \mathbf{e}''$ , we have that

$$\begin{cases} \mathbf{e}_i' = 0, & \text{if } i \in S^{t-1} \cup \{d\}, \\ \mathbf{e}_i' > 0, & \text{if } i \in S^t \setminus S^{t-1}. \end{cases}$$

Figure 22 illustrates the vectors  $\mathbf{e}^a$ ,  $\mathbf{e}^t$ ,  $\mathbf{e}''$ , and  $\mathbf{e}'$ .

We can get that  $\Delta\mathbf{r} = (\mathbf{I} - c\mathbf{P}')^{-1}\mathbf{e}'$ . The elements of  $c\mathbf{P}'$  are non-negative and  $\|c\mathbf{P}'\|_\infty < 1$ . We have that  $(\mathbf{I} - c\mathbf{P}')^{-1} = \sum_{l=0}^{\infty} c^l (\mathbf{P}')^l$ , where  $(\mathbf{P}')^l$  represents the matrix  $\mathbf{P}'$  to the power of  $l$ . We have that  $\Delta\mathbf{r} = (\mathbf{I} + \sum_{l=1}^{\infty} c^l (\mathbf{P}')^l)\mathbf{e}' = \mathbf{e}' + \sum_{l=1}^{\infty} c^l (\mathbf{P}')^l \mathbf{e}'$ . Since all the elements in matrix  $\mathbf{P}'$  are non-negative, we have that  $[(\mathbf{P}')^l]_{i,j} \geq 0$  for any  $l$  and any nodes  $i, j \in S^t \cup \{d\}$ .

For any node  $i \in S^t \setminus S^{t-1}$ ,  $\Delta\mathbf{r}_i > 0$ , i.e.,  $\mathbf{r}_i^a = \bar{\mathbf{r}}_d^t > \mathbf{r}_i^b$ .

For any node  $i \in S^{t-1}$ ,  $i \rightsquigarrow S^t \setminus S^{t-1}$  in the transition graph corresponding to  $\mathbf{P}'$ . Thus,  $[(\mathbf{P}')^l]_{i,j} = 0$  for any  $l$  and any node  $j \in S^t \setminus S^{t-1}$ . So,  $\Delta\mathbf{r}_i = 0$ .  $\square$

**Lemma 7.** (Step 3) For any node  $i \in S^t$ , we have that

$$\begin{cases} \mathbf{r}_i^b = \bar{\mathbf{r}}_i^t, & \text{if } i \rightsquigarrow S^t \setminus S^{t-1}, \\ \mathbf{r}_i^b > \bar{\mathbf{r}}_i^t, & \text{if } i \rightsquigarrow S^t \setminus S^{t-1}, \end{cases}$$

where  $\rightsquigarrow$  and  $\rightsquigarrow$  represent the reachability in the transition graph at the  $t$ -th iteration.

*Proof:* In step 3, we add the transition probabilities  $\{p_{j,i}\}$  from nodes  $j \in \delta S^{t-1}$  to the newly added nodes  $i \in S^t \setminus S^{t-1}$ , and remove their correspondences in  $\{p_{j,d}\}$ .

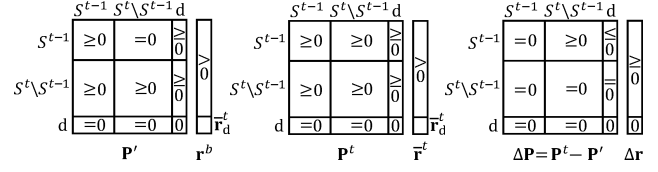
The sizes of matrices  $\mathbf{P}'$  and  $\mathbf{P}^t$  both are  $(|S^t| + 1) \times (|S^t| + 1)$ . The sizes of vectors  $\mathbf{r}^b$  and  $\bar{\mathbf{r}}^t$  both are  $(|S^t| + 1) \times 1$ . Figures 23(a) and 23(b) illustrate these matrices and vectors. We have the following two equations.

$$\begin{cases} \mathbf{r}^b = c\mathbf{P}'\mathbf{r}^b + \mathbf{e}^t, \\ \bar{\mathbf{r}}^t = c\mathbf{P}^t\bar{\mathbf{r}}^t + \mathbf{e}^t, \end{cases}$$

where  $\mathbf{e}_q^t = 1$ ,  $\mathbf{e}_d^t = \bar{\mathbf{r}}_d^t$ , and  $\mathbf{e}_i^t = 0$  if  $i \in S^t \setminus \{q\}$ .

Let  $\Delta\mathbf{r} = \mathbf{r}^b - \bar{\mathbf{r}}^t$  and  $\Delta\mathbf{P} = \mathbf{P}^t - \mathbf{P}'$ . Note that  $\Delta\mathbf{P}_{i,j} = 0$  if  $i \in (S^t \setminus S^{t-1}) \cup \{d\}$  or  $j \in S^{t-1}$ ,  $\Delta\mathbf{P}_{i,j} \geq 0$  if  $i \in S^{t-1}$  and  $j \in S^t \setminus S^{t-1}$ , and  $\Delta\mathbf{P}_{i,d} \leq 0$  if  $i \in S^{t-1}$ . The sum of elements in each row of  $\Delta\mathbf{P}$  equals 0. Figure 23(c) illustrates the matrix  $\Delta\mathbf{P}$  and vector  $\Delta\mathbf{r}$ . We have that  $\Delta\mathbf{r} = c\mathbf{P}'\mathbf{r}^b - c\mathbf{P}^t\bar{\mathbf{r}}^t = c\mathbf{P}^t\Delta\mathbf{r} - c\Delta\mathbf{P}\mathbf{r}^b = c\mathbf{P}^t\Delta\mathbf{r} + \mathbf{e}'$ , where  $\mathbf{e}' = -c\Delta\mathbf{P}\mathbf{r}^b$ .

Consider the set of nodes  $S' = \{i \in \delta S^{t-1} \text{ and } \mathbf{P}_{i,j}^t > 0 \text{ for some } j \in S^t \setminus S^{t-1}\}$ .  $S'$  comprises the set of nodes in  $\delta S^{t-1}$  that are adjacent to the newly added nodes. If node  $i \notin S'$ , the transition probability  $p_{i,j}$  does not change during step 3. Thus, we have that  $\Delta\mathbf{P}_{i,j} = 0$  and  $\mathbf{e}_i' = 0$  for any node  $i \notin S'$  and  $j \in S^t \cup \{d\}$ . If node  $i \in S'$ , the sum of elements in the  $i$ -th row of  $\Delta\mathbf{P}$  is equal to 0. By Lemma 6, we have that



(a) after step 2

(b) after step 3

(c) difference

Fig. 23. Transition probability matrices between two adjacent iterations (upper bound, step 3)

$\mathbf{r}_j^b < \bar{\mathbf{r}}_d^t$  for any node  $j \in S^t \setminus S^{t-1}$ . Thus, we have that  $\mathbf{e}_i' > 0$ . In conclusion, we have that

$$\begin{cases} \mathbf{e}_i' = 0, & \text{if } i \in (S^t \setminus S') \cup \{d\}, \\ \mathbf{e}_i' > 0, & \text{if } i \in S'. \end{cases}$$

We can get that  $\Delta\mathbf{r} = (\mathbf{I} - c\mathbf{P}^t)^{-1}\mathbf{e}'$ . The elements of  $c\mathbf{P}^t$  are non-negative and  $\|c\mathbf{P}^t\|_\infty < 1$ . We have that  $(\mathbf{I} - c\mathbf{P}^t)^{-1} = \sum_{l=0}^{\infty} c^l (\mathbf{P}^t)^l$ , where  $(\mathbf{P}^t)^l$  represents the matrix  $\mathbf{P}^t$  to the power of  $l$ . We have that  $\Delta\mathbf{r} = \sum_{l=0}^{\infty} c^l (\mathbf{P}^t)^l \mathbf{e}'$ . Since all the elements in matrix  $\mathbf{P}^t$  are non-negative, we have that  $[(\mathbf{P}^t)^l]_{i,j} \geq 0$  for any  $l$  and any nodes  $i, j \in S^t \cup \{d\}$ .

Consider a node  $i \in S^t$ . If  $i \rightsquigarrow S'$ ,  $[(\mathbf{P}^t)^l]_{i,j} > 0$  for some  $l$  and some node  $j \in S'$ . Thus,  $\Delta\mathbf{r}_i > 0$  and  $\mathbf{r}_i^b > \bar{\mathbf{r}}_i^t$ . If  $i \rightsquigarrow S'$ ,  $[(\mathbf{P}^t)^l]_{i,j} = 0$  for any  $l$  and any  $j \in S'$ . Thus,  $\Delta\mathbf{r}_i = 0$  and  $\mathbf{r}_i^b = \bar{\mathbf{r}}_i^t$ . Since each node in  $S'$  is adjacent to at least one node in  $S^t \setminus S^{t-1}$  and the graph is undirected,  $i \rightsquigarrow S'$  if and only if  $i \rightsquigarrow S^t \setminus S^{t-1}$ . This completes the proof.  $\square$

**Theorem 8.** For any node  $i \in \delta S^t$ , we have that  $\bar{\mathbf{r}}_d^t > \bar{\mathbf{r}}_i^t$ .

*Proof:* The theorem is proved by mathematical induction on the iteration. At the first iteration, for any node  $i \in \delta S^1$ , we have that  $\bar{\mathbf{r}}_i^1 < 1$  and  $\bar{\mathbf{r}}_d^1 = 1$ . So, the theorem is trivially satisfied.

Suppose that the theorem is satisfied at the  $(t-1)$ -th iteration with  $t \geq 2$ . That is,  $\bar{\mathbf{r}}_d^{t-1} > \bar{\mathbf{r}}_i^{t-1}$ , for any node  $i \in \delta S^{t-1}$ . Next, we will prove that this theorem is still satisfied at the  $t$ -th iteration.

Consider a node  $i \in \delta S^{t-1}$ . Note that  $i \rightsquigarrow d$  in the transition graph corresponding to  $\mathbf{P}^{t-1}$ . By Lemma 5,  $\bar{\mathbf{r}}_i^{t-1} > \mathbf{r}_i^a$ . By Lemma 6,  $\mathbf{r}_i^a = \mathbf{r}_i^b$ . By Lemma 7,  $\mathbf{r}_i^b \geq \bar{\mathbf{r}}_i^t$ . Thus,  $\bar{\mathbf{r}}_i^{t-1} > \bar{\mathbf{r}}_i^t$ . Since  $\bar{\mathbf{r}}_d^t = \max_{i \in \delta S^{t-1}} \bar{\mathbf{r}}_i^{t-1}$ , we have that  $\bar{\mathbf{r}}_d^t > \bar{\mathbf{r}}_i^t$ .

Consider a node  $i \in S^t \setminus S^{t-1}$ . By Lemma 6,  $\bar{\mathbf{r}}_d^t > \mathbf{r}_i^b$ . By Lemma 7,  $\mathbf{r}_i^b > \bar{\mathbf{r}}_i^t$ . Thus, we have that  $\bar{\mathbf{r}}_d^t > \bar{\mathbf{r}}_i^t$ .

Since  $\delta S^t \subset \delta S^{t-1} \cup (S^t \setminus S^{t-1})$ , we have that  $\bar{\mathbf{r}}_d^t > \bar{\mathbf{r}}_i^t$  for any node  $i \in \delta S^t$ . This completes the proof.  $\square$

The proof of Theorem 7 is as follows.

*Proof:* For any  $t$  ( $t \geq 1$ ) and any node  $i \in \delta S^t$ , we have that  $\bar{\mathbf{r}}_d^t > \bar{\mathbf{r}}_i^t$  by Theorem 8. Thus, the assumption for Lemmas 5, 6, and 7 is satisfied. Consider a node  $i \in S^{t-1}$ . Suppose that  $i \rightsquigarrow d$ . By Lemma 5,  $\bar{\mathbf{r}}_i^{t-1} > \mathbf{r}_i^a$ . By Lemma 6,  $\mathbf{r}_i^a = \mathbf{r}_i^b$ . By Lemma 7,  $\mathbf{r}_i^b \geq \bar{\mathbf{r}}_i^t$ . Thus, we have that  $\bar{\mathbf{r}}_i^{t-1} > \bar{\mathbf{r}}_i^t$ . Suppose that  $i \rightsquigarrow d$ . By Lemma 5,  $\bar{\mathbf{r}}_i^{t-1} = \mathbf{r}_i^a$ . By Lemma 6,  $\mathbf{r}_i^a = \mathbf{r}_i^b$ . By Lemma 7,  $\mathbf{r}_i^b = \bar{\mathbf{r}}_i^t$ . Thus, we have that  $\bar{\mathbf{r}}_i^{t-1} = \bar{\mathbf{r}}_i^t$ .  $\square$

## APPENDIX E

### TIGHTENING THE BOUNDS

The bounds used in FLoS can be further tightened by adding self-loop transition probabilities to the nodes in  $\delta S$ . We will still use PHP as an example to illustrate the process. We first define the star-to-mesh transformation on the transition graph, which is inspired by the star-mesh transformation in circuit theory [30].

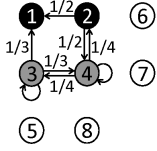
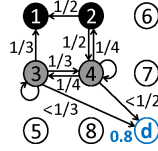
(a) adding self-loop for LB  $\underline{r}$ (b) adding self-loop for UB  $\bar{r}$ 

Fig. 24. Transition graphs with self-loops

**Definition 4.** [Star-to-mesh transformation] (1) Delete a node  $u \in V \setminus \{q\}$  and its incident transition probabilities; (2) For any pair of nodes  $i, j \in N_u$ , add transition probabilities  $p'_{i,j} = cp_{i,u}p_{u,j}$ .

Note that if  $i = j$ , it becomes the self-loop transition probability  $p'_{i,i} = cp_{i,u}p_{u,i}$ . Applying the star-to-mesh transformation for a node will not change the PHP proximity values of the remaining nodes.

**Lemma 8.** Applying the star-to-mesh transformation of node  $u$  will not change the PHP proximity values of nodes  $V \setminus \{q, u\}$ .

*Proof:* This can be proved by plugging the equation  $\mathbf{r}_u = c \sum_{i \in N_u} p_{u,i} \mathbf{r}_i$  into the recursive equations of neighbor nodes  $i \in N_u$  in the original transition graph.  $\square$

The self-loop transition probabilities generated in the star-to-mesh transformation can be used to further tighten the lower and upper bounds. Lemmas 9 and 10 analyze the tighter lower and upper bounds respectively. Figure 24 shows the transition graphs with self-loop transition probabilities for computing the lower and upper bounds. The original ones are shown in Figure 3.

**Lemma 9.** Adding self-loop transition probability  $p_{i,i} = c \sum_{j \in N_i \cap \delta \bar{S}} p_{i,j} p_{j,i}$  ( $\forall i \in \delta \bar{S}$ ) will tighten the lower bound.

*Proof:* First, we show the new bound values are still lower bounds. For the nodes in  $\delta \bar{S}$ , we apply star-to-mesh transformation sequentially in any order. After the star-to-mesh transformation of one node  $j \in \delta \bar{S}$ , we delete all the newly added transition probabilities except the self-loop transition probabilities of nodes in  $\delta \bar{S}$ . After all the nodes in  $\delta \bar{S}$  have been deleted, the self-loop transition probability of a node  $i \in \delta \bar{S}$  is  $p_{i,i} = c \sum_{j \in N_i \cap \delta \bar{S}} p_{i,j} p_{j,i}$ . During this process, we only apply star-to-mesh transformation and transition probability deletion. Therefore, the new bound values are still lower bounds.

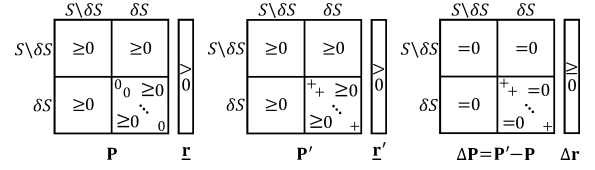
Next, we prove that the new lower bound is tighter.

Let  $\mathbf{P}$  be the transition probability matrix for computing the original lower bound vector  $\underline{r}$ . Since we add some self-loop transition probabilities to the nodes in  $\delta \bar{S}$ , we use  $\mathbf{P}'$  to denote the new transition probability matrix and  $\underline{r}'$  to denote the corresponding lower bound vector. Figures 25(a) and 25(b) illustrate the matrices and vectors. We have the following two equations

$$\begin{cases} \underline{r} = c\mathbf{P}\underline{r} + \mathbf{e}, \\ \underline{r}' = c\mathbf{P}'\underline{r}' + \mathbf{e}. \end{cases}$$

Let  $\Delta \mathbf{r} = \underline{r}' - \underline{r}$  and  $\Delta \mathbf{P} = \mathbf{P}' - \mathbf{P}$ . Note that  $\Delta \mathbf{P}_{i,i} > 0$  if  $i \in \delta \bar{S}$  and  $\Delta \mathbf{P}_{i,j} = 0$  otherwise. Figure 25(c) illustrates the matrix  $\Delta \mathbf{P}$  and vector  $\Delta \mathbf{r}$ . We have that  $\Delta \mathbf{r} = c\mathbf{P}'\underline{r}' - c\mathbf{P}\underline{r} = c\mathbf{P}'\Delta \mathbf{r} + c\Delta \mathbf{P}\underline{r} = c\mathbf{P}'\Delta \mathbf{r} + \mathbf{e}'$ , where  $\mathbf{e}' = c\Delta \mathbf{P}\underline{r}$ . We have that

$$\begin{cases} \mathbf{e}'_i = 0, & \text{if } i \in S \setminus \delta \bar{S}, \\ \mathbf{e}'_i > 0, & \text{if } i \in \delta \bar{S}. \end{cases}$$



(a) original

(b) with self-loop

(c) difference

Fig. 25. Transition probability matrices (lower bound)

We can get that  $\Delta \mathbf{r} = (\mathbf{I} - c\mathbf{P}')^{-1} \mathbf{e}'$ . The elements of  $c\mathbf{P}'$  are non-negative and  $\|c\mathbf{P}'\|_\infty < 1$ . We have that  $(\mathbf{I} - c\mathbf{P}')^{-1} = \sum_{l=0}^{\infty} c^l (\mathbf{P}')^l$ , where  $(\mathbf{P}')^l$  represents the matrix  $\mathbf{P}'$  to the power of  $l$ . Each element  $[(\mathbf{P}')^l]_{i,j}$  represents the probability of transitioning from node  $i$  to  $j$  in the  $l$ -th step.

Consider a node  $i \in S$ . If  $i \rightsquigarrow \delta \bar{S}$  in the transition graph corresponding to  $\mathbf{P}'$ ,  $[(\mathbf{P}')^l]_{i,j} > 0$  for some  $l$  and some node  $j \in \delta \bar{S}$ . So,  $[(\mathbf{P}')^l \mathbf{e}']_i > 0$  and  $\Delta \mathbf{r}_i > 0$ . If  $i \not\rightsquigarrow \delta \bar{S}$  in the transition graph corresponding to  $\mathbf{P}'$ ,  $[(\mathbf{P}')^l]_{i,j} = 0$  for any  $l$  and any node  $j \in \delta \bar{S}$ . So,  $[(\mathbf{P}')^l \mathbf{e}']_i = 0$  and  $\Delta \mathbf{r}_i = 0$ . In conclusion, we have that  $\underline{r}_i \leq \underline{r}'_i$  for any node  $i \in S$ . This completes the proof.  $\square$

**Lemma 10.** Adding self-loop transition probability  $p_{i,i} = c \sum_{j \in N_i \cap \delta \bar{S}} p_{i,j} p_{j,i}$  and setting the transition probability to dummy node as  $p_{i,d} = c \sum_{j \in N_i \cap \delta \bar{S}} p_{i,j} (1 - p_{j,i})$  ( $\forall i \in \delta \bar{S}$ ) will tighten the upper bound.

*Proof:* First, we show the new bound values are still upper bounds. For the nodes in  $\delta \bar{S}$ , we apply star-to-mesh transformation sequentially in any order. After the star-to-mesh transformation of one node  $j \in \delta \bar{S}$ , we change the destination of all the newly added transition probabilities except the self-loop transition probabilities of nodes in  $\delta \bar{S}$  to the dummy node  $d$ . After all the nodes in  $\delta \bar{S}$  have been deleted, the self-loop transition probability of a node  $i \in \delta \bar{S}$  is  $p_{i,i} = c \sum_{j \in N_i \cap \delta \bar{S}} p_{i,j} p_{j,i}$ , and the transition probability from a node  $i \in \delta \bar{S}$  to the dummy node is  $p_{i,d} = c \sum_{j \in N_i \cap \delta \bar{S}} p_{i,j} (1 - p_{j,i})$ . During this process, we only apply star-to-mesh transformation and change the destination of transition probabilities from a node  $i \in \delta \bar{S}$  to the dummy node. The proximity value of the dummy node is set as the maximum upper bound value of the nodes in the boundary node, i.e.,  $\bar{r}_d = \max_{i \in \delta \bar{S}} \bar{r}_i$ . Thus, we have that  $\bar{r}_d > \bar{r}_i$  for any node  $i \in \delta \bar{S}$ . Therefore, the new bound values are still upper bounds.

Next, we prove that the new upper bound is tighter.

Let  $\mathbf{P}$  be the transition probability matrix for computing the original upper bound vector  $\bar{r}$ . Let  $\mathbf{P}'$  denote the new transition probability matrix and  $\bar{r}'$  denote the corresponding upper bound vector. Figures 26(a) and 26(b) illustrate the matrices and vectors. We have the following two equations

$$\begin{cases} \bar{r} = c\mathbf{P}\bar{r} + \mathbf{e}', \\ \bar{r}' = c\mathbf{P}'\bar{r}' + \mathbf{e}', \end{cases}$$

where  $\mathbf{e}'_q = 1$ ,  $\mathbf{e}'_d = \bar{r}_d$ , and  $\mathbf{e}'_i = 0$  if  $i \in S \setminus \{q\}$ .

Let  $\Delta \mathbf{r} = \bar{r}' - \bar{r}$  and  $\Delta \mathbf{P} = \mathbf{P}' - \mathbf{P}$ . Comparing  $\mathbf{P}$  and  $\mathbf{P}'$ , we add some self-loop transition probabilities to the nodes in  $\delta \bar{S}$  and change the destination of transition probabilities from node  $i$  ( $i \in \delta \bar{S}$ ) to the dummy node  $d$ . Note that  $\Delta \mathbf{P}_{i,i} = p_{i,i}$  if  $i \in \delta \bar{S}$ ,  $\Delta \mathbf{P}_{i,d} < 0$  if  $i \in \delta \bar{S}$  and  $\Delta \mathbf{P}_{i,j} = 0$  otherwise. Moreover,  $-\Delta \mathbf{P}_{i,d} = \sum_{j \in N_i \cap \delta \bar{S}} p_{i,j} - c \sum_{j \in N_i \cap \delta \bar{S}} p_{i,j} (1 - p_{j,i}) > p_{i,i}$ . Thus, for any node  $i \in \delta \bar{S}$ ,  $\Delta \mathbf{P}_{i,i} + \Delta \mathbf{P}_{i,d} < 0$ . Figure 26(c) illustrates the matrix  $\Delta \mathbf{P}$  and vector  $\Delta \mathbf{r}$ .



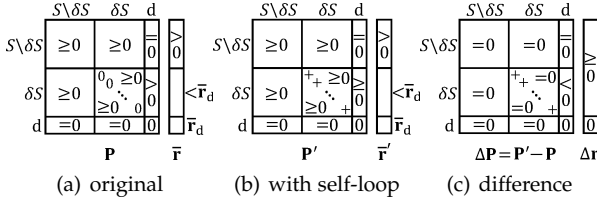


Fig. 26. Transition probability matrices (upper bound)

We have that  $\Delta \mathbf{r} = c\mathbf{P}\bar{\mathbf{r}} - c\mathbf{P}'\bar{\mathbf{r}}' = c\mathbf{P}'\Delta \mathbf{r} - c\Delta \mathbf{P}\bar{\mathbf{r}} = c\mathbf{P}'\Delta \mathbf{r} + \mathbf{e}''$ , where  $\mathbf{e}'' = -c\Delta \mathbf{P}\bar{\mathbf{r}}$ . Based on Theorem 8, we have that for any node  $i \in \delta S^t$ ,  $\bar{\mathbf{r}}_d^t > \bar{\mathbf{r}}_i^t$ . We get that

$$\begin{cases} \mathbf{e}''_i = 0, & \text{if } i \in (S \setminus \delta S) \cup \{d\}, \\ \mathbf{e}''_i > 0, & \text{if } i \in \delta S. \end{cases}$$

We can get that  $\Delta \mathbf{r} = (\mathbf{I} - c\mathbf{P}')^{-1}\mathbf{e}''$ . The elements of  $c\mathbf{P}'$  are non-negative and  $\|c\mathbf{P}'\|_\infty < 1$ . We have that  $(\mathbf{I} - c\mathbf{P}')^{-1} = \sum_{l=0}^{\infty} c^l (\mathbf{P}')^l$ , where  $(\mathbf{P}')^l$  represents the matrix  $\mathbf{P}'$  to the power of  $l$ . Each element  $[(\mathbf{P}')^l]_{i,j}$  represents the probability of transiting from node  $i$  to  $j$  in the  $l$ -th step.

Consider a node  $i \in S$ . If  $i \rightsquigarrow \delta S$  in the transition graph corresponding to  $\mathbf{P}'$ ,  $[(\mathbf{P}')^l]_{i,j} > 0$  for some  $l$  and some node  $j \in \delta S$ . So,  $[(\mathbf{P}')^l \mathbf{e}'']_i > 0$  and  $\Delta \mathbf{r}_i > 0$ . If  $i \not\rightsquigarrow \delta S$  in the transition graph corresponding to  $\mathbf{P}'$ ,  $[(\mathbf{P}')^l]_{i,j} = 0$  for any  $l$  and any node  $j \in \delta S$ . So,  $[(\mathbf{P}')^l \mathbf{e}'']_i = 0$  and  $\Delta \mathbf{r}_i = 0$ . In conclusion, we have that  $\bar{\mathbf{r}}_i \geq \bar{\mathbf{r}}'_i$  for any node  $i \in S$ . This completes the proof.  $\square$

## APPENDIX F

### ANALYSIS ON THE NUMBER OF VISITED NODES

In this subsection, we analyze the number of visited nodes. Let  $h$  be the average number of neighbors of a node. Suppose that the nodes in the boundary of visited nodes are  $\rho$  hops away from the query node  $q$ . We assume that the number of visited nodes equals  $h^\rho$ . Note that this is an upper bound of the number of visited nodes.

The proximity of one node will decrease by a factor of  $c$  when it is one hop farther away from the query node. The nodes in the boundary of the visited nodes have proximity less than  $c^\rho$  because they are  $\rho$  hops away from the query.

What is the distribution of the gaps between the upper and lower bounds? The upper bound  $\bar{\mathbf{r}}$  is computed based on the recursive equation  $\bar{\mathbf{r}} = c\mathbf{P}\bar{\mathbf{r}} + \mathbf{e}''$ , where vector  $\mathbf{e}''$  has only two non-zero elements  $\mathbf{e}''_q = 1$  and  $\mathbf{e}''_d = \mathbf{r}_d$ . When we set the proximity value of dummy node to 0, i.e.,  $\mathbf{e}''_d = 0$ , we will get the recursive equation  $\underline{\mathbf{r}} = c\mathbf{P}\underline{\mathbf{r}} + \mathbf{e}$  for the lower bound  $\underline{\mathbf{r}}$ . Let  $\mathbf{r}' = \bar{\mathbf{r}} - \underline{\mathbf{r}}$  be the gaps between the upper and lower bounds. We have that  $\mathbf{r}' = c\mathbf{P}\mathbf{r}' + \mathbf{e}'$ , where vector  $\mathbf{e}'$  has only one non-zero element  $\mathbf{e}'_d = \mathbf{r}_d$ . Thus the gaps  $\mathbf{r}'$  can be interpreted as the PHP proximity values when the query node is the dummy node  $d$  with constant proximity value  $\mathbf{r}_d$ .

Based on this observation, the gap  $\mathbf{r}'_u$  of one node  $u$  will decrease by a factor of  $c$  when it is one hop farther away from the boundary. Let  $\rho_u$  denote the number of hops that  $u$  is away from the query node. Node  $u$  is  $\rho - \rho_u$  hops away from the boundary. So, the gap  $\mathbf{r}'_u$  is less than  $c^{\rho - \rho_u} \cdot c^\rho$ , i.e.,  $\mathbf{r}'_u < c^{2\rho - \rho_u}$ .

For any two nodes  $u$  and  $v$ , we can distinguish their rankings if  $\mathbf{r}'_u + \mathbf{r}'_v < \epsilon$ , where  $\epsilon$  is the difference of the exact proximity values of nodes  $u$  and  $v$ . We already derive

that  $\mathbf{r}'_u < c^{2\rho - \rho_u}$  and  $\mathbf{r}'_v < c^{2\rho - \rho_v}$ . Thus if we have that  $\rho > \frac{1}{2} \log_c \epsilon - \frac{1}{2} \log_c (c^{-\rho_u} + c^{-\rho_v})$ , then the pair of nodes  $u$  and  $v$  can be distinguished.

Now we consider the case when  $u$  and  $v$  are the  $k$ -th and  $(k+1)$ -th nodes in the exact ranking list respectively. We have  $h^{\rho_u} \geq k$  and  $h^{\rho_v} \geq k+1$ . Thus, the inequality  $\rho > \frac{1}{2} \log_c \frac{\epsilon}{2} + \frac{1}{2} \log_h k$  should be satisfied. Therefore, we need to visit  $h^\rho = O((kh^{\log_c(\epsilon/2)})^{\frac{1}{2}})$  nodes to distinguish the  $k$ -th and  $(k+1)$ -th nodes.

## APPENDIX G

### EXTENSIONS OF FLOS TO THE PROXIMITY MEASURES HAVING LOCAL MAXIMUM

This subsection shows the details about how to extend the FLoS method to random walk with restart, RoundTripRank, Katz score, and absorption probability.

#### G.1 Extension to Random Walk with Restart

In this subsection, we discuss how to apply the FLoS method to random walk with restart (RWR) [8] which has local optimum. The key idea is to utilize the relationship between RWR and PHP. Utilizing such relationship, we can easily derive the lower and upper bounds for RWR based on the lower and upper bounds for PHP.

Random walk with restart (also known as personalized PageRank) [8] is a widely used proximity measure. RWR can be described as follows. From a node  $i$ , the random walker can walk to its neighbors with probabilities proportional to the edge weights. In each step, it has a probability of  $(1-c)$  to return to the query node  $q$ , where  $c$  is a constant. The proximity of node  $i$  w.r.t.  $q$  is defined as the stationary probability that the random walker will finally stay at  $i$ . RWR can be defined recursively as

$$\mathbf{r}_i = \begin{cases} c \sum_{j \in N_i} p_{j,i} \mathbf{r}_j + (1-c), & \text{if } i = q, \\ c \sum_{j \in N_i} p_{j,i} \mathbf{r}_j, & \text{if } i \neq q, \end{cases}$$

where  $c$  ( $0 < c < 1$ ) is a constant decay factor, and  $p_{j,i} = \frac{w_{j,i}}{w_j}$  is the transition probability from node  $j$  to  $i$ .

**Lemma 11.** RWR has local maximum.

*Proof:* Examples can be constructed to show that RWR has local maximum, which are omitted.  $\square$

Let  $\text{RWR}(i)$  and  $\text{PHP}(i)$  represent the RWR and PHP proximity values of node  $i$  respectively. We first show that RWR and PHP have the following relationship.

**Theorem 9.**  $\text{RWR}(i) \propto w_i \cdot \text{PHP}(i)$

*Proof:* Based on the recursive definition of RWR, we have  $\frac{\text{RWR}(i)}{w_i} = c \sum_{j \in N_i} \frac{w_{j,i}}{w_j} \cdot \frac{\text{RWR}(j)}{w_j} = c \sum_{j \in N_i} \frac{w_{j,i}}{w_i} \cdot \frac{\text{RWR}(j)}{w_j}$ , for any node  $i \neq q$ . Thus, for any node  $i \neq q$ , we have that  $\frac{\text{RWR}(i)}{w_i} = c \sum_{j \in N_i} p_{j,i} \cdot \frac{\text{RWR}(j)}{w_j}$ . This degree normalized RWR,  $\frac{\text{RWR}(i)}{w_i}$ , has the same recursive equation as PHP with decay factor  $c$ . So we have that  $\frac{\text{RWR}(i)}{w_i \cdot \text{PHP}(i)} = \frac{\text{RWR}(q)}{w_q \cdot \text{PHP}(q)}$ . When the query node  $q$  is fixed, we have that  $\text{RWR}(i) \propto w_i \cdot \text{PHP}(i)$ .  $\square$

How to extend FLoS to RWR by using this relationship is already discussed in Section 6.

Here, we provide more details about how to maintain the maximum degree of unvisited nodes. In the implementation of the algorithm, we use the vector container in C++. The

program can access any element in the vector in  $O(1)$  time, and can delete the last element in  $O(1)$  time. Each element in the vector is of structure type and contains two members, one represents the node degree and the other represents the node color. If a node is visited, its node color is black; otherwise, its node color is white.

When we pre-process the graph, we sort the elements in the vector by the node degree in ascending order, and set the initial node color of each element to white.

The last element always contains the node with the maximum degree of unvisited nodes. When we newly visit a node, we set its node color to black. If this newly visited node corresponds to the last element in the vector, we remove this element from the vector. And then we start removing the elements from the end of the vector until the last element's node color is white.

To maintain the maximum degree of unvisited nodes, the amortized time complexity is  $O(1)$  for each newly visited node.

## G.2 Extension to RoundTripRank

RoundTripRank (RT) [10] considers round trip paths between nodes. RT of node  $i$  is defined as the probability that a round trip from  $q$  contains node  $i$ . The round trip can be decoupled into two types of trips, i.e., the forward trip from node  $q$  to  $i$  and the backward trip from node  $i$  to  $q$ . Let  $\mathbf{r}_i$  denote the RT proximity value of node  $i$  with regard to the query node  $q$ . Let  $\mathbf{x}_i$  denote the proximity value defined in the forward direction from  $q$  to  $i$ , and  $\mathbf{y}_i$  denote the proximity value defined in the backward direction from  $i$  to  $q$ .  $\mathbf{x}_i$  and  $\mathbf{y}_i$  can be defined as

$$\mathbf{x}_i = \begin{cases} c \sum_{j \in N_i} p_{j,i} \mathbf{x}_j + (1-c), & \text{if } i=q, \\ c \sum_{j \in N_i} p_{j,i} \mathbf{x}_j, & \text{if } i \neq q, \end{cases}$$

$$\mathbf{y}_i = \begin{cases} c \sum_{j \in N_i} p_{i,j} \mathbf{y}_j + (1-c), & \text{if } i=q, \\ c \sum_{j \in N_i} p_{i,j} \mathbf{y}_j, & \text{if } i \neq q, \end{cases}$$

where  $c$  ( $0 < c < 1$ ) is a constant decay factor, and  $p_{i,j} = \frac{w_{i,j}}{w_i}$  is the transition probability from node  $i$  to  $j$ .

RT can be decomposed as the multiplication of these two values, i.e.,

$$\mathbf{r}_i \propto \mathbf{x}_i^\beta \cdot \mathbf{y}_i^{1-\beta},$$

where  $\beta$  ( $0 \leq \beta \leq 1$ ) is a constant and used to tune the trade-off between the forward and backward directions.

The relationship between RT and PHP is shown in the following theorem.

**Theorem 10.**  $\text{RT}(i) \propto w_i^\beta \cdot \text{PHP}(i)$

*Proof:* We can see that  $\mathbf{x}_i$  is exactly the RWR proximity value, and  $\mathbf{y}_i$  has a linear relationship with the EI proximity value, i.e.,  $\mathbf{y}_i = w_q \cdot \text{EI}(i)$ . Thus, we have that  $\text{RT}(i) \propto \text{RWR}(i)^\beta \cdot \text{EI}(i)^{1-\beta}$ . From Theorem 2, we have that  $\text{EI}(i) \propto \text{PHP}(i)$ . From Theorem 9, we have that  $\text{RWR}(i) \propto w_i \cdot \text{PHP}(i)$ . Thus,  $\text{RT}(i) \propto w_i^\beta \cdot \text{PHP}(i)$ .  $\square$

**Lemma 12.** RT has no local maximum when  $\beta=0$ , and has local maximum when  $0 < \beta \leq 1$ .

*Proof:* Based on Theorem 10, when  $\beta=0$ , we have that  $\text{RT}(i) \propto \text{PHP}(i)$ . Since PHP has no local maximum, RT has no local maximum. Examples can be constructed to show that

RT has local maximum when  $0 < \beta \leq 1$ , which are omitted.  $\square$

Suppose that node  $v \in \delta S^t$  has the largest PHP proximity value. Based on Theorem 1, for any  $i \in \bar{S}^t$ , we have that  $\text{PHP}(i) \leq \text{PHP}(v)$ . Let  $w(\bar{S}^t)$  denote the maximum degree of unvisited nodes in  $\bar{S}^t$ . We have that  $w_i^\beta \cdot \text{PHP}(i) \leq w(\bar{S}^t)^\beta \cdot \text{PHP}(i) \leq w(\bar{S}^t)^\beta \cdot \text{PHP}(v)$ . Therefore, if we maintain the maximum degree of unvisited nodes, we can develop the upper bound for the proximity values of unvisited nodes.

Specifically, we can apply FLoS to RT as follows. In Algorithm 3, we can change line 1 to the following line.

1:  $u \leftarrow \arg\max_{i \in \delta S^{t-1}} w_i^\beta \cdot (\mathbf{r}_i^{t-1} + \bar{\mathbf{r}}_i^{t-1})$ ;

In Algorithm 6, we can change line 2 and 3 to the following two lines.

2:  $K \leftarrow k$  nodes in  $S^t \setminus (\delta S^t \cup \{q\})$  with largest  $w_i^\beta \cdot \mathbf{r}_i^t$  values;  
3: **if**  $\min_{i \in K} w_i^\beta \cdot \mathbf{r}_i^t \geq \max_{i \in S^t \setminus (K \cup \{q\})} w_i^\beta \cdot \bar{\mathbf{r}}_i^t$  **and**  $\min_{i \in K} w_i^\beta \cdot \mathbf{r}_i^t \geq w(\bar{S}^t)^\beta \cdot \max_{i \in \delta S^t} \bar{\mathbf{r}}_i^t$  **then**  $\text{bStop} \leftarrow \text{true}$ ;

All other processes remain the same.

## G.3 Extension to the Katz Score

The Katz score (KZ) [18] measures the proximity between nodes via a weighted sum of the length of paths between them. Let  $\mathbf{R}$  denote the matrix containing all pairwise KZ proximity values (with  $\mathbf{R}_{i,j}$  denoting the KZ proximity value between nodes  $i$  and  $j$ ). Let  $\mathbf{W}$  be the symmetric adjacency matrix of an undirected and connected graph. We have that

$$\mathbf{R} = \kappa \mathbf{W} + \kappa^2 \mathbf{W}^2 + \dots = (\mathbf{I} - \kappa \mathbf{W})^{-1} - \mathbf{I},$$

where  $\kappa$  ( $0 < \kappa < 1$ ) is the decay factor.

In general, KZ has local maximum. But when  $\kappa$  is small, it does not have local maximum. To prove this, we first define a new proximity measure  $\text{PHP}'$  as follows

$$\mathbf{r}_i = \begin{cases} 1, & \text{if } i=q, \\ c \sum_{j \in N_i} p_{i,j} \mathbf{r}_j, & \text{if } i \neq q, \end{cases}$$

where  $c$  ( $0 < c < 1$ ) is the decay factor in the random walk process,  $p_{i,j} = w_{i,j}/w_{\max}$  denotes the transition probability from node  $i$  to  $j$ , and  $w_{\max}$  denotes the maximum degree among all the nodes. Different from PHP, the transition probability in  $\text{PHP}'$  is normalized by the maximum degree.

**Theorem 11.** If  $\kappa < 1/w_{\max}$ ,  $\text{KZ}(i) \propto \text{PHP}'(i)$ .

*Proof:* For the KZ proximity measure, we already have that  $\mathbf{R} = (\mathbf{I} - \kappa \mathbf{W})^{-1} - \mathbf{I}$ . The elements in the  $q$ -th column of  $\mathbf{R}$  represent the KZ proximity values when the query node is  $q$ . Thus, the KZ proximity vector is  $\mathbf{r} = (\mathbf{I} - \kappa \mathbf{W})^{-1} \mathbf{e} - \mathbf{e}$ . Then, we have that

$$\mathbf{r} + \mathbf{e} = \kappa \mathbf{W}(\mathbf{r} + \mathbf{e}) + \mathbf{e}.$$

To derive the relationship between KZ and  $\text{PHP}'$ , we define a new proximity measure  $\text{KZ}'$  as

$$\text{KZ}'(i) = \begin{cases} \frac{\text{KZ}(i)+1}{w_{\max}}, & \text{if } i=q, \\ \frac{\text{KZ}(i)}{w_{\max}}, & \text{if } i \neq q. \end{cases}$$

Let  $\mathbf{r}'$  denote the proximity vector of  $\text{KZ}'$  when the query node is  $q$ . The relationship between KZ and  $\text{KZ}'$  can be expressed in a matrix form as

$$\mathbf{r}' = (\mathbf{r} + \mathbf{e})/w_{\max}.$$

Then, the proximity vector of  $\text{KZ}'$  satisfies that

$$\mathbf{r}' = \kappa \mathbf{W} \mathbf{r}' + \mathbf{e}/w_{\max}.$$

Based on this matrix form,  $\mathbf{KZ}'$  has the following recursive definition

$$\mathbf{r}'_i = \begin{cases} \kappa \cdot w_{\max} \sum_{j \in N_i} p_{i,j} \mathbf{r}'_j + \frac{1}{w_{\max}}, & \text{if } i = q; \\ \kappa \cdot w_{\max} \sum_{j \in N_i} p_{i,j} \mathbf{r}'_j, & \text{if } i \neq q. \end{cases}$$

Since  $\kappa < 1/w_{\max}$ , we have that  $\kappa \cdot w_{\max} < 1$ . If we set the decay factor in  $\text{PHP}'$  as  $c = \kappa \cdot w_{\max}$ ,  $\mathbf{KZ}'$  and  $\text{PHP}'$  have the same recursive definition for any node  $i \neq q$ , and we have that  $\frac{\text{PHP}'(i)}{\mathbf{KZ}'(i)} = \frac{\text{PHP}'(q)}{\mathbf{KZ}'(q)}$ . Thus,  $\frac{\text{PHP}'(i)}{\mathbf{KZ}(i)} = \frac{\text{PHP}'(q)}{\mathbf{KZ}(q)+1}$ . When the query node  $q$  is fixed, we have that  $\mathbf{KZ}(i) \propto \text{PHP}'(i)$ .  $\square$

**Lemma 13.**  $\mathbf{KZ}$  does not have local maximum when  $\kappa < 1/w_{\max}$ , and has local maximum when  $\kappa \geq 1/w_{\max}$ .

*Proof:* Based on Theorem 11, when  $\kappa < 1/w_{\max}$ , we have that  $\mathbf{KZ}(i) \propto \text{PHP}'(i)$ . We can prove that  $\text{PHP}'$  has no local maximum. The proof is similar to that in Lemma 1. Thus  $\mathbf{KZ}$  has no local maximum. Examples can be constructed to show that  $\mathbf{KZ}$  has local maximum when  $\kappa \geq 1/w_{\max}$ , which are omitted.  $\square$

Theorem 11 says that  $\mathbf{KZ}(i)$  is proportional to  $\text{PHP}'(i)$  when  $\kappa$  is small. Therefore, ranking by  $\mathbf{KZ}$  is equivalent to ranking by  $\text{PHP}'$ . Comparing with  $\text{PHP}$ ,  $\text{PHP}'$  only has different transition probabilities, which are normalized by the maximum degree instead of the degree of each node as in  $\text{PHP}$ . Therefore, the FLoS method is readily applicable to  $\text{PHP}'$ .

#### G.4 Extension to Absorption Probability

In the absorption probability (AP) [11], with probability  $p_{i,i}$ , the random walker will be absorbed at node  $i$ , and with probability  $1 - p_{i,i}$ , the random walker will follow a random edge out of it. Let  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$  be a diagonal matrix, where  $\lambda_i > 0$  is a constant for any node  $i = 1, \dots, n$ . The transition probability in AP is defined as

$$p_{i,j} = \begin{cases} \frac{\lambda_i}{\lambda_i + w_i}, & \text{if } i = j, \\ \frac{w_{i,j}}{\lambda_i + w_i}, & \text{if } i \neq j. \end{cases}$$

The AP proximity  $r_{i,j}$  of node  $j$  with regard to node  $i$  is defined as the probability that a random walker starting from node  $i$  is absorbed at node  $j$ . AP can be defined recursively as

$$r_{i,j} = \begin{cases} \sum_{k \in N_i} p_{i,k} r_{k,j} + p_{i,i}, & \text{if } i = j, \\ \sum_{k \in N_i} p_{i,k} r_{k,j}, & \text{if } i \neq j. \end{cases}$$

Let  $\mathbf{R}$  be the matrix of AP proximity values, i.e.,  $\mathbf{R}_{i,j} = r_{i,j}$ . Let  $\mathbf{W}$  be the adjacency matrix, and  $\mathbf{D}$  be the diagonal matrix where each element  $\mathbf{D}_{i,i}$  equals the degree  $w_i$  of node  $i$ . The above equations can be expressed in a matrix form

$$\mathbf{R} = (\mathbf{\Lambda} + \mathbf{D})^{-1} \mathbf{W} \mathbf{R} + (\mathbf{\Lambda} + \mathbf{D})^{-1} \mathbf{\Lambda}.$$

Thus, we have that

$$\mathbf{R} = (\mathbf{\Lambda} + \mathbf{D} - \mathbf{W})^{-1} \mathbf{\Lambda}.$$

$\mathbf{R}$  is a non-negative matrix with each row summing up to 1 [11]. The elements in the  $q$ -th row of  $\mathbf{R}$  represent the AP proximity values with regard to the query  $q$ .

To extend FLoS to AP, we consider a variant of  $\text{PHP}$ , which is referred to as  $\text{PHP}''$ .  $\text{PHP}''$  is defined as

$$\mathbf{r}_i = \begin{cases} 1, & \text{if } i = q, \\ \sum_{j \in N_i} p_{i,j} \mathbf{r}_j, & \text{if } i \neq q, \end{cases}$$

where  $p_{i,j}$  is the transition probability defined in AP.

Compared with  $\text{PHP}$ , the transition probability in  $\text{PHP}''$  is changed. In  $\text{PHP}$ , the transition probability is normalized by degree  $w_i$ . In  $\text{PHP}''$ , the transition probability is normalized by  $\lambda_i + w_i$ . Another difference is that there is no decay factor in  $\text{PHP}''$ . Even though we do not have decay factor in  $\text{PHP}''$ , FLoS is still applicable to  $\text{PHP}''$ . Some key observations are shown as follows.

**Lemma 14.**  $\text{PHP}''$  has no local maximum.

*Proof:* Suppose that node  $i$  is a local maximum. We have that  $\mathbf{r}_i = \sum_{j \in N_i} p_{i,j} \mathbf{r}_j \leq \sum_{j \in N_i} p_{i,j} \mathbf{r}_i = \mathbf{r}_i \cdot \frac{w_i}{\lambda_i + w_i} < \mathbf{r}_i$ . We get a contradiction that  $\mathbf{r}_i < \mathbf{r}_i$ .  $\square$

We still use  $\mathbf{P}$  to denote the transition probability matrix with

$$\mathbf{P}_{i,j} = \begin{cases} 0, & \text{if } i = q \text{ or } i = j, \\ p_{i,j}, & \text{if } i \neq q \text{ and } i \neq j. \end{cases}$$

Then the proximity  $\mathbf{r}$  based on  $\text{PHP}''$  can be written in the following matrix form

$$\mathbf{r} = \mathbf{P} \mathbf{r} + \mathbf{e}.$$

We can see that the sum of each row of  $\mathbf{P}$  is smaller than 1, i.e.,  $\sum_j \mathbf{P}_{i,j} = \sum_j p_{i,j} = \frac{w_i}{\lambda_i + w_i} < 1$ . Thus, we have that  $\|\mathbf{P}\|_\infty < 1$ . Based on this, Theorems 3, 4 and 5 are still satisfied. Thus, we can develop the lower and upper bounds in a similar way as that for  $\text{PHP}$ . Therefore, the FLoS method is applicable to  $\text{PHP}''$ .

AP and  $\text{PHP}''$  have the following relationship.

**Theorem 12.**  $\text{AP}(i) \propto \lambda_i \cdot \text{PHP}''(i)$

*Proof:* The elements in the  $q$ -th row of  $\mathbf{R}$  represent the AP proximity values when the query node is  $q$ . Thus, the AP proximity vector is  $\mathbf{r} = \mathbf{R}^T \mathbf{e} = \mathbf{\Lambda}(\mathbf{\Lambda} + \mathbf{D} - \mathbf{W})^{-1} \mathbf{e}$ , where  $\mathbf{R}^T$  represents the transpose of  $\mathbf{R}$ . Then,  $(\mathbf{\Lambda} + \mathbf{D} - \mathbf{W}) \mathbf{\Lambda}^{-1} \mathbf{r} = \mathbf{e}$ . Thus, the proximity vector of AP satisfies that

$$\mathbf{\Lambda}^{-1} \mathbf{r} = (\mathbf{\Lambda} + \mathbf{D})^{-1} \mathbf{W} \mathbf{\Lambda}^{-1} \mathbf{r} + (\mathbf{\Lambda} + \mathbf{D})^{-1} \mathbf{e}.$$

We define a new proximity measure  $\text{AP}'$  as

$$\text{AP}'(i) = \frac{\text{AP}(i)}{\lambda_i},$$

for any node  $i \in V$ .

Let  $\mathbf{r}'$  denote the proximity vector of  $\text{AP}'$  when the query is  $q$ . The relationship between AP and  $\text{AP}'$  can be expressed in a matrix form as  $\mathbf{r}' = \mathbf{\Lambda}^{-1} \mathbf{r}$ . Thus, the proximity vector of  $\text{AP}'$  satisfies the following equation

$$\mathbf{r}' = (\mathbf{\Lambda} + \mathbf{D})^{-1} \mathbf{W} \mathbf{r}' + (\mathbf{\Lambda} + \mathbf{D})^{-1} \mathbf{e}.$$

Based on this matrix form,  $\text{AP}'$  has the following recursive definition

$$\mathbf{r}'_i = \begin{cases} \sum_{j \in N_i} p_{i,j} \mathbf{r}'_j + \frac{1}{\lambda_i + w_i}, & \text{if } i = q, \\ \sum_{j \in N_i} p_{i,j} \mathbf{r}'_j, & \text{if } i \neq q. \end{cases}$$

$\text{AP}'$  and  $\text{PHP}''$  have the same recursive definition for any node  $i \neq q$ . Thus,  $\frac{\text{AP}'(i)}{\text{PHP}''(i)} = \frac{\text{AP}'(q)}{\text{PHP}''(q)}$ . Since  $\text{AP}'(i) = \frac{\text{AP}(i)}{\lambda_i}$ , we have that  $\frac{\text{AP}(i)}{\lambda_i \cdot \text{PHP}''(i)} = \frac{\text{AP}(q)}{\lambda_q \cdot \text{PHP}''(q)}$ . When  $q$  is fixed, we have that  $\text{AP}(i) \propto \lambda_i \cdot \text{PHP}''(i)$ .  $\square$

**Lemma 15.** AP has local maximum.

*Proof:* Examples can be constructed to show that AP has local maximum and are omitted here.  $\square$

Suppose that node  $v \in \delta S^t$  has the largest  $\text{PHP}''$  proximity value. Based on Theorem 1, for any node  $i \in \bar{S}^t$ ,  $\text{PHP}''(i) \leq \text{PHP}''(v)$ . Let  $\lambda(\bar{S}^t)$  denote the maximum  $\lambda$  value of unvisited nodes in  $\bar{S}^t$ . We have that  $\lambda_i \cdot \text{PHP}''(i) \leq \lambda(\bar{S}^t) \cdot \text{PHP}''(i) \leq \lambda(\bar{S}^t) \cdot \text{PHP}''(v)$ . Therefore, if we maintain the

maximum  $\lambda$  value of unvisited nodes, we can develop the upper bound for the proximity values of unvisited nodes.

Specifically, we can apply FLoS to AP as follows. In Algorithm 3, we can change line 1 to the following line.

1:  $u \leftarrow \operatorname{argmax}_{i \in \delta S^{t-1}} \lambda_i \cdot (\mathbf{r}_i^{t-1} + \bar{\mathbf{r}}_i^{t-1});$

In Algorithm 6, we can change line 2 and 3 to the following two lines.

2:  $K \leftarrow k$  nodes in  $S^t \setminus (\delta S^t \cup \{q\})$  with largest  $\lambda_i \cdot \mathbf{r}_i^t$  values;

3: **if**  $\min_{i \in K} \lambda_i \cdot \mathbf{r}_i^t \geq \max_{i \in S^t \setminus (K \cup \{q\})} \lambda_i \cdot \bar{\mathbf{r}}_i^t$  **and**  $\min_{i \in K} \lambda_i \cdot \mathbf{r}_i^t \geq \lambda(\bar{S}^t) \cdot \max_{i \in \delta S^t} \bar{\mathbf{r}}_i^t$  **then**  $\text{bStop} \leftarrow \text{true};$

All other processes remain the same.

## APPENDIX H

### EXTENSIONS TO THE REVERSE PROXIMITY MEASURES

In this subsection, we show how to apply the FLoS method to solve the top- $k$  reverse-proximity query problem. Let  $\text{RWR}_i(j)$ ,  $\text{EI}_i(j)$ ,  $\text{PHP}_i(j)$ ,  $\text{DHT}_i(j)$ ,  $\text{RT}_i(j)$ , and  $\text{AP}_i(j)$  denote the RWR, EI, PHP, DHT, RT, and AP proximity values of node  $j$  when the query is node  $i$ .

#### H.1 Extension to Reverse RWR

In this subsection, we show that ranking by the reverse RWR proximity is the same as ranking by the PHP proximity. Thus the FLoS method for PHP can be directly applied to find the top- $k$  nodes for reverse RWR.

**Theorem 13.** PHP and reverse RWR give the same ranking results.

*Proof:* From Theorem 2, we have that PHP and EI give the same ranking results. Thus, it is equivalent to prove that EI and reverse RWR give the same ranking results.

EI is defined as the degree normalized RWR, i.e.,  $\text{EI}_q(i) = \frac{\text{RWR}_q(i)}{w_i}$ . Thus, we have that  $\text{RWR}_q(i) = w_i \cdot \text{EI}_q(i)$ . Therefore, we have that  $\text{RWR}_i(q) = w_q \cdot \text{EI}_i(q) = w_q \cdot \text{EI}_q(i)$ , where we use the fact that EI is symmetric, i.e.,  $\text{EI}_i(q) = \text{EI}_q(i)$ . Thus,  $\text{RWR}_i(q) \propto \text{EI}_q(i)$ , when the query  $q$  is fixed. This completes the proof.  $\square$

#### H.2 Extension to Reverse PHP and Reverse DHT

In this subsection, we first show how to solve the top- $k$  reverse PHP problem. Then, we prove that reverse PHP and reverse DHT give the same ranking results.

To find the top- $k$  nodes for reverse PHP, we use the following relationship between PHP and reverse PHP.

**Theorem 14.**  $\text{PHP}_i(q) \propto \frac{\text{PHP}_q(i)}{\text{EI}_i(i)}$

*Proof:* From Theorem 2, we have that  $\text{PHP}_q(i) = \frac{\text{EI}_q(i)}{\text{EI}_q(q)}$ . Thus,  $\text{PHP}_i(q) = \frac{\text{EI}_i(q)}{\text{EI}_i(i)} = \frac{\text{EI}_q(i)}{\text{EI}_i(i)} \cdot \frac{\text{PHP}_q(i)}{\text{EI}_q(i)}$ , where we use the fact that EI is symmetric, i.e.,  $\text{EI}_i(q) = \text{EI}_q(i)$ . Therefore, we have that  $\text{PHP}_i(q) \propto \frac{\text{PHP}_q(i)}{\text{EI}_i(i)}$  when the query node  $q$  is fixed.  $\square$

Suppose that we already pre-compute the values  $\text{EI}_i(i)$  for each node  $i$ , and node  $v \in \delta S^t$  has the largest PHP proximity value. Based on Theorem 1, for any node  $i \in \bar{S}^t$ ,  $\text{PHP}_q(i) \leq \text{PHP}_q(v)$ . Let  $\text{EI}(\bar{S}^t)$  denote the minimum  $\text{EI}_i(i)$  value of unvisited nodes  $i$  in  $\bar{S}^t$ . We have that  $\frac{\text{PHP}_q(i)}{\text{EI}_i(i)} \leq \frac{\text{PHP}_q(v)}{\text{EI}(\bar{S}^t)} \leq \frac{\text{PHP}_q(v)}{\text{EI}(\bar{S}^t)}$ . Therefore, if we maintain the minimum  $\text{EI}_i(i)$  value

of unvisited nodes, we can develop the upper bound for the proximity values of unvisited nodes.

Specifically, we can apply FLoS to reverse PHP as follows. In Algorithm 3, we can change line 1 to the following line.

1:  $u \leftarrow \operatorname{argmax}_{i \in \delta S^{t-1}} \frac{1}{\text{EI}_i(i)} \cdot (\mathbf{r}_i^{t-1} + \bar{\mathbf{r}}_i^{t-1});$

In Algorithm 6, we can change line 2 and 3 to the following two lines.

2:  $K \leftarrow k$  nodes in  $S^t \setminus (\delta S^t \cup \{q\})$  with largest  $\frac{\mathbf{r}_i^t}{\text{EI}_i(i)}$  values;

3: **if**  $\min_{i \in K} \frac{\mathbf{r}_i^t}{\text{EI}_i(i)} \geq \max_{i \in S^t \setminus (K \cup \{q\})} \frac{\bar{\mathbf{r}}_i^t}{\text{EI}_i(i)}$  **and**  $\min_{i \in K} \frac{\mathbf{r}_i^t}{\text{EI}_i(i)} \geq (\text{EI}(\bar{S}^t))^{-1} \cdot \max_{i \in \delta S^t} \bar{\mathbf{r}}_i^t$  **then**  $\text{bStop} \leftarrow \text{true};$

All other processes remain the same.

**Theorem 15.** Reverse PHP and reverse DHT give the same ranking results.

*Proof:* From Theorem 2, DHT and PHP have the relationship  $\text{DHT}_q(i) = \frac{1}{1-c}(1 - \text{PHP}_q(i))$ . Thus,  $\text{DHT}_i(q) = \frac{1}{1-c}(1 - \text{PHP}_i(q))$ . It means that reverse DHT is a linear function of reverse PHP. Thus reverse DHT and reverse PHP give the same ranking results.  $\square$

#### H.3 Extension to Reverse RT

In this subsection, we show how to find the top- $k$  nodes for reverse RT.

**Theorem 16.**  $\text{RT}_i(q) \propto w_i^{1-\beta} \cdot \text{PHP}_q(i)$

*Proof:* Let  $\text{RTF}_i(j)$  represent the forward proximity value defined in the direction from  $i$  to  $j$ . Thus, we have that  $\text{RTF}_q(i) = \mathbf{x}_i$ , where  $\mathbf{x}_i$  is defined in Appendix G.2.  $\text{RTF}_q(i)$  is exactly the RWR proximity value, i.e.,  $\text{RTF}_q(i) = \text{RWR}_q(i)$ . Let  $\text{RTB}_i(j)$  represent the backward proximity value defined in the direction from  $j$  to  $i$ . Thus, we have that  $\text{RTB}_q(i) = \mathbf{y}_i$ , where  $\mathbf{y}_i$  is defined in Appendix G.2.  $\text{RTB}_q(i)$  has a linear relationship with the EI proximity value, i.e.,  $\text{RTB}_q(i) = w_q \cdot \text{EI}_q(i)$ .

RT can be decomposed as

$$\text{RT}_q(i) \propto \text{RTF}_q(i)^\beta \cdot \text{RTB}_q(i)^{1-\beta}.$$

Reverse RT can be decomposed as

$$\begin{aligned} \text{RT}_i(q) &\propto \text{RTF}_i(q)^\beta \cdot \text{RTB}_i(q)^{1-\beta} \\ &\propto \text{RWR}_i(q)^\beta \cdot w_i^{1-\beta} \cdot \text{EI}_i(q)^{1-\beta}. \end{aligned}$$

From Theorem 13, we have that  $\text{RWR}_i(q) \propto \text{EI}_q(i)$ . Since EI is symmetric, we have that  $\text{EI}_i(q) = \text{EI}_q(i)$ . From Theorem 2, we have that  $\text{EI}_q(i) \propto \text{PHP}_q(i)$ . Thus, we have that  $\text{RT}_i(q) \propto w_i^{1-\beta} \cdot \text{PHP}_q(i)$ .  $\square$

The relationship between PHP and reverse RT in Theorem 16 is quite similar to the relationship between PHP and RT in Theorem 10. We only need to change  $\beta$  to  $(1 - \beta)$  in the FLoS method for RT to apply it for reverse RT.

#### H.4 Extension to Reverse AP

In this subsection, we show how to find the top- $k$  nodes for reverse AP.

**Theorem 17.**  $\text{PHP}''$  and reverse AP give the same ranking results.

*Proof:* In Appendix G.4, we have that the proximity matrix of AP is  $\mathbf{R} = (\mathbf{A} + \mathbf{D} - \mathbf{W})^{-1} \mathbf{A}$ . The elements in the  $q$ -th column of  $\mathbf{R}$  represent the reverse AP proximity values when the query node is  $q$ . Thus, the reverse AP proximity

vector is  $\mathbf{r} = \mathbf{R}\mathbf{e} = (\mathbf{A} + \mathbf{D} - \mathbf{W})^{-1} \mathbf{A}\mathbf{e}$ . Then,  $(\mathbf{A} + \mathbf{D} - \mathbf{W})\mathbf{r} = \mathbf{A}\mathbf{e}$ . Thus, the proximity vector of reverse AP satisfies the following equation

$$\mathbf{r} = (\mathbf{A} + \mathbf{D})^{-1} \mathbf{W}\mathbf{r} + (\mathbf{A} + \mathbf{D})^{-1} \mathbf{A}\mathbf{e}.$$

Based on this matrix form, reverse AP has the following recursive definition

$$\mathbf{r}_i = \begin{cases} \sum_{j \in N_i} p_{i,j} \mathbf{r}_j + \frac{\lambda_i}{\lambda_i + w_i}, & \text{if } i = q, \\ \sum_{j \in N_i} p_{i,j} \mathbf{r}_j, & \text{if } i \neq q, \end{cases}$$

where  $p_{i,j} = \frac{w_{i,j}}{\lambda_i + w_i}$  is the transition probability from node  $i$  to  $j$ .

PHP'' and reverse AP have the same recursive definition for any node  $i \neq q$ , thus  $\frac{AP_i(q)}{PHP''_q(i)} = \frac{AP_q(q)}{PHP''_q(q)}$ . When the query node  $q$  is fixed,  $AP_i(q) \propto PHP''_q(i)$ . This completes the proof.  $\square$

EI and KZ both are symmetric, thus the top- $k$  results for the reverse proximity are the same as the top- $k$  results for the original proximity.

## APPENDIX I

### COMPARISON OF CASTANET AND FLOS

Castanet [4] is an improved global iteration method, which needs to iterate over the entire graph and compute the exact proximity of each node. Castanet uses the developed bounds to identify the top- $k$  nodes earlier thus reduce the number of iterations required by the global iteration method.

The bounds in Castanet are based on the following probability distributions. Starting from the query node, the random surfer randomly walks to the neighbor nodes following the transition probabilities. Let  $p_i^t$  represent the probability that the random surfer is at node  $i$  at time  $t$ . We have that  $\sum_{i \in V} p_i^t = 1$  for any  $t$ . Initially,  $p_q^0 = 1$  and  $p_i^0 = 0$  if  $i \neq q$ .

At each iteration  $t$ , Castanet computes the probability  $p_i^t$  and uses  $p_i^t$  to refine the lower and upper bounds of node  $i$ . We can see that, at each iteration  $t$ , Castanet only uses the probability distribution at time  $t$  to refine the lower and upper bounds. However, at each iteration  $t$ , FLoS uses the probability distributions at all time points to refine the bounds. This is the key difference between the bounds in Castanet and FLoS and the reason why the bounds in FLoS are more effective than the bounds in Castanet.

## APPENDIX J

### SIMRANK

SimRank is a classic random walk based proximity measure [31]. It is based on the intuition that two nodes are similar if their neighbors are similar. The SimRank value between two nodes measures the expected number of steps required before two random surfers, one starting from each node, meet at the same node if they walk in lock-step.

**Lemma 16.** SimRank has local maximum.

*Proof:* Examples can be constructed to show that SimRank has local maximum, which are omitted.  $\square$

The random walk processes in SimRank and PHP are quite different. In SimRank, the random walk process involves two random surfers. In PHP, the random walk process only involves one random surfer. There is no direct relationship between SimRank and PHP. Thus, the FLoS algorithm is not applicable to SimRank at this time. We will investigate how to apply FLoS to SimRank in our future work.