# Using Social Networks to Raise HIV Awareness Among Homeless Youth

A. Yadav, H. Chan, A.X. Jiang, H. Xu, E. Rice, R. Petering, M. Tambe

## Abstract

Many homeless shelters conduct interventions to raise awareness about HIV (human immunodeficiency virus) among homeless youth. Due to human and financial resource shortages, these shelters need to choose intervention attendees strategically, in order to maximize awareness through the homeless youth social network. In this work, we propose HEALER (hierarchical ensembling based agent which plans for effective reduction in HIV spread), an agent that recommends sequential intervention plans for use by homeless shelters. HEALER's sequential plans (built using knowledge of homeless youth social networks) select intervention participants strategically to maximize influence spread, by solving POMDPs (partially observable Markov decision process) on social networks using heuristic ensemble methods. This paper explores the motivations behind HEALER's design, and analyzes HEALER's performance in simulations on real-world networks. First, we provide a theoretical analysis of the DIME (dynamic influence maximization under uncertainty) problem, the main computational problem that HEALER solves. HEALER relies on heuristic methods for solving the DIME problem due to its computational hardness. Second, we explain why heuristics used inside HEALER work well on real-world networks. Third, we present results comparing HEALER to baseline algorithms augmented by HEALER's heuristics. HEALER is currently being tested in real-world pilot studies with homeless youth in Los Angeles.

## Introduction

Homelessness has reached a crisis level, with over 565,000 homeless people in the US on any given night. Homeless youth (i.e., people below the age of 25) account for almost 34% of the total homeless population [1]. These homeless youth face significant difficulties, having to struggle for basic amenities such as healthcare, nutritious food, and primary education.

HIV has an extremely high incidence among homeless youth, as they are more likely to engage in high HIV-risk behaviors (e.g., unprotected sexual activity, injection drug use) than other sub-populations. In fact, previous studies show that homeless youth face a 10 times greater risk of HIV infection than stably housed populations [1].

To help prevent HIV infection among homeless youth, many homeless shelters implement social network based peer-leader intervention programs, where a select number of youth, called peer leaders, are taught strategies for reducing risk of contracting HIV. These intervention programs consist of day-long educational sessions in which these peer-leaders are provided with information about HIV prevention measures [2]. These "leaders" are then encouraged to share these messages with their peers in their social circles, in order to "lead" all their peers towards safer behaviors and practices.

These peer-leader-based intervention programs are motivated by the shelters limited financial resources, which prevents them from directly assisting on the entire homeless youth population. Therefore, they try to maximize the spread of awareness among the homeless youth population (via word-of-mouth influence) using the limited resources at their disposal. This leads to the well-known question from the field of influence maximization of how to select "influential" nodes (i.e., homeless youth) to maximize spread of awareness within a given social network? In the context of homeless shelters, this problem is further complicated by two factors. First, the social network structure is imperfectly known, which makes identifying "influential" nodes challenging [3]. Although some connections (friendships) are known, other connections may be uncertain. This is because homeless youth are a hard-to-reach population, and their social networks are harder to characterize than networks of stably housed youth [2]. Second, managing homeless youth (some of whom have emotional and behavioral problems) during an intervention, with the homeless shelter's limited personnel is not easy. As a result, the shelter officials can only manage small groups composed of three or four youths at one time. Therefore, the shelter officials prefer a series of small sized intervention camps organized sequentially (i.e., one after the other) to maximize the impact of their intervention [4]. In such camps, youth may reveal some additional information about the network; which can be used to inform future interventions.

The shelters need a plan to choose the participants (i.e., peer leaders) of their sequentially organized interventions. This plan must address four key points: (i) it must efficiently deal with uncertainties in the network structure, i.e., uncertainty about existence or absence of some friendships in the network; (ii) it needs to take into account new information uncovered during the interventions, which reduces the uncertainty in our understanding of the network; (iii) the plan needs to be deviation tolerant, as sometimes homeless youth may choose not to be a peer leader, thereby forcing the shelter to modify its plan; (iv) our approach should address the challenge of gathering information about social networks of homeless youth, which usually costs thousands of dollars and many months of time [4].

This paper presents three key contributions in addressing the sequential planning needs of homeless shelters. First, we model the shelters' sequential planning needs by introducing the dynamic influence maximization under uncertainty (or DIME) problem. The sequential selection of intervention participants under network uncertainty in DIME sets it apart from any other previous work on influence maximization, which mostly focuses on single shot decision problems (i.e., a set of nodes are selected just once, instead of selecting sets of nodes repeatedly) [5, 6, 7, 8]. We analyze several novel theoretical aspects of the DIME problem, which illustrates its computational hardness.

Second, we propose a new software agent, HEALER (hierarchical ensembling based agent which plans for effective reduction in HIV spread), to provide an end-to-end solution to the DIME problem. First, HEALER casts the DIME problem as a partially observable Markov decision process (POMDP) and solves it using HEAL (hierarchical ensembling algorithm for planning), a novel POMDP planner that quickly generates high-quality recommendations (of intervention

participants) for homeless shelter officials. In this paper, we discuss the design of HEALER and explain its method of gathering information about the homeless youth social network (at low cost) by interacting with youth via a network construction application. We also give a high-level overview of the HEAL algorithm and refer the reader to Yadav et. al. [9] for a more complete understanding.

Simulations presented in Yadav et. al. [9] show that even on small networks, HEAL achieves a 100-fold speed up and 70% improvement in solution quality over PSINET (POMDP based social interventions in networks for enhanced HIV testing) [10] (a baseline algorithm which uses POMDPs); and on larger networks *where PSINET is unable to run at all*, HEAL continues to provide high quality solutions quickly. This is in spite of several sub-optimal heuristics currently being used by HEALER. In order to better understand why HEALER's heuristics work so well on real-world networks, we analyzed the real-world networks used in HEALER's simulations to determine the reasons behind its outperforming its competitors.

Finally, to further explore and highlight HEALER's performance, we compare HEALER with baseline algorithms, including one augmented with the same heuristics. We then provide comparison results for HEALER and these augmented baseline algorithms.

HEALER has been tested in a real-world pilot study, in collaboration with a homeless shelter (called Safe Place for Youth), which provides food and lodging to homeless youth aged 12-25. They provide these facilities for ~55-60 homeless youth every day. They also operate an on-site medical clinic where free HIV and Hepatitis-C testing is provided. The recently completed pilot study enrolled 60 homeless youth, and then conducted three interventions on this population based on HEALER's recommended peer-leaders. To the best of our knowledge, this pilot study represents the first real-world evaluation of such sequential influence maximization algorithms, and it showed HEALER's effectiveness at spreading information in a social network effectively. Results from the pilot studies can be found in Yadav et. al. [11].

## Related work
There are three distinct areas of work related to the homeless shelter problem that we introduced above. The primary problem in computational influence maximization is to find optimal 'seed sets' of nodes in social networks, which can maximize the spread of information or influence in the social network (according to some a priori known influence model). While there are many algorithms for finding `seed sets' of nodes to maximize influence spread in networks [5, 6, 7, 8], most of these algorithms assume *no uncertainty in the network structure* and select a single seed set. In contrast, HEALER selects several seed sets sequentially in our work to select intervention participants for each successive training program, taking into account updates to the network structure revealed in past interventions. The DIME problem also incorporates uncertainty about the network structure and influence status of network nodes (i.e., whether a node is influenced or not). Finally, unlike [5, 6, 7, 8], HEALER uses a different diffusion model as we explain later in

this paper. Golovin et. al. [12] introduced adaptive submodularity and discussed adaptive sequential selection (similar to our problem), and they proved that a Greedy algorithm has a (1-1/ *e*) approximation guarantee. However, unlike the DIME problem, they assume no uncertainty in the network structure. We show that while the DIME problem can be cast into the adaptive stochastic optimization framework of [12], its influence function is not adaptive submodular (see section of paper titled "DIME Problem Statement") and because of this, their Greedy algorithm loses its approximation guarantees. Finally, Lei et. al. [13] use multi-armed bandit algorithms to pick influential nodes in social networks when influence probabilities are not known, but their approach requires lots of iterations to converge, thereby making it unsuitable for a real-world domain like ours.

Next, we discuss literature from *social work*. The general approach to these interventions is to use peer change agents (PCA) (i.e., peers who bring about change in attitudes) to engage homeless youth in interventions, but most studies do not use network characteristics to choose these PCAs [14]. A notable exception is Valente et. al. [15], who proposed selecting intervention participants with highest *degree centrality* (the most ties to other homeless youth). However, previous studies [10, 16] show that *degree centrality* performs poorly, as it does not account for potential overlaps in influence of two high degree centrality nodes.

The final field of related work is planning for reward and cost optimization. We only focus on the literature on Monte-Carlo (MC) sampling based online POMDP solvers since this approach allows significant scale-up [17]. The POMCP (Partially Observable Monte-Carlo Planning) solver [18]  uses Monte-Carlo UCT (upper confidence bound) tree search in online POMDP planning. Also, Somani et. al. [19] present the DESPOT (determinized sparse partially observable tree) algorithm, that improves the worst case performance of POMCP. Our initial experiments with POMCP and DESPOT showed that they run out of memory on even our small sized networks. A recent paper introduced PSINET-W [10], which is a MC sampling based online POMDP planner. We have discussed PSINET's shortcomings above, and how HEALER remedies them with the use of its heuristics. In particular, HEALER scales up whereas PSINET fails to do so. *HEALER's algorithmic approach also offers significant novelties in comparison with PSINET*.

## HEALER's Design

HEALER has a modular design [9], and consists of two major components. First, it has a network construction application for gathering information about social networks. Second, it has an algorithm called HEAL, which solves the DIME problem (introduced later) using heuristics. We first explain HEALER's components individually, and then explain how they are used inside HEALER's design.

### *Network Construction Application*

HEALER gathers information about social ties among homeless youth by interacting with these youth via its network construction application. Once a fixed number of homeless youth register

in its network application (which is hosted as a website to ensure ease of access for the youth), HEALER parses contact lists (on Facebook) of all the registered homeless youth and generates the social network that connects these youth. We choose Facebook for gathering information because previous studies [20] show that a large proportion (~80%) of homeless youth are regularly active on Facebook. Specifically, HEALER adds a link between two homeless youth, if and only if both youth are (i) friends on Facebook; and (ii) are registered in its application. Unfortunately, there is *uncertainty* in the generated network as friendship links between people who are only friends in real-life (and have not added each other as friends on Facebook) are not captured by HEALER's network construction application.

Previously, collecting accurate social network data on homeless youth was a technical and financial burden beyond the capacity of most agencies working with these youth [20]. Homeless shelters conducted tedious face-to-face interviews with homeless youth to infer ties between these youth, a process that costs thousands of dollars and many months of time. HEALER's network construction application enables homeless shelters to quickly generate a first approximation of the homeless youth social network at low cost. The HEAL algorithm (the second component in HEALER) subsequently corrects and improves the social network structure iteratively (as explained later), which is one of the major strengths of this approach. This network construction application has been tested multiple times by our collaborating homeless shelter with positive feedback.

### DIME Solver
The DIME Solver then takes the approximate social network (generated by HEALER's network construction application) as input and solves the DIME problem (formally defined later in the paper) using HEAL, the core algorithm running inside HEALER. The HEAL algorithm is an online POMDP solver, i.e., it interleaves planning and execution for each time step (explained later in the paper). The solution of the DIME problem generated by HEAL is provided as a series of recommendations (of intervention participants) to homeless shelter officials. Each recommendation would urge the officials to invite a particular set of youth for their intervention camp. For example, in Figure 1, HEALER would recommend inviting nodes *D* and *A* for the intervention.

### HEALER Design
HEALER's design begins with the network construction application constructing an *uncertain* network (as explained above). HEALER has a *sense-reason-act* cycle; where it repeats the following process for *T* interventions. It *reasons* about different long-term plans to solve the DIME problem, it *acts* by providing DIME's solution as a recommendation (of intervention participants) to homeless shelter officials. The officials may choose to not use HEALER's recommendation in selecting their intervention's participants. After finalizing the selection of participants, the shelter officials contact the chosen participants (via phone/email) and conduct the intervention with them. Upon the intervention's completion, HEALER *senses* feedback about the conducted intervention from the officials. This feedback includes new observations about the network, e.g., uncertainties in some links may be resolved as intervention participants are

interviewed by the shelter officials (explained more later). HEALER uses this feedback to update and improve its future recommendations.

## DIME Problem Statement

HEALER represents social networks as directed graphs (consisting of *nodes* and *directed edges*) where each *node* represents a person in the social network and a *directed edge* between two nodes *A* and *B* (say) represents that node *A considers* node *B* as their friend. *HEALER assumes directed-ness of edges as sometimes homeless shelters assess that the influence in a friendship is very much uni-directional; and to account for uni-directional follower links.* Otherwise friendships are encoded as two uni-directional links. In the following, we provide some background information that helps us define a precise problem statement for DIME. After that, we will show some hardness results about this problem statement.

*Uncertain Network*
The uncertain network is a directed graph $G = (V, E)$ with $|V| = N$ nodes and $|E| = M$ edges. The edges *E* in an uncertain network are of two distinct types: (i) the set of certain edges , that consists of friendships that we are certain about; and (ii) the set of uncertain edges ,which consists of friendships which we are uncertain about. Recall that uncertainties about friendships exist because HEALER's network construction application misses out on some links between people who are friends in real life, but not on Facebook.

To model the uncertainty about missing edges, every uncertain edge has an existence probability *u(e)* associated with it, which represents the likelihood of "existence" of that uncertain edge in the real-world. For example, if there is an uncertain edge *(A,B)* (i.e., we are unsure whether node *B* is node *A*'s friend), then *u(A,B)* = 0.75 implies that *B* is *A*'s friend with a 0.75 chance. This existence probability allows us to measure the potential value of influencing a given node. For example, if node A is connected to many uncertain edges with low *u(e)* values, then it is unlikely that node A is highly influential (as most of his supposed friendships may not exist in reality).

In addition, every edge in the network (both certain and uncertain) has a propagation probability *p(e)* associated with it. A propagation probability of 0.5 on directed edge *(A,B)* denotes that if node *A* is influenced (i.e., has information about HIV prevention), it influences node *B* (i.e., gives information to node *B*) with a 0.5 probability in each subsequent time step (our full influence model is defined below). This graph *G* with all relevant *p(e)* and *u(e)* values represents an uncertain network and serves as an input to the DIME problem. **Figure 1** shows an example of an uncertain network, where the dotted edges represent uncertain edges. We now explain how HEALER generates an uncertain social network.

First, HEALER uses its network construction application to generate a network with no uncertain edges. Next, we use well known link prediction techniques such as KronEM [21] to infer existence probabilities *u(e)* for additional friendships that might have been missed by the

network construction application. This process gives us an *uncertain network* which is then used by HEALER to generate recommendations, as we explain next.

Given the *uncertain network* as input, HEALER runs for *T rounds* (corresponding to the number of interventions organized by the homeless shelter). In each round, HEALER chooses *K nodes* (youth) as intervention participants. These participants are assumed to be influenced after the intervention (i.e., our intervention deterministically influences the participants). Upon influencing the chosen nodes, HEALER `observes' the true state of the *uncertain edges* (friendships) out-going from the selected nodes. This translates to asking intervention participants about their 1-hop social circles, which is within the capabilities of the homeless shelter [2].

After each round, influence spreads in the network according to our influence model (explained below) for *L time steps*, before we begin the next round. This *L* is the time duration in between two successive intervention camps. *In between rounds, HEALER does not observe the nodes that get influenced during L time steps*. Thus, while HEALER knows the influence model, it does not observe the random samples from the influence model which led to some nodes getting influenced. HEALER only knows that explicitly chosen nodes (our intervention participants in all past rounds) are influenced. Informally then, given an uncertain network and integers *T, K,* and *L* (as defined above), HEALER finds an online policy for choosing *exactly K* nodes for *T* successive rounds (interventions) that maximizes influence spread in the network at the end of *T* rounds.

***Influence Model***
Unlike most previous work in influence maximization [5, 6, 7, 8], HEALER uses a variation of the independent cascade model [22]. In the standard independent cascade model, all nodes that get influenced at time *t* get a *single* chance to influence their un-influenced neighbors at time *t+1*. If they fail to spread influence in this *single* chance, they don't spread influence to their neighbors in future rounds. On the other hand, HEALER's model assumes that nodes get *multiple* chances to influence their un-influenced neighbors. If they succeed in influencing a neighbor at a given time step *t'*, they stop influencing that neighbor for all future time steps. Otherwise, if they fail in step *t'*, they try to influence again with the same propagation probability in the next time step. This variant of independent cascade has been shown to empirically provide a better approximation to real influence spread than the standard independent cascade model [22, 23]. Further, we assume that nodes that get influenced at a certain time step remain influenced for all future time steps.

We now provide notation for defining HEALER's policy formally. Let  denote the set of *K* sized subsets of *V*, which represents the set of possible choices that HEALER can make at every time step . Let  denote HEALER's choice in the time step. Upon making choice , HEALER `observes' uncertain edges adjacent to nodes in , which updates its understanding of the network. Let  denote the uncertain network resulting from  with *observed* (additional edge) information from .

7

Formally, we define a history of length i as a tuple of past choices and observations . Denote by the set of all possible histories of length less than or equal to *i*. Finally, we define an *i*-step policy as a function that takes in histories of length less than or equal to *i* and outputs a *K* node choice for the current time step. We now provide an explicit problem statement for DIME.

***Problem Statement***

Given as input an uncertain network and integers *T, K* and *L* (as defined above). Denote by the *expected total number of influenced nodes at the end of round T*, given the *T*-length history of previous observations and actions , along with , the action chosen at time T. Let denote the expectation over the random variables and influence of , where are chosen according to , and are drawn according to the distribution over uncertain edges of that are revealed by . The objective of DIME is to find an optimal T-step policy .

Next, we show hardness results about the DIME problem. First, we analyze the value of having complete information in DIME. Then, we characterize the computational hardness of DIME.

***The Value of Information***

We characterize the impact of insufficient information (about the uncertain edges) on the achieved solution value. We show that no algorithm for DIME is able to provide a sufficiently good approximation to the *full-information solution value* (i.e., the best solution achieved w.r.t. the underlying ground-truth network), even with infinite computational power.

**Theorem 1** Given an uncertain network with *n* nodes, for any , there is no algorithm for the DIME problem that can guarantee a approximation to, the *full-information solution value.*
**Proof** We prove this statement by providing a counter-example in the form of a specific (ground truth) network for which there can exist no algorithm that can guarantee a approximation to . Consider an input to the DIME problem, an *uncertain network* with n nodes with uncertain edges between the *n* nodes, i.e., it is a completely connected uncertain network consisting of *only* uncertain edges (an example with *n=3* is shown in Figure 1). Let *p(e)=*1 and *u(e)=*0.5 on all edges in the *uncertain network*, i.e., all edges have the same propagation and existence probability. Let *K=*1*, L=*1 and *T=*1, i.e., we just select a single node in one shot (in a single round).

Further, consider a star graph (as the ground truth network) with n nodes such that propagation probability *p(e)* = 1 on all edges of the star graph (shown in Figure 1). Now, any algorithm for the DIME problem would select a single node in the *uncertain network* uniformly at random with equal probability of *1/n* (as information about all nodes is symmetrical). In expectation, the algorithm will achieve an expected reward . However, given the ground truth network, we get , because we always select the star node. As n goes to infinity, we can at best achieve a approximation to . Thus, no algorithm can achieve a approximation to for any

8

*Computational Hardness*

We now analyze the hardness of computation in the DIME problem in the next two theorems.

**Theorem 2** The DIME problem is NP-Hard.

**Proof** Consider the case where  and . This degenerates to the classical influence maximization problem which is known to be NP-hard. Thus, the DIME problem is also NP-hard.

Some NP-Hard problems exhibit nice properties that enable approximation guarantees for them. Golovin et. al. [11] introduced adaptive submodularity, an analog of submodularity for adaptive settings. Intuitively, adaptive submodularity deals with cases in which actions/items are to be picked in multiple stages, and newer information is revealed every time an action is picked. Adaptive submodularity requires that the expected marginal gain of picking an action can only decrease as more actions are picked and more information is revealed. Formally, adaptive submodularity requires that , where  represents the marginal gain/benefit of picking action *A*, conditioned on getting information . This makes the adaptive submodularity framework a natural fit for the DIME problem. Presence of adaptive submodularity ensures that a simply greedy algorithm provides a *(1-1/e)* approximation guarantee w.r.t. the optimal solution defined on the *uncertain network*. However, as we show next, while DIME can be cast into the adaptive stochastic optimization framework of [11], our influence function is not adaptive submodular, because of which their Greedy algorithm does not have a *(1-1/e)* approximation guarantee.

**Theorem 3** The influence function of DIME is not adaptive submodular.

**Proof** The definition of adaptive submodularity requires that the expected marginal increase of influence by picking an additional node is more when we have less observation. Here the expectation is taken over the random states that are consistent with current observation. We show that this is not the case in DIME problem. Consider a path with 3 nodes *A, B and C* and two directed edges  and . Let  i.e., propagation probability is 1; and  for some small enough  to be set. Thus, the only uncertainty comes from incomplete knowledge of the existence of edges.

Let us assume that we pick node *A*. After picking node *A*, the expected marginal benefit of picking node *C* is . However, after picking node *B*, if we get information , then the expected marginal benefit of picking node *C* goes to 1 (up from ). Since the expected marginal benefit of picking node *C* increased from  to 1 upon receiving more information and picking more actions, this contradicts the definition of adaptive submodularity. This shows that the influence function of DIME is not adaptive submodular.

## HEAL: DIME PROBLEM SOLVER

The above theorems show that DIME is a hard problem as it is difficult to even obtain any reasonable approximations. HEALER models DIME as a partially observable Markov decision process (POMDP) [24], which is a logical fit for the problem because of two reasons. First, several interventions are conducted sequentially, similar to sequential POMDP actions. Second,

there is *partial observability* (similar to POMDPs) due to uncertainties in network structure and influence status of nodes. We now provide a high level overview of HEALER's POMDP model.

## POMDP Model

A *state* in this model includes the influence status of all network nodes (i.e., which nodes are influenced and which nodes are not) and the true state of the uncertain edges (i.e., whether each uncertain edge exists or not in the real world). Thus, there are   possible POMDP states in a network with *N* nodes and *M* uncertain edges. Similarly, an *action* in this model is any possible subset of *K* network nodes, which can be called for an intervention. Thus, if *K* nodes are being selected in every intervention on a network with *N* nodes, there are possible POMDP actions. Finally, an *observation* in this model is based on the assumption that when a set of *K* nodes (i.e., *K* distinct homeless youth) are called in for intervention, the shelter officials can talk to these nodes (or youth) and resolve the status of the uncertain edges in their local neighborhood. Specifically, the shelter official observes the true state of each uncertain edge (i.e., whether it exists in the real world or not) outgoing from the *K* nodes chosen in that action. The observation of the true state of uncertain edge *(A,B)* leads to resetting of *u(A,B)* to either 1 or 0 (depending on whether edge *(A,B)* actually exists or not). Thus, when *M* uncertain edges are outgoing from the *K* nodes chosen in a POMDP action, there are   possible POMDP observations. Finally, the *rewards* in this model keep track of the number of new nodes that get influenced upon taking a POMDP action. Refer to Yadav et al. [9] for the full POMDP model.

## HEAL

HEAL is a heuristic based online POMDP planner for solving the DIME problem. HEAL solves the *original POMDP* using a novel *hierarchical ensembling heuristic*: it creates ensembles of imperfect (and smaller) POMDPs at *two* different layers, in a hierarchical manner (see **Figure 2**). HEAL's *top layer* creates an ensemble of smaller sized *intermediate POMDPs* by subdividing the original *uncertain network* into several smaller sized *partitioned networks* by using graph partitioning techniques [25]. Each of these partitioned networks is then mapped onto a POMDP, and these *intermediate POMDPs* form the *top layer* ensemble of POMDP solvers.

In the bottom layer, each *intermediate POMDP* is solved using TASP (**t**ree **a**ggregation for **s**equential **p**lanning), HEAL's POMDP planner, which subdivides the POMDP into another ensemble of smaller sized *sampled POMDPs*. Each member of this *bottom layer* ensemble is created by randomly sampling uncertain edges of the partitioned network to get a sampled network having no uncertain edges, and this sampled network is then mapped onto a *sampled POMDP*. Finally, the solutions of POMDPs in both the *bottom* and *top layer* ensembles are aggregated using novel techniques to get the solution for HEAL's original POMDP.

These heuristics enable scale up to real-world sizes (at the expense of sacrificing performance guarantees), as instead of solving one huge problem, HEAL solve several smaller problems. The primary difference between HEAL and PSINET (the previous state-of-the-art) is in the top layer

of HEAL, which uses the graph partitioning heuristic. This heuristic divides up the network into different partitions, with each partition corresponding to an intermediate POMDP (Figure 2). The partitions are chosen in a way which minimizes the number of cross-edges going across the partitions (while ensuring that the partitions have similar sizes. Since these partitions are almost disconnected, we solve each partition separately without accounting for influence going across the partitions. Simulations show that even on smaller settings, HEAL achieves a 100-fold speed up over PSINET, while providing a 70% improvement in solution quality; and on larger problems, *where PSINET is unable to run at all*, HEAL continues to provide high solution quality. This raises the following question: Why does the graph partitioning heuristic (a seemingly counter-intuitive heuristic) work so well in simulation? To provide an answer, we next analyze the structure of the real-world social networks of homeless youth that were used to simulate performance of HEAL and PSINET.

## Small-World Nature of Real World Networks

The small-world network model is known to mimic many properties of real-world networks. In fact, these small-world networks are observed in many different domains, ranging from biological, social and technological networks. While there is no rigid definition that classifies a social network as small-world or not, there are two widely accepted network characteristics that small-world networks should possess. First, the average path length (i.e., the average distance between any two nodes in the network) in a small-world network should be somewhat comparable to the average path length of a random network. Second, the average clustering coefficient (i.e., the average *cliquish*-ness of the network) should be significantly higher than the average clustering coefficient of a random network. This means that small-world networks have relatively low average path lengths (as the average path length in a random network is low), and consist of lots of network cliques (as the average clustering coefficient is higher than that in a random network).

The two networks of homeless youth used in our simulations are small-world networks, and hence well suited to HEAL's graph partitioning heuristic. HEAL's graph partitioning technique takes advantage of the clustering in these small-world networks by identifying the cliques in them, and considering their influence independently.

**Table 1** compares the average clustering coefficient and the average path length with the random clustering coefficient and random average path length for two different real-world networks of homeless youth. This table shows that in both networks, the average path length is comparable to that in a random network, whereas the average clustering coefficient is 8-fold higher on average. This explains why HEAL's graph partitioning heuristic works well, as it allows HEAL to find the different cliques in these small-world networks, which are then influenced independently.

Having established that the graph partitioning heuristic is the primary reason behind HEAL's good performance, we now try to see if we can augment other baseline algorithms with the same heuristic, and provide comparison results.

11

## Results comparing HEAL with Augmented Baselines

We use three algorithms as baselines: Greedy, Greedy with partitioning and Degree Centrality. We use the modified Greedy algorithm used by Yadav et. al. [9] (referred to as *Greedy* in **Figure 3**), their closest competitor to HEAL. Recall that even though Greedy has no theoretical guarantees in our domain (Theorem 3), we still want to test its empirical performance. Further, we augment this Greedy algorithm with the graph partitioning heuristic, i.e., we partition the network into different cliques, and then use the same Greedy algorithm inside each clique independently to find which nodes should be influenced in the network. This augmented Greedy algorithm is referred to as "*Greedy + Partition*" in Figure 3. We do not report results of PSINET as even with the graph partitioning heuristic, PSINET runs out of memory on our case study networks. Finally, we also compare against Degree Centrality (i.e., picking the nodes with the highest degree), the current modus operandi of homeless shelters.

This experiment was run on a 2.33 GHz 12-core Intel machine having 48 GB of RAM, and was averaged over 100 runs. We use a metric of "*Indirect Influence*" for comparison between different algorithms, which is number of nodes "*indirectly*" influenced by intervention participants. For example, on the homeless youth network having 170 nodes, by selecting 2 nodes (i.e., $K=2$) each for 10 interventions (horizon) (i.e., $T=10$), 20 nodes (a lower bound for any strategy) are influenced with certainty. However, the total number of influenced nodes might be 26 (say) and thus, the *Indirect Influence* is 26-20 = 6. In all experiments, the propagation and existence probability values on all network edges were uniformly set to 0.1 and 0.6, respectively. This was done based on findings in Kelly et. al. [26]. This comparison result is statistically significant under bootstrap-t (.

Figure 3 compares HEAL, Greedy, the "Greedy + Partition" approach and Degree Centrality (DC) on the two real-world networks of homeless youth in simulation (we see similar results on many other networks [9]). Each of these networks had around 170 nodes and 250 edges. The x-axis shows the two different networks, and the y-axis shows the indirect influence achieved. First, this figure shows that a static approach like DC performs very poorly as compared to HEAL, an adaptive POMDP based solution. This is in part due to HEAL's policy being responsive and flexible enough to incorporate any new observations that are seen during execution of HEAL's policy, in order to improve future decisions taken by HEAL. In contrast, DC does not change its node selection in future time steps, regardless of the observations received. In order to compare the flexibility of HEAL's policy with DC's policy, we analyzed the number of times HEAL chose different actions (i.e. chose different nodes) in future time steps, based on receiving different observations in earlier time steps. On one of the networks, HEAL chose an action in the first round, which led to four possible different observations. Corresponding to each of these four different observations (which we simulated), HEAL chose a different action in the next round. This is in comparison to DC which never modifies its chosen actions in response to getting different observations.

Second, this figure also shows that graph partitioning improves the performance of Greedy on only one network (due to the community structure of the networks) and not the other. Moreover, HEAL still outperforms the nearest competitor by ~32%. This is a secondary confirmation that even though the graph partitioning heuristic plays a major role in the determination of the solution quality, it is by no means the only reason for HEAL's superior performance, thereby illustrating the importance of the TASP solver (see [9] for details) in HEAL as well.

Third, the superiority of HEAL over Greedy and "Greedy + Partition" illustrates the importance of look-ahead search done by HEAL, which leads to higher solution qualities (i.e., indirect influence spread). On the other hand, Greedy does not do any look ahead search, thereby leading to lesser solution qualities.

## Conclusion

In this paper, we explored the reasons and motivations behind the design and superior performance of HEALER, an adaptive software agent which recommends intervention attendees to homeless shelter officials. HEALER solves a POMDP on a social network to come up with recommendations for which homeless youth in a social network should be chosen as intervention attendees. We first formally characterized the computational problem (called DIME) solved by HEALER, and showed that it is an NP-Hard problem. Moreover, well-known algorithms such as Greedy lose their approximation guarantees in the DIME problem due to the feedback about network structure received during interventions. We inferred that DIME's computational hardness forces HEALER to rely on heuristic methods for solving DIME. Further, we analyzed these heuristic methods, and showed that the primary reason behind the superior performance of HEALER is its graph partitioning heuristic, which works well due to the small-world nature of the real-world networks of homeless youth. However, we showed that the graph partitioning heuristic is not the only reason for HEALER's superior performance, as other baseline algorithms augmented with the graph partitioning heuristic don't perform as well. HEALER was recently tested in the real-world with 60 homeless youth, and it outperformed other baselines [11]. To the best of our knowledge, this is the first such evaluation of an influence maximization algorithm in the field.

## References

1. National HCH Council, "HIV/AIDS among Persons Experiencing Homelessness: Risk Factors, Predictors of Testing, and Promising Testing Strategies". [Online]. Available: www.nhchc.org/wp-content/uploads/2011/09/InFocus_Dec2012.pdf .
2. E. Rice, E. Tulbert, J. Cederbaum, A. B. Adhikari, N. G. Milburn, "Mobilizing Homeless Youth for HIV Prevention: a Social Network Analysis of the Acceptability of a face-to-face and Online Social Networking Intervention", *Health Education Research*, vol. 27, no. 2, pp. 226, 2012.

3. E. Rice, "The Positive Role of Social Networks and Social Networking Technology in the Condom-using Behaviors of Homeless Young People". *Public Health Reports*. vol. 125, no. 4, pp. 588, 2010.

4. E. Rice, A. Fulginti, H. Winetrobe, J. Montoya, A. Plant, T. Kordic, "Sexuality and Homelessness in Los Angeles public schools", *American Journal of Public Health*, vol. 102, 2012.

5. C. Borgs, M. Brautbar, J. Chayes, B. Lucier, "Maximizing Social Influence in Nearly Optimal Time", *Proc. 25th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 946-957, 2014.

6. Y. Tang, X. Xiao, Y. Shi, "Influence Maximization: Near-Optimal Time Complexity meets Practical Efficiency". *Proc. 2014 ACM SIGMOD International Conference on Management of Data*, pp. 75-86, 2014.

7. D. Kempe, J. Kleinberg, E. Tardos, "Maximizing the Spread of Influence through a Social Network", *Proc. 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp.* 137-146, 2003.

8. J. Leskovec, A. Krause, C. Guestrin, C. Faloutsos, J. VanBriesen, N. Glance, "Cost Effective Outbreak Detection in Networks", *Proc. 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp.* 420-429, 2007.

9. A. Yadav, H. Chan, A.X. Jiang, H. Xu, E. Rice, M. Tambe, "Using Social Networks to Aid Homeless Shelters: Dynamic Influence Maximization Under Uncertainty", (pp. 740-748) *Proc. International Conference on Autonomous Agents and Multiagent Systems (AAMAS) 2016*, Singapore.

10. A. Yadav, L. Marcolino, E. Rice, R. Petering, H. Winetrobe, H. Rhoades, M. Tambe, H. Carmichael, "Preventing HIV Spread in Homeless Populations Using PSINET" in *Proc. 27th Conference on Innovative Applications of Artificial Intelligence (IAAI)*, Austin, 2015.

11. A. Yadav, B. Wilder, E. Rice, R. Petering, J. Craddock, A. Yoshioka-Maxwell, M. Hemler, L. Onasch-Vera, M. Tambe, D. Woo, "Influence Maximization in the Field: The Arduous Journey from Emerging to Deployed Application" in *International Conference on Autonomous Agents and Multiagent Systems (AAMAS) 2017*, Sau Paulo, Brazil.

12. D. Golovin, A. Krause, "Adaptive Submodularity: Theory and Applications in Active Learning and Stochastic Optimization". *Journal of Artificial Intelligence Research*. vol. 42, pp. 427-486, 2011.

13. Lei, S., Maniu, S., Mo, L., Cheng, R., & Senellart, P. (2015, August). Online influence maximization. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 645-654). ACM.

14. J. Schneider, A. Ning. Zhou, E. O. Laumann, "A new HIV Prevention Network Approach: Sociometric Peer Change Agent Selection", *Social Science & Medicine*, vol. 125, pp. 192-202, 2015.

15. T. W. Valente, P. Pumpuang, "Identifying Opinion Leaders to Promote Behavior Change", *Health Education and Behavior*, pp.881-896, 2007.

16. E. Cohen, D. Delling, T. Pajor, R. F. Werneck, "Sketch-based Influence Maximization and Computation: Scaling up with guarantees". *Proc. 23rd ACM International Conference on Information and Knowledge Management, pp. 629-638, 2014.*

17. S. Ross, J. Pineau, S. Paquet, B. Chaib-Draa, "Online Planning Algorithms for POMDPs", *Journal of Artificial Intelligence Research*, pp. 663-704, 2008.

18. D. Silver, J. Veness, "Monte-Carlo Planning in large POMDPs", *Proc. Advances in Neural Information Processing Systems (NIPS)*, pp. 2164-2172, 2010.

19. A. Somani, N. Ye, D. Hsu, W. S. Lee, "DESPOT: Online POMDP Planning with Regularization". *Proc. Advances in Neural Information Processing Systems (NIPS)*, pp. 1772-1780, 2013.

20. S. D. Young, E. Rice, "Online Social Networking Technologies, HIV knowledge, and Sexual Risk and Testing Behaviors among Homeless Youth", *AIDS and Behavior*, vol. 15, no. 2, pp. 253-260, 2011.

21. Kim, M., & Leskovec, J. (2011, April). The network completion problem: Inferring missing nodes and edges in networks. In *Proceedings of the 2011 SIAM International Conference on Data Mining* (pp. 47-58). Society for Industrial and Applied Mathematics.

22. Q. Yan, S. Guo, D. Yang, "Influence Maximizing and Local Influenced Community Detection based on Multiple Spread Model", *Advanced Data Mining and Applications*, pp. 82-95, 2011.

23. J. P. Cointet, C. Roth, "How Realistic Should Knowledge Diffusion Models Be?", *Journal of Artificial Societies and Social Simulation*, vol. 10, no. 3, pp. 5, 2007.

24. Kaelbling, L. P., Littman, M. L., & Cassandra, A. R. (1998). Planning and acting in partially observable stochastic domains. *Artificial intelligence*, *101*(1), 99-134.

25. D. LaSalle, G. Karypis, "Multi-threaded Graph Partitioning", *Proc. 27th IEEE International Symposium on Parallel & Distributed Processing (IPDPS)*, pp. 225-236, 2013.

26. J. A. Kelly, D. A. Murphy, K. J. Sikkema, T. L. McAuliffe, R. A. Roffman, L. J. Solomon, R. A. Winett, and S. C. Kalichman. "Randomised, Controlled, Community-Level HIV-Prevention Intervention for Sexual-Risk Behaviour among Homosexual men in US cities.", *The Lancet*, vol. 350(9090). pp. 1500, 1997.

## Bios

**Amulya Yadav** *USC Center for Artificial Intelligence in Society, Computer Science Department, University of Southern California (USC), Los Angeles 90089 USA (amulyaya@usc.edu).* Mr. Yadav is a Ph.D. Candidate in the Computer Science Department of the USC Viterbi School of Engineering. He holds a Bachelors in Technology in Computer Science and Engineering from Indian Institute of Technology Patna. His research interests include influence maximization on social networks, multi-agent sequential decision making (MSDM) problems and game theory/ mechanism design.

**Hau Chan** *Trinity University, Computer Science Department, San Antonio, 78212 USA (hchan@trinity.edu).* Dr. Chan is currently a postdoctoral research associate (postdoc) at Trinity University. He obtained his Ph.D. in Computer Science at Stony Brook University in 2015. He

works mainly in the area of computational game theory (CGT) and AI. More specifically, his interests lie in modeling and representing social science problems compactly, and finding efficient algorithms and heuristics to compute Nash and Stackelberg equilibrium efficiently (by leveraging the compact representation).

**Albert Xin Jiang** *Trinity University, Computer Science Department, San Antonio, 78212 USA (xjiang@trinity.edu).* Dr. Jiang is an assistant professor in the Department of Computer Science at Trinity University. He received his PhD from the Department of Computer Science at the University of British Columbia, and was a postdoctoral research associate in the TEAMCORE research group at the Department of Computer Science at the University of Southern California. Much of his research is addressing computational problems arising in game theory, including the efficient computation of solution concepts such as Nash equilibrium, Stackelberg equilibrium and correlated equilibrium, as well as applications of game-theoretic computation to real-world domains such as large-scale infrastructure security and electronic commerce.

**Haifeng Xu** *USC Center for Artificial Intelligence in Society, Computer Science Department, University of Southern California (USC), Los Angeles 90089 USA (haifeng.ustc@gmail.com).* Mr. Xu is a Ph.D. Candidate in the Computer Science Department of the USC Viterbi School of Engineering. His research interests include computational game theory and mechanism design, auction theory, and signaling schemes in Stackelberg security games.
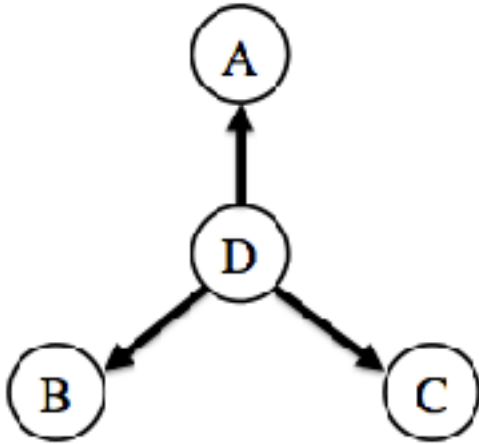
**Eric Rice** *USC Center for Artificial Intelligence in Society, School of Social Work, University of Southern California (USC), Los Angeles 90089 USA (ericr@usc.edu).* Dr. Rice is an Associate Professor at the University of Southern California's School of Social Work. He has been working on issues of youth homelessness since 2003. His work focuses on social networks, HIV prevention, and housing issues. He is committed to working with communities to end homelessness for youth.

**Robin Petering** *USC Center for Artificial Intelligence in Society, School of Social Work, University of Southern California (USC), Los Angeles 90089 USA (petering@usc.edu).* Ms. Petering is currently a Ph.D. Candidate at the USC School of Social Work. Ms. Petering holds a Masters in Social Work from the University of California, Los Angeles. Her current research interests include youth homelessness, gang involved youth and related issues of violence.
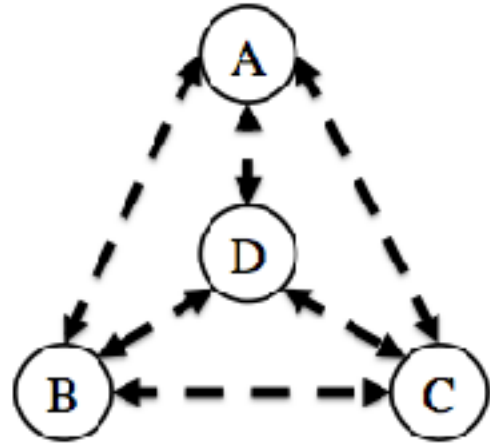
**Milind Tambe** *USC Center for Artificial Intelligence in Society, Computer Science Department, University of Southern California (USC), Los Angeles 90089 USA (tambe@usc.edu).* Dr. Tambe is the Helen N. and Emmett H. Jones Professor in Engineering at the University of Southern California. He is a fellow of AAAI (Association for the Advancement of Artificial Intelligence) and ACM (Association for Computing Machinery), and recipient of the ACM/SIGART (ACM Special Interest Group on Artificial Intelligence) Autonomous Agents Research Award, INFORMS (Institute for Operations Research and Management Sciences) Wagner prize for excellence in Operations Research practice, Rist Prize of the Military Operations Research Society, Christopher Columbus Fellowship Foundation Homeland security award, International

Foundation for Agents and Multiagent Systems influential paper award, Meritorious Team Commendation from the Commandant of the US Coast Guard and LAX (Los Angeles International Airport) Police, and Certificate of Appreciation from US Federal Air Marshals Service.

**Figures and Tables**



Ground Truth Network                    Uncertain Network

**Figure 1** Illustration of the value of information in the DIME Problem. A, B, C and D represent nodes, and the edges between them represent friendships. There are two kinds of edges: certain edges (denoted by solid edges as shown in the left figure); and uncertain edges (denoted by dotted edges as shown in the right figure). The propagation and existence probabilities on all the edges is assumed to be fixed.
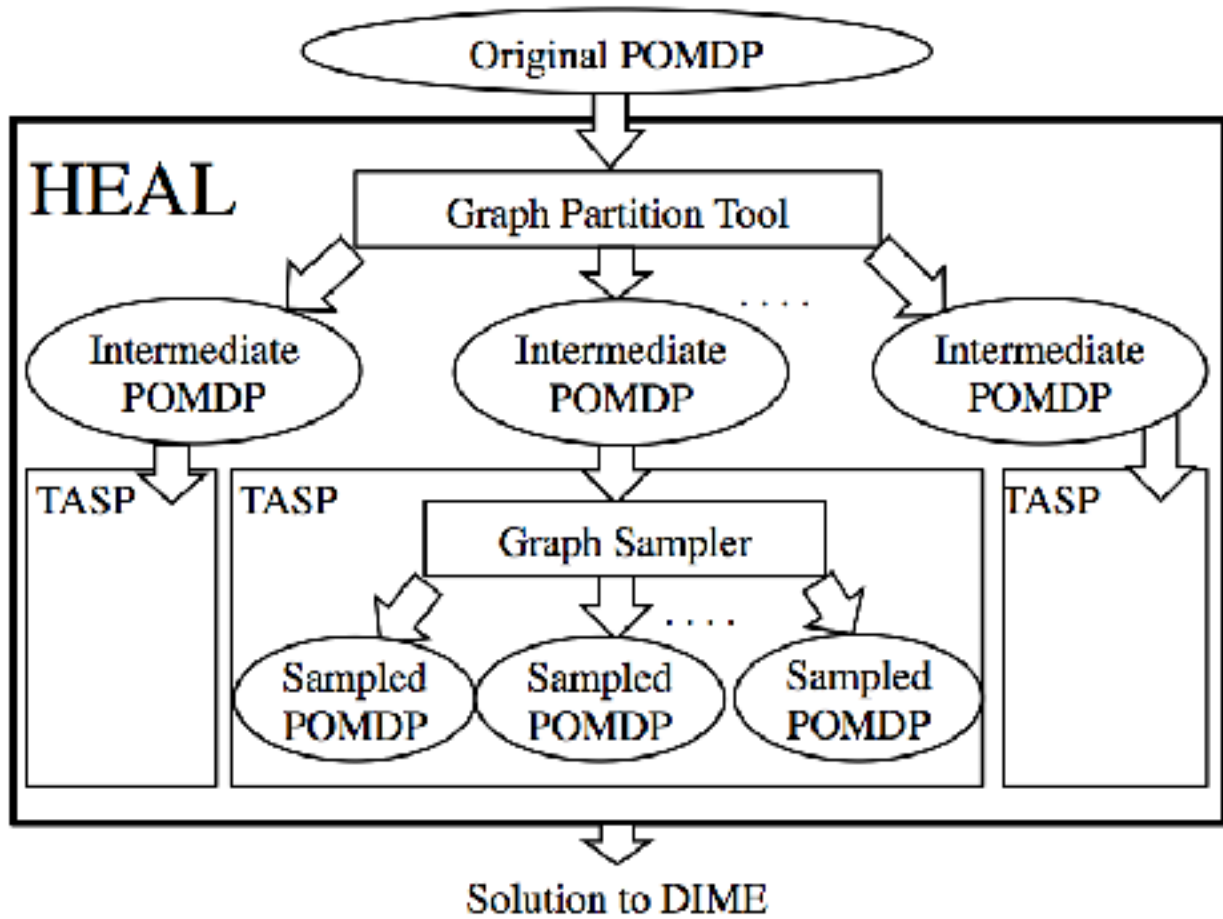
**Figure 2** Hierarchical decomposition in HEAL (From Yadav et. al. [9])

| Networks | Avg.<br>Path Length In<br>Real Nets | Avg.<br>Path Length<br>In Random Nets | Avg.<br>Clustering Coefficient<br>In Real Nets | Random<br>Clustering Coefficient<br>In Random Nets |
|---|---|---|---|---|
| Net 1 | 4.7758 | 3.1464 | 0.1752 | 0.0340 |
| Net 2 | 5.8369 | 3.8368 | 0.1746 | 0.0208 |

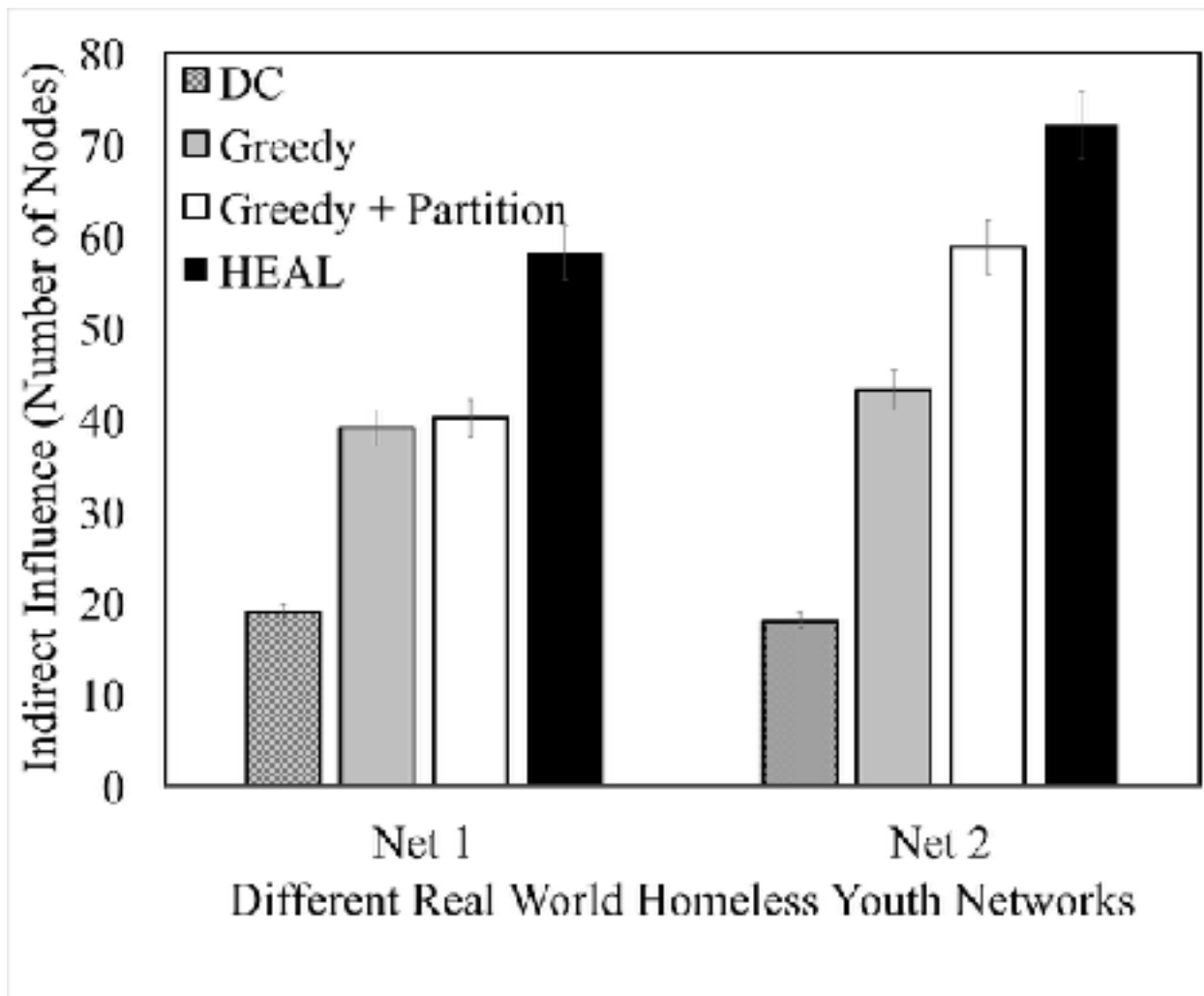**Table 1** Small World Characteristics of Real World Networks of Homeless Youth

**Figure 3** Influence Spread Comparison of HEAL with Augmented Baseline Algorithms