

# On Maximizing Diffusion Speed over Social Networks with Strategic Users

Jungseul Ok, *Student Member, IEEE*, Youngmi Jin, *Member, IEEE*, Jinwoo Shin, *Member, IEEE*,  
and Yung Yi, *Member, IEEE*

**Abstract**—A variety of models have been proposed and analyzed to understand how a new innovation (e.g., a technology, a product, or even a behavior) diffuses over a social network, broadly classified into either of epidemic-based or game-based ones. In this paper, we consider a game-based model, where each individual makes a selfish, rational choice in terms of its payoff in adopting the new innovation, but with some noise. We address the following two questions on the diffusion speed of a new innovation under the game-based model: (i) what is a good subset of individuals to seed for reducing the diffusion time significantly, i.e., convincing them to pre-adopt a new innovation, and (ii) how much diffusion time can be reduced by such a good seeding. For (i), we design near-optimal polynomial-time seeding algorithms for three representative classes of social network models, Erdős-Rényi, planted partition and geometrically structured graphs, and provide their performance guarantees in terms of approximation and complexity. For (ii), we asymptotically quantify the diffusion time for these graph topologies, further derive the seed budget threshold above which the diffusion time is dramatically reduced, i.e., phase transition of diffusion time. Furthermore, based on our theoretical findings, we propose a practical seeding algorithm, called PrPaS (Practical Partitioning and Seeding) and demonstrate that PrPaS outperforms other baseline algorithms in terms of the diffusion speed over a real social network topology. We believe that our results provide new insights on how to seed over a social network depending on its connectivity structure, where individuals rationally adopt a new innovation.

**Index Terms**—Influence Maximization, Clustering, Random Seeding

## I. INTRODUCTION

People are actively using social networks to get new information, exchange new ideas or behaviors, and adopt new innovations. Clearly, it is of significant importance to understand how such information diffuses over time, where diffusion by local interaction is the most prominent feature. Various fields including computer science, economics, and sociology have expressed their interests in understanding diffusion, e.g., [10], [45], [47]. People have first started to propose diffusion models in social network with close relevance to studies with long history on raging epidemic, e.g., SIRS model [28] or interacting particle system, e.g., Ising model [20]. Examples

of such epidemic-based diffusion model also include [15] and [6], often referred to as independent cascade or linear threshold models [26].

Different from epidemic-based models, people often make strategic choices, i.e., an individual adopts a new technology only if the new technology provides sufficient utility, which increases with the number of neighbors who adopt the same technology (i.e., coordination effect) [14], [19], [36], [38]. This is called game-based diffusion model, which is the main focus of this paper. A recent work by Montanari and Saberi [36] addressed the question of the equilibrium behavior as well as the impact of topological properties on diffusion speed. Under the assumption that individuals behave with bounded rationality (i.e., noisy best response dynamic), it has been proved that the number of innovation adopters increases and the innovation finally becomes widespread. However, the diffusion time can be significantly long so that in practice the innovation often diffuses within only a small number of individuals or even become extinct in practice. One of the approaches to reduce the diffusion time is to *seed* some individuals, i.e., convince a subset of individuals to pre-adopt the new innovation, e.g., by providing some incentives to those users.

The problem of maximizing the “degree of diffusion” by properly selecting seeds has been popularly studied in epidemic-based models, often referred to as *influence maximization*, whose major goal is to maximize the number of infected individuals. However, in game-based models, as in e.g., [36], the problem becomes completely different mainly because diffusion is widespread at the equilibrium. Thus, we study how to choose a constrained set of individuals to accelerate the speed of diffusion, which we call *diffusion speed maximization*.

Our main contribution is to (i) propose near-optimal seeding algorithms depending on network structures, (ii) quantify how much the diffusion time can be reduced by the algorithm asymptotically, and (iii) develop a practical seeding algorithm that works for real-world social networks. To this end, we first formulate a diffusion speed maximization problem, say **PI**, as minimizing the notion of typical hitting time which measures the time when every individual adopts the innovation. We discuss its computational challenges mainly stemming from (i) MCMC (Markov Chain Monte Carlo) based estimation and (ii) probabilistic feature of a typical hitting time, which is neither algebraic nor combinatorial (see Section III-C). Therefore, we transform the original problem **PI** into a combinatorial optimization, say **P2**, using the theory of meta-

J. Ok, J. Shin, and Y. Yi are with the Department of Electrical Engineering, KAIST, Daejeon 305-701, Korea (e-mail: {ockjs; jinwoos; yiyung}@kaist.ac.kr). Y. Jin is with KDDI R&D Labs., Saitama, Japan (e-mail: yo-jin@kddilabs.jp).

This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIP) (NRF-2013R1A2A2A01067633) and Institute for Information & communications Technology Promotion (IITP) grant funded by the Korea government (MSIP) (No.B0717-16-0034, Versatile Network System Architecture for Multi-dimensional Diversity).

Part of this work has been presented in the ACM Sigmetrics 2014.

stability of Markov chains [43], which, however, turns out to be computationally intractable as well as difficult to be reduced to a classical NP-hard problem amenable to approximation. For example, the influence maximization in epidemic-based models becomes the submodular maximization in most cases, whose greedy algorithm guarantees constant approximation [26]. However, we found that the optimization  $P2$  is not a submodular problem (see our discussion in Section III-C).

Despite this hardness of  $P2$ , we propose polynomial-time near-optimal algorithms for three representative classes of social network models, Erdős-Rényi, planted partition and geometrically structured graphs, and obtain their provable performance guarantees in terms of approximation ratio as well as complexity. We also analytically quantify the diffusion time taken by the proposed algorithms, where more details are elaborated in what follows:

- **Erdős-Rényi and planted partition graphs.** We show that an arbitrary seeding and a simple seeding proportional to the size of clusters are close to an optimal one with high probability for the dense Erdős-Rényi and planted partition graphs, respectively (see Theorems IV.1 and IV.2). The main technical ingredient for this result is on our concentration inequalities on the so-called ‘energy function’ (see Lemma A.1), which provides the exact approximation qualities of the random seeding via a solution of certain quartic equations. Then it is provably almost optimal via obtaining its approximate close-form solution.
- **Geometrically structured graphs.** For this graph class, including planar and  $d$ -dimensional graphs, we design an algorithmic framework, called PaS (Partitioning and Seeding), and provide a condition, which, if met, provably guarantees good approximation with polynomial complexity (see Theorem IV.3). PaS consists of two phases: (i) partitioning the graph into multiple clusters, and (ii) seeding within each cluster. The proposed PaS framework relies on our finding that the diffusion process in a graph is dominated by the slowest diffusion process among the underlying clusters. Thus, in the partitioning phase, a given graph should be smartly partitioned into the clusters in which a seeding problem becomes tractable (via seeding the “border individuals” among clusters). Then, to minimize the diffusion time, our focus simply becomes a good seed budget allocation to each cluster that minimizes the overall diffusion time. A greedy algorithm is run to achieve the desired budget allocation in the seeding phase.

The practical implications from our theoretical findings are summarized in what follows: Erdős-Rényi, planted partition and geometrically structured graphs represent (a) globally well-connected, (b) locally well-connected with big clusters, and (c) locally well-connected with small clusters, respectively. First, for globally well-connected graphs like Erdős-Rényi graphs, careful seeding is not highly required, because the underlying topological structure such as high symmetry and connectivity does not change significantly even after seeding with a small budget. However, for locally well-connected graphs, it is necessary to intelligently exploit their clustering characteristics, where the network-wide diffusion

time is governed by both intra-cluster diffusion and inter-cluster correlation. As is in sharp contrast to epidemic-based models, in game-based ones, it turns out that in (b) intra-cluster diffusion becomes the dominant factor, as opposed to in (c) where inter-cluster correlation dominantly determines the network-wide diffusion speed. Thus, as described in Sections IV-C and IV-D, for planted partition graphs, we focus only on how to distribute the seed budget to each (big) cluster, while for geometrically structured graphs, the seeds are mainly selected from the border individuals to remove inter-cluster correlation.

Using the new insights from our analysis, we develop a practical seeding algorithm, called PrPaS (Practical Partitioning and Seeding), and demonstrate that PrPaS outperforms other algorithms such as degree-based and random seeding for a real-world social network graph made by a facebook ego network having 4039 nodes and 88234 edges. Interestingly, degree-based seeding, which generally works well in epidemic-based models, performs worst out of all tested algorithms, which shows that smart seeding should be designed depending on how information diffuses over a given network.

## II. RELATED WORK

As discussed earlier, diffusion models in literature can be broadly classified into: (i) epidemic-based [2]–[4], [13], [17], [26], [28] and (ii) game-based [5], [14], [25], [48], depending on how diffusion occurs, i.e., just like a contagious disease or individuals’ strategic choices. In particular, game-based diffusion models [5], [14], [25], [48] adopt a networked coordination game where the payoff matrix appropriately models the value of accepting new technology for the neighbors’ selections, and studied the equilibrium and the dynamics. Especially, Kandori et al. [25] proved that the noisy best response dynamic converges to the equilibrium that the innovation becomes widespread. In [21] and [22], the authors also studied the stationary distribution using a mean field approximation of the game model with finite rationality, called graphical evolutionary game model. Recently, significant attention has been paid to the study of convergence time. In [36], [42], it was shown that in highly connected graph, the convergence becomes slower as opposed to in epidemic models. In [23], the authors showed that the external information such as advertisement on a new technology may slow down diffusion, again on the contrary to in epidemic models [4]. In practice, a small set of influential nodes, called *seeds*, can be convinced to pre-adopt a new technology, which can increase the effect of diffusion. See [12] for motivation in viral marketing, [40] in graph detection, and [29] in computer virus vaccine dissemination. The problem of how to maximize the diffusion effect for both diffusion models are summarized next, where depending on the adopted diffusion model, different problems can be formulated.

*Epidemic-based model.* In [26], [27], the authors addressed the so-called influence maximization problem in linear threshold (LT) and independent cascade (IC) models. In both LT and IC models, each individual has only one chance to infect its neighbors right after its infection. Thus, a main goal is

to maximize the influence spread, i.e., maximize the number of infected individuals. In [26], [27], it was first discussed that the problem is computationally intractable because of #P-completeness in measuring influence spread for a given seed set and NP-completeness in finding the optimal seed set that maximizes influence spread. Using the technique on the submodular set function maximization in [39], they showed that a greedy algorithm achieves at least  $(1 - 1/e - \varepsilon)$  of the optimal influence spread where  $\varepsilon$  represents the inaccuracy of Monte Carlo simulation for measuring the influence spread. Since the Monte-Carlo based measurement does not tend to scale with the network size, the authors in [9] proposed a scalable method called MIA using a tree structure. In [18], a clustering concept is proposed to reduce the computational complexity in measuring the influence spread. In [8], Chen et al. proposed modified LT and IC models by adding contact process, which delays infection chance of the infected individual from its infection. Using the modified models, the authors formulated an influence maximization with time deadline and proposed a greedy algorithm motivated by [26], [27]. In [16], Goyal et al. generalized the influence maximization problem in LT and IC models as an optimization problem with three dimensions: influence spread, seed budget, and time deadline.

*Game-based model.* In [11], [25], [30], [38], the authors considered only the best-response dynamics and studied the conditions (of network topology and the payoff difference between old and new technologies) on the existence of a small seed set, referred as the so-called ‘‘contagion set,’’ under which all individuals adopt new technology. In [32], a noisy best response was considered with objective of maximizing the influence spread by choosing a seed set assuming that there exists a set of ‘‘negative individuals,’’ and a greedy algorithm was proposed with simulation-based evaluations. As discussed in [36], without negative seeding, it is guaranteed to converge to a state where all individuals adopt the new technology. This paper studies a problem of minimizing the convergence time to such an equilibrium under a noisy best response dynamic.

### III. MODEL AND FORMULATION

#### A. Network Model and Coordination Game

**Network model.** We consider a social network as an undirected graph  $G = (V, E)$ , where  $V$  is the set of  $n$  nodes and  $E$  is the set of edges. Each node represents an individual (or a user) and each edge represents a social relationship between two individuals. We let  $N(i)$  be the set of node  $i$ 's neighbors, i.e.,  $N(i) = \{j \in V \mid (i, j) \in E\}$ . We simply use +1 and -1 to refer to new and old technologies, respectively. We are interested in how a new technology diffuses over the network.

**Networked coordination game.** We first consider the famous two-person coordination game whose payoff matrix is given by Table I, where an individual can choose one of new or old technologies, +1 and -1. We make the following practical assumptions on the payoffs. First, there always exists coordination gain, i.e.,  $a > d$  and  $b > c$ . Second, coordination gain becomes larger for the new technology, i.e.,  $a - d > b - c$ .

The two-person coordination game is extended to an  $n$ -person game over  $G$ . We let  $\mathbf{x} = (x_j \in \{-1, +1\} : j \in V)$ ,

TABLE I  
TWO-PERSON COORDINATION GAME

$P$	+1	-1
+1	$(a, a)$	$(c, d)$
-1	$(d, c)$	$(b, b)$

and  $\mathbf{x}_{-i} = (x_j : j \in V \setminus \{i\})$  be the states (i.e., a strategy vector chosen by the entire nodes) of all and those except for  $i$ , respectively. Then, in  $n$ -person game over  $G$ , node  $i$ 's payoff  $P_i(x_i, \mathbf{x}_{-i})$  for the state  $\mathbf{x}$  is modeled to be the aggregate payoff against all of  $i$ 's neighbors, i.e.,

$$P_i(x_i, \mathbf{x}_{-i}) = \sum_{j \in N(i)} P(x_i, x_j), \quad (1)$$

where  $P(x_i, x_j)$  is the payoff from the two-person coordination game, as in Table I. For notational convenience, let  $-1 =$  (resp.  $+1$ ) denote the state where every user adopts  $-1$  (resp.  $+1$ ).

#### B. Diffusion Dynamics

**Seed set.** We consider a continuous time model, where each node updates its strategy whenever its own independent Poisson clock with unit rate ticks. Let  $\mathbf{x}(t) = (x_i(t) : i \in V) \in \{+1, -1\}^V$  be the network state at time  $t$ , representing the strategies of all nodes at time  $t$ . We introduce the notion of *seed set*  $C \subset V$ , where each node in  $C$  is initialized by +1 and does not change its strategy over all time, i.e., for any  $i \in C$ ,  $x_i(t) = +1$  for all  $t \geq 0$ . Next, we describe how each non-seed individual updates its strategy.

**Best response.** As is well-known in game theory, in the best response dynamics, each (non-seed) individual selects a strategy that maximizes its own payoff: a node  $i$  chooses +1, if

$$(a - d)|N^+(i)| \geq (b - c)|N^-(i)| \quad (2)$$

where  $N^+(i)$  and  $N^-(i)$  denote the sets of node  $i$ 's neighbors adopting +1 and -1, respectively. Noting that for a given state  $\mathbf{x}$ ,  $P_i(+1, \mathbf{x}_{-i}) - P_i(-1, \mathbf{x}_{-i})$  represents the payoff difference between when node  $i$  chooses +1 and -1, the best response of node  $i$  is  $\text{sign}(P_i(+1, \mathbf{x}_{-i}) - P_i(-1, \mathbf{x}_{-i}))$ , simply expressed as:

$$\text{sign} \left( h_i + \sum_{j \in N(i)} x_j \right), \quad (3)$$

where  $h_i = h|N(i)|$  and  $h = \frac{a-d-b+c}{a-d+b-c}$

**Noisy best response: Logit dynamics.** In practice, individuals do not always make the ‘‘best’’ decision. We model such behavior by introducing small mutation probability that non-optimal strategy is chosen, often called noisy best response. A version of the noisy best response we focus on in this paper is *logit dynamics* [5], [34], [35], [37] that individuals adopt a strategy according to a distribution of the logit form which allocates larger probability to those strategies delivering larger payoffs. More formally, for the given state  $\mathbf{x}$ , non-seeded node

$i$  chooses the strategy  $y_i \in \{-1, +1\}$  with the following probability:

$$\mathbb{P}_\beta(y_i|\mathbf{x}) = \frac{\exp(\beta y_i K_i(\mathbf{x}))}{\exp(\beta K_i(\mathbf{x})) + \exp(-\beta K_i(\mathbf{x}))}. \quad (4)$$

where

$$K_i(\mathbf{x}) = \frac{1}{2} \left( h_i + \sum_{j \in N(i)} x_j \right).$$

Note that  $(a - d + b - c)y_i K_i(\mathbf{x})$  is the payoff gain for the strategy  $y_i$  instead of  $-y_i$  from (3) and  $(a - d + b - c)$  is removed just for convenient handling of other quantities later. Here, the parameter  $\beta$  represents the degree of user rationality, where  $\beta = \infty$  corresponds to the best response and  $\beta = 0$  lets users update their strategies uniformly at random. When the state changes according to the probability (4) and nodes' independent Poisson clock ticks, the system can be viewed as a continuous Markov chain with the state space  $\mathcal{S}_C = \{z \in \{-1, +1\}^V \mid z_i = 1 \text{ if } i \in C\}$ , recall  $C$  is a given seed set. The dynamics here is also called the Glauber dynamics in the "truncated" Ising model [41], where the truncation occurs due to the existence of hard-coded nodes (i.e., the nodes in the seed set  $C$ ). Then, it is not hard to see that this chain is time-reversible with the following stationary distribution  $\mu_\beta$ :

$$\mu_\beta(\mathbf{x}) \propto \exp(-\beta H(\mathbf{x})),$$

where

$$H(\mathbf{x}) = -\frac{1}{2} \left\{ \sum_{(i,j) \in E} x_i x_j + \sum_{i \in V} h_i x_i \right\} + (1 + 2h)|E|. \quad (5)$$

In the above, the constant term  $(1 + 2h)|E|$  is not necessarily needed to characterize the stationary distribution, but we add due to notational convenience in our proofs. We note that  $-H$  is often referred to as a *potential* function of the  $n$ -person game described in Section III-A and  $H$  is called the *energy* function in literature. Note that from the assumptions on the payoff matrix  $P$ ,  $h$  is strictly positive. Thus  $H$  has the global minimum at all  $+1$  state and the stationary distribution concentrates on all  $+1$  state.

### C. Problem Formulation

Our objective is to find a seed set  $C$  (within some budget constraint) which maximizes the speed of diffusion. To this end, we define a couple of related concepts.

First, a random variable called the hitting time (to the state where all users adopt  $+1$ ) of our system with a seed set  $C$  starting from the initial state  $\mathbf{y} \in \mathcal{S}_C$  defined by:

$$T_+(C, \mathbf{y}) = \inf\{t \geq 0 \mid \mathbf{x}(t) = +\mathbf{1}, \mathbf{x}(0) = \mathbf{y}\}.$$

Using this, we next define the *typical hitting time* to be:

$$\tau_+(C) = \sup_{\mathbf{y} \in \mathcal{S}_C} \inf \left\{ t \geq 0 \mid \mathbb{P}_\beta\{T_+(C, \mathbf{y}) \geq t\} \leq e^{-1} \right\}.$$

This means that with probability  $1 - 1/e$  ( $> 1/2$ ), every node adopts the innovation  $+1$  within time  $\tau_+(C)$ . This typical hitting time has also been used to measure the diffusion speed for a similar model via close relation between hitting and

mixing of the Markov chain, e.g., see [36]. Our goal is to solve the following optimization problem:

$$\begin{aligned} \mathbf{P1.} \quad & \min_{C \subset V} \tau_+(C) \\ & \text{subject to } |C| \leq k, \end{aligned}$$

where  $k$  is the given seed budget.

**Computational challenges of P1.** First, given a seed set  $C$ , the computation of the typical hitting time  $\tau_+(C)$  is a highly non-trivial task, primarily because the hitting time  $T_+(C, \cdot)$  is a random variable decided by the Markov chain of the logit dynamics whose underlying space is exponentially large, i.e.,  $|\mathcal{S}_C|$ . One can use the Markov Chain Monte Carlo (MCMC) method for estimating  $\tau_+(C)$ , which, however, takes at least the mixing time of the Markov chain of the logit dynamic that is typically exponentially large [36]. Even worse, a naive exhaustive search for the optimization **P1** requires computing the typical hitting time  $2^{\Omega(n)}$  times for  $k = \Omega(n)$ . Second, the hardness of the optimization **P1** also comes from the probabilistic definition of the minimizing objective  $\tau_+(C)$ , which is neither algebraic nor combinatorial. Due to these reasons, at a first glance, the optimization **P1** is a highly challenging computational task, similarly to other influence maximization problems in epidemic-based diffusion models, e.g., see [26]. It is not even clear whether the decision version of the optimization **P1** is in the computational class NP.

**Problem formulation via a combinatorial optimization.** To overcome such difficulties, we use the known combinatorial characterization of the typical hitting time  $\tau_+(C)$  from the theory of meta-stability [36], [43], where it was proved that for a given seed set  $C \subset V$ ,

$$\tau_+(C) = \exp(\beta \Gamma^*(C) + o(\beta)), \quad \text{as } \beta \rightarrow \infty, \quad (6)$$

where we refer to  $\Gamma^*(C)$  as the *diffusion exponent* with respect to the seed set  $C$ . In the above,  $\Gamma^*(C)$  is defined as

$$\Gamma^*(C) = \max_{w_0 \in \mathcal{S}_C} \min_{\underline{w}: w_0 \rightarrow +\mathbf{1}} \max_{t < |\underline{w}|} [H(w_t) - H(w_0)]. \quad (7)$$

where the minimization is taken over every possible path  $\underline{w} = (w_0, w_1, \dots, w_T = +\mathbf{1})$  such that for each  $t$ ,  $w_t$  and  $w_{t+1}$  are same except for one coordinate. This implies that  $\Gamma^*$  dominates the exponent of diffusion time  $\tau_+(C)$  for large  $\beta$ . Also,  $\Gamma^*$  can be interpreted as the "energy barrier" along the most probable path to  $+1$ . Two maximums in (7) choose the largest energy difference along a path toward  $+1$ . Then the (middle) minimum in (7) finds a path that has the smallest energy barrier to the ground state  $+1$  so that it is the most probable. In [36], it is known that the minimization of (7) is achieved just at a *monotone* path  $w_0 \prec w_2 \dots \prec w_T$ , i.e., a user is not allowed to take back from  $+1$  to  $-1$ .

The formula (6) provides a tractable approach for bounding  $\tau_+(C)$  through  $\Gamma^*(C)$  and motivated by this, we will focus on the following optimization instead of **P1**:

$$\begin{aligned} \mathbf{P2.} \quad & \min_{C \subset V} \Gamma^*(C) \\ & \text{subject to } |C| \leq k, \end{aligned}$$

where it becomes identical to **P1** as  $\beta \rightarrow \infty$  from (6).

**Further challenges of P2.** Note that it is still challenging to compute  $\Gamma^*(C)$  for a given seed set  $C$  for the following two reasons.

- First, there exist exponentially many monotone paths to consider for the minimization in (7). Characterizations of  $\Gamma^*(C)$  using ‘tilted cut’ and ‘tilted cut-width’ are known, but they are also computationally intractable, e.g., see Section 4.2 of [36]. Nevertheless,  $\Gamma^*(C)$  is defined as a form of combinatorial optimization and potentially more amenable to theoretical analysis than  $\tau_+(C)$ .
- Second, in epidemic-based diffusion models, the influence maximization problem [26], which maximizes the number of infected individuals, could enjoy an algorithmic convenience because of the key feature the objective function turns out to be submodular. Similar convenient features may also be applied to our case, which, if so, would facilitate our analysis significantly. However, unfortunately our objective function  $\Gamma^*(\cdot)$  is neither supermodular nor submodular, as proved by a counter-example in the supplemental material, which motivates our study of a different kind of approximation techniques.

#### IV. MAIN RESULT

In this section, we describe our polynomial-time approximation algorithms for the seeding problem **P2**. Each algorithm provides the guideline on which nodes should be seeded for fast diffusion over a game-based diffusion model for each of three graph classes, which is classified by the criterion on how globally and locally well-connected nodes are. To this end, we first introduce the following notion of ‘‘approximate solution’’.

##### A. $(\gamma, \delta)$ -Approximate Solution

**Definition IV.1.** A seed set  $C \subset V$  with  $|C| \leq k$  is called a  $(\gamma, \delta)$ -approximate solution of the seeding problem **P2** if

$$\Gamma^*(C) \leq \gamma \cdot \min_{C': |C'| \leq \delta k} \Gamma^*(C'),$$

where  $\gamma \geq 1$  and  $\delta \leq 1$ .

The parameters  $\gamma$  and  $\delta$  measure the quality of an approximate solution, quantifying the degrees of suboptimality in *objective value* and *budget*, respectively. One can observe that the solution with  $(\gamma, \delta) = (1, 1)$  corresponds to an optimal solution. Thus the distance between  $(\gamma, \delta)$  and  $(1, 1)$  quantifies the performance loss of  $(\gamma, \delta)$ -approximate solution comparing to the optimal solution. In what follows, we present the characteristics of approximate solutions in three graph classes which have different topological structures in terms of connectivity and the degree of clustering.

##### B. Erdős-Rényi Graphs

We first consider the popular *Erdős-Rényi (ER) graph*, denoted by  $G_{ER}(n, p)$ , which is a random graph on  $n$  nodes such that every node pair has an edge with probability  $p$ . Let  $\lambda = np$ , roughly corresponding to the average number of neighbors per node. For ER graphs, we obtain the following result, whose proof is presented in Appendix A.

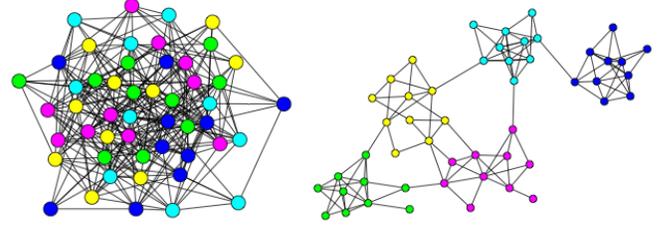


Fig. 1. An instance of ER-graph (left) and planted partition graph (right). Source: Lecture note of the network analysis and modeling course in Santa Fe Institute [1].

**Theorem IV.1.** Consider an ER graph  $G_{ER}(n, p)$  with  $\lambda = \Omega(1)$ . For the seed budget  $k = \kappa n$  with  $\kappa < \left(\frac{1-h}{2} - \frac{h}{\sqrt{\lambda}}\right)$ , every  $C \subset V$  with  $|C| = k$  is almost surely a  $(\gamma, \delta)$ -approximate solution as  $n \rightarrow \infty$ , where

$$\delta = 1, \quad \gamma = 1 + \frac{2}{\frac{\sqrt{\lambda}}{2(1-h^2)} \left(\frac{1-h}{2} - \kappa\right)^2 - 1} \quad (8)$$

and

$$\Gamma^*(C) = pn^2 \left[\frac{1-h}{2} - \kappa\right]_+^2 + o(pn^2)^1. \quad (9)$$

Three interpretations from Theorem IV.1 are in order. First, as in (8), for the relatively dense and (globally) well-connected ER graph, formally for the case  $\lambda = \omega(1)$ , an *arbitrary* seed set  $C$  is, somewhat surprisingly, an almost optimal solution, i.e.,  $(\gamma, \delta) \rightarrow (1, 1)$  as  $n$  grows. The near optimality of an arbitrary seeding in the dense ER graphs mainly comes from globally symmetric connectivities of nodes which makes the influencing effect by each node indistinguishable. Therefore no careful seeding mechanism is necessary for this globally well-connected graph. Second,  $\Gamma^*(\cdot)$  in (9) when  $\kappa < \frac{1-h}{2}$  implies *diminishing return* of adding more seed budget. Third, one needs a seed budget larger than  $\left(\frac{1-h}{2}\right)n$  in order to have an order-wise reduction in  $\Gamma^*$ .

##### C. Planted Partition Graphs

Second, we consider a generalized version of ER graphs and study the so-called *planted partition graph*<sup>2</sup>, which we denote by  $G_{PP}(n, p, q, \omega)$ . It is a popular model, e.g., [7], for social networks with big communities (also called clusters); Given a disjoint partition of the clusters  $\{V_1, \dots, V_m\}$ , with  $\bigcup_{l=1}^m V_l = V$ , let the fraction of nodes in the graph that belongs to a cluster  $l$  be  $\omega_l = |V_l|/n$  where  $\omega = (\omega_1, \dots, \omega_m) \in (0, 1)^m$ . For a pair of  $i, j \in V$ , an edge  $(i, j)$  exists between them with probability  $p$  for the nodes  $i$  and  $j$  if  $i, j$  belong to a same cluster, and with probability  $q < p$ , otherwise. We obtain the following result, whose proof is presented in Appendix B.

**Theorem IV.2.** Consider a planted partition graph  $G_{PP}(n, p, q, \omega)$  with  $q < p = \Theta(1)$ . For the seed budget  $k = \kappa n$  with  $\kappa < \frac{1-h}{2}$  and any small constant  $\varepsilon > 0$ , every  $C \subset V$  such that

$$C \in \arg \min_{\{C': |C'| \leq k\}} \max_{1 \leq l \leq m} \left( \frac{1-h}{2} |V_l| - |C' \cap V_l| \right) \quad (10)$$

<sup>1</sup>Here  $[x]_+ = \max\{x, 0\}$ . We note that the quantification of  $\Gamma^*(C)$  in (9) holds for all  $\kappa \in [0, 1]$ .

<sup>2</sup>This is often referred to as the stochastic block model.

is almost surely a  $(\gamma, \delta)$ -approximate solution as  $n \rightarrow \infty$ , where

$$\delta = 1, \quad \gamma = 1 + \frac{2}{p\xi^2/(q + \varepsilon) - 3} \quad (11)$$

and

$$\Gamma^*(C) = pn^2\xi^2 + o(pn^2)^3 \quad (12)$$

with

$$\xi = \min_{\{\nu \in [0,1]^m: |\nu|_1 \leq \kappa\}} \max_{1 \leq l \leq m} \left[ \frac{1-h}{2} \omega_l - \nu_l \right]_+.$$

In particular, for the homogeneous cluster size, i.e.,  $\omega = (\frac{1}{m}, \dots, \frac{1}{m})$ ,

$$\xi = \frac{1}{m} \left[ \frac{1-h}{2} - \kappa \right]_+.$$

Theorem IV.2 provides a guideline on how to allocate seeds, coming from solving a “simple” min-max optimization (10) whose computational complexity is  $O(1)$  ( $m$  is a given constant and only cardinality of  $C' \cap V_l$  is necessary in computing the min-max solution). Intuitively the resulting seed set  $C$  in (10) allocates more seeds to bigger clusters, and intra-cluster seeding does not have to be carefully chosen. More formally, any seed set  $C$  with such an allocation is an almost optimal solution, regardless of how to seed inside each cluster if the graph is locally well-connected with big clusters whose sizes scales with respect to  $n$  and the number of inter-cluster edges is ignorable comparing to intra-cluster ones, i.e.,  $|V_l| = \Omega(n)$  and  $p/q = \omega(1)$ . For locally well-connected graphs with clusters, it is necessary to intelligently exploit their clustering characteristics, where the network-wide diffusion time is governed by both (a) intra-cluster diffusion and (b) inter-cluster correlation. In locally well-connected with big clusters such as  $G_{pp}(n, p, q, \omega)$ , the intra-cluster diffusion  $\Gamma^*$  in each  $V_l$  dominates the inter-cluster correlation between  $V_l$  and  $V_{l'}$  with  $l \neq l'$ . Hence it suffices to focus on how much seed budget is distributed to each (big) cluster depending on its size. As in ER graphs, we obtain the quantification of  $\Gamma^*$  in (12), which implies that the minimum seed budget to have the order-wise reduction of  $\Gamma^*$  is  $\frac{1-h}{2}$ , and we have the diminishing return of adding seed budget.

#### D. Geometrically Structured Graphs

Third, we consider locally well-connected graphs with *small* clusters. Those graphs include geometrically structured graphs such as planar and  $d$ -dimensional graphs. In these graphs, the *inter-cluster correlation* dominantly determines the network-wide diffusion speed, and hence seeds should be selected with goal of removing the correlation. Different from the earlier two types of graphs, we here take an approach that rather than studying a particular type of graph, we first propose an algorithm and then study a sufficient condition that ensures good diffusion performance and is satisfied in the well-known geometrically structured graphs such as planar and  $d$ -dimensional graphs.

<sup>3</sup>We note that the quantification of  $\Gamma^*(C)$  in (12) holds for all  $\kappa \in [0, 1]$ .

---

**Input:** Graph  $G = (V, E)$  and seed budget  $k$

**Output:** Seed set  $C^{\text{PaS}}$

---

#### 1. Partitioning phase.

Construct a partition  $\{V_l : l = 0, 1, \dots, m\}$ , where  $V_0$  separates others, i.e., there is no edge between  $V_l$  and  $V_{l'}$  for all  $l \neq l' \geq 1$ ,

$$\bigcup_{l=0}^m V_l = V \quad \text{and} \quad V_l \cap V_{l'} = \emptyset \quad \text{for all } l \neq l' \geq 0.$$

Each component  $V_l$  becomes a cluster, i.e.,  $m + 1$  is the number of clusters found in this phase.

#### 2. Seeding phase.

**2-1.** Seed  $V_0$ , i.e.,  $C \leftarrow V_0$ .

**2-2.** Cluster selection.

Find the slowest cluster  $1 \leq l^* \leq m$  such that

$$l^* \in \arg \max_{1 \leq l \leq m: |C_l| < |V_l|} \Gamma^*(G_l, C_l \cup V_0),$$

where  $G_l$  is the subgraph induced by  $V_l \cup V_0$  and  $C_l$  is the set of seeds in  $V_l$ , i.e.,  $C_l = C \cap V_l$ .

**2-3.** Seed selection in the selected cluster.

Find an optimal seed set  $D$  in  $V_{l^*}$  with increased seed budget such that

$$D \in \arg \min_{D' \subset V_{l^*}: |D'| = |C_{l^*}| + 1} \Gamma^*(G_{l^*}, D' \cup V_0).$$

**2-4.** Update  $C \leftarrow (C \setminus C_{l^*}) \cup D$ , and repeat the steps **2-2**, **2-3**, and **2-4** whenever  $|C| < k$ .

#### 3. Terminate. Output $C$ .

---

**Algorithm 1:** PaS (Partitioning and Seeding) Algorithm

One of achieving the goal of removing inter-cluster correlation would be to seed the *border nodes* among small clusters. Motivated by this, we design a generic algorithm, called PaS (Partitioning and Seeding) (see **Algorithm 1** for a formal description) for finding good seeds. As the name implies, PaS has two phases: (i) partitioning and (ii) seeding, as elaborated in what follows.

(i) *Partitioning phase:* In this phase, PaS finds a partitioning with, a finite number of node clusters, where the number of clusters are chosen appropriately, depending on the underlying graph topologies. We call  $V_0$  *separator cluster* since after removing  $V_0$ , no edge exists between different clusters  $V_l, V_{l'}$  for all  $l \neq l' \geq 1$ . Except for the separator cluster  $V_0$ , which will be used as the initial seed set, PaS will find the seeds contained in each cluster by the seeding phase.

(ii) *Seeding phase:* In this phase, PaS runs in multiple rounds, where it starts from the initial seed set  $V_0$  (step **2-1**) and the seed set  $C$  increases by one in each round, until the entire seed set size becomes the target budget  $k$ . Let  $G_l$  and  $C_l$  be the subgraph induced and the seed contained, by  $l$ -th cluster  $V_l$ , respectively. The seeding phase consists of two sub-phases (a) partition selection and (b) seed selection. In (a), PaS finds the partition  $l^*$  that has the slowest diffusion time with the current seed set  $C_l$  (step **2-1**). In (b), for the chosen partition  $l^*$ , we

replace the existing seeds  $C_{l^*}$  by completely new set of seeds whose size increases by one. The new seed set is chosen such that the diffusion time in cluster  $l^*$  is minimized (step 2-2). Finally, the temporary seed  $C$  is updated by a new seed set in cluster  $l^*$ , which is repeated until  $|C| = k$  (steps 2-3 and 2-4). The choices of partition  $\{V_0, V_1, \dots, V_m\}$  in step 1 determines the performance and complexity of the PaS algorithm, where we will consider different choices for different social networks for rigorous analysis.

Now, we are ready to present the performance guarantees of the PaS algorithm. To that end, we introduce a notation:  $E_l$  is the edge set of the subgraph induced by  $V_l \cup V_0$ , where  $V_l$  is the  $l$ -th cluster resulting from the partitioning phase.

**Theorem IV.3.** *For given graph  $G = (V, E)$  and seeding budget  $k = \kappa n$  with  $\kappa \in (0, 1)$ , suppose that  $\{V_l : l = 0, 1, \dots, m\}$  in the partitioning phase of the PaS algorithm has the following condition:*

---

For some  $\varepsilon \in (0, 1)$ ,

$$|V_0| \leq \varepsilon n \quad \text{and} \quad |V_l| = O(1), \quad \text{for all } l = 1, \dots, m. \quad (13)$$


---

Then, the PaS algorithm outputs a  $(1, 1 - \frac{\varepsilon}{\kappa})$ -approximation solution  $C$  such that

$$\Gamma^*(C) = O(1), \quad (14)$$

and its seeding phase takes  $O(n^2)$  time.

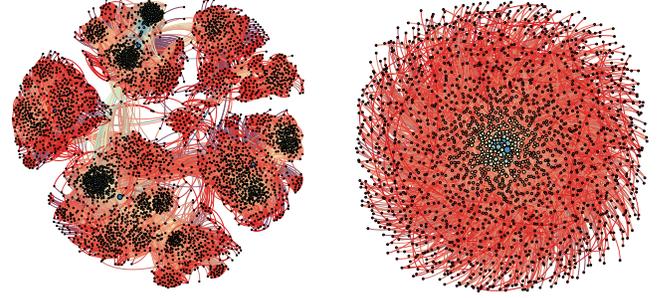
The proof of Theorem IV.3 is presented in Appendix C. Theorem IV.3 implies that if there exists an algorithm finding a ‘good’ partition (i.e.,  $|V_0|/n \leq \varepsilon$  for some small  $\varepsilon > 0$ ) with small clusters (i.e.,  $|V_l| = O(1)$ ), as specified in the condition (13), the PaS algorithm outputs an almost optimal solution. Note that  $V_0$  corresponds to the set of border nodes among clusters. This condition (13) does not always hold. However, for the following classes of social networks, polynomial-time algorithms are known for computing such a partition satisfying the condition for any  $\varepsilon = \Omega(1)$  [24].<sup>4</sup>

- ***d-dimensional Graph.*** A graph is called a  $d$ -dimensional graph, denoted by  $G_{dD}(n, d, D, R)$ , if each node  $i$  can be embedded to a position  $\pi_i$  in  $\mathbb{R}^d$  such that  $(i, j) \in E$  implies that the Euclidean distance between  $\pi_i$  and  $\pi_j$  is less than  $R$  and any cube of volume of  $B$  contains at most  $D \cdot B$  nodes, where  $d, D, R = O(1)$ .
- ***Planar Graph.*** A planar graph, denoted by  $G_{PL}(n, \Delta)$ , can be drawn on the plane without intersection of edges except nodes which is endpoints of edges and its maximum degree  $\Delta = O(1)$ .

Therefore, we can state the following corollary of Theorem IV.3.

**Corollary IV.1.** *For a  $d$ -dimensional graph  $G_{dD}(n, d, D, R)$  or planar graph  $G_{PL}(n, \Delta)$  and seeding budget  $k = \kappa n$  with*

<sup>4</sup>In fact, the author [24] considers polynomially-growing graphs and minor-excluded graphs, where  $d$ -dimensional graphs and planar graphs are their special cases, respectively.



(a) PPfacebook consisting of 4,039 users and 88,234 edges and having average clustering coefficient 0.6055 and degree distribution fit into power law distribution with exponent 1.18. (b) PLfacebook consisting of 1,899 users and 20,296 edges and having average clustering coefficient 0.1385 and degree distribution fit into power law distribution with exponent 1.334.

Fig. 2. Blueprints of PPfacebook [33] and PLfacebook [44].

$\kappa \in (0, 1)$ , there exists a polynomial-time<sup>5</sup> algorithm such that it outputs a  $(1, 1 - \varepsilon)$ -approximation solution  $C$  such that  $\Gamma^*(C) = O(1)$  for any  $\varepsilon \in (0, 1)$ .

We note that even if a geometrically structured graph satisfying (13) has extremely slow diffusion without seeding, where the diffusion time is exponentially increasing with respect to graph size, i.e.,  $\Gamma^*(G) = \omega(1)$ , the diffusion time can be significantly reduced by seed set  $C$  from PaS algorithm, i.e.,  $\Gamma^*(C) = O(1)$ . Further, if the graph is a  $d$ -dimensional graph or planar graph, the amount of seeds for having  $\Gamma^*(C) = O(1)$  is arbitrarily small, i.e.,  $|C| = \varepsilon n$  for any given  $\varepsilon \in (0, 1)$ . For example, consider a star-like graph with a center node surrounded by  $(n-1)$  nodes with the center node having degree  $(n-1)$  and the other  $n-1$  having degree 1. Then, it is a planar graph which can be partitioned into  $n$  partitions consisting of a node by letting the center node be separator cluster  $V_0$ . Hence, by seeding the center node only, we have  $\Gamma^*(C) = 0 = O(1)$ , while without seeding, we have  $\Gamma^*(G) = \frac{1-h}{2}(n-1) = \Omega(n)$ . In Section V, we will show that PaS algorithm shows indeed a good performance for a real social graph, showing its practical value.

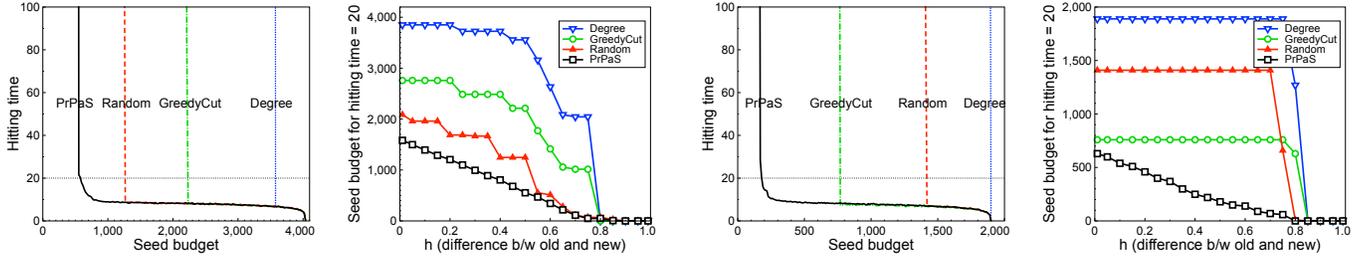
## V. PRACTICAL SEEDING AND SIMULATION RESULTS

In this section, we perform simulations using a real social network graph and show how our theoretical findings can be applied to the diffusion speed maximization in practice. Guided by the implications drawn from the analytic results based on three graph classes, we propose a practical, heuristic seeding algorithm and show how it performs, compared to other seeding algorithms.

### A. Setup

**Real-world social networks.** We use two topology data sets extracted from the social network among Facebook users originally obtained in [33] and [44]. Each data set forms an undirected graph where each node corresponds to a Facebook account and an edge corresponds to a social relationship

<sup>5</sup>It is a polynomial with respect to  $n$ , but may be exponential with respect to  $1/\varepsilon$ .



(a) Hitting time with varying seed budget and  $h = 0.5$  in PPfacebook. (b) Threshold(20) with varying  $h$  in PPfacebook. (c) Hitting time with varying seed budget and  $h = 0.5$  in PLfacebook. (d) Threshold(20) with varying  $h$  in PLfacebook.

Fig. 3. Simulation results on the performance of different algorithms in PPfacebook and PLfacebook.

(called “FriendList”) in Facebook. We name the graph from [33] PPfacebook, and the graph from [44] PLfacebook, whose graphical presentations are given in Figures 2(a) and 2(b), respectively. Namely, a clustering structure is more observed in PPfacebook, whereas a power law degree distribution is prominent in PLfacebook<sup>6</sup>.

**Parameters.** We use  $\beta = 10$  for the degree of rationality and vary  $h$  from 0 to 1 to investigate the impact of the difference between new and old technologies. We are interested in the regime of users are sufficiently rational and hence we tested various values of  $\beta$  larger than 10. They resulted in a similar trend and thus we just report the case of  $\beta = 10$  in this paper due to space limitation.

**Tested seeding algorithms.** We compare the performance of the following four algorithms, each of which is described in what follows.

- **Degree.** This choose  $k$  nodes in the order of their degrees.
- **GreedyCut.** This runs  $k$  iterations where at each iteration a node with the maximum number of edges is selected, and then the node and its edges are removed from the temporal graph.
- **Random.** This selects  $k$  nodes uniformly at random.
- **PrPaS.** This first identifies the partition, say  $\{V_1, \dots, V_m\}$ , from the given graph using the random-walk based approach [46], and then generates a seed set  $C$  whose per-cluster portion is kept equal, i.e.,  $|C \cap V_l|/|V_l| = k/n$  for  $l = 1, \dots, m$ . In each cluster, seeds are selected uniformly at random.

Inspired by our theoretical findings, we design PrPaS (Practical PaS) which can work in general graph without any prior information. and we will show its superiority by comparing it to the first three baseline algorithms. According to our analysis, we prefer a “good” partition consisting of locally well-connected clusters. We employ the random walk based partitioning scheme, borrowed from [46]. Then, with the resulting partition, we just balance the fraction of seeds in each cluster, so that the entire seed budget is allocated in proportion

<sup>6</sup>Our calculation states that the clustering coefficients of PPfacebook and PLfacebook are 0.606 and 0.139, respectively, and the degree distributions of those two graphs are fit into power law distributions with exponent 1.18 and 1.34, respectively.

to the cluster size. This can be regarded as a practical version of PaS in Section IV-D in the sense that (i) it works without explicit knowledge of  $h$ , which may be hard to be quantified in practice, and (ii) partitioning based on simple random walks is scalable and applicable to large-scale social networks. We assume the case when  $h$  is unknown, thus exact computation of  $\Gamma^*$  inside each cluster is infeasible, which is reason why we use per-cluster random seeding.

## B. Results

We compare the algorithms by the minimum seed budget with which the system hits the state +1 in a reasonable time. For convenience, we call this minimum seed budget for a given hitting time  $x$ ,  $\text{Threshold}(x)$ .

We first understand how hitting time changes with varying seed budgets. As shown in Figures 3(a) and 3(c), we observe that there exists a phase transition that the hitting time blows up after some seed budget, which differs across the algorithms. Due to space limitations we omit the results for other  $h$  values, where we observe a similar behavior with different seed budget leading to the hitting time blow-up. This phase transition is due to the existence of “bottleneck clusters”, without which diffusion would become fast. Hence, the seeding quality can be evaluated by how efficiently such bottleneck clusters are removed by the seeding. In our setting, we see that time 20 (a horizontal line in Figures 3(a) and 3(c)) can be a reasonable required hitting time to differentiate the tested algorithms. Hitting time 20 may or may not be the required time by seeders, because the absolute time should be computed by the duration of unit time and unit time can be different how actively individuals interact with each other over the given social network.

To investigate how the tested algorithms perform, we choose the time 20 as a given target hitting time, and compare  $\text{Threshold}(20)$  for all tested algorithms with varying  $h$ , whose results are shown in Figures 3(b) and 3(d). We first observe that across all ranges of  $h$ , PrPaS has the lowest threshold budget, performing significantly better than others. It is natural that for significantly high  $h$  (e.g., larger than 0.7) the performance difference is marginal because diffusion should occur very fast irrespective of the quality of seeding. In addition, PrPaS has linear curves of  $\text{Threshold}(20)$  with respect to  $h$ . This coincides with the analysis of  $\Gamma^*(C)$  in Theorems IV.1 and IV.2 where an order-wise reduction of diffusion time requires seed budget of  $\frac{1-h}{2}n$  at least.

In PPfacebook having a cluster structure, Random outperforms Degree and GreedyCut, because uniformly random seed selection allocates more seeds in larger clusters in the *average* sense. PrPaS performs much better than Random because PrPaS performs further optimization by considering the clustering and connectivity structure of the underlying graph. Conversely, in PLfacebook, seeding separator cluster becomes more important rather than the balanced seeding over clusters due to the skewed degree distribution. Hence GreedyCut, which prioritizes selecting seeds who separates graph, significantly outperforms Degree and Random. However PrPaS is superior to GreedyCut since PrPaS not only finds separator cluster but also balances the portion of seeds in each cluster. We provide additional experimental result with a larger data set in the supplemental material due to the limited space. The result also shows that PrPaS outperforms others.

## VI. CONCLUSION

In this paper, we have addressed the following two questions on the diffusion speed of a new innovation under a noisy game-based model: (i) what is a good subset of individuals to seed for reducing the diffusion time significantly, and (ii) how much diffusion time can be reduced by such a good seeding. For (i), we design near-optimal polynomial-time seeding algorithms for three representative classes of social network models, Erdős-Rényi, planted partition and geometrically structured graphs. Our analysis first implies that for globally well-connected graphs, a careful seeding is not necessary. However, for locally well-connected graphs, their clustering characteristics should be appropriately utilized for strengthening the seeding effect, where seeding inside and across clusters are of critical importance for the graphs having a mixture of big and small clusters, respectively. For (ii), we asymptotically quantify the diffusion time for these graph topologies, further derive the seed budget threshold above which the diffusion time experiences the phase transition of diffusion time.

## REFERENCES

- [1] <http://tuvalu.santafe.edu/~aaronc/courses/5352/>.
- [2] R. M. Anderson and R. M. May. *Infectious Diseases of Humans*. Oxford University Press, 1991.
- [3] N. T. J. Bailey. *The Mathematical Theory of Infectious Diseases and Its Applications*. Hafner Press, 1975.
- [4] S. Banerjee, A. Gopalan, A. K. Das, and S. Shakkottai. Epidemic spreading with external agents. *arXiv preprint arXiv:1206.3599*, 2012.
- [5] L. E. Blume. The statistical mechanics of strategic interaction. *Games and Economic Behavior*, 5(3):387–424, 1993.
- [6] D. Chakrabarti, Y. Wang, C. Wang, J. Leskovec, and C. Faloutsos. Epidemic thresholds in real networks. *ACM Transactions on Information and System Security*, 10(4):13:1–13:25, 2005.
- [7] K. Chaudhuri, F. C. Graham, and A. Tsiatas. Spectral clustering of graphs with general degrees in the extended planted partition model. *Journal of Machine Learning Research-Proceedings Track*, 23:35–1, 2012.
- [8] W. Chen, W. Lu, and N. Zhang. Time-critical influence maximization in social networks with time-delayed diffusion process. In *Proc. of AAAI*, 2012.
- [9] W. Chen, C. Wang, and Y. Wang. Scalable influence maximization for prevalent viral marketing in large-scale social networks. In *Proc. of ACM SIGKDD*, 2010.
- [10] J. S. Coleman, E. Katz, H. Menzel, et al. *Medical innovation: A diffusion study*. Bobbs-Merrill Company New York, NY, 1966.
- [11] E. Coupechoux and M. Lelarge. Impact of clustering on diffusions and contagions in random networks. In *Proc. of IEEE NetGCooP*, 2011.
- [12] P. Domingos and M. Richardson. Mining the network value of customers. In *Proc. of ACM SIGKDD*, 2001.
- [13] M. Draief and A. Ganesh. A random walk model for infection on graphs: spread of epidemics rumours with mobile agents. *Discrete Event Dynamic Systems*, 21(1):41–61, Mar. 2011.
- [14] G. Ellison. Learning, local action, and coordination. *Econometrica*, 61(5):1047–1071, Sep. 1993.
- [15] A. Ganesh, L. Massouli, and D. Towsley. The effect of network topology on the spread of epidemics. In *Proc. of IEEE Infocom*, 2003.
- [16] A. Goyal, F. Bonchi, L. V. S. Lakshmanan, and S. Venkatasubramanian. On minimizing budget and time in influence propagation over social networks. *Social Network Analysis and Mining*, pages 1–14, 2013.
- [17] H. W. Hethcote. The mathematics of infectious diseases. *SIAM review*, 42(4):599–653, 2000.
- [18] J. Hu, K. Meng, X. Chen, C. Lin, and J. Huang. Analysis of influence maximization in large-scale social networks. In *Proc. of ACM SIGMETRICS*, 2013.
- [19] N. Immorlica, J. Kleinberg, M. Mahdian, and T. Wexler. The role of compatibility in the diffusion of technologies through social networks. In *Proc. of ACM EC*, 2007.
- [20] E. Ising. Beitrag zur theorie des ferromagnetismus. *Zeitschrift für Physik A Hadrons and Nuclei*, 31(1):253–258, 1925.
- [21] C. Jiang, Y. Chen, and K. J. R. Liu. Evolutionary dynamics of information diffusion over social networks. *IEEE Transactions on Signal Processing*, 62(17):4573–4586, 2014.
- [22] C. Jiang, Y. Chen, and K. J. R. Liu. Graphical evolutionary game for information diffusion over social networks. *IEEE Journal of Selected Topics in Signal Processing*, 8(4):524–536, 2014.
- [23] Y. Jin, J. Ok, Y. Yi, and J. Shin. On the impact of global information on diffusion of innovations over social networks. In *Proc. of IEEE Infocom NetSci*, 2013.
- [24] K. Jung. *Approximate inference: decomposition methods with applications to networks*. PhD thesis, Massachusetts Institute of Technology, 2009.
- [25] M. Kandori, G. J. Mailath, and R. Rob. Learning, mutation, and long run equilibria in games. *Econometrica*, 61(1):29–56, Jan. 1993.
- [26] D. Kempe, J. Kleinberg, and E. Tardos. Maximizing the spread of influence through a social network. In *Proc. of ACM SIGKDD*, 2003.
- [27] D. Kempe, J. Kleinberg, and E. Tardos. Influential nodes in a diffusion model for social networks. In *Proc. of Intl. Colloq. on Automata, Languages and Programming*, pages 1127–1138. Springer, 2005.
- [28] W. O. Kermack and A. G. McKendrick. Contributions to the mathematical theory of epidemics. ii. the problem of endemicity. In *Proc. of the Royal society of London. Series A*, 1932.
- [29] M. Lelarge. Coordination in network security games. In *Proc. of INFOCOM*, 2012.
- [30] M. Lelarge. Diffusion and cascading behavior in random networks. *Games and Economic Behavior*, 75(2):752–775, 2012.
- [31] J. Leskovec, J. Kleinberg, and C. Faloutsos. Graph evolution: Densification and shrinking diameters. *ACM Transactions on Knowledge Discovery from Data*, 1(1):2, 2007.
- [32] S. Liu, L. Ying, and S. Shakkottai. Influence maximization in social networks: An ising-model-based approach. In *Proc. of IEEE Allerton*, 2010.
- [33] J. McAuley and J. Leskovec. Learning to discover social circles in ego networks. In *Proc. of NIPS*, 2012.
- [34] D. McFadden. *Chapter 4: Conditional logit analysis of qualitative choice behavior in Frontiers in Econometrics*. Academic Press, New York, 1973.
- [35] R. D. McKelvey and T. R. Palfrey. Quantal response equilibria for normal form games. *Games and Economic Behavior*, 10(1):6–38, 1995.
- [36] A. Montanari and A. Saberi. Convergence to equilibrium in local interaction games. In *Proc. of IEEE FOCS*, 2009.
- [37] D. Mookherjee and B. Sopher. Learning behavior in an experimental matching pennies game. *Games and Economic Behavior*, 7(1):62–91, 1994.
- [38] S. Morris. Contagion. *The Review of Economic Studies*, 67(1):57–78, 2000.
- [39] G. L. Nemhauser, L. A. Wolsey, and M. L. Fisher. An analysis of approximations for maximizing submodular set functions. *Mathematical Programming*, 14(1):265–294, 1978.
- [40] P. Netrapalli and S. Sanghavi. Learning the graph of epidemic cascades. In *Proc. of ACM SIGMETRICS*, 2012.

- [41] E. J. Neves and R. H. Schonmann. Behavior of droplets for a class of glauber dynamics at very low temperature. *Probability theory and related fields*, 91(3-4):331–354, 1992.
- [42] J. Ok, J. Shin, and Y. Yi. On the progressive spread over strategic diffusion: Asymptotic and computation. In *Proc. of INFOCOM*, 2012.
- [43] E. Olivieri and M. E. Vares. *Large Deviations and Metastability*. Cambridge University Press, 2005.
- [44] T. Opsahl and P. Panzarasa. Clustering in weighted networks. *Social networks*, 31(2):155–163, 2009.
- [45] E. M. Rogers. *Diffusion of innovations*. Free Press, 1962.
- [46] M. Rosvall and C. T. Bergstrom. Maps of random walks on complex networks reveal community structure. *PNAS*, 105(4):1118–1123, 2008.
- [47] D. Strang and S. A. Soule. Diffusion in organizations and social movements: From hybrid corn to poison pills. *Annual review of sociology*, 24:265–290, 1998.
- [48] H. P. Young. *Individual strategy and social structure: An evolutionary theory of institutions*. Princeton University Press, 1998.

## APPENDIX

### A. Proof of Theorem IV.1

We present the proof of Theorem IV.1 in this section. Consider Erdős-Rényi graph  $G_{ER}(n, p)$  and seed budget  $k = \kappa n$ . We will first show that for  $\kappa < \left(\frac{1-h}{2} - \frac{h}{\sqrt{\lambda}}\right)$  and  $\lambda = \Omega(1)$ , the following event occurs almost surely as  $n \rightarrow \infty$ :

$$\mathcal{L} \leq \frac{\Gamma^*(C)}{\lambda n} \leq \mathcal{U}, \quad \text{for all } C \text{ with } |C| = k, \quad (15)$$

where

$$\mathcal{L} = \left(\frac{1-h}{2} - \kappa\right)^2 - \frac{2(1-h^2)}{\sqrt{\lambda}},$$

$$\mathcal{U} = \left(\frac{1-h}{2} - \kappa\right)^2 + \frac{2(1-h^2)}{\sqrt{\lambda}}.$$

The above inequality (15) implies that  $\Gamma^*(C)$  is highly concentrated on the interval  $[\mathcal{L}, \mathcal{U}]$  for any arbitrary seed set  $C$  such that  $|C| = k$ . From Definition IV.1 we should have  $\gamma = \mathcal{U}/\mathcal{L}$  which directly implies (8) and (9) in Theorem IV.1 for  $\lambda = \Omega(1)$ . Hence we first focus on the proof of (15).

To begin with, recall the energy function  $H(\mathbf{x})$  in (7). For convenience, we abuse the terminology and define the energy function  $H(S)$  for a set  $S \subset V$  (not for a state  $\mathbf{x}$  as in (7)) as:

$$H(S) = \text{cut}(S, V \setminus S) - \sum_{i \in S} h|N(i)|$$

where  $\text{cut}(A, B)$  is the cardinality of the set  $\{(i, j) \in E \mid i \in A, j \in B\}$  for two disjoint subsets  $A, B \subset V$ . Note that the above definition coincides with the original definition (5) by setting  $x_i = 1$  if and only if  $i \in S$ . Using this energy function, one can express the function  $\Gamma^*(C)$  in (7) by:

$$\Gamma^*(C) = \max_{C \subset S_0 \subset V} \min_{\underline{S}: S_0 \rightarrow V} \max_{t < |\underline{S}|} [H(S_t) - H(S_0)], \quad (16)$$

where for  $A \subset V$ ,  $\underline{S}: A \rightarrow V$  is a monotone sequence of sets,  $A = S_0, S_1, \dots, S_{|\underline{S}|} = V$  such that  $S_{t-1} \subset S_t$  and  $S_t \setminus S_{t-1}$  is a vertex in  $V \setminus A$  for  $1 \leq t \leq |\underline{S}|$ .

To show the *concentration* of  $\Gamma^*$ , we first show the concentration of the energy function  $H$ , as stated in the next lemma whose proof is presented in Appendix D.

**Lemma A.1.** *Consider Erdős-Rényi graph  $G_{ER}(n, p)$  with  $\lambda = np = \Omega(1)$ . The following events occurs almost surely as  $n \rightarrow \infty$ :*

$$|H(S) - a(|S|)| \leq \eta(|S|),$$

where

$$a(s) = (1-h)s(n-s)p - hs(s-1)p,$$

$$\eta(s) = (1-h)\sqrt{2\lambda s(n-s)} + 2h\sqrt{\lambda s(s-1)}.$$

In Lemma A.1,  $H(S)$  is bounded by  $a(|S|) \pm \eta(|S|)$  which depends only cardinality of  $|S|$ . Thus, the paths, which are taken in min of  $\Gamma^*$ , have same bounds if they have same start  $S_0$ . Hence we have following:

$$\frac{\Gamma^*(C)}{\lambda n} = \frac{1}{\lambda n} \max_{C \subset S_0 \subset V} \min_{\underline{S}: S_0 \rightarrow V} \max_{t < |\underline{S}|} [H(S_t) - H(S_0)]$$

$$\leq \frac{1}{\lambda n} \max_{|C| \leq s_1 \leq s_2} a(s_2) + \eta(s_2) - a(s_1) + \eta(s_1) \quad (17)$$

$$= O\left(\frac{1}{n}\right) + \max_{\kappa \leq \sigma_1 \leq \sigma_2} \hat{a}(\sigma_2) + \hat{\eta}(\sigma_2) - \hat{a}(\sigma_1) + \hat{\eta}(\sigma_1), \quad (18)$$

where

$$\hat{a}(\sigma) = (1-h)\sigma(1-\sigma) - h\sigma^2,$$

$$\hat{\eta}(\sigma) = \frac{1-h}{\sqrt{\lambda}} + \frac{2h}{\sqrt{\lambda}}\sigma.$$

In (17), we have max over  $|C| \leq s_1 \leq s_2$  since  $C \subset S_0 \subset S_t$  for  $t < |\underline{S}|$ . Also, in (18), the  $O(\frac{1}{n})$  term is from  $O(\frac{1}{n} + \frac{1}{\lambda n})$  since we have  $\lambda = \Omega(1)$ .

We bound  $a, \eta$  by  $\hat{a}, \hat{\eta}$  for achieving an upper bound of a succinct close-form for  $\frac{\Gamma^*(C)}{\lambda n}$ . However, we note that one can directly consider (17) and obtain a tighter (but of a complicated form) upper bound for  $\frac{\Gamma^*(C)}{\lambda n}$ . Now it is not hard to check the maximum in (18) is

$$\left(\kappa - \left(\frac{1-h}{2} - \frac{h}{\sqrt{\lambda}}\right)\right)^2 + \frac{2(1-h^2)}{\sqrt{\lambda}}$$

at  $\sigma_1 = \kappa$  and  $\sigma_2 = \left(\frac{1-h}{2} + \frac{h}{\sqrt{\lambda}}\right)$  if  $\kappa \leq \left(\frac{1-h}{2} - \frac{h}{\sqrt{\lambda}}\right)$ . This implies that  $\frac{\Gamma^*(C)}{\lambda n} \leq \mathcal{U}$ . The proof of the lower bound  $\frac{\Gamma^*(C)}{\lambda n} \geq \mathcal{L}$  can be obtained similarly. This completes the proof of (15) and hence that of Theorem IV.1.

### B. Proof of Theorem IV.2

In this section, we present the proof of Theorem IV.2. Consider a planted partition graph  $G_{PP}(n, p, q, \omega)$  with  $p/q = \Omega(1)$ , and a seed set  $C'$  with budget  $k < \frac{1-h}{2}n$  satisfying the condition (10) in Theorem IV.2. Then, the quantification of  $\Gamma^*(G', C')$  is directly derived from the following lemma stating that where  $\Gamma^*(G', C')$  and  $\min_C \Gamma^*(G', C)$  is located, where the proof is provided in Appendix E.

**Lemma A.2.** *For every  $C'$  satisfying the conditions in Theorem IV.2, the following holds almost surely as  $n \rightarrow \infty$ ,*

$$\left| \frac{\Gamma^*(G', C')}{n^2} - \xi^2 p \right| \leq \frac{1}{2} n^{-0.4}$$

$$\left| \min_{C: |C| \leq k} \frac{\Gamma^*(G', C)}{n^2} - \xi^2 p \right| \leq \frac{1}{2} n^{-0.4}. \quad (19)$$

Thus, we focus on the proof of (11). To do so, it suffices to show that the following events occur almost surely as  $n \rightarrow \infty$ :

$$\frac{\Gamma^*(C') - \Gamma^*(C^*)}{\Gamma^*(C^*)} \leq \frac{2}{\left(\frac{p}{q+n^{-0.4}}\right)\xi^2 - 3} \quad (20)$$

where

$$\xi = \min_{\{\nu \in [0,1]^m: \nu_1 \leq \kappa\}} \max_{1 \leq l \leq m} \left( \frac{1-h}{2} \omega_l - \nu_l \right),$$

and  $C^*$  is an optimal seed set, i.e.,  $C^* \in \arg \min_{C: |C| \leq k = \kappa n} \Gamma^*(C)$ . This is because when  $p/q = \omega(1)$ , i.e.,  $q = o(1)$ , we have  $n^{-0.4}$  becomes arbitrarily small as  $n \rightarrow \infty$ , thus the result follows.

We first let  $G_l$  be the subgraph induced by each  $l$ -th cluster  $V_l$ , and  $E_l$  be the edges of  $G_l$ . We also let  $E_0 = E \setminus \cup_{l=1}^m E_l$ , which corresponds to the set of inter-cluster edges. Consider the ‘‘split graph’’  $G' = (V, E' = E \setminus E_0)$ , i.e.,  $G'$  is a graph removing the inter-cluster edges from  $G$ .

It is easy to have the following, which states that the difference of  $\Gamma^*$  between  $G$  and  $G'$  is bounded by the number of inter-cluster edges: For every  $C \subset V$ ,

$$|\Gamma^*(G, C) - \Gamma^*(G', C)| \leq 2|E_0|. \quad (21)$$

To check the above, for  $A, B$  such that  $A \subset B \subset V$ , we calculate  $H(B) - H(A)$  as below:

$$\begin{aligned} H(B) - H(A) &= (1-h) \cdot \text{cut}(B \setminus A, V \setminus B) - (1-3h) \cdot \text{cut}(A, B \setminus A) \\ &\quad + 2h \cdot \text{edge}(B \setminus A) \end{aligned} \quad (22)$$

where  $\text{edge}(S)$  is number of edges among nodes in  $S$ , i.e.,  $\text{edge}(S) = |\{(i, j) \in E | i, j \in S\}|$ . Note that in (22), three edge sets counted by cut and edge are disjoint. Thus, from removing an edge, change in value of (22) is at most  $\max(1-h, |1-3h|, 2h) \leq 2$  because of  $0 < h < 1$ . Also, we have  $S_0 \subset S_t$  in the expression of  $\Gamma^*$  in (16). Hence we have (21) since  $G'$  is the graph where  $E_0$  is removed from  $G$ .

Since the number of inter-cluster edges are stochastically dominated by a random variable with the binomial distribution  $B\left(\frac{n(n-1)}{2}, q\right)$ , we have:

$$\mathbb{P}\left[\frac{|E_0|}{n^2} \leq \frac{q}{2} + \frac{1}{4}n^{-0.4}\right] \rightarrow 1 \quad \text{as } n \rightarrow \infty, \quad (23)$$

where note that  $\mathbb{E}[|E_0|] = q \frac{n(n-1)}{2}$ .

Now, combining (21), (23), and Lemma A.2, leads to:

$$\left| \frac{\Gamma^*(C')}{n^2} - \xi^2 p \right| \leq (q + n^{-0.4}) \quad (24)$$

Furthermore, the following occurs almost surely as  $n \rightarrow \infty$ :

$$\begin{aligned} &\frac{\Gamma^*(G, C) - \Gamma^*(G, C^*)}{n^2} \\ &\stackrel{(a)}{\leq} \frac{\Gamma^*(G, C') - \Gamma^*(G', C^*)}{n^2} + \frac{2|E_0|}{n^2} \\ &\stackrel{(b)}{\leq} \frac{\Gamma^*(G', C')}{n^2} - \min_{C: |C| \leq k} \frac{\Gamma^*(G', C)}{n^2} + \frac{4|E_0|}{n^2} \\ &\stackrel{(c)}{\leq} n^{-0.4} + \frac{4|E_0|}{n^2} \stackrel{(d)}{\leq} 2(q + n^{-0.4}), \end{aligned} \quad (25)$$

where (a) is from (21), (b) is from (21) and the inequality:  $\min_C \Gamma^*(G', C) \leq \Gamma^*(G', C^*)$ , (c) is from Lemma A.2, and finally (d) is from (23). Then, noting the the bound  $\frac{\Gamma^*(C')}{\Gamma^*(C^*)} \leq \frac{\xi^2 p - 2(q + n^{-0.4})}{\xi^2 p - 3(q + n^{-0.4})}$ , (20) is a direct implication of (24) and (25). This completes the proof of (11).

### C. Proof of Theorem IV.3

This section provides the proof of Theorem IV.3. It is not hard to check the complexity of the seeding phase is  $O(n^2)$  for the following reason: In the seeding phase, we have total  $k = O(n)$  iterations. In each iteration, the number of clusters in the partition satisfying **P1** is  $O(n)(=m)$ . Further, the subphases of partition selection take  $O(m)$  and  $O(1)$  times, respectively, because using  $|V_l| = O(1)$ ,  $l = 1, \dots, m$ , we can compute the value  $\Gamma^*$  in each subgraph  $G_l$  in  $O(1)$  time (note that the nodes in  $V_0$  are already seeded).

We henceforth focus on the approximation quality of the output from the PaS algorithm and the quantity of  $\Gamma^*$  which the output has. To this end, we will use the following lemma whose proof is given in Appendix F.

**Lemma A.3.** *For every seed set  $C$  such that  $V_0 \subset C \subset V$ ,*

$$\Gamma^*(C) = \max_{l=1, \dots, m} \Gamma^*(G_l, C_l \cup V_0),$$

where  $C_l = C \cap V_l$ .

In addition, due to  $|V_l| = O(1)$ , it is not hard to check

$$\Gamma^*(G_l, C_l^{\text{PaS}} \cup V_0) = O(1)$$

which implies (14) with Lemma A.3. Hence, we will focus only on the quality of  $C^{\text{PaS}}$ .

To begin with, one can observe that the output  $C^{\text{PaS}}$  of the PaS algorithm minimizes  $\Gamma^*$  in each subgraph  $G_l$  for the budget allocation  $v_l^{\text{PaS}} = |C^{\text{PaS}} \cap V_l|$ , i.e.,

$$C_l^{\text{PaS}} \in \arg \min_{\{C_l \subset V_l: |C_l| \leq |C_l^{\text{PaS}}|\}} \Gamma^*(G_l, C_l \cup V_0), \quad (26)$$

where  $C_l^{\text{PaS}} = C^{\text{PaS}} \cap V_l$ . Recall that  $G_l$  is the subgraph induced by  $V_l \cup V_0$ . In addition,

From Lemma A.3 and (26), we have that

$$\Gamma^*(C^{\text{PaS}}) = \max_{1 \leq l \leq m} \min_{\{C_l \subset V_l: |C_l| \leq |C_l^{\text{PaS}}|\}} \Gamma^*(G_l, C_l \cup V_0). \quad (27)$$

Now we state the following key lemma, where its proof uses the above characterization of  $\Gamma^*(G, C^{\text{PaS}})$  and is presented in Appendix H.

**Lemma A.4.** *Given graph  $G = (V, E)$  and budget  $k$ , the output  $C^{\text{PaS}}$  of the PaS algorithm satisfies that*

$$C^{\text{PaS}} \in \arg \min_{C: |C| \leq k, V_0 \subset C} \Gamma^*(C).$$

From Lemma A.4, it follows that  $C^{\text{PaS}}$  is a  $(1, 1 - \frac{\xi}{\kappa})$ -approximation solution, since

$$\begin{aligned} \Gamma^*(C^{\text{PaS}}) &= \min_{C: |C| \leq k, V_0 \subset C} \Gamma^*(C) \\ &\leq \min_{C: |C| \leq k - |V_0|} \Gamma^*(C) \\ &\leq \min_{C: |C| \leq k(1 - \frac{\xi}{\kappa})} \Gamma^*(C), \end{aligned}$$

where we use  $|V_0| \leq \varepsilon n$ ,  $k = \kappa n$  and the monotone property of  $\Gamma^*$ , i.e., for all  $A, B$  such that  $A \subset B \subset V$ ,  $\Gamma^*(B) \leq \Gamma^*(A)$ . This completes the proof of Theorem IV.3.

#### D. Proof of Lemma A.1

Consider a subset  $S \subset V$ , where let  $s = |S|$ . For  $i \in S$ , we can split  $N(i)$  into two disjoint sets as  $N(i) = \left(N(i) \setminus S\right) \cup \left(N(i) \cap S\right)$ . Using this separation,  $H(S)$  in (16) can be written as:

$$H(S) = (1-h)\text{cut}(S, V \setminus S) - h \sum_{i \in S} |N(i) \cap S|. \quad (28)$$

In the ER graph, note that  $\text{cut}(S, V \setminus S)$  and  $\frac{1}{2} \sum_{i \in S} |N(i) \cap S|$  follows the binomial distributions  $B(s(n-s), p)$  and  $B(s(s-1)/2, p)$ , respectively. Then, from the Chernoff's bound, we have

$$\mathbb{P} \left[ \left| \text{cut}(S, V \setminus S) - ps(n-s) \right| \geq \sqrt{2\lambda s(n-s)} \right] \leq 2 \exp(-n), \quad (29)$$

$$\mathbb{P} \left[ \left| \frac{1}{2} \sum_{i \in S} |N(i) \cap S| - ps(s-1)/2 \right| \geq \sqrt{\lambda s(s-1)} \right] \leq 2 \exp(-n). \quad (30)$$

Thus, by applying the union bound to (29) and (30) and using (28), it follows that

$$\mathbb{P} \left[ |H(S) - a(s)| \geq \eta(s) \right] \leq 4 \exp(-n), \quad (31)$$

where  $a(s)$  and  $\eta(s)$  are defined in Lemma A.1. Finally, we complete the proof using the above inequality:

$$\mathbb{P} \left[ \bigcap_{S \subset V} [|H(S) - a(|S|)| \leq \eta(|S|)] \right] \geq 1 - 4 \exp(-n) \cdot 2^n \rightarrow 1 \quad \text{as } n \rightarrow \infty,$$

where we use the union bound and (31) for the first inequality.

#### E. Proof of Lemma A.2

We first note that each subgraph  $G_l$  is an ER graph  $G_{\text{ER}}(\omega_l n, p)$  where its  $\Gamma^*(G_l, \cdot)$  was already studied in Appendix A. Hence, from (18) with  $p = \Theta(1)$ , we have  $\hat{\eta}(\sigma) = O(n^{-0.5}) = o(n^{-0.4})$ .<sup>7</sup> Thus, for any  $C_l \subset V_l$  we have almost surely as  $n \rightarrow \infty$ :

$$\frac{\Gamma^*(G_{\text{ER}}(\omega_l n, p), C_l)}{n^2} = \begin{cases} \left(\frac{1-h}{2}\omega_l - \nu_l\right)^2 p + \frac{1}{2}n^{-0.4} & \text{if } \nu_l \leq \frac{1-h}{2} \\ \frac{1}{2}n^{-0.4} & \text{otherwise,} \end{cases}$$

where  $\nu_l = \frac{|C_l|}{n}$ . Also, we note that  $|V_l| = \omega_l n = \Omega(n)$ . Using the above, we have that almost surely as  $n \rightarrow \infty$ , for every  $C_l \subset V_l$ ,

$$\frac{\Gamma^*(G_l, C_l)}{n^2} = \left( \max \left( \frac{1-h}{2}\omega_l - \nu_l, 0 \right) \right)^2 p + \frac{1}{2}n^{-0.4}. \quad (32)$$

Since  $G'$  consists of disconnected subgraphs  $G_1, \dots, G_m$ , we provide the following which implies that the value  $\Gamma^*$  in the

<sup>7</sup>Here we have  $\lambda = np = \Theta(n)$ .

entire graph is decided by the maximum of the corresponding values in subgraphs: for every seed set  $C \subset V$ ,

$$\Gamma^*(G', C) = \max_{l=1, \dots, m} \Gamma^*(G_l, C_l). \quad (33)$$

The proof of (33) is almost identical to that of Lemma A.3, and we omit it for brevity.

Now observe that for every  $C \subset V$  with  $\frac{|C|}{n} \leq \kappa \leq \frac{1-h}{2}$ , there exists  $l$  such that  $\frac{|C_l|}{n} = \nu_l \leq \frac{1-h}{2}\omega_l$ . Thus, from (32) and (33), it follows that for every  $C \subset V$  such that  $|C| \leq k \leq \frac{1-h}{2}n$ ,

$$\frac{\Gamma^*(G', C)}{n^2} = \left( \max_{1 \leq l \leq m} \left( \frac{1-h}{2}\omega_l - \nu_l \right) \right)^2 p + \frac{1}{2}n^{-0.4}, \quad (34)$$

where  $\nu_l = \frac{|C_l|}{n}$ .

Therefore, it suffices to show the following:

$$\left| \max_{1 \leq l \leq m} \left( \frac{1-h}{2}\omega_l - \nu'_l \right) - \xi \right| \leq \frac{1}{2}n^{-0.4} \quad (35)$$

where  $\nu'_l = \frac{|C'_l|}{n}$ .

Since we consider  $C'$  satisfying (10),  $\max_{1 \leq l \leq m} \left( \frac{1-h}{2}\omega_l - \nu'_l \right)$  and  $\xi$  are the same except that the min is taken over  $\nu$  consisting of continuous  $\nu_l$  in  $\xi$  but we have the discreteness of  $\nu'_l = \frac{|C'_l \cap V_l|}{n}$ . Due to this discreteness,  $\xi$  and  $\max_{l=1, \dots, m} \frac{f(G_l, C'_l)}{n}$  have at most  $\frac{1}{n}$  difference which is less than  $n^{-0.4}$  as  $n \rightarrow \infty$ . This completes the proof.

#### F. Proof of Lemma A.3

We use proof by induction with respect to the number of clusters, i.e.  $m$ . The following claim states formally the base case  $m = 2$ , where its proof is presented in Appendix G.

**Proposition A.1.** *For given  $G = (V, E)$ , consider a partition  $\{V_l : l = 0, 1, 2\}$ , where there exists no edge between  $V_1$  and  $V_2$ ,*

$$\bigcup_{l \in \{0, 1, 2\}} V_l = V \quad \text{and} \quad V_l \cap V_{l'} = \emptyset, \quad \text{for all } l \neq l' \geq 0.$$

*Then, it follows that for any seed set  $C$  such that  $V_0 \subset C \subset V$ ,*

$$\Gamma^*(C) = \max_{l=1, 2} \Gamma^*(G_l, C_l \cup V_0),$$

*where  $G_l = (V_l \cup V_0, E_l)$  is the induced subgraph by  $V_l \cup V_0$  and  $C_l = C \cap V_l$ .*

We now consider two subgraphs  $G_1 = (V_1 \cup V_0, E_1)$  and  $G_{-1} = (V_{-1} \cup V_0, E_{-1})$  where

$$V_{-1} = \bigcup_{l=2}^m V_l \quad \text{and} \quad E_{-1} = \bigcup_{l=2}^m E_l.$$

Note that the separator  $V_0$  also partitions  $G$  into  $G_1$  and  $G_{-1}$  which are the subgraphs induced by  $V_1$  and  $V_{-1}$ , respectively. Then, from the construction of  $G_{-1}$  and Proposition A.1, for any seed set  $C$  such that  $V_0 \subset C \subset V$ , we have

$$\Gamma^*(C) = \max \{ \Gamma^*(G_1, C_1 \cup V_0), \Gamma^*(G_{-1}, C_{-1} \cup V_0) \},$$

where  $C_{-1} = C \cap V_{-1}$ .

Observe that  $V_0$  also partitions  $G_{-1} = (V_{-1} \cup V_0, E_{-1})$  into two subgraphs  $G_2 = (V_2 \cup V_0, E_2)$ ,  $G_{-2} = (V_{-2} \cup V_0, E_{-2})$

where  $V_{-2} = \cup_{l=3}^m V_l$  and  $E_{-2} = \cup_{l=3}^m E_l$ . Then, one can also apply Proposition A.1 to  $G_{-1}$  again: for any seed set  $C$  such that  $V_0 \subset C \subset V_{-1}$ ,

$$\Gamma^*(G_{-1}, C_{-1} \cup V_0) = \max \{ \Gamma^*(G_2, C_2 \cup V_0), \Gamma^*(G_{-2}, C_{-2} \cup V_0) \},$$

where  $C_{-2} = C \cap V_{-2}$ . Thus, we have, for any seed set  $C$  such that  $V_0 \subset C \subset V$ ,

$$\Gamma^*(C) = \max \left\{ \Gamma^*(G_2, C_2 \cup V_0), \max_{l=1,2} \Gamma^*(G_l, C_l \cup V_0) \right\}.$$

This provides the proof of Lemma A.3 for the case  $m = 3$ . One can repeat this procedure to complete the proof of Lemma A.3.

### G. Proof of Proposition A.1

For notational convenience, we will use the following definitions: for subset  $S_0 \subset V$  and monotone sequence of set  $\underline{S} \in S_0 \rightarrow V$ , we define

$$\Gamma(G, \underline{S}) = \max_{t \leq |\underline{S}|} [H(G, S_t) - H(G, S_0)] \quad (36)$$

$$\tilde{\Gamma}(G, S_0) = \min_{\underline{S}: S_0 \rightarrow V} \Gamma(G, \underline{S}).$$

Then, from the definition of  $\Gamma^*$ , we can write

$$\Gamma^*(G, C) = \max_{C \subset S_0 \subset V} \tilde{\Gamma}(G, S_0) = \max_{C \subset S_0 \subset V} \min_{\underline{S}: S_0 \rightarrow V} \Gamma(G, \underline{S}).$$

With the given partition, the following simple equality can be derived using (28) for any subset  $S$  such that  $V_0 \subset S$ ,

$$H(S) = H(G_1, S \cap W_1) + H(G_2, S \cap W_2). \quad (37)$$

where we let  $W_1 = V_1 \cup V_0$  and  $W_2 = V_2 \cup V_0$ .

Let  $C$  denote a seed set such that  $V_0 \subset C \subset V$ . Also, let  $C_1 = C \cap V_1$  and  $C_2 = C \cap V_2$ . To complete the proof of this proposition, we will show that the followings hold:

$$\max_{l=1,2} \Gamma^*(G_l, C_l \cup V_0) \leq \Gamma^*(C), \quad (38)$$

$$\Gamma^*(C) \leq \max_{l=1,2} \Gamma^*(G_l, C_l \cup V_0). \quad (39)$$

**Proof of (38).** For a subset  $X \subset V$  such that  $C \subset X$ , define

$$\begin{aligned} \mathcal{P}_l(X) &= \{ \underline{S}' : X \cap W_l \rightarrow W_l \}, \\ \mathcal{Q}_l(X) &= \{ \underline{S} : X \cup W_l \rightarrow V \}. \end{aligned}$$

Then, we have

$$\begin{aligned} &\tilde{\Gamma}(X \cup W_2) \\ &= \min_{\underline{S} \in \mathcal{Q}_2(X)} \max_{t \leq |\underline{S}|} [(H(S_t) - H(S_0))] \\ &= \min_{\underline{S} \in \mathcal{Q}_2(X)} \max_{t \leq |\underline{S}|} [(H(G_1, S_t \cap W_1) + H(G_2, S_t \cap W_2)) \\ &\quad - [H(G_1, S_0 \cap W_1) + H(G_2, S_0 \cap W_2)]] \quad (\because (37)) \\ &\stackrel{(a)}{=} \min_{\underline{S} \in \mathcal{Q}_2(X)} \max_{t \leq |\underline{S}|} [H(G_1, S_t \cap W_1) - H(G_1, S_0 \cap W_1)] \\ &\stackrel{(b)}{=} \min_{S' \in \mathcal{P}_1(X)} \max_{t \leq |\underline{S}'|} [H(G_1, S'_t) - H(G_1, S'_0)] \\ &= \tilde{\Gamma}(G_1, X \cap W_1) \end{aligned} \quad (40)$$

In the above, (a) holds since  $H(G_2, S_t \cap W_2) = H(G_2, W_2)$  for all  $t$ , which comes from the fact that  $V_2 \cup V_0 \subset S_t$ . (b) holds since there is a one-to-one correspondence between  $\mathcal{P}_1(X)$

and  $\mathcal{Q}_2(X)$ ; i.e.,  $\underline{S}'$  can be induced from  $\underline{S}$  by  $\underline{S}' = (S_0 - V_2, \dots, S_t - V_2, \dots, V - V_2 (= W_1))$  and vice versa. Similarly, one can show that

$$\tilde{\Gamma}(X \cup W_1) = \tilde{\Gamma}(G_2, X \cap W_2). \quad (41)$$

Since  $C \subset X \subset V$ , it follows that

$$\begin{aligned} \Gamma^*(C) &= \max_{C \subset S_0 \subset V} \tilde{\Gamma}(S_0) \\ &\geq \max_{l=1,2} \tilde{\Gamma}(X \cup W_l) = \max_{l=1,2} \tilde{\Gamma}(G_l, X \cap W_l) \end{aligned}$$

where the last equality holds from (40) and (41).

Now by taking the maximum of  $\max_{l=1,2} \tilde{\Gamma}(G_l, X \cap W_l)$  over all  $X$  such that  $C \subset X \subset V$ , we conclude that

$$\begin{aligned} \Gamma^*(C) &\geq \max_{C \subset X \subset V} \max_{l=1,2} \tilde{\Gamma}(G_l, X \cap W_l) \\ &= \max_{l=1,2} \max_{C \subset X \subset V} \tilde{\Gamma}(G_l, X \cap W_l) \\ &= \max_{l=1,2} \Gamma^*(G_l, C_l \cup V_0). \end{aligned}$$

This completes the proof of (38).

**Proof of (39).** Let  $S_0^*$  and  $\underline{S}^*$  be an optimal subset of  $V$  and an optimal monotone sequence of sets for  $G$ , i.e.,  $C \subset S_0^* \subset V$ ,  $\underline{S}^* : S_0^* \rightarrow V$ , and

$$\Gamma^*(C) = \tilde{\Gamma}(S_0^*) = \Gamma(\underline{S}^*).$$

In addition, let  $\underline{S}^1 : S_0^* \cap W_1 \rightarrow W_1$  and  $\underline{S}^2 : S_0^* \cap W_2 \rightarrow W_2$  be an optimal monotone sequences of sets for  $G_1, G_2$ , respectively. Then we have

$$\begin{aligned} \Gamma^*(G_1, C_1) &= \tilde{\Gamma}(G_1, S_0^* \cap W_1) = \Gamma(G_1, \underline{S}^1), \\ \Gamma^*(G_2, C_2) &= \tilde{\Gamma}(G_2, S_0^* \cap W_2) = \Gamma(G_2, \underline{S}^2). \end{aligned}$$

Now, construct  $\underline{S}^1 \cup S_0^* : S_0^* \rightarrow S_0^* \cup V_1$  and  $\underline{S}^1 \cup S_0^* : S_0^* \cup V_1 \rightarrow V$  such that

$$\begin{aligned} \underline{S}^1 \cup S_0^* &= (S_0^1 \cup S_0^*, \dots, S_t^1 \cup S_0^*, \dots, S_{|\underline{S}^1|}^1 \cup S_0^*), \\ \underline{S}^2 \cup V_1 &= (S_0^2 \cup V_1, \dots, S_t^2 \cup V_1, \dots, S_{|\underline{S}^2|}^2 \cup V_1). \end{aligned}$$

Since the end of  $\underline{S}^1 \cup S_0^*$  and the start of  $\underline{S}^2 \cup V_1$  are the same (note that  $S_0^1 \cup S_0^* = S_0^*$ ,  $S_{|\underline{S}^1|}^1 \cup S_0^* = S_0^* \cup V_1 = S_0^2 \cup V_1$  and  $S_{|\underline{S}^2|}^2 \cup V_1 = W_2 \cup V_1 = V$ ), and  $V_0 \subset S_0^*$ , we can construct a new monotone sequence of sets  $\underline{T} : S_0^* \rightarrow V$  by concatenating  $\underline{S}^1 \cup S_0^*$  and  $\underline{S}^2 \cup V_1$ :

$$\begin{aligned} \underline{T} &= (S_0^*, S_1^1 \cup S_0^*, S_2^1 \cup S_0^*, \dots, S_{|\underline{S}^1|}^1 \cup S_0^*, \\ &\quad S_1^2 \cup V_1, S_2^2 \cup V_1, \dots, S_{|\underline{S}^2|}^2 \cup V_1, V). \end{aligned}$$

Thus, we have

$$\begin{aligned} \Gamma(\underline{T}) &= \max \left( \max_{t \leq |\underline{S}^1|} H(S_t^1 \cup S_0^*), \max_{t \leq |\underline{S}^2|} H(S_t^2 \cup V_1) \right) \\ &\quad - H(S_0^*). \end{aligned}$$

Using the construction of  $\underline{T}$  with (36) and (37), it is not hard to check that

$$\max_{t \leq |\underline{S}^1|} H(S_t^1 \cup S_0^*) = \Gamma(G_1, \underline{S}^1) + H(S_0^*) \quad (42)$$

$$\max_{t \leq |\underline{S}^2|} H(S_t^2 \cup V_1) =$$

$$\Gamma(G_2, \underline{S}^2) + H(G_1, W_1) + H(G_2, S_0^* \cap W_2). \quad (43)$$

Furthermore, using (42), (43) and (37), we have

$$\max_{t \leq |\underline{S}^1|} H(S_t^1 \cup S_0^*) - H(S_0^*) = \Gamma(G_1, \underline{S}_1^1) \quad (44)$$

$$\begin{aligned} & \max_{t \leq |\underline{S}^2|} H(S_t^2 \cup V_1) - H(S_0^*) \\ &= \Gamma(G_2, \underline{S}^2) + H(G_1, W_1) - H(G_1, S_0^* \cap W_1). \end{aligned} \quad (45)$$

Recall that the state that all players choose +1 has the minimum of  $H(\cdot)$ . Hence on the subgraph  $G_1$ ,  $H(G_1, \cdot)$  has the minimum at  $W_1$ , i.e.,  $H(G_1, W_1) = \min_{S \subset W_1} H(G_1, S)$ . Thus, we have

$$H(G_1, W_1) - H(G_1, S_0^* \cap W_1) < 0.$$

Combining (44) and (45) leads us to:

$$\begin{aligned} \Gamma(\underline{T}) &\leq \max(\Gamma(G_1, \underline{S}^1), \Gamma(G_2, \underline{S}^2)) \\ &= \max(\tilde{\Gamma}(G_1, S_0^* \cap W_1), \tilde{\Gamma}(G_2, S_0^* \cap W_2)) \\ &\leq \max(\Gamma^*(G_1, C_1 \cup V_0), \Gamma^*(G_2, C_2 \cup V_0)), \end{aligned}$$

where the last inequality is due to  $C \subset S_0^*$ . Since  $\underline{T}$  and  $\underline{S}^*$  are monotone sequences of sets from  $S_0^* \rightarrow V$ , we have

$$\Gamma(\underline{S}^*) \leq \Gamma(\underline{T}) \leq \max_{l=1,2} \Gamma^*(G_l, C_l \cup V_0),$$

where the first equality holds by the definition of  $\underline{S}^*$ . This completes the proof of (39) and hence completes the proof of Proposition A.1.

#### H. Proof of Lemma A.4

We use proof by contradiction. To this end, suppose that there exists  $C^* \neq C^{\text{PaS}}$  such that,  $|C^*| = k$ ,  $V_0 \subset C^* \subset V$  and

$$\Gamma^*(C^{\text{PaS}}) > \Gamma^*(C^*). \quad (46)$$

Let  $C_l^{\text{PaS}} = C^{\text{PaS}} \cap V_l$  and  $C_l^* = C^* \cap V_l$ . Then, from  $C^* \neq C^{\text{PaS}}$ , there must exist  $l'$  such that

$$|C_{l'}^{\text{PaS}}| > |C_{l'}^*|. \quad (47)$$

The above inequality implies that the PaS algorithm selects the cluster  $l'$  (in step 2-3) more than  $|C_{l'}^*|$  times, where we say that it does for the  $|C_{l'}^*| + 1$  time at the  $t$ -th iteration of the seeding phase. This means that at the end of the  $(t-1)$ -th iteration, the set of seeds in the cluster  $l'$  has cardinality  $|C_{l'}^*|$  and the largest  $\Gamma^*$  among clusters, i.e.,  $|C_{l'}^{\text{PaS}}(t-1)| = |C_{l'}^*|$ , and

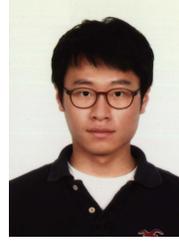
$$\begin{aligned} \Gamma^*(C^{\text{PaS}}(t-1)) &= \Gamma^*(G_{l'}, C_{l'}^{\text{PaS}}(t-1)) \\ &= \min_{C_{l'} \subset V_{l'}: |C_{l'}| \leq |C_{l'}^{\text{PaS}}(t-1)|} \Gamma^*(G_{l'}, C_{l'} \cup V_0) \\ &= \min_{C_{l'} \subset V_{l'}: |C_{l'}| \leq |C_{l'}^*|} \Gamma^*(G_{l'}, C_{l'} \cup V_0) \end{aligned} \quad (48)$$

where  $C^{\text{PaS}}(t-1)$  denotes the intermediate seed set at the end of the  $(t-1)$ -th iteration of the seeding phase. Therefore, it follows that

$$\begin{aligned} \Gamma^*(C^{\text{PaS}}) &\stackrel{(a)}{\leq} \Gamma^*(C^{\text{PaS}}(t'-1)) \\ &\stackrel{(b)}{=} \min_{C_{l'} \subset V_{l'}: |C_{l'}| \leq |C_{l'}^*|} \Gamma^*(G_{l'}, C_{l'} \cup V_0) \\ &\leq \Gamma^*(G_{l'}, C_{l'}^* \cup V_0) \end{aligned}$$

$$\begin{aligned} &\leq \max_{1 \leq l \leq m} \Gamma^*(G_l, C_l^* \cup V_0) \\ &\stackrel{(c)}{=} \Gamma^*(C^*), \end{aligned}$$

where (a) is from the fact that the PaS algorithm keeps reducing  $\Gamma^*$  at every iteration, (b) is due to (48), and (c) uses Lemma A.3. This conflicts to (46), and completes the proof of Lemma A.4.



**Jungseul Ok** Jungseul Ok received the B.S. in electrical engineering from the KAIST, Daejeon, Republic of Korea, in 2011. He is currently working toward the Ph.D degree in school of electrical engineering at KAIST. His research interests social networks, data mining, machine learning, and future wireless communication systems.



**Youngmi Jin** Youngmi Jin (S'00M'05) received the B.S. and M.S. degrees in mathematics from the Korea Advanced Institute of Science and Technology (KAIST), Daejeon, Korea, in 1991 and 1993, respectively, and the M.S. and Ph.D. degrees in electrical engineering from The Pennsylvania State University, University Park, PA, USA, in 2000 and 2005, respectively. She worked with University of Pennsylvania, Philadelphia, PA, USA, and KT, Sogang University, and KAIST in Korea. She is currently with KDDI R&D Labs, Saitama, Japan, as a Research Engineer.

Her research interests include cloud computing, social networks, Internet economics, and wireless networks.

**Jinwoo Shin** Jinwoo Shin is currently an assistant professor at the Department of Electrical Engineering at KAIST, South Korea. He obtained his B.S. degrees in Computer Science and Mathematics from Seoul National University in 2001 and his Ph.D. degree in Mathematics from Massachusetts Institute of Technology in 2010. After spending two years at Algorithms and Randomness Center, Georgia Institute of Technology, one year (2012-2013) at Business Analytics and Mathematical Sciences Department and IBM T. J. Watson Research, he



joined the KAIST department in Fall 2013. He received the best student paper award at ACM SIGMETRICS 2009, the best MIT CS doctoral thesis (George M. Sprowls) award 2010, the best paper award at ACM MOBIHOC 2013, the best publication award from INFORMS applied probability society 2013 and Bloomberg scientific research award 2015.



**Yung Yi** Yung Yi received his B.S. and the M.S. in the School of Computer Science and Engineering from Seoul National University, South Korea in 1997 and 1999, respectively, and his Ph.D. in the Department of Electrical and Computer Engineering at the University of Texas at Austin in 2006. From 2006 to 2008, he was a post-doctoral research associate in the Department of Electrical Engineering at Princeton University. Now, he is an associate professor at the Department of Electrical Engineering at KAIST, South Korea. His current

research interests include the design and analysis of computer networking and wireless communication systems, especially congestion control, scheduling, and interference management, with applications in wireless ad hoc networks, broadband access networks, economic aspects of communication networks, and green networking systems. He received the best paper awards at IEEE SECON 2013, ACM MOBIHOC 2013, and IEEE William R. Bennett Award 2016.