

# Approximating Probabilistic Group Steiner Trees in Graphs: Supplemental Materials

This supplement is available online [1]. The road map of this supplement is as follows.

- In Section S1, we prove the approximation guarantee of GRE-PATH.
- In Sections S2-S4, we discuss the time complexities of DUAL, GRE-TREE and GRE-PATH.
- In Section S5, we incorporate ImprovAPP with DUAL and GRE-TREE.
- In Section S6, we conduct additional case studies.
- In Section S7, we discuss the methods of selecting vertex groups in experiments.
- In Section S8, we report the memory consumption of algorithms in experiments.
- In Section S9, we show the sizes of  $|V_g|$  and  $|g_{min}|$  in experiments.
- In Section S10, we conduct experiments where edge weights are pairwise Jaccard distances.

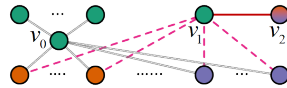


Fig. S1. The sharpness of  $\max\{1, \sum_{g \in \Gamma} \xi_g - 1\}$ .

## S1. THE APPROXIMATION GUARANTEE OF GRE-PATH

**Theorem 3.** GRE-PATH has a sharp approximation guarantee of

$$\max\{1, \sum_{g \in \Gamma} \xi_g - 1\}$$

for solving the probabilistic group Steiner tree problem.

*Proof.* Let  $\Theta_{opt}$  be an optimal solution. Suppose that  $\Theta_{opt}$  contains  $v \in g_{min}$ . Let  $\Theta_v$  be the feasible solution produced by GRE-PATH in the loop for  $v$  (Lines 3-18). Line 17 guarantees that

$$c(\Theta_v) \leq c(\Theta_{opt}). \quad (S1)$$

To satisfactorily cover each vertex group  $g$ , GRE-PATH merges at most  $\xi_g$  paths from  $v$  to vertices in  $g$ . Suppose that GRE-PATH merges  $P(v, u_1), P(v, u_2), \dots, P(v, u_x)$  for satisfactorily covering  $g$ , i.e.,  $u_1, u_2, \dots, u_x \in g$ , and

$$c(P(v, u_1)) \leq c(P(v, u_2)) \leq \dots \leq c(P(v, u_x)). \quad (S2)$$

Since  $\Theta_{opt}$  contains  $v$  and satisfactorily covers  $g$ , there are two possible scenarios: (i)  $\Theta_{opt}$  contains  $u_1, u_2, \dots, u_x$ , and does not contain any other vertex in  $g$ ; (ii)  $\Theta_{opt}$  contains  $u_y \in g$  such that  $u_y \notin \{u_1, u_2, \dots, u_x\}$ . Since GRE-PATH merges shortest paths in an increasing order of the weights of these paths, in both scenarios,

$$c(P(v, u_x)) \leq c(\Theta_{opt}). \quad (S3)$$

Thus, the weight of each path that is merged into  $\Theta_v$  is not larger than  $c(\Theta_{opt})$ . To satisfactorily cover every vertex group, GRE-PATH merges at most  $\sum_{g \in \Gamma} \xi_g$  paths into  $\Theta_v$ . Moreover, the single-vertex path  $P(v, v)$  can be seen as merged into  $\Theta_v$  for satisfactorily covering  $g_{min}$ . Since  $c(P(v, v)) = 0$ , we have

$$c(\Theta_v) \leq \left( \sum_{g \in \Gamma} \xi_g - 1 \right) \cdot c(\Theta_{opt}). \quad (S4)$$

Therefore,  $\sum_{g \in \Gamma} \xi_g - 1$  is an approximation guarantee of GRE-PATH. We prove the sharpness of this guarantee via the instance in Figure S1, where  $g_{min}$  is the set of green vertices. There are three special vertices:  $v_0 \in g_{min}$ ,  $v_1 \in g_{min}$  and  $v_2 \notin g_{min}$ . There is an edge between  $v_0$  and any other vertex except  $v_1$  and  $v_2$ , and the weight of each of these edges (i.e., each gray edge) is 1. There is an edge between  $v_1$  and any other vertex that is not in  $g_{min}$ , and the weight of each of these edges except edge  $(v_1, v_2)$  (i.e., each pink edge) is  $1 + \delta$ , while the weight of  $(v_1, v_2)$  (i.e., the red edge) is  $1 + 2\delta$ , where  $\delta$  is a positive value. There is no other edge in this graph.  $p_{g_{min}}(v)$  are equal small

values for every  $v \in g_{min} \setminus v_1$ , while  $p_{g_{min}}(v_1) = 1$ .  $v_2$  is in every vertex group  $g \in \Gamma \setminus g_{min}$  (e.g., the set of orange vertices and  $v_2$  is a group, and the set of purple vertices and  $v_2$  is another group).  $p_g(v_2) = 1$  for every  $g \in \Gamma \setminus g_{min}$ . Furthermore, the set of gray edges corresponds to a feasible solution. In the loop of Lines 3-18, GRE-PATH produces a feasible solution  $\Theta_{v_0}$ , which contains and only contains all the gray edges, and

$$c(\Theta_{v_0}) = \sum_{g \in \Gamma} \xi_g - 1. \quad (S5)$$

Moreover, GRE-PATH also produces a feasible solution  $\Theta_{v_1}$ , which contains and only contains all the pink edges, and

$$c(\Theta_{v_1}) = \sum_{g \in \Gamma \setminus g_{min}} \xi_g \cdot (1 + \delta). \quad (S6)$$

Suppose that  $c(\Theta_{v_0}) \leq c(\Theta_{v_1})$ , which means that

$$\sum_{g \in \Gamma} \xi_g - 1 \leq \sum_{g \in \Gamma \setminus g_{min}} \xi_g + \sum_{g \in \Gamma \setminus g_{min}} \xi_g \cdot \delta, \quad (S7)$$

$$\delta \geq \frac{\xi_{g_{min}} - 1}{\sum_{g \in \Gamma \setminus g_{min}} \xi_g}. \quad (S8)$$

Thus, when  $|\Gamma|$  is large,  $\delta$  can be a tiny value. When  $\delta$  is tiny, the optimal solution is  $\Theta_{opt} = (v_1, v_2)$ , and

$$c(\Theta_{opt}) = 1 + 2\delta. \quad (S9)$$

In the above case where  $c(\Theta_{v_0}) \leq c(\Theta_{v_1})$ , the solution of GRE-PATH is  $\Theta_3 = \Theta_{v_0}$ . The approximation ratio of GRE-PATH is

$$\lim_{\delta \rightarrow 0} \frac{c(\Theta_3)}{c(\Theta_{opt})} = \frac{\sum_{g \in \Gamma} \xi_g - 1}{1 + 2\delta} = \sum_{g \in \Gamma} \xi_g - 1. \quad (S10)$$

Moreover, in scenarios where  $\Gamma = \{g\}$  and  $\xi_g = 1$ , GRE-PATH can find optimal solutions, *i.e.*, the approximation ratio of GRE-PATH is 1. Hence,  $\max\{1, \sum_{g \in \Gamma} \xi_g - 1\}$  is a sharp approximation guarantee of GRE-PATH. This theorem holds.  $\square$

## S2. THE TIME COMPLEXITY OF DUAL

DUAL has a time complexity of

$$O\left(|\Gamma| \xi^2 |V|^{2\xi} \cdot \left(3^{2\xi} |V| + 2^{2\xi} |V| \cdot (2\xi |V| + |E| + |V| \log |V|)\right)\right),$$

where  $\xi$  is the smallest natural number that is larger than or equal to  $\log_{(1-p_{min})}(1-b)$ , and  $p_{min}$  is the minimum value of  $p_g(v)$  for every  $v \in g \in \Gamma$ . The details are as follows. DUAL initializes a tree (line 1) in  $O(1)$  time. It finds  $\Phi_g$  for every  $g \in \Gamma$  (Line 2) in  $O(|\Gamma| \cdot \xi^2 \cdot |V|^\xi)$  time. The reason is as follows. For any group  $g \in \Gamma$ , let  $V_g$  be an essential cover of  $g$ . Since

$$1 - (1 - p_{min})^\xi \geq 1 - (1 - p_{min})^{\log_{(1-p_{min})}(1-b)} \geq b, \quad (S11)$$

we have

$$|V_g| \leq \min\{|g|, \xi\}. \quad (S12)$$

That is to say, the size of any essential cover of any vertex group  $g$  is not larger than  $\min\{|g|, \xi\}$ . To find  $\Phi_g$ , we find every set of  $k$  vertices in  $g$  for every  $k \in [1, \min\{|g|, \xi\}]$ , and check whether these sets of vertices are essential covers of  $g$ . Since checking whether a set of vertices  $V'$  is an essential cover of  $g$  takes  $O(|V'|)$  time and enumerating every set of  $k$  vertices in  $g$  takes  $O(|g|^k)$  time, DUAL finds  $\Phi_g$  for every  $g \in \Gamma$  (Line 2) in

$$O\left(|\Gamma| \cdot (|V|^1 + \dots + |V|^\xi)\right) = O\left(|\Gamma| \cdot \xi^2 \cdot |V|^\xi\right) \quad (S13)$$

time. It finds  $\Phi_{g_x}$  in  $O(|\Gamma|)$  time. Since  $|V_{g_x}| \leq \min\{|g_x|, \xi\}$ ,

$$O(|\Phi_{g_x}|) = O\left(\sum_{j \in [1, \min\{|g_x|, \xi\}]} \binom{|g_x|}{j}\right) = O(\xi \cdot |V|^\xi). \quad (S14)$$

It processes each  $V'$  in  $\Phi_{g_x}$  as follows (Lines 4-18). If  $|\Gamma| = 1$  (Line 4), it uses PrunedDP++ to find  $\Theta_{V'}$  (Line 5) in

$$O(3^\xi |V| + 2^\xi |V| \cdot (\xi |V| + |E| + |V| \log |V|))$$

time, where  $O(3^\xi |V|)$  corresponds to a tree merging process in PrunedDP++, and  $O(2^\xi |V| \cdot (\xi |V| + |E| + |V| \log |V|))$  corresponds to merging  $O(\xi)$  shortest paths together and finding an MST of the merged graph for producing a feasible solution tree  $O(2^\xi |V|)$  times (details in [3]; notably, PrunedDP++ improves the previous DPBF algorithm [2] on practical efficiency, but has a larger time complexity than DPBF). If  $|\Gamma| > 1$ , it builds  $\Theta_{V'}$  as follows (Lines 7-16). It initializes  $G'$  in  $O(1)$  time (Line 7). For each  $g \in \Gamma \setminus g_x$  and each  $V_j \in \Phi_g$ , it finds  $\Theta(V', V_j)$  (Line 11) in

$$O(3^{2\xi} |V| + 2^{2\xi} |V| \cdot (2\xi |V| + |E| + |V| \log |V|))$$

time. Thus, the cost of Line 11 throughout the loop of Lines 8-15 is

$$O(|\Gamma| \xi |V|^\xi \cdot (3^{2\xi} |V| + 2^{2\xi} |V| \cdot (2\xi |V| + |E| + |V| \log |V|))).$$

It uses each computed  $\Theta(V', V_j)$  to update  $\Theta_{ST}(V', \Phi_g)$  (Line 12), and each update takes  $O(|V|)$  time. For each  $g \in \Gamma \setminus g_x$ , it merges  $\Theta_{ST}(V', \Phi_g)$  into  $G'$  (Line 14) in  $O(|V|)$  time. After the loop of Lines 8-15, it finds an MST as  $\Theta_{V'}$  (Line 16) in  $O(|E| + |V| \log |V|)$  time. It updates  $\Theta_1$  (Line 18) in  $O(|V|)$  time. At last, it returns  $\Theta_1$  (Line 20) in  $O(|V|)$  time.

### S3. THE TIME COMPLEXITY OF GRE-TREE

GRE-TREE has a time complexity of

$$O(\xi \cdot |g_{min}| \cdot (3^{|\Gamma|} |V| + 2^{|\Gamma|} |V| \cdot (|\Gamma| |V| + |E| + |V| \log |V|))).$$

The details are as follows. The algorithm initializes  $\Theta_2$  (Line 1) in  $O(1)$  time. Then, it finds  $g_{min}$  (Line 2) in  $O(|\Gamma|)$  time, and processes each vertex  $v \in g_{min}$  as follows (Lines 3-12). First, it initializes  $G' = \{v\}$  (Line 4) in  $O(1)$  time. It iteratively concatenates trees into  $G'$  through a loop (Lines 5-9). To efficiently check whether  $G'$  satisfactorily covers every group or not (Line 5), we record the probability that  $G'$  does not satisfactorily cover each group, and update these probabilities whenever a new vertex is added into  $G'$  in Line 8. By doing this, the cost of checking whether  $G'$  satisfactorily covers every group or not (Line 5) throughout the loop of Lines 5-9 is  $O(|\Gamma| |V|)$ . We also record  $\Gamma'$  in Line 6, and update  $\Gamma'$  whenever a new vertex is added into  $G'$  in Line 8. By doing this, the cost of constructing  $\Gamma'$  in Line 6 throughout the loop of Lines 5-9 is also  $O(|\Gamma| |V|)$ . There are  $O(\xi)$  iterations in the above loop. In each iteration, it uses PrunedDP++ to produce  $\Theta'$  (Line 7) at a cost of

$$O(3^{|\Gamma|} |V| + 2^{|\Gamma|} |V| \cdot (|\Gamma| |V| + |E| + |V| \log |V|)).$$

Then, it merges  $\Theta'$  into  $G'$  (Line 8) in  $O(|V|)$  time. After the above loop, it finds an MST as  $\Theta_v$  (Line 10) in  $O(|E| + |V| \log |V|)$  time, and uses  $\Theta_v$  to update  $\Theta_2$  (Line 11) at a cost of  $O(|V|)$ . After enumerating every  $v \in g_{min}$ , it returns  $\Theta_2$  (Line 13) at a cost of  $O(|V|)$ .

### S4. THE TIME COMPLEXITY OF GRE-PATH

GRE-PATH has a time complexity of

$$O(|g_{min}| \cdot \sum_{g \in \Gamma} \xi_g \cdot L |V| + |E| + |V| \log |V|),$$

where  $L$  is the average number of hub labels associated with each vertex (details in [5-7]). The details are as follows. First, the algorithm initializes an empty tree in  $O(1)$  time (Line 1). Then, it

finds  $g_{min}$  at a cost of  $O(|\Gamma|)$  (Line 2), and processes each vertex  $v \in g_{min}$  as follows (Lines 3-15). It initializes and popularizes  $|\Gamma|$  heaps (Lines 4-6) in  $O(|\Gamma| \cdot |V| \cdot L)$  time, since the time complexity of checking the cost of  $P(v, u)$  in Line 5 is  $O(L)$ , and we can build a Binary heap with  $n$  elements in  $O(n)$  time in a bottom-up manner. Subsequently, it initializes  $\Theta_v$  in  $O(1)$  time (Line 7). To efficiently check the probability that  $\Theta_v$  satisfactorily covers each group, we record the probability that  $\Theta_v$  does not satisfactorily cover each group, and update these probabilities whenever a new vertex is added into  $\Theta_v$  in Line 11. There are  $O(|V|)$  vertices added into  $\Theta_v$ . Whenever a new vertex is added into  $\Theta_v$ , we update the recorded probabilities in  $O(|\Gamma|)$  time. Thus, the cost of updating the recorded probabilities throughout the loop of Lines 8-13 is  $O(|\Gamma||V|)$ . The cost of checking the condition in Line 9 using the recorded probabilities is  $O(1)$ . There are  $O(\sum_{g \in \Gamma} \xi_g)$  while iterations (Lines 10-11). In each iteration, it pops out the top element of  $Q_g$  (Line 10) in  $O(\log |V|)$  time, and merges a path into  $\Theta_v$  (Line 11) in  $O(L|V|)$  time. After the loop, it updates  $\Theta_3$  using  $\Theta_v$  (Line 14) in  $O(|V|)$  time. In the end, it updates  $\Theta_3$  to be an MST that spans the vertices in  $\Theta_3$  (Line 16) in  $O(|E| + |V| \log |V|)$  time.

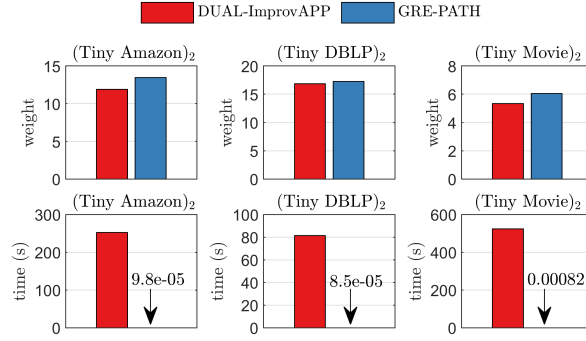


Fig. S2. DUAL-ImprovAPP is significantly slower than GRE-PATH.

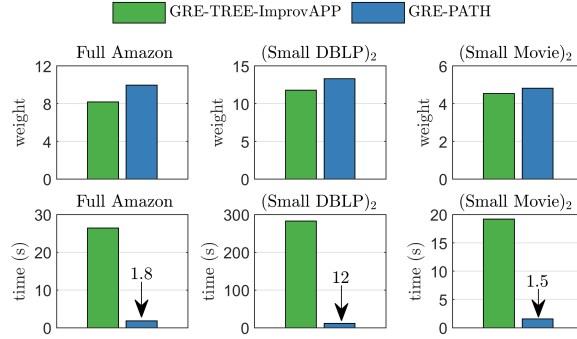


Fig. S3. GRE-TREE-ImprovAPP is significantly slower than GRE-PATH.

## S5. INCORPORATING IMPROVAPP WITH DUAL AND GRE-TREE

DUAL and GRE-TREE incorporate a state-of-the-art classical group Steiner tree algorithm: PrunedDP++ [3], which is based on a dynamic programming approach. We can also replace PrunedDP++ in DUAL and GRE-TREE with another state-of-the-art classical group Steiner tree algorithm: ImprovAPP [8], which is based on a greedy concatenation approach. Specifically, we can replace PrunedDP++ with ImprovAPP in Lines 5 and 11 of DUAL, as well as in Line 7 of GRE-TREE. We refer to the two algorithms after the replacement as DUAL-ImprovAPP and GRE-TREE-ImprovAPP, respectively.

In Figure S2, we compare DUAL-ImprovAPP with GRE-PATH, where  $|V| = 300$  for "Tiny Amazon", "Tiny DBLP" and "Tiny Movie", and the other parameter settings are the same with those in the main experiments. It can be seen that the solution weights of DUAL-ImprovAPP are smaller than those of GRE-PATH, at the cost of a significantly lower speed than GRE-PATH, since a large number of essential covers are enumerated in Lines 3 and 10 of Algorithm 1. In Figure S3, we compare GRE-TREE-ImprovAPP with GRE-PATH, where  $|V| = 548,552$  for "Amazon",  $|V| = 897,782$  for "Small DBLP" and  $|V| = 4,000$  for "Small Movie". Like the above experiment results, the solution weights of GRE-TREE-ImprovAPP are smaller than those of GRE-PATH, at the cost of a significantly

lower speed than GRE-PATH, since, for each vertex  $v \in g_{min}$ , GRE-TREE-ImprovAPP implements ImprovAPP multiple times in the while loop in Lines 5-9 of Algorithm 2.

Since ImprovAPP has an approximation guarantee of  $\max\{1, |\Gamma| - 1\}$  for solving the classical group Steiner tree problem, DUAL-ImprovAPP has an approximation guarantee of  $(\max\{1, |\Gamma| - 1\})^2$  (by changing  $\tau$  in the approximation guarantee of DUAL to  $\max\{1, |\Gamma| - 1\}$ ), while GRE-TREE-ImprovAPP has an approximation guarantee of  $\max\{1, |\Gamma| - 1\} \cdot \xi$  (by changing  $\tau$  in the approximation guarantee of GRE-TREE to  $\max\{1, |\Gamma| - 1\}$ ). Therefore, DUAL and GRE-TREE can achieve tighter guarantees than DUAL-ImprovAPP and GRE-TREE-ImprovAPP, respectively, by setting  $\tau \geq 1$  smaller than  $\max\{1, |\Gamma| - 1\}$  when  $|\Gamma| > 2$ .

Given that, in comparison with GRE-PATH, DUAL-ImprovAPP and GRE-TREE-ImprovAPP have a significantly low efficiency and thus mainly of theoretical interests, we choose to make these theoretical interests as strong as possible. As a result, we incorporate PrunedDP++, but not ImprovAPP, with DUAL and GRE-TREE, for achieving tight and tunable approximation guarantees.

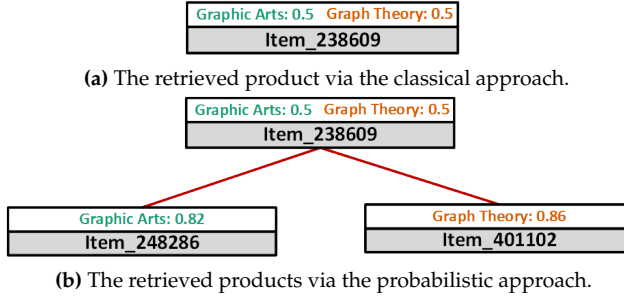


Fig. S4. A case study based on the Amazon dataset.

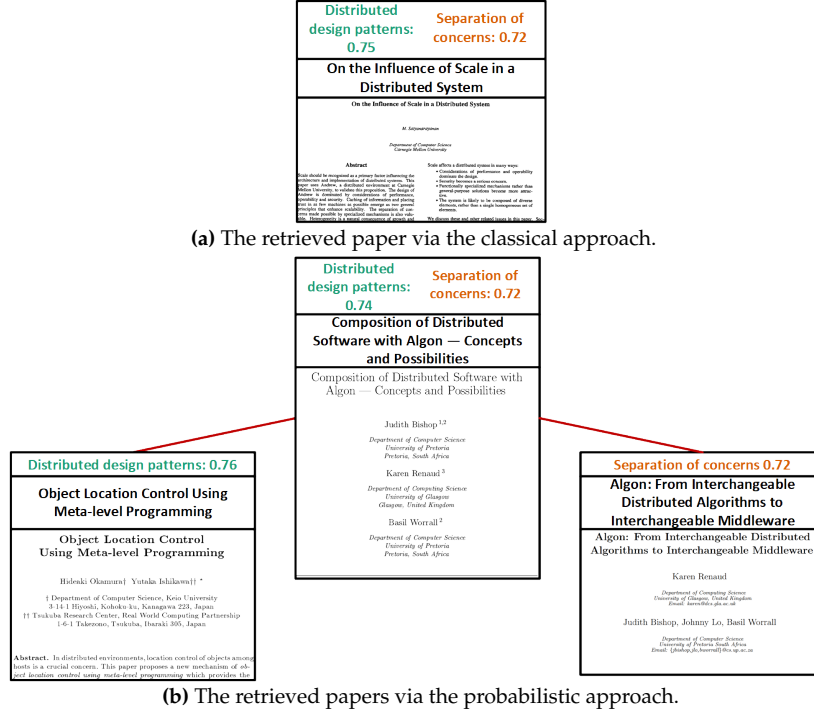


Fig. S5. A case study based on the DBLP dataset.

## S6. ADDITIONAL CASE STUDIES

In the main contents, we conduct a case study using the Movie dataset. In addition, here, we conduct similar case studies using the Amazon and DBLP datasets.

For Amazon, suppose that a user queries two keywords: {Graphic Arts, Graph Theory}. To help this user find related products, we input two corresponding vertex groups, and apply the existing

ImprovAPP and the proposed GRE-PATH to solve the classical and the probabilistic group Steiner tree problems, respectively. The existing ImprovAPP does not consider the probabilities of vertices for covering groups, and returns the single product "Item\_238609" in Figure S4a. This product has a probability of 0.5 of satisfying the user for each input keyword. This probability is smaller than the default threshold value  $b = 0.9$ . In comparison, the proposed GRE-PATH considers the probabilities of vertices for covering groups, and returns three related products in Figure S4b. These three products collectively have a probability of  $1 - (1 - 0.82) \cdot (1 - 0.5) = 0.91$  of satisfying the user for the keyword "Graphic Arts", and a probability of  $1 - (1 - 0.86) \cdot (1 - 0.5) = 0.93$  of satisfying the user for the keyword "Graph Theory". These two probabilities are larger than  $b$ . Thus, GRE-PATH retrieves closely related products that are more likely to satisfy the user.

For DBLP, suppose that a user queries two topics: {Distributed design patterns, Separation of concerns}. To help this user find related papers, we input two corresponding vertex groups, and apply the existing ImprovAPP and the proposed GRE-PATH to solve the classical and the probabilistic group Steiner tree problems, respectively. The existing ImprovAPP returns the single paper "On the Influence of Scale in a Distributed System" in Figure S5a. This paper has a probability of 0.75 of being in the field of "Distributed design patterns", and a probability of 0.72 of being in the field of "Separation of concerns". These two probabilities are smaller than  $b$ . In comparison, the proposed GRE-PATH returns three related papers in Figure S5b. These three papers collectively have a probability of  $1 - (1 - 0.76) \cdot (1 - 0.74) = 0.9376$  of being in the field of "Distributed design patterns", and a probability of  $1 - (1 - 0.72) \cdot (1 - 0.72) = 0.9216$  of being in the field of "Separation of concerns". These two probabilities are larger than  $b$ . Thus, GRE-PATH retrieves closely related papers that are more likely to be in the queried fields.

Like the Movie case study in the main contents, the above Amazon and DBLP case studies show that, in probabilistic scenarios, we could retrieve information that is more likely to be favorable via the probabilistic group Steiner tree approach than via the classical group Steiner tree approach.

## S7. THE METHODS OF SELECTING VERTEX GROUPS IN EXPERIMENTS

In the main experiments, we select  $|\Gamma|$  groups in such a way that the PoIs corresponding to the selected groups are often related and may appear together in practice. Since a Breadth First Search is conducted in this selection method, we refer to this selection method as the BFS selection method. A naive comparison of this method is the Simple Random selection method, *i.e.*, selecting  $|\Gamma|$  vertex groups uniformly at random. Notably, each group corresponds to a PoI. Here, we show some examples of the selected PoIs using the BFS and the Simple Random selection methods.

Specifically, in Tables S1 and S2, we show some examples of the selected PoIs using the BFS and the Simple Random selection methods, respectively, for Amazon. It can be seen from Table S1 that, when using the BFS selection method, the selected PoIs are often related and may occur together in practice, such as the five PoIs in the No.8 row in Table S1: "Thiebaud, Wayne | Realism | Painting | Art | Arts & Photography" (notably, Wayne Thiebaud is an American painter; [https://en.wikipedia.org/wiki/Wayne\\_Thiebaud](https://en.wikipedia.org/wiki/Wayne_Thiebaud)). Comparatively, it can be seen from Table S2 that, when using the Simple Random selection method, the selected PoIs are often unrelated and may rarely occur together in practice, *e.g.*, the two PoIs "Iranian Cinema" and "Composition & Creative Writing" in the No.9 row in Table S2 are barely related and rarely occur together in practice. We can get similar comparison results by analyzing the examples of the selected PoIs for DBLP in Tables S3-S4. Such comparison results are not so obvious for Movie in Tables S5-S6, due to the reason that there are only 19 PoIs in the Movie dataset, and most pairs of these PoIs may occur together in practice. Nevertheless, the above examples of the selected PoIs indicate that the BFS selection method is more practical than the Simple Random selection method, which is the reason why we use the BFS selection method in the experiments in the main contents.

## S8. THE MEMORY CONSUMPTION OF ALGORITHMS IN EXPERIMENTS

In Figures S6-S7, we report the memory consumption of algorithms in the main experiments. The reported memory consumption does not contain the memory consumed by the common inputs of algorithms, *i.e.*, the input graph  $G$  and the input set of vertex groups  $\Gamma$ , but contains all the other memory consumed in the computing process.

First, we report the memory consumption of algorithms in graphs with different sizes in Figure S6, which corresponds to Figure 4 in the main contents. We also report the memory consumption of hub labels of shortest paths in graphs in Figure S6. Both GRE-PATH and the baseline algorithms use these hub labels to retrieve the shortest paths to be merged. In the full Amazon and DBLP

graphs, the memories consumed by hub labels are much higher than the memories consumed by different algorithms, while in the full Movie graph, the memory consumed by hub labels is slightly smaller than that consumed by ENSteiner. The reason is that the Movie graph is much denser than the Amazon and DBLP graphs, and contains a large number of edges. These edges do not enlarge the size of hub labels significantly, due to the pruned landmark nature of the applied Pruned Landmark Labeling algorithm [6] for generating these labels. In comparison, these edges enlarge the memories consumed by ENSteiner significantly, since ENSteiner builds and records a graph by adding dummy vertices and edges into the input graph  $G$  in the computing process. Notably, as discussed in the main contents, we can also use different methods [9, 10] to generate hub labels with smaller sizes, at the cost of lower efficiencies of querying shortest paths using hub labels. We refer to [6, 9, 10] for more details of generating hub labels, and [11, 12] for the methods of updating hub labels for querying shortest paths in dynamic graphs.

We further observe that the memories consumed by DUAL and GRE-TREE are often larger than the memories consumed by PrunedDP++. This indicates that the implementations of PrunedDP++ in Line 11 of DUAL and Line 7 of GRE-TREE often consume larger amounts of memories than the baseline implementation of PrunedDP++, since PrunedDP++ is often used in Line 11 of DUAL and Line 7 of GRE-TREE to connect more groups, *e.g.*, in the initial implementation of PrunedDP++ in Line 7 of GRE-TREE,  $|\Gamma'|$  often equals  $|\Gamma| + 1$  and thus is often larger than  $|\Gamma|$ .

We also observe that, the memories consumed by GRE-PATH are always smaller than the memories consumed by the other algorithms, since it has a simple process of pushing  $O(|\Gamma||V|)$  values into priority queues and then popping out these values. In comparison, the other algorithms have more complex processes. For example, DPBF, PrunedDP++, DUAL and GRE-TREE incorporate dynamic programming processes to obtain classical group Steiner trees, and may record an exponential number of trees in the dynamic programming processes, and thus consume more memories than GRE-PATH. We further note that ENSteiner may consume more memories than the other algorithms in some cases, *e.g.*, in Figure S6c (3). The reason is that, different from the other algorithms, ENSteiner builds and records a graph by adding dummy vertices and edges into the input graph  $G$ . When  $G$  is large, such as the case in Figure S6c (3), the built graph in ENSteiner consumes a lot of memory.

Subsequently, we report the memory consumption of algorithms when varying parameters in Figure S7 (notably, the memory consumption of hub labels of shortest paths in graphs do not change with parameters; thus, we do not report the memory consumption of hub labels here). We observe that the memories consumed by DPBF and PrunedDP++ may increase significantly with  $|\Gamma|$ , *e.g.*, in Figure S7a (3). The reason is that the number of recorded trees in the dynamic programming processes of these two algorithms are exponential to  $|\Gamma|$ . Nevertheless, in Figure S7a (1), the memories consumed by PrunedDP++ and GRE-TREE, which incorporates PrunedDP++, do not increase with  $|\Gamma|$  as significantly as DPBF. This shows that the branch and bound techniques in PrunedDP++ effectively accelerate the dynamic programming process. On the other hand, we note that the memories consumed by different algorithms do not change with  $b$ ,  $\tau$ ,  $P_{min}$  and  $P_{max}$ . Moreover, we observe that the memories consumed by DPBF may increase with  $k$ , since more feasible solutions are computed with the increase of  $k$ .

## S9. THE SIZES OF $|V_g|$ AND $|g_{min}|$ IN EXPERIMENTS

We visualize the sizes of  $|V_g|$  and  $|g_{min}|$  in the main experiments in Figures S8-S9. In the full graphs in Figure S8c,  $|V_g|$  is an order of magnitude smaller than  $|V|$  for Amazon and DBLP, and is similar to  $|V|$  for Movie ( $|V|$  is 548552, 2497782 and 62423 for Amazon, DBLP and Movie, respectively). The reason is that the sizes of candidate groups are often small for Amazon and DBLP, but are often large for Movie, as shown in Figure 3 in the main contents. For a similar reason,  $|g_{min}|$  is larger for Movie than for Amazon and DBLP. On the other hand, in Figure S9,  $|V_g|$  and  $|g_{min}|$  do not change much with  $b$ ,  $\tau$ ,  $P_{min}$ ,  $P_{max}$  and  $k$ . In comparison, in Figure S9a,  $|V_g|$  increases with  $|\Gamma|$ , while  $|g_{min}|$  decreases with  $|\Gamma|$ . The reason is that, as  $|\Gamma|$  increases, more vertices are in the selected groups, and smaller groups are more likely to be selected.

## S10. ADDITIONAL EXPERIMENT RESULTS WITH PAIRWISE JACCARD DISTANCES

In the main experiments, we set edge weights to 1. Here, we conduct additional experiments by setting edge weights to pairwise Jaccard distances (*e.g.*, [8, 13]), *i.e.*, for edge  $e$  between vertices  $u$  and  $v$ , set the weight of  $e$  as  $c(e) = 1 - \frac{|V_u \cap V_v|}{|V_u \cup V_v|}$ , where  $V_u$  and  $V_v$  are the sets of vertices adjacent to  $u$  and  $v$ , respectively. We show that the key observations in the main experiments still hold here.



**DUAL is mainly of theoretical interests.** First, we show that DUAL can only be used in tiny graphs with dozens of vertices in Figure S10a, where  $|V| = 45$  for "Tiny Amazon",  $|V| = 90$  for "Tiny DBLP", and  $|V| = 70$  for "Tiny Movie". We observe that, in Figures S10a (4-6), DUAL is significantly slower than the other algorithms. These experiments show that DUAL can only be used in tiny graphs with dozens of vertices, and is mainly of theoretical interests.

**GRE-TREE is useful when group sizes are small.** We evaluate the performance of GRE-TREE in Figure S10b, where  $|V| = 188,552$  for "Small Amazon",  $|V| = 448,891$  for "Small DBLP", and  $|V| = 2,423$  for "Small Movie". We observe that GRE-TREE is significantly slower than the other algorithms for Movie, since group sizes are large for Movie. Like the main experiments, GRE-TREE can produce better solutions than the other algorithms. Thus, it may be preferable to use GRE-TREE when group sizes are small and graph sizes are not extremely large, *e.g.*, for Amazon.

**Experiment results in full graphs.** We evaluate the solution quality and speed of algorithms using the full datasets in Figure S10c. We observe that, in Figures S10c (1-3), the solution weights of GRE-TREE and GRE-PATH are significantly lower than those of the baseline algorithms. This shows the effectiveness of GRE-TREE and GRE-PATH for finding probabilistic group Steiner trees.

In conclusion, the key observations in the main experiments, *i.e.*, (i) DUAL can only be used in tiny graphs with dozens of vertices; (ii) GRE-TREE produces better solutions than the other algorithms in practice and is efficient when group sizes are small; and (iii) GRE-PATH produces considerably better solutions than baselines and scales well to large graphs, still hold here.

## REFERENCES FOR THE SUPPLEMENT

1. "Supplement," (2022). <https://github.com/rucdatascience/PGST/blob/main/Supplement.pdf>.
2. B. Ding, J. X. Yu, S. Wang, L. Qin, X. Zhang, and X. Lin, "Finding top-k min-cost connected trees in databases," in *IEEE International Conference on Data Engineering*, (IEEE, 2007), pp. 836–845.
3. R.-H. Li, L. Qin, J. X. Yu, and R. Mao, "Efficient and progressive group Steiner tree search," in *Proceedings of the 2016 International Conference on Management of Data*, (ACM, 2016), pp. 91–106.
4. R. M. Karp, "Reducibility among combinatorial problems," in *Complexity of computer computations*, (Springer, 1972), pp. 85–103.
5. E. Cohen, E. Halperin, H. Kaplan, and U. Zwick, "Reachability and distance queries via 2-hop labels," *SIAM Journal on Computing* **32**, 1338–1355 (2003).
6. T. Akiba, Y. Iwata, and Y. Yoshida, "Fast exact shortest-path distance queries on large networks by pruned landmark labeling," in *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, (2013), pp. 349–360.
7. Y. Li, L. H. U, M. L. Yiu, and N. M. Kou, "An experimental study on hub labeling based shortest path algorithms," *Proc. VLDB Endow.* **11**, 445–457 (2017).
8. Y. Sun, X. Xiao, B. Cui, S. Halgamuge, T. Lappas, and J. Luo, "Finding group steiner trees in graphs with both vertex and edge weights," *Proc. VLDB Endow.* **14**, 1137–1149 (2021).
9. W. Li, M. Qiao, L. Qin, Y. Zhang, L. Chang, and X. Lin, "Scaling distance labeling on small-world networks," in *Proceedings of the 2019 International Conference on Management of Data*, (2019), pp. 1060–1077.
10. W. Li, M. Qiao, L. Qin, Y. Zhang, L. Chang, and X. Lin, "Scaling up distance labeling on graphs with core-periphery properties," in *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, (2020), pp. 1367–1381.
11. D. Ouyang, L. Yuan, L. Qin, L. Chang, Y. Zhang, and X. Lin, "Efficient shortest path index maintenance on dynamic road networks with theoretical guarantees," *Proc. VLDB Endow.* **13**, 602–615 (2020).
12. M. Zhang, L. Li, W. Hua, and X. Zhou, "Efficient 2-hop labeling maintenance in dynamic small-world networks," in *2021 IEEE 37th International Conference on Data Engineering*, (IEEE, 2021), pp. 133–144.
13. T. Lappas, K. Liu, and E. Terzi, "Finding a team of experts in social networks," in *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, (ACM, 2009), pp. 467–476.

**Table S1.** Some examples of the selected PoIs (Amazon; the BFS selection method).

No.	The selected PoIs ( $ G  = 5$ )
1	Middle Eastern   Ancient   Figure Drawing   History of Religion   Bible
2	Toymaking   Projects   Crafts & Hobbies   Europe   Grandparenting
3	Pop Culture   South Atlantic   Virgin Islands   Boating   Hiking
4	Pulver, Liselotte   Satire   Specialty Stores   Wilder, Billy   Billy Wilder
5	Karaoke   Girl Groups   Music Video & Concerts   General Christmas   Music Outlet
6	Insurance Law   Administration & Policy   Business & Investing Books   Business   Nutshell
7	Rukeyser, Muriel   Millay, Edna St Vincent   Poetry, Drama & Short Stories   Poets, A-Z   Books on Tape
8	Thiebaud, Wayne   Realism   Painting   Art   Arts & Photography
9	Home Care   Medical   Hospice Care   Pharmacy   Spirituality
10	Rio de Janeiro   Travel   History   Social Sciences   Latin America



**Table S2.** Some examples of the selected PoIs (Amazon; the Simple Random selection method).

No.	The selected PoIs ( $ \Gamma  = 5$ )
1	Superman   Johnston, Aaron Kim   Stamper, J. B.   Hillenbrand, Will   Languages & Tools
2	Epistemology   McGuire, Christine   Mcquary, Chuck   Kober, Jeff   Gilded Age
3	Equipment   Balaban, Bob   Chicago Symphony Orchestra   Sirtis, Marina   Solti, Georg
4	MacDonald, Betty   Brooke, Hillary   Anatomy   Winston, Don   Parent Participation
5	Dorsey, Jimmy   Graduate Preparation   Lowry, Morton   Schrader, Paul   Data Structures
6	For Seniors   Finch, Jon   Languages & Tools   Sutras   Crawford, Broderick
7	Networking & System Administration   Yourcenar, Marguerite   Talbot, Nita   Kaiser, Walter   Overman, Lynne
8	Database Storage & Design   Yanne, Jean   Clásicos   Test Guides - High School   Millay, Edna St Vincent
9	Iranian Cinema   Composition & Creative Writing   Redding, Otis   Mutter, Scott   Browne, Roscoe Lee
10	Fertility Books   Ziglar, Zig   Morwood, Peter   Smith, Dan Warry   Torres, Liz

**Table S3.** Some examples of the selected PoIs (DBLP; the BFS selection method).

No.	The selected PoIs ( $ \Gamma  = 5$ )
1	Prefix   Binary prefix   Binary search algorithm   Encoding (memory)   Parallel computing
2	Solar battery   Distance-vector routing protocol   High voltage   Collinear antenna array   Renewable energy
3	Writing style   Personal development   Comprehension   Structure (mathematical logic)   Sentence
4	Fading   Electromagnetic reverberation chamber   Channel capacity   Reverberation room   Antenna (radio)
5	Psychological intervention   Forensic psychiatry   Anger   Computer security   Artificial intelligence
6	Good clinical practice   Health care   Medicine   Terminology   Torture
7	Unimodular lattice   Ring of integers   Lattice (group)   Leech lattice   Mass formula
8	Performance indicator   Self-organizing network   Integer programming   Supply chain network   Maintainability
9	Hardware obfuscation   Gate array   Identifier   Matrix decomposition   Lock (computer science)
10	Dissolved organic carbon   Downwelling   Biochemical oxygen demand   Satellite imagery   Artificial neural network

**Table S4.** Some examples of the selected PoIs (DBLP; the Simple Random selection method).

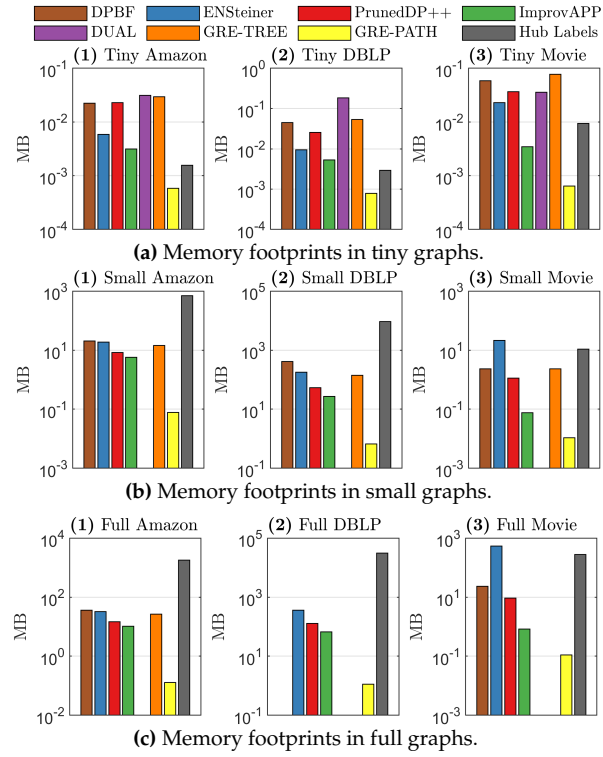
No.	The selected PoIs ( $ \Gamma  = 5$ )
1	Flexural strength   Documentary evidence   Radiation mode   Voltage compensation   Spectrum auction
2	Polytropic process   Guanidine   Environmental full-cost accounting   The Imaginary   Selective sweep
3	Food group   Frontal lobe   Psychological contract   Tax reform   Ammonia
4	Financialization   Virtual work   Organic search   Fibre optic gyroscope   Prime number theorem
5	Threading (protein sequence)   Junior school   Exome   Ultra-large-scale systems   Fuzzy classification
6	Hybrid security   Zonal polynomial   Pharmacology   Gas chromatography   Fermi level
7	CAAT box   Color filter array   Weather and climate   Individual mobility   Germanium
8	Inertial measurement unit   Streaking   Traditional education   Bispectrum   Honeycomb structure
9	Clicker   Japanese chess   Token passing   T/TCP   Data center network architectures
10	Toothbrush   Managed care   Audio power   Contraction mapping   Index set

**Table S5.** Some examples of the selected PoIs (Movie; the BFS selection method).

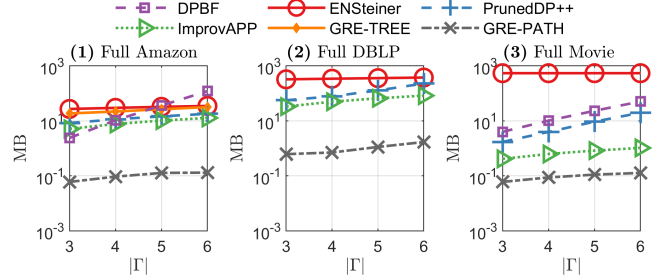
No.	The selected PoIs ( $ \Gamma  = 5$ )
1	Romance   Thriller   Western   Musical   Drama
2	Action   Comedy   Drama   Animation   Crime
3	Film-Noir   Action   Animation   Thriller   Drama
4	Musical   Drama   Romance   Action   Western
5	Crime   Comedy   Thriller   Children   Action
6	Film-Noir   Action   Musical   Crime   Western
7	Drama   Comedy   Western   Film-Noir   Crime
8	Animation   Thriller   Comedy   Romance   Drama
9	Crime   Drama   Film-Noir   Romance   Animation
10	Children   Horror   Action   Drama   Comedy

**Table S6.** Some examples of the selected PoIs (Movie; the Simple Random selection method).

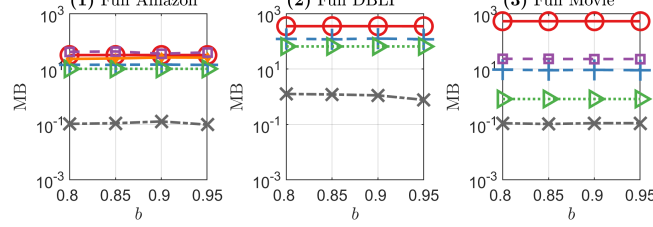
No.	The selected PoIs ( $ \Gamma  = 5$ )
1	Crime   Thriller   Animation   Western   Children
2	Drama   Horror   Film-Noir   Crime   Western
3	Musical   Western   Action   Children   Horror
4	Horror   Thriller   Animation   Western   Drama
5	Thriller   Western   Comedy   Crime   Children
6	Horror   Film-Noir   Comedy   Drama   Thriller
7	Drama   Animation   Thriller   Film-Noir   Horror
8	Thriller   Animation   Western   Romance   Action
9	Children   Horror   Animation   Thriller   Crime
10	Action   Comedy   Romance   Musical   Animation



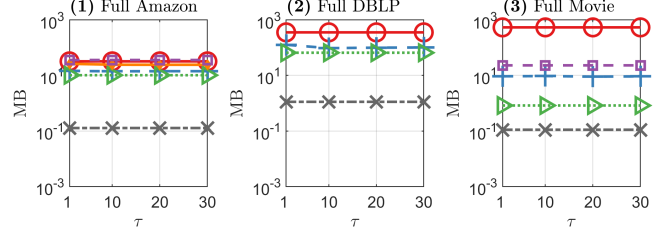
**Fig. S6.** Memory footprints in graphs with different sizes.



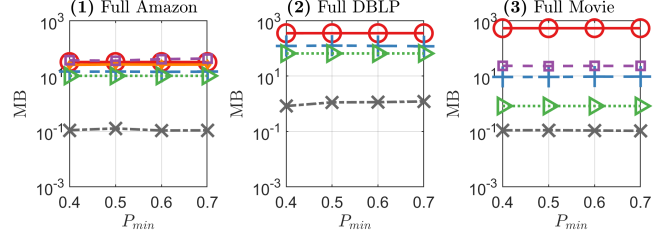
(a) Variation of the number of vertex groups:  $|\Gamma|$ .



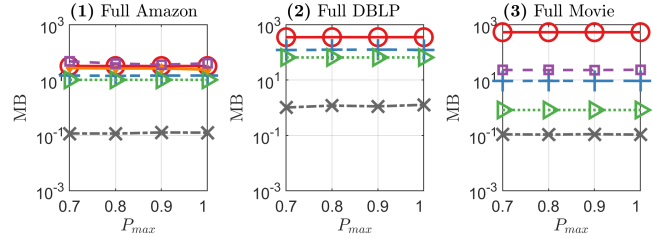
(b) Variation of the threshold value:  $b$ .



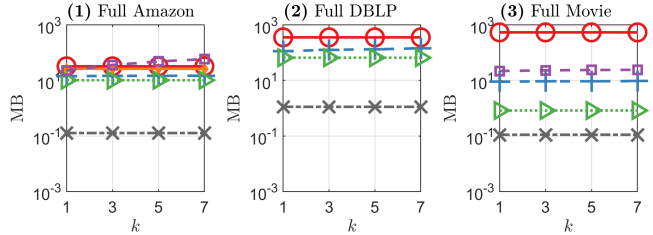
(c) Variation of the approximation parameter:  $\tau$ .



(d) Variation of the minimum positive probability value:  $P_{min}$ .

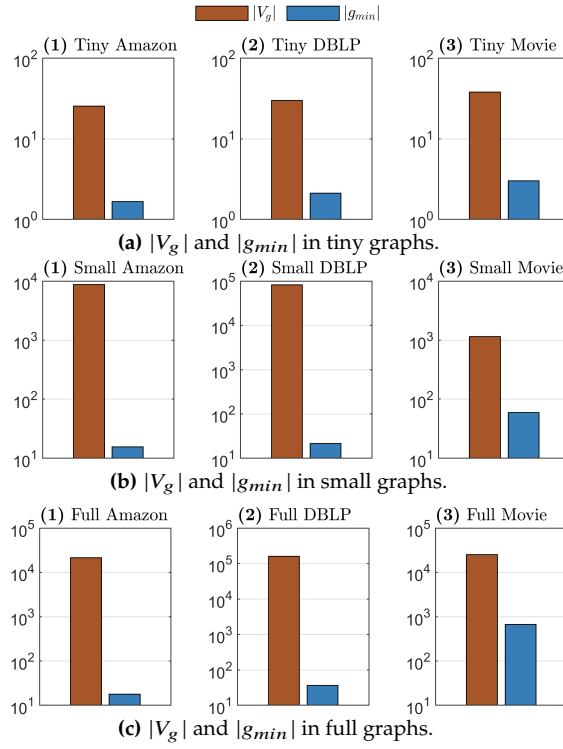


(e) Variation of the maximum positive probability value:  $P_{max}$ .

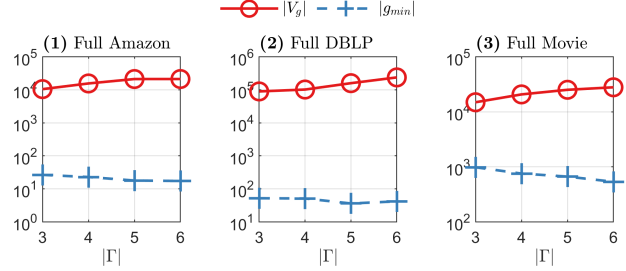


(f) Variation of the number of feasible solutions in baselines:  $k$ .

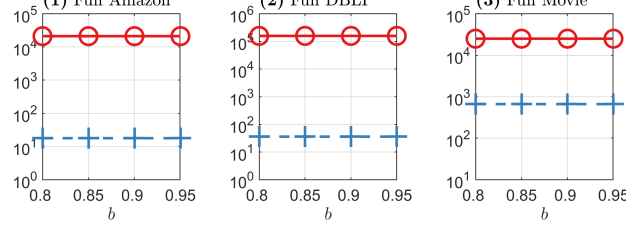
**Fig. S7.** Memory footprints of varying parameters.



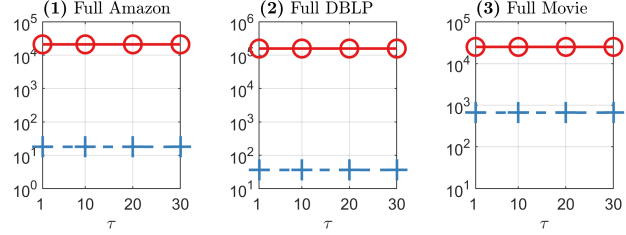
**Fig. S8.**  $|V_g|$  and  $|g_{min}|$  in graphs with different sizes.



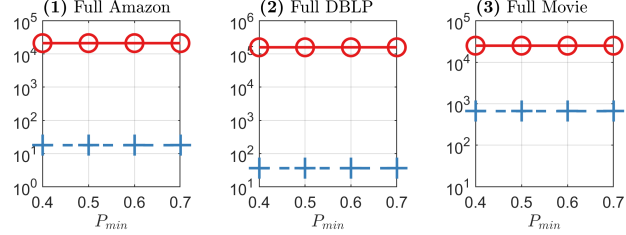
(a) Variation of the number of vertex groups:  $|\Gamma|$ .



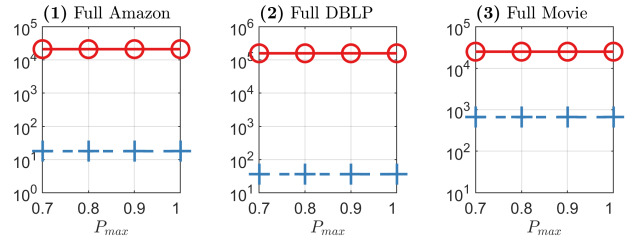
(b) Variation of the threshold value:  $b$ .



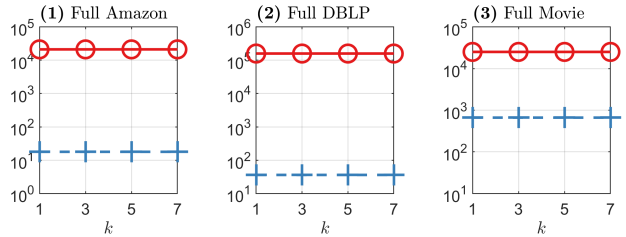
(c) Variation of the approximation parameter:  $\tau$ .



(d) Variation of the minimum positive probability value:  $P_{min}$ .

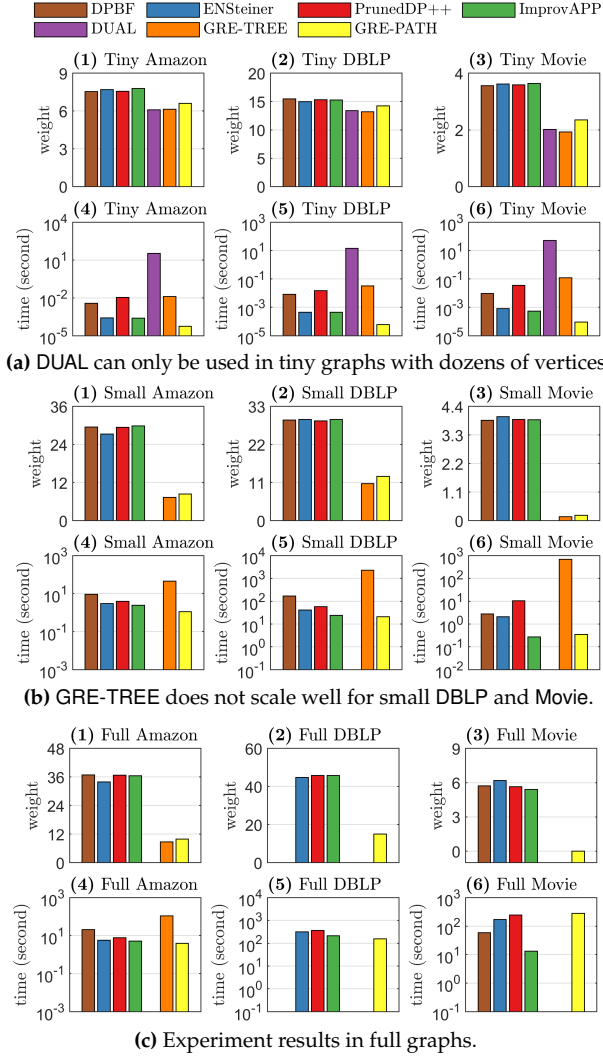


(e) Variation of the maximum positive probability value:  $P_{max}$ .



(f) Variation of the number of feasible solutions in baselines:  $k$ .

Fig. S9.  $|V_g|$  and  $|g_{min}|$  of varying parameters.



**Fig. S10.** Experiment results in graphs with different sizes (pairwise Jaccard distances).