

# Diffusion Archaeology for Diffusion Progression History Reconstruction

Emre Sefer  
Carnegie Mellon University  
Email: esefer@cs.cmu.edu

Carl Kingsford  
Carnegie Mellon University  
Email: carlk@cs.cmu.edu

**Abstract**—Diffusion through graphs can be used to model many real-world process, such as the spread of diseases, social network memes, computer viruses, or water contaminants. Often, a real-world diffusion cannot be directly observed while it is occurring — perhaps it is not noticed until some time has passed, continuous monitoring is too costly, or privacy concerns limit data access. This leads to the need to reconstruct how the present state of the diffusion came to be from partial diffusion data. Here, we tackle the problem of reconstructing a diffusion history from one or more snapshots of the diffusion state. This ability can be invaluable to learn when certain computer nodes are infected or which people are the initial disease spreaders to control future diffusions. We formulate this problem over discrete-time SEIRS-type diffusion models in terms of maximum likelihood. We design methods that are based on submodularity and a novel prize-collecting dominating-set vertex cover (PCDSVC) relaxation that can identify likely diffusion steps with some provable performance guarantees. Our methods are the first to be able to reconstruct complete diffusion histories accurately in real and simulated situations. As a special case, they can also identify the initial spreaders better than existing methods for that problem. Our results for both meme and contaminant diffusion show that the partial diffusion data problem can be overcome with proper modeling and methods, and that hidden temporal characteristics of diffusion can be predicted from limited data.

## I. INTRODUCTION

Dynamic processes over networks are used to model and analyze properties of various social and biological systems. Diffusion is special case of those processes in which a spread (e.g., an infection) starts from some part of the graph and spreads to other portions over time via edges of the graph. Some examples are virus propagation in computer networks [1], and idea and gossip spreading in social networks [2]. A diffusion model, such as the commonly studied SIRS and SEIRS models, defines the set of possible states that the nodes of the graph can be in and rules for probabilistically switching between states. Recently, [3] introduced the *VPM* (Virus Propagation Model) that generalizes all those Markovian diffusion models and defines the hierarchical relationships between them.

It is not always easy to know the whole diffusion progression, initial diffusion conditions, or the time it started due to several limitations. For example, existence of a computer virus diffusion over the computer network may only be noticed after a significant number of computers stop operating. A similar problem exists in detecting influenza diffusion [4]. We may also not track the diffusion of a virus in email networks and a contaminant in a water distribution network due to privacy

and physical limitations, respectively. In all these cases, it is essential to learn more about the past to take precautions to prevent future epidemics, to learn more about the true diffusion mechanics, to provide safer water, to break privacy and so on. However, given present-day diffusion data, it is not trivial to search for the most likely diffusion progression in the past as there will be many valid histories leading to the observed data.

In this paper, we tackle this problem of inferring a complete diffusion history from one or more diffusion snapshots for discrete-time SEIRS-type diffusion models that include SI, SIS, SIR, SIRS, SEIR, SEIS, SEIRS. Those models with their abstract states and independent cascade (IC) assumption [5] have been used to model many forms of diffusion in many different domains [2], [4], [6].

Complete diffusion history reconstruction has not been previously studied but similar problems exist in the literature. The most relevant such problem is *Initial Spreader Identification* where we want to identify the most probable initial infected nodes that started a diffusion. Among approaches for this problem, *Keffectors* [7] identifies the  $k$  best-possible initial spreaders. However, it requires an estimate of the number of initial spreaders to be given as input. *Rumor* [8] finds the most probable spreader by estimating the rumor centrality but it assumes a single initial spreader, and it is only defined for SI model. Lastly, *NetSleuth* [9] does not require the number of initial spreaders as an input. However, it works only for the restricted cases of SI model, and it is based on MDL principle without any provable performance guarantee. None of these methods infer the whole diffusion progression, as our approaches do. Another related problem is *Graph Inference* where we want to reconstruct the unknown graph from observed multiple diffusion traces over it. This problem is fundamentally different than the history reconstruction problem as graph inference methods [10], [11] search in the graph space assuming full observability of multiple traces whereas our methods search in temporal diffusion progression space as they try to complete the missing history of a single trace.

We formulate the diffusion history reconstruction problem as that of determining the maximum likelihood (ML) history given diffusion snapshots that may come from multiple time points. We designed an algorithm called *DHR-sub* (submodular history reconstruction on discrete dynamics) that reconstructs the history before the earliest measured time point by greedily maximizing the non-monotone submodular log-likelihood at each previous time step. It further reconstructs the history between the consecutive diffusion data time points by solving the problem as non-monotone submodular maximization under

matroid base constraints.

Though accurate and practical for smaller graphs, *DHR-sub* can take some time to solve. To reconstruct diffusion history faster, we designed *DHR-pcdsvc* that solves the first-order Taylor approximation relaxation of the log-likelihood. We define this new problem as *Prize-Collecting Dominating-Set Vertex Cover*, and show that it can be approximated within a factor of  $O(\log(|V|))$ . This problem can be further relaxed by removing the covering constraints; it becomes *Prize-Collecting Vertex Cover*, and we design *DHR-pcvc* that approximates it by a factor of 2 for non-bipartite models and solves this newer relaxation optimally by transforming it to s-t mincut for bipartite models. We also design ensemble approaches for all of our methods that estimate the robust set of initial spreaders from multiple runs of the algorithm.

In summary, our main contributions are:

- Our methods reconstruct the whole diffusion history nonparametrically for all SEIRS-type models whereas the existing methods only identify the initial spreaders for certain models;
- Our methods formulate the problem in terms of diffusion likelihood, and we give some performance guarantees on the quality of the obtained solutions;
- Our relaxation methods *DHR-pcdsvc* and *DHR-pcvc* scale well to history reconstruction over tens of thousands of nodes with provable performance guarantees;
- Our methods perform better by using the diffusion information from all the nodes (not just from infected nodes), and from multiple time points if available;
- We use reconstructed histories to predict several diffusion features such as speed and acceleration that are not apparent in the observed portion of the diffusion.

Our methods more accurately identify initial spreader sites on a water distribution network and on simulated networks. In terms of history reconstruction, we compared our methods with a baseline heuristic since there is no previous method. All our methods can accurately reconstruct several meme diffusion histories on blog networks. They also perform better on synthetic networks under different models. In general, all our methods reconstruct the diffusion history reasonably fast and accurately compared to the hardness of the problem (see Section VII-F). In many cases, relaxations of the original problem can reconstruct the diffusion history almost as good as the original formulations in a far shorter amount of time. Lastly, we also estimate the speed and acceleration dynamics of several memes over blog network from their reconstructed histories. In this case, estimated dynamics from quite a few diffusion snapshots match the true dynamics almost perfectly producing decent whole history reconstruction performance. Overall, our results for different types of diffusion show that many characteristics of complete diffusion history can be inferred with proper modeling and methods.

## II. SEIRS DIFFUSION DYNAMICS

SEIRS diffusion dynamics over directed graph  $G = (V, E)$  with possible state transitions is shown in Figure 1. The SEIRS

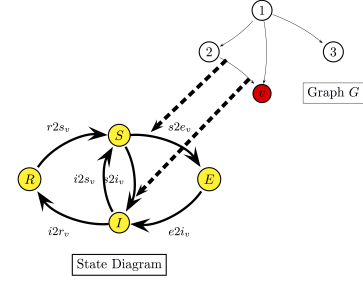


Fig. 1. SEIRS State Transition Diagram

Symbol	Definition and Description
$G = (V, E)$	directed graph $G$
$P(v), S(v)$	set of predecessors, successors of node $v$
$S$	SEIRS states ( $S, E, I, R$ )
$\mathcal{M}$	diffusion model of SEIRS-type models
$p_{uv}$	probability of diffusion from infected node $u$ to susceptible node $v$
$e2i_v, i2r_v, r2s_v, s2e_v, s2i_v, i2s_v$	probability of ( $E \rightarrow I, I \rightarrow R, R \rightarrow S, S \rightarrow E, S \rightarrow I, I \rightarrow S$ ) transition for $v$
$t_v^s, t_v^e, t_v^i, t_v^r$	time $v$ transitions into ( $S, E, I, R$ )
$S_t, E_t, I_t, R_t$	set of nodes that are in ( $S, E, I, R$ ) at time $t$
$D_t$	diffusion snapshot at time $t$
$\mathcal{D}$	given diffusion data
$l_D$	length of diffusion $D$
$T_D$	ordered set of time points of $\mathcal{D}$
$t_{min}, t_{max}$	$\min(T_D), \max(T_D)$
$f_{min}, f_{max}$	$\min(T_D)/l_D, \max(T_D)/l_D$

TABLE I. TABLE OF SYMBOLS

states are Susceptible ( $S$ ), Exposed but not contagious ( $E$ ), Infected and contagious ( $I$ ), and previously infected but now Recovered (or immune to the infection) ( $R$ ). Those states are general enough abstractions to model various forms of diffusion in different contexts [6], [4]. For instance, the infected state models people having influenza symptoms in influenza diffusion over humans, and it represents creation of a blog entry about a topic in idea diffusion. Similarly, the recovered state could represent recovery of a person from influenza or the decontamination of a water tower from chemical contaminants depending on the context.

In SEIRS model, diffusion starts at time  $t = 0$  from set of initially infected nodes and progresses over  $G$  in discrete time steps. Let  $S_t, E_t, I_t, R_t$  be the set of  $S, E, I, R$  nodes at time  $t$  respectively. At each time step, infected nodes spread the infection to the susceptible nodes with certain probability. This  $S \rightarrow E$  transition is exogenous; it is affected by  $G$  and probability of exogenous transition for susceptible node  $v$  at time  $t$  is  $1 - \prod_{u \in P(v) \cap I_t} (1 - p_{uv})$ , where  $P(v)$  is the set of nodes with edges into  $v$  and  $p_{uv}$  is the probability of transmission of the agent over edge  $(u, v)$ . The remaining  $E \rightarrow I, I \rightarrow R, R \rightarrow S$  transitions are endogenous; their transition probabilities are  $e2i_v, i2r_v, r2s_v$  respectively, and they are not affected by  $G$ . For every node at each time step, if a transition succeeds, the node transitions to a new state. Otherwise, it follows similar procedure at next time step, independent of the previous trials. SEIRS-type models are Markovian since state of a node at time  $t$  depends on its state and its neighbors' states at previous time steps, and it obeys independent cascade (IC) [5] assumption which states that a diffusion from one of nodes predecessor is enough for node to become exposed/infected. The symbols used in this

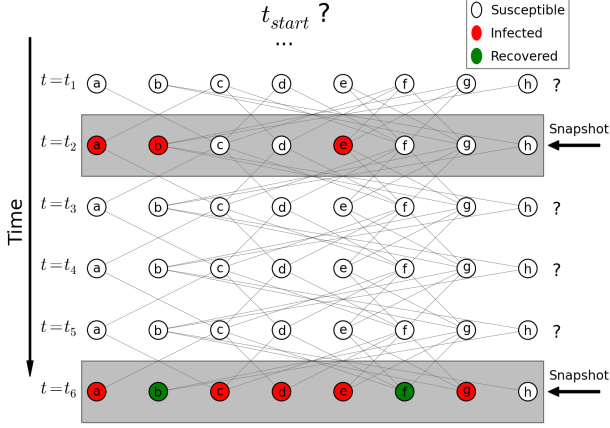


Fig. 2. Example Problem: SIR Diffusion over 8 node graph where we can only observe  $t_2$  and  $t_6$  without knowing the initial diffusion time  $t_{start}$ . We want to reconstruct the missing diffusion snapshots from  $t_{start}$  onwards

text are given in Table I for reference.

SEIRS-type models include the well-known SI, SIR, SIS, SIRS, SEIR, SEIRS models [12]. SEIRS is the most general model among these models, and some of its transitions disappear or change slightly in other models. For instance, in SIR, there is no exposed state; the exogenous transition is  $S \rightarrow I$  since nodes proceed directly to the infected state, and there is no  $R \rightarrow S$  transition. We can classify SEIRS-type models in various ways. SIRS, SEIRS are *loopy* models where  $R \rightarrow S$  transition is available whereas SI, SIR, SEIR are *non-loopy* models. We can also split SEIRS-type models into *bipartite* and *non-bipartite* models: a node that gets the infection directly transitions into infected state in non-bipartite models such as SI, SIR, SIRS, SIS whereas it goes through the exposed state for bipartite models such as SEIR, SEIRS. The model may also be either *uniform* in which case all of the transition probabilities are the same for each edge and node, or *non-uniform* in which case the probabilities may vary over the edges and nodes. We discuss the general non-uniform case here; the uniform case is a simple specialization.

### III. DIFFUSION HISTORY RECONSTRUCTION PROBLEM

For diffusion  $D$ , let  $D_t = (S_t, E_t, I_t, R_t)$  be the state of the nodes at the time  $t$ , where  $S_t$  is the set of susceptible nodes, etc.  $D_t$  is a *diffusion snapshot*. We define Problem 1 to reconstruct the diffusion history when the diffusion length is unknown:

**Problem 1.** We are given: a graph  $G = (V, E)$ , state transition probabilities  $(p_{uv}, e2i_v, i2s_v, i2r_v, r2s_v)$  that define an SEIRS-type model, a collection of time points  $T_D$  at which snapshots were taken, and a collection of diffusion snapshots  $\mathcal{D} = \{D_t\}$  for  $t \in T_D$ . Each snapshot records the state of every node at a single time point, partitioning them into  $V = S_t \cup E_t \cup I_t \cup R_t$ .

Our goal is to infer the past states (susceptible, exposed, infected and recovered) of every node at every time  $t \notin T_D$ .

Figure 2 illustrates the history reconstruction problem for the SIR model. Each subsequent layer shows the progression of time, and we want to reconstruct the diffusion progression from

unknown initial time  $t_{start}$  onwards given full state knowledge at subset of time points. The *initial spreader identification problem* is special case of Problem 1 where we want to identify only the initial *infected* nodes.

### IV. NON-MONOTONE SUBMODULAR HISTORY RECONSTRUCTION (DHR-sub)

Let  $t_{max} = \max(T_D)$ ,  $t_{min} = \min(T_D)$ ,  $t_{start}$  be the unknown initial diffusion time, and  $\mathcal{M}$  be the specific SEIRS-type model that diffusion snapshots are collected over, SEIRS-type models are Markovian so the probability of diffusion  $D = \{D_{t_{start}}, \dots, D_{t_{max}}\}$  that starts at  $t_{start}$  and progresses until  $t_{max}$  can be written as the multiplication of the probability of each time step in terms of previous time steps as in (1):

$$P(D) = \prod_{j=t_{start}+1}^{t_{max}} P(D_j | D_{j-1}, \dots, D_{t_{start}}) P(D_{t_{start}}) \quad (1)$$

We assume the state transition probabilities  $(p_{uv}, e2i_v, i2r_v)$  to be same at each time step, so the overall diffusion probability in (1) simplifies to (2) under this memoryless property:

$$P(D) = \prod_{j=t_{start}+1}^{t_{max}} P(D_j | D_{j-1}) P(D_{t_{start}}) \quad (2)$$

Let  $T_D$  be the ordered set of observed time points,  $X_j = (S_j, E_j, I_j, R_j)$  be the unknown state knowledge at time  $\forall j \notin T_D$ , and  $X = \{X_j : j \in \{t_{start}, \dots, t_{max}\} \setminus T_D\}$ . Given a collection of diffusion snapshots  $\mathcal{D}$  for  $T_D$ , our goal is to reconstruct the most probable diffusion progression ( $X$ ) by maximizing the log-likelihood as in Equation (3)–(5):

$$\operatorname{argmax}_X \log(\mathcal{L}(X | \mathcal{D})) = \underbrace{\log(\mathcal{L}^{pre})}_{\text{DHR-sub-early}} + \sum_{(j,k) \in \mathcal{P}(T_D)} \underbrace{\log(\mathcal{L}_{j,k}^{in})}_{\text{DHR-sub-between}} \quad (3)$$

$$\text{s.t. } \mathbf{IntraConsistent}(X_j, \mathcal{M}), j \in t_{start}, \dots, t_{max} - 1 \quad (4)$$

$$\mathbf{InterConsistent}(X_j, X_{j+1}, \mathcal{M}), j \in t_{start}, \dots, t_{max} - 1 \quad (5)$$

where  $\mathcal{L}^{pre}$  and  $\mathcal{L}_{j,k}^{in}$  are defined in (6)–(7),  $\mathcal{P}(T_D) = \{(t_j, t_{j+1}), j \in 1, \dots, |T_D| - 1\}$ , and maximum log-likelihood estimate is the same as the maximum likelihood estimate since the logarithm is a monotonically increasing function:

$$\mathcal{L}_{j,k}^{in} = P(X_{j+1} | D_j) P(D_k | X_{k-1}) \prod_{t \in \{j+1, \dots, k-2\}} P(X_{t+1} | X_t) \quad (6)$$

$$\mathcal{L}^{pre} = \prod_{j=1}^{t_{min} - t_{start}} P(X_{t_{min}-j+1} | X_{t_{min}-j}) P(X_{t_{start}}) \quad (7)$$

There are two types of constraints: **IntraConsistency** constraints (4) make sure that the variable assignments at each time step are valid under  $\mathcal{M}$ : every node belongs to a single state at each  $j$ , and **InterConsistency** constraints (5) make sure that the diffusion between each consecutive time steps is valid according to rules of  $\mathcal{M}$ : every node that got the infection at  $j$  has at least one infected predecessors at  $j - 1$ , and node transitions are valid according to  $\mathcal{M}$ . For instance, recovered nodes cannot become susceptible if  $\mathcal{M}$  is not loopy. These constraints are described in more detail below.

The diffusion history between consecutive observed  $T_D$  pairs is independent of each other since Eq. 2 is memoryless, and each  $D_j$  completely describes states of all nodes at time  $j$ . Thus, maximizing (3)–(5) can be partitioned into multiple independent subproblems of two types that can be optimized independently. The first type (**DHR-sub-early**) maximizes  $\log(\mathcal{L}^{pre})$  under the consistency constraints to reconstruct the history before  $t_{min}$ . The second type (**DHR-sub-between**) maximizes  $\log(\mathcal{L}_{j,k}^{in})$  to reconstruct the history between the snapshots from time  $j$  and time  $k$  under the consistency constraints. We define algorithms for both types of subproblems below. In the text, we use  $D_j$  and  $X_j$  interchangeably for  $\forall j \in T_D$ .

#### A. History reconstruction before the earliest observed snapshot (DHR-sub-early)

To find the most likely diffusion history before  $t_{min}$ , we solve the problem:

$$\operatorname{argmax}_X \log(\mathcal{L}^{pre}) = \sum_{j=t_{start}}^{t_{min}-1} \log\left(P(X_{j+1}|X_j)\right) + \log\left(P(X_{t_{start}})\right) \quad (8)$$

$$\text{s.t. } \mathbf{IntraConsistent}(X_j, \mathcal{M}), j \in t_{start}, \dots, t_{min} - 1 \quad (9)$$

$$\mathbf{InterConsistent}(X_j, X_{j+1}, \mathcal{M}), j \in t_{start}, \dots, t_{min} - 1 \quad (10)$$

We assume a uniform prior  $P(X_{t_{start}})$  over set of initially infected nodes since we do not have any extra information about them. We now discuss how to formulate the objective function and constraints above in terms of binary variables representing each node's state.

1) *Expressing the objective function* (8): Given  $X_j$ , the probability of observing the diffusion snapshot  $X_{j-1}$  at time  $j-1$  can be expressed as:

$$P(X_j|X_{j-1}) = \mathcal{L}(X_{j-1}|X_j) = \mathcal{L}_s^j \mathcal{L}_e^j \mathcal{L}_i^j \mathcal{L}_r^j \quad (11)$$

where  $\mathcal{L}_s^j, \mathcal{L}_e^j, \mathcal{L}_i^j, \mathcal{L}_r^j$  are the likelihoods of the nodes in  $S_j, E_j, I_j, R_j$  states respectively in terms of  $X_{j-1}$ .

To define  $\mathcal{L}_s^j, \mathcal{L}_e^j, \mathcal{L}_i^j, \mathcal{L}_r^j$ , we introduce a single binary variable for each node to define its state at time  $j-1$  given its state at time  $j$ . A binary variable is sufficient because there are only two possibilities for a node at time  $j-1$  given its state at time  $j$ : either the node is in the same state as time  $j$ , or the node has made a state transition at time  $j$ , and when computing  $\mathcal{L}(X_{j-1}|X_j)$ , the state at time  $j$  is known.

We define a variable  $s_{v,j-1}$  for every node  $v \in E_j$ , and  $r_{v,j-1}$  for every node  $v \in R_j$ . For every node  $v \in I_j$ , we define variable for the incoming state of  $I$  ( $e_{v,j-1}$  for bipartite  $\mathcal{M}$  and  $s_{v,j-1}$  for non-bipartite  $\mathcal{M}$ ). Similarly, for every node  $v \in S_j$ , we define a variable for the incoming state of  $S$  if  $\mathcal{M}$  is loopy, otherwise we do not need to define the variable, since we know that if a node is in  $S_j$  it must be in  $S_{j-1}$ .

With these variables, the likelihoods in (11) are explicitly

defined as in (12)–(15) for SEIRS model as:

$$\mathcal{L}_e^j = \prod_{v \in E_j} \left( \mathcal{L}_{e2e}^{v,j-1} \left( 1 - \mathcal{L}_{s2e}^{v,j,I} \mathcal{L}_{s2e}^{v,j,R} \right)^{s_{v,j-1}} \right) \quad (12)$$

$$\mathcal{L}_s^j = \prod_{v \in S_j} \left( \left( \mathcal{L}_{s2e}^{v,j,I} \mathcal{L}_{s2e}^{v,j,R} \right)^{s_{v,j-1}} (r2s_v)^{1-s_{v,j-1}} \right) \quad (13)$$

$$\mathcal{L}_r^j = \prod_{v \in R_j} \left( (i2r_v)^{1-r_{v,j-1}} (1 - r2s_v)^{r_{v,j-1}} \right) \quad (14)$$

$$\mathcal{L}_i^j = \prod_{v \in I_j} \left( ((e2i_v)^{e_{v,j-1}} (1 - i2r_v)^{1-e_{v,j-1}}) \right) \quad (15)$$

where the sub-terms are defined as:

$$\mathcal{L}_{s2e}^{v,j,I} = \prod_{u \in P(v) \cap I_j} (1 - p_{uv})^{1-e_{u,j-1}}$$

$$\mathcal{L}_{s2e}^{v,j,R} = \prod_{u \in P(v) \cap R_j} (1 - p_{uv})^{1-r_{u,j-1}}$$

$$\mathcal{L}_{e2e}^{v,j} = \prod_{v \in E_j} (1 - e2i_v)^{1-s_{v,j-1}}$$

Each likelihood above for a given state ( $\mathcal{L}_s^j, \mathcal{L}_e^j, \mathcal{L}_i^j, \mathcal{L}_r^j$ ) has two parts: the likelihood of the nodes staying at the given state, and the likelihood of the nodes transitioning towards the given state. For example,  $\mathcal{L}_e^j$  is the likelihood for nodes  $v \in E_{j-1}$  not to transition to infected state at time  $j$ , and nodes  $v \in S_{j-1}$  to become exposed at time  $j$ . This gives us an explicit definition of objective function (8) in terms of a collection of binary variables.

Likelihoods (12)–(15) are defined for the most general model SEIRS, some of the likelihood terms disappear, or change slightly for models that are missing some of the states. For instance, parts including  $r2s_v$  in  $\mathcal{L}_s^j$  and  $\mathcal{L}_r^j$  disappear for non-loopy  $\mathcal{M}$ , the likelihood representing the exogenous transition  $\mathcal{L}_{s2e}$  is replaced by the similarly defined  $\mathcal{L}_{s2i}$  for non-bipartite models, etc.

2) *Expressing the constraints in equations* (9) and (10):

The intra-consistency constraints (9) that require every node have a single state at each time step are already implied by the objective function since there is only a single variable for every node modeling the two possibilities. The inter-consistency constraints (10) can be modeled as packing constraints:

$$\sum_{u \in P(v) \cap I_j} e_{u,j-1} + \sum_{u \in P(v) \cap R_j} r_{u,j-1} + s_{v,j-1} \leq d_v, \forall v \in E_j \quad (16)$$

These constraints make sure every node that became exposed at time  $t$  ( $v \in E_j, S_{j-1}$ ) has at least one incoming edge from node  $u \in I_{j-1}$  ( $u \in I_j \cup R_j$ ). However, these constraints (16) are already represented in the objective function (11) since the higher order term  $\log(\mathcal{L}_e^j)$  takes the lowest possible value  $\log(0) = -\infty$  when any of them are not satisfied.

3) *Optimizing the likelihood under these constraints:*

Since  $t_{start}$  is unknown, we reconstruct the history using the above likelihoods and constraints by iteratively maximizing the likelihood at each time step  $t_{start} \leq j < t_{min}$  backwards starting from  $t_{min} - 1$ , where the state is known. In each

iteration, given  $X_j$ , we reconstruct the states at the previous time step  $j - 1$  ( $X_{j-1}$ ) by maximizing:

$$\max F = \log(\mathcal{L}_s^j) + \log(\mathcal{L}_e^j) + \log(\mathcal{L}_i^j) + \log(\mathcal{L}_r^j) \quad (17)$$

Objective  $F$  for single step reconstruction is submodular as proven in Theorem IV.1 for all SEIRS-type models except SIS. For SIS model, history reconstruction can still be expressed as submodular maximization under packing and partition matroid constraints by modifying  $F$  as in Theorem IV.2. Proofs of all theorems can be found in the extended version of this paper.

**Theorem IV.1.**  $F$  in Equation (17) is non-monotone submodular for all SEIRS-type models except SIS.

**Theorem IV.2.** Program (17)–(16) for SIS can be expressed as submodular maximization under both packing and partition matroid constraints.

Therefore, optimizing (17) is a non-monotone submodular maximization problem. Non-monotone submodular maximization is NP-hard since its special cases such as *MAX DICUT* is NP-hard [13]. To solve this problem, we apply the deterministic non-monotone submodular maximization method by [14] repeatedly between adjacent time steps and iterate until the estimates between the consecutive time steps are same, indicating that we have reached the initial  $t_{start}$  state. At each step, [14] maximizes a normalized  $F_n$  at every step that is obtained by adding  $-F(\emptyset)$  to every  $S \subset 2^N$  so  $F_n(\emptyset) = 0$ , where  $F(\emptyset)$  is the value of the objective if no nodes change states between adjacent time points. As applied here, this is done by starting with an initial solution  $X_j' = \emptyset$  that represents the same state assignments between the consecutive time steps  $j$  and  $j - 1$ . For each time step, we add the node with the most increase in  $F_n$  to the set of nodes that have changed state. Algorithm 1 gives a schematic outline of the procedure.

---

**Algorithm 1** *DHR-sub-early*

---

```

1:  $j \leftarrow t_{min} - 1$ 
2: repeat
3:    $\{X_j'$  is the set of nodes that changed state at time  $j\}$ 
4:    $X_j' \leftarrow \emptyset$ 
5:   repeat
6:     Add nodes to  $X_j'$  according to the rule for
       non-monotone submodular maximization approxima-
       tion [14]
7:   until no node can be added that increases the score
8:    $j \leftarrow j - 1$ 
9: until  $X_j' = X_{j+1}'$ 

```

---

The method by [14] has a  $\frac{1}{3}$  approximation ratio for normalized submodular functions, and we found it to perform in practice better than the randomized algorithm [15] with approximation ratio 0.5 due to the structure of our problem. The  $\frac{1}{3}$  ratio for normalized  $F_n$  implies a data-dependent bound for  $F$  as proven in Theorem IV.3 where  $F(\emptyset) = -S_0$ ,  $X_{opt}$  is the set of elements maximizing  $F$  and  $F(X_{opt}) = -O$ .

**Theorem IV.3.** Algorithm 1 has approximation guarantee of  $k + \frac{S_0}{O}(1 - k)$  for  $k = \frac{1}{3}$  in terms of minimization of supermodular  $-F$  for each of its iteration.

## B. History Reconstruction Between Consecutive Snapshots (DHR-sub-between)

History reconstruction for every interval between by consecutive, observed  $T_D$  pairs is independent of other intervals. Therefore, we can solve each independently by solving the following problem:

$$\begin{aligned} \operatorname{argmax}_X T = & \log(P(X_{j+1}|D_j)) + \log(P(D_k|X_{k-1})) \\ & + \sum_{t \in j+1, \dots, k-2} \log(P(X_{t+1}|X_t)) \end{aligned} \quad (18)$$

$$\text{s.t. } \mathbf{IntraConsistent}(X_t, \mathcal{M}), t \in j + 1, \dots, k - 1 \quad (19)$$

$$\mathbf{InterConsistent}(X_t, X_{t+1}, \mathcal{M}), t \in j, \dots, k - 1 \quad (20)$$

where we only consider  $X_t$  that lay within the  $[j, k]$  interval bracketed by diffusion snapshot observations  $D_j$  and  $D_k$

1) *Expressing objective* (18): Objective (18) has three parts:  $\log(P(D_k|X_{k-1}))$  is same as the single step backwards reconstruction of *DHR-sub-early*, and  $\log(P(X_{j+1}|D_j))$  is a trivial forward diffusion expression with unknown  $X_{j+1}$  and known  $D_j$ . On the other hand, both  $X_{t+1}$  and  $X_t$  are unknown in  $\log(P(X_{t+1}|X_t))$ , and it can be written explicitly as:

$$\log(P(X_{t+1}|X_t)) = \log(\mathcal{L}_s^{t+1}) + \log(\mathcal{L}_e^{t+1}) + \log(\mathcal{L}_i^{t+1}) + \log(\mathcal{L}_r^{t+1}) \quad (21)$$

where the likelihoods are defined as:

$$\mathcal{L}_e^{t+1} = \prod_{v \in V} ((1 - e2i_v)^{e_{v,t} e_{v,t+1}} (1.0 - \mathcal{L}_{exo})^{e_{v,t+1}}) \quad (22)$$

$$\mathcal{L}_s^{t+1} = \prod_{v \in V} (\mathcal{L}_{exo}^{s_{v,t+1}} (r2s_v)^{r_{v,t} s_{v,t+1}}) \quad (23)$$

$$\mathcal{L}_r^{t+1} = \prod_{v \in V} ((i2r_v)^{i_{v,t} r_{v,t+1}} (1.0 - r2s_v)^{r_{v,t} r_{v,t+1}}) \quad (24)$$

$$\mathcal{L}_i^{t+1} = \prod_{v \in V} ((e2i_v)^{e_{v,t} i_{v,t+1}} (1.0 - i2r_v)^{i_{v,t} i_{v,t+1}}) \quad (25)$$

and the term in (22) is:

$$\mathcal{L}_{exo} = \prod_{u \in P(v)} (1 - p_{uv})^{i_{u,t} s_{v,t}}$$

Objective (18) is non-monotone submodular as proven in Theorem IV.4.

**Theorem IV.4.**  $T$  in (18) is non-monotone submodular for all SEIRS-type models.

2) *Expressing the inter- and intra-consistency constraints:* The inter-consistency constraints (20) ensure the validity of diffusion, and they are explicitly written, using binary state variables  $s_{v,t}$ ,  $e_{v,t}$ ,  $i_{v,t}$ ,  $r_{v,t}$  for every  $v \in V$  and  $t \in j, \dots, k$ , as:

$$s_{v,t} + i_{v,t+1} + r_{v,t+1} \leq 1, \quad (26)$$

$$e_{v,t} + s_{v,t+1} + r_{v,t+1} \leq 1, \quad (27)$$

$$i_{v,t} + s_{v,t+1} + e_{v,t+1} \leq 1, \quad (28)$$

$$r_{v,t} + e_{v,t+1} + i_{v,t+1} \leq 1, t \in j, \dots, k - 1, v \in V \quad (29)$$

$$e_{v,t+1} - s_{v,t} \leq \sum_{u \in P(v)} i_{u,t}, t \in j, \dots, k - 1, v \in V \quad (30)$$

Constraints (26)–(29) ensure that state transitions obey SEIRS dynamics rules such as a node infected at  $t$  cannot be susceptible or exposed at  $t+1$ . The remaining constraint (30) ensures that a newly exposed node must have at least one infected predecessor at previous time step. The constraints (30) are already represented in the objective function, since  $(1.0 - \mathcal{L}_{exo})^{e_{v,t+1}}$  in  $\mathcal{L}_e^{t+1}$ , and (21) takes the lowest possible value  $\log(0) = -\infty$  when any of them is not satisfied. So, we can remove the constraints (30) without affecting the results. (Some of these constraints are modified accordingly for subset of models. For example, (28) becomes  $i_{v,t} + s_{v,t+1} \leq 1$  for SI.)

The intra-consistency constraints (19) ensure that every node belongs to a single state at each time step:

$$s_{v,t} + e_{v,t} + i_{v,t} + r_{v,t} = 1, t \in j+1, \dots, k-1, v \in V \quad (31)$$

3) *Optimizing (18) in practice:* Let  $E = \{s_{v,t}, e_{v,t}, i_{v,t}, r_{v,t} | v \in V, t \in j, \dots, k\}$ . Then the intra-consistency constraints (31) define base of a partition matroid over the ground set  $E$  [16], and constraints (26)–(29) define 2-independence system over the same ground set: Its rank quotient is 2 since the ratio of cardinality of the largest base (maximal independent set) to the cardinality of the smallest base is at most 2. For more detailed information on matroid and independence system, see [17], [18].

Combining (21) and (18) with the discussion above, the history reconstruction between time steps  $j$  and  $k$  can be then written as optimizing

$$\max T = \sum_{t=j+1}^k \log(\mathcal{L}_s^t) + \log(\mathcal{L}_e^t) + \log(\mathcal{L}_i^t) + \log(\mathcal{L}_r^t) \quad (32)$$

subject to inter-consistency constraints (26)–(29) and intra-consistency constraints (31). (31) cannot be removed as in Section IV-A since each node may belong to any state at time  $t$  as we do not know the node states at  $t-1$  or  $t+1$  (except for boundary times  $j$  and  $k$ ).

When considered together, constraints (26)–(29) and (31) are base of a new matroid defined by the intersection of the partition matroid and 2-independence system. Proof is as follows: Intersection of constraints (26)–(29) and (31) relaxed to  $\leq$  define a matroid  $\mathcal{M}_p = (E_p, \mathcal{I}_p)$  where  $E_p = E$ , and independent set  $\mathcal{I}_p$  is subset of  $E$  satisfying (26)–(29) and relaxed (31).  $\mathcal{M}_p$  defines a matroid since all its bases (maximal independent set) have the same cardinality  $(k-j-1)|V|$ ; we can always find a state assignment for every node and every time step that satisfies the constraints, and we cannot assign multiple states to each node at each time step. Then, equality constraints in the original equations (31) force independent sets in  $\mathcal{I}_p$  to be bases of  $\mathcal{M}_p$ , as cardinality of an independent set in  $\mathcal{I}_p$  will now always be  $(k-j-1)|V|$ .

This problem becomes non-monotone submodular maximization under matroid base constraints. It is NP-hard [16], and its normalized version can be approximated by  $\frac{1}{6}$  by modified local search [16]. We run this method by [16] in *DHR-sub-between* to reconstruct the history between  $j$  and  $k$ .

*DHR-sub-between* has three main steps: In the first step, it starts with a base of  $\mathcal{M}_p$  that also satisfies (30), and it finds a base  $B_1 \subseteq \mathcal{M}_p$  that is optimal under swap operations. In

the second step, it removes  $B_1$  from  $\mathcal{M}_p$  and greedily finds independent set  $X_2$  that is locally optimal under addition and deletion operations. In the third step, it contracts independent set  $X_2$  from  $\mathcal{M}_p$  and finds two disjoint bases  $B_a, B_b$  that are guaranteed to exist when the original matroid  $\mathcal{M}_p$  has two disjoint bases. Lastly, it returns the best of three bases  $B_1, X_2 \cap B_a$  or  $X_2 \cap B_b$ . The resulting solution always satisfies (30) without explicitly checking for them: Local search will not replace the current solution with a low-score invalid solution as the objective (32) takes the lowest possible value  $-\infty$  if any of (30) are not satisfied.

## V. PRIZE COLLECTING (DOMINATING SET) VERTEX COVER RELAXATIONS (*DHR-pcdsvc*, *DHR-pcvc*)

Although accurate and practical for smaller graphs, *DHR-sub* may take some time to solve for larger graphs. We can reconstruct the history before the earliest observed snapshot faster by relaxing *DHR-sub-early*. For a relaxed version of the problem, we define variables differently than above. When reconstructing the history at previous time  $j-1$ , we define  $i_{v,j-1}, \forall v \in I_j \cup R_j, s_{v,j-1}, \forall v \in S_j$ , and  $e_{v,j-1}, \forall v \in E_j$ . After this transformation,  $\mathcal{L}_e^j$  (12) turns into

$$\begin{aligned} \mathcal{L}_{s2e}^{v,j} &= \prod_{u \in P(v) \cap (I_j \cup R_j)} (1 - p_{uv})^{i_{u,j-1}} \\ \mathcal{L}_{e2e}^{v,j} &= \prod_{v \in E_j} (1 - e2i_v)^{e_{v,j-1}} \\ \mathcal{L}_e^j &= \prod_{v \in E_j} \left( \mathcal{L}_{e2e}^{v,j-1} \left(1 - \mathcal{L}_{s2e}^{v,j}\right)^{1-e_{v,j-1}} \right) \end{aligned} \quad (33)$$

The other likelihoods are transformed similarly.

The hardness of *DHR-sub-early* comes from higher-order terms  $\left(1 - \mathcal{L}_{s2e}^{v,j}\right)^{1-e_{v,j-1}}$  in (33), so we replace them with their first-order Taylor expansion  $\mathcal{T}_{s2e}^{v,j}$  at the point  $(i_{u,j-1} = 1, \forall u \in P(v) \cap \{I_j \cup R_j\} \cup e_{v,j-1} = 1)$  as in (34).

$$\mathcal{T}_{s2e}^{v,j} = \log(K_{s2e}) + \frac{1}{K_{s2e}} \sum_{u \in I_j \cup R_j} \frac{\partial \mathcal{L}_{s2e}^{v,j,(I,R)}}{\partial i_{u,j-1}} (i_{u,j-1} - 1) \quad (34)$$

In (34),  $K_{s2e} = \mathcal{L}_{s2e}^{v,j}(1, \dots, 1) \approx 1$ , so the original reconstruction Problem (17)–(16) for single time step turns into minimizing  $-F_r$  as in (35)–(36):

$$\begin{aligned} \min -F_r &= \sum_{(u,v) \in E^*} w_{uv} \bar{i}_{u,j-1} \bar{e}_{v,j-1} + \sum_{u \in I_j \cup R_j} w_u i_{u,j-1} + \\ &\quad \sum_{v \in E_j} w_v e_{v,j-1} + \sum_{v \in S_j} w_v s_{v,j-1} \quad (35) \\ \text{s. t.} &\quad \sum_{u \in P(v) \cap \{I_j \cup R_j\}} i_{u,j-1} + e_{v,j-1} \geq 1, v \in E_j \quad (36) \end{aligned}$$

where  $\bar{i}_{u,j} = 1 - i_{u,j}$  and  $\bar{e}_{v,j} = 1 - e_{v,j}$ . The covering constraints (36) are inter-consistency constraints ensuring the validity of the diffusion. Similar to *DHR-sub-early*, we do not need intra-consistency constraints since we are reconstructing the history step by step. This problem is *Prize Collecting Dominating Set Vertex Cover* (PCDSVC) over the graph



$G^* = (V^*, E^*)$  where  $V^* = V$  with weights  $w_v$  and directed edge from  $u$  to  $v$  with weight  $w_{uv} = -\log(1 - p_{uv})$  exists when  $v \in E_j, u \in P(v) \cap \{I_j \cup R_j\}$  for bipartite  $\mathcal{M}$  and  $v \in I_j, u \in P(v) \cap \{I_j \cup R_j\}$  for non-bipartite  $\mathcal{M}$ .

PCDSVC is different than *Vertex Cover* because (1) We may not cover an edge  $(u, v)$  if we pay its price  $w_{uv}$ , and (2) A feasible solution is a vertex dominating set. This problem has not been studied before, it is NP-hard, and it can be approximated by  $O(\log(|V^*|))$  by formulating it as *Minimum Hitting Set* and running the greedy method for *Set Cover* as proven in Theorem V.1.

**Theorem V.1.** *Prize Collecting Dominating Set Vertex Cover (PCDSVC) is NP-hard, and it can be approximated by  $O(\log(|V^*|))$ .*

We can relax this problem further by removing (36) and it becomes *Prize Collecting Vertex Cover (PCVC)*. PCVC can be approximated by a factor of 2 using the LP relaxation [19], and it can be solved optimally for *bipartite diffusion models* by expressing it as s-t mincut as proven in Theorem V.2.

**Theorem V.2.** *The Taylor expansion relaxation of (17) for bipartite diffusion models can be expressed as s-t mincut.*

The algorithms for these relaxed versions, *DHR-pcdsvc* and *DHR-pcvc*, are similar to *DHR-sub* except they run PCDSVC and PCVC respectively instead of submodular maximization at each iteration.

## VI. ENSEMBLE INITIAL SPREADER IDENTIFICATION

We define *DHR-sub-ens*, *DHR-pcdsvc-ens* and *DHR-pcvc-ens* for the ensemble versions of our methods: they estimate the most likely subset of nodes that explains the diffusion data from multiple runs. For each initial time point seen in the multiple runs, we greedily select the subset of nodes seen in that time point that best explains  $\mathcal{D}$  in terms of minimum absolute difference  $F_{dif}$  from Equation 37:

$$F_{dif} = |S_t^e - S_t^t| + |E_t^e - E_t^t| + |I_t^e - I_t^t| + |R_t^e - R_t^t| \quad (37)$$

where  $S_t^e, E_t^e, I_t^e, R_t^e$  are the set of estimated nodes whereas  $S_t^t, E_t^t, I_t^t, R_t^t$  are the set of true nodes for  $S, E, I, R$  states at time  $t$  respectively. We keep adding the node that improves  $F_{dif}$  the most until there is no improvement. Lastly, we return set of nodes that has the minimum score among the all possible initial time points as our initial spreader prediction.

## VII. EXPERIMENTAL RESULTS

### A. Comparison and Evaluation

We compared our methods with *NetSleuth*, *Keffectors* and *Rumor* in identifying the initial spreaders. *Keffectors* and *Rumor* require estimates of the number of initial spreaders, so we provide them an estimate of the initial spreader count by the number of clusters in  $G$  estimated by modularity [20]. We compared our methods with the baseline heuristic *Greedy-Forward* for history reconstruction that reconstructs the history in each interval by simulating a forward trace starting from the interval's earlier time. We return the topmost  $k$  spreaders from *Rumor* sorted by its rumor centrality metric where  $k$  is the number of clusters in  $G$ .

We validated the history reconstruction performance by Kendall Tau-b statistic [21] ( $\tau_B$ ) that measures the similarity between true and estimated node orderings defined in terms of infection times by also adjusting for ties:

$$\tau_B(T, O) = \frac{n_c - n_d}{\sqrt{(n_0 - n_1)(n_0 - n_2)}} \quad (38)$$

Here,  $T$  and  $O$  are true and inferred node orderings respectively in terms of given state (such as *infected*). Let  $V_T$  be set of nodes seen in true ordering  $T$ , then  $\tau_B$ ,  $n_c$  are  $n_d$  are concordant and discordant pairs respectively,  $n_0 = \frac{|V_T|(|V_T|-1)}{2}$ ,  $n_1$  and  $n_2$  are sum of tied quantities in the true and observed orderings respectively. Kendall tau-b adjusts for ties by subtracting  $n_1$  and  $n_2$  from  $n_0$  in the denominator.

We validated the initial spreaders identification performance by graph-based average matching score ( $\bar{M}_G$ ). Let  $\hat{V}_t$  and  $\hat{V}_o$  be true and estimated initial nodes respectively, and  $G_b = (\hat{V}_t \cup \hat{V}_o, \hat{V}_t \times \hat{V}_o)$  be a weighted bipartite graph with weights  $w_{ab} = \frac{1}{1+d_{ab}}$  for every  $a \in \hat{V}_t, b \in \hat{V}_o$  where  $d_{ab}$  is the distance between  $a$  and  $b$  in  $G$ .  $\bar{M}_G$  estimates the maximum bipartite matching score in  $G_b$ , and returns the average. When  $|\hat{V}_o| \neq |\hat{V}_t|$ ,  $M_G$  is modified to account for the unmatched vertices by matching them independently to the best ones.

Both  $\tau_B$  and  $\bar{M}_G$  are normalized, and higher score means better performance in both. We implemented all our methods, synthetic trace generator and existing methods *NetSleuth*, *Rumor*, *Keffectors* in Python, solved LP relaxations by CPLEX [22], and modified and used C++ maximum flow code from [23]. We run all our experiments on Macbook Pro with 2.5 Ghz CPU and 8 GB memory. All our code, data and supplementary text including proofs of the theorems are available on the web<sup>1</sup>.

### B. Reconstruction Performance on Synthetic Data

We generated 5 networks of 500 nodes and 5000 edges that are grown by *Erdős-Reyni* [24], *Forest Fire* [25], *linear preferential attachment* [26] network growth models. We generated each synthetic trace by choosing the given number of source nodes randomly, making them infected, and running the diffusion over the network until either all nodes become recovered (or infected under the SI model) or until the spread dies out. When multiple snapshots are given, we sample them uniformly in the range  $(t_{min}, t_{max})$ .

We test our methods on SI, SIR, SEIR by modeling the transition distributions from a geometric distribution with different parameters in each model in order to assess performance under various conditions. In SI, we selected  $p_{uv}$  for every  $(u, v) \in E$  uniformly between 0.1 and 0.4. In SIR, we selected  $p_{uv}, \forall (u, v) \in E$  uniformly in the range (0.2, 0.6) and  $i2r_v, \forall v \in V$  uniformly in the range (0.5, 0.6). In SEIR, we selected  $p_{uv}, \forall (u, v) \in E, e2i_v, \forall v \in V, i2r_v, \forall v \in V$  each uniformly in the range (0.4, 0.8).

*DHR-sub* and its ensemble version *DHR-sub-ens* perform the best on all the models in terms of identifying the initial spreaders as in Table II. Its relaxations *DHR-pcdsvc* and *DHR-pcvc* also perform better than the existing methods, and they

<sup>1</sup><http://www.cs.cmu.edu/~ckingsf/software/dhrec>

	Initial Spreader			History			
	FF			LPA		RDS	
	SI	SIR	SEIR	SI	SIR	SI	SIR
<i>DHR-sub</i>	0.8	0.83	0.81	0.97	0.88	0.69	0.77
<i>DHR-sub-ens</i>	0.87	0.88	0.89	-	-	-	-
<i>DHR-pcdsvc</i>	0.78	0.8	0.81	0.9	0.82	0.64	0.73
<i>DHR-pcvc</i>	0.76	0.76	0.79	0.88	0.77	0.59	0.72
<i>Rumor</i>	0.74	0.7	0.6	-	-	-	-
<i>NetSleuth</i>	0.75	0.8	0.64	-	-	-	-
<i>Keffectors</i>	0.77	0.74	0.7	-	-	-	-
<i>GreedyForward</i>	-	-	-	0.34	0.28	0.31	0.23

TABLE II.  $\bar{M}_G, \tau_B$  vs. GROWTH AND DIFFUSION MODELS FOR SPREADER IDENTIFICATION (5 TRUE SPREADERS) AND HISTORY RECONSTRUCTION FROM  $|T_D| = 2$  SNAPSHOTS

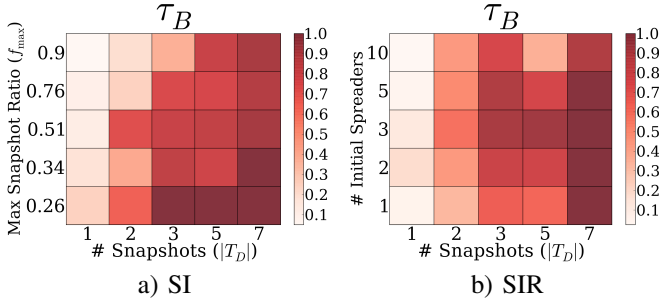


Fig. 3.  $\tau_B$  vs. number of snapshots ( $x$  axis) and max snapshot ratio ( $y$ -axis), number of true initial spreaders ( $y$ -axis) for history reconstruction over *Forest Fire* for *DHR-sub* a) SI b) SIR

are good alternatives to *DHR-sub* considering their faster running times. The performance difference between our methods and the existing methods become more apparent especially for SIR and SEIR models.

In terms of history reconstruction, all our methods perform much better than the greedy baseline *GreedyForward*. All our methods perform better when multiple snapshots are available as seen in Figure 3 for *DHR-sub* for both SI and SIR. *DHR-sub* reconstructs the histories more precisely when the interval to be reconstructed has lower maximum snapshot ratio ( $f_{max} = t_{max}/l_D$ ) where  $l_D$  is the diffusion length, and its performance is not significantly affected by the number of initial spreaders given the same number of snapshots as in Figure 3. Lower reconstruction performance for higher  $f_{max}$  intervals is due to increasing number of similar quality diffusion histories. In its extreme,  $\tau_B$  may become close to 0 when reconstructing histories of longer intervals from a single snapshot.

### C. Reconstructing Meme Diffusion History From Blog Data

We used our methods to extract the diffusion history of memes that are defined as short textual phrases that travel through the Web. We inferred the diffusion progression of several memes in two blog networks under SI using the true diffusion data from [11]: *Top-Blog* has 5000 nodes and 30072 edges and it shows the connection between the topmost 5000 blogs, *Rand-Blog* has 250 nodes and 3342 edges, and it shows the connection between random 250 blogs. In both networks, nodes represent either personal blogs or mass media, and edges represent hyperlinks from one blog to the another one. We do not know the true  $p_{uv}$ , so we estimate them by a geometric distribution with  $p$  being 0.3 between mass media, 0.25 from mass media to bloggers, 0.15 between bloggers, and 0.05 from bloggers to media. Traces for several topics were obtained from

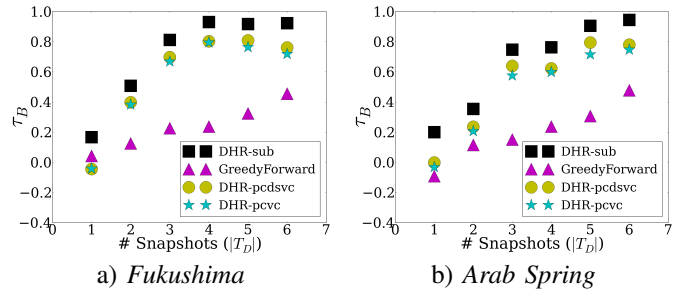


Fig. 4.  $\tau_B$  vs. number of snapshots for a) *Fukushima* (1 : 5), b) *Arab Spring* (1 : 5) on *Top-Blog*

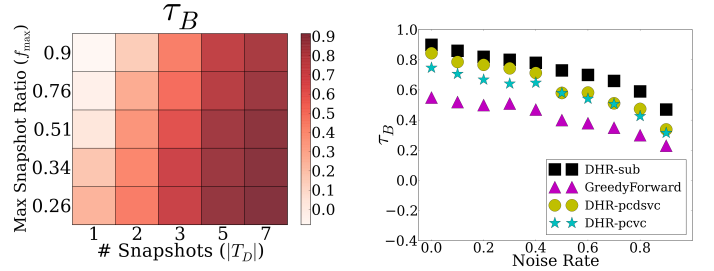


Fig. 5.  $\tau_B$  vs.  $|T_D|$  and  $f_{max}$  for *DHR-sub* of *Nba* (1 : 10) on *Rand-Blog*

Fig. 6.  $\tau_B$  vs noise ratio ( $p$ ) over *Water-sm*

the same source [11]. When tracking the diffusion of a topic, if a blog publishes about it at multiple time points, we assume blog is infected at the earliest time point.

We reconstructed the diffusion history of the memes *Fukushima*, *Arab Spring* and *Nba* on both *Top-Blog* and *Rand-blog* as in Figures (4)–(5). Values inside the parentheses define the time scale for the meme progression (1 : 5 = 1 time unit for 5 days).  $\tau_B$  are lower than the synthetic case especially when fewer than 2 snapshots are available but they are still reasonable since the true diffusion parameters are unknown. In Figure 4, *DHR-sub* performs the best, and all methods reconstruct the diffusion history better when more diffusion data is available. When run with multiple snapshots, *DHR-sub* better captures the diffusion direction and performs almost close to 1 whereas heuristic method *GreedyForward*'s  $\tau_B$  never exceeds 0.5. Although *Fukushima* and *Arab Spring* have different diffusion dynamics [11], both trajectories can be reconstructed precisely by *DHR-sub*. Similar to the synthetic case, performance of *DHR-sub* increases if more snapshots are available, and it decreases as  $f_{max}$  increases as in Figure 5. Overall, both *DHR-sub* and *DHR-pcdsvc* can nicely fill in the missing gaps of the meme diffusion history.

In another example, the order of diffusion estimated by *DHR-sub* matches the true order of the meme *Occupy* reasonably well ( $\tau_B = 0.77$ ). In this case, most of the initial diffusion of *Occupy* happens between mass media, and diffusion at personal blogs start to show up later. However, the speed of the predicted diffusion trajectory is more uniform than true *Occupy* trajectory.



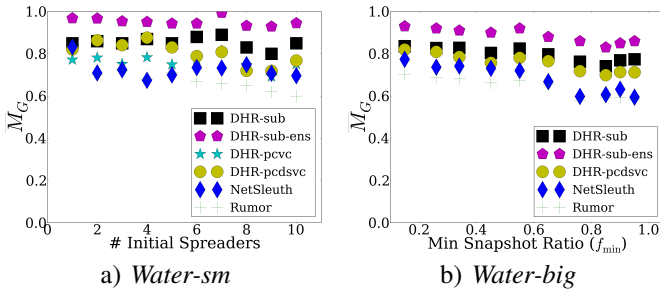


Fig. 7. a)  $\bar{M}_G$  vs. number of initial spreaders for *Water-sm*, b)  $\bar{M}_G$  vs.  $f_{min}$  for *Water-big* (5 initial sites)

#### D. Identifying Initial Water Contamination Sites

We inferred the initial contaminant locations over two water distribution networks [27] where nodes are water demand-supply locations, and the edges represent the water pipes: *Water-sm* has 130 nodes and 173 edges, *Water-big* has 12527 nodes and 14595 edges. We used contaminant diffusion data generated by the water distribution simulator EPANET [28].

We identified the initial contamination sites in *Water-sm* and *Water-big* under SIR where the *recovered* state models the dilution of the contaminant. We approximate the true hydraulic water diffusion dynamics by SIR as follows: we assume that  $p_{uv} = K_1/l_{uv}$  and  $i2r_{uv} = K_2/l_{uv}$  where  $K_1, K_2$  are constants, and  $l_{uv}$  is the length of pipe  $(u, v)$ . Ensemble methods perform the best as in Figure 7 on *Water-sm*, and *DHR-sub* (without ensemble) also performs better than the existing methods. Our methods' performance is consistent across different numbers of initial contamination sites whereas the existing methods' performance is affected by the number of initial sites. Our methods are nonparametric as they do not require number of initial spreaders as input, and our methods' performance consistency makes them the topmost candidates for application domains with multiple but unknown number of initial spreaders.

Performance of both *DHR-sub-ens* and *DHR-pcdsvc* decreases for higher  $f_{min}$  as in Figure 7 on *Water-big*, but they still perform at least 10% better than the best performing *NetSleuth*. This lower performance is due to both difficulty of differentiating between the initial spreader candidates with similar scores, and the decreasing ability to estimate the correct number of initial spreaders. Our methods may miss the true initial spreaders, but their estimates are within close distance to the original spreaders as reflected by higher performance in various cases.

#### E. Predicting temporal diffusion features

We may answer questions related to temporal diffusion features from the reconstructed histories such as How quickly did it spread over time?, Did it spread faster at the beginning slowing down at later time steps?, etc. Here, we compared the speed (first-order) and acceleration (second-order) dynamics of *Unemployment* and *Fukushima* estimated from *DHR-sub* with the true ones from [11] as in Figures (8)–(9). We define speed of a meme as the number of blogs that publishes about the meme for the first time per time unit, and acceleration as the diffusion speed change per time unit.

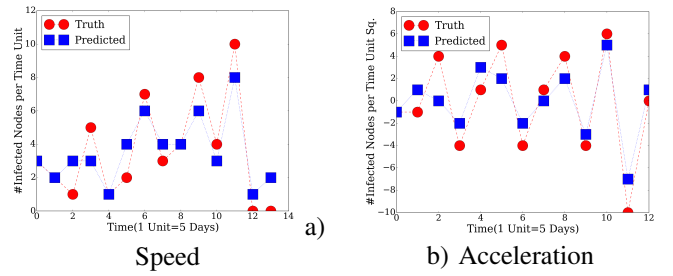


Fig. 8. a) Speed, b) Acceleration Dynamics of True and Predicted Diffusion of *Unemployment* over Time from 3 snapshots

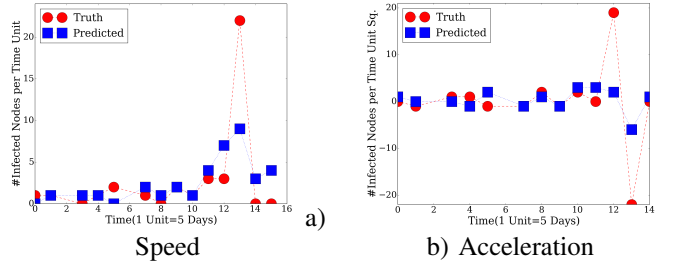


Fig. 9. a) Speed, b) Acceleration Dynamics of True and Predicted Diffusion of *Fukushima* over Time from 3 snapshots

*Unemployment* is a more commonly-used meme than *Fukushima*, and such difference is reflected in their diffusion dynamics: *Unemployment*'s diffusion speed is more uniform over time whereas *Fukushima* shows more bursty dynamics. Diffusion speed of *Unemployment* has multiple local optima for the time points it peaks in news cycle whereas the speed of *Fukushima* has a single peak when it takes attention of the main media sites. However, such difference in diffusion dynamics does not make a difficulty for *DHR-sub* as *DHR-sub* predicted speed of both memes closely approximate their true ones even from 3 snapshots.

*DHR-sub* reconstructed histories of both memes also mimic closely their true acceleration dynamics. The change of diffusion speed for *Unemployment* is more uniform than the one for *Fukushima*, and its uniform dynamics are predicted almost perfectly by *DHR-sub* whereas the main peak of *Fukushima*'s acceleration dynamics was missed by *DHR-sub* except precise approximation at remaining time points. Overall, *DHR-sub* reconstructed histories from only 3 snapshots mimic the true speed and acceleration dynamics of both memes quite precisely even though the original prediction scores are below 0.75 ( $\tau_B = 0.74$  for *Unemployment*,  $\tau_B = 0.65$  for *Fukushima*).

#### F. Scalability and Robustness of History Reconstruction

All our methods reconstruct the history on *Top-Blog* in less than 2 minutes, and our relaxation methods *DHR-pcdsvc*, *DHR-pcvc* reconstruct the history in less than 10 minutes on a large 2D grid graph having 90000 nodes and 179400 edges, with reasonable performance ( $\tau_B = 0.71, 0.63$ ) as in Table III whereas *DHR-sub* takes more than an hour on a personal laptop. When combined with previous sections' results, running times in Table III suggest that *DHR-pcdsvc* and *DHR-pcvc* are nice alternatives to *DHR-sub* for scalable history reconstruction on large graphs. However, we still need faster methods for scalable reconstruction on million-node graphs.

	Top-Blog		2D-GRID	
	$ T_D  = 1$	$ T_D  = 3$	$ T_D  = 1$	$ T_D  = 3$
DHR-sub	112.5	48.9	-	-
DHR-pcdsvc	53.9	28.2	592.1	199.2
DHR-pcvc	49.9	14.3	351.7	82.7

TABLE III. HISTORY RECONSTRUCTION TIME (IN SECONDS) FOR Top-Blog AND A 2D GRID GRAPH FOR DIFFERENT NUMBERS OF DIFFUSION SNAPSHOTS.

Figure 6 shows the performance of contaminant diffusion history reconstruction over *Water-sm* under increasing noise levels. Let  $p$  be the noise ratio between 0.0 and 1.0, and we added the synthetic noise  $p$  as follows: For each node and each state, we randomly select a value  $m$  between 0 and  $pl_D$  where  $l_D$  is length of the diffusion and flip a coin to either add  $m$  to the current state transition time  $t_v$ , or subtract it from  $t_v$ . If modified transition time ( $t_v + m$ ) is less than 0, we make it 0.

Our methods do not show a sudden performance drop by increasing noise levels, as *DHR-sub* can still reconstruct histories with performance over  $\tau_B = 0.7$  even when the noise levels are 0.5. Similarly, *DHR-sub-ens* achieves  $\bar{M}_G = 0.72$  in identifying the initial contaminant locations over *Water-sm* for  $p = 0.5$  (results are not shown). In general, all our methods are robust to the noise in the diffusion data.

## VIII. CONCLUSIONS

We designed several methods for estimating diffusion histories that either optimize the likelihood or its relaxations with provable performance guarantees for local steps. Our methods do not require the number of initial spreaders and diffusion length as parameters. They identify the initial spreaders better than the existing methods specially designed for this task. They reconstruct the history accurately in a number of scenarios. We also accurately estimated temporal diffusion characteristics of several semantically different memes from partial data. These findings suggest the reconstructability of diffusion history from partial data under several settings. Partial diffusion data is not an unsolvable bottleneck as missing diffusion history can be completed by our methods accurately.

**Acknowledgements.** This work has been partially funded by the US National Science Foundation (CCF-1256087, CCF-1053918) and US National Institutes of Health (R21HG006913 and R01HG007104). C.K. received support as an Alfred P. Sloan Research Fellow.

## REFERENCES

- [1] G. Serazzi and S. Zanero, "Computer virus propagation models," in *In Tutorials of the 11th IEEE/ACM International Symposium on Modeling, Analysis and Simulation of Computer and Telecommunications Systems (MASCOTS03)*. Springer-Verlag, 2003.
- [2] J. Yang and J. Leskovec, "Patterns of temporal variation in online media," in *ACM International Conference on Web Search and Data Mining (WSDM)*, 2011, pp. 177–186.
- [3] B. A. Prakash, D. Chakrabarti, M. Faloutsos, N. Valler, and C. Faloutsos, "Threshold conditions for arbitrary cascade models on arbitrary networks," in *Proceedings of the 2011 IEEE 11th International Conference on Data Mining*. Washington, DC, USA: IEEE Computer Society, 2011, pp. 537–546.
- [4] M. Salathé, M. Kazandjieva, J. W. Lee, P. Levis, M. W. Feldman, and J. H. Jones, "A high-resolution human contact network for infectious disease transmission." *Proc. Natl. Acad. Sci. USA*, vol. 107, no. 51, pp. 22 020–5, 2010.

- [5] D. Kempe, J. Kleinberg, and E. Tardos, "Maximizing the spread of influence through a social network," in *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. New York, NY, USA: ACM, 2003, pp. 137–146.
- [6] J. Leskovec, L. A. Adamic, and B. A. Huberman, "The dynamics of viral marketing," *ACM Trans. on the Web*, vol. 1, no. 1, p. 5, 2007.
- [7] T. Lappas, E. Terzi, D. Gunopulos, and H. Mannila, "Finding effectors in social networks," in *KDD '10*. New York, New York, USA: ACM Press, Jul. 2010, p. 1059.
- [8] D. Shah and T. Zaman, "Finding Rumor Sources on Random Graphs," *arXiv*, p. 1110.6230, 2011.
- [9] B. A. Prakash, J. Vreeken, and C. Faloutsos, "Spotting culprits in epidemics: How many and which ones?" in *ICDM*, 2012, pp. 11–20.
- [10] M. Gomez-Rodriguez, J. Leskovec, and A. Krause, "Inferring Networks of Diffusion and Influence," *ACM Transactions on Knowledge Discovery from Data*, vol. 5, no. 4, pp. 1–37, Feb. 2012.
- [11] M. Gomez-Rodriguez, J. Leskovec, and B. Schölkopf, "Structure and dynamics of information pathways in online media," in *WSDM '13*. New York, NY, USA: ACM, 2013, pp. 23–32.
- [12] H. W. Hethcote, "The mathematics of infectious diseases," *SIAM Rev.*, vol. 42, no. 4, pp. 599–653, 2000.
- [13] U. Feige and M. Goemans, "Approximating the value of two power proof systems, with applications to max 2sat and max dicut," in *Theory of Computing and Systems, 1995. Proceedings., Third Israel Symposium on the*, Jan 1995, pp. 182–189.
- [14] U. Feige, V. S. Mirrokni, and J. Vondrak, "Maximizing non-monotone submodular functions," in *Proceedings of the 48th Annual IEEE Symposium on Foundations of Computer Science*, ser. FOCS '07. Washington, DC, USA: IEEE Computer Society, 2007, pp. 461–471.
- [15] N. Buchbinder, M. Feldman, J. S. Naor, and R. Schwartz, "A Tight Linear Time (1/2)-Approximation for Unconstrained Submodular Maximization," in *2012 IEEE 53rd Annual Symposium on Foundations of Computer Science*. IEEE, Oct. 2012, pp. 649–658.
- [16] J. Lee, V. S. Mirrokni, V. Nagarajan, and M. Sviridenko, "Non-monotone submodular maximization under matroid and knapsack constraints," in *Proceedings of the Forty-first Annual ACM Symposium on Theory of Computing*. New York, NY, USA: ACM, 2009, pp. 323–332.
- [17] A. Schrijver, *Combinatorial Optimization - Polyhedra and Efficiency*. Springer, 2003.
- [18] A. Gupta, A. Roth, G. Schoenebeck, and K. Talwar, "Constrained non-monotone submodular maximization: Offline and secretary algorithms," *CoRR*, vol. abs/1003.1517, 2010.
- [19] D. S. Hochbaum, "Instant recognition of polynomial time solvability, half integrality and 2-approximations," in *APPROX '00*. Springer-Verlag, 2000, pp. 2–14.
- [20] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," *J. of Statistical Mechanics: Theory and Experiment*, vol. 2008, no. 10, pp. P10008+, 2008.
- [21] A. Agresti, *Categorical Data Analysis*, 2nd ed., ser. Wiley Series in Probability and Statistics. Wiley-Interscience, 2002.
- [22] "IBM ILOG CPLEX Optimizer," 2010. [Online]. Available: <http://www.ilog.com/products/cplex/>
- [23] Y. Boykov and V. Kolmogorov, "An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 26, no. 9, pp. 1124–1137, 2004.
- [24] P. Erdős and A. Rnyi, "On the evolution of random graphs," in *Publication of the Mathematical Institute of the Hungarian Academy of Sciences*, 1960, pp. 17–61.
- [25] J. Leskovec, J. Kleinberg, and C. Faloutsos, "Graphs over time: densification laws, shrinking diameters and possible explanations," in *KDD '05*. New York, NY, USA: ACM, 2005, pp. 177–187.
- [26] A.-L. Barabási and R. Albert, "Emergence of scaling in random networks," *Science*, vol. 286, no. 5439, pp. 509–512, Oct. 1999.
- [27] Avi Ostfeld et al., "The battle of water sensor networks (bwsn): A design challenge for engineers and algorithms," *Journal of Water Resources Planning and Management*, vol. 134, no. 6, pp. 556–568, 2008.
- [28] L. Rossman, "The epanet programmer's toolkit for analysis of water distribution systems," in *WRPMD '99*, 1999, pp. 1–10.

## IX. APPENDIX

**Theorem IV.1.** *Log-likelihood  $F$  in Equation (17) is non-monotone submodular for all SEIRS-type models except SIS.*

*Proof:*

$F$  has three types of terms; higher order terms from  $\log(\mathcal{L}_e^j)$ , quadratic or linear terms from  $\log(\mathcal{L}_s^j)$  depending on  $\mathcal{M}$  and linear terms from  $\log(\mathcal{L}_i^j)$  and  $\log(\mathcal{L}_r^j)$ .  $F$  is non-monotone since linear and quadratic terms are either positive or negative depending on  $\mathcal{M}$ , transition distribution parameters and the terms from  $\log(\mathcal{L}_s^j)$  that model the probability of susceptible nodes not being infected/exposed.

$F$  is submodular when  $F(A+x) - F(A) \geq F(B+x) - F(B)$  for every  $A \subset B$  and for every  $x \in U \setminus (A \cup B)$ . To prove submodularity of  $F$ , we prove the submodularity of each term in  $F$  since summation of submodular functions is also submodular. Linear terms of  $F$  are unimodular, so they are submodular. Quadratic terms show up in  $\log(\mathcal{L}_s^j)$  when  $\mathcal{M}$  is loopy and when the model is not SIS, each quadratic term is one of the following:  $Q(r_{v,j-1}, r_{u,j-1}) = \log(1 - p_{uv})(1 - r_{v,j-1})(1 - r_{u,j-1})$ ,  $Q(r_{v,j-1}, e_{u,j-1}) = \log(1 - p_{uv})(1 - r_{v,j-1})(1 - e_{u,j-1})$  or  $Q(r_{v,j-1}, i_{u,j-1}) = \log(1 - p_{uv})(1 - r_{v,j-1})(1 - i_{u,j-1})$ . All those terms are submodular since they satisfy the inequality  $Q(0,0) + Q(1,1) \leq Q(0,1) + Q(1,0)$ .

Then, we need to prove the submodularity of the higher-order terms that depend on  $G$  to prove submodularity of  $F$ . Higher-order terms appear in either  $\log(\mathcal{L}_e^j)$  for bipartite models or  $\log(\mathcal{L}_i^j)$  for non-bipartite diffusion models. Depending on  $\mathcal{M}$ , we need to prove either  $T = s_{v,j-1} \log\left(1.0 - \mathcal{L}_{s2e}^{v,j,I} \mathcal{L}_{s2e}^{v,j,R}\right)$  or  $T = s_{v,j-1} \log\left(1.0 - \mathcal{L}_{s2i}^{v,j,I} \mathcal{L}_{s2i}^{v,j,R}\right)$ . Each variable might appear at two positions of  $T$ ; either inside or outside the logarithm. When  $\mathcal{M}$  is bipartite, each variable can only appear in one of those positions whereas it can appear in both positions for non-bipartite  $\mathcal{M}$ . Let  $V_e = \bigcup_{u \in P(v) \cap I_j} e_{u,j-1}$ ,  $V_r = \bigcup_{u \in P(v) \cap E_j} r_{u,j-1}$ ,  $x$  be the variable to be added,  $X$  be the current set of added variables,  $K = \prod_{V_e \cup V_r} (1 - p_{uv})$  and  $P_t = (1 - p_{tv})^t$  for every  $t \in V_e \cup V_r$ ,  $T$  is submodular as proven below.

- If  $x$  is outside the logarithm, let  $A = \{a, b\}$  and  $B = \{a, b, c\}$ . Then,  $T(A+x) = \log\left(1 - \frac{K}{P_a P_b P_x}\right)$ ,  $T(B+x) = \log\left(1 - \frac{K}{P_a P_b P_c P_x}\right)$  and  $T(A+x) - T(A) \geq T(B+x) - T(B)$  will be satisfied since  $T(A+x) \geq T(B+x)$  and  $T(A) = T(B) = 0$ .
- If  $x$  is inside the logarithm, when  $s_{v,j-1} \notin X$ , submodularity is trivially satisfied since  $T(A) = T(A+x) = T(B) = T(B+x) = 0$ . When  $s_{v,j-1} \in X$ , let  $A = \{a\}$  and  $B = \{a, c\}$  ( $A \subset B$ ), submodularity is satisfied as shown in Equation (39)–(41).

$$T(A+x) - T(A) \geq T(B+x) - T(B) \quad (39)$$

$$\log\left(\frac{1 - \frac{K}{P_a P_x}}{1 - \frac{K}{P_a}}\right) \geq \log\left(\frac{1 - \frac{K}{P_a P_b P_x}}{1 - \frac{K}{P_a P_b}}\right) \quad (40)$$

$$K P_a P_b (1 - P_b)(1 - P_a) \geq 0 \quad (41)$$

Then,  $F$  is submodular since each summation term including the higher-order ones is submodular. ■

**Theorem IV.2.** *History reconstruction from log-likelihood for SIS model can be expressed as submodular maximization under both packing and partition matroid constraints.*

*Proof:*

Quadratic terms  $Q(s_{v,j-1}, s_{u,j-1})$  from  $\mathcal{L}_s^j$  are supermodular for SIS but they can be turned into submodular ones as follows: We define new variable  $i_{v,j-1}$  for every node  $v \in \{S_j \cup I_j\}$  to represent whether  $v$  is infected at time  $j-1$ . Then, we obtain the new objective function  $F^*$  by replacing each supermodular  $Q(s_{v,j-1}, s_{u,j-1}) = \log(1 - p_{uv})s_{v,j-1}(1 - s_{u,j-1})$  with  $Q^*(s_{v,j-1}, s_{u,j-1}) = \log(1 - p_{uv})s_{v,j-1}i_{v,j-1}$ . We also add assignment constraints of  $s_{v,j-1} + i_{v,j-1} = 1$  for every node  $v \in \{S_j \cup I_j\}$  to make sure node  $v$  is either infected or susceptible at  $j-1$ . Each  $Q^*(s_{v,j-1}, s_{u,j-1})$  in  $F^*$  is submodular since it satisfies the inequality  $Q^*(0,0) + Q^*(1,1) \leq Q^*(0,1) + Q^*(1,0)$ . Then,  $F^*$  is submodular since the rest of the higher-order terms are submodular as proven in Theorem IV.1. Assignment constraints define partition matroid and the problem of reconstructing history at time  $j-1$  becomes submodular maximization under both partition matroid and existing packing constraints for SIS model. ■

**Theorem IV.3.** *Algorithm 1 has approximation guarantee of  $k + \frac{S_0}{O}(1 - k)$  for  $k = \frac{1}{3}$  in terms of minimization of supermodular  $-F$  for each of its iteration.*

*Proof:*

Let  $X$  be the set of elements returned by the non-monotone submodular maximization algorithm and  $F(X) = -M$ . We are interested in upper-bounding the supermodular minimization ratio ( $\frac{M}{O}$ ) for  $-F$ . Since  $F_n$  is obtained by adding  $S_0$  to each set in  $F$ ,  $\frac{F_n(X)}{F_n(X_{opt})} = \frac{S_0 - M}{S_0 - O} \geq k$  and we obtain  $\frac{M}{O} \leq k + \frac{S_0}{O}(1 - k)$ . Here,  $\frac{S_0}{O}(1 - k)$  makes the approximation ratio data-dependent and this ratio is the best we can achieve when  $k$  is tight for non-monotone submodular maximization. This data-dependent bound is also the best we can achieve in terms of supermodular minimization perspective since non-negative supermodular minimization problem cannot be approximated in constant factor unless  $P = NP$  [1]. ■

**Theorem IV.4.**  *$\log(\mathcal{L}_{j,k}^{in})$  in Equation 18 is non-monotone submodular for all SEIRS-type models.*

*Proof:*

We prove the submodularity of  $\log(\mathcal{L}_{j,k}^{in})$  by proving the submodularity of each of its summation terms.  $\log(P(X_{j+1}|D_j))$  estimates the most probable diffusion snapshot at  $j+1$  given  $D_j$ . It is a forward estimate and if we use the same variable naming as in Section IV-A, it becomes a linear function of  $X_{j+1}$  and thus submodular.

$\log(P(D_k|X_{k-1}))$  is same as  $F$  (17) in Section IV-A and it is submodular as proven in Theorem IV.1.

Every  $\log(P(X_{t+1}|X_t))$  involves the variables from both time steps  $t$  and  $t+1$ . Here, we do not know the exact node states at both time steps so we define all possible state variables for every node for both time steps  $(s_{v,t}, e_{v,t}, i_{v,t}, r_{v,t}, s_{v,t+1}, e_{v,t+1}, i_{v,t+1}, r_{v,t+1}, \forall v \in V)$ .  $\log(P(X_{t+1}|X_t))$  can be expressed as in Equation 42 where the likelihoods are defined as in Equation (43)–(46). Each term in  $\log(P(X_{t+1}|X_t))$  is additive and log-likelihood terms of endogenous transitions are submodular since they are quadratic terms with negative coefficient. Log-Likelihood terms of exogenous transitions are also submodular by following the submodularity proof of the higher-order terms from Theorem IV.1.

$$\log(P(X_{t+1}|X_t)) = \log(\mathcal{L}_s^{t+1}) + \log(\mathcal{L}_e^{t+1}) + \log(\mathcal{L}_i^{t+1}) + \log(\mathcal{L}_r^{t+1}) \quad (42)$$

$$\mathcal{L}_{exo} = \prod_{u \in P(v)} (1 - p_{uv})^{i_{u,t} s_{v,t}}$$

$$\mathcal{L}_e^{t+1} = \prod_{v \in V} ((1 - e_{2i_v})^{e_{v,t} e_{v,t+1}} (1.0 - \mathcal{L}_{exo})^{e_{v,t+1}}) \quad (43)$$

$$\mathcal{L}_s^{t+1} = \prod_{v \in V} (\mathcal{L}_{exo}^{s_{v,t+1}} (r_{2s_v})^{r_{v,t} s_{v,t+1}}) \quad (44)$$

$$\mathcal{L}_r^{t+1} = \prod_{v \in V} ((i_{2r_v})^{i_{v,t} r_{v,t+1}} (1.0 - r_{2s_v})^{r_{v,t} r_{v,t+1}}) \quad (45)$$

$$\mathcal{L}_i^{t+1} = \prod_{v \in V} ((e_{2i_v})^{e_{v,t} i_{v,t+1}} (1.0 - i_{2r_v})^{i_{v,t} i_{v,t+1}}) \quad (46)$$

■

**Theorem V.1.** *Prize Collecting Dominating Set Vertex Cover (PCDSVC) is NP-hard, and it can be approximated by  $O(\log(|V^*|))$ .*

*Proof:*

PCDSVC is NP-hard since its special case *Dominating Set* is NP-hard that is obtained when all edge weights are 0 ( $w_{uv} = 0$ ).

Given PCDSVC problem over graph  $G^* = (V^*, E^*)$ , we construct *Minimum Hitting Set* instance  $(S, C)$  as follows: We define the set of elements as  $S = \{v \in V^*\} \cup \{e \in E^*\}$  where the cost of each element in  $S$  is  $w_u$  for every  $u \in V^*$  and  $w_{uv}$  for every  $(u, v) \in E^*$ . Subsets  $C = C_1 \cup C_2$  of  $S$  are defined as:  $C_1 = \{e_u, e_v, e_{uv}\}, \forall (u, v) \in E^*$  and  $C_2 = \{e_u, u \in P(v) \cup \{v\}\}, \forall v \in V^*$ . This reduction is linear time, approximation preserving and the solution of this *Minimum Hitting Set* gives us the solution for PCDSVC. Here  $|S| = |E^*| + |V^*|$  and **Greedy** method for *Set Cover* approximates this problem by  $\log(|S|) + 1 \approx O(\log(|E^*| + |V^*|)) + 1 \approx O(\log(|V^*|)) + 1$ .

One can also easily show that each *Minimum Hitting Set* instance can be reduced to PCDSVC and this reduction is also approximation preserving. Then, *Minimum Hitting Set* and PCDSVC are *equivalent* under linear reduction and this approximation ratio for PCDSVC is the best we can achieve unless  $P=NP$  [2].

■

**Theorem V.2.** *The Taylor expansion relaxation of (17) for bipartite diffusion models can be expressed as s-t mincut.*

*Proof:*

Minimization problem for bipartite  $\mathcal{M}$  has objective  $F_{bi}$  as seen in Equation 47.  $F_{bi}$  is a regular function [3]: when expressed as the summation of first and second-order terms as in Equation 48, each second order term  $E^{u,v}(s_{v,j-1}, i_{u,j-1})$  satisfies  $E^{u,v}(0, 0) + E^{u,v}(1, 1) \leq E^{u,v}(0, 1) + E^{u,v}(1, 0)$  in regular functions. Regular functions can be solved optimally by transforming it into s-t mincut [3]. Transformation is as follows:

$$\min -F_{bi} = \sum_{(u,v) \in E^*} \frac{1}{\log(1 - p_{uv})} (1 - i_{u,j-1}) s_{v,j-1} + \sum_{v \in E_j \cup S_j} w_v s_{v,j-1} + \sum_{v \in I_j \cup R_j} w_v i_{v,j-1} \quad (47)$$

$$-F_{bi} = \sum_{u \in I_j \cup R_j, v \in E_j \cup S_j} E^{u,v}(i_{u,j-1}, s_{v,j-1}) + \sum_{v \in I_j \cup R_j} E^v(i_{v,j-1}) + \sum_{v \in S_j \cup E_j} E^v(s_{v,j-1}) \quad (48)$$

We define new directed graph  $G' = (V', E')$  where  $V' = V^* \cup \{s\} \cup \{t\}$ . For every  $v \in V^*$ , we add edge  $(s, v)$  with weight  $E^v(1)$  if  $E^v(1) > 0$  and add edge  $(v, t)$  with weight  $-E^v(1)$  if  $-E^v(1) < 0$ . For every  $u \in I_j \cup R_j$  and  $v \in S_j \cup E_j$ , we add edge  $(u, v)$  with weight  $E^{u,v}(0, 1)$ . s-t mincut solution of this graph gives us the resulting node partition; after the cut edges removed, variables of the nodes that are reachable from  $s$  are assigned 1 and the variables of the nodes that have a path to  $t$  are assigned 0.

■

## REFERENCES

- [1] L. A. Wolsey and G. L. Nemhauser, *Integer and Combinatorial Optimization*. Wiley-Interscience, 1999.
- [2] U. Feige, "A threshold of  $\ln n$  for approximating set cover," *Journal of the ACM*, vol. 45, no. 4, pp. 634–652, Jul. 1998.
- [3] V. Kolmogorov and R. Zabih, "What energy functions can be minimized via graph cuts," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 26, pp. 65–81, 2004.