

# Efficient and Accurate SimRank-based Similarity Joins: Experiments, Analysis, and Improvement [Experiment, Analysis & Benchmark]

Qian Ge<sup>†</sup>  
Peking University  
geqian@pku.edu.cn

Lei Zou  
Peking University  
zoulei@pku.edu.cn

Yu Liu<sup>†§</sup>  
Yinghao Zhao  
Yuetian Sun  
Beijing Jiaotong University  
{yul,yinghaozhao,yuetiansun}@bjtu.edu.cn

Yuxing Chen  
Anqun Pan  
Tencent Inc.  
{axingguchen,aaronpan}@tencent.com

## ABSTRACT

SimRank-based similarity joins, which mainly include threshold-based and top- $k$  similarity joins, are important types of all-pair SimRank queries. Although a line of related algorithms have been proposed recently, they still fall short of providing approximation guarantee and suffer from scalability issues on medium and large graphs. Meanwhile, we also lack an extensive analysis of existing techniques in terms of accuracy and efficiency. Motivated by these challenges, we first conduct detailed analysis of state-of-the-art algorithms and provide additional theoretical results. Second, to address the limitations of existing techniques, we propose simple yet effective algorithm frameworks for both queries to theoretically guarantee the approximation bound, and present a more efficient all-pair algorithm inspired by randomized local push of Personalized PageRank. Next, we analyze the algorithmic complexity of threshold-based and top- $k$  similarity joins by leveraging a reasonable assumption of SimRank distribution. Through extensive experiments, we find that our proposed methods far exceed existing ones with respect to query efficiency, approximation guarantee and practical accuracy, while our theoretical analysis nicely matches the empirical study.

### PVLDB Reference Format:

Qian Ge<sup>†</sup>, Yu Liu<sup>†§</sup>, Yinghao Zhao, Yuetian Sun, Lei Zou, Yuxing Chen, and Anqun Pan. Efficient and Accurate SimRank-based Similarity Joins: Experiments, Analysis, and Improvement [Experiment, Analysis & Benchmark]. PVLDB, 14(1): XXX-XXX, 2020.  
doi:XX.XX/XXX.XX

### PVLDB Artifact Availability:

The source code and data have been made available at <https://github.com/xinghun0525/R2LP>.

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 14, No. 1 ISSN 2150-8097.  
doi:XX.XX/XXX.XX

<sup>†</sup>These authors contributed equally to this work.

<sup>§</sup>Corresponding author.

## 1 INTRODUCTION

SimRank [12] is one of the most important measures for pairwise node similarity, which finds its applications in recommender systems [19], link prediction [26], and graph embeddings [33]. The SimRank similarity of node  $u$  and  $v$  is defined by the following recursive equation, with the intuition that “two nodes in a graph are similar if their in-neighbors are similar, and a node is most similar to itself”:

$$s(u, v) = \begin{cases} 1, & \text{if } u = v \\ \frac{c}{|I(u)||I(v)|} \sum_{x \in I(u)} \sum_{y \in I(v)} s(x, y), & \text{otherwise.} \end{cases} \quad (1)$$

Here,  $I(u)$  denote the in-neighbors of node  $u$ , and  $c \in (0, 1)$  is the decay factor, which is typically set to 0.6 [32, 35, 38] or 0.8 [12]. Due to its computational complexity, efficient computation of SimRank has been extensively studied in the past two decades [8, 16, 18, 21, 23, 25, 27, 31], and remains to be a hot topic in the very recent years [22, 30, 35, 38, 46].

According to the query inputs and outputs, SimRank queries can be broadly categorized as *single-pair* queries, *single-source* queries, and *all-pair* queries (see Section 7 for more details). In particular, the all-pair query computes the estimated SimRank values for each pair of nodes (with some additive error  $\varepsilon$ ). On small graphs, e.g., with tens of thousands of nodes, the *Power Method* [12] is sufficient to compute the *ground truth* SimRank values<sup>1</sup>, which guarantee a very small absolute error  $\varepsilon_{min}$  such as  $10^{-7}$  [35] or  $10^{-11}$  [32]. However, this is infeasible for larger graphs because of excessive computational and memory cost. Among existing works [20, 25, 27, 28, 31, 38, 44–46] that tackle all-pair queries, state-of-the-art algorithms [38, 46] can hardly guarantee an additive error of  $\varepsilon = 0.01$  on million-node graphs. Besides, a latest work [35] also indicates that it is “essentially hopeless” to compute all-pair queries on large graphs, as the output size can be as large as  $O(n^2)$ . Fortunately, two modified versions of all-pair queries have been proposed and studied, namely, *threshold-based* [28, 44, 45] and *top- $k$*  [20, 31, 46] *similarity join*, of which the output sizes can be significantly smaller. Both queries

<sup>1</sup>We distinguish between the *exact* and the *ground truth* SimRank values, where the former is the exact solution to Equation 1 and the latter is an accurate approximation.

are more useful in practice because we are only interested in the node pairs with non-negligible similarities [46].

**DEFINITION 1 (THRESHOLD-BASED SIMILARITY JOIN).** *Given a graph  $G = (V, E)$  and a threshold  $\theta > 0$ , return the set of node pairs  $R(\theta)$  with  $u \neq v$  and SimRank value  $s(u, v) \geq \theta$ .*

**DEFINITION 2 (TOP- $k$  SIMILARITY JOIN).** *Given a graph  $G = (V, E)$  and a parameter  $k > 0$ , return a  $k$ -sized set  $R(k)$  containing the node pairs with the top- $k$  largest SimRank values among all pairs  $\{(u, v) | u, v \in V, u \neq v\}$ .*

Following previous works [12, 35], we allow a small error parameter  $\varepsilon_{min}$  (e.g.,  $10^{-7}$ ) in SimRank estimation. Considering the hardness of the problem, we additionally adopt a parameter  $\rho \in (0, 1)$  named *approximation bound*.

**DEFINITION 3 (APPROXIMATION BOUND).** *Given an approximation parameter  $\rho \in (0, 1)$ , we say an algorithm  $\mathcal{A}$  has approximation bound  $\rho$  for threshold-based (resp. top- $k$ ) similarity join of SimRank, if for any input graph  $G$ , the returned set of node pairs contains at least  $\rho$  fraction of the ground truth answer.*

## 1.1 Motivations

We note that existing algorithms for threshold-based and top- $k$  similarity joins significantly fall short of both query efficiency and approximation guarantee. We identify the following issues.

**Motivation 1. Lack of experimental analysis of efficiency and accuracy, especially on medium and large graphs.** Even for state-of-the-art algorithms [38, 46] of general all-pair queries, the result accuracy is not evaluated except for small graphs due to lack of ground truth. Other representative algorithms [20, 44] even do not evaluate accuracy on small graphs.

**Motivation 2. Lack of approximation guarantee for threshold-based and top- $k$  similarity joins.** Surprisingly, there does not exist an algorithm that guarantees the approximation bound for threshold-based or top- $k$  similarity join. Algorithms with absolute error guarantee cannot be directly extended to achieve this goal.

**Motivation 3. Lack of understanding of the computational complexity.** To the best of our knowledge, it is unclear how the input parameters and graph properties impact the algorithmic complexity of both queries. The complexity analysis for general all-pair algorithms does not reveal the intrinsic hardness of threshold-based and top- $k$  similarity joins.

## 1.2 Contributions

For both types of all-pair SimRank queries, we analyze state-of-the-art algorithms in detail, provide improved algorithms with approximation guarantee and lower theoretical complexity, and discuss the problem complexity w.r.t. the input parameters and graph properties. Our main contributions are summarized as follows.

**Detailed analysis of state-of-the-art algorithms.** We choose four representative algorithms, *UISim* [46], *FLP & Opt-LP* [38], *H-go SRJ* [44], and *KSimJoin* [20] for comparison. *UISim* and *FLP & Opt-LP* are state-of-the-art algorithms for general all-pair queries, which are directly extended to answer SimRank-based similarity joins [46]. *H-go SRJ* and *KSimJoin* are the state-of-the-art algorithms for threshold-based and top- $k$  similarity joins, respectively. For each algorithm, we discuss its basic idea, theoretical guarantee,

complexity analysis and empirical study. Specifically, we propose additional theoretical results (w.r.t. complexity or error guarantee) of these algorithms that are incorrect or missed in the original papers, and discuss their limitations in answering similarity join queries.

**Improved algorithms for SimRank-based similarity joins.** We first propose two simple yet effective algorithm frameworks for both queries respectively with approximation bound. The basic idea is to invoke an efficient algorithm for all-pair queries that estimates SimRank values within  $\varepsilon$  error, and to gradually shrink  $\varepsilon$  until the approximation bound is theoretically guaranteed. Next, we utilize the equivalence of all-pair SimRank on the input graph  $G$  and single-target Personalized PageRank (PPR) on the *SimRank graph*  $G^S$  [12, 38], and devise Randomized Reverse Local Push ( $R^2LP$ ) inspired by randomized backward push for PPR [34]. It improves the all-pair query complexities from  $O(\frac{\sum_{u,v \in V} d_{in}(u)d_{in}(v)s(u,v)}{\varepsilon})$  (the best known algorithm [38]) to  $\tilde{O}(\frac{\sum_{u,v \in V} \sqrt{d_{in}(u)d_{in}(v)s(u,v)}}{\varepsilon})$ . We further propose a pruning strategy to dramatically enhance practical efficiency while retaining error guarantee. Our empirical study shows that  $R^2LP$  significantly improves query efficiency for general all-pair queries and similarity joins compared to existing methods.

**Complexity analysis by utilizing SimRank distribution.** We find that it is almost impossible to give a non-trivial complexity bound for both similarity join queries without any assumption of SimRank distribution. On the other hand, our empirical analysis validates that for most real-world graphs, the ground truth SimRank values follow some *generalized* version of power-law distribution [2, 5], which is consistent with previous studies [22, 35]. We leverage this property to propose verifiable assumptions of SimRank distribution and present complexity analysis for our algorithms, which serves as the upper bound of the algorithmic complexity of SimRank-based similarity joins. Our theoretical results nicely match the empirical study in Section 6.

**Extensive experiments on medium and large graphs.** We conduct extensive experiments on small, medium, and large sized graphs and compare our algorithms with state-of-the-art methods. As far as we know, we first employ medium and large datasets to systematically evaluate threshold-based and top- $k$  similarity joins in terms of running time, absolute error and query accuracy, and demonstrate the necessities of holding approximation bound as well as devising more efficient algorithms. Experimental study demonstrates that our algorithms outperform existing ones in most cases, achieving better accuracy while being faster by up to an order of magnitude. More importantly, our empirical findings validate that it is the skewness of SimRank distribution that mainly determines the problem hardness. This results verifies the effectiveness of our theoretical analysis.

## 2 PRELIMINARIES

The definition in Equation 1 can be interpreted via the *node-pairs graph*  $G^2$  [12]. Each node in  $G^2$  represents a pair of nodes  $(u, v)$  in  $G$ . If there exist edges  $(u, v)$  and  $(x, y)$  in  $G$ , edge  $((u, x), (v, y))$  is included in  $G^2$ . Specifically, every node  $(v, v)$  is called the *singleton node*. To this end, SimRank can be thought of as propagating the similarity among nodes in  $G^2$ , starting from all the singleton nodes.

**Matrix formation of SimRank.** The recursive definition can be represented in matrix formation [43], i.e.,  $S = cP^\top SP \vee I$ , where  $S$  is an  $n \times n$ -dimensional similarity matrix with  $S[u, v] = s(u, v)$ ,

**Table 1: Table of notations.**

Notation	Description
$G = (V, E)$	A graph with node set $V$ and edge set $E$ , where $n =  V $ and $m =  E $
$G^2, G^s$	The node-pairs graph [12] and the SimRank graph [38]
$u, v$	Nodes in $G$
$I(u), O(u)$	The incoming and outgoing neighbors of $u$ , and we have $d_{in}(u) =  I(u) , d_{out}(u) =  O(u) $
$s(u, v)$	The SimRank value, i.e., the solution of Equation 1
$s^*(u, v)$	The ground truth SimRank value, with $ s^*(u, v) - s(u, v)  \leq \varepsilon_{min}$ (with high probability)
$\hat{s}(u, v)$	The estimated SimRank value
$\varepsilon, \delta$	The additive error and the failure probability
$\theta$	The input parameter of threshold-based similarity join
$k$	The input parameter of top- $k$ similarity join
$\rho$	The approximation bound

$P$  is the column-normalized adjacency matrix of  $G$ ,  $I$  is the identity matrix, and operator  $\vee$  denotes element-wise maximum. The all-pair ground truth of SimRank can be computed by the *Power Method* [12], which is essentially a fixed-point iteration process on  $G^2$ . Recent literature [18, 25, 38] further derives various close-form solutions based on this formation.

**The random walk interpretation.** Jeh and Widom explain the SimRank definition from random walk perspective in their original paper [12]. We adopt  $\sqrt{c}$ -walk, its revised version following [13, 21, 32]. Let  $W_{\sqrt{c}}(u)$  denote a random walk starting from node  $u$  and following incoming edges at each step with walking probability  $\sqrt{c}$ . More precisely, at each step, with probability  $\sqrt{c}$ , it randomly chooses an in-neighbor and continues the walk, otherwise the walk terminates at the current node. [32] proves that the SimRank of  $u$  and  $v$  equals the probability that two  $\sqrt{c}$ -walks  $W_{\sqrt{c}}(u)$  and  $W_{\sqrt{c}}(v)$  meet.

Table 1 lists the frequently used notations throughout the paper.

### 3 ANALYSIS OF STATE OF THE ART

We present detailed comparison and analysis of state-of-the-art algorithms for threshold-based and top- $k$  similarity joins. Specifically, we choose *UISim* [46], *FLP & Opt-LP* [38], *H-go SRJ* [44] and *KSimJoin* [20] because they have not been outperformed or evaluated by previous empirical analysis. We illustrate their basic ideas in Figure 1, and summarize the theoretical guarantee, time complexity, and previous empirical study in Table 2.

#### 3.1 UISim

The algorithm depends on the random walk interpretation of SimRank and adopts a *hub-based strategy* to prioritize all possible random walk pairs. First, the hub set  $H$  is constructed from the top- $|H|$  nodes with the largest *in-degrees*. We say an  $l$ -length random walk  $W(u) = (p_1 = u, p_2, \dots, p_l)$  from  $u$  has *hub length*  $k$  if  $|H \cap \{p_2, \dots, p_{l-1}\}| = k$ . Intuitively, large hub length implies small walking probability. Second, the SimRank value of two nodes  $u$  and  $v$  is decomposed according to the hub length:

$$\hat{s}^{(\eta)}(u, v) = \sum_{k=0}^{\eta} \sum_{\max(i,j) \leq k} R(P_u^i \bowtie P_v^j), \quad (2)$$

where  $\eta$  denotes the maximum hub length,  $P_u^i$  denotes the set of random walks that start from  $u$  and are of hub length  $i$ , and the join operator assembles pairs of random walks from  $P_u^i$  and  $P_v^j$  that have the same walk length and meet at the same node. Given the set of assembled walk pairs, the  $R(\cdot)$  operator returns the sum of their meeting probabilities. As demonstrated in Equation 2, *UISim* assembles random walks in ascending order of hub length to prioritize them roughly by the meeting probabilities.

In practice, *UISim* introduces the notion of *prime subgraphs* and leverages graph expansion to estimate SimRank values incrementally (Figure 1(a)). The prime in-subgraph (resp. out-subgraph) of node  $v$  contains all reachable nodes starting from  $v$  and following in-edges (resp. out-edges), until a hub node is encountered. For all-pair SimRank query, *UISim* first considers the prime out-subgraphs of all nodes with at least two out-neighbors (denoted as  $V_{\geq 2}^o$ ), and assembles all possible pairs of random walks following Equation 2 (with  $k = 0$ ). Then, each subgraph is expanded with the prime out-subgraphs of the encountered hub nodes, which corresponds to increasing  $k$  by one. The algorithm terminates when  $k$  reaches  $\eta$ .

**Theoretical Guarantee.** *UISim* relaxes the first-meeting constraints of SimRank definition<sup>2</sup>, and conducts analysis of the expected error bound based on this assumption. Generally, it does not guarantee an absolute error of  $\varepsilon$  for SimRank estimation  $\hat{s}(u, v)$ , which says  $|\hat{s}(u, v) - s(u, v)| \leq \varepsilon$ .

**Complexity Analysis.** According to [46], the time complexity of *UISim* for all-pair queries is bounded by  $O(|V_{\geq 2}^o||H|^{\eta}T)$  in expectation, where  $T = O\left(d^L\left(1 - \sum_{v \in H} d_v/2m\right)^L\right)$  stands for the expected number of random walks up to length  $L$ . It is assumed that every prime subgraph has  $O(|H|)$  hubs.

**Empirical Analysis.** *UISim* is evaluated for all-pair queries (in terms of absolute error) and for top- $k$  similarity joins (in terms of Precision@ $k$ ). However, the evaluation is only conducted on small graphs with tens of thousands nodes.

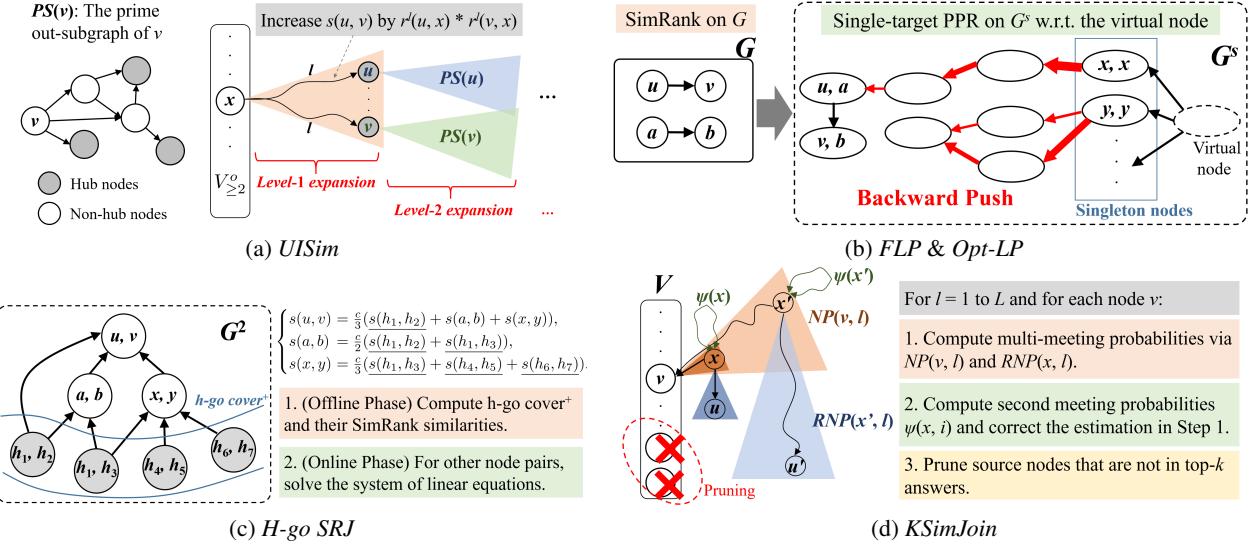
#### 3.2 FLP & Opt-LP

Motivated by a modified version of the node-pairs graph, namely, the *SimRank graph*, Wang et al. [38] propose a new linear system of SimRank with the Kronecker product and vector operators:

$$[I - c(I - E)(P^\top \otimes P^\top)]\vec{v}\vec{c}(S) = \vec{v}\vec{c}(I), \quad (3)$$

where  $E = \text{Diag}(\vec{v}\vec{c}(I))$ ,  $\otimes$  denotes the Kronecker product,  $S$  is the SimRank matrix with  $S[u, v] = s(u, v)$ ,  $\vec{v}\vec{c}(\cdot)$  transforms an  $N \times N$ -dimensional matrix to an  $N^2$ -dimensional vector by concatenating its columns, and  $\text{Diag}(\cdot)$  transforms an  $N'$ -dimensional vector to an  $N' \times N'$ -dimensional diagonal matrix. Based on that, a local push algorithm *ForwardLocalPush* (*FLP*) is proposed. The authors further devise a reduced linear system to avoid the redundant computation of both  $s(u, v)$  and  $s(v, u)$ , and to effectively speed up the computation involving self-loops. The optimized algorithm is referred to as *OptimizeLocalPush* (*Opt-LP*). We propose a slightly modified version of the SimRank graph [38].

<sup>2</sup>Note that [46] provides a correction method to exclude *multi-hop double-meeting probabilities*, i.e., the probability of two random walks meeting exactly twice. It is non-trivial to extend the method to *multi-meeting probabilities*.



**Figure 1: Overview of state-of-the-art algorithms for SimRank-based similarity joins.**

**Table 2: Comparison of state of the art. N/A means that the specific query is not studied in the original paper.**

Algorithm	Theoretical Guarantee			Time Complexity	Empirical Study		
	Absolute	Threshold	Top- $k$		Absolute	Threshold	Top- $k$
<i>UISim</i> [46]	✗	N/A	✗	$O( V_{\geq 2}^o H ^nT)$	✓	N/A	✓
<i>FLP &amp; Opt-LP</i> [38]	✓	N/A	N/A	$O(\frac{\sum_{u,v \in V} d_{in}(u)d_{in}(v)s(u,v)}{\epsilon})$	✓	N/A	N/A
<i>H-go SRJ</i> [44]	✓	✗	N/A	$\tilde{O}(nd^h + \frac{n H(G) }{\epsilon^2})$ (offline), $O(nd^{h+3})$ (online)	N/A	✓(only query time)	✓(only query time)
<i>KSimJoin</i> [20]	✓	N/A	✗	$O(nd^{\bar{d}\log\frac{1}{\epsilon}})$	N/A	N/A	✓(only query time)

**DEFINITION 4 (SIMRANK GRAPH).** Given a graph  $G$ , its SimRank graph  $G^s$  is obtained from the node-pairs graph  $G^2$  by removing all the in-edges of the singleton nodes. Besides, we modify  $G^s$  by adding a virtual node  $(v_r, v_r)$ , and for each singleton node  $(v, v)$ , we add one in-edge from it to the virtual node.

Noticing the close relation of SimRank to PPR, [38] claims that the algorithm is essentially the *Backward Push* [24] for single-target PPR queries, except that the PPR is defined following in-edges rather than out-edges, which is referred to as the *reverse PPR* (Figure 1(b)). In practice, the algorithm does not have to materialize  $G^s$ , and we demonstrate it in Figure 1(b) as the dashed rectangle. We reformalize the claim in [38] as follows. For space constraint, we refer all proofs to the full version of the paper [1].

**LEMMA 1.** The SimRank similarity of  $u$  and  $v$  is exactly  $\frac{1}{c(1-c)}$  times the reverse PPR value of node  $(u, v)$  w.r.t. the target node  $(v_r, v_r)$  on SimRank graph  $G^s$  and with stopping probability  $1 - c$ :

$$s(u, v) = \frac{1}{c(1-c)} \pi_{G^s}((u, v), (v_r, v_r)). \quad (4)$$

**Theoretical Guarantee.** FLP & Opt-LP achieve absolute error guarantee following *Backward Push*. The algorithm takes an error parameter  $\epsilon$  as input, and guarantees  $|\hat{s}(u, v) - s(u, v)| \leq \epsilon$  for each node pair  $(u, v)$ .

**Complexity Analysis.** We observe that the proposed average time complexity  $O(\frac{\bar{d}^2}{\epsilon})$  (Proposition 1 of [38]) is *incorrect* where  $\bar{d}$  denotes the average degree of  $G$ . As with *Backward Push* [24], the average complexity is achieved by considering all possible target nodes in the graph. However, as Lemma 1 shows, the target node is *fixed* as the virtual node. We fix the issue using the following lemma.

**LEMMA 2.** To achieve an absolute error of  $\epsilon$  for all-pair queries, following [24, 34], the time complexity of FLP and Opt-LP are bounded by  $O(\frac{\sum_{u,v \in V} d_{in}(u)d_{in}(v)s(u,v)}{\epsilon})$ .

**Empirical Analysis.** [38] conducts experiments for all-pair queries on both small and large graphs with respect to query time, but only evaluates estimation error on small graphs with thousands of nodes. Due to the deficiency of *Backward Push* and as our experiment shows, the algorithms have scalability issues on million-node graphs.

### 3.3 H-go SRJ

*H-go SRJ* [44] is an index-based algorithm for threshold-based similarity joins. Inspired by the distance-based upper bound of SimRank, it presents an index called *h-go cover*<sup>+</sup> for the node-pairs graph  $G^2$ , so that for any node pair, the SimRank value is either recorded by the index or can be recovered from the indexed neighbors with at most  $h$ -hops. As shown in Figure 1(c), in the offline phase, the

algorithm first computes an h-go cover<sup>+</sup> without materializing the node-pairs graph. For each indexed node pair, its SimRank value can be estimated by a single-pair query. In the online phase, for every non-indexed node pair, *H-go SRJ* first generates an extending tree rooted at this node pair, then builds a system of linear equations following the recursive definition of SimRank and with the help of the index. The equations can be solved using Gaussian elimination. **Theoretical Guarantee.** [44] does not discuss the error guarantee of *H-go SRJ*. We remedy this issue with the following analysis. Given an error parameter  $\varepsilon$ , to estimate each indexed node pair with absolute error guarantee, we invoke the Monte Carlo algorithm to sample  $\tilde{O}(\frac{1}{\varepsilon^2})$  pairs of random walks (the  $\tilde{O}$  notation ignores log factors). The following lemma bounds the estimation error of the online phase.

**LEMMA 3.** *If all indexed SimRank similarities have at most  $\varepsilon$  additive error, then this  $\varepsilon$  error guarantee also holds for the similarities of any node pair computed in the online phase.*

**Complexity Analysis.** According to [44], the h-go cover<sup>+</sup> selection takes  $O(n\bar{d}^h)$  time where  $\bar{d}$  is the average degree of  $G$ . We notice that computing SimRank estimation for the indexed node pairs takes an extra  $\tilde{O}(\frac{n|H(G)|}{\varepsilon^2})$  time, where  $H(G)$  denotes the h-go cover<sup>+</sup> of  $G$ . Therefore, the offline phase needs  $\tilde{O}(n\bar{d}^h + n\frac{|H(G)|}{\varepsilon^2})$  time. For the online phase, given a node pair, generating the extending tree needs  $O(\bar{d}^{2h})$  time, while solving the linear equations incurs  $O(\bar{d}^{2h+1})$  time with all optimizations [44]. Since the linear system contains  $O(\bar{d}^{2(h-1)})$  unknown values on average, the amortized cost for computing each non-indexed node pair is  $O(\bar{d}^3)$ . Assume there exist  $O(n\bar{d}^h)$  non-indexed node pairs after distance-based pruning, the online phase needs  $O(n\bar{d}^h \cdot \bar{d}^3) = O(n\bar{d}^{h+3})$  time.

**Empirical Analysis.** The authors of *H-go SRJ* conduct experiments for query efficiency on graphs with up to millions of nodes. However, they do not evaluate the accuracy of the proposed method.

### 3.4 KSimJoin

Li et al. [20] propose an incremental framework for top- $k$  similarity joins based on the random walk interpretation (Figure 1(d)). The key idea is to first compute all meeting probabilities and then exclude the multi-meeting cases. In particular, it decomposes the SimRank value by meeting steps:  $s(u, v) = \sum_{l=0}^{\infty} s^{(l)}(u, v)$ . Here,  $s^{(l)}(u, v)$  is the first meeting probability of two  $\sqrt{c}$ -walks from  $u$  and  $v$  at exact  $l$ -th step. To estimate  $s^{(l)}(u, v)$ , the following equation is adopted:

$$s^{(l)}(u, v) = \sum_{x \in V} \Pr((u, v) \xrightarrow{l} (x, x)) - \sum_{i \in [1, l-1], y \in V} \Pr((u, v) \xrightarrow{i} (y, y)) \psi(y, l-i), \quad (5)$$

where  $\psi(y, l-i)$  is the probability of two  $\sqrt{c}$ -walks starting from  $y$  first meeting at exact  $(l-i)$ -th step, which is referred to as the *second meeting probability*.

To implement the above idea, the *KSimJoin* algorithm estimates  $s^{(l)}(u, v)$  in ascending order of  $l$  for  $l \in [1, L]$ . In  $l$ -th iteration, for every node  $v$  in the remaining node set  $R$  (where  $R$  is initialized as  $V$ ), it first performs a breadth-first search (BFS) from  $v$  following in-edges. All visited node at step  $l$  is referred to as  $NP(v, l)$ . Then,

for each node  $x \in NP(v, l)$ , a BFS following out-edges is conducted to find all node  $u$  that is  $l$ -step from  $x$ , denoted as  $RNP(x, l)$ . The multi-meeting probabilities are computed from  $NP(v, l)$  and  $RNP(x, l)$  for all  $x \in NP(v, l)$ . By leveraging Equation 5, it excludes all multi-meeting cases, where the second meeting probabilities  $\psi(x, i)$  for each  $x$  and  $i \in [1, l-1]$  can be computed following a similar BFS-based and layer-by-layer iteration approach. For top- $k$  queries, after each iteration of  $l$ , *KSimJoin* uses a carefully designed upper bound to prune nodes in  $R$  that cannot be the top- $k$  answers. The upper bound is computed by relaxing Equation 5.

**Theoretical Guarantee.** As [20] does not analyze the error guarantee of *KSimJoin*, we formally present the following lemma, which states that *KSimJoin* is essentially equivalent to the *Power Method*.

**LEMMA 4.** *By setting  $L = O(\log \frac{1}{\varepsilon})$ , *KSimJoin* achieves  $O(\varepsilon)$  absolute error guarantee.*

**Complexity Analysis.** Computing the multi-meeting probabilities and the upper bound incurs  $O(|R|\bar{d}^L)$  time [20]. The cost of estimating the second meeting probabilities is asymptotically identical. Since  $L = O(\log \frac{1}{\varepsilon})$  and  $|R| = O(n)$ , the overall time complexity is  $O(n\bar{d}^{\log \frac{1}{\varepsilon}})$ .

**Empirical Analysis.** Empirical study of [20] only focuses on the efficiency of top- $k$  queries on small and medium graphs.

## 3.5 Comparison of State-of-the-art Algorithms

*H-go SRJ* and *FLP & Opt-LP* rely on the recursive SimRank definition while *UISim* and *KSimJoin* adopt the random walk interpretation. To compute the meeting probabilities, *KSimJoin* uses breadth-first search, *FLP* & *Opt-LP* depend on local push algorithms, while *UISim* is considered as a tradeoff between them. Moreover, we derive a few conclusions. First, their theoretical complexities are not comparable. Particularly, [44, 46] introduce parameters defined by the algorithms. Second, existing experimental study is quite limited, including lack of evaluation on large graphs [20, 38, 46] and for estimation error [20, 44]. Third, although some algorithms achieve absolute error guarantee or effectively prune a fraction of unpromising nodes, none of them provides approximation guarantees for both types of similarity joins. For example, as shown by our experiments, threshold-based similarity joins with parameter  $\theta$  cannot be directly answered by estimating node pairs with  $O(\theta)$  error.

## 4 OUR IMPROVED ALGORITHMS

### 4.1 Algorithm Framework

We adopt a simple yet effective idea to answer threshold-based and top- $k$  queries with approximation bound. It takes an all-pair algorithm  $\mathcal{AP}$  as subprocedure, and progressively decreases the input error parameter of  $\mathcal{AP}$  until the bound is satisfied. Though similar idea has been used to answer single-source top- $k$  queries for PPR [36] and SimRank [22], we present specific stopping conditions for our problems.

**4.1.1 APTres for Threshold Queries.** Recall that threshold-based similarity join finds (at least  $\rho$  fraction of) all node pairs with SimRank similarity above  $\theta$ . Given an all-pair algorithm that guarantees an additive error of  $\varepsilon$  (with high probability), i.e.,  $|\hat{s}(u, v) -$

$s(u, v) \leq \varepsilon$  for each node pair  $(u, v)$ , if the estimated SimRank value  $\hat{s}(u, v)$  satisfies that  $\hat{s}(u, v) \geq \theta + \varepsilon$ , then we have  $s(u, v) \geq \theta$ , that is,  $(u, v)$  is included in the result set. Similarly, for  $(u', v')$  with  $\hat{s}(u', v') < \theta - \varepsilon$ , we should exclude it from the result set. Intuitively,  $\varepsilon$  should be small enough to meet the approximation bound  $\rho$ . We iteratively decrease  $\varepsilon$  by half until the condition is satisfied.

**Algorithm Description.** The framework for threshold-based queries, denoted as *APThres*, is illustrated in Algorithm 1. It takes graph  $G$ , threshold  $\theta$ , approximation bound  $\rho$ , and an all-pair algorithm  $\mathcal{AP}$  with absolute error guarantee as input. The algorithm initializes the result set  $\mathcal{R}(\theta)$  and the candidate set  $\mathcal{C}$  as empty set, and set  $\varepsilon_1 = O(\theta)$  (e.g.,  $\frac{\theta}{2}$ ) for the first iteration (Line 1). Then, at  $i$ -th iteration, *APThres* iteratively invokes  $\mathcal{AP}$  with error parameter  $\varepsilon_i$  (Lines 2-3). Particularly, we use  $\hat{S}_i = \{\hat{s}_i(u, v) | u, v \in V, u \neq v\}$  to denote the all-pair SimRank estimations, where the subscript  $i$  specifies the iteration number. In practice, we do not store node pairs with zero estimated values. We use  $\hat{S}_i$  to update  $\mathcal{R}(\theta)$  (Line 4) and  $\mathcal{C}$  (Line 5). The algorithm terminates if  $\frac{|\mathcal{R}(\theta)|}{|\mathcal{R}(\theta)|+|\mathcal{C}|} \geq \rho$  holds, which implies that  $\mathcal{R}(\theta)$  contains at least  $\rho$  fraction of the result (Lines 6-7). Otherwise, we set the error parameter of the next iteration as the half of its current value (Line 8). Note that if our stopping condition is never satisfied, the algorithm will stop when the SimRank estimation of every node pair has no more than  $\varepsilon_{min}$  error. In both cases,  $\mathcal{R}(\theta)$  is returned (Line 9).

---

#### Algorithm 1: The *APThres* framework

---

**Input:** Directed graph  $G = (V, E)$ , threshold  $\theta$ , approximation bound  $\rho$ , all-pair algorithm  $\mathcal{AP}$

**Output:**  $\mathcal{R}(\theta)$ , which contains at least  $\rho$  fraction of all node pairs with similarities above  $\theta$

```

1  $\mathcal{R}(\theta) \leftarrow \emptyset, \mathcal{C} \leftarrow \emptyset, \varepsilon_1 \leftarrow O(\theta);$ 
2 for  $i = 1$  to  $\log_2 \frac{1}{\varepsilon_{min}}$  do
3    $\hat{S}_i \leftarrow \mathcal{AP}(G, \varepsilon_i);$ 
4    $\mathcal{R}(\theta) \leftarrow \mathcal{R}(\theta) \cup \{(u, v), \hat{s}_i(u, v) \in \hat{S}_i | \hat{s}_i(u, v) \geq \theta + \varepsilon\};$ 
5    $\mathcal{C} \leftarrow (\mathcal{C} \cup \{(u, v), \hat{s}_i(u, v) \in \hat{S}_i | \hat{s}_i(u, v) \in [\theta - \varepsilon, \theta + \varepsilon]\}) \setminus \mathcal{R}(\theta);$ 
6   if  $\frac{|\mathcal{R}(\theta)|}{|\mathcal{R}(\theta)|+|\mathcal{C}|} \geq \rho$  then
7     break;
8    $\varepsilon_{i+1} \leftarrow \varepsilon_i / 2;$ 
9 return  $\mathcal{R}(\theta);$ 
```

---

**Correctness.** The following lemma proves that *APThres* has approximation bound for threshold-based queries.

LEMMA 5. *APThres* with an all-pair algorithm  $\mathcal{AP}$  of absolute error guarantee holds approximation bound  $\rho$  for threshold-based similarity joins.

4.1.2 *APTop-k* for Top-k Queries. We provide a framework *APTop-k* for top-k similarity joins with approximation bound  $\rho$ , which is analogous to *APThres* but with different stopping condition.

**Algorithm Description.** We demonstrate the pseudocode in Algorithm 2. The algorithm sets the initial value of the error parameter as  $O(c)$  (Line 1), and iteratively invokes the all-pair algorithm (Lines 2-3). At  $i$ -th iteration, the estimation values  $\hat{S}_i$  are sorted in descending

---

#### Algorithm 2: The *APTop-k* framework

---

**Input:** Directed graph  $G = (V, E)$ , parameter  $k$ , approximation bound  $\rho$ , all-pair algorithm  $\mathcal{AP}$

**Output:**  $\mathcal{R}(k)$ , which contains at least  $\rho$  fraction of node pairs with top- $k$  largest SimRank values

```

1  $\varepsilon_1 \leftarrow O(c);$ 
2 for  $i = 1$  to  $\log_2 \frac{1}{\varepsilon_{min}}$  do
3    $\hat{S}_i \leftarrow \mathcal{AP}(G, \varepsilon_i);$ 
4    $\{\hat{s}_{i,1}, \hat{s}_{i,2}, \dots, \hat{s}_{i,k}, \dots, \hat{s}_{i,n(n-1)}\} \leftarrow \text{sort}(\hat{S}_i)$  (in descending order);
5   if  $\hat{s}_{i,\lceil \rho k \rceil} - \varepsilon_i \geq \hat{s}_{i,k+1} + \varepsilon_i$  then
6     break;
7    $\varepsilon_{i+1} \leftarrow \varepsilon_i / 2;$ 
8 Let  $\mathcal{R}(k)$  be the node pairs with top- $k$  largest estimations;
9 return  $\mathcal{R}(k);$ 
```

---

order, denoted as  $\{\hat{s}_{i,1}, \hat{s}_{i,2}, \dots, \hat{s}_{i,k}, \dots, \hat{s}_{i,n(n-1)}\}$  (Line 4). We set the stopping conditions as  $\hat{s}_{i,\lceil \rho k \rceil} - \varepsilon_i \geq \hat{s}_{i,k+1} + \varepsilon_i$  (Lines 5-6), which says that the lower bound of the  $\lceil \rho k \rceil$ -th largest estimation is no smaller than the upper bound of the  $(k+1)$ -th one. The algorithm decreases the error parameter by half if it does not terminate this round (Line 7). Finally, *APTop-k* returns the set of node pairs with the top- $k$  largest estimations (Lines 8-9).

**Correctness.** We have the following lemma for *APTop-k*.

LEMMA 6. *APTop-k* with an all-pair algorithm  $\mathcal{AP}$  of absolute error guarantee holds approximation bound  $\rho$  for top- $k$  queries.

**Remark.** One might speculate that iteratively invocation of the all-pair algorithm  $\mathcal{AP}$  in *APThres* and *APTop-k* is costly and can be improved by the following strategy. When the size of the candidate set  $\mathcal{C}$  is below some threshold (note that we can retain  $\mathcal{C}$  for *APTop-k* similar to *APThres*), or the estimated cost of handling  $\mathcal{C}$  is less than invoking  $\mathcal{AP}$ , we switch to estimate the SimRank values of node pairs in  $\mathcal{C}$  independently. Actually, we have tried this strategy but always find it having inferior query efficiency. We anticipate that it may be useful for partial-pair queries (especially when the query node set is small), however, a fully investigation of this setting is beyond the scope of this paper.

## 4.2 Randomized Reverse Local Push ( $R^2LP$ )

We present an all-pair algorithm with better theoretical complexity, which is inspired by the advantage of *Randomized Backward Search (RBS)* [34] over *Backward Push* for PPR. To guarantee the practical efficiency, we devise a pruning strategy with absolute error guarantee, which dramatically enhances query time.

**Algorithm Description.** We refer to the algorithm as *Randomized Reverse Local Push ( $R^2LP$ )* (see Algorithm 3), which naturally extends the *RBS* algorithm [34] for PPR to all-pair SimRank queries. Specifically, the SimRank value  $s(u, v)$  is split into  $s^{(l)}(u, v)$  for  $l \in [0, \infty)$ , where  $s^{(l)}(u, v)$  is the probability of two  $\sqrt{c}$ -walks  $W_{\sqrt{c}}(u)$  and  $W_{\sqrt{c}}(v)$  first meet after exact  $l$  steps. It is obvious that  $s(u, v) = \sum_{l=0}^{\infty} s^{(l)}(u, v)$ .  $R^2LP$  estimates  $s^{(l)}(u, v)$  for all node pairs in ascending order of  $l$ . To guarantee  $\varepsilon$  additive error, it is sufficient to consider  $l \leq L = O(\log \frac{1}{\varepsilon})$  (Line 1). We initialize

the algorithm in Lines 2-3, where  $s^{(l)}(u, v)$  corresponds to the reserves of PPR [24, 34], and we omit the residues because it holds that  $r^{(l)}(u, v) = s^{(l)}(u, v)/(1 - c)$ . Then, for  $l = 0$  to  $L - 1$ , we compute the estimates  $\hat{s}^{(l+1)}(u, v)$  based on  $s^{(l)}(u, v)$  (Lines 4-10). For each  $s^{(l)}(u, v)$  at level  $l$ , we check if its value is large enough (Line 5). For now, we set  $f(\varepsilon) = 0$  to be the same as *RBS*. Similarly, our algorithm checks each out-neighbor  $(u', v')$  of  $(u, v)$  in the SimRank graph  $G^s$  (Lines 6&9). According to the values of  $s^{(l)}(u, v), d_{in}(u'), d_{in}(v')$ , and  $\varepsilon$ , the algorithm either pushes deterministically (Line 7) or randomly (Lines 8&10). The returned value  $\hat{s}(u, v)$  summarizes the estimations of every step  $l$  (Lines 11-12).

---

**Algorithm 3:** Randomized Reverse Local Push ( $R^2LP$ )

---

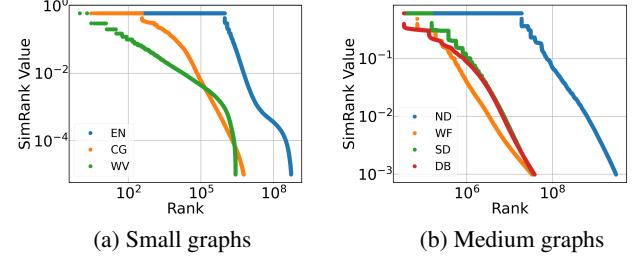
**Input:** Directed graph  $G = (V, E)$  with sorted adjacency lists, additive error parameter  $\varepsilon$   
**Output:**  $\hat{s}(u, v)$  for all  $u, v \in V$

- 1  $L \leftarrow \log_c(1 - c)\varepsilon;$
- 2  $\hat{s}^{(l)}(u, v) \leftarrow 0$  for all  $l \in [0, L]$  and  $u, v \in V$ ;
- 3  $\hat{s}^{(0)}(v, v) \leftarrow 1$  for all  $v \in V$ ;
- 4 **for**  $l = 0$  to  $L - 1$  **do**
- 5     **for each**  $(u, v)$  with  $\hat{s}^{(l)}(u, v) > f(\varepsilon)$  **do**
- 6         **for each**  $u' \in O(u), v' \in O(v), u' \neq v'$  and  
 $\sqrt{d_{in}(u')d_{in}(v')} \leq \frac{c\hat{s}^{(l)}(u, v)}{(1-c)\varepsilon}$  **do**
- 7              $\hat{s}^{(l+1)}(u', v') \leftarrow \hat{s}^{(l+1)}(u', v') + \frac{c\hat{s}^{(l)}(u, v)}{d_{in}(u')d_{in}(v')}$ ;
- 8          $r \leftarrow \text{rand}(0, 1)$ ;
- 9         **for each**  $u' \in O(u), v' \in O(v), u' \neq v'$  and  
 $\frac{c\hat{s}^{(l)}(u, v)}{(1-c)\varepsilon} < \sqrt{d_{in}(u')d_{in}(v')} \leq \frac{c\hat{s}^{(l)}(u, v)}{(1-c)\varepsilon r}$  **do**
- 10              $\hat{s}^{(l+1)}(u', v') \leftarrow \hat{s}^{(l+1)}(u', v') + \frac{(1-c)\varepsilon}{\sqrt{d_{in}(u')d_{in}(v')}};$
- 11      $\hat{s}(u, v) \leftarrow \sum_{l=0}^L \hat{s}^{(l)}(u, v)$ ;
- 12 **return** non-zero  $\hat{s}(u, v)$ s for  $u, v \in V, u \neq v$ ;

---

We discuss two implementation details. First, we do not need to materialize  $G^s$ . To get all out-neighbors of  $(u, v)$  in  $G^s$ , it is sufficient to traverse the out-neighbors of  $u$  and  $v$  in  $G$  simultaneously. Second, to guarantee the efficiency, the randomized push strategy [34] needs the out-neighbors of  $(u, v)$  to be sorted in *ascending* order of *in-degree*. This also can be achieved by the corresponding operations on  $G$ . Actually, we first preprocess the adjacency lists of  $G$  before invoking the algorithm. One may note that this strategy is not applicable to dynamic graphs. In this paper, we focus on the unsolved static problem setting and leave the problem on dynamic graphs for future work.

**Pruning Strategy.** The practical efficiency of  $R^2LP$  still needs to be improved to handle large graphs as we empirically observe many node pairs with negligible similarities during the push procedure. Recall that  $R^2LP$  saves computational cost by adopting the randomized push strategy. We additionally propose a deterministic pruning strategy to directly rule out estimations with small values at each step  $l$ . To be precise, we set  $f(\varepsilon) = \tilde{O}(\varepsilon)$  (i.e., we omit the logarithmic factor of  $\varepsilon$ ) in Line 5 of Algorithm 3.



**Figure 2: Distribution of all-pair SimRank values.**

**Correctness.** The following lemma shows that our  $R^2LP$  algorithm guarantees additive error  $\varepsilon$  with high probability. Particularly, the above pruning strategy does not affect the error guarantee.

**LEMMA 7.** *By conducting the randomized push procedure of  $R^2LP$  for  $O(\log \frac{n}{\delta})$  times and apply the Median-of-Mean trick [6], we have  $|\hat{s}(u, v) - s(u, v)| \leq \varepsilon$  with at least  $1 - \delta$  probability.*

**Complexity Analysis.** We use the following theorem to guarantee the efficiency of  $R^2LP$ .

**THEOREM 1.** *The expected time complexity of  $R^2LP$  is bounded by  $\tilde{O}\left(\frac{\sum_{u \in V} \sum_{v \in V} \sqrt{d_{in}(u)d_{in}(v)}s(u, v)}{\varepsilon}\right)$ .*

## 5 COMPUTATIONAL COMPLEXITY ANALYSIS

We analyze the computational complexity of threshold-based and top- $k$  similarity joins with approximation bound  $\rho$ . By empirical study we observe the skewed SimRank distribution of real-world graphs. Meanwhile, existing literature [22, 35] has also demonstrated that single-source SimRank distribution approximately follows power law. This motivates us to make reasonable assumption of the distribution and then derive non-trivial complexity bounds that match the practical performance of algorithms.

### 5.1 Modeling SimRank Distribution

We demonstrate in Figure 2 the all-pair SimRank distributions of three small graphs and four medium graphs used in our experiments (Section 6). We omit the base cases with  $s(v, v) = 1$ . Generally the distributions show different degrees of skewness with a fraction of the top ranking node pairs having a similarity of  $c$ . Inspired by [2, 5], we propose the following assumption of SimRank distribution.

**DEFINITION 5 (POWER LAW BOUNDED DISTRIBUTION OF SIMRANK).** *Assume we have  $nnz$  node pairs with non-zero SimRank values for a given graph  $G$ . Let  $s_j$  denote the  $j$ -th largest SimRank value, for  $j \in [1, nnz]$ . We say the SimRank distribution is power law bounded (PLB) if for any  $x = 0, \dots, \lfloor \log_2 nnz \rfloor$ , there exist parameters  $\beta > 0$  and  $b_2 > b_1 > 0$  such that*

$$\sum_{j=2^x}^{2^{x+1}-1} b_1 \cdot r(j)^{-\beta} \leq \sum_{j=2^x}^{2^{x+1}-1} s_j \leq \sum_{j=2^x}^{2^{x+1}-1} b_2 \cdot r(j)^{-\beta}, \quad (6)$$

where  $r(j) = \max(1, j-t)$ , and we set the shift  $t = |\{(u, v) | s(u, v) = c, u, v \in G\}|$ . Specifically,  $\beta$  is referred to as the PLB exponent of SimRank distribution, while  $b_1$  and  $b_2$  are small numbers that can be taken as constants.

Please refer to the full version of our paper [1] for empirical validation of the assumption.

## 5.2 Complexity of Threshold-based Queries

With the above assumption, we derive the *upper bound* of the complexity of threshold-based queries. We do this by bounding the complexity of our *APThres* (as well as *APTop-k*) framework which integrates an all-pair algorithm with absolute error guarantee.

**LEMMA 8.** *Given a graph  $G$  and an all-pair algorithm  $\mathcal{AP}$  that guarantees absolute error, assume that *APThres* (or *APTop-k*) terminates when the error parameter reaches  $\varepsilon_t$ , and we have  $\text{cost}(\mathcal{AP}(G, \varepsilon_{i+1})) \geq p \cdot \text{cost}(\mathcal{AP}(G, \varepsilon_i))$  for  $i \in [1, t]$  with  $p > 1$ . The time complexity of the algorithm is then bounded by  $O(\text{cost}(\mathcal{AP}(G, \varepsilon_t)))$ .*

**Remark.** We can choose  $\mathcal{AP}$  as *FLP*, *Opt-LP*, *KSimJoin*, or *R<sup>2</sup>LP*. The lemma cannot be applied to *UISim* because it has no absolute error guarantee, nor *H-go SRJ* because the complexity depends on an algorithm-specific parameter  $h$ .

Next, we bound the error parameter  $\varepsilon_t$  for threshold-based queries with the help of Definition 5.

**LEMMA 9.** *We are given the problem input as in Lemma 8. Moreover, the graph  $G$  has PLB SimRank distribution. For threshold-based queries with approximation bound  $\rho$ , we have*

$$\varepsilon_t = O\left(\max\left(\frac{(b_1 - \rho^\beta b_2)\theta}{b_1 + \rho^\beta b_2}, \varepsilon_{\min}\right)\right). \quad (7)$$

The above lemma implies that given a graph  $G$  and an approximation bound  $\rho$ , the approximation guarantee can be non-trivially satisfied only if  $b_1 - \rho^\beta b_2 > 0$ , which is determined by both graph properties and the approximation parameter. The following theorem bounds the computational complexity of threshold-based queries.

**THEOREM 2.** *Given a graph  $G$  with PLB SimRank distribution, an approximation bound  $\rho$  with  $b_1 - \rho^\beta b_2 > 0$ , and an all-pair algorithm  $\mathcal{AP}$  that guarantees absolute error, let  $\text{cost}(\mathcal{AP}(G, \varepsilon))$  denote the complexity of  $\mathcal{AP}$  given error parameter  $\varepsilon$ . The *APThres* framework with  $\mathcal{AP}$  answers  $\rho$ -approximation threshold-based queries in  $O(\text{cost}(\mathcal{AP}(G, \frac{(b_1 - \rho^\beta b_2)\theta}{b_1 + \rho^\beta b_2})))$  time.*

## 5.3 Complexity of Top-k Queries

For top- $k$  queries, the following theorem is provided, with the intuition that the complexity is mainly determined by the gap between the  $\lceil \rho k \rceil$ -th and the  $(k+1)$ -th largest values.

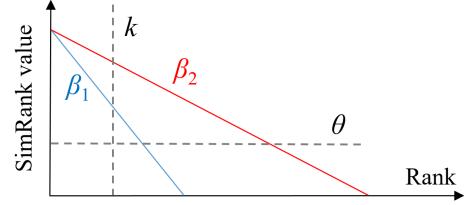
**THEOREM 3.** *We are given the problem input as in Lemma 9. For top- $k$  queries with approximation bound  $\rho$ , with the PLB assumption of SimRank distribution, let  $x_1 = 2^{\lceil \log_2 \lceil \rho k \rceil \rceil} - 1$  and  $x_2 = 2^{\lfloor \log_2 (k+1) \rfloor}$ . We have*

$$\varepsilon'_t = O\left(\max(b_1(x_1 - t)^{-\beta} - b_2(x_2 - t)^{-\beta}, \varepsilon_{\min})\right). \quad (8)$$

Furthermore, with the power law assumption of all-pair SimRank distribution, i.e.,  $s_j = \Theta((j-t)^{-\beta})$ , we have

$$\varepsilon'_t = O\left(\max((b_1(\lceil \rho k \rceil - t)^{-\beta} - b_2(k+1 - t)^{-\beta}), \varepsilon_{\min})\right). \quad (9)$$

The *APTop-k* framework with  $\mathcal{AP}$  answers  $\rho$ -approximation top- $k$  queries in  $O(\text{cost}(\mathcal{AP}(G, \varepsilon'_t)))$  time.



**Figure 3: Illustration of the complexity of threshold-based and top- $k$  similarity joins.**

Lastly, we present a qualitative analysis of the complexities w.r.t. the input parameters for threshold-based and top- $k$  queries. For both queries, the complexity upper bound increases as we set a higher approximation bound  $\rho$ . For threshold-based (resp. top- $k$ ) query, generally the problem tends to be harder as we decrease  $\theta$  (resp. increase  $k$ ). However, the PLB exponent  $\beta$  has different impacts on two types of queries. An illustrative demonstration is referred to Figure 3. For threshold-based queries, a smaller  $\beta$  (e.g.  $\beta_2$ ) leads to a smaller  $\varepsilon_t$  in Lemma 9, and the complexity upper bound increases. Intuitively, more node pairs are included in the answer, while given an error parameter  $\varepsilon$ , the size of node pairs to be carefully handled (e.g., with similarity values between  $\theta - \varepsilon$  and  $\theta + \varepsilon$ ) also increases. In contrast, as Theorem 3 shows, a larger  $\beta$  (e.g.  $\beta_1$ ) makes a smaller gap between the  $\lceil \rho k \rceil$ -th and the  $(k+1)$ -th SimRank values, so the algorithm needs more effort to separate them. This is in consistent with Figure 3, where the  $k$ -th SimRank value becomes smaller (also note that the figure is in log-log scale).

## 6 EXPERIMENTAL EVALUATION

### 6.1 Experimental Setup

**Datasets.** We use three small graphs (CG, WV, and EN), four medium graphs (SD, DB, WF, and ND), and three large graphs (CP, LJ, and WZ) for evaluation. We do not include larger datasets due to excessive memory consumption. All graphs are obtained from [3, 4, 15, 17] with their basic statistics listed in Table 3. In particular,  $\bar{d}$  is the average degree and  $\beta$  is the fitted PLB exponent of all-pair SimRank distribution[1].

**Table 3: Datasets and their statistics.**

Dataset	Type	$n$	$m$	$\bar{d}$	$\beta$
ca-GrQc (CG)	U	5.2K	14.5K	2.77	1.063
Wiki-Vote (WV)	D	7.1K	103.7K	14.57	0.495
email-Enron (EN)	U	36.7K	183.8K	5.01	1.303
Slashdot0922 (SD)	D	82.2K	948.5K	11.54	1.098
DBLP (DB)	U	317.1K	1.05M	3.31	0.866
Wikilinks-fy (WF)	D	65.6K	1.07M	16.35	1.098
Notre Dame (ND)	D	325.7K	1.5M	4.6	0.964
cit-Patents (CP)	D	3.77M	16.52M	4.38	-
LiveJournal (LJ)	D	4.85M	68.99M	14.23	-
Wikilinks-zh (WZ)	D	1.79M	72.61M	40.65	-

**Ground truths.** For small and medium graphs, we compute the all-pair ground truth using *Power Method* [12]. We set the number of iterations to 30 and achieve an absolute error about  $10^{-7}$ . For medium graphs, we only store SimRank values above  $10^{-3}$  due to

excessive disk consumption. For large graphs, since *Power Method* is computationally intractable, we randomly choose 100 nodes in each graph, take each node as query input and use *ExactSim* [35] to compute the single-source ground truth with an absolute error of  $10^{-7}$ . This partial-pair ground truth is then used for evaluation.

**Baselines.** We include *FLP* & *Opt-LP* [38], *UISim* [46], and *K-SimJoin* [20] as our baselines. We do not include [25, 27, 28, 31] as they are outperformed by the baselines according to their empirical evaluation. The implementation of [38, 46] is obtained from the authors. We implement other baselines and our improved algorithms in C++. Specifically, we do not consider *H-go SRJ* [44] for efficiency issues. We find that it incurs tremendous time and space cost for index construction and query execution even on small graphs, with significant inferior performance compared to other baselines.

We conduct all experiments on an Ubuntu server with an Hygon C86 7151 processor and 1TB memory, repeat each experimental settings for ten times and report the average measures.

## 6.2 Evaluation of Additive Error

We first evaluate the additive error of each algorithm in answering general all-pair queries, since they are the building blocks of threshold-based and top- $k$  algorithms.

**Experimental settings.** For algorithms with absolute error guarantee, including *FLP*, *Opt-LP*, *R<sup>2</sup>LP* and *KSimJoin*, we vary the error parameter  $\varepsilon$ . For small graphs (CG, WV, and EN) and medium graphs (SD, DB, WF, and ND), we set  $\varepsilon = \{0.05, 0.01, 0.005, 0.001\}$ . For large graphs (CP, LJ, and WZ), we vary  $\varepsilon = \{0.1, 0.05, 0.01\}$ . In particular, we slightly modify *KSimJoin* to answer all-pair queries by excluding the upper bound computation and the iterative pruning framework. For *UISim*, we fix the number of iterations  $\eta = 2$  and vary the number of hub nodes  $|H| = \{0.2n, 0.5n\}$ . As for its implementation, the authors adopt another parameter named stopping reachability (*stopRea*), where the graph expansion stops either when a hub node is encountered or the walking probability falls below *stopRea*. As we empirically find that *stopRea* has more impact on query efficiency and accuracy, we vary it in  $\{0.005, 0.001\}$  for small graphs and  $\{0.05, 0.01\}$  for medium graphs, and set it as 0.01 for large graphs. We fix  $L = 2$  for *KSimJoin*.

**Evaluation metric.** We use *AvgErr* for evaluation, which is defined as  $\text{AvgErr} = \frac{\sum_{(u,v) \in S(\theta)} |s^*(u,v) - \hat{s}(u,v)|}{|S(\theta)|}$ , where  $s^*(u,v)$  and  $\hat{s}(u,v)$  are the ground truth and estimated value respectively, and  $S(\theta) = \{(u,v) | s^*(u,v) \geq \theta, u, v \in V, u \neq v\}$ . We set  $\theta = 10^{-3}$  for the following reasons. Applications usually focus on non-negligible SimRank values, which only make up a small portion of the SimRank matrix. With most other similarity values approaching zero, it is not wise to take them into consideration, as an algorithm that simply returns zero for each node pair achieves very small average error. Moreover, we cannot afford to compute and store the ground truths for node pairs with very small similarities for medium and large graphs.

The results are demonstrated in Figure 4. Throughout the paper, all reported query times are in seconds. For *UISim*, the unfilled markers correspond to the default values of *stopRea* (0.001 and 0.01), while the filled markers refer to the increase of *stopRea* by 5 times (i.e., 0.005 and 0.05) on small and medium graphs. Surprisingly, although *UISim* does not guarantee the absolute error, it outperforms

*FLP* and *Opt-LP* on 6 of 10 graphs, and achieves comparable performance on another graph (DB). However, we note that varying the number of hubs has quite limited influence on both query efficiency or accuracy, that is, the two unfilled markers (as well as the two filled markers) are very close, especially for WV, EN, SD, WF, and LJ. By contrast, the filled and unfilled markers show very different performance. The reason is that most of graph expansions terminate because of stopping reachability rather than encountering a hub node. As we increase *stopRea*, the algorithm returns SimRank estimations of larger additive errors with less response time.

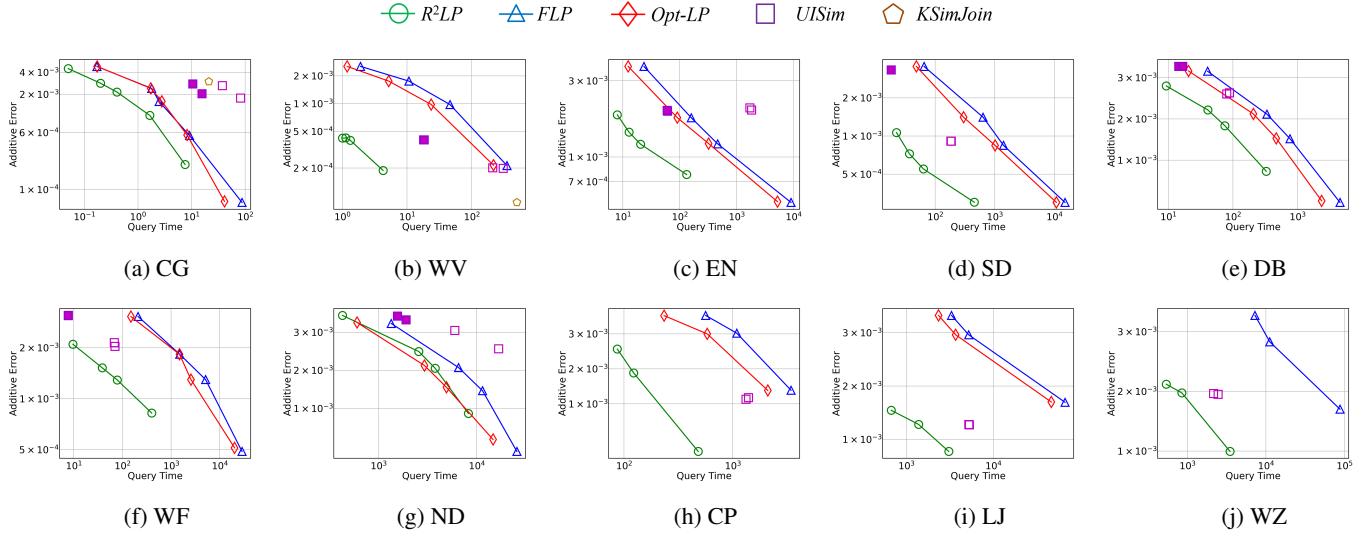
For algorithms with absolute error guarantees, *Opt-LP* outperforms *FLP* on most graphs thanks to the optimizations, except that it fails on the largest graph WZ (not terminated for over 10 hours). However, they are usually slower than *UISim*, because backward push on the SimRank graph is very time consuming. *KSimJoin* only manages to answer all-pair queries on the two smallest graphs and has similar performance as *UISim*, while the tremendous memory consumption makes it infeasible on larger graphs. For example, on WV and with  $L = 2$ , the *NP* structure for all nodes has about 1.8M items. On the contrary, with the randomized push strategy, *R<sup>2</sup>LP* achieves significantly better performance compared to all other baselines. On most graphs, it is faster than *FLP*, *Opt-LP*, and *KSimJoin* by up to an order of magnitude with similar additive errors. The only exception is ND, where *R<sup>2</sup>LP* and *Opt-LP* show comparable results. Recall that ND has numerous node pairs with very large similarity values (see Figure 2(b)), thus the randomized part of *R<sup>2</sup>LP* is seldomly invoked. In general, to achieve good results, we need both theoretical guarantee of the algorithm as well as strategies to ensure practical efficiency.

## 6.3 Evaluation of Threshold-based Queries

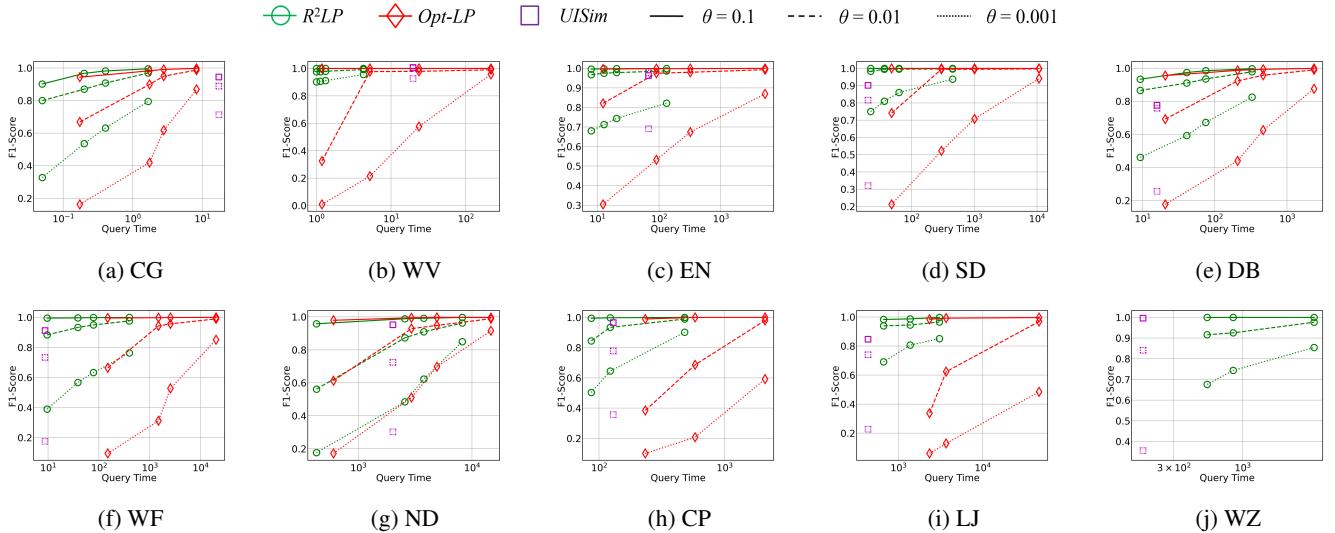
**Experimental settings and metrics.** For small and medium graphs, we directly evaluate the results of all-pair threshold-based queries. For three large graphs, as we only have the ground truth of 100 single-source queries, we ignore the algorithm outputs that are not contained in this ground truth. We adopt *Precision*, *Recall* and *F1-score* for evaluation. Given a threshold-based query with parameter  $\theta \in (0, c]$ , let  $R^*(\theta)$  denote the set of node pairs with ground truth SimRank values above or equal to  $\theta$  (with up to  $\varepsilon_{min} = 10^{-7}$  error), and we assume  $|R^*(\theta)| > 0$ . Let  $R_{\mathcal{A}}(\theta)$  denote the returned set by algorithm  $\mathcal{A}$ . We have  $\text{Precision} = |R_{\mathcal{A}}(\theta) \cap R^*(\theta)| / |R_{\mathcal{A}}(\theta)|$ ,  $\text{Recall} = |R_{\mathcal{A}}(\theta) \cap R^*(\theta)| / |R^*(\theta)|$ , and  $\text{F1-score} = 2 / (\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}})$ .

In our first set of experiments, we vary  $\theta = \{0.1, 0.01, 0.001\}$  and evaluate the all-pair algorithms used in Section 6.2. We only include three best algorithms *R<sup>2</sup>LP*, *Opt-LP*, and *UISim*. Note that none of these algorithms provides approximation bound for threshold-based queries, but we can still compare their empirical performance. The results are demonstrated in Figure 5<sup>3</sup> with the following conclusions derived. First, for each algorithm and each of its parameter configuration, the *F1-scores* vary significantly for different graphs and with different  $\theta$ . It is not recommended to answer the threshold-based query with a general all-pair algorithm by simply setting  $\varepsilon = O(\theta)$ . For example, on CG, EN, DB, and WF, the *F1-scores* with  $\varepsilon = \theta = 0.001$  are below 0.9. Second, *R<sup>2</sup>LP* generally achieves the best tradeoff between efficiency and accuracy

<sup>3</sup>For the results of *Precision* and *Recall*, please refer to the full version of our paper [1].



**Figure 4: Additive error of SimRank estimation.**



**Figure 5: F1-score for threshold-based queries. No approximation bound is provided.**

except on ND. Third, for *UISim*, we choose the best one among the four configurations used in Section 6.2. On a majority of graphs, it surpasses *Opt-LP* mainly due to the fast query time. However, as we decrease *stopRea* and increase  $\eta$ , the algorithm fails to compute an output in reasonable time (say 3 hours), which implies the hub-based strategy is less flexible than the probability-guided ones.

Next, we apply the *APThres* framework and integrate it with  $R^2LP$  and *Opt-LP*, the best algorithms with absolute error guarantee. By fixing  $\theta = 0.01$ , we test their actual performance with different approximation bounds, which is shown in Table 4. Consistent with our theoretical guarantee, both *APThres+Opt-LP* and *APThres+R<sup>2</sup>LP* hold the approximation bound on all graphs and all settings of  $\rho$ . *APThres+R<sup>2</sup>LP* achieves similar *F1-score* with much less query time except on ND, of which the reasons are explained in Section 6.2. Also note that the advantage of *Opt-LP* over  $R^2LP$  with an iterative

framework lies in that it can compute from the ‘‘borders’’ of the previous iteration rather than from scratch.

However, compared to Figure 5, both algorithms incur significantly more time, and fail on LJ for different reasons<sup>4</sup> while *APThres+Opt-LP* also fails on WZ. Recall that  $R^2LP$  is able to handle LJ with acceptable accuracy if we do not force an approximation bound. Meanwhile, the actual *F1-scores* shown in Table 4 far exceed the bound  $\rho$ . We believe this phenomenon is attributed to the pessimistic estimation of the absolute error, that is, for almost all node pairs, it turns out that  $|\hat{s}(u, v) - s(u, v)| \ll \varepsilon$ . We have tried other optimization strategies (see Section 4.1) and find that for now  $R^2LP$  is the most efficient algorithm to guarantee an absolute error. Hence, it

<sup>4</sup>We say an algorithm is time out if it takes over an hour in one iteration of *APThres*.

**Table 4: Threshold-based queries with  $\theta = 0.01$ . We vary  $(\rho_1, \rho_2) = (0.9, 0.99), (0.9, 0.95), (0.7, 0.9)$  for small, medium and large graphs respectively.**

	$Opt-LP(\rho_1)$		$Opt-LP(\rho_2)$		$R^2LP(\rho_1)$		$R^2LP(\rho_2)$	
	Time	F1	Time	F1	Time	F1	Time	F1
CG	5.55	0.991	11.99	0.997	2.85	0.978	8.28	0.993
WV	152.4	0.998	159.7	0.998	31.2	0.998	32.0	0.999
EN	2944	0.995	10625	0.999	201.8	0.989	1040	0.996
SD	9451	0.998	12011	0.998	611.6	0.997	1457	0.998
DB	1994	0.994	3480	0.997	501.6	0.985	1119	0.992
WF	13336	0.994	8938	0.994	649.9	0.984	1370	0.99
ND	10691	0.991	13129	0.996	6607	0.936	10269	0.958
CP	5487	0.999	8504	0.999	1903	0.999	4350	0.999
LJ	Time Out			Out of Memory (OOM)				
WZ	Time Out			3743	0.985	7411	0.99	

remains an open problem to answer threshold-based queries more efficiently and with theoretical approximation bound.

#### 6.4 Evaluation of Top- $k$ Queries

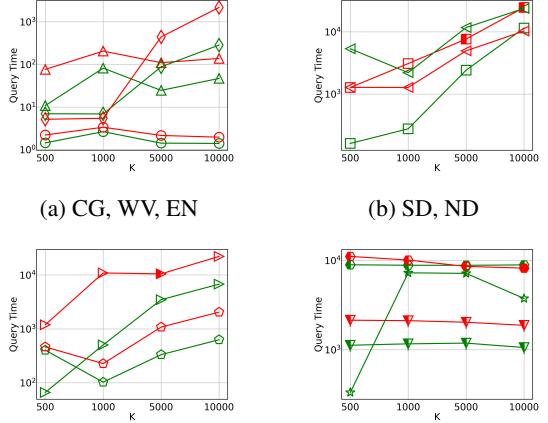
**Experimental settings and metrics.** For top- $k$  similarity joins,  $Precision@k$  is used for evaluation. Let  $R^*(k)$  and  $R_{\mathcal{A}}(k)$  denote the ground truth and the results of algorithm  $\mathcal{A}$ , respectively. We have  $Precision@k = |R_{\mathcal{A}}(k) \cap R^*(k)|/k$ . As shown in Figure 2, numerous node pairs have similarity value  $c$ , and the numbers of such node pairs vary significantly for different graphs. Instead of setting different  $k$  for our datasets, the following strategy is adopted. For medium and large graphs, we conduct *all-pair queries* but only evaluate the result accuracy on a pair of query node sets  $A$  and  $B$ , which is referred to as *partial-pair evaluation*. Specifically, we set  $A = B = V$  for CG and WV, while for EN, SD, WF, and ND we set  $|A| = |B| = 10^3$ . Since the all-pair ground truth of DB is very sparse, we set  $|A| = |B| = 10^4$ . For three large graphs, we fix  $A$  as the 100 query nodes in generating ground truth, and set  $|B| = 10^3$ . We randomly sample the query node sets and repeat the experiment for ten times. Note that we only use partial-pair evaluation to facilitate the comparison of different algorithms and do not consider the possible optimizations for partial-pair queries.

**Table 5: Precision@5000 and query time ( $\rho = 0.9$ ). The practical accuracy is underlined if it falls below  $\rho$ .**

	<i>UISim</i>		<i>Opt-LP</i>		<i>KSimJoin</i>		<i>R<sup>2</sup>LP</i>	
	Time	P@k	Time	P@k	Time	P@k	Time	P@k
CG	82.5	0.935	2.16	0.999	6.3	0.926	1.42	0.998
WV	308.2	0.992	110.3	0.999	405.9	0.994	24.6	0.998
EN	1189	0.993	443.7	0.997	OOM		88.0	0.996
SD	185.4	0.922	7682	0.996	OOM		2424	0.997
DB	84.0	0.777	1082	0.998	OOM		333.7	0.995
WF	56.0	0.776	10501	0.99	OOM		3488	0.995
ND	13448	0.959	4964	0.997	OOM		11749	0.993
CP	1416	0.95	2037	0.992	OOM		1190	0.994
LJ	5604	0.917	8586	0.964	OOM		8824	0.97
WZ	2375	0.904	Time Out		OOM		7171	1

We first compare the state-of-the-art algorithms by fixing  $k = 5000$  and  $\rho = 0.9$  with the results listed in Table 5. For *UISim*,



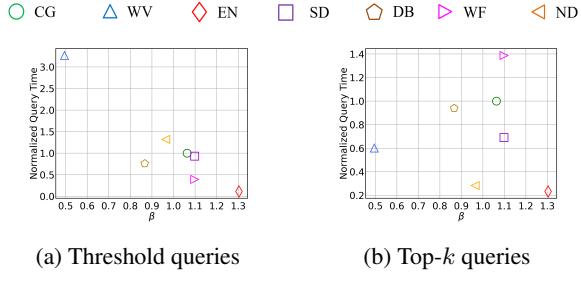


**Figure 6: Query time vs.  $k$ , with  $\rho = 0.9$ .**

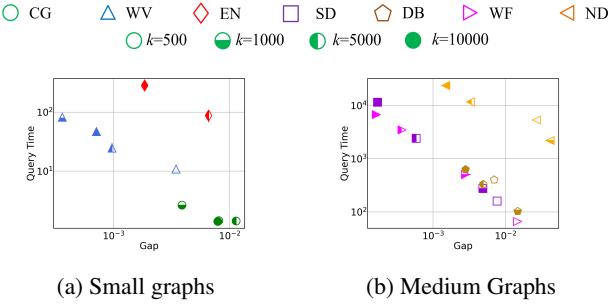
we fix  $\eta = 2$  and set  $stopRea$  as  $10^{-4}$  for CG and WV,  $10^{-3}$  for EN, and  $10^{-2}$  for all other datasets. For *KSimJoin*, it shares the same parameter settings as in Section 6.2. We integrate *Opt-LP* and *R<sup>2</sup>LP* with *APTop-k* to answer top- $k$  queries. Particularly, we terminate the algorithm early if it takes more than an hour in one iteration, and set it as time out if it costs over 30000 seconds. Although being fast, *UISim* mainly suffers from accuracy issues. On the other hand, *KSimJoin* has severe scalability problem on medium and large graphs in spite of its pruning strategy. For both *APTop-k+Opt-LP* and *APTop-k+R<sup>2</sup>LP*, the query accuracy is guaranteed. With theoretical bound, they also face the open problem as for the threshold-based queries. Nonetheless, *APTop-k+R<sup>2</sup>LP* is the fastest on 5 of 10 graphs while achieving much better practical accuracy than its counterparts.

Next, we evaluate the query efficiency w.r.t.  $k$  by setting  $k = \{500, 1000, 5000, 10000\}$ . Only *APTop-k+Opt-LP* and *APTop-k+R<sup>2</sup>LP* are considered because other methods fail to guarantee stable accuracy. As demonstrated in Figure 6, the red (resp. green) line denotes the result of *Opt-LP* (resp. *R<sup>2</sup>LP*), while different types of markers indicate whether the result has theoretical guarantee or even has inferior accuracy than the bound. Recall that if the algorithm terminates earlier due to excessive time cost, it is unable to guarantee the theoretical accuracy. In most cases, *APTop-k+R<sup>2</sup>LP* outperforms *APTop-k+Opt-LP* by a significant margin. On three large graphs CP, LJ, and WZ, both algorithms lose the theoretical guarantee due to early termination. Fortunately, the practical accuracy of *APTop-k+R<sup>2</sup>LP* is always above  $\rho$ , whereas *APTop-k+Opt-LP* has inferior accuracy on LJ with  $k = 10000$  and even fails on WZ.

Another interesting observation from Figure 6 is that the query time does not necessarily increase when we enlarge  $k$ . Generally, both algorithms demonstrate similar trends on all graphs. The explanation is given in the next subsection.



**Figure 7: Query time vs.  $\beta$ .**



**Figure 8: Query time vs.  $s_{\lceil \rho k \rceil} - s_{k+1}$  for top- $k$  queries.**

## 6.5 Empirical Study of Computational Complexity

We discuss the impact of approximation bound  $\rho$  and input parameters  $\theta$  and  $k$  on query efficiency, which can be concluded from the experiments in Section 6.3 & 6.4. For example, Figure 5 and Table 4 imply that a larger approximation bound  $\rho$  or a smaller parameter  $\theta$  leads to more time cost, while Figure 6 directly suggests the time complexity of top- $k$  queries generally increases as we enlarge  $k$ . For a deeper investigation of the algorithmic complexities, we further conduct two sets of experiments to validate our theoretical analysis w.r.t. graph properties. As  $R^2LP$  and  $Opt-LP$  demonstrate similar performance in the above experiments (except for  $R^2LP$  being faster), we only use  $R^2LP$  in the following analysis.

We conduct threshold queries with  $\theta = 0.01$  and top- $k$  queries with  $k = 1000$  following the settings in the above subsections. To illustrate the intrinsic problem complexities w.r.t.  $\beta$ , the PLB exponent of SimRank distribution, we first make sure the *practical accuracy* of  $R^2LP$  slightly surpasses  $\rho = 0.9$  by adjusting the error parameter  $\varepsilon$ . Next, since the graph size does have an important impact on query efficiency, say, the numerator of the computational complexity (see Theorem 1 for  $R^2LP$ ), we approximate it via computing  $\sum_{u \in V} \sum_{v \in V} \sqrt{d_{in}(u)d_{in}(v)}s(u, v)$  for each dataset using the all-pair ground truth. We then divide the actual query time by the approximate numerator and normalize the result by adjusting the value of CG to 1. Figure 7 plots the *normalized* query time w.r.t.  $\beta$ . The normalized query time only reflects the impact of  $\varepsilon$ , which is mainly affected by both  $\rho$  (fixed in our experiment) and graph properties, i.e.,  $\beta$ ,  $b_1$  and  $b_2$ . Generally speaking, the time is in negatively (resp. positively) correlated with  $\beta$  for threshold-based (resp. top- $k$ ) queries, which coincides with the derived computational complexity (Theorem 2 & 3). For top- $k$  queries, the data points are more

scattered because the impact of  $\beta$  is more complicated as shown by Theorem 3. Therefore, an additional analysis is conducted.

For top- $k$  queries, our theoretical analysis suggests that the time complexity heavily relies on the gap between the  $\lceil \rho k \rceil$ -th and the  $(k+1)$ -th largest SimRank values. Figure 8 directly shows query time w.r.t. the gap for small and medium graphs, where the all-pair ground truths are available. It turns out that query time nicely correlates with the gap on all datasets. Note that for a majority of studied graphs including CG, WV, DB and ND, the gap does not necessarily decrease as we increase  $k$ . This explains the phenomenon observed in Figure 6. It is worthy pointing out that the observed correlation also fits the sampling complexity of the top- $k$  identification problem [10, 14] in the multi-armed bandit setting, where the backward push procedure (such as  $R^2LP$  and  $Opt-LP$ ) is considered as generating an estimation for each arm (i.e., node pair) with  $\varepsilon$  confidence interval.

## 7 OTHER RELATED WORK

To accelerate the *Power Method*, early work [23, 43] improves the algorithmic complexity by pruning redundant or unnecessary computation. Another line of work [11, 18, 40] studies all-pair SimRank computation on dynamic graphs. We also note that several all-pair algorithms [9, 39–42] rely on a modified definition of SimRank, i.e.  $S = cP^TSP + (1 - c)I$ , which results in estimations without absolute error guarantee.

As for other types of SimRank queries, most of the existing work (e.g., [13, 16, 21, 22, 29, 30, 32, 35]) focuses on single-source settings. Given a query node  $u$ , the returned result can either be the estimated SimRank value w.r.t.  $u$  for each node  $v$  [13, 29, 30, 32, 35], the set of nodes with top- $k$  largest SimRank values (the *single-source top- $k$  query* [16, 21, 22]), or all nodes whose similarity w.r.t.  $u$  are above some threshold  $\theta$  (the *single-source thresholding query* [22]). Various algorithms have been proposed based on random walk sampling [13, 29], graph traversal [16, 21, 22, 32], and local push [30], while recent literature [21, 22, 30, 35] has significantly improved the query efficiency on large graphs. For example, SimPush [30] answers single-source queries on billion-edge graphs in seconds with additive error less than  $10^{-3}$ . Besides, a single-pair query takes two query nodes  $u$  and  $v$  as input and computes an estimation of  $s(u, v)$  with additive error up to  $\varepsilon$ . Existing methods such as the Monte Carlo (*MC*) algorithm [8] and its recent improvement [37] can answer single-pair queries within milliseconds for  $\varepsilon \geq 10^{-4}$ .

## 8 CONCLUSION

This paper presents detailed empirical studies, new algorithms with better efficiency and approximation bounds, and novel theoretical analysis for two types of all-pair SimRank queries, namely, threshold-based and top- $k$  similarity joins. We conduct a detailed comparison of state-of-the-art algorithms followed by the introduction of algorithm frameworks that theoretically guarantee the result accuracy for both queries, along with a more efficient all-pair algorithm inspired by the randomized local push. To gain a deep understanding of algorithm performance and the computational complexities of these problems, we propose theoretical analysis by leveraging the SimRank distribution and conduct extensive experimental studies. Our experiments show that the proposed algorithms significantly improve both query efficiency and accuracy, while our theoretical results nicely match the empirical analysis. For future work, we aim to study all-pair SimRank computation in dynamic settings.

## REFERENCES

- [1] [n.d.]. <https://github.com/xinghun0525/R2LP>.
- [2] Jonathan W Berry, Luke K Fostvedt, Daniel J Nordman, Cynthia A Phillips, C Seshadri, and Alyson G Wilson. 2014. Why do simple algorithms for triangle enumeration work in the real world?. In *Proceedings of the 5th conference on Innovations in theoretical computer science*. 225–234.
- [3] Paolo Boldi, Marco Rosa, Massimo Santini, and Sebastiano Vigna. 2011. Layered Label Propagation: A MultiResolution Coordinate-Free Ordering for Compressing Social Networks. In *Proceedings of the 20th international conference on World Wide Web*, Sadagopan Srinivasan, Krithi Ramamritham, Arun Kumar, M. P. Ravindra, Elisa Bertino, and Ravi Kumar (Eds.). ACM Press, 587–596.
- [4] Paolo Boldi and Sebastiano Vigna. 2004. The WebGraph Framework I: Compression Techniques. In *Proc. of the Thirteenth International World Wide Web Conference (WWW 2004)*. ACM Press, Manhattan, USA, 595–601.
- [5] Paweł Brach, Marek Cygan, Jakub Lkacki, and Piotr Sankowski. 2016. Algorithmic complexity of power law networks. In *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM, 1306–1325.
- [6] Moses Charikar, Kevin C. Chen, and Martin Farach-Colton. 2002. Finding Frequent Items in Data Streams. In *Automata, Languages and Programming, 29th International Colloquium, ICALP 2002, Malaga, Spain, July 8-13, 2002, Proceedings (Lecture Notes in Computer Science)*, Peter Widmayer, Francisco Triguero Ruiz, Rafael Morales Bueno, Matthew Hennessy, Stephan J. Eidenbenz, and Ricardo Conejo (Eds.), Vol. 2380. Springer, 693–703. [https://doi.org/10.1007/3-540-45465-9\\_59](https://doi.org/10.1007/3-540-45465-9_59)
- [7] Aaron Clauset, Cosma Rohilla Shalizi, and Mark EJ Newman. 2009. Power-law distributions in empirical data. *SIAM review* 51, 4 (2009), 661–703.
- [8] Dániel Fogaras and Balázs Rácz. 2005. Scaling link-based similarity search. In *Proceedings of the 14th international conference on World Wide Web*. 641–650.
- [9] Yasuhiro Fujiwara, Makoto Nakatsuji, Hiroaki Shiokawa, and Makoto Onizuka. 2013. Efficient search algorithm for SimRank. In *2013 IEEE 29th International Conference on Data Engineering (ICDE)*. IEEE, 589–600.
- [10] Victor Gabillon, Mohammad Ghavamzadeh, and Alessandro Lazaric. 2012. Best arm identification: A unified approach to fixed budget and fixed confidence. In *Advances in Neural Information Processing Systems*. 3212–3220.
- [11] Guoming He, Haijun Feng, Cuiping Li, and Hong Chen. 2010. Parallel SimRank computation on large graphs with iterative aggregation. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*. 543–552.
- [12] Glen Jeh and Jennifer Widom. 2002. SimRank: a measure of structural-context similarity. In *SIGKDD*. 538–543.
- [13] Minhao Jiang, Ada Wai-Chee Fu, and Raymond Chi-Wing Wong. 2017. READS: a random walk approach for efficient and accurate dynamic SimRank. *PVLDB* 10, 9 (2017), 937–948.
- [14] Shivaram Kalyanakrishnan, Ambuj Tewari, Peter Auer, and Peter Stone. 2012. PAC Subset Selection in Stochastic Multi-armed Bandits.. In *ICML*, Vol. 12. 655–662.
- [15] Jérôme Kunegis. 2013. Konect: the koblenz network collection. In *Proceedings of the 22nd international conference on world wide web*. 1343–1350.
- [16] Pei Lee, Laks VS Lakshmanan, and Jeffrey Xu Yu. 2012. On top-k structural similarity search. In *2012 IEEE 28th International Conference on Data Engineering*. IEEE, 774–785.
- [17] Jure Leskovec and Andrej Krevl. 2014. SNAP Datasets: Stanford Large Network Dataset Collection. <http://snap.stanford.edu/data>.
- [18] Cuiping Li, Jiawei Han, Guoming He, Xin Jin, Yizhou Sun, Yintao Yu, and Tianyi Wu. 2010. Fast computation of simrank for static and dynamic information networks. In *Proceedings of the 13th International Conference on Extending Database Technology*. 465–476.
- [19] Lina Li, Cuiping Li, Hong Chen, and Xiaoyong Du. 2013. Mapreduce-based SimRank computation and its application in social recommender system. In *2013 IEEE international congress on big data*. IEEE, 133–140.
- [20] Ruiqi Li, Xiang Zhao, Haichuan Shang, Yifan Chen, and Weidong Xiao. 2017. Fast top-k similarity join for SimRank. *Inf. Sci.* 381 (2017), 1–19. <https://doi.org/10.1016/j.ins.2016.10.042>
- [21] Yu Liu, Bolong Zheng, Xiaodong He, Zhewei Wei, Xiaokui Xiao, Kai Zheng, and Jiaheng Lu. 2017. Probesim: scalable single-source and top-k simrank computations on dynamic graphs. *arXiv preprint arXiv:1709.06955* (2017).
- [22] Yu Liu, Lei Zou, Qian Ge, and Zhewei Wei. 2020. SimTab: accuracy-guaranteed simrank queries through tighter confidence bounds and multi-armed bandits. *Proceedings of the VLDB Endowment* 13, 12 (2020), 2202–2214.
- [23] Dmitry Lizorkin, Pavel Velikhov, Maxim Grinev, and Denis Turdakov. 2008. Accuracy Estimate and Optimization Techniques for SimRank Computation. *Proc. VLDB Endow.* 1, 1 (aug 2008), 422C433. <https://doi.org/10.14778/1453856.1453904>
- [24] Peter Lofgren and Ashish Goel. 2013. Personalized pagerank to a target node. *arXiv preprint arXiv:1304.4658* (2013).
- [25] Juan Lu, Zhiguo Gong, and Xuemin Lin. 2016. A novel and fast simrank algorithm. *IEEE transactions on knowledge and data engineering* 29, 3 (2016), 572–585.
- [26] Linyuan Lü and Tao Zhou. 2011. Link prediction in complex networks: A survey. *Physica A: statistical mechanics and its applications* 390, 6 (2011), 1150–1170.
- [27] Takanori Maehara, Mitsu Kusumoto, and Ken-ichi Kawarabayashi. 2014. Efficient SimRank computation via linearization. In *The 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14, New York, NY, USA - August 24 - 27, 2014*, Sofus A. Macskassy, Claudia Perlich, Jure Leskovec, Wei Wang, and Rayid Ghani (Eds.). ACM, 1426–1435. <https://doi.org/10.1145/2623330.2623696>
- [28] Takanori Maehara, Mitsu Kusumoto, and Ken-ichi Kawarabayashi. 2015. Scalable SimRank join algorithm. In *31st IEEE International Conference on Data Engineering, ICDE 2015, Seoul, South Korea, April 13-17, 2015*, Johannes Gehrke, Wolfgang Lehner, Kyuseok Shim, Sang Kyun Cha, and Guy M. Lohman (Eds.). IEEE Computer Society, 603–614. <https://doi.org/10.1109/ICDE.2015.7113318>
- [29] Yingxia Shao, Bin Cui, Lei Chen, Mingming Liu, and Xing Xie. 2015. An efficient similarity search framework for SimRank over large dynamic graphs. *PVLDB* 8, 8 (2015), 838–849.
- [30] Jieming Shi, Tianyuan Jin, Renchi Yang, Xiaokui Xiao, and Yin Yang. 2020. Realtime Index-Free Single Source SimRank Processing on Web-Scale Graphs. *Proc. VLDB Endow.* 13, 7 (2020), 966–978. <https://doi.org/10.14778/3384345.3384347>
- [31] Wenbo Tao, Minghe Yu, and Guoliang Li. 2014. Efficient Top-K SimRank-based Similarity Join. *Proc. VLDB Endow.* 8, 3 (2014), 317–328. <https://doi.org/10.14778/2735508.2735520>
- [32] Boyu Tian and Xiaokui Xiao. 2016. SLING: A Near-Optimal Index Structure for SimRank. In *SIGMOD*. 1859–1874.
- [33] Anton Tsitsulin, Davide Mottin, Panagiotis Karras, and Emmanuel Müller. 2018. Verse: Versatile graph embeddings from similarity measures. In *Proceedings of the 2018 world wide web conference*. 539–548.
- [34] Hanzhi Wang, Zhewei Wei, Junhao Gan, Sibo Wang, and Zengfeng Huang. 2020. Personalized pagerank to a target node, revisited. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 657–667.
- [35] Hanzhi Wang, Zhewei Wei, Yu Liu, Ye Yuan, Xiaoyong Du, and Ji-Rong Wen. 2021. ExactSim: benchmarking single-source SimRank algorithms with high-precision ground truths. *The VLDB Journal* 30, 6 (2021), 989–1015.
- [36] Sibo Wang, Renchi Yang, Xiaokui Xiao, Zhewei Wei, and Yin Yang. 2017. FO-RA: simple and effective approximate single-source personalized pagerank. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 505–514.
- [37] Yue Wang, Lei Chen, Yulin Che, and Qiong Luo. 2019. Accelerating pairwise simrank estimation over static and dynamic graphs. *The VLDB Journal* 28, 1 (2019), 99–122.
- [38] Yue Wang, Xiang Lian, and Lei Chen. 2018. Efficient simrank tracking in dynamic graphs. In *2018 IEEE 34th international conference on data engineering (ICDE)*. IEEE, 545–556.
- [39] Weiren Yu, Xuemin Lin, and Wenjie Zhang. 2013. Towards efficient SimRank computation on large networks. In *2013 IEEE 29th international conference on data engineering (ICDE)*. IEEE, 601–612.
- [40] Weiren Yu, Xuemin Lin, and Wenjie Zhang. 2014. Fast incremental simrank on link-evolving graphs. In *2014 ieee 30th international conference on data engineering*. IEEE, 304–315.
- [41] Weiren Yu, Xuemin Lin, Wenjie Zhang, Lijun Chang, and Jian Pei. 2013. More is simpler: Effectively and efficiently assessing node-pair similarities based on hyperlinks. *Proceedings of the VLDB Endowment* 7, 1 (2013), 13–24.
- [42] Weiren Yu and Julie A McCann. 2015. Efficient partial-pairs simrank search on large networks. *Proceedings of the VLDB Endowment* 8, 5 (2015), 569–580.
- [43] Weiren Yu, Wenjie Zhang, Xuemin Lin, Qing Zhang, and Jiajin Le. 2012. A space and time efficient algorithm for SimRank computation. *World Wide Web* 15, 3 (2012), 327–353.
- [44] Weiguo Zheng, Lei Zou, Lei Chen, and Dongyan Zhao. 2017. Efficient SimRank-Based Similarity Join. *ACM Trans. Database Syst.* 42, 3 (2017), 16:1–16:37. <https://doi.org/10.1145/3083899>
- [45] Weiguo Zheng, Lei Zou, Yansong Feng, Lei Chen, and Dongyan Zhao. 2013. Efficient SimRank-based Similarity Join Over Large Graphs. *Proc. VLDB Endow.* 6, 7 (2013), 493–504. <https://doi.org/10.14778/2536349.2536350>
- [46] Fanwei Zhu, Yuan Fang, Kai Zhang, Kevin Chang, Hongtai Cao, Zhen Jiang, and Minghui Wu. 2021. Unified and Incremental SimRank: Index-free Approximation with Scheduled Principle. *IEEE Transactions on Knowledge and Data Engineering* (2021).

## 9 APPENDIX

### 9.1 Proofs

**Proof of Lemma 1.** Since SimRank has unique solution [12], we only need to prove that the RHS of Equation 4 is a solution of Equation 1. For each singleton node  $(v, v) \in G^s$ , we have  $\pi_{G^s}((v, v), (v_r, v_r)) = c(1 - c)$ . To be precise, since  $(v_r, v_r)$  is the only in-neighbor of  $(v, v)$ , the probability that a random walk (with stopping probability  $1 - c$ ) starts from  $(v, v)$  and ends at  $(v_r, v_r)$  is  $c(1 - c)$ . It follows that  $\frac{1}{c(1-c)}\pi_{G^s}((v, v), (v_r, v_r)) = 1$ . Next, we consider non-singleton node  $(u, v)$ . Recall that the reverse PPR is defined as  $\pi(s, t) = (1 - a) \sum_{v \in O(t)} \frac{\pi(s, v)}{d_{in}(v)} + a_{s=t}$ . Moreover, it holds that  $\pi(s, t) = \frac{1-a}{d_{in}(s)} \sum_{v \in I(s)} \pi(v, t) + a_{s=t}$  [24]. Therefore, for  $u \neq v$  we have

$$\pi_{G^s}((u, v), (v_r, v_r)) = \frac{c}{|I'(u, v)|} \sum_{(x, y) \in I'(u, v)} \pi_{G^s}((x, y), (v_r, v_r)), \quad (10)$$

where  $I'(u, v)$  is the in-neighbors of node  $(u, v)$ . This is exactly the definition of SimRank for non-singleton nodes. Hence, the RHS of Equation 4 satisfies Equation 1, and the lemma follows.  $\square$

**Proof of Lemma 2.** Recall the time complexity of *Backward Push* [24] for PPR computation w.r.t. the target node  $t$ . With an error parameter  $\varepsilon$ , each push of node  $v$  converts at least  $\alpha\varepsilon$  residue to  $v$ 's reserve, where  $\alpha$  is the stopping probability. Since each push has  $O(d_{in}(v))$  cost, and the number of pushes for  $v$  is bounded by  $\frac{\pi(v, t)}{\alpha\varepsilon}$ . Therefore, the complexity of *Backward Push* is bounded by  $O(\frac{\sum_{v \in V} \pi(v, t) d_{in}(v)}{\varepsilon})$ . Note that we assume the target node  $t$  is given (i.e., fixed), so we do not compute the expected complexity by randomly choosing the target. According to the similar analysis in [34] which considers the increment of residues for  $v$ 's in-neighbors, the complexity is also bounded by  $O(\frac{\sum_{v \in V} \pi(v, t) d_{out}(v)}{\varepsilon})$ . Since all-pair SimRank on  $G$  is equivalent to reverse PPR on  $G^s$  w.r.t. the virtual node  $(v_r, v_r)$ , the time complexity is bounded by  $O(\frac{\sum_{(u, v) \in V_{G^s}} \pi_{G^s}((u, v), (v_r, v_r)) d_{out}(u, v, G^s)}{\varepsilon})$ , and is equivalent to  $O(\frac{\sum_{u, v \in V} d_{out}(u) d_{out}(v) s(u, v)}{\varepsilon})$ . Following [34], it is also bounded by  $O(\frac{\sum_{u, v \in V} d_{in}(u) d_{in}(v) s(u, v)}{\varepsilon})$ .  $\square$

**Proof of Lemma 3.** Suppose that we solve a system of linear equations with  $k$  unknown values  $x_1, \dots, x_k$ , and the h-go cover<sup>+</sup> is represented as  $s_1, \dots, s_{k'}$ . We have  $x_i = \frac{c}{d(i)} (\sum_j A_{ij} x_j + \sum_{j'} A_{ij'} s_{j'})$  where the summation has  $d(i)$  terms. In the offline phase, we have an estimated value of  $s_i$ , denoted as  $\hat{s}_i$ , with  $|\hat{s}_i - s_i| \leq \varepsilon$  for  $i \in [1, k']$ . In the online phase, let  $i = \text{argmax}_j |x_j - \hat{x}_j|$ , where the  $\hat{x}_j$ s are computed from the  $\hat{s}_j$ s. We have  $|x_i - \hat{x}_i| \leq \frac{c}{d(i)} (\sum_j A_{ij} |x_j - \hat{x}_j| + \sum_{j'} A_{ij'} |s_{j'} - \hat{s}_{j'}|) \leq \frac{c}{d(i)} \cdot d(i) \varepsilon \leq \varepsilon$ , and the lemma follows.  $\square$

**Proof of Lemma 4.** By eliminating  $s^{(l)}(u, v)$  for  $l \in [L+1, \infty)$ , the estimation error of *KSimJoin* is bounded by  $c^{L+1}$  [23]. Thus it is sufficient to set  $L = O(\log \frac{1}{\varepsilon})$  to achieve  $O(\varepsilon)$  error.  $\square$

**Proof of Lemma 5.** First, note that all node pairs included in  $\mathcal{R}(\theta)$  belong to the actual result. Second, if some node pair is missed, it must be in the candidate set  $C$ . Therefore,  $\frac{|\mathcal{R}(\theta)|}{|\mathcal{R}(\theta)| + |C|}$  is a lower bound of the actual approximation guarantee. The above claims also hold for an algorithm  $\mathcal{AP}$  with a small failure probability. Assume that we set the failure probability of *APThres* as  $\delta$  (e.g.  $o(\frac{1}{n})$ ). It is

sufficient to set the failure probability of  $\mathcal{AP}$  as  $\delta / \log_2 \frac{1}{\varepsilon_{min}}$ . Then the lemma holds.  $\square$

**Proof of Lemma 6.** We claim when the stopping condition is satisfied, for each  $j \leq \lceil \rho k \rceil$  and each  $j' \geq k+1$ , the SimRank similarity of the  $j$ -th node pair is no smaller than that of the  $j'$ -th one. This is obvious as the lower bound of the former is no smaller than the upper bound of the latter. It means that the first  $\lceil \rho k \rceil$  node pairs are included in the exact result. We also include the next  $k - \lceil \rho k \rceil$  node pairs because the error parameter is typically a loose upper bound of the confidence interval. Similar to Lemma 5, the above discussion holds for algorithm  $\mathcal{AP}$  with a small failure probability.  $\square$

**Proof of Lemma 7.** Since  $R^2LP$  is considered as an extension of *Randomized Backward Search (RBS)* [34] to all-pair SimRank computation, we follow the proofs in [34] for the ease of understanding. Specifically, we first prove the properties of the expectation and variance with our pruning strategy for single-target PPR estimation, which corresponds to Lemma 4.3 & 4.4 of [34]. Then our lemma follows for the SimRank problem.

**LEMMA 10 (BOUNDED BIAS FOR EXPECTATION).** *By applying our pruning strategy to RBS, it holds that  $\mathbb{E}[\hat{\pi}_{l+1}(u, t)] \geq \pi_{l+1}(u, t) - \frac{1-\alpha}{\alpha} \varepsilon$  for  $\forall u \in V$  and  $l \in [0, L]$ .*

**PROOF.** First note that by setting  $f(\varepsilon) = 0$ , which is the same as *RBS*, we have  $\mathbb{E}[\hat{\pi}_{l+1}(u, t)] = \pi_{l+1}(u, t)$ . In particular, let  $X_{l+1}(u, v)$  and  $X'_{l+1}(u, v)$  denote the increment of  $\hat{\pi}_{l+1}(u, t)$  by conducting *Push*( $v$ ) without and with pruning, respectively. According to Lemma 4.3 of [34], we have  $\mathbb{E}[X_{l+1}(u, v) | \{\hat{\pi}_l\}] = \frac{(1-\alpha)\hat{\pi}_l(v, t)}{d_{out}(u)}$ . Next, with the pruning strategy, we have  $\mathbb{E}[X'_{l+1}(u, v) | \{\hat{\pi}_l\}] \geq \frac{(1-\alpha)(\hat{\pi}_l(v, t) - \varepsilon)}{d_{out}(u)}$ . We prove by the following induction procedure. For  $l = 0$ ,  $\hat{\pi}_0(t, t) = \alpha$  and  $\hat{\pi}_0(v, t) = 0$  for all other  $vs$ . For  $l = 1$ , we have  $\mathbb{E}[\hat{\pi}'_1(v, t)] = \pi_1(v, t)$  (we assume  $\varepsilon < \alpha$ ). Similarly,  $\mathbb{E}[\hat{\pi}'_1(v, t)] \geq \pi_1(v, t) - (1 - \alpha)\varepsilon$ . Suppose  $\mathbb{E}[\hat{\pi}'_l(v, t)] \geq \pi_l(v, t) - \sum_{i=1}^{l-1} (1 - \alpha)^i \varepsilon$  holds, then it follows that  $\hat{\pi}'_{l+1}(u, t) = \sum_{v \in N_{out}(u)} X'_{l+1}(u, v)$  and

$$\begin{aligned} \mathbb{E}[\hat{\pi}'_{l+1}(u, t) | \{\hat{\pi}_l\}] &= \sum_{v \in N_{out}(u)} \mathbb{E}[X'_{l+1}(u, v) | \{\hat{\pi}_l\}] \\ &\geq \frac{(1 - \alpha)\hat{\pi}'_l(v, t)}{d_{out}(u)} - (1 - \alpha)\varepsilon. \end{aligned} \quad (11)$$

Hence, we can derive that

$$\begin{aligned} \mathbb{E}[\hat{\pi}'_{l+1}(u, t)] &\geq \frac{(1 - \alpha)\mathbb{E}[\hat{\pi}'_l(v, t)]}{d_{out}(u)} - (1 - \alpha)\varepsilon \\ &\geq \frac{(1 - \alpha)(\pi_l(v, t) - \sum_{i=1}^{l-1} (1 - \alpha)^i \varepsilon)}{d_{out}(u)} - (1 - \alpha)\varepsilon \\ &= \pi_{l+1}(u, t) - \sum_{i=1}^l (1 - \alpha)^i \varepsilon \\ &\geq \pi_{l+1}(u, t) - \frac{1 - \alpha}{\alpha} \varepsilon. \end{aligned} \quad (12)$$

$\square$

**LEMMA 11 (BOUNDED VARIANCE WITH PRUNING).** *With our pruning strategy, the variance of every  $\hat{\pi}_l(v, t)$  is still bounded by Lemma 4.4 of [34]. Specifically, for additive error guarantee, it holds that  $\text{Var}[\hat{\pi}_l(v, t)] \leq \alpha \varepsilon^2$ .*

**PROOF.** The key observation is that (1)the variance comes from the randomized push [34] and (2)via pruning, we have

$$\text{Var}[\hat{\pi}'_{l+1}(u, t)|\{\hat{\pi}_l\}] \leq \text{Var}[\hat{\pi}_{l+1}(u, t)|\{\hat{\pi}_l\}]. \quad (13)$$

Then the lemma holds following the proof of Lemma 4.4 in [34].  $\square$

Per the above lemmas, we have  $\mathbb{E}[\hat{\pi}(u, t)] \geq \pi(u, t) - L \frac{1-\alpha}{\alpha} \varepsilon$ , and the variance of  $\hat{\pi}(u, t)$  is also bounded by  $L \alpha \varepsilon^2$ . By setting  $L = O(\log \frac{1}{\varepsilon})$ , the truncation error is bounded by  $O(\varepsilon)$ . To be precise, we can set  $L = \log_{1-\alpha} \frac{\varepsilon}{3}$  to achieve a truncation error  $\varepsilon_t$  of  $\frac{\varepsilon}{3}$ . The pruning error parameter  $\varepsilon_p$  is set to  $\frac{\varepsilon}{3} \cdot \frac{\alpha}{1-\alpha}$  so that  $L \frac{1-\alpha}{\alpha} \varepsilon_p \leq \frac{2\varepsilon}{3}$ . The total estimation error is thus bounded by  $\varepsilon$ . By adopting the Median-of-Mean trick, we can bound the failure probability as in [34]. The above conclusion also holds for all-pair SimRank estimation because of the equivalence to single-target PPR (Lemma 1).  $\square$

**Proof of Theorem 1.** It follows directly from Lemma 4.5 of [34].  $\square$

**Proof of Lemma 8.** According to the *APThres* (or *APTop-k*) framework and with  $p > 1$ , the computational complexity is bounded by  $\sum_{i=1}^t \text{cost}(\mathcal{AP}(G, \varepsilon_i)) \leq \text{cost}(\mathcal{AP}(G, \varepsilon_t)) \cdot (1 + \frac{1}{p} + \dots + \frac{1}{p^{t-1}}) = O(\text{cost}(\mathcal{AP}(G, \varepsilon_t)))$ .  $\square$

**Proof of Lemma 9.** Recall that *APThres* stops when  $\frac{|\mathcal{R}(\theta)|}{|\mathcal{R}(\theta)| + |\mathcal{C}|} \geq \rho$ . Assume that at  $t$ -th iteration the condition is satisfied, we develop the lower bound of  $|\mathcal{R}(\theta)|$  and the upper bound of  $|\mathcal{R}(\theta) \cup \mathcal{C}|$  (note that  $\mathcal{R}(\theta)$  and  $\mathcal{C}$  are disjoint). With high probability, the following claims hold. We have  $\{(u, v) | s(u, v) \geq \theta + 2\varepsilon_t\} \subseteq \mathcal{R}(\theta)$ , as any such node pair  $(u, v)$  has  $s_t(u, v) \geq \theta + \varepsilon_t$ . We find the maximum  $j$  and  $x$  (denoted as  $j^*$  and  $x^*$  respectively) such that  $b_1 r(j^*)^{-\beta} \geq \theta + 2\varepsilon_t$  and  $j^* = 2^{x^*+1} - 1$ . It follows that  $j^* \leq (\frac{b_1}{\theta + 2\varepsilon_t})^{\frac{1}{\beta}} + t$ . Similarly, we have  $\mathcal{R}(\theta) \cup \mathcal{C} \subseteq \{(u, v) | s(u, v) \geq \theta - 2\varepsilon_t\}$ . Then we find the minimum  $j$  and  $x$  (denoted as  $j'$  and  $x'$  respectively) with  $b_2 r(j')^{-\beta} \leq \theta - 2\varepsilon_t$  and  $j' = 2^{x'+1} - 1$ . It follows that  $j' \geq (\frac{b_2}{\theta - 2\varepsilon_t})^{\frac{1}{\beta}} + t$ . If  $j^* \geq \rho j'$ , the algorithm must have been stopped at  $t$ -th iteration. It is sufficient to have  $(\frac{b_1}{\theta + 2\varepsilon_t})^{\frac{1}{\beta}} \geq \rho (\frac{b_2}{\theta - 2\varepsilon_t})^{\frac{1}{\beta}}$ , which is equivalent to  $\varepsilon_t \leq \frac{(b_1 - \rho^\beta b_2)\theta}{2(b_1 + \rho^\beta b_2)}$ , on condition that  $b_1 - \rho^\beta b_2 > 0$ . By setting  $\varepsilon_t = O(\max(\frac{(b_1 - \rho^\beta b_2)\theta}{b_1 + \rho^\beta b_2}, \varepsilon_{min}))$ , the lemma follows.  $\square$

**Proof of Theorem 2.** It follows directly from Lemma 8 & 9.  $\square$

**Proof of Theorem 3.** To achieve an approximation bound of  $\rho$ , it is sufficient to set  $s_{\lceil \rho k \rceil} - s_{k+1} \geq 2\varepsilon'_t$  for *APTop-k*. By Equation 6, let  $x_1 = 2^{\lceil \log_2 \lceil \rho k \rceil \rceil} - 1$  and  $x_2 = 2^{\lfloor \log_2(k+1) \rfloor}$ , we have  $s_{\lceil \rho k \rceil} \geq b_1(x_1 - t)^{-\beta}$  and  $s_{k+1} \leq b_2(x_2 - t)^{-\beta}$ . By setting  $\varepsilon'_t = \max(\frac{1}{2}(b_1(x_1 - t)^{-\beta} - b_2(x_2 - t)^{-\beta}), \varepsilon_{min})$ , *APTop-k* guarantees the approximation bound with high probability. (Note that the first term may be less than zero, so we need the second term  $\varepsilon_{min}$ .) By assuming the power law distribution of SimRank and using similar derivation, we have  $\varepsilon'_t = \max(\frac{1}{2}(b_1(\lceil \rho k \rceil - t)^{-\beta} - b_2(k+1 - t)^{-\beta}), \varepsilon_{min})$ .  $\square$

## 9.2 Additional Experiments

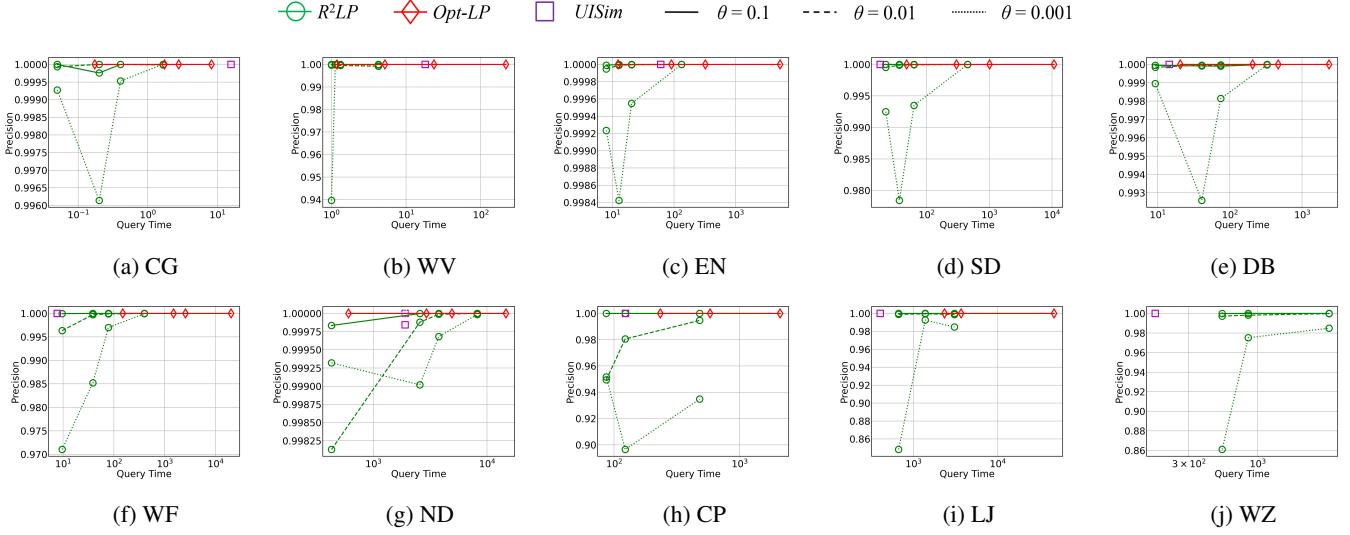
**Additional experiments for threshold-based similarity joins.** We show the *Precision* (Figure 9) and *Recall* (Figure 10) separately for

the experiments in Figure 5. For *Precision*, first note that *Opt-LP* has *Precision* equal to 1 on all datasets because it always makes underestimations for SimRank values. Surprisingly, *UISim* also underestimates all node similarities on almost all queries (except on ND). The algorithms only considers random walks with limited length in practice, therefore the truncation error outweighs the overestimation by allowing multi-meeting probabilities. As for *R<sup>2</sup>LP*, the *Precision* is close to 1 on most small and medium graphs. In our experiments, following the setting of the baselines, we return all node pairs with SimRank estimations above  $\theta$  rather than  $\theta + \varepsilon$ , where the latter also achieves *Precision* close to 1 in practice. As randomness of the algorithm introduces both underestimation as well as overestimation, it also has false positives. On the other hand, *R<sup>2</sup>LP* significantly outperforms *Opt-LP* and *UISim* in terms of *Recall*, since the baselines sacrifice *Recall* for *Precision*. Overall, *R<sup>2</sup>LP* still has a clear advantage over the baselines as demonstrated in Figure 5.

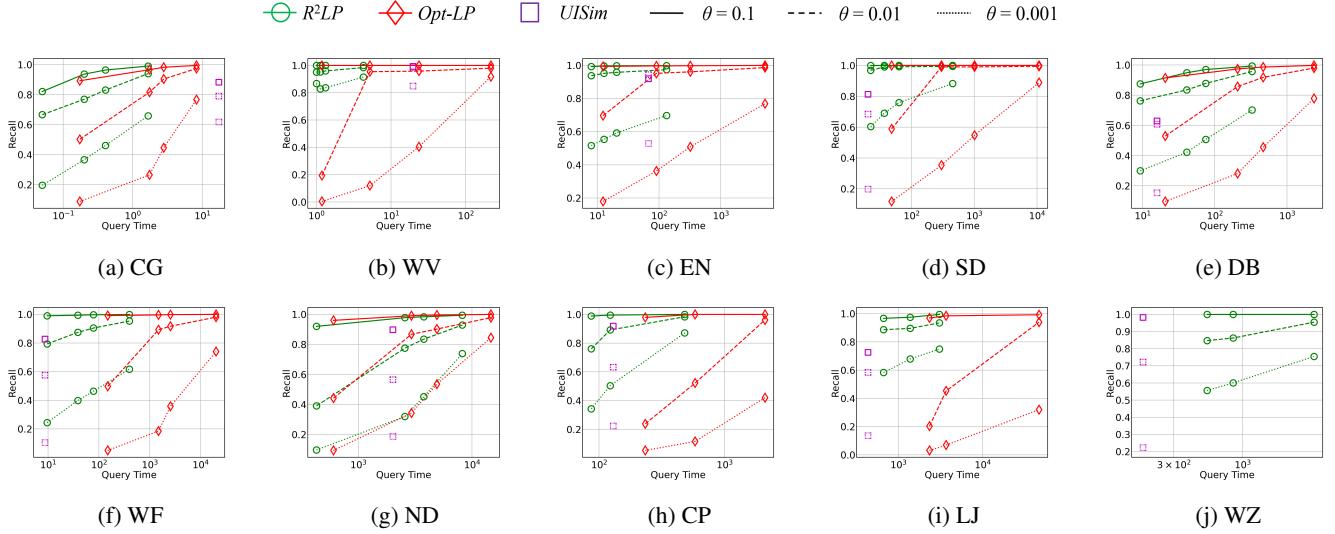
**Additional experiments for top-k similarity joins.** Table 6 demonstrates the query time and Precision@ $k$  of the experiments in Figure 6. Generally speaking, *APTop-k*+*R<sup>2</sup>LP* is more efficient than *APTop-k*+*Opt-LP* to achieve comparable accuracy. On 7 of 10 datasets, the speedup is by at least 2×.

**Empirical validation for the PLB distribution of SimRank.** We empirically validate that all the small and medium graphs used in our experiments have ground truth SimRank distribution follow PLB. The three large datasets are ignored since we only have the single-source ground truth of 100 query nodes. Note that the definition of PLB is different from the power law distribution, namely, we do not assume the concrete distribution form for the SimRank values. We instead suppose that the distribution is both upper and lower bounded by two power law distributions of the same exponent but with different scaling factors. Therefore, we do not employ the maximum likelihood estimation (MLE) [7] and directly minimize our fitting objective by considering all logarithmic bins.

The pseudocode is illustrated in Algorithm 4. In particular, we empirically observe that in almost all datasets, the all-pair SimRank distribution approximately follow power law for node pairs except those with similarity values significantly above 0.1. (Actually, those node pairs with similarity above 0.1 also follow power law approximately but with a different slope. Besides, they are easy to estimate for most algorithms since two random walks from such a node pair tend to meet after only one or two steps.) Therefore, we focus on the modeling of these node pairs by only considering the bins with  $x \in [x_{min}, x_{max}]$ , where we set  $x_{min}$  as the smallest bin with SimRank value below 0.1, and set  $x_{max} = \lfloor \log_2 nnz \rfloor$ . In particular, we first compute the average SimRank of each bin (Lines 1-2). Next, we iterate over the possible values of  $\beta \in (0, 5]$  with step equal to 0.001 (Lines 4-8). For each value of  $\beta$ , we find the maximum  $b_1$  and minimum  $b_2$  with  $b_1, b_2 \in [0.001, 50]$  and  $b_2 > b_1$ , such that  $b_1 \sum_j j^{-\beta}$  (resp.  $b_2 \sum_j j^{-\beta}$ ) is the lower (resp. upper) bound of  $S_x$ . Our object is to minimize the ratio between  $b_2$  and  $b_1$ . The fitted values are demonstrated in Table 6.



**Figure 9: Precision of the experiments in Figure 5 for threshold-based queries. No approximation bound is provided.**



**Figure 10: Recall of the experiments in Figure 5 for threshold-based queries. No approximation bound is provided.**

---

**Algorithm 4:** Fitting the parameters of PLB distribution

---

**Input:** The all-pair ground truth SimRank values  $\{s_1, \dots, s_{nnz}\}$  sorted in descending order

**Output:**  $\beta$ ,  $b_1$ , and  $b_2$

- 1  $b_1^* \leftarrow 0, b_2^* \leftarrow \infty;$
- 2 **for**  $x = x_{min}$  to  $x_{max}$  **do**
- 3   | Compute  $S_x = \sum_{j=2^x}^{\min(2^{x+1}-1, nnz)} s_j$ ;
- 4 **for**  $\beta = 0.001$  to  $5$  **do**
- 5   | Find maximum  $b_1$  and minimum  $b_2$  with  $b_1, b_2 \in [0.001, 50]$  and  $b_2 > b_1$ , such that  $b_1 \sum_j j^{-\beta}$  (resp.  $b_2 \sum_j j^{-\beta}$ ) is the lower (resp. upper) bound of  $S_x$ ;
- 6   | **if**  $b_2/b_1 < b_2^*/b_1^*$  **then**
- 7     | |  $\beta^* \leftarrow \beta, b_1^* \leftarrow b_1, b_2^* \leftarrow b_2$ ;
- 8   | Increase  $\beta$  by  $0.001$ ;
- 9 **return**  $\beta^*, b_1^*$ , and  $b_2^*$ ;

---

**Table 6: Query time and Precision@ $k$  for  $k = \{500, 1000, 5000, 10000\}$  with  $\rho = 0.9$ . This corresponds to the experiments shown in Figure 6.**

	Opt-LP								R <sup>2</sup> LP							
	k = 500		k = 1000		k = 5000		k = 10000		k = 500		k = 1000		k = 5000		k = 10000	
	Time	P@ $k$	Time	P@ $k$	Time	P@ $k$	Time	P@ $k$	Time	P@ $k$	Time	P@ $k$	Time	P@ $k$	Time	P@ $k$
CG	2.21	0.994	3.36	0.998	2.16	0.999	1.96	0.997	1.45	0.998	2.65	0.997	1.42	0.998	1.40	0.995
WV	76.6	1	206.0	0.999	110.3	0.999	138.6	0.998	10.8	1	82.1	0.999	24.6	0.998	47.0	0.997
EN	5.2	1	5.5	1	443.7	0.997	2185.5	0.998	7.0	1	6.9	1	88.0	0.996	284.9	0.996
SD	1277	0.984	3095	0.997	7682	0.996	24798	0.990	161	0.984	276	0.998	2424	0.997	11564	0.997
DB	465	0.998	228	0.996	1082	0.998	2052	0.998	404	0.996	103	0.994	334	0.995	631	0.996
WF	1197	0.988	10927	0.998	10501	0.99	22157	0.989	66	0.99	506	0.997	3488	0.995	6777	0.994
ND	1281	0.99	1278	0.998	4964	0.997	10336	0.999	5360	0.996	2211	0.996	11749	0.993	23895	0.996
CP	2147	0.986	2122	1	2037	0.992	1882	0.994	1126	1	1168	0.999	1190	0.994	1063	0.979
LJ	11107	0.948	10109	0.973	8586	0.964	8210	0.639	8953	0.988	8845	0.986	8824	0.970	8926	0.947
WZ	OOM								334	1	7277	1	7171	1	3762	0.998