

# Efficient and Provable Effective Resistance Computation on Large Graphs: an Index-based Approach

## ABSTRACT

Effective resistance (ER) is a fundamental metric for measuring node similarities in a graph, and it finds applications in various domains including graph clustering, recommendation systems, link prediction, and graph neural networks. The state-of-the-art algorithm for computing effective resistance relies on a landmark technique, which involves selecting a node that is easy to reach by all the other nodes as a landmark. The performance of this technique heavily depends on the chosen landmark node. However, in many real-life graphs, it is not always possible to find an easily reachable landmark node, which can significantly hinder the algorithm's efficiency. To overcome this problem, we propose a novel multiple landmarks technique which involves selecting a set of landmark nodes  $\mathcal{V}_l$  such that the other nodes in the graph can easily reach any one of a landmark node in  $\mathcal{V}_l$ . Specifically, we first propose several new formulas to compute ER with multiple landmarks, utilizing the concept of Schur complement. These new formulas allow us to pre-compute and maintain several small-sized matrices related to  $\mathcal{V}_l$  as a compact index. With this powerful index technique, we demonstrate that both single-pair and single-source ER queries can be efficiently answered using a newly-developed  $\mathcal{V}_l$ -absorbed random walk sampling or  $\mathcal{V}_l$ -absorbed push technique. Comprehensive theoretical analysis shows that all proposed index-based algorithms achieve provable performance guarantees for both single-pair and single-source ER queries. Extensive experiments on 5 real-life datasets demonstrate the high efficiency of our multiple landmarks-based index techniques. For instance, our algorithms, with a 1.5 GB index size, can be up to 4 orders of magnitude faster than the state-of-the-art algorithms while achieving the same accuracy on a large road network.

## ACM Reference Format:

. 2022. Efficient and Provable Effective Resistance Computation on Large Graphs: an Index-based Approach. In *Proceedings of SIGMOD '23: International Conference on Management of Data (SIGMOD '23)*. ACM, New York, NY, USA, 20 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Effective resistance is a fundamental metric to measure node similarities of a graph. Given an undirected graph  $\mathcal{G}$  and two nodes  $s, t$ , the effective resistance between  $s$  and  $t$ , denoted by  $r(s, t)$ , is equivalent to (up to a constant factor) the expected number of steps taken by a random walk starting from  $s$ , visiting  $t$ , and coming back to  $s$  [38]. Intuitively, a small  $r(s, t)$  indicates a high similarity between  $s$  and  $t$ .

This is because a small effective resistance implies that it is relatively easy for a random walk starting from  $s$  to reach  $t$ , and vice versa. In other words, the nodes  $s$  and  $t$  are more likely to be well-connected or have similar neighbors in the graph, suggesting a high degree of similarity between them. Compared to the classic shortest-path distance metric, effective resistance is often more robust to noise, such as the deletion or insertion of a small number of edges in a graph, as it considers all possible paths [20, 38].

The effective resistance metric has found wide applications in graph data management and mining, including clustering in geo-social networks [34], long-tail recommendation systems [49], link prediction in social networks [32], and anomaly detection in time-varying graphs [37]. Recently, it has also been employed to analyze the over-squashing problem in graph neural networks [23, 39]. These applications leverage the power of effective resistance to capture and quantify node similarities in various graph-based scenarios.

Despite many efforts have been made to compute effective resistance in various studies [20, 29, 36, 47], there is still a lack of an efficient and provable algorithm for computing effective resistance on large graphs. Recently, Peng et al. [29] proposed several efficient algorithms to compute single-pair effective resistance using an  $L$ -step random-walk sampling technique. They proved that their algorithms achieves sub-linear time complexity with respect to the size of the graph. Building upon this, Yang et al. [47] further improved the algorithm by reducing the required length of  $L$  while maintaining accuracy. However, a major limitation of these algorithms is that when the effective resistance  $r(s, t)$  is large, the random walk length  $L$  can still be very large. Additionally, these algorithms can only compute single-pair effective resistance and would need to be invoked  $n - 1$  times to compute the single-source effective resistance, resulting in significant computational cost.

In addition to the work by Yang et al. [47], Liao et al. [20] proposed an alternative random walk sampling algorithm for effective resistance computation, which utilizes a landmark node. This algorithm is considered as the state-of-the-art for both single-pair and single-source effective resistance calculations on large real-life graphs, as demonstrated in our experiments. The key idea behind their algorithm is as follows: rather than performing random walks from the source node  $s$  to the target node  $t$ , random walks are conducted from both  $s$  and  $t$  towards an easily reachable node (e.g., the highest-degree node), referred to as the landmark node. By terminating the random walk once it hits the landmark node, the random walk sampling process can be executed quickly. The estimated value of  $r(s, t)$  is then derived based on these random walks. Additionally, the landmark idea is also extended to develop a random spanning forest sampling technique to compute both single-pair and single-source effective resistance. A notable advantage of their algorithm for computing single-pair effective resistance is that it only requires exploration of a small portion of the graph, resulting in high efficiency. However, a limitation of this approach is that its performance heavily relies on the selection of an appropriate landmark node.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

SIGMOD '23, June 03–05, 2022, Woodstock, NY

© 2022 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

In cases where finding an easily reachable landmark node is challenging, such as in road networks, their algorithm may yield poor performance, as observed in our experiments.

To overcome this problem, in this work, we propose a novel multiple landmarks technique based on several newly-developed effective resistance formulas. This technique allows us to strategically select a set of multiple landmark nodes, denoted by  $\mathcal{V}_l$ , which in turn enables faster random walk sampling. By leveraging the multiple landmarks technique, we develop a novel, efficient and provable index-based approach to efficiently compute both the single-pair and single-source effective resistance queries. More specifically, we first propose three new effective resistance formulas with multiple landmarks, using a classic concept of Schur complement [7]. Based on these new formulas, we can pre-compute and maintain several small-sized matrices, which are related the landmark node set  $\mathcal{V}_l$ , as a compact index. To construct the index, we develop two efficient and provable Monte Carlo algorithms based on several interesting and newly-discovered probability interpretations of the Schur complement. Armed with our indexing technique, we propose new and provable  $\mathcal{V}_l$ -absorbed random walk and  $\mathcal{V}_l$ -absorbed push algorithms to efficiently answer both the single-pair and single-source effective resistance queries. We show that the state-of-the-art algorithm based on a single landmark node is a special case of our algorithms. In addition, we also present comprehensive theoretical analysis of our algorithms, and the results demonstrate that all our algorithms can achieve an  $\epsilon$ -absolute error guarantee while using only sub-linear time. Finally, we conduct extensive experiments using 5 large real-life graphs to evaluate our algorithms, and the results show that our algorithms are substantially faster than the state-of-the-art algorithms while maintaining the same accuracy level. To summarize, the main contributions of this paper are as follows.

**New theoretical results.** We propose three new formulas to compute the effective resistance with multiple landmarks  $\mathcal{V}_l$ . We propose novel and interesting probability interpretations of the Schur complement of a graph w.r.t. the landmark nodes set  $\mathcal{V}_l$ . We believe that these novel probability interpretations could be of independent interest. In addition, we also establish several novel connections among effective resistance and three new concepts called  $\mathcal{V}_l$ -absorbed random walk,  $\mathcal{V}_l$ -absorbed push, and  $\mathcal{V}_l$ -rooted random spanning forest. We present detailed theoretical analysis for all our algorithms and the results show that all of them can achieve provable accuracy guarantee and take sub-linear running time.

**Novel index-based algorithms.** We propose novel index-based approaches to both single-pair and single-source effective resistance computations. Specifically, we first propose several novel techniques to construct the index, including  $\mathcal{V}_l$ -absorbed random walk sampling,  $\mathcal{V}_l$ -rooted random spanning forest sampling, and loop-erased random walk sampling. Then, armed with our index, we develop two new  $\mathcal{V}_l$ -absorbed random walk and  $\mathcal{V}_l$ -absorbed push algorithms to efficiently answer both the single-pair and single-source effective resistance queries. We show that the time complexity of our query processing algorithms is strictly lower than the state-of-the-art (SOTA) algorithms.

**Extensive experiments.** We conduct comprehensive experiments on 5 large real-life graphs to evaluate our algorithms. The results

show that by selecting a small set of landmark nodes, our algorithms can be up to 4 orders magnitude faster than the SOTA algorithms for both single-pair and single-source effective resistance computations, when achieve the same accuracy. For example, on the Road-PA dataset (1.09 million nodes, 1.54 million edges), our algorithms takes 0.19 and 30 seconds to achieve a relative error 0.05 for single-pair and single-source queries respectively, while the SOTA algorithm consumes 184 and  $2 \times 10^4$  seconds respectively. Additionally, the results also show that our index can be constructed very efficient and uses acceptable space. For instance, on Road-PA dataset, our index can be built within 788 seconds using 1493 MB spaces. For reproducible purpose, the source code of the paper can be found at <https://anonymous.4open.science/r/ML-EED5>.

## 2 PRELIMINARIES

**Notations and definitions.** Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be an undirected graph with  $|\mathcal{V}| = n$  nodes and  $|\mathcal{E}| = m$  edges. For a node  $u$ , denote by  $\mathcal{N}(u)$  the set of neighbor nodes of  $u$ . Let  $\mathbf{A}$  be the adjacency matrix of  $\mathcal{G}$  and  $\mathbf{D}$  be the diagonal degree matrix with  $\mathbf{D}_{uu} = d(u) = |\mathcal{N}(u)|$ . Let  $\mathbf{L}$  be the Laplacian matrix of  $\mathcal{G}$  defined as  $\mathbf{L} = \mathbf{D} - \mathbf{A}$ . Without loss of generality, we assume that the graph is connected (if the graph is disconnected, we can consider each connected component separately). For any  $n$ -nodes connected graph, the Laplacian matrix  $\mathbf{L}$  has a rank  $n - 1$ , thus the inverse of  $\mathbf{L}$  does not exist. Let  $0 = \mu_1 \leq \dots \leq \mu_n$  be the eigenvalues of  $\mathbf{L}$ , and  $\mathbf{u}_1, \dots, \mathbf{u}_n$  be the corresponding eigenvectors. The eigen-decomposition of  $\mathbf{L}$  is  $\mathbf{L} = \sum_{i=2}^n \mu_i \mathbf{u}_i \mathbf{u}_i^T$ , because  $\mu_1 = 0$ . Based on this, the classic Moore-Penrose pseudo-inverse of  $\mathbf{L}$  can be defined as  $\mathbf{L}^\dagger = \sum_{i=2}^n \frac{1}{\mu_i} \mathbf{u}_i \mathbf{u}_i^T$ . Given a graph  $\mathcal{G}$  and two nodes  $s, t$ , the effective resistance (ER) between  $s$  and  $t$  is defined as

$$r(s, t) = (\mathbf{L}^\dagger)_{ss} + (\mathbf{L}^\dagger)_{tt} - 2(\mathbf{L}^\dagger)_{st}. \quad (1)$$

Eq. (1) is often used for computing the ER. However, as shown in [7], the ER can also be computed by the so-called  $g$ -inverse of  $\mathbf{L}$ . Specifically, the  $g$ -inverse of  $\mathbf{L}$ , denoted by  $\mathbf{H}$ , is a symmetric matrix that satisfies  $\mathbf{LH} = \mathbf{L}$ . Then,  $r(s, t)$  can be determined by the following formula:

$$r(s, t) = (\mathbf{H})_{ss} + (\mathbf{H})_{tt} - 2(\mathbf{H})_{st}. \quad (2)$$

Random walk is a random process on graphs. In each step, a random surfer jumps to the neighbor of the current node  $u$  with probability  $\frac{1}{d(u)}$ . The probability transition matrix is defined as  $\mathbf{P} = \mathbf{D}^{-1}\mathbf{A}$ . The commute time  $c(s, t)$  between two nodes  $s, t$  is the expected number of steps of the random walk starting from  $s$ , visiting  $t$ , and then coming back to  $s$ . It is well known that the commute time is closely related to ER, i.e.,  $c(s, t) = 2m \times r(s, t)$  [38]. The hitting time  $h(s, t)$  between two nodes  $s, t$  is the expected number of steps of the random walk starting from  $s$  and visiting  $t$  for the first time. Similarly, the hitting time  $h(s, \mathcal{V}_l)$  is the expected number of steps of the random walk starting from  $s$  and visiting a node set  $\mathcal{V}_l$  for the first time.

Let  $\mathbf{L}_v$  be a sub-matrix of  $\mathbf{L}$  which is obtained by deleting the  $v$ -th row and the  $v$ -th column of  $\mathbf{L}$ . Unlike  $\mathbf{L}$ , the inverse of  $\mathbf{L}_v$  exists for any node  $v$ . Recently, Liao et al. [20] shows that ER can also be computed based on the matrix  $\mathbf{L}_v$ :

$$r(s, t) = (\mathbf{L}_v^{-1})_{ss} + (\mathbf{L}_v^{-1})_{tt} - 2(\mathbf{L}_v^{-1})_{st}, \quad u_1, u_2 \neq v; \quad (3)$$

$$r(u, v) = (\mathbf{L}_v^{-1})_{uu}. \quad (4)$$

Note that both Eq. (3) and Eq. (4) always hold for any node  $v$ .

**Problem formulation.** Based on the definition of ER, we formulate two ER computation problems as follows.

*Definition 2.1 (Single-pair ER computation).* Given an undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and two nodes  $s, t$  with  $s \neq t$ , the single-pair ER computation problem is to calculate  $r(s, t)$ .

*Definition 2.2 (Single-source ER computation).* Given an undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and a source node  $s$ , the single-source ER computation problem is to calculate  $r(s, t)$  for every  $t \in \mathcal{V}$ .

Note that by definition, computing the exact effective resistance requires to calculate the matrix inverse (or pseudo-inverse) which is typically intractable for large graphs. As a result, existing solutions often focus on deriving an  $\epsilon$ -estimation of the effective resistance [20, 29, 47], i.e., outputting an approximate effective resistance  $\hat{r}(s, t)$  that satisfies  $|\hat{r}(s, t) - r(s, t)| \leq \epsilon$  with a high probability (e.g., the failure probability is smaller than 0.01). In this work, we also focus on developing efficient  $\epsilon$ -estimation algorithms for both single-pair and single-source ER computation. Below, we first establish several new effective resistance formulas with multiple landmark nodes, based on which we will develop a powerful index-based approach to compute single-pair and single-source ER queries.

### 3 NEW EFFECTIVE RESISTANCE FORMULAS

Recall that Liao et al. [20] proposed an effective resistance (ER) formula based on the matrix  $\mathbf{L}_v$  that involves only one landmark node  $v$ . A natural question arises: can we extend their formula to handle scenarios with multiple landmark nodes? Note that solving this problem is quite nontrivial as it needs to incorporate the concept of Schur complement. Specifically, let  $\mathcal{V}_l \subset \mathcal{V}$  be a small set of *landmark nodes* (e.g.,  $|\mathcal{V}_l| \leq 100$ ), and  $\mathcal{U} = \mathcal{V} \setminus \mathcal{V}_l$  be the set of remaining nodes. Denote by  $\mathbf{L}_{\mathcal{U}\mathcal{U}}$  the sub-matrix obtained by deleting the rows and columns indexed by the landmark node set  $\mathcal{V}_l$ .

Then,  $\mathbf{L}$  can be represented as  $\mathbf{L} = \begin{bmatrix} \mathbf{L}_{\mathcal{U}\mathcal{U}} & \mathbf{L}_{\mathcal{U}\mathcal{V}_l} \\ \mathbf{L}_{\mathcal{V}_l\mathcal{U}} & \mathbf{L}_{\mathcal{V}_l\mathcal{V}_l} \end{bmatrix}$ . Based on this block representation, we formally define the Schur complement of the landmark node set  $\mathcal{V}_l$ , denoted by  $\mathbf{L}/\mathcal{V}_l$ , as follows.

*Definition 3.1.* For a node set  $\mathcal{V}_l$ , the Schur complement of  $\mathcal{V}_l$  is  $\mathbf{L}/\mathcal{V}_l \triangleq \mathbf{L}_{\mathcal{V}_l\mathcal{V}_l} - \mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l}$ .

**New ER formulas.** It is well-known that the Schur complement  $\mathbf{L}/\mathcal{V}_l$  is also the Laplacian matrix of a weighted graph  $\mathcal{G}/\mathcal{V}_l$  [9]. Let  $(\mathbf{L}/\mathcal{V}_l)^\dagger$  be the pseudo-inverse of  $\mathbf{L}/\mathcal{V}_l$ . Similar to  $\mathbf{L}_v$ , the inverse of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}$  also exists. Below, we construct a new  $g$ -inverse of  $\mathbf{L}$  with the pseudo-inverse of the Schur complement  $\mathbf{L}/\mathcal{V}_l$ .

**LEMMA 3.2.** *Given a landmark node set  $\mathcal{V}_l$  and the remaining node set  $\mathcal{U} = \mathcal{V} \setminus \mathcal{V}_l$ , a  $g$ -inverse of  $\mathbf{L}$ , denoted by  $\mathbf{H}$ , can be constructed as:*

$$\begin{aligned} \mathbf{H} &= \begin{bmatrix} \mathbf{I} & \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l} \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} & 0 \\ 0 & (\mathbf{L}/\mathcal{V}_l)^\dagger \end{bmatrix} \begin{bmatrix} \mathbf{I} & 0 \\ \mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} & \mathbf{I} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} + \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l}(\mathbf{L}/\mathcal{V}_l)^\dagger\mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} & -\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l}(\mathbf{L}/\mathcal{V}_l)^\dagger \\ -(\mathbf{L}/\mathcal{V}_l)^\dagger\mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} & (\mathbf{L}/\mathcal{V}_l)^\dagger \end{bmatrix}. \end{aligned} \quad (5)$$

**PROOF.** By definition, if a symmetric matrix  $\mathbf{H}$  satisfies the property  $\mathbf{L}\mathbf{H}\mathbf{L} = \mathbf{L}$ ,  $\mathbf{H}$  is a  $g$ -inverse of  $\mathbf{L}$ . Clearly,  $\mathbf{H}$  defined in Eq. (5) is symmetric. Since  $\begin{bmatrix} \mathbf{I} & \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l} \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{I} & -\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l} \\ 0 & \mathbf{I} \end{bmatrix} = \mathbf{I}$ ,

and  $\begin{bmatrix} \mathbf{I} & 0 \\ -\mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{I} & 0 \\ \mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} & \mathbf{I} \end{bmatrix} = \mathbf{I}$ , we have:

$$\begin{aligned} \mathbf{L}\mathbf{L}^\dagger\mathbf{L} &= \begin{bmatrix} \mathbf{I} & 0 \\ \mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{L}_{\mathcal{U}\mathcal{U}} & 0 \\ 0 & \mathbf{L}/\mathcal{V}_l \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l} \\ 0 & \mathbf{I} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{I} & -\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l} \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{L}_{\mathcal{U}\mathcal{U}} & 0 \\ 0 & (\mathbf{L}/\mathcal{V}_l)^\dagger \end{bmatrix} \begin{bmatrix} \mathbf{I} & 0 \\ -\mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} & \mathbf{I} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{I} & 0 \\ \mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{L}_{\mathcal{U}\mathcal{U}} & 0 \\ 0 & \mathbf{L}/\mathcal{V}_l \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l} \\ 0 & \mathbf{I} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{I} & 0 \\ \mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{L}_{\mathcal{U}\mathcal{U}} & 0 \\ 0 & \mathbf{L}/\mathcal{V}_l \end{bmatrix} \begin{bmatrix} \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} & 0 \\ 0 & (\mathbf{L}/\mathcal{V}_l)^\dagger \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{L}_{\mathcal{U}\mathcal{U}} & 0 \\ 0 & \mathbf{L}/\mathcal{V}_l \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l} \\ 0 & \mathbf{I} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{I} & 0 \\ \mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{L}_{\mathcal{U}\mathcal{U}} & 0 \\ 0 & \mathbf{L}/\mathcal{V}_l \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l} \\ 0 & \mathbf{I} \end{bmatrix} \\ &= \mathbf{L}, \end{aligned}$$

Thus, the Lemma is established.  $\square$

Based on Lemma 3.2 and Eq. (1), we present new formulas for ER computation as follows:

**LEMMA 3.3.** *(New effective resistance formulas) Given a landmark node set  $\mathcal{V}_l$  and the remaining node set  $\mathcal{U} = \mathcal{V} \setminus \mathcal{V}_l$ , the effective resistance between any two nodes can be computed by the following formulas.*

(1) For  $u_1, u_2 \in \mathcal{U}$ , we have

$$\begin{aligned} r(u_1, u_2) &= (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{u_1 u_1} + (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{u_2 u_2} - 2(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{u_1 u_2} \\ &\quad + (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l}(\mathbf{L}/\mathcal{V}_l)^\dagger\mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{u_1 u_1} \\ &\quad + (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l}(\mathbf{L}/\mathcal{V}_l)^\dagger\mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{u_2 u_2} \\ &\quad - 2(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l}(\mathbf{L}/\mathcal{V}_l)^\dagger\mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{u_1 u_2}; \end{aligned} \quad (6)$$

(2) For  $u \in \mathcal{U}, v \in \mathcal{V}_l$ , we have

$$\begin{aligned} r(u, v) &= (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} + ((\mathbf{L}/\mathcal{V}_l)^\dagger)_{vv} \\ &\quad + (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l}(\mathbf{L}/\mathcal{V}_l)^\dagger\mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} \\ &\quad - 2(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l}(\mathbf{L}/\mathcal{V}_l)^\dagger)_{uv}; \end{aligned} \quad (7)$$

(3) For  $v_1, v_2 \in \mathcal{V}_l$ , we have

$$\begin{aligned} r(v_1, v_2) &= ((\mathbf{L}/\mathcal{V}_l)^\dagger)_{v_1 v_1} + ((\mathbf{L}/\mathcal{V}_l)^\dagger)_{v_2 v_2} \\ &\quad - 2((\mathbf{L}/\mathcal{V}_l)^\dagger)_{v_1 v_2}. \end{aligned} \quad (8)$$

**PROOF.** Based on Lemma 3.2, since  $r(s, t) = (\mathbf{H})_{ss} + (\mathbf{H})_{tt} - 2(\mathbf{H})_{st} = (\mathbf{e}_s - \mathbf{e}_t)^T \mathbf{H} (\mathbf{e}_s - \mathbf{e}_t)$ , the formulas can be derived in three cases: (1) For  $u_1, u_2 \in \mathcal{U}$ , we can derive  $r(u_1, u_2) = (\mathbf{e}_{u_1} - \mathbf{e}_{u_2})^T \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}(\mathbf{e}_{u_1} - \mathbf{e}_{u_2}) + (\mathbf{e}_{u_1} - \mathbf{e}_{u_2})^T \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l}(\mathbf{L}/\mathcal{V}_l)^\dagger\mathbf{L}_{\mathcal{V}_l\mathcal{U}}\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}(\mathbf{e}_{u_1} - \mathbf{e}_{u_2})$  by Eq. (5). By expanding the bilinear form, we can derive Eq. (6); (2) For  $u \in \mathcal{U}, v \in \mathcal{V}_l$ , we have  $r(u, v) = (\mathbf{e}_u - \mathbf{e}_v)^T \mathbf{H} (\mathbf{e}_u - \mathbf{e}_v)$ , by Eq. (5), we can derive Eq. (7); (3) For  $v_1, v_2 \in \mathcal{V}_l$ , we have  $r(v_1, v_2) = (\mathbf{e}_{v_1} - \mathbf{e}_{v_2})^T (\mathbf{L}/\mathcal{V}_l)^\dagger (\mathbf{e}_{v_1} - \mathbf{e}_{v_2})$ . By expanding the bilinear form, we can obtain Eq. (8). The proof is completed.  $\square$

By Lemma 3.3, to compute the effective resistance for any two nodes, it is sufficient to calculate (i) the matrix  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}\mathbf{L}_{\mathcal{U}\mathcal{V}_l}$  (note that

it is no need to repeatedly compute  $\mathbf{L}_{\mathcal{V}_l \mathcal{U}} \mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1}$ , because  $\mathbf{L}_{\mathcal{V}_l \mathcal{U}} \mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} = (\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l})^T$ ; (ii) the pseudo-inverse of the Schur complement:  $(\mathbf{L}/\mathcal{V}_l)^\dagger$ ; and (iii) the inverse of the Laplacian submatrix:  $\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1}$ . By Definition 3.1, the Schur complement  $\mathbf{L}/\mathcal{V}_l$  is closely related to the matrix  $\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l}$ . Note the Schur complement is a  $|\mathcal{V}_l| \times |\mathcal{V}_l|$  matrix and  $|\mathcal{V}_l|$  is often small (e.g.,  $|\mathcal{V}_l| \leq 100$ ), thus  $(\mathbf{L}/\mathcal{V}_l)^\dagger$  can be easily computed after obtaining the matrix  $\mathbf{L}/\mathcal{V}_l$ . As a consequence, for both (i) and (ii), the key is to determine the matrix  $\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l}$ . Below, we propose several interesting probability interpretations of the matrix  $\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l}$ , based on which we will develop several novel and efficient sampling techniques to estimate such a matrix.

**Novel probability interpretations of  $\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l}$ .** We first give a random walk interpretation for each element of the matrix  $\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l}$ . For a given set of landmark nodes  $\mathcal{V}_l$ , we define a  $\mathcal{V}_l$ -absorbed random walk as a random walk that terminates once it reaches any node in the set  $\mathcal{V}_l$ . That is to say, for a  $\mathcal{V}_l$ -absorbed random walk starting from a node  $u \in \mathcal{U}$ , it must stop at a node  $v \in \mathcal{V}_l$ . Let  $\tau_{\mathcal{V}_l}[u_1, u_2]$  be the expected number of visits to node  $u_2$  in a  $\mathcal{V}_l$ -absorbed walk starting from node  $u_1$ . Denote by  $\mathbf{P}_{\mathcal{U} \mathcal{U}}$  a sub-matrix of the probability transition matrix  $\mathbf{P} = \mathbf{D}^{-1} \mathbf{A}$  obtained by removing the rows and columns indexed by the set  $\mathcal{V}_l$ . Then, by definition, we have  $\tau_{\mathcal{V}_l}[u_1, u_2] = \sum_{k=0}^{\infty} (\mathbf{P}_{\mathcal{U} \mathcal{U}}^k)_{u_1 u_2}$ . Specifically,  $(\mathbf{P}_{\mathcal{U} \mathcal{U}}^k)_{u_1 u_2}$  is the probability that a  $\mathcal{V}_l$ -random walk starting from  $u_1$  visits  $u_2$  at the  $k$ -th step. Based on this, we have the following results.

LEMMA 3.4.  $(\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1})_{u_1 u_2} = \frac{\tau_{\mathcal{V}_l}[u_1, u_2]}{d(u_2)}$ .

PROOF. Let  $\mathbf{D}_{\mathcal{U} \mathcal{U}}$  be a sub-matrix of  $\mathbf{D}$  obtained by removing the rows and columns indexed by the set  $\mathcal{V}_l$ . Since  $\mathbf{D}_{\mathcal{U} \mathcal{U}}(\mathbf{I} - \mathbf{P}_{\mathcal{U} \mathcal{U}}) = \mathbf{L}_{\mathcal{U} \mathcal{U}}$ , we have  $(\mathbf{I} - \mathbf{P}_{\mathcal{U} \mathcal{U}})^{-1} \mathbf{D}_{\mathcal{U} \mathcal{U}} = \mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1}$ . Note that  $\tau_{\mathcal{V}_l}[u_1, u_2] = \sum_{k=0}^{\infty} (\mathbf{P}_{\mathcal{U} \mathcal{U}}^k)_{u_1 u_2} = (\mathbf{I} - \mathbf{P}_{\mathcal{U} \mathcal{U}})^{-1} \mathbf{D}_{\mathcal{U} \mathcal{U}}^{-1}$ . Thus, the lemma is established.  $\square$

Then, by Lemma 3.4, we can derive an interesting probability interpretation of the matrix  $\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l}$  as follows.

LEMMA 3.5. ( $\mathcal{V}_l$ -absorbed random walk interpretation) Given a node  $u \in \mathcal{U}$  and a node  $v \in \mathcal{V}_l$ , the  $u, v$ -th element of the matrix  $-(\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l})_{uv}$  is the probability  $p_{uv}^r$  that a  $\mathcal{V}_l$ -absorbed random walk starts from  $u$  and terminates at the node  $v \in \mathcal{V}_l$ .

PROOF. Recall that  $\tau_{\mathcal{V}_l}[u, w] = \sum_{k=0}^{\infty} (\mathbf{P}_{\mathcal{U} \mathcal{U}}^k)_{uw}$ , where  $(\mathbf{P}_{\mathcal{U} \mathcal{U}}^k)_{uw}$  is the probability that a  $\mathcal{V}_l$ -random walk starting from  $u$  visits  $w$  at the  $k$ -th step. If  $w \in \mathcal{N}(v)$ , then with probability  $\frac{1}{d(w)}$ , the random walk will reach  $v \in \mathcal{V}_l$  in the next step (i.e., it will terminate at  $v$ ). Thus, by adding up all such probabilities, we have

$$\begin{aligned} p_{uv}^r &= \sum_{w \in \mathcal{N}(v) \cap \mathcal{U}} \frac{\sum_{k=0}^{\infty} (\mathbf{P}_{\mathcal{U} \mathcal{U}}^k)_{uw}}{d_w} = \sum_{w \in \mathcal{N}(v) \cap \mathcal{U}} \frac{\tau_{\mathcal{V}_l}[u, w]}{d_w} \\ &= \sum_{w \in \mathcal{N}(v) \cap \mathcal{U}} (\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1})_{uw}. \end{aligned}$$

The last equality is due to Lemma 3.4. Since  $\mathbf{D}_{\mathcal{U} \mathcal{V}_l}$  is a zero matrix by definition, we have  $\mathbf{L}_{\mathcal{U} \mathcal{V}_l} = -\mathbf{A}_{\mathcal{U} \mathcal{V}_l}$ . Therefore,  $\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l} = -\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{A}_{\mathcal{U} \mathcal{V}_l}$ . As a consequence,  $p_{uv}^r = \sum_{w \in \mathcal{N}(v) \cap \mathcal{U}} (\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1})_{uw} = \mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{A}_{\mathcal{U} \mathcal{V}_l} = -(\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l})_{uv}$ . This completes the proof.  $\square$

Lemma 3.5 indicates that the matrix  $\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l}$  can be efficiently estimated by using  $\mathcal{V}_l$ -random walk sampling technique. Below, we further propose a different probability interpretation for  $\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l}$  based on a concept of random spanning forest.

A spanning forest is a subgraph of  $\mathcal{G}$  with no cycles. For each connected component of a spanning forest, we specify a node as the root node in a rooted spanning forest. The set of all root nodes is called the root set  $\mathcal{R}$ . For convenience, we let the root set be the landmark node set  $\mathcal{V}_l$ , i.e.,  $\mathcal{R} = \mathcal{V}_l$ . For a tree  $\mathcal{T}$  with a root node  $v$ , we call that “ $u$  is rooted at  $v$ ” for each node  $u \in \mathcal{T}$ . A random spanning forest with a prescribed root set  $\mathcal{V}_l$  is a spanning forest uniformly sampled from all spanning forests of  $\mathcal{G}$  with the root set  $\mathcal{V}_l$ . With this concepts, we can derive the following results.

LEMMA 3.6. ( $\mathcal{V}_l$ -rooted random spanning forest interpretation) Given a node  $u \in \mathcal{U}$  and a node  $v \in \mathcal{V}_l$ , the  $u, v$ -th element  $-(\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l})_{uv}$  is the probability  $p_{uv}^f$  that in a random spanning forest  $\mathcal{F}$  with root set  $\mathcal{V}_l$ ,  $u$  is rooted at  $v$  in a tree of  $\mathcal{F}$ .

PROOF. According to the all-minors matrix forest theorem [10], we have  $(\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1})_{uw} = \frac{|\mathcal{F}_{\mathcal{V}_l|u,w}|}{|\mathcal{F}_{\mathcal{V}_l}|}$ , where  $\mathcal{F}_{\mathcal{V}_l}$  is the set of spanning forests with root set  $\mathcal{V}_l$ ,  $\mathcal{F}_{\mathcal{V}_l|u,w}$  is the set of spanning forests with root set  $\mathcal{V}_l \cup \{u\}$ , and  $w$  belongs to the tree component with root  $u$ . Let  $\mathcal{F}_{\mathcal{V}_l:u \leadsto v}$  be the set of spanning forests with root set  $\mathcal{V}_l$ ,  $u$  is rooted at  $v$ . There is a bijection between  $\mathcal{F}_{\mathcal{V}_l:u \leadsto v}$  and  $\bigcup_{w \in \mathcal{N}(v) \cap \mathcal{U}} \mathcal{F}_{\mathcal{V}_l|u,w}$ . To see this, if we remove the last edge on the path from  $u$  to  $\mathcal{V}_l$  in each spanning forest in  $\mathcal{F}_{\mathcal{V}_l:u \leadsto v}$ , we will exactly obtain a spanning forest in  $\bigcup_{w \in \mathcal{N}(v) \cap \mathcal{U}} \mathcal{F}_{\mathcal{V}_l|u,w}$ . And it is easy to verify that the reverse is also true. Thus, we have:

$$\begin{aligned} p_{uv}^f &= -\frac{\sum_{w \in \mathcal{N}(v) \cap \mathcal{U}} |\mathcal{F}_{\mathcal{V}_l|u,w}|}{|\mathcal{F}_{\mathcal{V}_l}|} = \sum_{w \in \mathcal{N}(v) \cap \mathcal{U}} (\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1})_{uw} \\ &= -(\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l})_{uv}. \end{aligned}$$

As a result, the element  $-(\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l})_{uv}$  is exactly the probability that in a random spanning forest with root set  $\mathcal{V}_l$ , the node  $u$  is rooted at  $v$ .  $\square$

Note that each row of the matrix  $-\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l}$  sums up to 1, this is because a random walk starts from a node  $u \in \mathcal{U}$  must stop when hitting a node  $v \in \mathcal{V}_l$ , thus the probability  $\sum_{v \in \mathcal{V}_l} p_{uv}^r = 1$ . Similarly, in a random spanning forest with root set  $\mathcal{V}_l$ , each node  $u \in \mathcal{U}$  must be root at a certain landmark node  $v \in \mathcal{V}_l$ , the probability  $\sum_{v \in \mathcal{V}_l} p_{uv}^f = 1$ . We use  $\mathbf{P}_r \in \mathbb{R}^{|\mathcal{U}| \times |\mathcal{V}_l|}$  to denote the random walk probability matrix where  $(\mathbf{P}_r \in \mathbb{R}^{|\mathcal{U}| \times |\mathcal{V}_l|})_{uv} = p_{uv}^r$ , and  $\mathbf{P}_f \in \mathbb{R}^{|\mathcal{U}| \times |\mathcal{V}_l|}$  to denote the random spanning forest probability matrix where  $(\mathbf{P}_f \in \mathbb{R}^{|\mathcal{U}| \times |\mathcal{V}_l|})_{uv} = p_{uv}^f$ . Therefore, we have  $\mathbf{P}_r = \mathbf{P}_f = -\mathbf{L}_{\mathcal{U} \mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U} \mathcal{V}_l}$  by Lemma 3.5 and Lemma 3.6. Then, we can further obtain three new ER formulas based on the matrices  $\mathbf{P}_r$  and  $\mathbf{P}_f$ .

LEMMA 3.7. Let  $\mathbf{p}_u$  be the  $u$ -th row of the matrix  $\mathbf{P}_r$  ( $\mathbf{P}_f$ ) for  $u \in \mathcal{U}$ , where  $\mathbf{p}_u(v)$  is the probability that a random walk from  $u$  hits  $v \in \mathcal{V}_l$  (the probability that in a random spanning forest with root set  $\mathcal{V}_l$ ,  $u$  is rooted at  $v$ ). Then, we have:

(1) For  $u_1, u_2 \in \mathcal{U}$ , we have

$$r(u_1, u_2) = (\mathbf{e}_{u_1} - \mathbf{e}_{u_2})^T (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}) (\mathbf{e}_{u_1} - \mathbf{e}_{u_2}) + (\mathbf{p}_{u_1} - \mathbf{p}_{u_2})^T (\mathbf{L}/\mathcal{V}_l)^\dagger (\mathbf{p}_{u_1} - \mathbf{p}_{u_2}); \quad (9)$$

(2) For  $u \in \mathcal{U}, v \in \mathcal{V}_l$ , we have

$$r(u, v) = \mathbf{e}_u^T \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} \mathbf{e}_u + (\mathbf{p}_u - \mathbf{e}_v)^T (\mathbf{L}/\mathcal{V}_l)^\dagger (\mathbf{p}_u - \mathbf{e}_v); \quad (10)$$

(3) For  $v_1, v_2 \in \mathcal{V}_l$ , we have

$$r(v_1, v_2) = (\mathbf{e}_{v_1} - \mathbf{e}_{v_2})^T (\mathbf{L}/\mathcal{V}_l)^\dagger (\mathbf{e}_{v_1} - \mathbf{e}_{v_2}). \quad (11)$$

PROOF. Based on Lemma 3.3 and Lemma 3.5, Lemma 3.6, the  $u$ -th row of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U}\mathcal{V}_l}$  is a probability vector, which represents the probability that a  $\mathcal{V}_l$ -absorbed walk starts from  $u$  and hits the node  $v \in \mathcal{V}_l$  (the probability that in a random spanning forest with root set  $\mathcal{V}_l$ ,  $u$  is rooted at  $v$ ). We also have  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U}\mathcal{V}_l})^T = \mathbf{L}_{\mathcal{V}_l\mathcal{U}} \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$ . Thus, by substituting the matrix  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U}\mathcal{V}_l}$  by  $\mathbf{P}_r$  ( $\mathbf{P}_f$ ) and combining the terms into bilinear forms, the proof is established.  $\square$

**Novel probability interpretations of the Schur complement  $\mathbf{L}/\mathcal{V}_l$ .** Based on Lemma 3.5 and Lemma 3.6, we can derive probability interpretations for the Schur complement  $\mathbf{L}/\mathcal{V}_l = \mathbf{L}_{\mathcal{V}_l\mathcal{V}_l} - \mathbf{L}_{\mathcal{V}_l\mathcal{U}} \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U}\mathcal{V}_l}$ . Note that  $\mathbf{L}_{\mathcal{V}_l\mathcal{U}} \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U}\mathcal{V}_l}$  is a  $|\mathcal{V}_l| \times |\mathcal{V}_l|$  matrix. Then, we have  $(\mathbf{L}_{\mathcal{V}_l\mathcal{U}} \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U}\mathcal{V}_l})_{v_i v_j} = -(\mathbf{L}_{\mathcal{V}_l\mathcal{U}} \mathbf{P}_r)_{v_i v_j} = \sum_{u \in \mathcal{U}, u \in \mathcal{N}(v_i)} -p_{uv_j}^r$  for any  $v_i, v_j \in \mathcal{V}_l$ . In other words, the  $v_i, v_j$ -th element of the matrix  $-\mathbf{L}_{\mathcal{V}_l\mathcal{U}} \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U}\mathcal{V}_l}$  is the sum of the probability of a  $\mathcal{V}_l$ -absorbed random walk starting from  $u$  with  $u \in \mathcal{U}$  and  $u \in \mathcal{N}(v_i)$  that hits the landmark node  $v_j \in \mathcal{V}_l$ . Since  $\sum_{v \in \mathcal{V}_l} p_{uv}^r = 1$  for  $u \in \mathcal{U}$ , we can derive that the  $v_i$ -th row of the matrix sums up to  $-d_{\mathcal{V}_l}(v_i)$ , where  $d_{\mathcal{V}_l}(u) = d(u) - d_{\mathcal{U}}(u)$  denotes the number of nodes  $u$  that satisfy  $u \in \mathcal{V}_l$  and  $u \in \mathcal{N}(v_i)$ , and  $d_{\mathcal{U}}(u)$  denotes the number of nodes  $u$  that meet  $u \in \mathcal{U}$  and  $u \in \mathcal{N}(v_i)$ . The Schur complement  $\mathbf{L}/\mathcal{V}_l$  is also the Laplacian matrix of a Schur complement graph, which we denote by  $\mathcal{G}/\mathcal{V}_l$ . Similarly, the same statement also holds for the spanning forest probability  $p_{uv}^f$ . Thus, together with the matrix  $\mathbf{L}_{\mathcal{V}_l\mathcal{V}_l}$ , we can derive probability interpretations of the Schur complement  $\mathbf{L}/\mathcal{V}_l$ .

LEMMA 3.8. (*Probability interpretations of Schur complement*) The Schur complement  $\mathbf{L}/\mathcal{V}_l$  is the Laplacian matrix of a weighted graph  $\mathcal{G}/\mathcal{V}_l$  with node set  $\mathcal{V}_l$ , the degree of each node  $v \in \mathcal{V}_l$  is  $d(v) - \sum_{u \in \mathcal{N}(v) \cap \mathcal{U}} p_{uv}^r$ , the weight of each edge  $(v_i, v_j)$  for  $v_i, v_j \in \mathcal{V}_l$  is  $\mathbf{1}\{v_i \sim v_j\} + \sum_{u \in \mathcal{N}(v_i) \cap \mathcal{U}} p_{uv_j}^r$ , where  $\mathbf{1}\{v_i \sim v_j\}$  is an indicator variable that equals 1 when  $v_i$  and  $v_j$  are connected in the original graph, equals 0 otherwise. The statement still holds when we replace  $p_{uv}^r$  with  $p_{uv}^f$ .

PROOF. Since  $\mathbf{L}/\mathcal{V}_l = \mathbf{L}_{\mathcal{V}_l\mathcal{V}_l} - \mathbf{L}_{\mathcal{V}_l\mathcal{U}} \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U}\mathcal{V}_l}$ , based on the interpretations presented in Lemma 3.5 and Lemma 3.6, the Lemma can be established.  $\square$

Note that Lemma 3.8 shows that we can easily construct the Schur complement  $\mathbf{L}/\mathcal{V}_l$  based on the probability matrix  $\mathbf{P}_r$  (or  $\mathbf{P}_f$ ). In the following section, we will apply such probability interpretations to design an efficient sampling algorithm to compute the Schur complement.

**Discussions.** If there is only one landmark node  $v$  in  $\mathcal{V}_l$ , then since for each random walk that hits  $\mathcal{V}_l$ , it must hit  $v$ , the probability  $p_{uv}^r = 1$ . According to Lemma 3.8, the Schur complement  $\mathbf{L}/\mathcal{V}_l$  is 0. Therefore, the pseudo-inverse  $(\mathbf{L}/\mathcal{V}_l)^\dagger$  is also 0. As a result, Eq. (5) degrades to:

$$\mathbf{H} = \begin{bmatrix} \mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} & \mathbf{0} \\ \mathbf{0}^T & 0 \end{bmatrix}. \quad (12)$$

Then, the ER formula Eq. (3) and Eq. (4) can be recovered easily by Eq. (12), indicating that the result established in [20] is a special case of our results.

## 4 SINGLE-PAIR ER COMPUTATION

In this section, we propose a novel index-based approach to compute single-pair ER based on the theoretical results established in Section 3. Recall that the key to compute the ER is to compute three matrices: (i)  $\mathbf{P}_r$  (or  $\mathbf{P}_f$ ), (ii)  $(\mathbf{L}/\mathcal{V}_l)^\dagger$ , and (iii)  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$ . The basic idea of our approach is as follows. First, we pre-compute both  $\mathbf{P}_r$  (or  $\mathbf{P}_f$ ) and  $(\mathbf{L}/\mathcal{V}_l)^\dagger$  and maintain these two matrices as an index. Clearly, such an index consumes  $O(|\mathcal{V}_l|n)$  space. Since  $|\mathcal{V}_l|$  is often a small constant, the index only requires a small amount of additional space. Second, armed with the index, we can efficiently process any single-pair ER query (e.g.,  $r(s, t)$ ) by computing at most three elements of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$  (e.g.,  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{ss}$ ,  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{tt}$ , and  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$ ) based on Lemma 3.7. Below, we first propose two efficient index construction algorithms, followed by the query processing algorithms.

### 4.1 Index Construction Algorithms

Our index includes two matrices  $\mathbf{P}_r$  (or  $\mathbf{P}_f$ ) and  $(\mathbf{L}/\mathcal{V}_l)^\dagger$ . To compute the matrix  $\mathbf{P}_r$  ( $\mathbf{P}_f$ ), it is identical to calculate  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1} \mathbf{L}_{\mathcal{U}\mathcal{V}_l}$  by Lemma 3.5 and Lemma 3.6. Let  $\mathbf{p}_{v_i}$  be the  $v_i$ -th column of the matrix  $\mathbf{P}_r$ , in which an element denotes the probability that a  $\mathcal{V}_l$ -absorbed random walk starting from a node  $u_i \in \mathcal{U}$  hits  $v_i \in \mathcal{V}_l$ . We can compute  $\mathbf{p}_{v_i}$  by solving a linear system  $\mathbf{L}_{\mathcal{U}\mathcal{U}} \mathbf{p}_{v_i} = \mathbf{L}_{\mathcal{U}\mathcal{V}_l} \mathbf{e}_{v_i}$ , where  $\mathbf{e}_{v_i}$  is a  $|\mathcal{V}_l|$ -dimensional one-hot vector such that the element indexed by  $v_i$  is 1 and other elements are 0. Thus, a basic method to compute  $\mathbf{P}_r$  ( $\mathbf{P}_f$ ) requires solving  $|\mathcal{V}_l|$  linear systems. After determining  $\mathbf{P}_r$ , we can construct the Schur complement  $\mathbf{L}/\mathcal{V}_l$  by Definition 3.1, and then compute its pseudo-inverse by eigen-decomposition. Clearly, such a basic method is costly for large graphs, because it needs to solve  $|\mathcal{V}_l|$  linear systems. Below, we propose two new and efficient Monte Carlo approaches to estimate the probability matrix  $\mathbf{P}_r$  ( $\mathbf{P}_f$ ).

**Index building by  $\mathcal{V}_l$ -absorbed random walk sampling.** Based on the  $\mathcal{V}_l$ -absorbed random walk interpretation (Lemma 3.5), we can construct the index matrices  $\mathbf{P}_r$  and  $(\mathbf{L}/\mathcal{V}_l)^\dagger$  by  $\mathcal{V}_l$ -absorbed random walk sampling. The detailed algorithm is outlined in Algorithm 1. Specifically, Algorithm 1 samples a number of  $\mathcal{V}_l$ -absorbed random walk from each node  $u \in \mathcal{U}$ , and uses the proportion of walks that hit  $v \in \mathcal{V}_l$  as an estimation of  $(\mathbf{P}_r)_{uv}$  (Lines 4-6). Clearly, such an estimation is unbiased by Lemma 3.5. Similarly, by Lemma 3.8, Algorithm 1 constructs an unbiased estimator for each element of the Schur complement  $\mathbf{L}/\mathcal{V}_l$  (Lines 7-11), resulting in an estimated Laplacian matrix  $\mathbf{L}_{\mathcal{H}}$ . After that, the algorithm computes the pseudo-inverse of the matrix  $\mathbf{L}_{\mathcal{H}}$  as the final estimation of  $(\mathbf{L}/\mathcal{V}_l)^\dagger$ . Note that although  $\mathbf{L}_{\mathcal{H}}$  is an unbiased estimator of  $\mathbf{L}/\mathcal{V}_l$ ,  $\mathbf{L}_{\mathcal{H}}^\dagger$  is not necessarily

an unbiased estimator of  $(\mathbf{L}/\mathcal{V}_l)^\dagger$ . However, this does not affect the theoretical guarantee of our algorithm as analyzed below.

**Theoretical analysis of Algorithm 1.** To analyze the accuracy and complexity of our algorithm, we need an additional concept, called  $\epsilon$ -spectral sparsifier, which is defined as follows:

**Definition 4.1.** ( $\epsilon$ -spectral sparsifier) Given two undirected, un-weighted graphs  $\mathcal{G}, \mathcal{H}$  with the same number of nodes, let  $\mathbf{L}_\mathcal{H}$  and  $\mathbf{L}_\mathcal{G}$  be the Laplacian matrix of  $\mathcal{G}$  and  $\mathcal{H}$ . We call  $\mathcal{H}$  is an  $\epsilon$ -spectral sparsifier if for any vector  $\mathbf{x}$ ,  $(1 - \epsilon)\mathbf{x}^T \mathbf{L}_\mathcal{G} \mathbf{x} \leq \mathbf{x}^T \mathbf{L}_\mathcal{H} \mathbf{x} \leq (1 + \epsilon)\mathbf{x}^T \mathbf{L}_\mathcal{G} \mathbf{x}$ .

We first analyze the sample size needed for Algorithm 1 to construct  $\mathcal{H}$  as an  $\epsilon$ -spectral sparsifier of the Schur complement graph  $\mathcal{G}/\mathcal{V}_l$ , i.e. the matrix  $\mathbf{L}_\mathcal{H}$  satisfies  $(1 - \epsilon)\mathbf{x}^T \mathbf{L}_\mathcal{G} \mathbf{x} \leq \mathbf{x}^T \mathbf{L}_\mathcal{H} \mathbf{x} \leq (1 + \epsilon)\mathbf{x}^T \mathbf{L}_\mathcal{G} \mathbf{x}$ . We need a variant of matrix Chernoff bound presented in [36].

**THEOREM 4.2.** [36] Let  $\mathcal{H}_1^{(m)}, \dots, \mathcal{H}_\omega^{(m)}$  be random multi-graphs consisting random multi-edges (a multi-edge can appear in a multi-graph more than one time) over the graph  $\mathcal{G}$ . Let  $\mathcal{H}_1, \dots, \mathcal{H}_\omega$  be the weighted graphs that the weight of an edge  $(e_1, e_2)$ ,  $w_\mathcal{H}$  is the sum of the weight of all multi-edges  $(e_1, e_2)$ . Suppose that the following properties is satisfied: (i) The expectation sums to the graph  $\mathcal{G}$ , i.e.  $\sum_{i=1}^\omega E[\mathcal{H}_i] = \mathcal{G}$ ; (ii) For each  $\mathcal{H}_i$ , any edge  $e = (e_1, e_2) \in \mathcal{H}_i$  has low effective resistance in  $\mathcal{G}$ , i.e.  $w_\mathcal{H}(e)r_\mathcal{G}(e_1, e_2) \leq \frac{\epsilon^2}{O(\log n)}$ . Then, with high probability,  $\mathcal{H} = \sum_{i=1}^\omega \mathcal{H}_i$  is an  $\epsilon$ -spectral sparsifier of  $\mathcal{G}$ .

Algorithm 1 constructs the approximated Schur complement graph  $\mathcal{H}$  by adding edges to the induced graph with node set  $\mathcal{V}_l$  when the  $\mathcal{V}_l$ -absorbed random walk sampled from a node  $u \in \mathcal{U}$  and  $u$  is a neighbor of a node  $v \in \mathcal{V}_l$ . According to Lemma 3.8,  $\mathcal{H}$  is an unbiased estimator of  $\mathcal{G}/\mathcal{V}_l$ . Thus, it remains to prove that each random edge added has low effective resistance in the graph  $\mathcal{G}/\mathcal{V}_l$  that we want to approximate. We have the following Lemma:

**LEMMA 4.3.** Let  $\Delta_\mathcal{G}$  be the diameter of the graph  $\mathcal{G}$ . When the sample size  $\omega$  is larger than  $O(\frac{\Delta_\mathcal{G} |\mathcal{V}_l| \log n}{\epsilon^2})$ , the approximated Schur complement graph  $\mathcal{H}$  in Algorithm 1 is an  $\epsilon$ -spectral sparsifier of  $\mathcal{G}/\mathcal{V}_l$ .

**PROOF.** Let  $\mathcal{H}_i$  denote the random multi-graph (i.e. graph consisting several random edges that can be repeated) including the edges as follows: in the  $i$ -th round of random walks from all  $u \in \mathcal{U}$ , for any  $v_1, v_2 \in \mathcal{V}_l$ , suppose that the random walk starts from the node  $u \in N(v_i) \cap \mathcal{U}$  hits  $\mathcal{V}_l$  by  $v_2$ , then there is an edge  $(v_1, v_2)$  in  $\mathcal{H}_i$  with weight  $\frac{1}{\omega}$ . The rounds of random walks are independent of each other. Since the weight of the edge  $(v_i, v_j)$  for  $v_i, v_j \in \mathcal{V}_l$  in the Schur complement  $\mathcal{G}/\mathcal{V}_l$  is  $\sum_{u \in N(v_i) \cap \mathcal{U}} p_{uw}^r$ , suppose that  $\mathcal{H} = \sum_{i=1}^\omega \mathcal{H}_i$ , we have  $E[\mathcal{H}] = \mathcal{G}/\mathcal{V}_l$  according to Lemma 3.8. The condition (i) in Theorem 4.2 is satisfied. Moreover, for an arbitrary node  $u$ , it is the neighbor of at most  $|\mathcal{V}_l|$  nodes  $v_i \in \mathcal{V}_l$ . Thus, the weight of an edge (when summing up all the repeated multi-edges) in a random subgraph that  $u$  contributes is at most  $\frac{|\mathcal{V}_l|}{\omega}$ , when all the nodes in  $\mathcal{V}_l$  have the neighbor  $u$ . As a result, for each edge in the approximate graph  $\mathcal{H}_i$ , the weight of the edge  $w_\mathcal{H}(e)$  is at most  $\frac{|\mathcal{V}_l|}{\omega}$ , we have  $w_\mathcal{H}(e)r_\mathcal{G}(e_1, e_2) \leq \frac{|\mathcal{V}_l| \Delta_\mathcal{G}}{\omega}$ , since  $\Delta_\mathcal{G}$  is

### Algorithm 1: RwIndex

---

**Input:** Graph  $\mathcal{G}$ , landmark node set  $\mathcal{V}_l$ , sample size  $\omega$   
**Output:** The estimated index matrices  $\tilde{\mathbf{P}}_r$  and  $\mathbf{L}_\mathcal{H}^\dagger$

- 1  $\tilde{\mathbf{P}}_r \leftarrow \mathbf{0}; \mathcal{U} \leftarrow \mathcal{V} \setminus \mathcal{V}_l;$
- 2 Let  $\mathcal{H}$  be a subgraph of  $\mathcal{G}$  induced by the node set  $\mathcal{V}_l$ ;
- 3 For each edge  $(u, v) \in \mathcal{H}$ , let  $w(u, v) = 1$  be the weight of  $(u, v)$ ;
- 4 **for**  $i = 1 : \omega$  **do**
- 5     **for each node**  $u \in \mathcal{U}$  **do**
- 6         Sample a  $\mathcal{V}_l$ -absorbed random walk from  $u$ ; Suppose the terminated node is  $v \in \mathcal{V}_l$ , then  $(\tilde{\mathbf{P}}_r)_{uv} \leftarrow (\tilde{\mathbf{P}}_r)_{uv} + \frac{1}{\omega}$ ;
- 7     **for each node**  $v \in \mathcal{V}_l$  **do**
- 8         **for each neighbor node**  $u \in N(v) \cap \mathcal{U}$  **do**
- 9             **if the**  $\mathcal{V}_l$ -absorbed random walk sampled from  $u$  terminates at  $v$  **then**
- 10                 **if the edge**  $(u, v) \in \mathcal{H}$ ,  $w(u, v) \leftarrow w(u, v) + \frac{1}{\omega}$ ; **Otherwise,**  
                        create a new edge  $(u, v)$  with weight  $w(u, v) \leftarrow \frac{1}{\omega}$ ;
- 11 Let  $\mathbf{L}_\mathcal{H}$  be the Laplacian matrix of the weighted graph  $\mathcal{H}$ ;
- 12 Compute the pseudo-inverse  $\mathbf{L}_\mathcal{H}^\dagger$  by eigen-decomposition;
- 13 **return**  $\tilde{\mathbf{P}}_r, \mathbf{L}_\mathcal{H}^\dagger$

---

a natural bound of the effective resistance in  $\mathcal{G}$  [7]. Thus, when  $\omega \geq O(\frac{|\mathcal{V}_l| \Delta_\mathcal{G} \log n}{\epsilon^2})$ , the condition (ii) in Theorem 4.2 is satisfied. According to Theorem 4.2,  $\mathcal{H}$  is an  $\epsilon$ -spectral sparsifier of  $\mathcal{G}$ .  $\square$

Based on Lemma 4.3, when the sample size  $\omega \geq O(\frac{\Delta_\mathcal{G} |\mathcal{V}_l| \log n}{\epsilon^2})$ ,  $\mathcal{H}$  is an  $\epsilon$ -spectral sparsifier of  $\mathcal{G}/\mathcal{V}_l$ . We will show later in Section 4.2 that the computed  $\mathbf{P}_r$  and  $\mathbf{L}_\mathcal{H}^\dagger$  is sufficient to support an  $\epsilon$ -estimation single-pair query efficiently if  $\mathcal{H}$  is an  $\epsilon$ -spectral sparsifier of  $\mathcal{G}/\mathcal{V}_l$ . Thus, we can obtain the overall time complexity of Algorithm 1.

**LEMMA 4.4.** The expected time complexity of Algorithm 1 to obtain an  $\epsilon$ -spectral sparsifier of  $\mathcal{G}/\mathcal{V}_l$  is  $O(\frac{\tilde{\mathbf{I}}^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \tilde{\mathbf{I}} \Delta_\mathcal{G} |\mathcal{V}_l| \log n}{\epsilon^2})$ .

**PROOF.** Since each sample of random walks from all nodes  $u \in \mathcal{U}$  in Algorithm 1 costs expected time  $\sum_{u \in \mathcal{U}} h(u, \mathcal{V}_l) = \tilde{\mathbf{I}}^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \tilde{\mathbf{I}}$ . The equation holds because  $h(u, \mathcal{V}_l)$ , the hitting time from node  $u$  to the set  $\mathcal{V}_l$  can be written as  $\mathbf{e}_u^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \tilde{\mathbf{I}}$  [38], the expected running time for sampling  $\mathcal{V}_l$ -absorbed random walks separately from  $u \in \mathcal{U}$  is the sum of all such hitting times, which is  $\tilde{\mathbf{I}}^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \tilde{\mathbf{I}}$ . According to Lemma 4.3, to obtain an  $\epsilon$ -spectral sparsifier of  $\mathcal{G}/\mathcal{V}_l$ , the sample size should be larger than  $O(\frac{\Delta_\mathcal{G} |\mathcal{V}_l| \log n}{\epsilon^2})$ . Thus, the expected total time complexity is  $\omega \sum_{u \in \mathcal{U}} h(u, \mathcal{V}_l) = O(\frac{\tilde{\mathbf{I}}^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \tilde{\mathbf{I}} \Delta_\mathcal{G} |\mathcal{V}_l| \log n}{\epsilon^2})$ , which establishes the Theorem.  $\square$

In Lemma 4.4,  $\tilde{\mathbf{I}}^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \tilde{\mathbf{I}}$  is the expected running time of random walks from  $\mathcal{U}$  to hit  $\mathcal{V}_l$ . This quantity is not very large when  $\mathcal{V}_l$  is properly selected. For example, on a real-life road network PowerGrid [18],  $\tilde{\mathbf{I}}^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \tilde{\mathbf{I}}$  is  $71 \times n$  when choosing the 10 highest degree nodes as  $\mathcal{V}_l$ . In practice,  $|\mathcal{V}_l|$  is often selected as a small number (i.e.,  $|\mathcal{V}_l| \leq 100$ ). The diameter of the graph  $\mathcal{G}$  is also often not very large in real-world graphs. Thus, Algorithm 1 is efficient in practice.

### Index building by $\mathcal{V}_l$ -rooted random spanning forest sampling.

Likewise, based on the  $\mathcal{V}_l$ -rooted random spanning forest interpretation (Lemma 3.6), we can derive the index matrices  $\mathbf{P}_f$  and  $(\mathbf{L}/\mathcal{V}_l)^\dagger$  by  $\mathcal{V}_l$ -rooted random spanning forest sampling. The detailed implementation is shown in Algorithm 2. First, Algorithm 2 uniformly

draws a set of  $\mathcal{V}_l$ -rooted random spanning forests using the Wilson algorithm [45] (Lines 4-5). Note that here we need to slightly modify the original Wilson algorithm to obtain the  $\mathcal{V}_l$ -rooted random spanning forests. The Wilson algorithm is based on the concept of loop-erased random walk (LERW), in which all loops in the random walk trajectory are removed. Specifically, the Wilson algorithm first fixes a node ordering. Then, the algorithm initializes a tree  $\mathcal{T} = \{v\}$  with a root node  $v$  and performs an LERW from the first node until it hits the tree  $\mathcal{T}$  [20]. Once the LERW hits any node in  $\mathcal{T}$ , the LERW trajectory will be added into  $\mathcal{T}$ . The algorithm processes the nodes following the fixed node ordering, until all nodes are visited, a rooted spanning tree  $\mathcal{T}$  is generated uniformly [2, 20, 33, 45]. To generate a  $\mathcal{V}_l$ -rooted random spanning forest, we only need to initialize the tree as  $\mathcal{T} = \mathcal{V}_l$  and set the root set as  $\mathcal{V}_l$ , and then invoke the same LERW procedure as the Wilson algorithm. By the results established in [1, 5, 30], such a modified Wilson algorithm can generate a uniformly  $\mathcal{V}_l$ -rooted random spanning forest.

Second, for each  $u \in \mathcal{U}$ , Algorithm 2 utilizes the proportion of  $\mathcal{V}_l$ -rooted random spanning forests in which  $u$  is rooted at  $v$  as an unbiased estimator of  $(P_f)_{uv}$  based on Lemma 3.6 (Lines 4-6). Similar to Algorithm 1, Algorithm 2 creates an unbiased estimator for each element of the Schur complement  $L/\mathcal{V}_l$  based on Lemma 3.8 (Lines 7-11), and then computes the pseudo-inverse as the estimator of  $(L/\mathcal{V}_l)^\dagger$  using eigen-decomposition (Line 12). In the following, we analyze the accuracy guarantee and time complexity of Algorithm 2.

**Theoretical analysis of Algorithm 2.** First, similar to Algorithm 1, the sample size for Algorithm 2 to obtain an  $\epsilon$ -spectral sparsifier can also be derived based on Theorem 4.2.

**LEMMA 4.5.** *Let  $\Delta_G$  be the diameter of the graph  $\mathcal{G}$ , let  $d_{\mathcal{V}_l}^{max}$  be the maximum degree of the nodes in  $\mathcal{V}_l$  in the original graph  $\mathcal{G}$ . When the sample size  $\omega$  is larger than  $O(\frac{\Delta_G d_{\mathcal{V}_l}^{max} \log n}{\epsilon^2})$ , the approximate Schur complement graph  $\mathcal{H}$  in Algorithm 2 is an  $\epsilon$ -spectral sparsifier of  $\mathcal{G}/\mathcal{V}_l$ .*

**PROOF.** Let  $\mathcal{H}_i$  denote the random multi-graph including the edges as follows: in the  $i$ -th spanning forest instance  $F_i$ , for any  $v_1, v_2 \in \mathcal{V}_l$ , suppose that  $n(\mathcal{N}(v_1), v_2)$  is the number of the nodes  $u \in \mathcal{N}(v_1) \cap \mathcal{U}$  which are rooted in  $v_2$ , then there is an edge  $(v_1, v_2)$  in  $\mathcal{H}_i$  with weight  $\frac{n(\mathcal{N}(v_1), v_2)}{\omega}$ . The spanning forests are independent of each other. Since the weight of the edge  $(v_i, v_j)$  for  $v_i, v_j \in \mathcal{V}_l$  in the Schur complement  $\mathcal{G}/\mathcal{V}_l$  is  $\sum_{u \in \mathcal{N}(v_i) \cap \mathcal{U}} P_{uv_j}$ , suppose that  $\mathcal{H} = \sum_{i=1}^{\omega} \mathcal{H}_i$ , we have  $E[\mathcal{H}] = \mathcal{G}/\mathcal{V}_l$ . Thus, the condition (i) in Theorem 4.2 is satisfied. Moreover, the weight of an edge in a random subgraph that  $u$  contributes is at most  $\frac{d_{\mathcal{V}_l}^{max}}{\omega}$ , when all the neighbors of the node  $v_1$  have the root  $v_2$ . As a result, for each edge in the approximate graph  $\mathcal{H}_i$ , the weight of the edge  $w_{\mathcal{H}_i}(e)$  is at most  $\frac{d_{\mathcal{V}_l}^{max}}{\omega}$ , we have  $w_{\mathcal{H}}(e) r_{\mathcal{G}}(e_1, e_2) \leq \frac{d_{\mathcal{V}_l}^{max} \Delta_G}{\omega}$ , since  $\Delta_G$  is a natural bound of the effective resistance in  $\mathcal{G}$  [7]. Thus, when  $\omega \geq O(\frac{d_{\mathcal{V}_l}^{max} \Delta_G \log n}{\epsilon^2})$ , the condition (ii) in Theorem 4.2 is satisfied. The proof is completed.  $\square$

Similar to Algorithm 1, Algorithm 2 can also support an  $\epsilon$ -estimation query of ER if it constructs an  $\epsilon$ -spectral sparsifier, which we will discuss later. We can give the overall time complexity as follows:

### Algorithm 2: RsIndex

---

**Input:** Graph  $\mathcal{G}$ , landmark node set  $\mathcal{V}_l$ , sample size  $\omega$   
**Output:** The estimated index matrices  $\tilde{P}_f$  and  $L_{\mathcal{H}}^\dagger$

- 1  $\tilde{P}_f \leftarrow 0; \mathcal{U} \leftarrow \mathcal{V} \setminus \mathcal{V}_l;$
- 2 Let  $\mathcal{H}$  be a subgraph of  $\mathcal{G}$  induced by the node set  $\mathcal{V}_l$ ;
- 3 For each edge  $(u, v) \in \mathcal{H}$ , let  $w(u, v) = 1$  be the weight of  $(u, v)$ ;
- 4 **for**  $i = 1 : \omega$  **do**
- 5     Sample a  $\mathcal{V}_l$ -rooted random spanning forest  $F \in \mathcal{F}_{\mathcal{V}_l}$  using the Wilson algorithm  
        (with an initial root set  $\mathcal{V}_l$ ); **for each node**  $u \in \mathcal{U}$  **do**
- 6         Suppose  $v \in \mathcal{V}_l$  is the root of  $u$ ,  $(\tilde{P}_f)_{uv} \leftarrow (\tilde{P}_f)_{uv} + \frac{1}{\omega};$
- 7     **for each node**  $v \in \mathcal{V}_l$  **do**
- 8         **for each neighbor node**  $u \in \mathcal{N}(v) \cap \mathcal{U}$  **do**
- 9             **if**  $v \in \mathcal{V}_l$  is the root of  $u$  in  $F$  **then**
- 10             If the edge  $(u, v) \in \mathcal{H}$ ,  $w(u, v) \leftarrow w(u, v) + \frac{1}{\omega}$ ; Otherwise,  
                create a new edge  $(u, v)$  with weight  $w(u, v) \leftarrow \frac{1}{\omega};$
- 11 Let  $L_{\mathcal{H}}$  be the Laplacian matrix of the weighted graph  $\mathcal{H}$ ;
- 12 Compute the pseudo-inverse  $L_{\mathcal{H}}^\dagger$  by eigen-decomposition;
- 13 **return**  $\tilde{P}_f, L_{\mathcal{H}}^\dagger$

---

**LEMMA 4.6.** *The expected time complexity of Algorithm 2 to obtain an  $\epsilon$ -spectral sparsifier of  $O(\frac{\text{Tr}((I - P_{\mathcal{U}\mathcal{U}})^{-1}) \Delta_G d_{\mathcal{V}_l}^{max} (\mathcal{V}_l) \log n}{\epsilon^2})$ .*

**PROOF.** Based on the results in [27, 31], the expected running time of the loop-erased walk procedure in Algorithm 2 is  $\text{Tr}((I - P_{\mathcal{U}\mathcal{U}})^{-1})$ . According to Lemma 3.5,  $\tau_{\mathcal{V}_l}[u, u] = ((I - P_{\mathcal{U}\mathcal{U}})^{-1})_{uu}$ . Thus,  $\text{Tr}((I - P_{\mathcal{U}\mathcal{U}})^{-1}) = \sum_{u \in \mathcal{U}} \tau_{\mathcal{V}_l}[u, u]$ , which is strictly smaller than  $\tilde{1}^T (I - P_{\mathcal{U}\mathcal{U}})^{-1} \tilde{1} = \sum_{u \in \mathcal{U}} h(u, \mathcal{V}_l)$ . In other words, the expected running time of a loop-erased walk with root set  $\mathcal{V}_l$  is the sum of all expected number of passes to a node in a  $\mathcal{V}_l$ -absorbed random walk starts from  $u$  for all  $u \in \mathcal{U}$ . Thus, the total time complexity is  $\omega \sum_{u \in \mathcal{U}} \tau_{\mathcal{V}_l}[u, u] = O(\frac{\text{Tr}((I - P_{\mathcal{U}\mathcal{U}})^{-1}) \Delta_G d_{\mathcal{V}_l}^{max} \log n}{\epsilon^2})$ .  $\square$

In Lemma 4.6,  $t_{\mathcal{V}_l}^e = \text{Tr}((I - P_{\mathcal{U}\mathcal{U}})^{-1})$  is the expected running time of a loop-erased walk with root set  $\mathcal{V}_l$ , which is strictly smaller than the expected running time of  $|\mathcal{U}|$  random walks  $(\tilde{1}^T (I - P_{\mathcal{U}\mathcal{U}})^{-1} \tilde{1})$  in Algorithm 1. According to Lemma 3.5,  $\tau_{\mathcal{V}_l}[u, u] = ((I - P_{\mathcal{U}\mathcal{U}})^{-1})_{uu}$ . Thus,  $t_{\mathcal{V}_l}^e$  can also be written as  $\sum_{u \in \mathcal{U}} \tau_{\mathcal{V}_l}[u, u]$ . Since  $\tau_{\mathcal{V}_l}[u, u]$  is the expected number of passes to a node  $u$  in a random walk from  $u$  before hitting  $\mathcal{V}_l$ . In real-world networks, if  $\mathcal{V}_l$  is properly chosen, there is little probability that a node will pass itself twice. As a result,  $t_{\mathcal{V}_l}^e$  is nearly  $O(n)$ . For example, if we choose  $\mathcal{V}_l$  as the top-100 highest degree nodes,  $t_{\mathcal{V}_l}^e = 3.2 \times n$  in a real-life graph PowerGrid [18]. Thus, sampling a spanning forest is much faster than sampling random walks. For the sample size, it is hard to compare the bound of Algorithm 1 ( $\frac{|\mathcal{V}_l| \Delta_G \log n}{\epsilon^2}$ ) and the bound of Algorithm 2 ( $\frac{d_{\mathcal{V}_l}^{max} \Delta_G \log n}{\epsilon^2}$ ). In general,  $d_{\mathcal{V}_l}^{max}$  is larger than  $|\mathcal{V}_l|$ , since we prefer a small set of nodes with high degrees. However, these bounds are worst-case theoretical bounds, we find that in practice, Algorithm 2 can also obtain comparable index quality compared to Algorithm 1 when the sample size is the same. For the overall performance, we will show in experiments that Algorithm 2 is significantly better than Algorithm 1.

## 4.2 Query processing algorithms

In this section, we propose three novel algorithms to efficiently process any single-pair ER query based on our index. Note that equipped with our index, it is sufficient to compute at most three elements of the matrix  $L_{\mathcal{U}\mathcal{U}}^{-1}$  (i.e.,  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{ss}$ ,  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{tt}$ ,  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{st}$ ) to process a ER query  $r(s, t)$  based on Lemma 3.7. Below, we first propose three new algorithms to estimate an element of  $L_{\mathcal{U}\mathcal{U}}^{-1}$ , followed by the query processing algorithm.

**Estimating  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  by  $\mathcal{V}_l$ -absorbed random walk.** By Lemma 3.4, we have already shown that each element of  $L_{\mathcal{U}\mathcal{U}}^{-1}$  can be determined by the probability of a  $\mathcal{V}_l$ -absorbed random walk. Based on this, we can easily devise a  $\mathcal{V}_l$ -absorbed random walk sampling algorithm to estimate the element of  $L_{\mathcal{U}\mathcal{U}}^{-1}$ . Specifically, to estimate  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  ( $(L_{\mathcal{U}\mathcal{U}}^{-1})_{ss}$  and  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{tt}$  can be estimated similarly), we first simulate a number of  $\mathcal{V}_l$ -absorbed random walks from  $s$ , and then take the average number of visits to the node  $t$  as an estimator of  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{st}$ . Clearly, by Lemma 3.4, such an estimator is unbiased. Note that such an algorithm can also simultaneously estimate  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  for each node  $t \in \mathcal{V}$  by simulating  $\mathcal{V}_l$ -absorbed random walks from  $s$ . This is because for each  $\mathcal{V}_l$ -absorbed random walk, we can simultaneously count the number of visits for all nodes in this walk. Thus, for each node, we can derive an unbiased estimator by taking the average visits across all  $\mathcal{V}_l$ -absorbed random walks. Such a  $\mathcal{V}_l$ -absorbed random walk sampling approach that simultaneously estimates  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  for every node  $t \in \mathcal{U}$  will be used in computing the single-source ER query.

Additionally, it is worth mention that if the landmark node set  $\mathcal{V}_l$  contains only one node  $v$ , then our algorithm is exactly equivalent to the  $v$ -absorbed random walk algorithm proposed in [20]. In this sense, the  $v$ -absorbed random walk algorithm [20] is a special case of our algorithm. However, compared to the  $v$ -absorbed random walk algorithm, our algorithm is often much more efficient, because the  $\mathcal{V}_l$ -absorbed random walk (with multiple landmark nodes) typically terminates much faster than  $v$ -absorbed random walk (with only one landmark node).

**Estimating  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  by  $\mathcal{V}_l$ -absorbed push.** Push is a powerful technique to compute the personalized PageRank vector [3, 4, 8]. Recently, Liao et al. [20] propose a *push-style* algorithm, called  $v$ -absorbed push, to compute the effective resistance. However, their push algorithm only relies on a single landmark nodes. Here, we propose a new *push-style* algorithm to compute the elements of  $L_{\mathcal{U}\mathcal{U}}^{-1}$ , called  $\mathcal{V}_l$ -absorbed push, by extending their technique with multiple landmark nodes. Such a  $\mathcal{V}_l$ -absorbed push can be regarded as the deterministic version of our  $\mathcal{V}_l$ -absorbed random walk algorithm. Algorithm 3 details our  $\mathcal{V}_l$ -absorbed push algorithm.

Let  $\tau_s$  be the  $s$ -th row of  $(I - P_{\mathcal{U}\mathcal{U}})^{-1}$ . Since  $(I - P_{\mathcal{U}\mathcal{U}})^{-1}D_{\mathcal{U}\mathcal{U}}^{-1} = L_{\mathcal{U}\mathcal{U}}^{-1}$ , we can easily obtain that the  $t$ -th element of  $\tau_s$  is  $\tau_s(t) = \tau_{\mathcal{V}_l}[s, t]$  by Lemma 3.4.

As described in Algorithm 3, a  $\mathcal{V}_l$ -absorbed push procedure maintains two  $|\mathcal{U}|$ -dimensional vectors: (i)  $\tilde{\tau}_s$ , which is the estimated vector of  $\tau_s$ ; and (ii) "res", which is the residual vector. Initially,  $\tau_s$  is set as  $\mathbf{0}$  and "res" is set as  $\mathbf{e}_s$  (Line 1). Then, every time there exists a node  $u \in \mathcal{U}$  such that  $\text{res}(u) > r_{\max}$ , we will conduct push operation on the node  $u$  (Lines 2-6). Specifically, we add the residual of  $u$  to the estimated value of  $u$  (Line 3), and uniformly distribute

---

### Algorithm 3: The $\mathcal{V}_l$ -absorbed push algorithm

---

**Input:** Graph  $\mathcal{G}$ , landmark node set  $\mathcal{V}_l$ , a source node  $s$ , a threshold  $r_{\max}$   
**Output:**  $\tilde{\tau}_s$  as an approximation of  $\tau_s$

```

1  $\tilde{\tau}_s \leftarrow \mathbf{0}$ ,  $\text{res} \leftarrow \mathbf{e}_s$ ;
2 while  $\exists u \in \mathcal{U}$  such that  $\text{res}(u) > r_{\max}$  do
3    $\tilde{\tau}_s(u) \leftarrow \tilde{\tau}_s(u) + \text{res}(u)$ ;
4   for each  $w \in \mathcal{N}(u) \cap \mathcal{U}$  do
5      $\text{res}(w) \leftarrow \text{res}(w) + \frac{\text{res}(u)}{d(u)}$ ;
6    $\text{res}(u) \leftarrow 0$ ;
7 return  $\tilde{\tau}_s$ ;

```

---

the residual of  $u$  to its neighbors (Lines 4-6). If the neighbor node belongs to  $\mathcal{V}_l$ , the residual will vanish. It is easy to show (by induction) that the following invariant holds during the whole push process:

$$\tau_{\mathcal{V}_l}[s, u] = \tilde{\tau}_s(u) + \sum_{w \in \mathcal{U}} \tau_{\mathcal{V}_l}[w, u] \text{res}(w). \quad (13)$$

This is because each push operation in Algorithm 3 (Lines 3-6) does not violate the invariant equation. Note that when the landmark set  $\mathcal{V}_l$  contains only one node  $v$ , Algorithm 3 degrades to the  $v$ -absorbed push proposed by Liao et al. [20]. Likewise, our  $\mathcal{V}_l$ -absorbed push is often much faster than the  $v$ -absorbed push, as the residuals can be vanished faster with multiple landmark nodes compared to the case that has only one landmark node. Intuitively, by Eq. (13), the vector  $\tilde{\tau}_s$  output by Algorithm 3 is a good approximation of the *true* vector  $\tau_s$  if the threshold  $r_{\max}$  is small, because the final residual vector  $\text{res}$  satisfies  $\text{res}(w) \leq r_{\max}$  when the algorithm terminates. As a result, by Lemma 3.4, we can easily obtain an approximation of  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  for all  $t \in \mathcal{V}$ . Later, we will present detailed error guarantee and time complexity analysis of the  $\mathcal{V}_l$ -absorbed push algorithm.

**Combining  $\mathcal{V}_l$ -absorbed random walk and push.** To further improve the efficiency, we propose a bidirectional algorithm by combining the techniques of  $\mathcal{V}_l$ -absorbed random walk and  $\mathcal{V}_l$ -absorbed push based on the invariant equation Eq. (13). Specifically, the bidirectional algorithm first invokes a  $\mathcal{V}_l$ -absorbed push with a *relatively-large* parameter  $r_{\max}$  to obtain a *rough* estimation of  $\tilde{\tau}_s$ . Then, by Eq. (13), we can easily show that the additive error of the estimation is bounded by  $\sum_{w \in \mathcal{U}} \tau_{\mathcal{V}_l}[w, u] \text{res}(w)$ . Here the "res" vector is the residual vector output by  $\mathcal{V}_l$ -absorbed push. Subsequently, to make the estimation more accurate, we can apply  $\mathcal{V}_l$ -absorbed random walk to estimate the second term of the right hand side of Eq. (13). To achieve this, we draw a source node  $s'$  from the probability distribution  $\frac{\text{res}}{\|\text{res}\|_1}$ , and sample a  $\mathcal{V}_l$ -absorbed random walk from  $s'$ . For each node  $w$  the  $\mathcal{V}_l$ -absorbed random walk visits, we add  $\|\text{res}\|_1$  to the estimation  $\tilde{\tau}_s(w)$ . It can be shown that this is an unbiased estimation of the second term of the right hand side of Eq. (13). By properly setting the threshold  $r_{\max}$  and the sample size  $\omega$ , we can obtain an algorithm that is better than both of  $\mathcal{V}_l$ -absorbed random walk and  $\mathcal{V}_l$ -absorbed push. We will present a detailed analysis of this bidirectional algorithm in the following.

**The single-pair ER query processing algorithm.** Equipped with the above three different techniques to estimate the elements of  $L_{\mathcal{U}\mathcal{U}}^{-1}$ , we propose three algorithms for processing the single-pair ER query based on Lemma 3.7.

The detailed description of our algorithm is shown in Algorithm 4. Algorithm 4 takes two indexed matrices  $\mathbf{P}_r$  ( $\mathbf{P}_f$ ) and  $L_{\mathcal{H}}^\dagger$  as input. Let



**Algorithm 4:** The single-pair ER query processing algorithm

---

**Input:** Graph  $\mathcal{G}$ , landmark node set  $\mathcal{V}_l$ , a source node  $s$ , a target node  $t$ , two indexed matrices  $\mathbf{P}_r$  and  $\mathbf{L}_{\mathcal{H}}^\dagger$

**Output:**  $\tilde{r}(s, t)$  as an estimation of  $r(s, t)$

```

1 if  $s, t \in \mathcal{V}_l$  then
2    $\tilde{r}(s, t) \leftarrow (\mathbf{e}_s - \mathbf{e}_t)^T \mathbf{L}_{\mathcal{H}}^\dagger (\mathbf{e}_s - \mathbf{e}_t)$ ;
3 else if  $s \in \mathcal{U}, t \in \mathcal{V}_l$  then
4   Invoking an algorithm to obtain an estimation  $(\tilde{\mathbf{L}}_{\mathcal{U}\mathcal{U}}^{-1})_{ss}$  of
       $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{ss}$ ;
5    $\tilde{r}(s, t) \leftarrow (\tilde{\mathbf{L}}_{\mathcal{U}\mathcal{U}}^{-1})_{ss} + (\mathbf{p}_s - \mathbf{e}_t)^T \mathcal{L}_{\mathcal{H}}^\dagger (\mathbf{p}_s - \mathbf{e}_t)$ ;
6 else if  $s \in \mathcal{V}_l, t \in \mathcal{U}$  then
7   Invoking an algorithm to obtain an estimation  $(\tilde{\mathbf{L}}_{\mathcal{U}\mathcal{U}}^{-1})_{tt}$  of
       $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{tt}$ ;
8    $\tilde{r}(s, t) \leftarrow (\tilde{\mathbf{L}}_{\mathcal{U}\mathcal{U}}^{-1})_{tt} + (\mathbf{p}_t - \mathbf{e}_s)^T \mathcal{L}_{\mathcal{H}}^\dagger (\mathbf{p}_t - \mathbf{e}_s)$ ;
9 else if  $s, t \in \mathcal{U}$  then
10  Invoking an algorithm to obtain an estimation  $(\tilde{\mathbf{L}}_{\mathcal{U}\mathcal{U}}^{-1})_{ss}$ ,
       $(\tilde{\mathbf{L}}_{\mathcal{U}\mathcal{U}}^{-1})_{tt}$  and  $(\tilde{\mathbf{L}}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  of  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{ss}$ ,  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{tt}$  and  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$ 
      respectively;
11   $\tilde{r}(s, t) \leftarrow (\mathbf{e}_s - \mathbf{e}_t)^T \tilde{\mathbf{L}}_{\mathcal{U}\mathcal{U}}^{-1} (\mathbf{e}_s - \mathbf{e}_t) + (\mathbf{p}_s - \mathbf{p}_t)^T \mathcal{L}_{\mathcal{H}}^\dagger (\mathbf{p}_s - \mathbf{p}_t)$ ;
12 return  $\tilde{r}(s, t)$ 

```

---

$\mathbf{p}_{u_i}$  be the  $u_i$ -th row of  $\mathbf{P}_r$ . By Lemma 3.7, there are four cases depending on which set  $s$  and  $t$  belongs to. When  $s, t \in \mathcal{V}_l$ , i.e., both of them belong to the landmark node set, the algorithm computes an estimator  $\tilde{r}(s, t)$  using the indexed matrix  $\mathbf{L}_{\mathcal{H}}^\dagger$  within  $O(1)$  time (Lines 1-2). When one node ( $s$  or  $t$ ) belongs to  $\mathcal{U}$ , the other node is located in  $\mathcal{V}_l$ . Algorithm 4 first invokes an algorithm devised above (e.g.,  $\mathcal{V}_l$ -absorbed random walk, or  $\mathcal{V}_l$ -absorbed push, or the bidirectional algorithm) to estimate only one element of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$ , and then compute  $\tilde{r}(s, t)$  by Lemma 3.7 (Lines 3-8). When  $s, t \in \mathcal{U}$ , the situation becomes more complicated. Specifically, in this case, Algorithm 4 needs to estimate three elements of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$  by Lemma 3.7 (Lines 9-11). For convenience, we refer to Algorithm 4 that is equipped with  $\mathcal{V}_l$ -absorbed random walk,  $\mathcal{V}_l$ -absorbed push, and the bidirectional algorithm as RW, Push and Bipush respectively. Below, we present comprehensive theoretical analyses of these algorithms.

**Theoretical analysis of the query processing algorithms.** First, we give the time complexity for the three techniques to achieve an  $\epsilon$ -estimation of  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$ . We will utilize the following Hoeffding lemma in our analysis:

LEMMA 4.7. (Hoeffding lemma [11]) Let  $X_1, \dots, X_\omega$  be random variables with expectation  $E[X]$ . Suppose that  $X$  can be bounded by  $a \leq X \leq b$ . Then, with high probability (e.g., the failure probability is smaller than 0.01), when  $\omega \geq \frac{(a-b)^2}{\epsilon^2}$ , we have  $|\frac{1}{\omega} \sum_{i=1}^\omega -E[X]| \leq \epsilon$ .

Notice that for our  $\mathcal{V}_l$ -absorbed sampling algorithm, although it usually terminates very fast, in the worst case the random walk length can be infinite. Thus, it is hard to directly analyze the algorithm. However, if we truncate the random walks with length  $\sigma$ , we can give a rigorous analysis of the modified algorithm, by separately analyzing the sampling error and the truncated error.

LEMMA 4.8. Let  $\lambda_{\mathcal{V}_l}$  be the spectral radius of the probability transition matrix  $\mathbf{P}_{\mathcal{U}\mathcal{U}}$ ,  $\sigma_{\mathcal{V}_l} = O(\frac{\log(1/(\epsilon(1-\lambda_{\mathcal{V}_l})))}{\log(1/\lambda_{\mathcal{V}_l})})$ , the  $\mathcal{V}_l$ -absorbed random walk sampling algorithm can estimate  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  with  $\epsilon$  absolute error in time  $O(\frac{\sigma_{\mathcal{V}_l}^3}{\epsilon^2})$ .

PROOF. For an arbitrary element  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$ , recall that  $\mathbf{L}_{\mathcal{U}\mathcal{U}} = (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \mathbf{D}_{\mathcal{U}\mathcal{U}}$ . We have  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st} = \frac{((\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1})_{st}}{d_t}$ . By sampling random walks with length at most  $l$  from  $s$  to hit  $\mathcal{V}_l$ , suppose that  $X$  is the random variable that represents the number of passes to the node  $t$ , we have  $E[\frac{X}{d_t}] = (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$ . The total error consists two parts: (i) the sampling error  $\epsilon_1$  and (ii) the truncation error  $\epsilon_2$ . For the (i) sampling error, since  $X \leq \sigma$ , according to Hoeffding bound, when  $\omega \geq O(\frac{\sigma^2}{\epsilon_1^2})$ , with high probability, the first part error is less than  $\epsilon_1$ . For the (ii) truncated error, we bound the error  $\epsilon_2 = \mathbf{e}_s^T \sum_{k=\sigma+1}^\infty \mathbf{P}_{\mathcal{U}\mathcal{U}}^k \mathbf{e}_t$ . We have  $\|\mathbf{e}_s^T \sum_{k=\sigma+1}^\infty \mathbf{P}_{\mathcal{U}\mathcal{U}}^k \mathbf{e}_t\|_2 \leq \|\mathbf{e}_s\|_2 \frac{\lambda_{\mathcal{V}_l}^{\sigma+1}}{1-\lambda_{\mathcal{V}_l}} \|\mathbf{e}_t\|_2 = \frac{\lambda_{\mathcal{V}_l}^{\sigma+1}}{1-\lambda_{\mathcal{V}_l}}$ . Thus, when  $\sigma \geq O(\frac{\log(1/(\epsilon_2(1-\lambda_{\mathcal{V}_l})))}{\log(1/\lambda_{\mathcal{V}_l})})$ , the second part error can be bound by  $\epsilon_2$ . Let  $\epsilon_1 = \epsilon_2 = \frac{\epsilon}{2}$ , the total time complexity for the  $\mathcal{V}_l$ -absorbed random walk sampling algorithm to achieve an  $\epsilon$  absolute error of  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  is  $\omega\sigma = O(\frac{\log^3(1/(\epsilon(1-\lambda_{\mathcal{V}_l})))}{\epsilon^2 \log^3(1/\lambda_{\mathcal{V}_l})}) = O(\frac{\sigma_{\mathcal{V}_l}^3}{\epsilon^2})$ , which establishes the Lemma.  $\square$

Note that  $\lambda_{\mathcal{V}_l}$  is often near to 1 [29]. Thus, by Lemma 4.8,  $\sigma_{\mathcal{V}_l}$  is a small number, which results in an efficient sublinear algorithm. Recall that the  $v$ -absorbed walk sampling algorithm proposed in [20] is a special case of the  $\mathcal{V}_l$ -absorbed walk sampling algorithm when  $\mathcal{V}_l = \{v\}$ . Let  $\lambda_v$  denote the spectral radius of  $\mathbf{P}_v$ ,  $\sigma_v = O(\frac{\log(1/(\epsilon(1-\lambda_v)))}{\log(1/\lambda_v)})$ . The time complexity for the  $v$ -absorbed walk sampling algorithm is  $O(\frac{\sigma_v^3}{\epsilon^2})$ . According to the Cauchy Interlacing theorem [7],  $\lambda_{\mathcal{V}_l^{(1)}} \leq \lambda_{\mathcal{V}_l^{(2)}}$  if  $\mathcal{V}_l^{(1)} \subseteq \mathcal{V}_l^{(2)}$ . Thus, we have  $\sigma_{\mathcal{V}_l} \leq \sigma_v$ , the theoretical bound of the algorithm is strictly better than [20]. For example, on a real-life road network PowerGrid,  $\sigma_v = 1.5 \times 10^5$  while  $\sigma_{\mathcal{V}_l} = 1321$  if  $\mathcal{V}_l$  is the top-100 highest degree nodes. Next, we derive a theoretical bound for the  $\mathcal{V}_l$ -absorbed push algorithm.

LEMMA 4.9. Let  $\bar{d} = \frac{2m}{n}$  be the average degree, the  $\mathcal{V}_l$ -absorbed push algorithm can estimate  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  with  $\epsilon$  absolute error in time  $O(\frac{\bar{d}h(s, \mathcal{V}_l)h(t, \mathcal{V}_l)}{\epsilon})$ .

PROOF. For an arbitrary element  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$ , the  $\mathcal{V}_l$ -absorbed push algorithm with source node  $s$  estimates the  $s$ -th row of the matrix  $(\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1}$ , then divided the elements by degree. In each push operation, the algorithm will add at most  $r_{\max}$  to the estimation vector  $\tau$ . During the whole  $\mathcal{V}_l$ -absorbed push algorithm, it will add at most  $\mathbf{e}_s^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \mathbf{1}$  to the vector  $\tau$ . Thus, there are at most  $O(\frac{\mathbf{e}_s^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \mathbf{1}}{r_{\max}})$  push operations. The amortized time of each push operation is  $\bar{d}$  [25], which results in the time complexity  $O(\frac{\bar{d} \mathbf{e}_s^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \mathbf{1}}{r_{\max}}) = O(\frac{\bar{d}h(s, \mathcal{V}_l)}{r_{\max}})$ , where  $h(s, \mathcal{V}_l)$  is the hitting time from  $s$  to  $\mathcal{V}_l$ . At the same time, according to Eq. (13), the additive error of the  $t$ -th element is bounded by  $r_{\max} \mathbf{1}^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \mathbf{e}_t$ . Since

we have  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st} = \frac{((\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1})_{st}}{d_t}$ , the additive error of  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  can be bounded by  $r_{\max} \frac{\tilde{\mathbf{1}}^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \mathbf{e}_t}{d_t} \leq r_{\max} \mathbf{e}_t^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \tilde{\mathbf{1}} = r_{\max} h(t, \mathcal{V}_l)$ . By setting  $r_{\max} = \frac{\epsilon}{h(t, \mathcal{V}_l)}$ , the total time complexity to achieve an  $\epsilon$  absolute error is  $O(\frac{dh(s, \mathcal{V}_l)h(t, \mathcal{V}_l)}{\epsilon})$  for estimating  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$ . Thus, the lemma is established.  $\square$

Based on the results of Lemma 4.9, the time complexity of the bidirectional algorithm is closely related to the hitting time from the query node to  $\mathcal{V}_l$ . Notice that the push algorithm proposed in [20] is  $\frac{dh(s, v)h(t, v)}{\epsilon}$  based on our results. The hitting time  $h(u, \mathcal{V}_l)$  is intuitively smaller than  $h(u, v)$  for  $u \in \mathcal{U}$ , since there are more nodes to hit. Thus, the time complexity of our algorithm is strictly lower than that of the algorithm in [20]. For example, on a real-life network PowerGrid,  $\bar{h}_v = \sum_{u \in \mathcal{V} \setminus \{v\}} h(u, v)$  is  $1.3 \times 10^4$  while  $\bar{h}_{\mathcal{V}_l} = \sum_{u \in \mathcal{U}} h(u, \mathcal{V}_l)$  is 71 if  $\mathcal{V}_l$  is the top-100 highest degree nodes. Next, we present theoretical analysis for the bidirectional algorithm.

LEMMA 4.10. *The bidirectional algorithm can estimate  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  with  $\epsilon$  absolute error in time  $O(\frac{(ndh(s, \mathcal{V}_l))^{\frac{2}{3}} \sigma_{\mathcal{V}_l}}{\epsilon^{\frac{2}{3}}})$ .*

PROOF. The bidirectional algorithm first applies a  $\mathcal{V}_l$ -absorbed push procedure, then performs random walks. Similar to Lemma 4.8, we can analyze the algorithm when the random walks are truncated by a length  $\sigma$ , and the total error of the algorithm is the sum of the sampling error and the truncated error. When the push procedure terminates, all the elements of the residual vector "res" are less than  $r_{\max}$ . According to Eq. (13), let  $X$  be the random variable that a  $\mathcal{V}_l$ -absorbed walk starts from a node  $s'$  sampled from the distribution  $\frac{\text{res}}{\|\text{res}\|_1}$  passes node  $t$  for  $X$  times, we have  $E[\|\text{res}\|_1 X] = \tau_{\mathcal{V}_l}[s, t] - \bar{\tau}_s(t)$ . Since there are less than  $n$  elements in the vector "res",  $\|\text{res}\|_1$  can be bounded by  $nr_{\max}$ . Similar to the proof of Lemma 4.8,  $X$  can be bounded by  $\sigma$ ,  $\|\text{res}\|_1 X$  can be bounded by  $nr_{\max} \sigma$ , we have that when  $\omega \geq \frac{n^2 r_{\max}^2 \sigma^2}{\epsilon_1^2}$ , the sampling error is smaller than  $\epsilon_1$ . Likewise, the truncated error  $\epsilon_2 = (\frac{\text{res}}{\|\text{res}\|_1})^T \sum_{k=\sigma+1}^{\infty} \mathbf{p}_{\mathcal{U}\mathcal{U}}^k \mathbf{e}_t$ . We have  $\|(\frac{\text{res}}{\|\text{res}\|_1})^T \sum_{k=\sigma+1}^{\infty} \mathbf{p}_{\mathcal{U}\mathcal{U}}^k \mathbf{e}_t\|_2 \leq \|(\frac{\text{res}}{\|\text{res}\|_1})\|_2 \frac{\lambda_{\mathcal{V}_l}^{\sigma+1}}{1 - \lambda_{\mathcal{V}_l}} \|\mathbf{e}_t\|_2 = \frac{\lambda_{\mathcal{V}_l}^{\sigma+1}}{1 - \lambda_{\mathcal{V}_l}}$ . Thus, when  $\sigma \geq O(\frac{\log(1/(\epsilon_2(1 - \lambda_{\mathcal{V}_l})))}{\log(1/\lambda_{\mathcal{V}_l})}) = \sigma_{\mathcal{V}_l}$ , the second part error can be bound by  $\epsilon_2$ . Let  $\epsilon_1 = \epsilon_2 = \frac{\epsilon}{2}$ , the total absolute error is smaller than  $\epsilon$  with high probability. According to Lemma 4.9, the push procedure costs time  $O(\frac{dh(s, \mathcal{V}_l)}{r_{\max}})$ . The sampling procedure costs time  $O(\frac{n^2 r_{\max}^2 \sigma_{\mathcal{V}_l}^3}{\epsilon^2})$ . Thus, the total time complexity is minimized when  $r_{\max} = \sqrt[3]{\frac{dh(s, \mathcal{V}_l)}{2n^2 r_{\max}^2 \sigma_{\mathcal{V}_l}^3}}$ , the total time complexity is thereby  $O(\frac{(ndh(s, \mathcal{V}_l))^{\frac{2}{3}} \sigma_{\mathcal{V}_l}}{\epsilon^{\frac{2}{3}}})$ . The proof is completed.  $\square$

As the bidirectional algorithm is a combination of the  $\mathcal{V}_l$ -absorbed walk sampling algorithm and the  $\mathcal{V}_l$ -absorbed push algorithm, the dependency of the bound on  $\sigma_{\mathcal{V}_l}$ ,  $h(u, \mathcal{V}_l)$  and  $\epsilon$  is better than both of the algorithms. Although there is an additional  $n$  term, the bound presented in Lemma 4.10 is only a worst-case theoretical bound. We find that in experiments, the bidirectional algorithm performs

much better than both  $\mathcal{V}_l$ -absorbed random walk sampling and  $\mathcal{V}_l$ -absorbed push. Moreover, as discussed above, the time complexity bound is also strictly better than the Bipush algorithm proposed in [20], since  $\sigma_{\mathcal{V}_l}$ ,  $h(u, \mathcal{V}_l)$  are both smaller.

Recall that in Algorithm 4, the error of the single-pair query algorithm includes both the error of building index and the error of estimating  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$ . Next, we first show the time complexity for Algorithm 1 and Algorithm 2 to achieve an  $\epsilon$ -absolute error of index building. Then, we combine the error of estimating  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$  to give the total error of the query processing algorithms. Here, the error of index building is different in four cases, depending on the query node  $s$  and  $t$ : (i) If  $s, t \in \mathcal{V}_l$ , it is the error of estimating  $(\mathbf{e}_s - \mathbf{e}_t)^T \mathbf{L}_{\mathcal{H}}^{\dagger} (\mathbf{e}_s - \mathbf{e}_t)$ ; (ii) If  $s \in \mathcal{U}$ ,  $t \in \mathcal{V}_l$ , it is the error of estimating  $(\mathbf{p}_s - \mathbf{e}_t)^T \mathbf{L}_{\mathcal{H}}^{\dagger} (\mathbf{p}_s - \mathbf{e}_t)$ ; (iii) If  $s \in \mathcal{V}_l$ ,  $t \in \mathcal{U}$ , it is the error of estimating  $(\mathbf{p}_t - \mathbf{e}_s)^T \mathbf{L}_{\mathcal{H}}^{\dagger} (\mathbf{p}_t - \mathbf{e}_s)$ ; (iv) If  $s, t \in \mathcal{V}_l$ , it is the error of estimating  $(\mathbf{p}_s - \mathbf{p}_t)^T \mathbf{L}_{\mathcal{H}}^{\dagger} (\mathbf{p}_s - \mathbf{p}_t)$ . Since all the four cases deal with the bilinear form of the pseudo-inverse of the Schur complement  $\mathbf{L}/\mathcal{V}_l$ , and we have shown the time complexity for Algorithm 2 and Algorithm 2 to output an  $\epsilon$ -spectral sparsifier, we will utilize the following lemma proposed in [14].

LEMMA 4.11. *(Property of an  $\epsilon$ -spectral sparsifier [14]) If  $\mathcal{H}$  is an  $\epsilon$ -spectral sparsifier of  $\mathcal{G}$ , let  $\mathbf{L}_{\mathcal{G}}$  and  $\mathbf{L}_{\mathcal{H}}$  be the Laplacian matrix of  $\mathcal{G}$  and  $\mathcal{H}$ , respectively. Then, for any vector  $\mathbf{x}$ , we have:  $(1 - \epsilon) \mathbf{x}^T \mathbf{L}_{\mathcal{G}}^{\dagger} \mathbf{x} \leq \mathbf{x}^T \mathbf{L}_{\mathcal{H}}^{\dagger} \mathbf{x} \leq (1 + \epsilon) \mathbf{x}^T \mathbf{L}_{\mathcal{G}}^{\dagger} \mathbf{x}$ .*

Notice that Lemma 4.11 is not enough for the error of index building, since in cases (ii)-(iv), the probability vector  $\mathbf{p}_u$  for  $u \in \mathcal{U}$  is also an estimated value. Below, we show the time complexity for Algorithm 1 and Algorithm 2 to achieve an  $\epsilon$ -absolute error of index building.

LEMMA 4.12. *Algorithm 1 can achieve an  $\epsilon$ -absolute error for index building in time  $O(\frac{\tilde{\mathbf{1}}^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \tilde{\mathbf{1}} \Delta_{\mathcal{G}}^3 |\mathcal{V}_l| \log n}{\epsilon^2})$ .*

PROOF. As we have proved in Lemma 4.4, Algorithm 1 can output an  $\frac{\epsilon}{\Delta_{\mathcal{G}}}$ -spectral sparsifier  $\mathcal{H}$  in time  $O(\frac{\tilde{\mathbf{1}}^T (\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1} \tilde{\mathbf{1}} \Delta_{\mathcal{G}}^3 |\mathcal{V}_l| \log n}{\epsilon^2})$ . Based on Theorem 4.11, for any  $\mathbf{x}$ , we have:  $(1 - \epsilon) \mathbf{x}^T (\mathbf{L}/\mathcal{V}_l)^{\dagger} \mathbf{x} \leq \mathbf{x}^T \mathbf{L}_{\mathcal{H}}^{\dagger} \mathbf{x} \leq (1 + \epsilon) \mathbf{x}^T (\mathbf{L}/\mathcal{V}_l)^{\dagger} \mathbf{x}$ . Below, we discuss the four cases of the error index building. (i) If  $s, t \in \mathcal{V}_l$ , the error of index building is  $(\mathbf{e}_s - \mathbf{e}_t)^T \mathbf{L}_{\mathcal{H}}^{\dagger} (\mathbf{e}_s - \mathbf{e}_t)$ . The output of the ER value in Algorithm 4 is also  $(\mathbf{e}_s - \mathbf{e}_t)^T \mathbf{L}_{\mathcal{H}}^{\dagger} (\mathbf{e}_s - \mathbf{e}_t)$ . Thus, let  $\mathbf{x} = \mathbf{e}_s - \mathbf{e}_t$ , it directly holds that  $|\hat{r}(s, t) - r(s, t)| \leq \frac{\epsilon}{\Delta_{\mathcal{G}}} r(s, t) \leq \frac{\epsilon}{\Delta_{\mathcal{G}}} \Delta_{\mathcal{G}} = \epsilon$ . Thus, the absolute error of index building in case (i) is smaller than  $\epsilon$ . (ii) When  $s \in \mathcal{U}$ ,  $t \in \mathcal{V}_l$ ,  $r(s, t) = (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{ss} + (\mathbf{e}_t - \mathbf{p}_s)^T \mathbf{L}_{\mathcal{H}}^{\dagger} (\mathbf{e}_t - \mathbf{p}_s)$ . Since  $\mathbf{p}_s$  is the  $|\mathcal{V}_l|$ -dimensional probability vector that the  $v$ -th element represents the probability of a random walk starts from  $s$  hits  $\mathcal{V}_l$  by  $v$ . Thus, suppose that the random walk starts from  $s$  hits  $\mathcal{V}_l$  by  $v$ , let  $X_v$  denote the random variable that  $X_v = (\mathbf{e}_t - \mathbf{e}_v)^T \mathbf{L}_{\mathcal{H}}^{\dagger} (\mathbf{e}_t - \mathbf{e}_v) = r_{\mathcal{H}}(t, v)$ , then we have  $E[X_v] = \sum_{v \in \mathcal{V}_l} \mathbf{p}(v) r_{\mathcal{H}}(t, v) = \sum_{v \in \mathcal{V}_l} \mathbf{p}(v) r(t, v) = (\mathbf{e}_t - \mathbf{p}_s)^T \mathbf{L}_{\mathcal{H}}^{\dagger} (\mathbf{e}_t - \mathbf{p}_s)$ . Since  $\mathcal{H}$  is an  $\epsilon$ -spectral sparsifier of  $\mathcal{G}/\mathcal{V}_l$ , we have  $r_{\mathcal{H}}(t, v) = (\mathbf{e}_t - \mathbf{e}_v)^T \mathbf{L}_{\mathcal{H}}^{\dagger} (\mathbf{e}_t - \mathbf{e}_v) \leq (1 + \epsilon) r(t, v) \leq (1 + \epsilon) \Delta_{\mathcal{G}}$ . In other words, the random variable  $X_v$  is bounded by  $(1 + \epsilon) \Delta_{\mathcal{G}}$ . Thus, according to the Hoeffding lemma, with high probability, when

the sampling size in Algorithm 1 satisfies  $\omega \geq (1 + \epsilon)^2 \Delta_G^2 \frac{\log n}{\epsilon^2}$ , we have  $|(e_t - p_s)^T L_{\mathcal{H}}^\dagger (e_t - p_s) - (e_t - p_s)^T L_{\mathcal{V}_l}^\dagger (e_t - p_s)| \leq \epsilon$ . Thus, for case (ii), we can also achieve an  $\epsilon$ -absolute error in the desired complexity. (iii) When  $s \in \mathcal{V}_l$ ,  $t \in \mathcal{U}$ , the situation is similar to (ii). (iv) When  $s, t \in \mathcal{U}$ ,  $r(s, t) = (e_s - e_t)^T L_{\mathcal{U}\mathcal{U}}^{-1} (e_s - e_t) + (p_s - p_t)^T L_{\mathcal{H}}^\dagger (p_s - p_t)$ . Again, suppose that the random walk starts from  $s$  hits  $\mathcal{V}_l$  by  $v_1$ , the random walk starts from  $t$  hits  $\mathcal{V}_l$  by  $v_2$ , let  $X$  denote the random variable that  $X = (e_{v_1} - e_{v_2})^T L_{\mathcal{H}}^\dagger (e_{v_1} - e_{v_2}) = r_{\mathcal{H}}(v_1, v_2)$ , then we have  $E[X] = \sum_{v_1, v_2 \in \mathcal{V}_l} p(v_1) p(v_2) r_{\mathcal{H}}(v_1, v_2) = (p_{v_1} - p_{v_2})^T (\mathcal{G}/\mathcal{V}_l)^\dagger (p_{v_1} - p_{v_2})$ . With high probability,  $r_{\mathcal{H}}(v_1, v_2) \leq (1 + \epsilon) r(v_1, v_2) \leq (1 + \epsilon) \Delta_G$ . The random variable  $X$  can also be bounded by  $(1 + \epsilon) \Delta_G$ . Thus, according to the Hoeffding lemma, with high probability, when  $\omega \geq (1 + \epsilon)^2 \Delta_G^2 \frac{\log n}{\epsilon^2} = O(\Delta_G^2 \frac{\log n}{\epsilon^2})$ , we have  $|(p_s - p_t)^T L_{\mathcal{H}}^\dagger (p_s - p_t) - (e_s - p_t)^T L_{\mathcal{V}_l}^\dagger (p_s - p_t)| \leq \epsilon$ . Thus, for case (iv), we can also achieve an  $\epsilon$ -absolute error in the desired complexity. The proof is completed.  $\square$

LEMMA 4.13. *Algorithm 2 can achieve an  $\epsilon$ -absolute error for index building in time  $O(\frac{\text{Tr}((I - P_{\mathcal{U}\mathcal{U}})^{-1}) d_{\mathcal{V}_l}^{\max} \Delta_G^2 |\mathcal{V}_l| \log n}{\epsilon^2})$ .*

PROOF. According to Lemma 4.6, Algorithm 2 can output an  $\epsilon$ -spectral sparsifier in time  $O(\frac{\text{Tr}((I - P_{\mathcal{U}\mathcal{U}})^{-1}) d_{\mathcal{V}_l}^{\max} |\mathcal{V}_l| \log n}{\epsilon^2})$ . Similar to the proof of Lemma 4.12, Algorithm 2 can achieve an  $\epsilon$ -absolute error of index building in time  $O(\frac{\text{Tr}((I - P_{\mathcal{U}\mathcal{U}})^{-1}) d_{\mathcal{V}_l}^{\max} \Delta_G^2 |\mathcal{V}_l| \log n}{\epsilon^2})$ . The proof is completed.  $\square$

Based on Lemma 4.12 and Lemma 4.13, we can give the query time complexity for all the three algorithms to achieve an  $\epsilon$ -absolute error of  $r(s, t)$  based on the fact that we have built the index with an  $\epsilon$ -absolute error.

LEMMA 4.14. *The query complexity of RW to achieve an  $\epsilon$ -absolute error for  $r(s, t)$  is  $O(\frac{\sigma_{\mathcal{V}_l}^3}{\epsilon^2})$ .*

PROOF. Based on Lemma 4.8, RW can output an estimation of  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  with  $\epsilon$ -absolute error in time  $O(\frac{\sigma_{\mathcal{V}_l}^3}{\epsilon^2})$ . Since there are at most four elements of  $L_{\mathcal{U}\mathcal{U}}^{-1}$  to estimate, the query complexity of RW is  $O(\frac{\sigma_{\mathcal{V}_l}^3}{\epsilon^2})$  if we use the  $O(\cdot)$  notation to ignore the constant terms.  $\square$

LEMMA 4.15. *Let  $h_{\mathcal{V}_l} = \max\{h(s, \mathcal{V}_l), h(t, \mathcal{V}_l)\}$ . The query complexity of Push to achieve an  $\epsilon$ -absolute error for  $r(s, t)$  is  $O(\frac{\bar{d} h_{\mathcal{V}_l}^2}{\epsilon})$ .*

PROOF. Based on Lemma 4.9, Push can output an estimation of  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  with  $\epsilon$ -absolute error in time  $O(\frac{\bar{d} h(s, \mathcal{V}_l) h(t, \mathcal{V}_l)}{\epsilon})$ . Since there are at most four elements of  $L_{\mathcal{U}\mathcal{U}}^{-1}$  to estimate, the total time complexity can also be bounded by  $O(\frac{\bar{d} h_{\mathcal{V}_l}^2}{\epsilon})$  if we use the  $O(\cdot)$  notation to ignore the constant terms.  $\square$

LEMMA 4.16. *The query complexity of Bipush algorithm to achieve an  $\epsilon$ -absolute error for  $r(s, t)$  is  $O(\frac{(n \bar{d} h_{\mathcal{V}_l})^{\frac{2}{3}} \sigma_{\mathcal{V}_l}}{\epsilon^{\frac{2}{3}}})$ .*

PROOF. Based on Lemma 4.9, Push can output an estimation of  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  with  $\epsilon$ -absolute error in time  $O(\frac{(n \bar{d} h(s, \mathcal{V}_l))^{\frac{2}{3}} \sigma_{\mathcal{V}_l}}{\epsilon^{\frac{2}{3}}})$ . Since there are at most four elements of  $L_{\mathcal{U}\mathcal{U}}^{-1}$  to estimate, the total time complexity is  $O(\frac{(n \bar{d} h_{\mathcal{V}_l})^{\frac{2}{3}} \sigma_{\mathcal{V}_l}}{\epsilon^{\frac{2}{3}}})$  if we use the  $O(\cdot)$  notation to ignore the constant terms.  $\square$

As discussed above, the bounds of the three query algorithms are strictly better than the RW, Push and Bipush algorithm proposed in [20], which are all special cases of our algorithms. Moreover, the bound is also better than the other SOTA method GEER [47]. Let  $\lambda_1 \leq \dots \leq \lambda_n = 1$  be the eigenvalues of the probability transition matrix  $P$ ,  $\lambda = \max\{|\lambda_1|, \lambda_{n-1}\}$ . Then, suppose that  $\sigma = O(\frac{\log(1/\epsilon(1-\lambda))}{\log(1/\lambda)})$ , the time complexity bound for GEER to achieve an  $\epsilon$ -absolute error is  $O(\frac{\sigma^3}{\epsilon^2})$ , which is similar to the bound  $O(\frac{\sigma_{\mathcal{V}_l}^3}{\epsilon^2})$  of RW. Although it is hard to compare these bounds in theory. We find that in practice, when  $\mathcal{V}_l$  is properly selected,  $\sigma_{\mathcal{V}_l}$  can be significantly smaller than  $\sigma$ . For example, on a real-life road graph PowerGrid [18],  $\sigma_{\mathcal{V}_l} = 1321$  if we select  $\mathcal{V}_l$  as the top-100 highest degree nodes, while  $\sigma$  is as large as  $3.3 \times 10^4$ . Thus, the theoretical bound of our query processing algorithms RW to achieve an  $\epsilon$ -estimation of effective resistance is strictly better than the SOTA methods. Push and Bipush can be better than RW further as we have discussed. We will show in experiments that the empirical performance of our query algorithms are significantly better than the SOTA methods.

**Heuristic landmark nodes set selection.** Recall that the efficiency of our algorithms is closely related to the landmark set nodes. Intuitively, if every node in  $\mathcal{U}$  can easily hit any one landmark node in  $\mathcal{V}_l$  by a random walk, then our algorithms can be very efficient, because both the  $\mathcal{V}_l$ -absorbed random walk and  $\mathcal{V}_l$ -absorbed push can terminate very quickly in this case. Based on this intuition, we present two heuristic landmark node selection methods as follows. The first approach is based on the degree of the nodes. Specifically, we first sort the nodes in a non-increasing ordering by their degrees. Then, we add the highest-degree node into  $\mathcal{V}_l$  and delete all the neighbors of this selected node from the graph. In the remaining graph, we iteratively select the highest-degree nodes according to this rule, until we get  $|\mathcal{V}_l|$  nodes. Since we always select the highest-degree nodes, this method can often obtain a good landmark node set such that other nodes can easily hit the landmark nodes. Instead of using the degrees, our second method is based on the PageRank values. The selection procedure is same as the degree-based method, i.e., it iteratively selects the highest PageRank node from the remaining graph (when a node is selected, all of its neighbors will be deleted). Our experiments confirm that these methods are very effective in practice.

## 5 SINGLE-SOURCE ER COMPUTATION

A straightforward single-source ER computation approach is to invoke the single-pair algorithm for  $n$  times. However, even equipped with the index proposed in Section 4.1, such a straightforward method may still need to compute  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{ss}$ ,  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{st}$ ,  $(L_{\mathcal{U}\mathcal{U}}^{-1})_{tt}$  for all  $t \in \mathcal{V}$ , which is very costly for large graphs.

**Algorithm 5: RsflIndex-SS**


---

**Input:** Graph  $\mathcal{G}$ , landmark node set  $\mathcal{V}_l$ , sample size  $\omega$

**Output:**  $\tilde{\mathbf{P}}_f, \mathbf{L}_{\mathcal{H}}^\dagger, (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  for each  $u \in \mathcal{U}$

- 1  $\tilde{\mathbf{P}}_f, \mathbf{L}_{\mathcal{H}}^\dagger \leftarrow \text{RsflIndex}(\mathcal{G}, \mathcal{V}_l, \omega)$ ;
- 2 Let  $\mathcal{F}_{\mathcal{V}_l}$  be the set of  $\mathcal{V}_l$ -rooted random spanning forest sampled by RsflIndex;
- 3  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} \leftarrow 0$  for  $u \in \mathcal{U}$ ;
- 4 Fix a path  $\mathcal{P}_{u \rightsquigarrow \mathcal{V}_l}$  from node  $u$  to  $\mathcal{V}_l$  for each  $u \in \mathcal{U}$  in  $\mathcal{G}$ ;
- 5 **for** each  $F \in \mathcal{F}_{\mathcal{V}_l}$  **do**
- 6     **for** each node  $u \in \mathcal{U}$  **do**
- 7         Let  $\mathcal{P}_{F:u \rightsquigarrow \mathcal{V}_l}$  denote the unique path from  $u$  to  $\mathcal{V}_l$  in  $F$ ;
- 8         **for** each edge  $(e_1, e_2) \in \mathcal{P}_{F:u \rightsquigarrow \mathcal{V}_l}$  **do**
- 9             **if**  $(e_1, e_2) \in \mathcal{P}_{F:u \rightsquigarrow \mathcal{V}_l}$  **then**
- 10                  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} \leftarrow (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} + \frac{1}{\omega}$
- 11             **if**  $(e_2, e_1) \in \mathcal{P}_{F:u \rightsquigarrow \mathcal{V}_l}$  **then**
- 12                  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} \leftarrow (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} - \frac{1}{\omega}$
- 13 **return**  $\tilde{\mathbf{P}}_f, \mathbf{L}_{\mathcal{H}}^\dagger, (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  for each  $u \in \mathcal{U}$

---

Recall that for a given source node  $s$ , we can simultaneously compute  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  for all  $t \in \mathcal{V}$  by invoking the  $\mathcal{V}_l$ -absorbed random walk or  $\mathcal{V}_l$ -absorbed push algorithms proposed in Section 4.2. Thus, for the single-source ER computation problem, the most challenging step is to compute  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{tt}$  for all  $t \in \mathcal{V}$ , i.e., the diagonal elements of the matrix  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$ . Since the diagonal elements of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$  is independent on the source node, a natural question is that can we pre-compute all such diagonal elements and maintain them as an index? A basic method to compute all diagonal elements requires to invoke the  $\mathcal{V}_l$ -absorbed random walk or  $\mathcal{V}_l$ -absorbed push algorithms for  $n$  times, which is costly for large graphs. In this section, we first propose two efficient Monte Carlo algorithms to estimate  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{tt}$  for all  $t \in \mathcal{V}$ . Together with the two index matrices  $\tilde{\mathbf{P}}_f$  and  $\mathbf{L}_{\mathcal{H}}^\dagger$  constructed in Section 4.1, we obtain a new index, including all diagonal elements of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$ , for handling single-source ER query. Based on the new index, we then propose two efficient query processing algorithms for handling single-source ER queries.

## 5.1 Index Construction Algorithms

To construct our index, the key is to efficiently estimate all diagonal elements of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$ . To achieve this goal, we propose two novel algorithms based on the Wilson algorithm with a root set  $\mathcal{V}_l$ . We show that all diagonal elements of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$  can be approximated at the same time when building the index matrices  $\tilde{\mathbf{P}}_f$  and  $\mathbf{L}_{\mathcal{H}}^\dagger$  using Algorithm 2.

### Index building by $\mathcal{V}_l$ -rooted random spanning forest sampling.

Recall that Algorithm 2 can obtain two index matrices  $\tilde{\mathbf{P}}_f$  and  $\mathbf{L}_{\mathcal{H}}^\dagger$  by sampling a set of  $\mathcal{V}_l$ -rooted random spanning forests. Let  $\mathcal{F}_{\mathcal{V}_l}$  be the set of  $\mathcal{V}_l$ -rooted random spanning forests sampled by Algorithm 2. Interestingly, we find that all the diagonal elements of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$  can also be estimated by  $\mathcal{F}_{\mathcal{V}_l}$ . Specifically, we first establish a useful lemma which shows a close connection between the number of  $\mathcal{V}_l$ -rooted random spanning forests to the diagonal elements of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$ .

Let  $\mathcal{F}_{\mathcal{V}_l}$  be the set of  $\mathcal{V}_l$ -rooted random spanning forests. Denote by  $\mathcal{P}_{u \rightsquigarrow \mathcal{V}_l}$  the path from  $u$  to the node set  $\mathcal{V}_l$ . Here, a path from a node  $u$  to a set  $\mathcal{V}_l$  means that the path comes from  $u$  and terminates

at any one node in  $\mathcal{V}_l$ . Further, we let  $\mathcal{F}_{\mathcal{V}_l:u \rightsquigarrow \mathcal{V}_l}^{(e_1, e_2)}$  be the set of  $\mathcal{V}_l$ -rooted random spanning forests that consist of all  $F \in \mathcal{F}_{\mathcal{V}_l}$  and the unique path from node  $u$  to  $\mathcal{V}_l$  passes through the edge  $(e_1, e_2)$ . Then, we have the following results.

LEMMA 5.1.

$$(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} = \sum_{(e_1, e_2) \in \mathcal{P}_{u \rightsquigarrow \mathcal{V}_l}} \frac{|\mathcal{F}_{\mathcal{V}_l:u \rightsquigarrow \mathcal{V}_l}^{(e_1, e_2)}| - |\mathcal{F}_{\mathcal{V}_l:u \rightsquigarrow \mathcal{V}_l}^{(e_2, e_1)}|}{|\mathcal{F}_{\mathcal{V}_l}|},$$

where the sum is taking over edges in the path  $\mathcal{P}_{u \rightsquigarrow \mathcal{V}_l}$ .

PROOF. According to the all-minors matrix forest theorem [10], we have:  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{ue_1} = \frac{|\mathcal{F}_{\mathcal{V}_l|u, e_1}|}{|\mathcal{F}_{\mathcal{V}_l}|}$ , where  $\mathcal{F}_{\mathcal{V}_l|u, e_1}$  denotes the set of spanning forests including all  $F$  that has  $|\mathcal{V}_l| + 1$  connected components, with  $|\mathcal{V}_l|$  connected components having a node  $v_i \in \mathcal{V}_l$  as the root, and one connected component which includes  $u$  and  $e_1$ , with  $e_1$  being the root. For each edge  $(e_1, e_2) \in \mathcal{E}_{\mathcal{U}}$ ,  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{ue_1} - (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{ue_2} = \frac{|\mathcal{F}_{\mathcal{V}_l|u, e_1}| - |\mathcal{F}_{\mathcal{V}_l|u, e_2}|}{|\mathcal{F}_{\mathcal{V}_l}|}$ . Let  $\mathcal{F}_{u, e_1|e_2, \mathcal{V}_l}$  denote the set of spanning forests with  $|\mathcal{V}_l| + 1$  connected components, with  $|\mathcal{V}_l|$  connected components having a node  $v_i \in \mathcal{V}_l$  as the root, and one connected component which includes  $e_1$  and  $u$ , with  $e_1$  being the root. Specifically,  $e_1$  and  $e_2$  are in different connected components. Depending on whether  $e_1, e_2$  are in the connected component of  $u$ , we can partition the spanning forests, which results in  $|\mathcal{F}_{\mathcal{V}_l|u, e_1}| - |\mathcal{F}_{\mathcal{V}_l|u, e_2}| = |\mathcal{F}_{u, e_1|e_2, \mathcal{V}_l}| - |\mathcal{F}_{u, e_2|e_1, \mathcal{V}_l}|$ . For each edge  $(e_1, e_2) \in \mathcal{E}_{\mathcal{U}}$ , there is a bijection between  $\mathcal{F}_{\mathcal{V}_l:u \rightsquigarrow \mathcal{V}_l}^{(e_1, e_2)}$  and  $\mathcal{F}_{u, e_1|e_2, \mathcal{V}_l}$ , since if we remove the edge  $(e_1, e_2)$  from a spanning forest in  $\mathcal{F}_{\mathcal{V}_l:u \rightsquigarrow \mathcal{V}_l}^{(e_1, e_2)}$ , we will exactly obtain a spanning forest in  $\mathcal{F}_{u, e_1|e_2, \mathcal{V}_l}$ . If we delete the edge  $(e_1, e_2)$  from a spanning forest in  $\mathcal{F}_{u, e_1|e_2, \mathcal{V}_l}$ , we will exactly obtain a spanning forest in  $\mathcal{F}_{\mathcal{V}_l:u \rightsquigarrow \mathcal{V}_l}^{(e_1, e_2)}$ . Furthermore, we have  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} = \sum_{(e_1, e_2) \in \mathcal{P}_{u \rightsquigarrow \mathcal{V}_l}} (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{ue_1} - (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{ue_2}$ . Thus, the Lemma is established.  $\square$

Based on Lemma 5.1, we can fix a set of paths from each node  $u \in \mathcal{U}$  to  $\mathcal{V}_l$ . When sampling a spanning forest  $F \in \mathcal{F}_{\mathcal{V}_l}$ , we examine whether the paths from  $u$  to  $\mathcal{V}_l$  in  $F$  pass through the edges in the pre-fixed paths to obtain an unbiased estimation of  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  for  $u \in \mathcal{U}$ . Note that this can be easily implemented by a depth-first search (DFS) procedure. Specifically, we start the DFS procedure by regarding the landmark node set  $\mathcal{V}_l$  as a start node, and exploit the graph in a DFS manner. During the DFS procedure, we record the visit time and the finish time for each node in  $\mathcal{U}$  that we visit [12]. After the DFS procedure terminates, we can examine whether the edge  $(e_1, e_2)$  passes the path  $\mathcal{P}_{F:u \rightsquigarrow \mathcal{V}_l}$  by comparing the visit time and finish time of  $e_1$  and  $u$ . The benefit of the implementation is that we can finish the DFS procedure in  $O(m + n)$  time and examine each node in  $O(1)$  time [12], which is very efficient. Combining with RsflIndex to approximate  $\mathbf{P}_f$  and  $(\mathbf{L}/\mathcal{V}_l)^\dagger$ , the resulting index building algorithm is shown in Algorithm 5. In Lines 5-12, the algorithm estimates the diagonal elements by examining a set of fixed paths in a sampled spanning forest. When the algorithm terminates, all  $\tilde{\mathbf{P}}_f, \mathbf{L}_{\mathcal{H}}^\dagger, (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  for  $u \in \mathcal{U}$  are stored as the index. The size of the index is still  $O(|\mathcal{V}_l| \cdot n)$ . The following

**Algorithm 6:** LeWalkIndex-SS**Input:** Graph  $\mathcal{G}$ , landmark node set  $\mathcal{V}_l$ , sample size  $\omega$ **Output:**  $\tilde{\mathbf{P}}_f, \mathbf{L}_{\mathcal{H}}^\dagger, (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  for each  $u \in \mathcal{U}$ 


---

```

1  $\tilde{\mathbf{P}}_f, \mathbf{L}_{\mathcal{H}}^\dagger \leftarrow \text{RsflIndex}(\mathcal{G}, \mathcal{V}_l, \omega);$ 
2  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} \leftarrow 0$  for  $u \in \mathcal{U};$ 
3 for  $i = 1 : \omega$  do
4   In the loop-erased walk of the Wilson algorithm in the  $i$ -th iteration of RsflIndex, if a
     node is visited for  $n_u$  times, then  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} \leftarrow (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} + \frac{n_u}{\omega};$ 
5 return  $\tilde{\mathbf{P}}_f, \mathbf{L}_{\mathcal{H}}^\dagger, (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  for  $u \in \mathcal{U};$ 

```

---

theorem shows the accuracy and time complexity of Algorithm 5 to build our index.

LEMMA 5.2. *Algorithm 5 can output an estimation of  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  with absolute error  $\epsilon$  for all  $u \in \mathcal{U}$  with the time complexity  $O(\frac{\text{Tr}((\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1}) d_{\mathcal{V}_l}^{\max} \Delta_{\mathcal{G}}^2 |\mathcal{V}_l| \log n}{\epsilon^2})$ .*

PROOF. Let  $X$  be the random variable added to  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  in Algorithm 5, that is, for each edge in  $\mathcal{P}_{u \rightsquigarrow \mathcal{V}_l}$ , if the edge appears in the path from  $u$  to  $\mathcal{V}_l$  in the sampled spanning forest  $F$ , then add  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  by 1. If it appears in the opposite direction, then add  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  by  $-1$ . For all such edges, since the path has at most  $\Delta_{\mathcal{G}}$  edges, we have  $|X| \leq \Delta_{\mathcal{G}}$ . Thus, according to Hoeffding bound, when  $\omega \geq \frac{4\Delta_{\mathcal{G}}^2 \log n}{\epsilon}$ , the algorithm can achieve an  $\epsilon$  absolute error for  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  with high probability. Sampling the spanning forest costs time  $\text{Tr}((\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1})$ . The traverse operation in Line 5-12 costs time  $|\mathcal{U}| \Delta_{\mathcal{G}}$ . Thus, the total time complexity is dominated by Algorithm 2, which is  $O(\frac{\Delta_{\mathcal{G}}^2 \log n \text{Tr}((\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1})}{\epsilon^2})$ . The proof is completed.  $\square$

**Index building by loop-erased random walk sampling.** Note that Algorithm 5 requires matching the paths in the random spanning forests with a pre-fixed path, which introduces additional time overheads. Here we propose an alternative algorithm based on the loop-erased random walk (LERW) sampling. Below, we first present a novel loop-erased random walk interpretation of  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  for  $u \in \mathcal{U}$ , and based on this interpretation we can develop our algorithm.

LEMMA 5.3. *Let  $\tau_{\mathcal{V}_l}^{\text{LE}}[u, u]$  be the expected number of passes to a node  $u \in \mathcal{U}$  in a loop-erased random walk with root set  $\mathcal{V}_l$ . Then, we have  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} = \frac{\tau_{\mathcal{V}_l}^{\text{LE}}[u, u]}{d_u}$ .*

PROOF. In a loop-erased random walk with root set  $\mathcal{V}_l$ , the output distribution is independent on the node ordering according to the analysis in [45]. When setting the node  $u$  as the first node, the probability distribution of the  $\mathcal{V}_l$ -absorbed random walk starts from  $u$  and the loop-erased random walk with root set  $\mathcal{V}_l$  are just the same. Note that when the first node  $u$  hits a node in  $\mathcal{V}_l$ , the loop-erased walk will never pass through  $u$  again. Thus, by Lemma 3.4, we have  $\tau_{\mathcal{V}_l}^{\text{LE}}[u, u] = \tau_{\mathcal{V}_l}[u, u] = (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} d_u$ . As a result, the Lemma is established.  $\square$

The detailed description of this index-building algorithm is shown in Algorithm 6. According to Lemma 5.3, we can utilize the loop-erased random walks with root set  $\mathcal{V}_l$  for sampling random spanning forests in Algorithm 2 to simultaneously estimate the diagonal elements  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  for all  $u \in \mathcal{U}$ . Specifically, by Lemma 5.3, we use the degree-normalized average number of visits to each node  $u \in \mathcal{U}$  as an unbiased estimation of  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  for all  $u \in \mathcal{U}$  (Line 4). Below, we analyze the accuracy and time complexity of Algorithm 6.

LEMMA 5.4. *Algorithm 6 can output an estimation of  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  for all  $u \in \mathcal{U}$  in time  $O(\frac{(\sigma_{\mathcal{V}_l}^2 + d_{\mathcal{V}_l}^{\max}) \log n \text{Tr}((\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1})}{\epsilon^2})$ .*

PROOF. Let  $X$  be the random variable which represents the number of visits to a node  $u$  in a  $\mathcal{V}_l$ -absorbed loop-erased walk. The distribution of the walk is identical to a  $v$ -absorbed walk. Thus, the distribution of the number of visits to a node  $u$  in a loop-erased walk is equal to the distribution of the number of passes to the node  $u$  in a  $\mathcal{V}_l$ -absorbed walk starts from  $u$ . According to the proof of Lemma 4.9, when the sample size satisfies  $\omega \geq \frac{\sigma_{\mathcal{V}_l}^2 \log n}{\epsilon^2}$ , the algorithm can achieve an  $\epsilon$  absolute error for  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$ . Here, the random walk steps can be computed via utilizing LERW to sample spanning forests, with little additional overhead. Combined with Lemma 4.6, the overall time complexity is  $O(\frac{(\sigma_{\mathcal{V}_l}^2 + d_{\mathcal{V}_l}^{\max}) \log n \text{Tr}((\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1})}{\epsilon^2})$ , which establishes the lemma.  $\square$

**Discussions.** The sample size of Algorithm 5 and Algorithm 6 are  $\frac{d_{\mathcal{V}_l}^{\max} \Delta_{\mathcal{G}}^2 |\mathcal{V}_l| \log n}{\epsilon^2}$  and  $\frac{(\sigma_{\mathcal{V}_l}^2 + d_{\mathcal{V}_l}^{\max}) \log n}{\epsilon^2}$  respectively. As discussed before,  $\Delta_{\mathcal{G}}, \sigma_{\mathcal{V}_l}$  are typically small numbers. We also find in experiments that our algorithms can achieve low errors with very small sample size. Note that in [20] there are also two index building algorithms based on sampling spanning trees and loop-erased walks. The time complexity of our algorithm for drawing a sample  $(\text{Tr}((\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1}))$  is strictly smaller than [20]  $(\text{Tr}((\mathbf{I} - \mathbf{P}_o)^{-1}))$ , since there are more nodes to hit. For example, on PowerGrid [18],  $\text{Tr}((\mathbf{I} - \mathbf{P}_o)^{-1}) = 4.3 \times 10^5$ , while  $\text{Tr}((\mathbf{I} - \mathbf{P}_{\mathcal{U}\mathcal{U}})^{-1}) = 1.5 \times 10^5$  when  $\mathcal{V}_l$  is the top-100 highest degree nodes. We will also show significantly empirical improvement of our index building algorithms compared to [20].

## 5.2 Query processing algorithms

Armed with our index, we propose two efficient query processing algorithms based on  $\mathcal{V}_l$ -absorbed walk sampling and  $\mathcal{V}_l$ -absorbed push respectively. Algorithm 7 shows the detailed framework of our algorithm. If  $s \in \mathcal{V}_l$  (Lines 1-5), there are two cases to compute  $r(s, u)$  by Lemma 3.7. For  $u \in \mathcal{U}$ ,  $r(s, u) = (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} + (\mathbf{p}_s - \mathbf{e}_u)^T \mathcal{L}_{\mathcal{H}}^\dagger (\mathbf{p}_s - \mathbf{e}_u)$  (Lines 2-3). Since all the diagonal element  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$ , the probability vector  $\mathbf{p}_s$  and the pseudo-inverse  $\mathcal{L}_{\mathcal{H}}^\dagger$  can be directly obtained from the index, the computation is very fast. For  $u \in \mathcal{V}_l$ , two nodes are both located in the landmark node set,  $r(s, u) = (\mathbf{e}_s - \mathbf{e}_u)^T \mathcal{L}_{\mathcal{H}}^\dagger (\mathbf{e}_s - \mathbf{e}_u)$  (Lines 4-5). Again, all the elements can be directly obtained from the index. When  $s \in \mathcal{U}$ , the situation is more complicated (Lines 6-11). According to Lemma 3.3, there are also two cases to compute  $r(s, u)$ . For  $u \in \mathcal{U}$ , both of the two nodes are located in the node set  $\mathcal{U}$ , we have

**Algorithm 7:** The single-source ER query processing algorithm

---

**Input:** Graph  $\mathcal{G}$ , landmark set  $\mathcal{V}_l$ , a source node  $s$ ,  $\bar{\mathbf{p}}_f$ ,  $\mathcal{L}_{\mathcal{H}}^\dagger$ ,  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  for all  $u \in \mathcal{U}$   
**Output:**  $\tilde{r}(s, u)$  for all  $u \in \mathcal{V}$

---

```

1 if  $s \in \mathcal{V}_l$  then
2   for each node  $u \in \mathcal{U}$  do
3      $\tilde{r}(s, u) \leftarrow (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} + (\mathbf{p}_s - \mathbf{e}_u)^T \mathcal{L}_{\mathcal{H}}^\dagger (\mathbf{p}_s - \mathbf{e}_u)$ ;
4   for each node  $u \in \mathcal{V}_l$  do
5      $\tilde{r}(s, u) \leftarrow (\mathbf{e}_s - \mathbf{e}_u)^T \mathcal{L}_{\mathcal{H}}^\dagger (\mathbf{e}_s - \mathbf{e}_u)$ ;
6 else if  $s \in \mathcal{U}$  then
7   Invoke the  $\mathcal{V}_l$ -absorbed random walk sampling or the  $\mathcal{V}_l$ -absorbed push algorithms
   to obtain an estimation  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{st}$  for all  $t \in \mathcal{V}$ ;
8   for each node  $u \in \mathcal{U}$  do
9      $\tilde{r}(s, u) \leftarrow$ 
       $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{ss} + (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} - 2(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{su} + (\mathbf{p}_s - \mathbf{p}_u)^T \mathcal{L}_{\mathcal{H}}^\dagger (\mathbf{p}_s - \mathbf{p}_u)$ ;
10  for each node  $u \in \mathcal{V}_l$  do
11     $\tilde{r}(s, u) \leftarrow (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} + (\mathbf{p}_s - \mathbf{e}_u)^T \mathcal{L}_{\mathcal{H}}^\dagger (\mathbf{p}_s - \mathbf{e}_u)$ ;
12 return  $\tilde{r}(s, u)$  for  $u \in \mathcal{V}$ ;
```

---

$r(s, u) = (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{ss} + (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} - 2(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{su} + (\mathbf{p}_s - \mathbf{e}_u)^T \mathcal{L}_{\mathcal{H}}^\dagger (\mathbf{p}_s - \mathbf{e}_u)$ . Except the elements  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{su}$  for all  $u \in \mathcal{U}$ , all other elements can be directly obtained from the index. Thus, we utilize  $\mathcal{V}_l$ -absorbed random walk sampling or  $\mathcal{V}_l$ -absorbed push to estimate a column of the matrix  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$  (Line 7). For convenience, we refer to these algorithms as RW-SS and Push-SS, respectively. For  $u \in \mathcal{V}_l$ ,  $r(s, u) = (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} + (\mathbf{p}_s - \mathbf{e}_u)^T \mathcal{L}_{\mathcal{H}}^\dagger (\mathbf{p}_s - \mathbf{e}_u)$  (Line 10-11). It can also be directly computed based on the index. Combining all these procedures together, it shows that with an  $O(n \times |\mathcal{V}_l|)$  size index, we can answer a single-source ER query in time of answering a single-pair ER query. Below, we present a detailed theoretical analysis of our algorithms.

**Theoretical analysis of the algorithms.** Based on the above results, we can give the query time complexity for RW-SS and Push-SS to output an estimation of  $r(s, u)$  with an  $\epsilon$ -absolute error for all  $u \in \mathcal{U}$ , equipped with the index we described before.

**LEMMA 5.5.** *The query time complexity of RW-SS is  $O(\frac{\sigma_{\mathcal{V}_l}^3}{\epsilon^2})$  for an  $\epsilon$ -estimation of  $r(s, u)$  for all  $u \in \mathcal{V}$ .*

**PROOF.** Based on Lemma 4.13, with the time complexity give therein, the error of index building is smaller than  $\epsilon$ . According to Lemma 5.2 and Lemma 5.4, the diagonal elements are also estimated within an  $\epsilon$ -absolute error. The results of Lemma 4.8 show that we can estimate a row of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$  simultaneously in time  $O(\frac{\sigma_{\mathcal{V}_l}^3}{\epsilon^2})$ . Thus, the proof is established.  $\square$

**LEMMA 5.6.** *Let  $h_{\mathcal{V}_l}^{\max}(t)$  be the maximum value of hitting time from all each nodes  $u \in \mathcal{U}$  to  $\mathcal{V}_l$ , the query time complexity of Push-SS is  $O(\frac{dh(s, \mathcal{V}_l)h_{\mathcal{V}_l}^{\max}(t)}{\epsilon})$  for an  $\epsilon$ -estimation of  $r(s, u)$  for all  $u \in \mathcal{V}$ .*

**PROOF.** Based on Lemma 4.13, with the time complexity give therein, the error of index building is smaller than  $\epsilon$ . According to Lemma 5.2 and Lemma 5.4, the diagonal elements are also estimated within an  $\epsilon$ -absolute error. Let  $h_{\mathcal{V}_l}^{\max}(t)$  denote the maximum value of hitting time from all each nodes  $u \in \mathcal{U}$  to  $\mathcal{V}_l$ , the results of

**Table 1: Datasets ( $\bar{d}$ : average degree;  $\Delta_G$ : diameter of the graph)**

Dataset	$n$	$m$	$\bar{d}$	$\Delta_G$
PowerGrid	4,941	6,594	2.67	46
Email-enron	33,696	180,811	10.73	11
Road-PA	1,087,562	1,541,514	2.83	786
Road-CA	1,957,027	2,760,388	2.82	849
Orkut	3,072,441	117,185,083	76.28	9

Lemma 4.9 show that we can estimate a row of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$  simultaneously in time  $O(\frac{dh(s, \mathcal{V}_l)h_{\mathcal{V}_l}^{\max}(t)}{\epsilon^2})$ . Thus, the proof is established.  $\square$

Based on the discussions before, the query time complexity of our algorithms is strictly better compared to the single-source query algorithms in [20]. Specifically, the query complexity RW-SS algorithm in [20] is  $O(\frac{\sigma_{\mathcal{V}_l}^3}{\epsilon^2})$  based on our analysis, and the Push-SS algorithm is  $O(\frac{dh(s, v)h_{\mathcal{V}_l}^{\max}(t)}{\epsilon})$ . We improve the  $\sigma_{\mathcal{V}_l}$  term into  $\sigma_{\mathcal{V}_l}$  and the  $h(s, v)$ ,  $h_{\mathcal{V}_l}^{\max}(t)$  terms into  $h(s, \mathcal{V}_l)$ ,  $h_{\mathcal{V}_l}^{\max}(t)$ . Similar to the single-pair algorithms, RW-SS and Push-SS perform much better than such worst-case bounds in real-life graphs.

## 6 EXPERIMENTS

### 6.1 Experimental Settings

**Datasets.** We employ five real-world datasets encompassing various types of graphs, mainly including social networks and road networks. The detailed statistics of the datasets is illustrated in Table 1. Among the five datasets, Email-enron and Orkut are social networks. Road-PA and Road-CA are road networks of America. For PowerGrid, it is an infrastructure network, which we find its statistics similar to the road networks. Road networks and social networks have very different structures. It is common to observe a comparatively smaller average degree combined with a larger diameter on road networks. The algorithm behavior is also different on the two types of networks. For single-pair query, we randomly generate 100 node pairs as the query set and report the average results. For single-source query, we randomly generate 50 source nodes as the query set and report the average performance. We can obtain ground truth for both single-pair and single-source effective resistance by applying deterministic eigen-decomposition to the small datasets. On large graphs, we obtain the ground-truth results for single-pair query by applying the state-of-the-art deterministic algorithm Push-v [20] with  $r_{\max} = 10^{-9}$  following [20, 47]. We apply LEwalk in [20] with  $\omega = 10^6$  to obtain the ground-truth results for single-source query following [20] on large datasets.

**Different algorithms.** For the index building algorithm, we implement two index building algorithms RwIndex (Algorithm 1) and RsflIndex (Algorithm 2). To support single-source query, we also implement two extended algorithms based on RsflIndex named RsflIndex-SS (Algorithm 5) and LeWalkIndex-SS (Algorithm 6). We compare the two single-source index building algorithms with the state-of-the-art index building algorithm UST and LEwalk proposed in [20] which are special cases of our index building algorithms when  $|\mathcal{V}_l| = 1$ . For single-pair query processing algorithms, we implement three algorithms RW, Push and Bipush which is described in Algorithm 4. We compare the three query processing algorithms with four state-of-the-art algorithms, GEER [47] and the single landmark

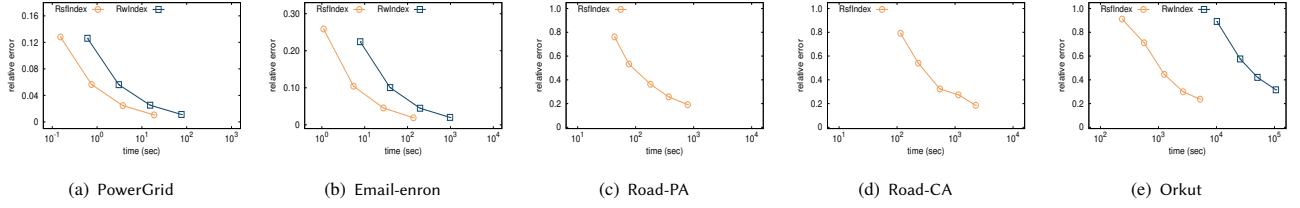


Figure 1: Performance of different index building algorithms for single-pair ER query

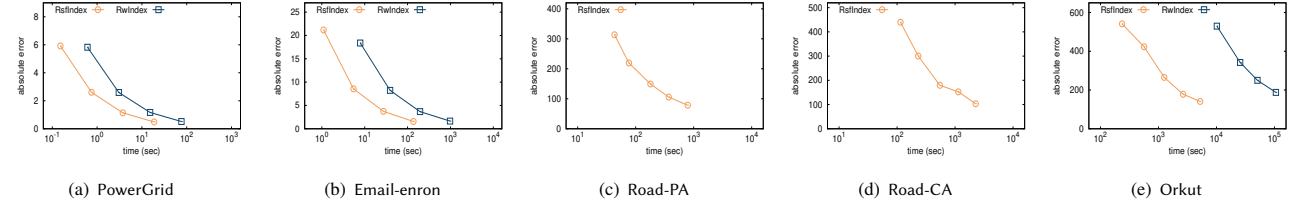


Figure 2: Performance of different index building algorithms for single-pair ER query (absolute error)

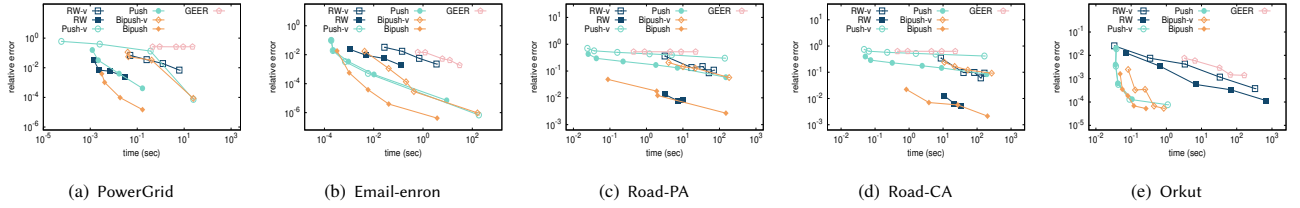


Figure 3: Performance of different query processing algorithms for single-pair ER query

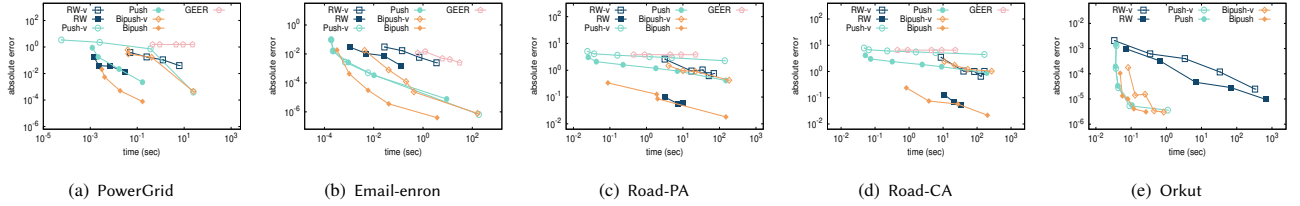


Figure 4: Performance of different query processing algorithms for single-pair ER query (absolute error)

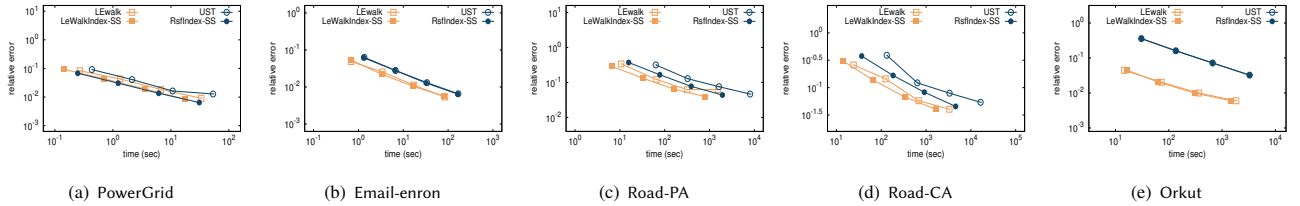
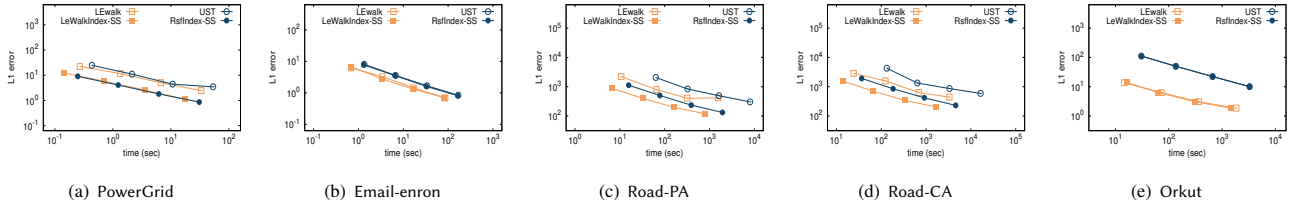


Figure 5: Performance of different index building algorithms for single-source ER query

Figure 6: Performance of different index building algorithms for single-source ER query ( $L_1$ -error)

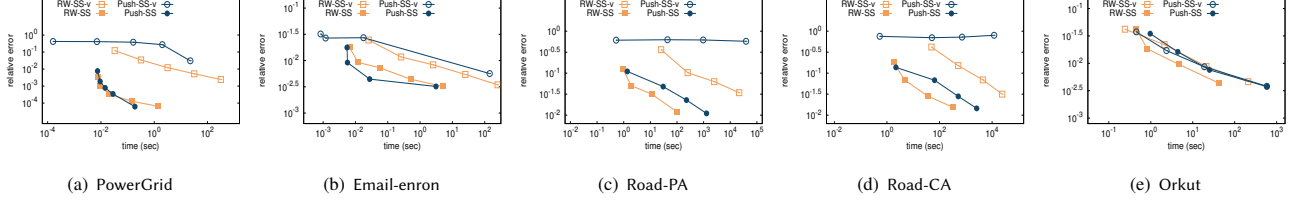


Figure 7: Performance of different query processing algorithms for single-source ER query

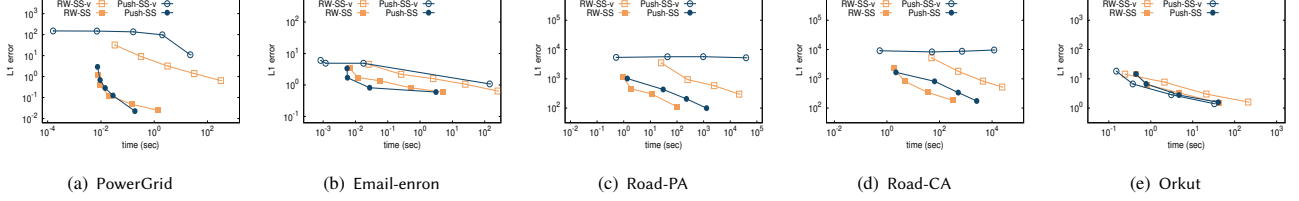


Figure 8: Performance of different query processing algorithms for single-source ER query (L1-error)

node special cases RW-v, Push-v and Bipush-v proposed in [20]. For single-source query, we implement two query processing algorithms RW-SS and Push-SS which are described in Algorithm 7. We compare them with the single landmark node special cases RW-SS-v and Push-SS-v proposed in [20].

**Different algorithms.** For the index building algorithm, we implement two algorithms Rwlndex (Algorithm 1) and Rsflndex (Algorithm 2). To support single-source query, we also implement two extended algorithms based on Rsflndex named Rsflndex-SS (Algorithm 5) and LeWalkIndex-SS (Algorithm 6). We compare the two single-source index building algorithms with the state-of-the-art index building algorithm UST and LEwalk proposed in [20] which are special cases of our index building algorithms when  $|\mathcal{V}_l| = 1$ . For single-pair query processing algorithms, we implement three algorithms RW, Push and Bipush which is described in Algorithm 4. We compare the three query processing algorithms with four state-of-the-art algorithms, GEER [47] and RW-v, Push-v and Bipush-v proposed in [20], which are the special cases of our algorithms with only one landmark node. For single-source query, we implement two query processing algorithms RW-SS and Push-SS which are described in Algorithm 7. We compare them with the state-of-the-art algorithms RW-SS-v and Push-SS-v proposed in [20], which are the special cases of our algorithms with one landmark node.

**Parameters.** For all the proposed algorithms including the index building algorithms and query processing algorithms for both single-pair query and single-source query, we have shown the sample size needed and the total time complexity for the algorithms to achieve an  $\epsilon$ -estimation of ER. However, the bounds are just worst-case bounds for comparing the algorithms theoretically. There will be too much unnecessary work if we use these bounds to determine sample sizes. Thus, following [20, 29, 47], we compare the algorithms by carefully setting the sample size  $\omega$  and the error threshold  $r_{max}$  to make the algorithms run in a reasonable time. Specifically, for the single-pair index building algorithms, we set  $\omega = 10^4$  by default for both Rwlndex and Rsflndex. For the single-pair query processing algorithms, we set  $\omega = 10^4$  for RW,  $r_{max} = 10^{-4}$  for Push, and  $\omega = 10^3$ ,  $r_{max} = 10^{-3}$  for Bipush by default. We also vary the

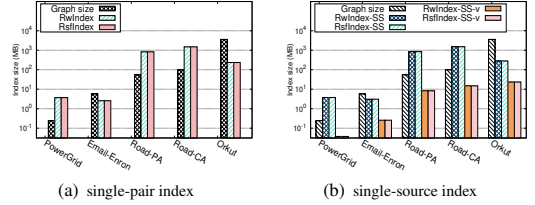


Figure 9: Index space consumption of different algorithms

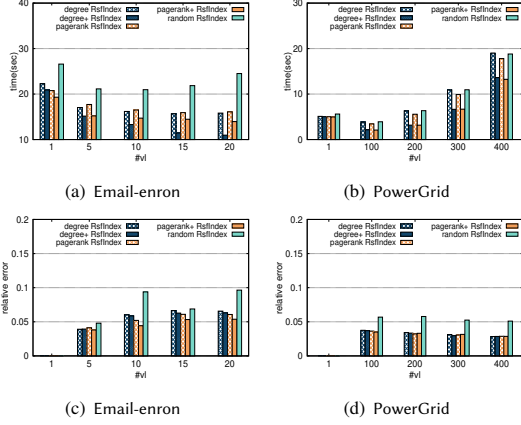
sample size  $\omega$  and the error threshold  $r_{max}$  to compare the empirical errors. For the compared algorithms, there is one parameter  $\epsilon$  which controls the accuracy for GEER, we set  $\epsilon$  as 0.01 by default following [47]. For the algorithms proposed in [20], we set  $\omega = 10^4$  for RW-v,  $r_{max} = 10^{-4}$  for Push-v, and  $\omega = 10^{-3}$ ,  $r_{max} = 10^{-3}$  for Bipush-v following [20]. For the single-source index building algorithms, we set  $\omega = 10^4$  by default for Rsflndex-SS and LeWalkIndex-SS. We also set  $\omega = 10^4$  for UST and LEwalk following [20]. For single-source query processing algorithms, we set  $\omega = 10^4$  for RW-SS and  $r_{max} = 10^{-4}$  for Push-SS by default. We also set  $\omega = 10^4$  for RW-SS-v and  $r_{max} = 10^{-4}$  for Push-SS-v following [20]. For the landmark node set  $\mathcal{V}_l$ , we set  $|\mathcal{V}_l| = 10$  in social networks and  $|\mathcal{V}_l| = 100$  in road networks and use degree+ to select the landmark node set in all our algorithms. We also vary  $|\mathcal{V}_l|$  and different landmark selection methods to compare the performance in Section 6.4. For the query processing algorithms, we use the index built by Rsflndex for single-pair query and the index built by LeWalkIndex-SS for single-source query by default. We also evaluate the performance of the query processing algorithms with the index built by Rwlndex and Rsflndex-SS in Section 6.4.

## 6.2 Results of Single-pair ER Computation

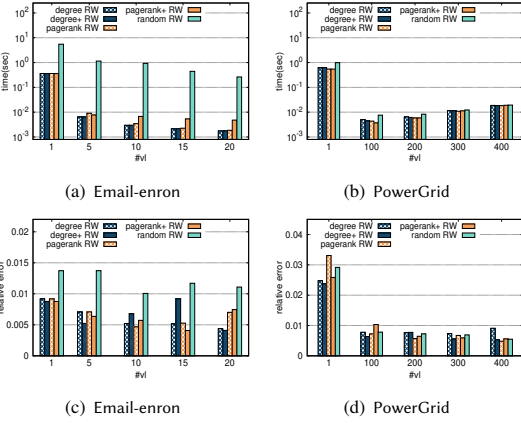
We first study the performance of the algorithms for single-pair effective resistance computation.

**Evaluation of index building algorithms.** We evaluate the index building time and the index size for different single-pair index building algorithms. We use the relative error of the norm of the related

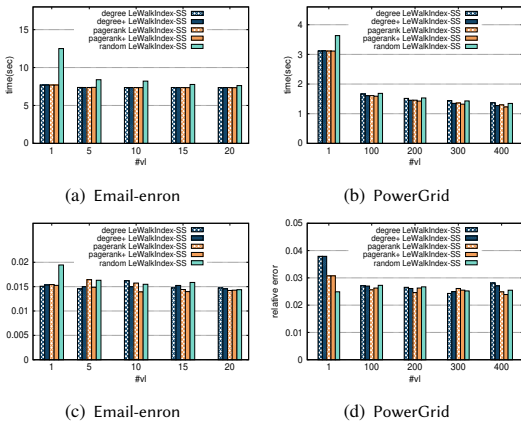




**Figure 10: Effect of the choice of  $V_l$  for single-pair index building algorithms**

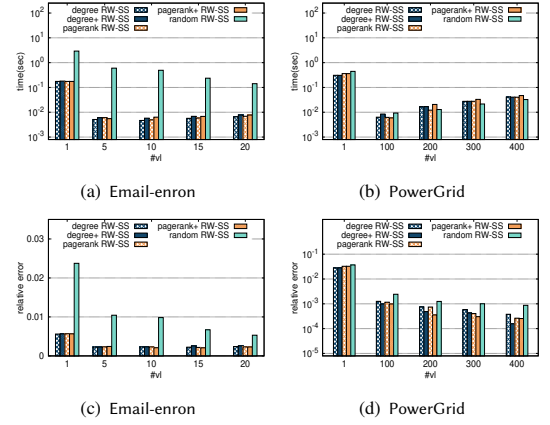


**Figure 11: Effect of the choice of  $V_l$  for single-pair query processing algorithms**



**Figure 12: Effect of the choice of  $V_l$  for single-source index building algorithms**

matrices to evaluate the index building quality. Suppose that  $\tilde{\mathbf{P}}$  and  $\mathbf{L}_{\mathcal{H}}^{\dagger}$  is the estimated matrices output by RwlIndex and RsflIndex, the



**Figure 13: Effect of the choice of  $V_l$  for single-source query processing algorithms**

relative error is defined as  $\frac{\|\tilde{\mathbf{P}} - \mathbf{L}_{\mathcal{H}}^{-1} \mathbf{L}_{\mathcal{U}} \mathbf{V}_l\|_F}{\|\mathbf{L}_{\mathcal{H}}^{-1} \mathbf{L}_{\mathcal{U}} \mathbf{V}_l\|_F} + \frac{\|\mathbf{L}_{\mathcal{H}}^{\dagger} - (\mathbf{L}/\mathbf{V}_l)^{\dagger}\|_F}{\|\mathbf{L}_{\mathcal{H}}^{\dagger}\|_F}$ , where  $\|\cdot\|_F$  is the Frobenius norm of a matrix. We also evaluate the error by the absolute error of the Frobenius norm. We vary the parameters of RwlIndex and RsflIndex and plot the running time and empirical error trade-offs, the results are illustrated in Fig. 1 and Fig. 2. On two large road networks Road-PA and Road-CA, RwlIndex cannot output a result within 24 hours. Thus, we omit the results of RwlIndex. As can be seen, both RwlIndex and RsflIndex can approximately build the index with a small relative error. Specifically, RsflIndex is much faster than RwlIndex. For example, on Orkut, it takes  $2.7 \times 10^3$  secs for RsflIndex to build the index with a relative error 0.3, while it takes  $1.06 \times 10^5$  secs for RwlIndex, which is more than 2 orders of magnitude slower. The reason is that RsflIndex uses loop-erased walks which are more efficient than sampling  $|\mathcal{U}|$   $V_l$ -absorbed random walks. The index size of different index building algorithms is illustrated in Fig. 9(a). As can be seen, the index size of RwlIndex and RsflIndex is smaller than the graph size on social networks, and the index size of RwlIndex and RsflIndex is larger than the graph size on road graphs. Recall that the space complexity of the index is  $O(|V_l| \times n)$ . However, even on the large road network Road-CA, it consumes only 1493 MB space to store the index when  $|V_l| = 100$ . Thus, the index size is acceptable and can be easily stored in memory.

**Evaluation of query processing algorithms.** We compare the proposed single-pair query processing algorithms RW, Push, Bipush with four SOTA methods, the RW-v, Push-v, Bipush-v algorithm proposed in [20], which is the special cases of our algorithms, as well as GEER proposed in [47]. For all the proposed algorithms, we use the index that we build with the algorithm RsflIndex in a similar degree of accuracy. We vary the parameters of the algorithms and plot the empirical error and query time trade-offs. Specifically, let  $\hat{r}(s, t)$  be the estimated ER value, we evaluate the errors by the relative error defined as  $\frac{|\hat{r}(s, t) - r(s, t)|}{r(s, t)}$ . We also evaluate the errors by the absolute error defined as  $|\hat{r}(s, t) - r(s, t)|$ . The results are illustrated in Fig. 3 and Fig. 4. We can find that on all five datasets, the query processing algorithms RW, Push and Bipush all improve significantly over their single landmark node special cases RW-v, Push-v and Bipush-v. They can also achieve lower errors compared

to the SOTA method GEER. This demonstrates the effectiveness of the multiple landmark nodes approaches. Among these algorithms, Bipush obtains the best performance, especially on road networks. For example, on Road-PA, Bipush can achieve a relative error 0.05 in 0.09 secs, while the fastest SOTA method Bipush-v takes 184 secs. Thus, our algorithm is up to 4-orders of magnitude faster than the SOTA methods. Push performs well on social networks, while RW performs better on road networks. This is perhaps because, on road networks, the hitting time from a node  $u$  to the landmark node set  $\mathcal{V}_l$  is relatively larger, which makes the push operation inefficient. In such cases, random walks exploit the graph more quickly. However, the Bipush algorithm takes advantage of both algorithms, which performs well on all the datasets.

### 6.3 Results of single-source ER Computation

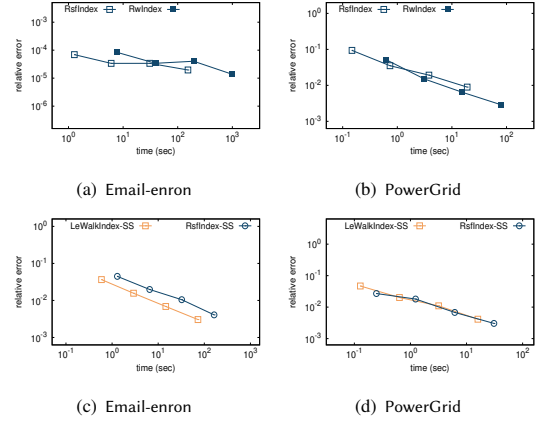
Next, we study the performance of the algorithms for single-source effective resistance computation.

**Evaluation of index building algorithms.** First, we evaluate the performance of the index building time algorithms. Notice that the performance of RsflIndex-SS and LeWalkIndex-SS for estimating the matrix  $\tilde{\mathbf{P}}$  and the pseudo-inverse  $\mathbf{L}_{\mathcal{H}}^{\dagger}$  is totally the same as RsflIndex. Thus, we focus on evaluating the performance of estimating the diagonal of  $\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1}$ . Suppose that  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$  is the estimation value of  $(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}$ , we evaluate the errors by the maximum relative error,

which is defined as  $\max_{u \in \mathcal{U}} \{ \frac{|(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} - (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}|}{(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}} \}$ . We also evaluate

the errors by the  $L1$ -error, which is defined as  $\sum_{u \in \mathcal{U}} |(\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu} - (\mathbf{L}_{\mathcal{U}\mathcal{U}}^{-1})_{uu}|$ . For the compared methods UST and LEwalk, we also evaluate the errors by the maximum relative error of the diagonal vector. We vary the parameters of all the algorithms and plot the error and time trade-offs. The results are illustrated in Fig. 5 and Fig. 6. As can be seen, both RsflIndex-SS and LeWalkIndex-SS can achieve similar errors with the same running time compared to UST and LEwalk on social networks. They are slightly faster than UST and LEwalk on road networks. RsflIndex-SS is the best algorithm. For example, on Road-PA, it takes 65 secs for LEwalk to achieve a relative error 0.12, while it takes 32 secs for LeWalkIndex-SS. Thus, our algorithm only takes 50% time. The index size of different index building algorithms is also illustrated in Fig 9(b). As can be seen, the index size of RsflIndex-SS and LeWalkIndex-SS is similar to the index size of RwlIndex and RsflIndex. Recall that the space complexity is still  $O(n \times |\mathcal{V}_l|)$ . On the largest road network Road-PA, it takes 1523 MB memory. Thus, the index size is still acceptable on large graphs.

**Evaluation of query processing algorithms.** We also study the performance of the query processing algorithms. We compare RW-SS and Push-SS with the SOTA methods RW-SS-v and Push-SS-v. We vary the parameters and plot the empirical error v.s. running time trade-offs. We evaluate the errors by the maximum relative error, we also evaluate the errors by the  $L1$ -error. The results are illustrated in Fig. 7 and Fig. 8. We can find that RW-SS and Push-SS can both answer a single-source query much faster than RW-SS-v and Push-SS-v. For example, on Road-PA, it takes 30 secs for RW-SS to achieve a relative error 0.05, while it takes  $2 \times 10^4$  secs. Thus,



**Figure 14: Effect of the accuracy of index for different query processing algorithms**

our algorithm is up to 4 orders of magnitude faster than the SOTA methods. This further confirms the effectiveness of our multiple landmark nodes approaches.

### 6.4 Effect of Parameters

In this subsection, we study the performance of our algorithms with different parameters, including (i) different settings of the landmark node set  $\mathcal{V}_l$ ; (ii) different index building algorithms to support query processing algorithms.

**Effect of  $\mathcal{V}_l$ .** First, we study the effect of the choice of the landmark node set  $\mathcal{V}_l$ . We vary both the selection strategies and  $|\mathcal{V}_l|$  (i.e. the size of the landmark node set  $\mathcal{V}_l$ ). Specifically, we vary  $|\mathcal{V}_l|$  from 1 to 20 for Email-enron, and vary  $|\mathcal{V}_l|$  from 1 to 500 for PowerGrid. Five heuristic landmark node set selection strategies including degree, pagerank, Random, degree+ and pagerank+ are evaluated. As discussed in Section 4, degree chooses the nodes with the highest degree, pagerank picks the node with the largest PageRank centrality value when  $\alpha = 0.15$ , and Random selects nodes randomly. Moreover, degree+ first picks the node with the highest degree. Then, it removes the node and its neighbors, picks the next node with the highest degree among the remaining nodes. This process repeats until  $|\mathcal{V}_l|$  nodes have been selected. pagerank+ is similar to degree+ but the node is selected by its PageRank centrality values. The results of Email-enron and PowerGrid are illustrated in Fig. 10-Fig. 13. Similar results can also be observed on the other datasets. For single-pair index building algorithms, we choose the algorithm RwlIndex as an example. The results of other index building algorithms are consistent. As can be seen in Fig. 10, the running time of RwlIndex decreases when  $|\mathcal{V}_l|$  grows larger for all the selection strategies, while the empirical error of all the algorithms is similar. Among all the landmark node selection methods, degree+ can achieve the lowest running time, followed by pagerank+, degree, pagerank. The running time of RwlIndex is very poor when using Random. The reason is that the random walk is hard to hit a small-size random chosen node set. For single-pair query processing algorithms, we choose RW as an example, the results of other query processing algorithms are consistent. The results can be found in Fig. 11, as expected, when  $|\mathcal{V}_l|$  increases, running time of the RW

decreases. When we use degree+ to select the landmark nodes, it has the smallest running time, followed by pagerank+, pagerank and degree. The empirical errors of all the algorithms are similar. For Email-enron, when  $|\mathcal{V}_l|$  exceeds 10, the running time will drop much slower. For PowerGrid, this number is around 100. After  $|\mathcal{V}_l| = 100$ , the running time of RW grows larger instead on PowerGrid. This is because the index query operations on the landmark node set related matrices also grow larger. Thus, we set  $|\mathcal{V}_l| = 10$  for social networks and  $|\mathcal{V}_l| = 100$  for road networks by default. For single-source index building algorithms, we select the results of LeWalkIndex-SS to illustrate. As shown in Fig. 12, we can find similar results that degree+ achieves the best running time while the errors are similar. For single-source query processing algorithms, we can find that in Fig. 13, if we choose  $\mathcal{V}_l$  by degree and pagerank, the performance of the algorithms will be much better than Random. degree+ and pagerank+ further improve the performance of the algorithms. For example, on Email-enron, RW-SS with degree+ is 2-orders of magnitude faster than RW-SS with Random. These results suggest that degree+ is a very good landmark node set choice heuristics in practice, which is also used in our previous experiments.

**Effect of different index building algorithms.** In former experiments, we study the query performance by using a pre-defined index structure built by RsflIndex and LeWalkIndex-SS. Here, we also study the effect of the index quality with different index building algorithms. We vary the sample size of the index building algorithms from  $10^2$  to  $10^5$ , and plot the relative error of the query processing algorithm v.s. the index building time of different algorithms. We show the results for Bipush on Email-enron and PowerGrid, the results of other algorithms and datasets are consistent. The results are illustrated in Fig. 14. As can be seen, the empirical error of the query processing algorithm Bipush decreases when the index building time increases. For single-pair query, the relative error of the query processing algorithm based on RwIndex is similar to the error based on RsflIndex when the sample size is the same, but the index building time is much longer than RsflIndex. For single-source query, the relative error of the query processing algorithm based on LeWalkIndex-SS is smaller than the error based on RsflIndex-SS, and the running time is also much smaller. Thus, we choose the index built by RsflIndex and LeWalkIndex-SS as the default index structure for single-pair and single-source query processing algorithms, respectively.

## 7 RELATED WORK

**Effectiveness resistance computation.** There exist several practical algorithms to compute single-pair (or single-source) effective resistance on large graphs [20, 29, 47]. For example, Peng et al. [29] proposed several sublinear local algorithms based on random walk sampling to compute single-pair effective resistance on the graphs with bounded mixing time. Recently, Yang et al. [47] further improved their algorithm by optimizing the lengths of random walks. However, both of these algorithms can be costly when the effective resistance between the two query nodes is large. Liao et al. [20] proposed an alternative random walk sampling approach that utilizes an easily reachable landmark node. The fundamental idea behind their algorithm is that by using the easily reachable landmark node, the random walk sampling process can be terminated more efficiently,

leading to significant improvements in efficiency. Nonetheless, their algorithm only employs a single landmark node, making it challenging to extend the method to use multiple landmark nodes. In our work, we address this challenge and also develop several efficient index algorithms for effective resistance computation.

In addition, there also exist several algorithms to compute the effective resistance of each edge (also called spanning edge centrality) in the graph [15, 28, 51]. Note that such a problem is often much easier than the problem of computing the effective resistance of any pair of query nodes (the two query nodes may not have an edge), and existing techniques for the spanning edge centrality computation cannot be used for single-pair (or single-source) effective resistance computation. It is worth mentioning that theoretically, all-pairs of effective resistance can be approximated in  $O(\frac{m \log n}{\epsilon^2} + n^2)$  time [36]. However, the algorithm proposed in [36] relies on several complicated techniques which is only with theoretical interests, and its practical performance is often very poor on large graphs [20, 29, 47].

**Personalized PageRank computation.** Our problem is also closely related to the personalized PageRank (PPR) computation problem. Random walk sampling algorithms [6, 19, 44] and push based algorithms [4, 8, 13, 25, 26, 40] are proposed for computing PPR. The state-of-the-art algorithms compute PPR by combining the two techniques [19, 21, 22, 24, 42–44, 46]. Index-based methods are also proposed by computing the information of the hub nodes [16, 35, 41, 50]. These techniques have also been generalized to compute heat Kernel PageRank [17, 48]. However, all these techniques are mainly designed for PPR computation, and it is unclear how to generalize these techniques for effective resistance computation. In this paper, we propose several novel random walk sampling and push based algorithms for computing effective resistance, which are fundamentally different from existing PPR computation approaches.

## 8 CONCLUSION

In this paper, we develop a novel index-based approach to efficiently answer both the single-pair and single-source effective resistance (ER) queries. We first propose three new ER formulas based on a newly-developed multiple landmarks technique and a concept of Schur complement. These new formulas enable us to pre-compute a compact index which contains several small-sized matrices related to the landmark nodes. We propose several efficient and provable Monte Carlo algorithms to construct the index based on the newly-established probability interpretations of the Schur complement. With this powerful index, we also develop several provable query processing algorithms to efficiently answer both the single-pair and single-source ER queries. Extensive experiments on 5 real-life graphs show that with an acceptable additional space, our algorithms can achieve up to 4 orders of magnitude faster than the state-of-the-art algorithms while maintaining the same accuracy.

## REFERENCES

- [1] Rafiq Agaev and Pavel Chebotarev. 2006. Spanning Forests of a Digraph and Their Applications. *CoRR* abs/math/0602061 (2006). arXiv:math/0602061
- [2] David J Aldous. 1990. The random walk construction of uniform spanning trees and uniform labelled trees. *SIAM Journal on Discrete Mathematics* 3, 4 (1990), 450–465.
- [3] Reid Andersen, Christian Borgs, Jennifer T. Chayes, John E. Hopcroft, Vahab S. Mirrokni, and Shang-Hua Teng. 2007. Local Computation of PageRank Contributions. In *WAW*. 150–165.

- [4] Reid Andersen, Fan R. K. Chung, and Kevin J. Lang. 2006. Local Graph Partitioning using PageRank Vectors. In *FOCS*. 475–486.
- [5] Luca Avena and Alexandre Gaudillière. 2018. Two applications of random spanning forests. *Journal of Theoretical Probability* 31, 4 (2018), 1975–2004.
- [6] Konstantin Avrachenkov, Nelly Litvak, Danil Nemirovsky, and Natalia Osipova. 2007. Monte Carlo Methods in PageRank Computation: When One Iteration is Sufficient. *SIAM J. Numer. Anal.* 45, 2 (2007), 890–904.
- [7] Ravindra B Bapat. 2010. *Graphs and matrices*. Vol. 27. Springer.
- [8] Pavel Berkhin. 2006. Bookmark-Coloring Approach to Personalized PageRank Computing. *Internet Math.* 3, 1 (2006), 41–62.
- [9] Béla Bollobás. 1998. *Modern graph theory*. Vol. 184. Springer Science & Business Media.
- [10] Seth Chaiken. 1982. A combinatorial proof of the all minors matrix tree theorem. *SIAM Journal on Algebraic Discrete Methods* 3, 3 (1982), 319–329.
- [11] Fan R. K. Chung and Lincoln Lu. 2006. Survey: Concentration Inequalities and Martingale Inequalities: A Survey. *Internet Math.* 3, 1 (2006), 79–127.
- [12] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. 2022. *Introduction to algorithms*. MIT press.
- [13] Mustafa Coskun, Ananth Grama, and Mehmet Koyutürk. 2018. Indexed Fast Network Proximity Querying. *VLDB* 11, 8 (2018), 840–852.
- [14] Gramoz Goranci, Monika Henzinger, and Pan Peng. 2018. Dynamic Effective Resistances and Approximate Schur Complement on Separable Graphs. In *26th Annual European Symposium on Algorithms, ESA 2018, August 20–22, 2018, Helsinki, Finland (LIPIcs, Vol. 112)*. 40:1–40:15.
- [15] Takanori Hayashi, Takuya Akiba, and Yuichi Yoshida. 2016. Efficient Algorithms for Spanning Tree Centrality. In *IJCAI*. 3733–3739.
- [16] Jinhong Jung, Namyoung Park, Lee Sael, and U Kang. 2017. BePI: Fast and Memory-Efficient Method for Billion-Scale Random Walk with Restart. In *SIGMOD*. 789–804.
- [17] Kyle Kloster and David F. Gleich. 2014. Heat kernel based community detection. In *KDD*. ACM, 1386–1395.
- [18] Jure Leskovec and Andrej Krevl. 2014. SNAP Datasets: Stanford Large Network Dataset Collection. <http://snap.stanford.edu/data>.
- [19] Meihao Liao, Rong-Hua Li, Qiangqiang Dai, Hongyang Chen, Hongchao Qin, and Guoren Wang. 2023. Efficient Personalized PageRank Computation: The Power of Variance-Reduced Monte Carlo Approaches. *Proc. ACM Manag. Data* 1, 2 (2023), 160:1–160:26.
- [20] Meihao Liao, Rong-Hua Li, Qiangqiang Dai, Hongyang Chen, Hongchao Qin, and Guoren Wang. 2023. Efficient Resistance Distance Computation: The Power of Landmark-based Approaches. *Proc. ACM Manag. Data* 1, 1 (2023), 68:1–68:27.
- [21] Meihao Liao, Rong-Hua Li, Qiangqiang Dai, and Guoren Wang. 2022. Efficient Personalized PageRank Computation: A Spanning Forests Sampling Based Approach. In *SIGMOD*. ACM, 2048–2061.
- [22] Dandan Lin, Raymond Chi-Wing Wong, Min Xie, and Victor Junqiu Wei. 2020. Index-Free Approach with Theoretical Guarantee for Efficient Random Walk with Restart Query. In *ICDE*. 913–924.
- [23] Yang Liu, Chuan Zhou, Shirui Pan, Jia Wu, Zhao Li, Hongyang Chen, and Peng Zhang. 2023. CurvDrop: A Ricci Curvature Based Approach to Prevent Graph Neural Networks from Over-Smoothing and Over-Squashing. In *WWW*. ACM, 221–230.
- [24] Peter Lofgren, Siddhartha Banerjee, and Ashish Goel. 2016. Personalized PageRank Estimation and Search: A Bidirectional Approach. In *WSDM*. 163–172.
- [25] Peter Lofgren and Ashish Goel. 2013. Personalized PageRank to a Target Node. *CoRR* abs/1304.4658 (2013). [arXiv:1304.4658](http://arxiv.org/abs/1304.4658) <http://arxiv.org/abs/1304.4658>
- [26] Takanori Maehara, Takuya Akiba, Yoichi Iwata, and Ken-ichi Kawarabayashi. 2014. Computing Personalized PageRank Quickly by Exploiting Graph Structures. *VLDB* 7, 12 (2014), 1023–1034.
- [27] Philippe Marchal. 2000. Loop-erased random walks, spanning trees and Hamiltonian cycles. *Electronic Communications in Probability* 5 (2000), 39–50.
- [28] Charalampos Mavroforakis, Richard Garcia-Lebron, Ioannis Koutis, and Evimaria Terzi. 2015. Spanning Edge Centrality: Large-scale Computation and Applications. In *WWW*. 732–742.
- [29] Pan Peng, Daniel Lopatta, Yuichi Yoshida, and Gramoz Goranci. 2021. Local Algorithms for Estimating Effective Resistance. In *KDD*. 1329–1338.
- [30] Jim Pitman and Wenpin Tang. 2018. Tree formulas, mean first passage times and Kemeny’s constant of a Markov chain. (2018).
- [31] Jim Pitman and Wenpin Tang. 2018. Tree formulas, mean first passage times and Kemeny’s constant of a Markov chain. *Bernoulli* 24, 3 (2018), 1942 – 1972.
- [32] Purnamrita Sarkar, Andrew W. Moore, and Amit Prakash. 2008. Fast incremental proximity search in large graphs. In *ICML*.
- [33] Aaron Schild. 2018. An almost-linear time algorithm for uniform random spanning tree generation. In *STOC*. 214–227.
- [34] Jieming Shi, Nikos Mamoulis, Dingming Wu, and David W. Cheung. 2014. Density-based place clustering in geo-social networks. In *SIGMOD*. ACM, 99–110.
- [35] Kijung Shin, Jinhong Jung, Lee Sael, and U Kang. 2015. BEAR: Block Elimination Approach for Random Walk with Restart on Large Graphs. In *SIGMOD*, Timos K. Sellis, Susan B. Davidson, and Zachary G. Ives (Eds.). 1571–1585.
- [36] Daniel A. Spielman and Nikhil Srivastava. 2008. Graph sparsification by effective resistances. In *STOC*. ACM, 563–568.
- [37] Kumar Sricharan and Kamalika Das. 2014. Localizing anomalous changes in time-evolving graphs. In *SIGMOD*. ACM, 1347–1358.
- [38] Prasad Tetali. 1991. Random walks and the effective resistance of networks. *Journal of Theoretical Probability* 4, 1 (1991), 101–109.
- [39] Jake Topping, Francesco Di Giovanni, Benjamin Paul Chamberlain, Xiaowen Dong, and Michael M. Bronstein. 2022. Understanding over-squashing and bottlenecks on graphs via curvature. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25–29, 2022*.
- [40] Hanzhi Wang, Zhewei Wei, Junhao Gan, Sibowang, and Zengfeng Huang. 2020. Personalized PageRank to a Target Node, Revisited. In *KDD*. 657–667.
- [41] Sibowang, Youze Tang, Xiaokui Xiao, Yin Yang, and Zengxiang Li. 2016. HubPPR: Effective Indexing for Approximate Personalized PageRank. *VLDB* 10, 3 (2016), 205–216.
- [42] Sibowang and Yufei Tao. 2018. Efficient Algorithms for Finding Approximate Heavy Hitters in Personalized PageRanks. In *SIGMOD*. 1113–1127.
- [43] Sibowang, Renchi Yang, Xiaokui Xiao, Zhewei Wei, and Yin Yang. 2017. FORA: Simple and Effective Approximate Single-Source Personalized PageRank. In *KDD*. 505–514.
- [44] Zhewei Wei, Xiaodong He, Xiaokui Xiao, Sibowang, Shuo Shang, and Ji-Rong Wen. 2018. TopPPR: Top-k Personalized PageRank Queries with Precision Guarantees on Large Graphs. In *SIGMOD*. 441–456.
- [45] David Bruce Wilson. 1996. Generating random spanning trees more quickly than the cover time. In *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*. 296–303.
- [46] Hao Wu, Junhao Gan, Zhewei Wei, and Rui Zhang. 2021. Unifying the Global and Local Approaches: An Efficient Power Iteration with Forward Push. In *SIGMOD*. 1996–2008.
- [47] Renchi Yang and Jing Tang. 2023. Efficient Estimation of Pairwise Effective Resistance. *Proc. ACM Manag. Data* 1, 1 (2023), 16:1–16:27.
- [48] Renchi Yang, Xiaokui Xiao, Zhewei Wei, Sourav S. Bhowmick, Jun Zhao, and Rong-Hua Li. 2019. Efficient Estimation of Heat Kernel PageRank for Local Clustering. In *SIGMOD*. ACM, 1339–1356.
- [49] Hongzhi Yin, Bin Cui, Jing Li, Junjie Yao, and Chen Chen. 2012. Challenging the Long Tail Recommendation. *VLDB* 5, 9 (2012), 896–907.
- [50] Minji Yoon, Jinhong Jung, and U Kang. 2018. TPA: Fast, Scalable, and Accurate Method for Approximate Random Walk with Restart on Billion Scale Graphs. In *ICDE*. 1132–1143.
- [51] Shiqi Zhang, Renchi Yang, Jing Tang, Xiaokui Xiao, and Bo Tang. 2023. Efficient Approximation Algorithms for Spanning Centrality. In *KDD*. ACM, 3386–3395.