

# gIM: GPU Accelerated RIS-based Influence Maximization Algorithm

Soheil Shahrouz, Saber Salehkaleybar, and Matin Hashemi

This article is published. Please cite as S. Shahrouz, S. Salehkaleybar, M. Hashemi, "gIM: GPU Accelerated RIS-based Influence Maximization Algorithm," IEEE Transactions on Parallel and Distributed Systems.

**Abstract**—Given a social network modeled as a weighted graph  $G$ , the influence maximization problem seeks  $k$  vertices to become initially influenced, to maximize the expected number of influenced nodes under a particular diffusion model. The influence maximization problem has been proven to be NP-hard, and most proposed solutions to the problem are approximate greedy algorithms, which can guarantee a tunable approximation ratio for their results with respect to the optimal solution. The state-of-the-art algorithms are based on Reverse Influence Sampling (RIS) technique, which can offer both computational efficiency and non-trivial  $(1 - 1/e - \epsilon)$ -approximation ratio guarantee for any  $\epsilon > 0$ . RIS-based algorithms, despite their lower computational cost compared to other methods, still require long running times to solve the problem in large-scale graphs with low values of  $\epsilon$ . In this paper, we present a novel and efficient parallel implementation of a RIS-based algorithm, namely IMM, on GPU. The proposed GPU-accelerated influence maximization algorithm, named gIM, can significantly reduce the running time on large-scale graphs with low values of  $\epsilon$ . Furthermore, we show that gIM algorithm can solve other variations of the IM problem, only by applying minor modifications. Experimental results show that the proposed solution reduces the runtime by a factor up to  $220\times$ . The source code of gIM is publicly available online.

**Index Terms**—CUDA, GPGPU, Graph Diffusion Process, Influence Maximization (IM), Parallel Processing, Reverse Influence Sampling (RIS).



## 1 INTRODUCTION

WITH the explosion of social network services in the last decades, hundreds of millions of people can easily interact with each other. The vast prevalence of social networks facilitates large-scale viral marketing through "word-of-mouth" effects, where each individual recommends a product to his/her friends, contributing to a large number of adoptions of the product. In order to succeed in a viral marketing campaign, it is required to choose a few *influential* individuals and provide incentives (e.g. free samples of the product and/or cash) to create a cascade of product adoptions as large as possible. Obtaining such a set of users is commonly known as the "influence maximization (IM) problem". In particular, for a given network  $G$  and a probability model representing diffusion mechanism in the network, IM aims to select a set of users (called seed set) which maximizes the expected number of users that are affected in the diffusion process.

In a seminal work, Kempe et al. [1] formulated the IM problem as a combinatorial optimization problem for two popular diffusion models, namely, independent cascade (IC) and linear threshold (LT). They showed that the IM problem is NP-hard in both models, and proposed a greedy framework that can obtain  $(1 - \frac{1}{e} - \epsilon)$ -approximate solutions for any  $\epsilon > 0$ . More specifically, they considered the influence spread (the expected number of affected users) as a function of seed set. Starting from an empty set, a user with the

maximum marginal gain to the influence spread is added to the seed set in an iterative manner until the desired number of users are selected. Almost all the proposed algorithms for the IM problem followed this greedy framework. Later, it has been shown that computing the influence spread as a function of seed set is #P-hard under both IC and LT models [2], [3]. These theoretical results have motivated many researchers to design efficient estimators for the influence spread. In this regard, existing algorithms for the IM problem can be classified into three main approaches: I) simulation-based approach, II) proxy-based approach, and III) sketch-based approach.

In the simulation-based approach, Monte Carlo (MC) simulations are executed to estimate the influence spread for a given seed set by averaging the number of affected users over all the generated instances. This approach was considered in [1] for the first time and subsequent works tried to reduce the number of MC simulations through *lazy evaluation* [4], [5]. The main advantage of this approach is that it can be utilized for any diffusion model. However, it is not scalable to very large graphs since it requires to generate too many instances to compute the influence spread with a desirable estimation error. In the proxy-based approach, the main idea is to use simple models, such as shortest path or PageRank, instead of complicated diffusion models. This reduces the runtime for computing the influence spread significantly. However, improving the time complexity of evaluating the influence spread comes at the expense of losing theoretical guarantees on estimation error bounds. In fact, it has been shown that the proxy-based approach might have unstable behavior in some graphs in the sense that the influence spread could change significantly by a small

Authors are with Learning and Intelligent Systems Laboratory, Department of Electrical Engineering, Sharif University of Technology, Tehran, Iran.  
E-mails: soheil.shahrouz@ee.sharif.edu, saleh@sharif.edu, matin@sharif.edu.  
Webpage: <http://lis.ee.sharif.edu>  
(Corresponding author: Saber Salehkaleybar)

modification in the graph structure [6]. The sketch-based approach has been proposed to overcome computational inefficiency issues of the simulation-based approach and instability problems of the proxy-based approach. The main idea is to first generate *sketches* for a given diffusion model. Afterwards, the influence spread is estimated based on the generated sketches. Due to the desired properties of this approach, many recent algorithms for the IM problem have focused on efficiently generating sketches (see Section 2).

Although the most recent algorithms for the IM problem are mainly based on constructing sketches, they still have computational issues for small estimation error bounds in large graphs. Moreover, this problem might become critical in some recent applications of IM such as multi-round IM [7], where it is required to solve the IM problem multiple times for a given graph.

The aim of this paper is to accelerate the state-of-the-art sketch-based algorithm (IMM [8]) using parallel processing on GPU. Hence, the proposed solution is named gIM. Our main contributions are as follows:

- One of the key parts in IMM algorithm is to find a reverse reachable set of nodes from a randomly selected node by executing a BFS algorithm. This part of the algorithm is repeated multiple times. Many common approaches for the acceleration of BFS on GPU run only one BFS at a time, and therefore, produce a single reverse reachable set. Here, we propose a method to simultaneously generate multiple reverse reachable sets, and at the same time, parallelize each task of generating a reverse reachable set as well. More specifically, several BFSs are executed in parallel, and in each one, adjacent edges of a node are processed in parallel, but nodes in the frontier queue are processed sequentially (Section 3.2).
- We also propose methods to judiciously store large frontier queues of the BFSs in limited GPU shared memory in order to further reduce the runtime (Section 3.3).
- We propose data structures and methods to optimally store the sketches in GPU global memory (Section 3.5) and to process them efficiently in order to select the seed set (Section 3.8). This helps to fit the required internal data structures into GPU global memory for large graphs with millions of nodes.
- We discuss how to optimize the required synchronization mechanism and also the number of threads and blocks in order to exploit all processing power of GPU (Section 3.6).
- The proposed parallelization of IMM algorithm on GPU can be utilized in other applications that use a variant of IMM as a subroutine. For instance, we show that our implementation can be adapted to the multi-round IM problem. The experiments show that running times are greatly reduced (Section 4.8).
- Experimental results on real social networks show that the proposed gIM algorithm can outperform state-of-the-art solutions by a factor of up to  $220\times$ .

In the following, we briefly review some of the most prominent works in the literature, which focused on parallel implementation of the IM algorithm. Liu et al. [9]

presented the IMGPU framework to accelerate the basic greedy algorithm with MC simulations on GPU. They first converted each instance graph to a DAG to avoid redundant graph traversals in the computations of each node’s influence spread. Next, they designed a bottom-up traversal algorithm to compute the marginal gain of all nodes in the resulting DAG on GPU hardware. IMGPU reduces the execution time of the greedy algorithm by a factor of  $60\times$  but it still needs tens of minutes to solve the problem in relatively large graphs.

Pal et al. [10] focused on GPU acceleration of the greedy algorithm with MC simulations under continuous-time diffusion models. They assumed that each node’s influence spread is limited to its local neighborhood and restricted the graph traversal to this small neighborhood. They adopted a node-level parallelization approach for calculating marginal gains. Despite achieving outstanding speedups, the limited graph traversal degrades the seed set quality, and hence, this method does not reach the same influence spread as the state-of-the-art algorithms.

Minutoli et al. [11] proposed Ripples, a framework for multi-core and distributed implementation of the IMM algorithms. Later, by introducing cuRipples [12] framework, they extended their work to provide support for multi-GPU systems and employed GPUs along with multi-core CPUs to further reduce the execution time.

Gokturk et al. [13] accelerated simulation-based algorithms on multi-core CPUs by employing a direction-oblivious hash-based sampling approach to fuse instance graph generation and marginal gain calculation. They also exploited SIMD vectorization to run multiple concurrent simulations on a single core. However, direction-oblivious sampling causes their algorithm to be applicable only on undirected graphs.

Nguyen et al. [14] proposed a CPU-GPU method to accelerate Spread Interdiction (SI) problem which is similar to the IM problem. Although they used concepts similar to reverse influence sampling, their proposed method is not directly applicable to the IM problem.

In a parallel line of research, there have been efforts to make large-scale graph processing more tractable by graph sampling, that is, producing a graph of much smaller size while preserving the desirable characteristics of the original graph. Pandey et al. [15] proposed C-SAW framework to implement various graph sampling methods and paralleled them on GPU. Zeng et al. [16] introduced GraphSAINT, a graph sampling method to train graph convolutional networks (GCN), which paves the way for training GCNs on large graphs and thus opens the opportunity to use GCN for solving the IM problem.

The rest of this paper is organized as follows: Section 2 reviews the preliminaries. Section 3 presents our proposed GPU-accelerated algorithm for the IM problem, named gIM. Section 4 presents experiment results and comparisons with related works. Section 5 concludes the paper.

## 2 PRELIMINARIES

In this section, we introduce some notations and investigate two common diffusion models, namely independent cascade and linear threshold models. Then we formally define

the influence maximization problem and review Kempe’s greedy algorithm [1] and methods based on reverse influence sampling.

## 2.1 Notation

Directed graph  $G = (V, E)$  models a social network, where  $V$  denotes the set of vertices (i.e., users) and  $E$  indicates the set of directed edges (i.e., relationships between users). For every edge  $e = (u, v) \in E$ , we define  $u$  as an incoming neighbor of  $v$ , and  $v$  as an outgoing neighbor of  $u$ . Also, we denote sets of **incoming** and **outgoing** neighbors of every node  $u$ , by  $N_I(u)$  and  $N_O(u)$ , respectively. Moreover, we assume that every edge  $e = (u, v) \in E$  is associated with an **influence probability**  $p_{uv} \in [0, 1]$ . We denote the total number of nodes and edges in  $G$  by  $n$  and  $m$ , respectively.

## 2.2 Diffusion Models

There is an immense amount of literature on models that capture the diffusion phenomenon’s behavior [17], [18], [19], [20]. Among these models, independent cascade (IC) and linear threshold (LT) models have been extensively studied. In the following, we briefly discuss these models.

### *Independent Cascade (IC) Diffusion Model*

Given a graph  $G = (V, E)$ , edge influence probabilities  $p_{uv}$  for every edge  $e = (u, v)$ , and a subset of vertices  $S \subseteq V$ , which is commonly called **seed set**, an instance of influence diffusion process under the IC model is generated as follows:

The influence propagates in discrete time steps. At time  $t_0$ , only the nodes in the seed set  $S$  are activated, and all the other nodes remain inactive. At every time step  $t_{i+1}$ , every node  $u$ , which has been activated at time step  $t_i$ , has a chance to activate its inactive outgoing neighbors. Node  $u$  succeeds to activate an inactive outgoing neighbor  $v$ , with probability  $p_{uv}$ . When a node is activated (either at  $t_0$  or by an incoming active neighbor), it remains active in the following time steps. The influence diffusion process terminates at a time step at which none of the previously activated nodes can activate a new inactive node.

### *Linear Threshold (LT) Diffusion Model*

The LT model imposes a new restriction on edge influence probabilities, by which the sum of probabilities of incoming edges for every node  $u$ , must be less than or equal to 1. Moreover, in the LT model, every node  $u$  is associated with a threshold  $\tau_u$ , which is a real value in the range  $[0, 1]$ . Given a seed set  $S$ , a single diffusion process instance under LT model can be described as follows:

Similar to the IC model, the influence propagates in discrete time steps. At time step  $t_0$ , we assign a random number sampled uniformly from  $[0, 1]$  to each node  $u$ , as its threshold  $\tau_u$ . Furthermore, we initially activate every node  $u \in S$  and leave the other nodes as inactive. At time step  $t_{i+1}$ , every node that has been previously activated remains active, and every inactive node  $v$  becomes active when it satisfies the following condition

$$\sum_{u \in N_I(v)} p_{uv} \cdot \mathbb{1}_A(u) \geq \tau_v, \quad (1)$$

where  $A$  is the set of all activated nodes until some time step  $t_i$ , and  $\mathbb{1}_A(u)$  is indicator function which is equal to one if  $u \in A$ . The influence diffusion process stops at a time step at which no inactive node can become active.

## 2.3 Problem Statement

The **influence spread** of seed set  $S$  in a single influence diffusion process instance is defined as the total number of active nodes after the termination of the diffusion process, and is denoted by  $I(S)$ . It should be noted that  $I(S)$  is not a deterministic function of  $S$ . The influence diffusion procedure is a random process which depends on edge influence probabilities (and also node thresholds in LT model). Therefore,  $I(S)$  is a random variable.

Given graph  $G$ , edge influence probabilities  $p_{uv}$ , a particular diffusion model, and a constant positive integer  $k$ , the goal of the influence maximization problem is to find seed set  $S^* \subset V$  of size  $k$ , which has the maximum **expected influence spread**  $\mathbf{E}[I(S)]$ . The problem can also be expressed as follows:

$$S^* = \operatorname{argmax}_{S \subset V} \{\mathbf{E}[I(S)] \mid |S| = k\}, \quad (2)$$

where  $S^*$  is the optimal seed set with size  $k$ . It has been shown [1] that the IM problem is NP-Hard under both IC and LT diffusion models.

## 2.4 Basic Greedy Algorithm

Kempe et al. [1] proposed a greedy algorithm to solve the IM problem. This algorithm iteratively selects nodes for insertion into the seed set. First, the seed set  $S$  is initialized by an empty set. In every iteration, the node  $v$  whose insertion into the seed set  $S$  leads to the largest increase in  $\mathbf{E}[I(S)]$  is selected. The amount of increase is called **marginal gain**. The algorithm terminates after  $k$  iterations, i.e., when the size of seed set  $S$  reaches  $k$ .

The greedy algorithm has a simple structure, but its main drawback is that it requires  $\mathbf{E}[I(S \cup u)]$  to be computed in every iteration for all the remaining nodes. Unfortunately, evaluating  $\mathbf{E}[I(S)]$  has been proven to be #P-Hard [2]. To tackle this problem, Kempe et al. [1] proposed a method based on Monte Carlo (MC) simulations in order to estimate  $\mathbf{E}[I(S \cup u)]$ . In order to execute a single MC simulation under IC diffusion model, an **instance graph**  $g$  is created from  $G$ , by removing every edge  $(u, v) \in E$  by probability  $1 - p_{uv}$ . **Reachable set**  $R(S)$  is defined as the set of all nodes that are reachable from  $S$  in the resulted graph  $g$ . Kempe et al. have shown that  $\mathbf{E}[|R(S)|] = \mathbf{E}[I(S)]$ , and therefore, one can estimate  $\mathbf{E}[I(S)]$  by determining  $\mathbf{E}[|R(S)|]$  through running a large number of MC simulations and generating many instances of  $g$ . They also proved that by knowing the exact values of  $\mathbf{E}[I(S)]$ , greedy algorithm’s solution has an approximation ratio of  $(1 - \frac{1}{e})$ . However, in practice, the value of  $\mathbf{E}[I(S)]$  is estimated and the true approximation ratio is  $(1 - \frac{1}{e} - \epsilon)$ , in which  $\epsilon$  depends on both graph  $G$  and the total number of MC simulations. In [1], they did not provide any formal analysis on the number of required MC simulations in order to achieve  $(1 - \frac{1}{e} - \epsilon)$  approximation ratio and only suggested using 10,000 MC simulations. Several following works also used 10,000 as the number of MC

simulations without providing any formal analysis. Later, Chen et al. [21] conducted an analysis on the total number of MC simulations and relative error  $\epsilon$ . They showed that to achieve  $(1 - \frac{1}{e} - \epsilon)$  approximation ratio under both IC and LT diffusion models, with a probability of  $1 - \frac{1}{n}$ , one needs to execute  $\Theta(\epsilon^{-2}k^2n \log(n^2k))$  number of MC simulations. Therefore, lower values of  $\epsilon$  requires larger number of MC simulations.

## 2.5 RIS-based Algorithms

As stated above, in every iteration, the greedy algorithm evaluates the marginal gain (i.e., the improvement in expected influence spread) resulting from adding one node which is not currently in the seed set. Thus, the number of marginal gain evaluations is  $O(n)$ , and because there are  $k$  iterations in total, the greedy algorithm conducts  $O(nk)$  marginal gain evaluations. Every marginal gain evaluation involves calculating  $\mathbf{E}[|R(S)|]$ , which is a computationally expensive task. In addition, in every iteration, we are only interested in the node with the largest marginal gain, therefore, marginal gain evaluations for all the other nodes are virtually wasted and are not useful in subsequent iterations.

To address this problem, Borgs et al. [22] proposed a method called **Reverse Influence Sampling (RIS)**, by which there is no need to run many MC simulations for every remaining node in every iteration in order to evaluate the marginal gain.

In order to describe RIS-based algorithms, we need the following definitions. Given a graph  $G$  and edge influence probabilities  $p_{uv}$  for every edge  $e = (u, v) \in E$ , one can generate an instance graph  $g$  from  $G$ , by removing every edge  $e = (u, v)$  by probability  $1 - p_{uv}$ . Let  $v$  be a node in graph  $G$ , then **reverse reachable (RR) set** for  $v$  in  $g$  is defined as the set of all nodes in  $g$  that can reach  $v$ . An RR set is defined to be a **random RR set** if instance graph  $g$  is created from  $G$  as described before, and node  $v$  is selected randomly from the set of nodes  $V$ .

In other words, every node  $u$  that is present in an RR set generated for a particular node  $v$ , has an opportunity to activate  $v$ , if it has been activated already. Borgs et al. [22] showed that the probability by which a node  $u$  can appear in a random RR set is proportional to the expected influence spread of "that" node. Moreover, they proved that the probability by which a particular seed set  $S$ , overlaps a random RR set can be used as an unbiased estimator of expected influence spread of the seed set, as given in the following equation:

$$\mathbf{E}[I(S)] = n \times Pr[S \cap RR \neq \emptyset]. \quad (3)$$

According to above results, Borgs et al. [22] proposed an algorithm to solve the IM problem. As illustrated in Algorithm 1, it consists of two main steps. First, given  $\epsilon$ , the required number of random RR sets  $\theta$  to guarantee the approximation ratio  $(1 - \frac{1}{e} - \epsilon)$  is estimated, and then, a total number of  $\theta$  random RR sets are generated. Second, the problem of maximum coverage on all generated random RR sets is considered. The goal of the problem is to find a set of  $k$  nodes which covers the maximum number of RR sets. The problem is known to be NP-Hard, and it has been proven that the hill-climbing greedy algorithm yields

---

**Algorithm 1** RIS-based algorithm for the IM problem.

---

**Input:**  $G, k$

**Output:**  $S$

- 1:  $S = \emptyset, R = \emptyset$
  - 2: // Step 1: estimation and sampling
  - 3: Estimate  $\theta$ , i.e., the required number of RR sets
  - 4: Generate  $\theta$  random RR sets and insert them into  $R$
  - 5: // Step 2: maximum coverage
  - 6: **for**  $i = 1$  to  $k$  **do**
  - 7:    $v = \operatorname{argmax}_{u \in V \setminus S} \{ \sum_{RR \in R} \mathbb{1}(u \in RR) \}$
  - 8:    $S = S \cup \{v\}$
  - 9:   Remove all RR sets covered by  $v$  from  $R$
  - 10: **end for**
- 

a  $(1 - \frac{1}{e})$  approximate solution. After solving the maximum coverage problem, the resulting set of nodes (of size  $k$ ) is returned as the solution to the IM problem [22].

Many previous works on RIS-based influence maximization algorithms have focused on providing formal analysis on the required number of random RR sets to ensure  $(1 - \frac{1}{e} - \epsilon)$ -approximation guarantee [8], [23], [24]. Specifically, most studies have tried to reduce the number of required random RR sets, while preserving theoretical guarantees. Borgs et al. [22] terminated the random RR set generation procedure when the total number of traversed edges during random RR set generation exceeds a predefined threshold  $\tau$ . They showed that by setting  $\tau$  equal to  $O(\epsilon^{-2}k(m+n) \log(n))$ , a  $(1 - \frac{1}{e} - \epsilon)$ -approximation solution is guaranteed. The main issue with this method is the large constant factor in the equation that determines the threshold, which incurs heavy computational cost. Tang et al. [23] proposed TIM and TIM+ algorithms. They showed that  $\theta$  should be greater than  $\frac{\lambda}{OPT}$ , where parameter  $\lambda$  is a constant factor determined by  $n, \epsilon$  and  $k$ , and parameter  $OPT$  is the optimal solution's expected influence spread. Later, Tang et al. proposed IMM [8], which reduces the value of  $\lambda$ , and as a result, the value of  $\theta$  based on a martingale approach [25]. Nguyen et al. [24] proposed SSA/DSSA algorithms and reduced  $\theta$ , but later, Huang et al. [26] showed that there are some issues with SSA's analysis and it cannot preserve theoretical guarantees.

Both TIM and IMM algorithms need optimal solution's influence spread, i.e.,  $OPT$ , in order to estimate the required number of RR sets, but in reality there is no prior knowledge on the value of  $OPT$ . These two algorithms resolve this issue by providing a (preferably close) lower bound  $LB$  for the value of  $OPT$ , and therefore, a larger estimated value for  $\theta$ , i.e.,  $\theta = \frac{\lambda}{LB} > \frac{\lambda}{OPT}$ .

To obtain a reasonable lower bound for  $OPT$ , both TIM and IMM use bootstrapping estimation techniques to perform a hypothesis testing (see Algorithm 2). They first set the value of  $\theta$  with an initial estimation (line 4). We refer readers to [8], [23] for more details on the function  $f$  in line 4. Next, they sample random RR sets until the total number of RR sets reaches  $\theta$  (line 5). Then, they construct a seed set from the generated RR sets. The function SeedSelection in line 6 in Algorithm 2 is equivalent to some extent to lines 6 – 10 in Algorithm 1. After that, they estimate the influence spread of the generated seed set and compare it with the approximation bound of the estimator (line 7).

**Algorithm 2** Estimating the lower bound LB.

---

**Input:**  $G, k, \epsilon$   
**Output:**  $LB$

- 1:  $S = \emptyset, R = \emptyset$
- 2: **for**  $i = 1$  to  $\log_2(n) - 1$  **do**
- 3:    $x = n/2^i$
- 4:    $\theta = f(n, \epsilon, k)/x$
- 5:   Sample random RR sets until  $|R| < \theta$
- 6:    $S = \text{SeedSelection}(R, k)$
- 7:   **if**  $n \times F_R(S) \geq (1 + \sqrt{2}\epsilon) \times x$  **then**
- 8:      $LB = n \times F_R(S)/(1 + \sqrt{2}\epsilon)$
- 9:     **break**
- 10:   **end if**
- 11: **end for**

---

If the estimated influence spread is much higher than the approximation bound, they use a discounted value of that estimate as  $LB$  (line 8). Otherwise, they double the number of RR sets (lines 3 – 4) and repeat the described procedure.

### 3 PROPOSED RIS-BASED PARALLEL ALGORITHM

In this section, we present our solution to accelerate computational bottlenecks of RIS-based IM algorithms. Our GPU-accelerated IM algorithm is named gIM. Similar to IMM, our gIM algorithm is composed of two main steps: 1) sampling RR sets, and 2) generating the seed set.

#### 3.1 Baseline Parallel Methods for RR Set Sampling

Before discussing our proposed method, we first investigate two simpler approaches that can be used to parallelize random RR set generation in RIS-based algorithms. For each of these two approaches, we first explain the method and then mention its drawbacks. Our proposed method for RR set sampling will be discussed in Section 3.2.

##### *Thread-Level Parallelization:*

Since in RIS-based algorithms, every random RR set can be generated independently from other random RR sets, one simple approach to parallelize random RR set generation is to assign the task of creating every random RR set to one thread, and launch a large number of threads on GPU.

The main issue with this approach is its severe workload imbalance. Since size of a random RR set may range from a single node to a large portion of the graph, and node degrees vary as well, the amount of work to generate a random RR set varies drastically. This causes severe workload imbalance among the threads within a block, which makes a large number of threads inactive most of the time, and hence, highly degrades the overall performance.

Memory capacity is another impediment to this approach. Since every thread needs to separately maintain large arrays (e.g., visited nodes), memory capacity can easily become the bottleneck and prevent this approach to scale-up to large graphs with millions of nodes.

##### *Parallel Breadth First Search:*

Generating a random RR set is equivalent to running breadth-first search (BFS) algorithm starting from a randomly selected node on a reversed instance graph. One can avoid generating a whole reversed instance graph by incorporating instance graph generation into the BFS algorithm. In the standard BFS, once a node  $u$  is picked up from the front of the frontier queue, all of its outgoing edges are traversed. If instead of traversing all the outgoing edges, one traverses every outgoing edge of  $u$  with probability of  $p_{uv}$ , then the result is equivalent to running a BFS on a reversed instance graph.

Therefore, one approach towards accelerating random RR set generation is to employ parallel algorithms previously proposed for BFS on GPU [27], [28], [29]. Such methods parallelize the procedure of generating one random RR set on GPU.

Most previous solutions on parallelizing BFS on GPU are level-synchronous, that is each level can be processed in parallel, but consecutive levels must be processed sequentially. These methods often maintain a node frontier queue, which is produced by processing the previous level, and is used to produce the node frontier for the next level. In level-synchronous methods, parallelization is usually implemented among the nodes in the current frontier queue or their outgoing edges, therefore, if the size of current frontier is relatively small, then not much gain can be attained from parallelization. Even the overhead of synchronization between consecutive levels can exceed any gain achieved from parallelization. Unfortunately, the problem of small frontier is not rare. Under the LT model, each node has at most one active outgoing edge so the maximum number of nodes in the frontier does not exceed 1. Even under the IC model, where nodes may have more than one active edge, low values of influence probabilities in real data cause instance graphs to be usually much sparser than  $G$ . Therefore, the average number of active edges per node is low and the frontier cannot grow very much.

Gaihre et al. [30] proposed XBFS and implemented asynchronous BFS which allows nodes from different levels be visited in the same iteration. However, XBFS performs only one BFS at a time. Liu et al. [31] exploited the similarity between frontiers of BFS traversals started from different nodes and proposed an algorithm to run multiple concurrent BFSs on the same GPU. This algorithm still requires relatively large frontiers to achieve performance gains and is also heavily reliant on the assumption of high resemblance between frontiers, which may not be the case for generating random RR sets.

Another important issue with using previous parallel BFS algorithms is producing redundant nodes in the frontier queue. Since all nodes in the current frontier queue are processed in parallel, a node which is adjacent to more than one node in the current frontier might be placed in the next level's frontier more than once. In standard BFS, this issue may degrade speed, but does not affect the correctness of the algorithm. However, edge traversal in random RR set generation is randomized, and when a node is placed more than once in the frontier, its outgoing edges are processed multiple times, and this issue violates the correctness of RR set generation. For example, if node  $u$  is placed in

the frontier queue twice, each one of its outgoing edges  $e = (u, v)$  are also processed twice, which accordingly increases influence probability from  $p_{uv}$  to  $1 - (1 - p_{uv})^2$ .

### 3.2 Proposed Parallel Method for RR Set Sampling

The above two approaches are actually located at two extreme ends of the spectrum. One launches many random RR set generation tasks but executes each one of them sequentially, and the other generates a single random RR set at a time and strives to parallelize that single task using all the GPU resources. gIM employs a parallel method which is located somewhere between these two extremes.

#### Overview:

We run many random RR set generations in parallel, and also, parallelize the execution of every random RR set generation to some extent. In specific, we judiciously assign the task of generating every random RR set to one CUDA block. See Algorithm 3. Every block is associated with a fixed-size frontier queue in the shared memory which is denoted by  $Q_{shr}$ , a *Visited* array of length  $n$  to keep account of visited nodes during BFS, and a custom-designed data structure named  $RR_{tmp}$  which is used to temporarily hold the generated RR set. Both *Visited* and  $RR_{tmp}$  are located in GPU global memory.

The number of RR sets to which a particular node belongs, is an estimator of that node’s influence spread. Hence, we also need to maintain the number of occurrences of every node in all the generated RR sets. To do so, we employ array *Occur* of length  $n$ , whose elements are initialized to zero.

#### Data Structure $RR_{tmp}$ :

The number of nodes in an RR set may vary from one to all the nodes. Therefore,  $RR_{tmp}$  should be capable of storing an RR set as large as the set  $V$ . If  $RR_{tmp}$  is implemented as an array of length  $n$ , since every block is associated with a  $RR_{tmp}$ , the total number of blocks that can be executed in parallel would be restricted by the amount of memory required to store all these large arrays. If  $RR_{tmp}$  is implemented as a linked list, whose memory footprint can grow dynamically, the overhead of memory management would be too large. To tackle this issue, we devise a data structure that is similar to a linked list. However, the granularity for memory management is set larger than a single node. In specific, each element in the linked list is composed of a pointer to the next element and a fixed-length array to store some of the nodes.

#### Graph Representation:

We represent graph  $G$  in the compressed sparse row (CSR) format. As illustrated in Fig. 1, the CSR representation consists of two arrays, namely, the column indices array  $C$ , and the row offsets array  $R$ . For a graph  $G$  with  $n$  nodes and  $m$  edges, array  $C$  of size  $m$  is created by concatenating adjacency lists of all the nodes in  $G$ . Array  $R$  has  $n + 1$  elements, where element  $R[i]$  denotes the location of the adjacency list of node  $i$  in array  $C$ . One can find all outgoing neighbors of a particular node  $i$ , by reading all elements in  $C$  indexed from  $R[i]$  to  $R[i + 1]$ . In addition, array  $W$  of size  $m$  contains influence probability  $p_{uv}$  for every edge  $e = (u, v)$ .

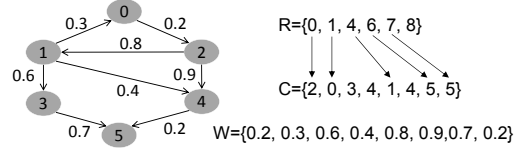


Fig. 1. Example: a graph and its corresponding CSR representation.

**Algorithm 3** Parallel kernel for random RR set generation (see Section 3.2).

**Input:**  $R, C, W, \theta$

**Output:**  $Occur, RR_{tmp}$

**# of blocks:**  $\theta$

**# of threads / block:**  $N_{th}$

```

1: if  $tx == 0$  then
2:    $RR_{tmp}.init()$ 
3:    $Q_{shr}.init()$ 
4:    $Q_{res}.init()$  // Section 3.3
5:    $u = randSelect(V)$ 
6:    $Q_{shr}.enqueue(u)$ 
7: end if
8: while  $!Q_{shr}.empty()$  do
9:    $u = Q_{shr}.front()$ 
10:  if  $tx == 0$  then
11:     $RR_{tmp}.insert(u)$ 
12:     $Q_{shr}.dequeue()$ 
13:     $atomicAdd(Occur[u], 1)$ 
14:  end if
15:   $s = R[u], e = R[u + 1]$ 
16:  for ( $i = tx; i < e - s; i = i + N_{th}$ ) do
17:     $ofloadQueue(Q_{shr}, Q_{res})$  // Section 3.3
18:     $v = C[s + i], p_{uv} = W[s + i], p = U(0, 1)$ 
19:    if  $p < p_{uv} \ \& \ Visited[v] == false$  then
20:       $Visited[v] = true$ 
21:       $Q_{shr}.atomic_enqueue(v)$ 
22:    end if
23:  end for
24:   $reloadQueue(Q_{shr}, Q_{res})$  // Section 3.3
25: end while

```

Note that in our CSR representation, all the nodes and edges are kept in the same ordering as  $G$ , i.e., we do not perform any sorting or pre-processing on the nodes or edges of the input graph  $G$ .

#### Algorithm Details:

Algorithm 3 shows our parallel algorithm for random RR set generation under IC diffusion model. The input graph is represented by  $R, C$  and  $W$  arrays. Parameter  $\theta$  is the number of required random RR sets to be generated.

The algorithm works as the following. Every thread within a block is distinguished by an index denoted as  $tx$ , which ranges from 0 to  $N_{th} - 1$ . First, thread zero ( $tx = 0$ ) performs the required initializations, randomly selects a node from  $V$  as the source node for BFS, and inserts this node at the front of  $Q_{shr}$  (lines 2 – 6). Next, all  $N_{th}$  threads start a parallel randomized BFS from the randomly selected node. In this parallel BFS, nodes are processed sequentially and parallelization is realized in evaluation of

adjacent edges to the current node. The BFS continues as long as the shared queue  $Q_{shr}$  is not empty (line 8).

If  $Q_{shr}$  is not empty, all threads extract the node at the front of  $Q_{shr}$  (line 9). Next, thread zero ( $tx = 0$ ) inserts the extracted node to  $RR_{tmp}$ , removes it from  $Q_{shr}$ , and increments its corresponding element in array  $Occur$  (lines 10 – 14). Then, all threads within the block find the range of the extracted node’s out-neighbors in array  $C$  (line 15). Next, they collaboratively read the out-neighbors of the current node from  $C$  and their corresponding influence propagation probabilities from  $W$  (line 18). While processing an edge, every thread produces a random number from the uniform distribution  $U(0, 1)$ , and compares the generated random number with influence probability of that edge, and based on this comparison, it decides whether to traverse that edge or not. When a thread decides to traverse a particular edge, it checks the visited flag of the destination node of that edge, and if the node has not been visited already, it is added to the frontier queue (lines 19 – 22). Since threads within a block are executed in parallel, more than one thread may decide to add a node to  $Q_{shr}$ . Therefore adding a node to  $Q_{shr}$  should be protected by atomic operations (line 21). Details of *offloadQueue()* in line 17 and *reloadQueue()* in line 24 are discussed in Section 3.3.

### 3.3 Avoiding Overflow of the Frontier Queue

Shared memory is one of the most limited resources in GPU. In the proposed parallel algorithm, every block uses shared memory to store its own fixed-size frontier queue. If we allocate small queues, the actual frontier may grow larger and cause the queue to overflow, in which case, the produced RR set is no longer valid. In order to guarantee that no queue overflow occurs, we could increase the queue size to the total number of nodes. However, by doing this, we are no longer able to maintain the queue in GPU shared memory. Moving the queue to GPU global memory incurs long memory access latencies. To resolve this issue, gIM employs the following solution. We define a threshold  $\tau_q$  as

$$\tau_q = |Q_{shr}| - N_{th}, \quad (4)$$

where  $|Q_{shr}|$  is queue capacity (the maximum number of nodes that every queue may hold), and  $N_{th}$  is the total number of threads within a block. When the number of nodes currently stored in the queue exceeds  $\tau_q$ , there is a chance for the queue to overflow. This is because each one of  $N_{th}$  threads within the block processes one edge, and therefore, has a chance to add one new node to the queue. In the worst case, when all  $N_{th}$  threads add a node to the queue, overflow occurs.

To resolve this issue, we associate a reservoir queue, to every block, which is denoted by  $Q_{res}$  and is actually implemented as a stack whose elements are arrays of length  $N_{th}$ . The reservoir queue is stored in GPU global memory and its size may grow dynamically, hence, it does not have shared queue’s limitations. Our approach to avoid possible overflows is the following. As shown in line 17 in Algorithm 3, *offloadQueue()* function is executed before all threads process their incident edges. As shown in Algorithm 4, this function checks the condition  $Q_{shr}.size() > \tau_q$ . If the condition is satisfied, thread zero allocates the required

---

#### Algorithm 4 *offloadQueue* function

---

```

1: function OFFLOADQUEUE( $Q_{shr}, Q_{res}$ )
2:   if  $Q_{shr}.size() > \tau_q$  then
3:     if  $tx == 0$  then
4:        $ptr = Q_{res}.alloc()$ 
5:     end if
6:      $ptr[tx] = Q_{shr}.deque\_block(tx)$ 
7:   end if
8: end function

```

---



---

#### Algorithm 5 *reloadQueue* function

---

```

1: function RELOADQUEUE( $Q_{shr}, Q_{res}$ )
2:   if  $Q_{shr}.empty() \ \& \ !Q_{res}.empty()$  then
3:     if  $tx == 0$  then
4:        $ptr = Q_{res}.pop()$ 
5:     end if
6:      $Q_{shr}.enqueue\_block(ptr[tx], tx)$ 
7:   end if
8: end function

```

---

memory to expand the reservoir queue. Then all threads collaboratively move  $N_{th}$  elements from the shared queue to the newly allocated space in reservoir queue. This reduces the number of nodes in the shared queue and prevents possible overflow.

The elements of  $Q_{shr}$  which have been moved to  $Q_{res}$ , need to be returned back to  $Q_{shr}$  when it has enough empty space. This procedure is implemented in *reloadQueue()* function, which is shown in Algorithm 5. In this function, threads check if  $Q_{shr}$  is empty and also if there are some nodes in the  $Q_{res}$  that can be brought back to  $Q_{shr}$  (line 2). If this condition is satisfied, threads collaboratively move  $N_{th}$  nodes from  $Q_{res}$  to  $Q_{shr}$  (line 6). Once a batch of nodes are transferred from  $Q_{shr}$  to  $Q_{res}$ , the order in which nodes are processed differs from the conventional BFS. However, this does not break the correctness of the parallel algorithm, as the order in which nodes are visited is not important.

### 3.4 An Illustrative Example

Fig. 2 illustrates the execution flow of a random RR set generation a simple graph by using one small CUDA block with  $N_{th} = 3$  threads. Since outgoing edges are processed via parallel threads, edges are marked with different colors according to their corresponding threads. Black edges are not processed throughout this example, because their source node is not visited and thus is not placed in the queue. Check marks and cross marks on the edges represent whether the corresponding edge is traversed or not.

The execution starts by thread zero which places node 0 into  $Q_{shr}$ . Next, node 0 is extracted by all threads and its outgoing edges are processed. Since node 0 has more than  $N_{th} = 3$  outgoing edges, processing these edges are performed in two consecutive steps and newly visited nodes are put in  $Q_{shr}$ . Then, node 1 is extracted from  $Q_{shr}$ , but threads find that  $N_q = 3 > \tau_q = 2$ , so  $N_{th}$  nodes are moved from  $Q_{shr}$  to  $Q_{res}$ . Next, outgoing edges of node 1 are processed, and recently visited node 6 is added to  $Q_{shr}$ . In the next step, node 6 is extracted from  $Q_{shr}$  and its outgoing edges are processed. Afterwards, threads find that

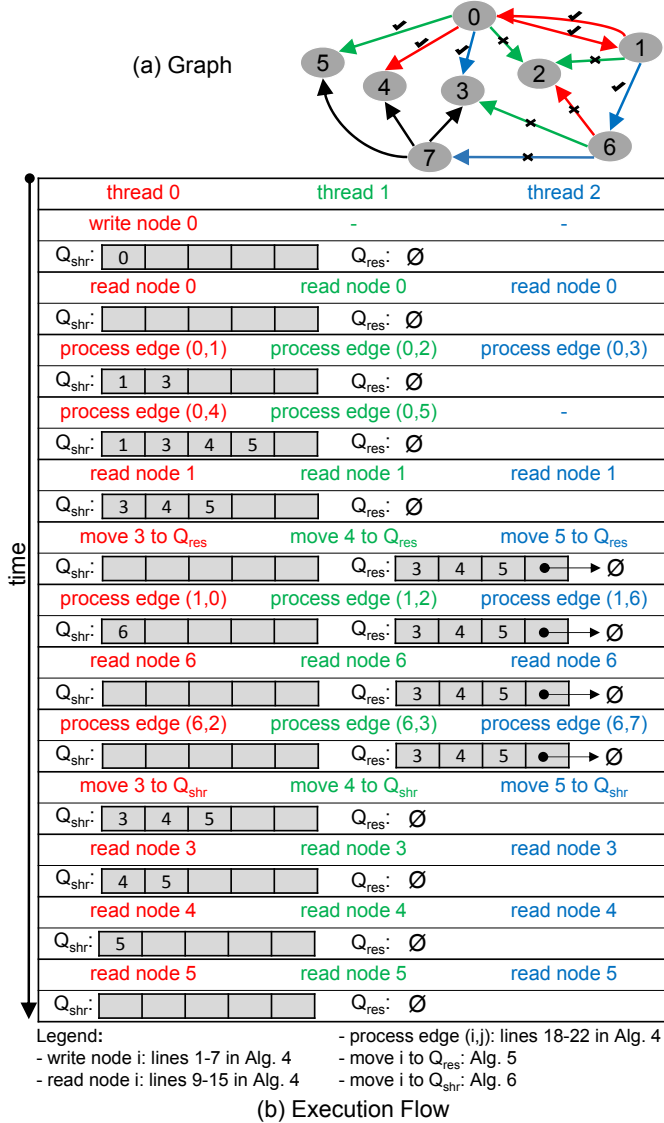


Fig. 2. (a) An example graph. (b) Execution flow of a random RR set generation starting from node 0 within a small block with three threads. Details are discussed in Section 3.4.

$Q_{shr}$  is empty and thus bring back  $N_{th}$  nodes from  $Q_{res}$  to  $Q_{shr}$ . The rest of the procedure is straightforward. Nodes are removed from the queue one by one until both  $Q_{shr}$  and  $Q_{res}$  become empty.

### 3.5 Optimizing Memory Usage

In most cases, size of the generated RR set is very small, and therefore,  $RR_{tmp}$  data structure has only few elements and the allocated array inside the last element of this data structure is underutilized. This under-utilization of GPU global memory should be avoided since the number of RR sets to be generated is large.

To this end, Algorithm 3 is modified as shown in Algorithm 6. Right after producing an RR set, it is copied from  $RR_{tmp}$  to a continuous chunk of memory called  $RR$ . This data structure is basically formed by concatenation of all the generated RR sets.

In addition, the total number of concurrent blocks is reduced from  $\theta$  (the number of RR sets to be generated)

**Algorithm 6** Memory-optimized parallel kernel for random RR set generation (see Section 3.5).

**Input:**  $R, C, W, \theta$

**Output:**  $Occur, RR$

# of blocks:  $N_b$

# of threads / block:  $N_{th}$

```

1: repeat
2:   Generate one random RR set as in Algorithm 3
3:   if  $tx == 0$  then
4:      $offset = tail_{RR}$ 
5:      $Offsets_{RR}[N_{RR}] = offset$ 
6:      $N_{RR} = N_{RR} + 1$ 
7:      $tail_{RR} = tail_{RR} + RR_{tmp}.size()$ 
8:   end if
9:   for ( $i = tx; i < RR_{tmp}.size(); i = i + N_{th}$ ) do
10:     $RR[offset + i] = RR_{tmp}[i]$ 
11:   end for
12: until  $N_{RR} \leq \theta$ 

```

to a constant number  $N_b$ . Parameter  $N_b$  is discussed later in Section 3.6. This is done because every block requires to store intermediate data (i.e.,  $Q_{res}$  and  $Visited$ ) in GPU global memory.

In Algorithm 6,  $N_{RR}$  and  $tail_{RR}$  are scalar variables which lie in GPU global memory and keep track of the number of generated random RR sets so far, and the sum of their sizes, respectively. Once a random RR set is generated (line 2), the location where it should be copied to in array  $RR$  is determined and stored in  $Offsets_{RR}$  (lines 4 – 5). This is required so that the boundaries between consecutive random RR sets in array  $RR$  can be distinguished later. In addition,  $N_{RR}$  is incremented by 1, and  $tail_{RR}$  is incremented by the size of the generated RR set (lines 6 – 7). Finally, all threads within the block collaboratively copy the generated RR set from  $RR_{tmp}$  to the specified location in array  $RR$  (lines 9 – 11). A block terminates its execution when the number of generated RR sets has reached the number of required RR sets (line 12).

Note that since  $N_{RR}$  and  $tail_{RR}$  are accessible by all concurrent blocks, lines 4 – 7 should be protected from concurrent execution using mutex and atomicCAS in CUDA API. For brevity, this is not shown in Algorithm 6.

Without the optimization presented in this section, the required data would not fit into GPU global memory. For example, storing one  $Visited$  array for a graph with 4 million nodes requires about 488 kB of GPU global memory, hence, if we were to launch  $\theta$  blocks instead of  $N_b$  blocks, launching  $\theta = 1$  million blocks would require 465 GB of GPU global memory only to store all the  $Visited$  arrays. This is about 10 to 60 times larger than the amount of memory available on even high-end GPUs today.

### 3.6 Tuning to the Underlying Hardware

*Tuning  $N_{th}$  and Synchronization Mechanism:*

To ensure correct functionality, we need to synchronize all the threads within a block after lines 7, 9 and 14 in Algorithm 3, line 5 in Algorithm 4 and line 5 in Algorithm 5. This is normally done by using `syncthreads()` API call



in CUDA. For example, a `syncthreads()` call needs to be inserted after line 7 because thread zero inserts a node in  $Q_{shr}$ , which affects the execution of other threads.

The `syncthreads()` API call is known to degrade performance. Hence, we propose the following optimization. In CUDA, every 32 threads inside a block form a warp. All threads within a warp execute in lock steps. Hence, we set  $N_{th} = 32$ , and as a result, no `syncthreads()` calls are necessary anymore. This is because every block consists of only one warp, i.e.,  $N_{th} = 32$  threads. Since all threads within a warp execute the same instruction, they are already in synchronization. This is the case in Nvidia Kepler, Maxwell, and Pascal GPU architectures.

In recent Volta, Turing, and Ampere GPU architectures, however, the SIMT execution model has experienced some modifications, and as a result, the assumption of implicit synchronization among threads in a warp is no longer valid. Consequently, to ensure correct functionality, we need to insert `syncwarp()` API calls after lines 7, 9 and 14 in Algorithm 3, and after line 5 in both Algorithm 4 and Algorithm 5. The `syncwarp()` API call only synchronizes the threads within a warp, and has lower overhead compared to `syncthreads()` API call.

Therefore, we avoid `syncthreads()` API call by selecting  $N_{th} = 32$ . Another reason for selecting  $N_{th} = 32$  is that average node degree is very low in many real-world graphs, and 32 threads are enough to process all incident edges in parallel in most cases. Using larger values for  $N_{th}$  causes a large number of threads to remain idle during the process of incident edges of low-degree nodes. See Section 4.7 for a discussion on the effect of larger values of  $N_{th}$  on performance.

*Tuning  $N_b$  and  $|Q_{shr}|$ :*

In order to fully utilize parallel processing capacity of GPU, we determine the total number of blocks to be launched as the following equation:

$$N_b = \#SM \times MaxResBlockPerSM, \quad (5)$$

where  $\#SM$  is the total number of streaming multiprocessors (SM) in the target GPU device, and  $MaxResBlockPerSM$  is the maximum number of resident blocks per SM. The above formula means that we launch just enough blocks in Algorithm 6 to keep all SMs busy, while avoiding the overhead of launching too many blocks, i.e.,  $\theta$  blocks as in Algorithm 3.

To ensure that the number of concurrent blocks is not limited by the amount of available shared memory,  $|Q_{shr}|$  is determined such that:

$$\lfloor \frac{|ShMem|}{|Q_{shr}|} \rfloor \geq MaxResBlockPerSM, \quad (6)$$

where  $|ShMem|$  is the amount of shared memory per SM, and  $|Q_{shr}|$  is the capacity of the fixed-size frontier queue. For a GPU with  $|ShMem| = 48$  KB and  $MaxResBlockPerSM = 32$ , we must select  $|Q_{shr}| < 384$ , i.e., 1.5 KB.

### 3.7 Modifications under the LT Model

Our proposed parallel algorithm for random RR set sampling can be modified to support LT diffusion model as

---

**Algorithm 7** Maximum Coverage parallel kernel.

---

**Input:**  $RR, Occur, u$

**Output:**  $Occur, Covered$

**# of blocks:**  $N_b$

**# of threads / block:**  $N_{th}$

```

1: for ( $i = bx; i < N_{RR}; i = i + N_b$ ) do
2:   if  $Covered[i] == false$  then
3:      $found = false$ 
4:      $offset = Offsets_{RR}[i]$ 
5:      $length = Offsets_{RR}[i + 1] - offset$ 
6:     for ( $j = tx; j < length \ \& \ !found; j = j + N_{th}$ ) do
7:       if  $RR[offset + j] == u$  then
8:          $found = true$ 
9:       end if
10:    end for
11:    if  $found == true$  then
12:       $Covered[i] = true$ 
13:      for ( $j = tx; j < length; j = j + N_{th}$ ) do
14:         $atomicAdd(Occur[RR[offset + j]], -1)$ 
15:      end for
16:    end if
17:  end if
18: end for

```

---

well. Under this model, every node has at most one active incoming edge, which is selected randomly according to the edge weights. First, thread zero samples a random number from  $U(0, 1)$ . Then all threads within the block run a parallel scan algorithm [32] on the edge weights, and instead of comparing the generated random number with edge weights, they compare it with two consecutive numbers that resulted from the parallel scan. The first thread that finds the generated random number between two consecutive outputs of scan, marks its edge as active and broadcasts an early termination signal to all other threads within the block.

In addition, since every node visits at most one other node during random BFS, the size of the frontier queue never exceeds one, and therefore, the fixed-size shared frontier queue never overflows and the reservoir queue is not required anymore.

### 3.8 Proposed Parallel Method for Seed Set Generation

The final step is to solve the Maximum Coverage problem on the generated random RR sets, in order to generate the seed set. IMM [8] used the standard greedy algorithm to solve this problem, which is illustrated in lines 6 – 10 of Algorithm 1.

Our proposed solution, however, employs a more efficient approach. As mentioned in Section 3.2, every node is associated with a counter which keeps track of the number of its occurrences across all the generated random RR sets. All such counters are stored in array  $Occur$ . By running a parallel reduction algorithm [33] with max operator on this array, the node with maximum coverage is found.

Once a node is selected and added to the seed set, all RR sets which contain that node must be removed, and also, all counters whose corresponding nodes are present in those RR sets must be decremented. To do this, we associate

every RR set with a *Covered* flag which indicates whether or not the RR set is covered by the seed set so far. In order to mark the covered RR sets and update the counters, we employ the parallel kernel shown in Algorithm 7. Here, the generated RR sets are distributed among  $N_b$  blocks (line 1). When a block processes an RR set, it first checks that the RR set has not already been covered by any node in the seed set (line 2). Then, all threads within the block collaboratively search for all occurrences of node  $u$ , i.e., the newly added node in the seed set (lines 3 – 10). If the RR set contains node  $u$ , it is flagged as covered, and all threads within the block collaboratively decrement the counters whose corresponding nodes are present in this RR set (lines 11 – 16). Note that all the threads within a block need to be synchronized after lines 2 and 10, according to the discussion in Section 3.6.

## 4 EXPERIMENTAL EVALUATION

### 4.1 gIM Source Code

The proposed gIM algorithm is implemented in C++ language using CUDA parallel programming framework. The source code of gIM is available online [34].

### 4.2 Experiment Setup

Benchmark datasets are graphs extracted from real social networks. All of these datasets are publicly available for download from [35], and have been widely employed in the literature on the IM problem to assess the influence spread and running time of different algorithms [8], [11], [12], [36]. Their main characteristics are shown in Table 1.

We compare our proposed gIM algorithm with IMM [8], Ripples [11], cuRipples [12] and SKIM [36]. IMM is the state-of-the-art algorithm which provides theoretical guarantee on the approximation ratio of influence spread. Ripples is a framework which provides multi-core implementation of IMM, and cuRipples is the extended version of Ripples that adds support for multi-GPU systems. SKIM is among the best algorithms which does not provide the theoretical guarantee but runs faster than IMM.

To run the experiments, we employed a Linux machine with Intel Xeon Scalable 4110 CPU operating at 2.50 GHz, and Tesla V100 GPU operating at 1245 MHz and equipped with 16GB of global memory. We employ CUDA version 10.2. Since IMM [8] and SKIM [36] are sequential methods, they are executed on a single CPU core, while Ripples and cuRipples use all the available CPU cores. In gIM, the parallel kernels are executed on GPU, and the rest of the program is executed on a single CPU core.

Although some efforts [37], [38], [39] have been made to learn influence probabilities based on user actions extracted from social networks, most previous works used heuristics to assign influence probabilities to the edges. Weighted Cascade (WC) [1] is the most commonly used schemes in previous works [8], [23], [36], [40]. In this scheme, the influence probability of edge  $e = (u, v)$  is set to  $1/d_v^{in}$ , where  $d_v^{in}$  is the in-degree of node  $v$ . By using this scheme, the probability of a node being influenced is heuristically less dependent on the number of incoming edges. Also, WC scheme, in which the sum of incoming edges' probabilities

TABLE 1  
Benchmark datasets.

Dataset	Type	# of nodes ( $n$ )	# of edges ( $m$ )
soc-Epinions1	Directed	75,879	508,837
soc-Slashdot0922	Directed	77,360	905,468
higgs-twitter	Directed	456,631	14,855,875
soc-Pokec	Directed	1,632,803	30,622,564
soc-LiveJournal1	Directed	4,847,571	68,993,773
com-Orkut	Undirected	3,072,441	117,185,083

for each node is exactly equal to 1, can be used under the LT model. In our experiments, we use WC scheme to assign influence probabilities to edges, under both IC and LT models.

### 4.3 Performance Comparison

Table 2 illustrates the runtime of different methods with  $k = 50$  and  $\epsilon = 0.05$  under the IC model in different datasets. In every algorithm, and for every dataset, the algorithm is executed ten times, runtime values are measured, and the average runtime is reported. The observed standard deviation is about 2% of the average value. Note that the entire runtime is measured, which includes the time required to perform memory allocations on GPU and transfer data between CPU and GPU.

We used the default parameters in the evaluation of SKIM. It is noteworthy to mention that the solution quality, i.e., the influence spread of the resulting seed set, is the same for IMM and gIM, however, the solution quality of SKIM is slightly lower (about 1% to 3% lower).

The runtime of IMM varies from 4.42 seconds to 813.09 seconds, and it increases as graph size grows. SKIM runs faster than IMM. The runtime varies from 0.29 to 36.53 seconds. gIM, however, attains much smaller runtime compared to both IMM and SKIM. Our solution's runtime ranges from 131 milliseconds to 3.689 seconds. Note that the measured runtime values include all overheads such as CPU to GPU data transfer latency. The speedup ratio of gIM over IMM ranges from  $33.74\times$  to  $220.41\times$ .

The runtime of Ripples' multi-core implementation of IMM, with 16 parallel threads, ranges from 11.09 to 586.52 seconds. cuRipples has lower runtimes compared to Ripples. However, despite using the computational power of both CPU and GPU, cuRipples is slower than gIM, which only employs the computational power of GPU. cuRipples moves the generated RR sets back and forth between the device and host memories, and therefore, is not limited to the size of GPU memory. cuRipples used two different methods for GPU implementation of random BFS: under the IC model, it employed a parallel BFS which generates one RR set at a time, and under the LT model, it adopted thread-level parallelism on GPU. Both of these methods are very similar to what discussed in Section 3.1.

It is noteworthy to mention that Minutoli et al. [11], [12] used a uniform distribution  $U(0, 1)$  to assign weights to the edges, while similar to other previous works such as [8], [23], [36], [40] we use WC scheme. By using uniform distribution, half of the edges become active during the generation of a random RR set, because the mean of  $U(0, 1)$  is 0.5. While in WC scheme, since the in-degree of most nodes is larger than two, most edges are assigned with

TABLE 2  
Runtime of different methods under IC diffusion model.

		soc-Epinions1	soc-Slashdot0922	higgs-twitter	soc-Pokec	soc-LiveJournal1	com-Orkut
Runtime (sec.)	IMM [8]	4.42	12.69	92.94	122.91	251.86	813.09
	Ripples [11]	11.09	18.35	56.86	65.23	133.71	586.52
	cuRipples [12]	1.25	1.68	13.08	21.42	87.53	96.98
	SKIM [36]	0.29	0.58	2.49	4.43	7.71	36.53
	gIM (proposed)	0.131	0.143	0.802	1.084	2.276	3.689
Speedup ratio over IMM	Ripples [11]	0.40	0.69	1.63	1.88	1.88	1.39
	cuRipples [12]	3.54	7.55	7.11	5.74	2.88	8.38
	SKIM [36]	15.24	21.87	37.33	27.74	32.67	7.82
	gIM (proposed)	33.74	88.74	115.89	113.39	110.66	220.41

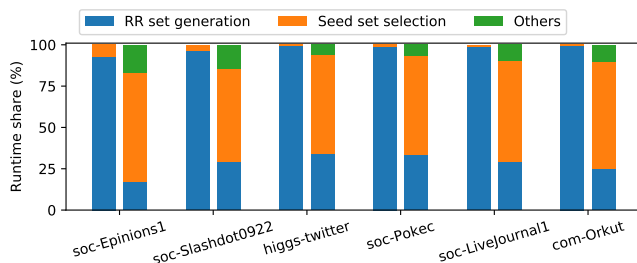


Fig. 3. Runtime breakdown of gIM (the proposed solution) compared to IMM. In each dataset, the left bar represents IMM and the right bar gIM.

weights less than 0.5, and therefore, we expect to encounter less active edges during the generation of a random RR set in comparison with the uniform distribution model. Therefore, the difference between our experimental results for cuRipples and the ones in [12] is due to the usage of different methods for assigning weights.

#### 4.4 Runtime Breakdown

RIS-based algorithms can be decomposed into two main parts: random RR set sampling and seed selection, each of which has different computational characteristics. Fig. 3 illustrates the share of each part in the total runtime, for both IMM and gIM. Note that the time required for memory allocation and moving data between CPU and GPU are also taken into account in gIM results.

In IMM, a large portion of the runtime is spent on generation of random RR sets, while in gIM, seed set selection has a larger share of the runtime. This is because the speed-up ratio for our GPU acceleration is different for the two parts, in specific, the speedup ratio is higher for random RR set sampling in comparison with seed set selection.

#### 4.5 Impact of Parameters $k$ and $\epsilon$

The number of required random RR sets to guarantee  $(1 - \frac{1}{e} - \epsilon)$ -approximation ratio, which directly affects runtime, is dependent on the values of  $k$  and  $\epsilon$ . Fig. 4 and Fig. 5 compare the runtime of gIM against IMM with  $\epsilon = 0.1$  and for different values of  $k$ . As can be seen, the speedup factor is relatively consistent for different values of  $k$ . It can be seen that the runtime of IMM is monotonically increasing, while the runtime of gIM sometimes drops when  $k$  increases. This is because by increasing the value of  $k$ , the amount of influence of the seed set also increases and it might cause the condition in line 7 of Algorithm 2 to be satisfied one iteration earlier. In IMM, calling greedy seed set generation

procedure for one fewer iteration does not significantly affect the runtime, because as it is evident in Fig. 3, random RR set generation is the dominant contributor to the runtime. While in gIM, random RR generation experiences much higher speedup in comparison with greedy seed selection. Therefore, the execution time of greedy seed selection procedure calls become comparable to that of random RR set generation, and one fewer iteration to this procedure can decrease the runtime in spite of the increased required number of RR sets.

We also studied the impact of  $\epsilon$  on runtime. Fig. 6 illustrates the running times for different values of  $\epsilon$  with  $k = 50$ . The number of required RR samples (i.e.,  $\theta$ ) has an inverse quadratic relation with the value of  $\epsilon$ . As can be seen, the runtime values follow the same trend, and therefore the amount of speedup is almost preserved for different values of  $\epsilon$ .

#### 4.6 Scalability

In this section, we evaluate the scalability of gIM. Specifically, we measure the runtime of gIM against IMM on synthetic graphs with varying densities and show that the speedup increases, as the density of the input graph grows.

In order to generate synthetic graphs, we use Barabasi-Albert model, which is an algorithm for generating undirected random scale-free graphs by using preferential attachment mechanism. This algorithm starts with a clique of  $r_0$  nodes, and then, new nodes are added one by one. Every newly added node is randomly connected to  $r$  nodes that are already in the graph ( $r \leq r_0$ ). The constant parameter  $r$  is the graph density. The probability that the newly added node is connected to node  $i$  is determined by  $p_i = \frac{d_i}{\sum_j d_j}$ , where  $d_i$  denotes the current degree of node  $i$ , and the sum is over the nodes that already exist in the graph. It has been proven that Barabasi-Albert model generates a graph with a power law degree distribution [41].

We generated various graphs with  $n = 10^6$  nodes and different densities ranging from  $r = 2$  to  $r = 32$ , by using Barabasi-Albert algorithm. We measure the runtime of gIM and IMM over the generated graphs, with parameters  $k = 50$  and  $\epsilon = 0.05$ . Fig. 7 illustrates the resulting runtimes and speedup. As can be seen, the speedup increases, as  $r$  grows. This observation can be attributed to the way by which edges are processed in gIM:  $N_{th}$  threads within a warp process all outgoing edges of a single node. Therefore, if the number of outgoing edges of a node is less than  $N_{th}$ , some threads remain idle and perform no useful computation. When  $r$  grows, the average degree also increases, and

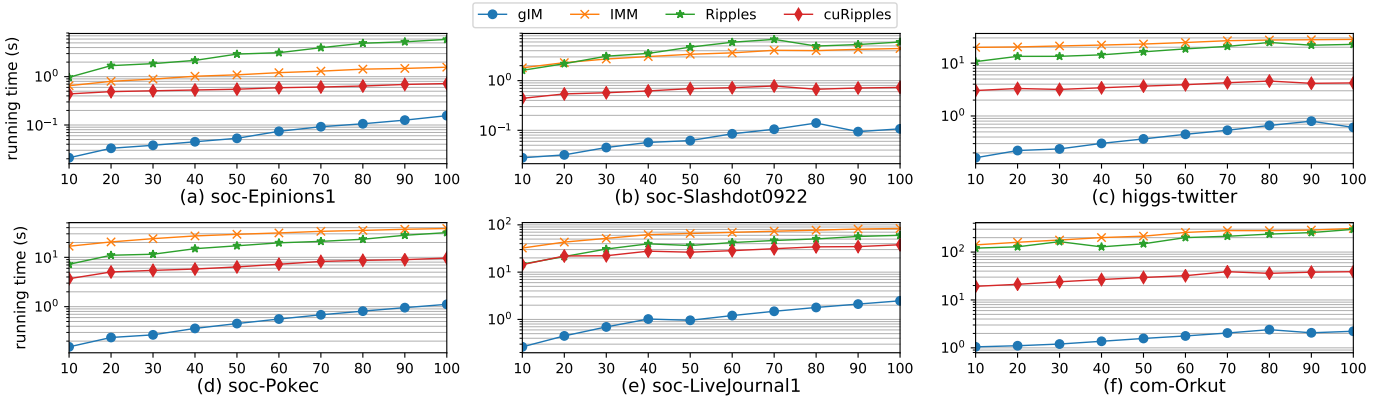


Fig. 4. Runtime of different algorithms (y axis) versus  $k$  (x axis) under IC model and  $\epsilon = 0.1$ .

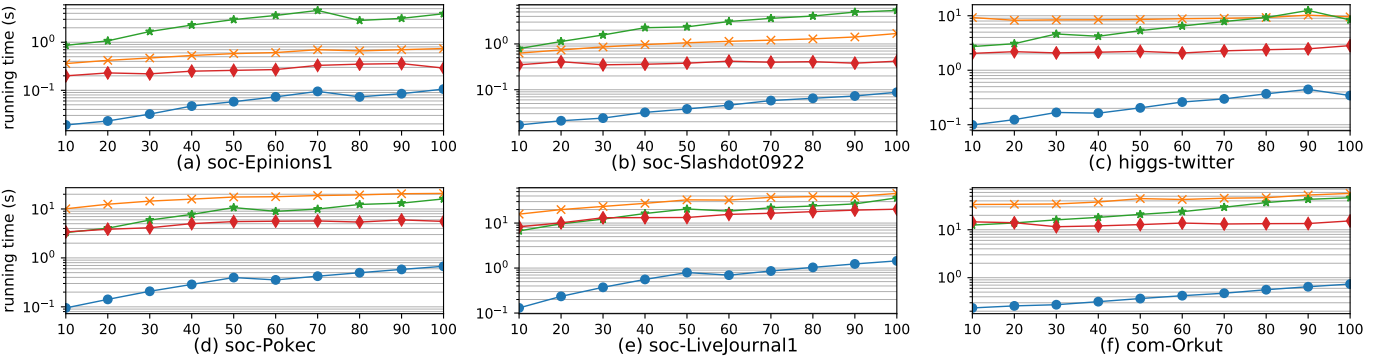


Fig. 5. Runtime of different algorithms (y axis) versus  $k$  (x axis) under LT model and  $\epsilon = 0.1$ .

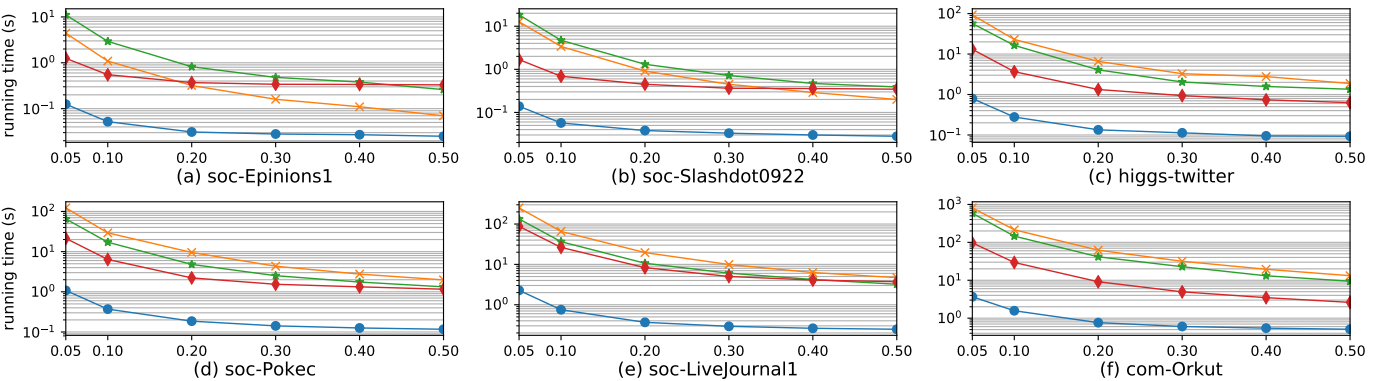


Fig. 6. Runtime of different algorithms (y axis) versus  $\epsilon$  (x axis) under IC model and  $k = 50$ .

as a result, fewer threads remain idle and more potential parallelism can be exploited.

#### 4.7 Impact of Block and Grid Size

As discussed in Section 3.6, we set the number of threads per block to  $N_{th} = 32$  to avoid having too many idle threads while processing low-degree nodes. In order to evaluate the impact of  $N_{th}$  on the performance of gIM, we replaced `syncwarp()` API calls with `syncthreads()` API calls and set  $N_{th} = 64, 96$ , and  $128$ . Fig. 8a illustrates the runtime with different values of  $N_{th}$ , normalized to the runtime of  $N_{th} = 32$ . As it can be seen, the performance degrades by increasing  $N_{th}$ . One reason is that the average degree in social graphs is relatively low, and having a larger block causes a higher percentage of the threads to remain idle while processing low-degree nodes. In addition, since the maximum number of resident threads per SM is constant, having larger blocks means fewer number of blocks may run concurrently.

To efficiently utilize GPU parallel computing capabilities, we set the number of blocks, i.e. grid size, according to Equation 5 to  $N_b = 2560$  for the employed V100 GPU device. In order to analyze the effect of grid size on performance, we varied  $N_b$  and measured the runtime of gIM on different datasets. The resulting execution times are shown in Fig. 8b. As can be seen, the smallest execution times are achieved when  $N_b$  is set to the value specified by Equation 5. According to this equation,  $N_b = 2560$  blocks are enough to fully utilize the available hardware resources. Increasing  $N_b$  results in more block management overheads, and hence, slightly increases the runtime. Reducing  $N_b$  under-utilizes the GPU, and hence, increases the runtime.

#### 4.8 Multi-round Influence Maximization

Although some previously-proposed algorithms for solving the standard influence maximization problem have successfully been used to resolve some real-world problems [42], [43], it is evident that the standard IM problem cannot

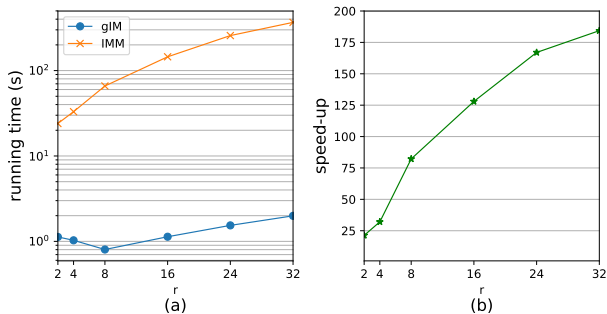


Fig. 7. (a) Runtime of gIM (the proposed solution) and IMM on random scale-free graphs with different densities. (b) The resulting speedup, i.e., the ratio of the two curves in (a).

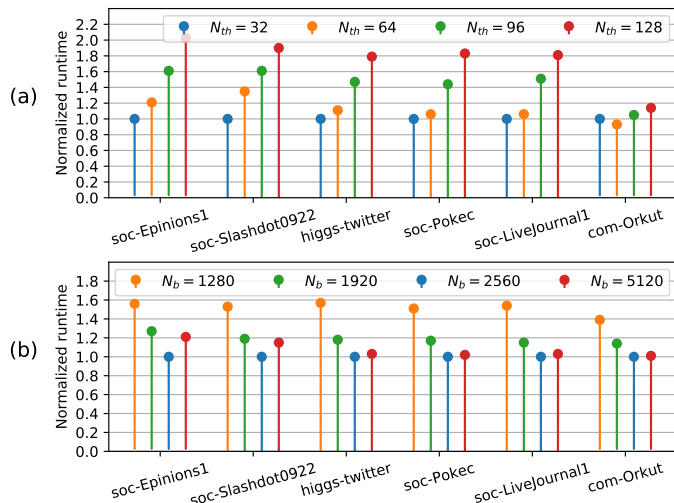


Fig. 8. (a) Runtime of gIM for  $k = 50$  and  $\epsilon = 0.05$ , with different values of  $N_{th}$ , normalized to the runtime of  $N_{th} = 32$  (the selected block size). (b) Runtime of gIM for different values of  $N_b$ , normalized to the runtime of  $N_b = 2560$  (the selected grid size).

grasp all aspects of real-world phenomena. Therefore, many researchers have attempted to propose new variations of the standard IM problem, in order to incorporate new details such as time restriction [18], location awareness [44] and topic awareness [45] into the standard IM problem and make it more realistic for real-world situations. Interestingly, most of the variations to the IM problem can be easily solved by making some subtle modifications to RIS-based algorithms. Therefore, our parallel algorithm not only is able to solve the standard IM problem in a very short time but also is capable of solving different variations of it by applying simple modifications. In order to illustrate this, we introduce a specific variation of the IM problem called multi-round influence maximization [7] and show that our algorithm is able to solve this problem.

In the multi-round influence maximization problem (MRIM), influence propagates in multiple rounds, and the goal is to find a seed set for each round in order to maximize the total number of nodes that have been influenced at least once. Sun et al. [7] proposed three different algorithms to solve the MRIM problem, all of which can be implemented efficiently by using reverse influence sampling. We slightly modify our algorithm to parallelize CR-NAIMM algorithm in [7], which is more computationally expensive than the others. Specifically, after selecting a random node, we ini-

TABLE 3  
Running time of CR-NAIMM algorithm in seconds.

Dataset	GPU	CPU	Speed-up
soc-Epinions1	0.087	2.4	27.59
soc-Slashdot0922	0.066	8.97	135.91
higgs-twitter	0.74	56	75.68
soc-Pokec	2.055	60	29.13
soc-LiveJournal1	10.755	113	10.51
com-Orkut	5.118	591	115.47

tiate a random BFS originating from the selected node as many times as the number of rounds. Also, each element in a random RR set is a tuple of node-id and round number. The rest of our algorithm remains almost intact. We performed experiments using real-world datasets. We set the size of the seed set to  $k = 10$ , the number of rounds to  $T = 5$ , and  $\epsilon = 0.1$ . The runtimes of both serial and parallel algorithms on different datasets are shown in Table 3. We achieve a speed-up of up to  $136\times$  and an average of  $65.72\times$  on six datasets.

## 5 CONCLUSION

In this paper, we proposed gIM, an efficient parallel implementation of IMM algorithm on GPU. Extensive experiments on real social network graphs demonstrated that gIM significantly reduces the runtime of IMM. Moreover, we show that our algorithm is able to solve other variants of the IM problem, only by applying minor modifications.

## REFERENCES

- [1] D. Kempe et al., "Maximizing the spread of influence through a social network," in *ACM SIGKDD international conference on Knowledge discovery and data mining*, 2003, pp. 137–146.
- [2] W. Chen, Y. Yuan, and L. Zhang, "Scalable influence maximization in social networks under the linear threshold model," in *IEEE international conference on data mining*, 2010, pp. 88–97.
- [3] W. Chen, C. Wang, and Y. Wang, "Scalable influence maximization for prevalent viral marketing in large-scale social networks," in *ACM SIGKDD international conference on Knowledge discovery and data mining*, 2010, pp. 1029–1038.
- [4] J. Leskovec et al., "Cost-effective outbreak detection in networks," in *ACM SIGKDD international conference on Knowledge discovery and data mining*, 2007, pp. 420–429.
- [5] A. Goyal et al., "Celf++ optimizing the greedy algorithm for influence maximization in social networks," in *International conference companion on World wide web*, 2011, pp. 47–48.
- [6] X. He and D. Kempe, "Stability of influence maximization," in *ACM SIGKDD international conference on Knowledge discovery and data mining*, 2014, pp. 1256–1265.
- [7] L. Sun et al., "Multi-round influence maximization," in *ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 2249–2258.
- [8] Y. Tang et al., "Influence maximization in near-linear time: A martingale approach," in *ACM SIGMOD International Conference on Management of Data*, 2015, pp. 1539–1554.
- [9] X. Liu et al., "IMGPU: GPU-accelerated influence maximization in large-scale social networks," *IEEE Transactions on Parallel and Distributed Systems*, vol. 25, no. 1, pp. 136–145, 2013.
- [10] K. Pal et al., "Fast GPU-based influence maximization within finite deadlines via node-level parallelism," in *Industrial Conference on Data Mining*, 2017, pp. 151–165.
- [11] M. Minutoli, M. Halappanavar et al., "Fast and scalable implementations of influence maximization algorithms," in *IEEE International Conference on Cluster Computing*, 2019, pp. 1–12.
- [12] M. Minutoli, M. Drocco et al., "cuRipples: Influence maximization on multi-GPU systems," in *ACM International Conference on Supercomputing*, 2020, pp. 1–11.

- [13] G. Gökürtük and K. Kaya, "Boosting parallel influence-maximization kernels for undirected networks with fusing and vectorization," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 5, pp. 1001–1013, 2021.
- [14] H. T. Nguyen *et al.*, "Blocking self-avoiding walks stops cyber-epidemics: a scalable GPU-based approach," *IEEE Transactions on Knowledge and Data Engineering*, vol. 32, no. 7, pp. 1263–1275, 2019.
- [15] S. Pandey *et al.*, "C-saw: a framework for graph sampling and random walk on gpus," *arXiv preprint arXiv:2009.09103*, 2020.
- [16] H. Zeng *et al.*, "Graphsaint: Graph sampling based inductive learning method," *arXiv preprint arXiv:1907.04931*, 2019.
- [17] S. Wen *et al.*, "A sword with two edges: Propagation studies on both positive and negative information in online social networks," *IEEE Transactions on Computers*, vol. 64, no. 3, pp. 640–653, 2015.
- [18] W. Chen *et al.*, "Time-critical influence maximization in social networks with time-delayed diffusion process," in *AAAI Conference on Artificial Intelligence*, 2012, pp. 592–598.
- [19] J. Kim *et al.*, "CT-IC: Continuously activated and time-restricted independent cascade model for viral marketing," *Knowledge-Based Systems*, vol. 62, no. 1, pp. 57–68, 2014.
- [20] M. Xie *et al.*, "Dynadiffuse: A dynamic diffusion model for continuous time constrained influence maximization," in *AAAI Conference on Artificial Intelligence*, 2015, pp. 346–352.
- [21] W. Chen *et al.*, "Information and influence propagation in social networks," *Synthesis Lectures on Data Management*, vol. 5, no. 4, pp. 1–177, 2013.
- [22] C. Borgs *et al.*, "Maximizing social influence in nearly optimal time," in *ACM-SIAM symposium on Discrete algorithms*, 2014, pp. 946–957.
- [23] Y. Tang *et al.*, "Influence maximization: Near-optimal time complexity meets practical efficiency," in *ACM SIGMOD international conference on Management of data*, 2014, pp. 75–86.
- [24] H. T. Nguyen *et al.*, "Stop-and-stare: Optimal sampling algorithms for viral marketing in billion-scale networks," in *International Conference on Management of Data*, 2016, pp. 695–710.
- [25] D. Williams, *Probability with martingales*. Cambridge university press, 1991.
- [26] K. Huang *et al.*, "Revisiting the stop-and-stare algorithms for influence maximization," *Proceedings of the VLDB Endowment*, vol. 10, no. 9, pp. 913–924, 2017.
- [27] F. Busato and N. Bombieri, "BFS-4K: an efficient implementation of BFS for Kepler GPU architectures," *IEEE Transactions on Parallel and Distributed Systems*, vol. 26, no. 7, pp. 1826–1838, 2014.
- [28] H. Liu and H. H. Huang, "Enterprise: breadth-first graph traversal on GPUs," in *International Conference for High Performance Computing, Networking, Storage and Analysis*, 2015, pp. 1–12.
- [29] D. Merrill *et al.*, "Scalable GPU graph traversal," *ACM SIGPLAN Notices*, vol. 47, no. 8, pp. 117–128, 2012.
- [30] A. Gaihre *et al.*, "Xbfs: exploring runtime optimizations for breadth-first search on gpus," in *International Symposium on High-Performance Parallel and Distributed Computing*, 2019, pp. 121–131.
- [31] H. Liu *et al.*, "iBFS: Concurrent breadth-first search on GPUs," in *International Conference on Management of Data*, 2016, pp. 403–416.
- [32] M. Harris *et al.*, "Parallel prefix sum (scan) with cuda," *GPU Gems*, vol. 3, no. 39, pp. 851–876, 2007.
- [33] M. Harris, "Optimizing parallel reduction in cuda," *Nvidia developer technology*, vol. 2, no. 4, p. 70, 2007.
- [34] gIM source code. [Online]. Available: <http://lis.ee.sharif.edu>
- [35] J. Leskovec and A. Krevl, "SNAP Datasets: Stanford large network dataset collection," Jun. 2014. [Online]. Available: <http://snap.stanford.edu/data5>
- [36] E. Cohen *et al.*, "Sketch-based influence maximization and computation: Scaling up with guarantees," in *ACM International Conference on Conference on Information and Knowledge Management*, 2014, pp. 629–638.
- [37] K. Saito *et al.*, "Prediction of information diffusion probabilities for independent cascade model," in *International conference on knowledge-based and intelligent information and engineering systems*, 2008, pp. 67–75.
- [38] A. Goyal *et al.*, "Learning influence probabilities in social networks," in *International conference on web search and data mining*, 2010, pp. 241–250.
- [39] K. Kutzkov *et al.*, "STRIP: stream learning of influence probabilities," in *ACM SIGKDD international conference on Knowledge discovery and data mining*, 2013, pp. 275–283.
- [40] W. Chen *et al.*, "Efficient influence maximization in social networks," in *ACM SIGKDD international conference on Knowledge discovery and data mining*, 2009, pp. 199–208.
- [41] A.-L. Barabási *et al.*, "Mean-field theory for scale-free random networks," *Physica A: Statistical Mechanics and its Applications*, vol. 272, no. 1–2, pp. 173–187, 1999.
- [42] A. Yadav *et al.*, "Using social networks to aid homeless shelters: Dynamic influence maximization under uncertainty," in *International Conference on Autonomous Agents & Multiagent Systems*, vol. 16, 2016, pp. 740–748.
- [43] D. L. Gibbs and I. Shmulevich, "Solving the influence maximization problem reveals regulatory organization of the yeast cell cycle," *PLoS computational biology*, vol. 13, no. 6, p. e1005591, 2017.
- [44] G. Li *et al.*, "Efficient location-aware influence maximization," in *ACM SIGMOD international conference on Management of data*, 2014, pp. 87–98.
- [45] W. Chen *et al.*, "Real-time topic-aware influence maximization using preprocessing," *Computational social networks*, vol. 3, no. 1, p. 8, 2016.



**Soheil Shahrouz** received the B.Sc. degree in electrical engineering from Amirkabir University of Technology, Tehran, Iran in 2018. He is currently working towards the M.Sc. degree in electrical engineering at Sharif University of Technology, Tehran, Iran. His research interests include parallel computing and hardware acceleration.



**Saber Salehkaleybar** received the B.Sc., M.Sc. and Ph.D. degrees in electrical engineering from Sharif University of Technology, Tehran, Iran, in 2009, 2011, and 2015, respectively. He was a postdoctoral researcher in Coordinated Science Lab. (CSL) at University of Illinois, Urbana-Champaign in 2016–2017. He is currently an assistant professor of electrical engineering at Sharif University of Technology, Tehran, Iran. His research interests include distributed systems, machine learning, and causal inference.



**Matin Hashemi** received the B.Sc. degree in electrical engineering from Sharif University of Technology, Tehran, Iran, in 2005, and the M.Sc. and Ph.D. degrees in computer engineering from University of California, Davis, in 2008 and 2011, respectively. He is currently an assistant professor of electrical engineering at Sharif University of Technology, Tehran, Iran. His research interests include algorithm design and hardware acceleration for machine learning and big data applications.