
Dynamic influence maximization

Binghui Peng
Columbia University
bp2601@columbia.edu

Abstract

We initiate a systematic study on *dynamic influence maximization* (DIM). In the DIM problem, one maintains a seed set S of at most k nodes in a dynamically involving social network, with the goal of maximizing the expected influence spread while minimizing the amortized updating cost. We consider two involution models. In the *incremental model*, the social network gets enlarged over time and one only introduces new users and establishes new social links, we design an algorithm that achieves $(1 - 1/e - \epsilon)$ -approximation to the optimal solution and has $k \cdot \text{poly}(\log n, \epsilon^{-1})$ amortized running time, which matches the state-of-art offline algorithm with only poly-logarithmic overhead. In the fully dynamic model, users join in and leave, influence propagation gets strengthened or weakened in real time, we prove that under the Strong Exponential Time Hypothesis (SETH), no algorithm can achieve $2^{-(\log n)^{1-o(1)}}$ -approximation unless the amortized running time is $n^{1-o(1)}$. On the technical side, we exploit novel adaptive sampling approaches that reduce DIM to the dynamic MAX-k coverage problem, and design an efficient $(1 - 1/e - \epsilon)$ -approximation algorithm for it. Our lower bound leverages the recent developed distributed PCP framework.

1 Introduction

Influence maximization (IM) is the algorithmic task of given a social network and a stochastic diffusion model, finding a seed set of at most k nodes with the largest expected influence spread over the network [29]. Influence maximization and its variants have been extensively studied in the literature over past decades [29, 12, 23, 39, 43, 32, 4, 20, 8], and it has applications in viral market, rumor control, advertising, etc.

Social influence can be highly dynamic and the propagation tendencies between users can alter dramatically over time. For example, in a Twitter network, new users join in and existing users drop out in real time, pop stars arise instantly for breaking news and trending topics; in a DBLP network, scientific co-authorship is built up and expands gradually over time. The classic IM algorithms make crucial assumptions on a stationary social network and they fail to capture the elastic nature of social networks. As a consequence, their seeding set could become outdated rapidly in a constantly involving social network. To mitigate the issue, one designs a dynamic influence maximization algorithm, which maintains a feasible seed set with high influence impact over time, and at the same time, saturates low average computation cost per update of the social network.

In this paper, we initiate a systematic study on the dynamic influence maximization (DIM) problem. Two types of involution models are of interests to us. In an *incremental model*¹, the social networks keep growing: new users join in and new social relationship is built up or strengthened over time. The motivating example is the DBLP network, where co-authorship gets expanded over time. Another justification is the preferential attachment involving model of social network [11, 35]. In a *fully*

¹These terms are standard notions in the dynamic algorithm literature, e.g., see [24, 1]

dynamic model, the social network involves over time, users can join in and leave out, social links emerge and disappear, influence impacts get strengthened or weakened. The motivating examples are social media networks like Twitter or Facebook, and the advertising market.

There are some pioneer works on the topic of dynamic influence maximization and various *heuristic* algorithms have been proposed in the literature [46, 38, 21]. These algorithms have different kinds of approximation guarantee on the seed set, but the crucial drawback is that the amortized running time per update is no better than $\Omega(n)$, here n is the total number of nodes in a social network. This is extremely unsatisfactory from a theoretical view, as it indicates these algorithm do no better than re-run the state-of-art IM algorithm upon every update. Hence, the central question over the field is that, *can we achieve the optimal approximation guarantee with no significant overhead on amortized running time?*

Our contribution We address the above questions and provide clear resolutions. In the incremental model, we prove it is possible to maintain a seed set with $(1 - 1/e - \epsilon)$ -approximation ratio in amortized running time of $k \cdot \text{poly}(\log n, \epsilon^{-1})$, which matches the state-of-art offline algorithm up to poly-logarithmic factors (Theorem 3.1). While in the fully dynamic algorithm, assuming the Strong Exponential Time Hypothesis (SETH), we prove no algorithm can achieve $2^{-(\log n)^{1-o(1)}}$ approximation unless the amortized running time is $n^{1-o(1)}$ (Theorem 4.3 and Theorem 4.4). This computational barrier draws a sharp separation between these two dynamic models and delivers the following surprising message: There is no hope to achieve any meaningful approximation guarantee for DIM, even one only aims for a slightly reduced running time than the naive approach of re-instantiating an IM algorithm upon every update.

On the technical side, our DIM algorithm in the incremental model uses novel adaptive sampling strategies and reduces the DIM problem to the dynamic MAX-k coverage problem. We provide an efficient dynamic algorithm to the later problem. Both the adaptive sampling strategy and the combinatorial algorithm could be of independent interests to the influence maximization community and the submodular maximization community. For the computaional hardness result, we exploit a novel application of the recent developed distribution PCP framework for fine grained complexity.

Related work Influence maximization as a discrete optimization task is first proposed in the seminal work of Kempe et al. [29], who propose the Independent Cascade (IC) model and the Linear Threshold (LT) model, prove the submodularity property and study the performance of greedy approximation algorithm. Upon then, influence maximization and its variant has been extensively studied in the literature, including scalable algorithms [19, 12, 45, 44], adaptive algorithms [23, 22, 39, 17, 43, 6], learning algorithms [20, 8, 33, 26], etc. For detailed coverage over the area, we refer interested readers to the survey of [16, 34].

Efficient and scalable influence maximization has been the central focus of the research. Influence maximization is known to be NP-hard to approximate within $(1 - 1/e)$ under the Independent Cascade model [29], and APX-hard under the Linear Threshold model [42], while at the same, the simple greedy algorithm achieves $(1 - 1/e)$ approximation under both diffusion models, but with running time $\Omega(mnk)$. The breakthrough comes from Borgs et al. [12], who propose the Reverse influence sampling (RIS) approach and their algorithm has running time $O((m + n)k\epsilon^{-3} \log^2 n)$. The running time has been improved to $O((m + n)k\epsilon^{-2} \log n)$ and made practical in [45, 44].

There are some initial efforts on the problem of dynamic influence maximization [46, 38, 21, 3, 36]. All of them provide heuristics and none of them has rigorous theoretical guarantee on both the approximation ratio and the amortized running time. Chen et al. [21] propose Upper Bound Interchange (UBI) method with $1/2$ -approximation ratio. Ohsaka et al. [38] design a RIS based algorithm that achieves $(1 - 1/e)$ -approximation. Inspired by recent advance on streaming submodular maximization, Wang et al. [46] give a practical efficient algorithm that maintains a constant approximated solution. However, none of the above mentioned algorithm has rigorous guarantees on the amortized running time, and they can be as worse as $\Omega(n)$. The failure on their theoretical guarantees are not accidental, as our hardness result indicates that $\Omega(n^{1-\epsilon})$ amortized running time is essential to achieve any meaningful approximation. Meanwhile, our algorithm for incremental model substantially deviates from these approaches and provides rigorous guarantee on amortized running time.

Influence maximization is closely related to submodular maximization [14, 5, 27, 10, 9, 7, 13, 28, 37, 31], a central area of discrete optimization. Ours is especially related to the dynamic submodular maximization problem, which has been recently studied in [37, 31]. The concurrent work of [37] and [31] gives $(1/2 - \epsilon)$ -approximation to the fully dynamic submodular under cardinality constraints, with amortized query complexity $O(k^2 \log^2 n \cdot \epsilon^{-3})$ and $O(\log^8 n \cdot \epsilon^{-6})$. The dynamic submodular maximization is studied under the query model and measured in query complexity, while we consider time complexity in dynamic influence maximization. These two are generally incomparable.

2 Preliminary

We consider the well-studied Independent Cascade (IC) model and the Linear Threshold (LT) model.

Independent Cascade model In the IC model, the social network is described by a directed influence graph $G = (V, E, p)$, where V ($|V| = n$) is the set of nodes and $E \subseteq V \times V$ ($|E| = m$) describes the set of directed edges. There is a probability $p_{u,v}$ associated with each directed edge $(u, v) \in E$. In the information diffusion process, each *activated* node u has one chance to activate its out-going neighbor v , with independent success probability $p_{u,v}$. The *live-edge* graph $L = (V, L(E))$ is a random subgraph of G , where each edge $(u, v) \in E$ is included in $L(E)$ with an independent probability $p_{u,v}$. The diffusion process can also be seen as follow. At time $\tau = 0$, a live-edge graph L is sampled and nodes in seed set $S \subseteq V$ are activated. At every discrete time $\tau = 1, 2, \dots$, if a node u was activated at time $\tau - 1$, then its out-going neighbors in L are activated at time τ . The propagation continues until no more activated nodes appears at a time step.

Linear Threshold model In the LT model, the social network is a directed graph $G = (V, E, w)$, with V denotes the set of nodes and E denotes the set of edge. There is a weight $w_{u,v} > 0$ associate with each edge $(u, v) \in E$ and the weight satisfies $\sum_{u \in N_{in}(v)} w_{u,v} \leq 1$ for every node v , where $N_{in}(v)$ contains all incoming neighbors of node v . In the information diffusion process, a threshold t_v is sampled uniformly and independently from $[0, 1]$ for each node v , and a node v becomes active if its active in-coming neighbors have their weights exceed the threshold t_v . In the LT model, the live-edge graph $L = (V, L(E))$ can be obtained as follow. Each node v samples an incoming neighbor from $N_{in}(v)$, where the node u is sampled with probability $w_{u,v}$ and the edge (u, v) is included in $L(E)$. With probability $1 - \sum_{u \in N_{in}(v)} w_{u,v}$, no edge is added. The diffusion process can also be described by the live-edge graph. That is, a live-edge graph L is sampled at the beginning and the influence is spread along the graph.

Influence maximization Given a seed set S , the influence spread of S , denoted as $\sigma(S)$, is the expected number of nodes activated from seed set S , i.e., $\sigma(S) = \mathbb{E}_{L \sim G} [|\Gamma(S, L)|]$, and $\Gamma(S, L)$ is the set of nodes reachable from S in graph L . In the influence maximization problem, our goal is to find a seed set S of size at most k that maximizes the expected influence, i.e., finding $S^* \in \arg \max_{S \subseteq V, |S| \leq k} \sigma(S)$.

The reverse reachable set [12] has been the key concept for all near-linear time IM algorithms [12, 45, 44].

Definition 2.1 (Reverse Reachable Set). *Under the IC model and the LT model, a reverse reachable (RR) set with a root node v , denoted R_v , is the random set of nodes that node v reaches in one reverse propagation. Concretely, R_v can be derived by randomly sampling a live-edge graph $L = (V, L(E))$ and include all nodes in V that can reach v under the live-edge graph L .*

Dynamics of network Social networks are subject to changes and involve over time. We consider the following dynamic model of a social network. At each time step t , one of the following four types of changes could happen. (1) Introduction of a new user. This corresponds to insert a new node to influence graph; (2) Establishment of a new relationship. It is equivalent to insert a new directed edge to the graph²; (3) Diminishing of an old relationship. It is equivalent to remove an existing edge of the graph; (4) Leave of an old user. This means to remove an existing node in the influence graph.

²In order to adding a new directed edge (u, v) , one also specifies the probability $p_{u,v}$ under IC model or the weight $w_{u,v}$ under LT model.

In the *incremental* model, we only allow the first two types of changes. That is, we assume the social network gets enlarged over time and we only consider the introduction of new users and new relations. In the *fully dynamic* model, we allow all four types of changes.

The theoretical results developed in this paper also adapts to other forms of change on the social network, including (a) strengthen of influence, which means the increase of the propagation probability p_e for some edge under IC model or the increase of the weight w_e under LT model; (b) weaken the influence. which means the decrease of propagation probability p_e under IC model or the decrease of the weight w_e under LT model.

Dynamic influence maximization Let $G_t = (V_t, S_t)$ be the social network at time t . Define the expected influence to be $\sigma_t(S) = \mathbb{E}_{L \sim G_t} [|\Gamma(S, L)|]$. The goal of the dynamic influence maximization problem is to maintain a seed set S_t of size at most k that maximizes the expected influence at *every* time step t . That is to say, we aim to find $S_t^* \in \arg \max_{S_t \subseteq V_t, |S_t| \leq k} \sigma(S_t)$ for all t . The (expected) amortized running time of an algorithm is defined as the (expected) average running time per update.

Dynamic MAX-k coverage The (dynamic) influence maximization problem is closely related to the (dynamic) MAX-k coverage problem. In a MAX-k coverage problem, there is a collection of n sets \mathcal{A} defined on the ground element $[m]$. The goal is to find k sets A_1, \dots, A_k such that their coverage is maximized, i.e. finding $\arg \max_{A_1, \dots, A_k \in \mathcal{A} \mid \cup_{i \in [k]} A_i} |\cup_{i \in [k]} A_i|$. The problem can also be formulated in terms of a bipartite graph $G = (V_L, V_R, E)$, named coverage graph, where V_L ($|V_L| = n$) corresponds to set in \mathcal{A} and V_R ($|V_R| = m$) corresponds to element in $[m]$, for any $i \in [n], j \in [m]$, there is an edge between $V_{L,i}$ and $V_{R,j}$ if and only if $j \in A_i$, here $V_{L,i}$ is the i -th node in V_L and $V_{R,j}$ is the j -th node in V_R . In the dynamic MAX-k coverage problem, nodes and edges arrive or leave one after another. Let $V \subseteq V_L$. At any time step t , we use $f_t(S)$ to denote the number of neighbors of nodes in S . Our goal is to maintain a set of node $S_t^* \subseteq V_L$ that maximizes the coverage for every time step t . That is, finding $S_t^* = \arg \max_{S_t \subseteq V_L, |S_t| \leq k} f_t(S_t)$.

Submodular functions Let V be a finite ground set and $f : 2^V \rightarrow \mathbb{R}$ be a set function. Given two sets $X, Y \subseteq V$, the marginal gain of Y with respect to X is defined as $f_X(Y) := f(X \cup Y) - f(X)$. The function f is *monotone*, if for any element $e \in V$ and any set $X \subseteq V$, it holds that $f_X(e) \geq 0$. The function f is *submodular*, if for any two sets X, Y satisfy $X \subseteq Y \subseteq V$ and any element $e \in V \setminus Y$, one has $f_X(e) \geq f_Y(e)$. The influence spread function σ is proved to be monotone and submodular under both the IC model and the LT model in the seminar work of [29].

3 Dynamic influence maximization on a growing social network

We study the dynamic influence maximization problem in the incremental model. Our main result is to show that it is possible to maintain an $(1 - 1/e - \epsilon)$ -approximate solution in $k \cdot \text{poly}(\epsilon^{-1}, \log n)$ amortized time. The amortized running time is near optimal and matches the state-of-art offline algorithm up to poly-logarithmic factors.

Theorem 3.1. *In the incremental model, there is a randomized algorithm for the dynamic influence maximization under both Independent Cascade model and the Linear Threshold model. With probability $1 - \delta$, the algorithm maintains an $(1 - 1/e - \epsilon)$ -approximately optimal solution in all time steps, with amortized running time $O(k\epsilon^{-3} \log^4(n/\delta))$.*

Technique overview We provide a streamlined technique overview over our approach. All existing near-linear time IM algorithms [12, 45, 44] make use of the reverse influence sampling (RIS) technique and adopt the following two-stage framework. In the first stage, the algorithm *uniformly* samples a few RR sets, with the guarantee that the total time steps is below some threshold. In the second stage, the algorithm solves a MAX-k coverage problem on the coverage graph induced by the RR sets. Both steps are challenging to be applied in the dynamic setting and require new ideas.

Following the RIS framework, we first reduce the dynamic influence maximization problem to a dynamic MAX-k coverage problem (see Section 3.1). The first obstacle in dynamic setting is that it is hard to dynamically maintain a uniform distribution over nodes, and hence, one can not uniformly sample the RR sets. We circumvent it with a simple yet elegant idea, instead of uniformly sampling RR sets, we include and sample each node's RR set independently with a uniformly fixed probability p . This step is easy to implement in a dynamic stream, but to make the idea work, it needs some

extra technique machinery. First of all, it is not obvious how to set the sampling probability p , as if p is too large, it brings large overhead on the running time; and if it is too small, the sampled coverage graph won't yield a good approximation to the influence spread function. We use a simple trick. We first uniformly sample RR sets within some suitable amount of steps (the ESTIMATE procedure), and by doing so, we have an accurate estimate on the sampling probability (Lemma 3.3). Note that we only execute this step at beginning and it is an one-time cost. Independently sampling is friendly to dynamic algorithm, but we still need to control the total number of steps to execute it. More importantly, the sampling probability should not be a fixed one throughout the dynamic stream, since the average cost to sample a RR set could go up and down, and if the cost is too large, we need to decrease the sample probability. A simple rescue is to restart the above procedure once the coverage graph induced by independently sampling is too large, but this brings along two additional issues. First, the stopping time is highly correlated with the realization of RR sets, and this is especially problematic for independent sampling, as each incremental update on RR sets does not follow identical distribution. Second, we certainly don't want to restart the procedure for too many times. After all, if we need to restart upon every few updates, we gain nothing than the naive approach. We circumvent the first issue by maintaining two pieces of coverage graph \mathcal{H}_{est} and \mathcal{H}_{cv} . Especially, we use \mathcal{H}_{est} to estimate the average cost for sampling a RR set, and use \mathcal{H}_{cv} to approximate the influence spread function. This decouples the correlation. By Martingale concentration, we can prove that (1) \mathcal{H}_{est} yields good estimation on the cost of sampling RR sets (Lemma 3.5), (2) \mathcal{H}_{cv} gives good approximation on the influence spread function (Lemma 3.9), and at the same time, (3) conditioning on the first event, the computation cost for \mathcal{H}_{cv} won't be too large (Lemma 3.6). We circumvent the second issue by setting suitable threshold for constructing \mathcal{H}_{est} (see Line 9 in INSERT-EDGE), and we periodically restart the algorithm each time the number of nodes or edges get doubled. The later step could lead to an extra $O(\log n)$ overhead, but it guarantees that each time the construction time of \mathcal{H}_{est} goes above the threshold, the average steps for sampling a RR set increase by a factor of 2 (see Lemma 3.11). Hence, we restart the entire procedure no more than $O(\log n)$ times.

We next give an efficient dynamic algorithm for dynamic MAX-k coverage (Section 3.2). The offline problem can be solved in linear time via greedy algorithm, but it requires several new ideas to make the dynamic variant have low amortized cost. The major issue with offline greedy algorithm is that we don't know which element has the maximum marginal contribution until we see all elements, and in the dynamic setting, we can't afford to wait till the end. We get around it with two adaptations. First, we divide the algorithm into $O(\epsilon^{-1} \log n)$ threads³ and for each thread, we make a guess on the optimal value (the i -th guess is $(1 + \epsilon)^{i-1}$). Second, instead of selecting element with the maximum margin, we maintain a threshold $(\text{OPT}_i - f(S_i))/k$ and add any element whose marginal lies above the threshold, where S_i is the current solution set maintained in the i -th thread. We prove this still gives $(1 - 1/e)$ -approximation (Lemma 3.13). The key challenge is that once we add a new element to the solution set, the threshold changes and we need to go over all existing element again. We carefully maintain a dynamic data structure and prove the amortized running time for each thread is $O(1)$ by a potential function based argument (Lemma 3.14).

3.1 Reducing DIM to dynamic MAX-k coverage

Lemma 3.2. *In the incremental model, suppose there is an algorithm that maintains an α -approximate solution to the dynamic MAX-k coverage, with amortized running time of \mathcal{T}_{mc} . Then there is an algorithm, with probability at least $1 - \delta$, maintains an $(\alpha - 2\epsilon)$ -approximate solution for dynamic influence maximization and has amortized running time at most $O(\mathcal{T}_{\text{mc}} k \epsilon^{-2} \log^3(n/\delta))$.*

The reduction is described in Algorithm 1-5, we provide some explanations here. From a high level, the execution of the algorithm is divided into *phases*. It restarts whenever the number of edges or the number of nodes gets doubled from the previous phase. Since the total number of steps is at most $(m + n)$, there are at most $O(\log n)$ phases. Within each phase, the algorithm is divided into *iterations*. At the beginning of each iteration, the algorithm runs the ESTIMATE procedure in order to estimate the total number of RR set it needs to sample. It then fixes a sampling probability p and samples two graphs, \mathcal{H}_{cv} and \mathcal{H}_{est} , dynamically. As explained before, \mathcal{H}_{est} is used to measure the number of steps one takes for sampling RR sets, and \mathcal{H}_{cv} is used to approximate the influence spread function. We stop sampling \mathcal{H}_{est} and \mathcal{H}_{cv} and start a new iteration whenever the number of step for constructing \mathcal{H}_{est} exceeds $16Rm_0$.

³We remark this approach has been used for submodular maximization under various setting [5, 9, 31].

Due to space limits, we defer detailed proof to the full version of this paper and outline a sketch below.

Algorithm 1 INITIALIZE

```

1:  $m \leftarrow 0, m_0 \leftarrow 0, n \leftarrow 0, n_0 \leftarrow 0.$ 
2:  $V \leftarrow \emptyset, E \leftarrow \emptyset, \mathcal{H}_{cv} \leftarrow \emptyset, \mathcal{H}_{est} \leftarrow \emptyset.$ 
3:  $R \leftarrow 0, p \leftarrow 0.$ 

```

Algorithm 2 ESTIMATE

```

1:  $K \leftarrow 0.$ 
2: repeat
3:   Sample a RR set uniformly over  $V.$ 
4:    $K \leftarrow K + 1.$ 
5: until  $Rm_0$  steps have been taken
6: return  $K$ 

```

Algorithm 3 BUILDGRAPH

```

1:  $\mathcal{H}_{cv} \leftarrow (V, \emptyset, \emptyset), \mathcal{H}_{est} \leftarrow (V, \emptyset, \emptyset),$ 
2:  $R \leftarrow c\epsilon^{-2}k \log(n_0/\delta).$ 
3:  $K \leftarrow \text{ESTIMATE}().$ 
4:  $p \leftarrow \frac{K}{n_0}.$ 
5: for node  $v \in V$  do
6:   Sample the RR set  $R_{v,est}$  with prob  $p.$ 
7:   Sample the RR set  $R_{v,cv}$  with prob  $p.$ 
8:   for node  $u \in R_{v,cv}$  do
9:     INSERT-EDGE-COV( $u, v$ ) to  $\mathcal{H}_{cv}.$ 
10:  end for
11: end for

```

Algorithm 4 INSERT-NODE(u)

```

1:  $V \leftarrow V \cup \{u\}, n \leftarrow n + 1.$ 
2: if  $n \geq 2n_0$  then
3:    $m_0 \leftarrow m, n_0 \leftarrow n.$ 
4:   BUILDGRAPH().
5: else
6:    $V_{L,est} \leftarrow V_{L,est} \cup \{u\}.$ 
7:    $V_{L,cv} \leftarrow V_{L,cv} \cup \{u\}.$ 
8:    $V_{R,est} \leftarrow V_{R,est} \cup \{u\}$  with prob  $p.$ 
9:    $V_{R,cv} \leftarrow V_{R,cv} \cup \{u\}$  with prob  $p.$ 
10: end if

```

Algorithm 5 INSERT-EDGE(u, v)

```

1:  $E \leftarrow E \cup \{(u, v)\}, m \leftarrow m + 1$ 
2: if  $m \geq 2m_0$  then
3:    $m_0 \leftarrow m, n_0 \leftarrow n.$ 
4:   BUILDGRAPH().
5: else
6:   for node  $v' \in V_{R,est}$  do
7:     Augment the RR set  $R_{v',est}$  to
       $R_{v',est}^{new}$ .
8:      $R_{v',est} \leftarrow R_{v',est}^{new}.$ 
9:     if the total steps of building  $\mathcal{H}_{est}$ 
      exceed  $16Rm_0$  then
10:       $m_0 \leftarrow m, n_0 \leftarrow n.$ 
11:      Restart and BUILDGRAPH().
      end if
12:   end for
13:   for node  $v' \in V_{R,cv}$  do
14:     Augment the RR set  $R_{v',cv}$  to  $R_{v',cv}^{new}.$ 
15:     INSERT-EDGE-COV( $u', v'$ ) to  $\mathcal{H}_{cv}$ 
      for all  $u' \in R_{v',cv}^{new} \setminus R_{v',cv}.$ 
16:      $R_{v',cv} \leftarrow R_{v',cv}^{new}.$ 
17:   end for
18: end if
19: end if

```

We first focus on one single iteration of the algorithm. Recall that in each iteration, the number of nodes satisfies $n \in [n_0, 2n_0)$ and the number of edge satisfies $m \in [m_0, 2m_0)$, and the number of steps for sampling \mathcal{H}_{est} never exceeds $16Rm_0$.

Bounds at initialization We prove the ESTIMATE procedure makes a good estimate on the sampling probability at the beginning of each iteration. Let $\text{AVG} \cdot m_0$ be the expected number of steps taken of sampling a uniformly random RR set, and recall K is the number of RR sets ESTIMATE samples.

Lemma 3.3. *With probability at least $1 - \frac{2\delta}{n^{ck/8}}$, the number of RR sets ESTIMATE samples satisfies $(1 - \epsilon) \cdot \frac{R}{\text{AVG}} \leq K \leq (1 + \epsilon) \cdot \frac{R}{\text{AVG}}.$*

The initial steps to built \mathcal{H}_{est} and \mathcal{H}_{cv} are within $[(1 - \epsilon)^2 Rm_0, (1 + \epsilon) Rm_0]$ with high probability.

Lemma 3.4. *Conditioning on the event of Lemma 3.3, with probability at least $1 - \frac{4\delta}{n^{ck/8}}$, the total number of initial steps for building \mathcal{H}_{est} and \mathcal{H}_{cv} are within $[(1 - \epsilon)^2 R, (1 + \epsilon)^2 R],$*

Bounds on adaptive sampling We dynamically sample and augment RR sets after we initiate \mathcal{H}_{est} and \mathcal{H}_{cv} . Consider the t -th edge to arrive, let $Z_{t,i} \cdot 2m_0$ be the number of edges checked by INSERT-EDGE when it augments the RR set of the i -th node. Note if INSERT-NODE and BUILDGRAPH

did not sample the i -th node and its RR set, then we simply take $Z_{t,i} = 0$. Slightly abuse of notation, we write Z_j to denote the j -th random variable in the sequence $Z_{0,1}, \dots, Z_{0,n_0}, Z_{1,1}, \dots$, and we know $\{Z_j\}_{j \geq 1}$ forms a martingale. Let J be the number of steps executed by INSERT-EDGE for \mathcal{H}_{est} , our next Lemma asserts that the total number of steps is around its expectation.

Lemma 3.5. *Let J_1 be the smallest number such that $\sum_{i=1}^{J_1} \mathbb{E}[Z_i] \geq (1 - \epsilon)8R$, and J_2 be the largest number that $\sum_{i=1}^{J_2} \mathbb{E}[Z_i] \leq (1 + \epsilon)8R$. With probability at least $1 - \frac{2\delta}{n^{ck/3}}$, we have $J_1 \leq J \leq J_2$.*

With high probability, the total number of edge checked for \mathcal{H}_{cv} is bounded by $16(1 + \epsilon)^2 Rm_0$.

Lemma 3.6. *Conditioning on the event of Lemma 3.5, with probability at least $1 - \frac{\delta}{n^{ck}}$, the total number of edges checked by INSERT-EDGE for \mathcal{H}_{cv} is at most $16(1 + \epsilon)^2 Rm_0$.*

Approximation guarantee We prove that the coverage function induced on \mathcal{H}_{cv} yields a good approximation on the influence spread function. At any time step t , let V_t be the set of node at time step t . For any node $v \in V_t$ and node set $S \subseteq V_t$, let $x_{t,v,S} = 1$ if $S \cap R_{v,\text{cv}} \neq \emptyset$, and $x_{t,v,S} = 0$ otherwise. That is, $x_{t,v,S} = 1$ iff S has non-zero intersection with the RR set of node v at time step t . We note that if a node v was not sampled by INSERT-NODE and BUILDGRAPH, and hence does not appear in $V_{R,\text{cv}}$, we set $x_{t,v,S} = 0$.

Definition 3.7 (Normalized coverage function). *Define $f_{\text{cv},t} : 2^{V_t} \rightarrow \mathbb{R}^+$ to be the normalized coverage function. $f_{\text{cv},t}(S) = \frac{n_0}{K} \sum_{v \in V_t} x_{t,v,S}$.*

The normalized coverage function $f_{\text{cv},t}$ is an unbiased estimator on the influence spread function and it achieves good approximation with high probability.

Lemma 3.8. *At any time step t and $S \subseteq V_t$, we have $\mathbb{E}[f_{\text{cv},t}(S)] = \sigma_t(S)$, where σ_t is the influence spread function at time t .*

Lemma 3.9. *After initializing \mathcal{H}_{est} and \mathcal{H}_{cv} , at any time step t . Let $S \subseteq V_t$, $|S| \leq k$. Assume the condition in Lemma 3.3 holds, then we have*

$$\Pr [f_{\text{cv},t}(S) > \mathbb{E}[f_{\text{cv},t}(S)] + \epsilon \text{OPT}_t] \leq \frac{\delta}{n^{ck/6}}, \text{ and,} \quad (1)$$

$$\Pr [f_{\text{cv},t}(S) < \mathbb{E}[f_{\text{cv},t}(S)] - \epsilon \text{OPT}_t] \leq \frac{\delta}{n^{ck}}. \quad (2)$$

The expectation is taken over the randomness of the construction of \mathcal{H}_{cv} , and OPT_t is defined as the maximum (expected) influence spread of a set of size k at time t .

The following Lemma indicates pointwise approximation carries over the approximation ratio.

Lemma 3.10. *Let $c \geq 2$ and let $f : 2^V \rightarrow \mathbb{R}^+$ be an arbitrary set function. Let \mathcal{D} be a distribution over functions g such that for all $S \subseteq V$, $\Pr_{g \sim \mathcal{D}}[|f(S) - g(S)| - \gamma] \leq \frac{\delta}{n^{ck}}$. Let S_g be an α -approximate solution for function g , i.e. $g(S_g) \geq \alpha \max_{S \subseteq V, |S| \leq k} g(S)$, then we have*

$$\Pr_{g \sim \mathcal{D}} \left[f(S_g) \leq \alpha \max_{S \subseteq V, |S| \leq k} f(S) - 2\gamma \right] \leq \frac{\delta}{n^{ck/2}}.$$

Bounds for amortized running time We next bound the total number of iterations within each phase. The key observation is that each time the algorithm restarts, the average steps of sampling a random RR set increases at least by a factor of 2, with high probability.

Lemma 3.11. *With probability at least $1 - \frac{4\delta}{n^{ck/16}}$, there are at most $O(\log n)$ iterations in a phase.*

Proof Sketch of Lemma 3.2. For correctness, by Lemma 3.8 and Lemma 3.9, the normalized coverage function $f_{\text{cv},t}$ (defined on \mathcal{H}_{cv}) guarantees $|f_{\text{cv},t}(S) - \mathbb{E}[f_{\text{cv},t}(S)]| = |f_{\text{cv},t}(S) - \sigma_t(S)| \leq \epsilon \text{OPT}_t$ for every time step t . Combining with Lemma 3.10, an α -approximate solution to the dynamic MAX- k coverage problem translates to a solution set S_t that satisfies $f(S) \geq (\alpha - 2\epsilon) \text{OPT}_t$.

For amortized running time. There are $O(\log n)$ phases, and by Lemma 3.11, there are at most $O(\log n)$ iterations in each phase. Within each iteration, the ESTIMATE procedure and the construction of \mathcal{H}_{est} takes $O(m_0 c k \epsilon^{-2} \log n)$ steps in total and $O(k \epsilon^{-2} \log n)$ per update. By Lemma 3.4 and Lemma 3.6, with high probability, the construction of \mathcal{H}_{cv} takes $16(1 + \epsilon)^2 Rm_0$ steps in total and $O(k \epsilon^{-2} \log n)$ steps per updates. Hence, there are $O(k m_0 \epsilon^{-2} \log n)$ updates to the dynamic MAX- k coverage problem on \mathcal{H}_{cv} . Taking an union bound, the overall amortized running time per update is bounded by $\log n \cdot \log n \cdot (k \epsilon^{-2} \log n + k \epsilon^{-2} \log n + k \epsilon^{-2} \log n \cdot \mathcal{T}_{\text{mc}}) \leq O(\mathcal{T}_{\text{mc}} k \epsilon^{-2} \log^3 n)$. \square

Algorithm 6 INITIALIZE

1: $\text{OPT}_i \leftarrow (1 + \epsilon)^i$, $S_i = \emptyset$, $\forall i \in I$

Algorithm 7 REVOKE(i)

1: **if** there exists a node u such that $f_{S_i}(u) \geq \frac{\text{OPT}_i - f(S_i)}{k}$ **and** $|S_i| < k$ **then**
2: $S_i \leftarrow S_i \cup \{u\}$.
3: REVOKE(i).
4: **end if**

Algorithm 8 INSERT-EDGE-COV(u, v)

1: **for** $i \in I$ **do**
2: **if** $f_{S_i}(u) \geq \frac{\text{OPT}_i - f(S_i)}{k}$ **and** $|S_i| < k$ **then**
3: $S_i \leftarrow S_i \cup \{u\}$.
4: REVOKE(i).
5: **end if**
6: **end for**

3.2 Solving dynamic MAX- k coverage in near linear time

For ease of presentation, we assume an upper bound on the value of n is known and we set $I = \{0, 1, \dots, \lceil \epsilon^{-1} \log n \rceil\}$. The algorithm maintains $|I|$ threads and for the i -th thread, INSERT-EDGE-COV and REVOKE augments the solution set S_i only when the marginal value of a node u exceeds the threshold, i.e. $f_{S_i}(u) \geq \frac{\text{OPT}_i - f(S_i)}{k}$. In particular, the threshold decreases over time and each time it decreases, REVOKE scans over all existing nodes in V_L . The algorithm returns the solution set with maximum value at each time step, i.e., return $\arg \max_{S_i} f_t(S_i)$.

Theorem 3.12. *In the incremental model, there is an algorithm for dynamic MAX- k coverage that maintains a solution set with $(1 - 1/e - \epsilon)$ -approximation at every time step and the amortized running time of the algorithm is at most $O(\epsilon^{-1} \log n)$.*

The approximation guarantee and the amortized running time of Algorithm 6-8 are analysed separately. We defer detailed proof to the full version of this paper and outline a sketch below.

Lemma 3.13. *The solution set is $(1 - 1/e - \epsilon)$ -approximate to the optimal one.*

Proof Sketch. Let the current optimum satisfies $(1 + \epsilon)^i < \text{OPT} \leq (1 + \epsilon)^{i+1}$ for some $i \in I$. For the i -th thread, if $|S_i| < k$, then one can prove $f(S_i) \geq (1 + \epsilon)^{-1} \text{OPT}$. On the other side, if $|S_i| = k$. Let $s_{i,j}$ be the j -th element added to the set S_i , and denote $S_{i,j} = \{s_{i,1}, \dots, s_{i,j}\}$. Our algorithm guarantees $f(S_{i,j+1}) - f(S_{i,j}) \geq \frac{1}{k}(\text{OPT}_i - f(S_{i,j}))$ for all $j \in [k]$. Unravel the recursion, one has $f(S_i) \geq (1 - \frac{1}{e} - \epsilon) \text{OPT}$. \square

Lemma 3.14. *The algorithm 6-8 has amortized running time $O(\epsilon^{-1} \log k)$.*

Proof Sketch. We prove that for each thread $i \in I$, the amortized running time is $O(1)$. Let $V_R^i \subseteq V_R$ be nodes covered by S_i . Let $P(t)$ denote the number of operations performed on the i -th thread up to time t . For each edge $e = (u, v)$, let X_e denote whether the edge e is covered by V_R^i , i.e. $X_e = 1$ if $v \in V_R^i$ and $X_e = 0$ otherwise. Similarly, for each node $v \in V_R$, let Y_v denote whether node v is included in V_R^i , i.e., $Y_v = 1$ if $v \in V_R^i$ and $Y_v = 0$ otherwise. For each node $u \in V_L$, let Z_u denote whether node u is included in S_i , i.e., $Y_v = 1$ if $v \in S_i$ and $Y_v = 0$ otherwise. Define the potential function $\Phi : t \rightarrow \mathbb{R}^+$:

$$\Phi(t) := 2|V_t| + 2|E_t| + 2 \sum_{e \in E_t} X_e + \sum_{v \in V_R} Y_v + 2 \sum_{v \in V_L} Z_u.$$

We prove $P(t) \leq \Phi(t)$ always holds. This suffices for our purpose as one can easily show $\Phi(t) \leq 5t$. The claim is executed by an induction showing $P(t) - P(t-1) \leq \Phi(t) - \Phi(t-1)$ holds for t . \square

Taking $\alpha = 1 - 1/e$ and $\mathcal{T}_{\text{mc}} = O(\epsilon^{-1} \log n)$, we wrap up the proof of Theorem 3.1.

4 Fully dynamic influence maximization

In the fully dynamic model, the social network involves over time and all four types of change exhibit. Our main delivery is a sharp computational hardness result. We prove that under the SETH, unless

the amortized running time is $n^{1-o(1)}$, the approximation ratio can not even be $2^{-(\log n)^{1-o(1)}}$. We first provide some background on fine-grain complexity. We reduce from the Strong Exponential Time Hypothesis (SETH), which is a pessimistic version of $P \neq NP$ that postulates that much better algorithms for k -SAT do not exist. The SETH is first proposed in the seminal work of [25] and it is widely believed in the computational complexity, see [47, 41] for detailed survey.

Conjecture 4.1 (Strong Exponential Time Hypothesis (SETH), [25]). *For any $\epsilon > 0$, there exists $k \geq 3$ such that k -SAT on variables can not be solved in time $O(2^{(1-\epsilon)n})$.*

Our starting point is the following SETH-based hardness of approximation result. The result is proven in [2, 15] using the distributed PCP framework [2, 40] for hardness of approximation results in P.

Theorem 4.2 ([2, 15, 1]). *Let $\epsilon > 0$, $m = n^{o(1)}$ and $t = 2^{(\log n)^{1-o(1)}}$. Given two collections of n sets \mathcal{A} and \mathcal{B} over universe $[m]$. Unless SETH is false, no algorithm can distinguish the following two cases in $O(n^{2-\epsilon})$:*

YES instance. *There exists two sets $A \in \mathcal{A}$, $B \in \mathcal{B}$ such that $B \subseteq A$;*

NO instance. *For every $A \in \mathcal{A}$, $B \in \mathcal{B}$ we have $|A \cap B| < |B|/t$.*

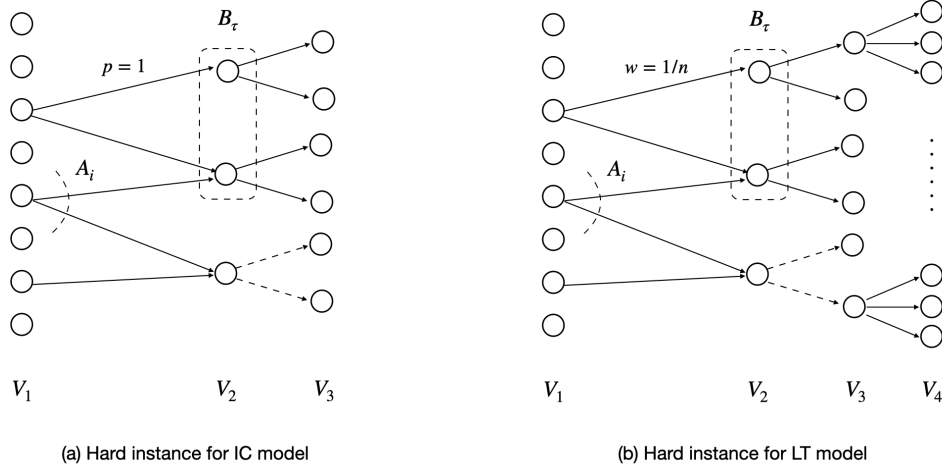


Figure 1: The hard instance for fully dynamic influence maximization.

Theorem 4.3. *Assuming SETH, in the fully dynamic influence maximization problem under IC model, no algorithm can achieve $2^{-(\log n)^{1-o(1)}}$ approximation unless the amortized running time is $n^{1-o(1)}$.*

Proof Sketch. Let $m = n^{o(1)}$, $t = 2^{(\log n)^{1-o(1)}}$, $k = 1$. Given an instance \mathcal{A}, \mathcal{B} of the problem in Theorem 4.2, we reduce it to DIM. We assume $|B_\tau| \geq t$ for all $\tau \in [n]$ as we can duplicate the ground element for t times. Consider the following influence graph $G = (V, E, p)$, where the influence probability on all edges are 1. The node set V is partitioned into $V = V_1 \cup V_2 \cup V_3$, $|V_1| = n$, $|V_2| = m$, $|V_3| = m^2 t$. Intuitively, the i -th node $v_{1,i}$ of V_1 corresponds to the set $A_i \in \mathcal{A}$, and the j -th node $v_{2,j}$ in V_1 correspond to the j -th element of ground set $[m]$. There is a directed edge from $v_{1,i}$ to $v_{2,j}$ iff the j -th element is contained in the set A_i . We write nodes in V_3 as $\{v_{3,j,\ell}\}_{j \in [m], \ell \in [mt]}$, and there is a directed edge node $v_{2,j}$ to $v_{3,j,\ell}$, $j \in [m], \ell \in [mt]$.

Consider the following update sequence of the DIM problem. The graph G is loaded first and then followed by n consecutive epochs. In the τ -th epoch, all edges between V_2 and V_3 are deleted, and for each $j \in B_\tau$, we add back the edge between $v_{2,j}$ and $v_{3,j,\ell}$ for all $\ell \in [mt]$.

One can show the total number of updates is at most $n^{1+o(1)}$. We prove by contrary and suppose there is an algorithm for DIM that achieves $2/t$ -approximation in $n^{1-\epsilon}$ time. Under the above reduction, we output YES, if for some epoch $\tau \in [n]$, the DIM algorithm outputs a solution with influence spread greater than $2m|B_\tau|$. We output NO otherwise. One can prove (1) if there exists $A_i \in \mathcal{A}$, $B_\tau \in \mathcal{B}$ such that $B_\tau \subseteq A_i$, then we output YES, and, (2) if $|A_i \cap B_\tau| < |B_\tau|/t$ for any $i, \tau \in [n]$,

the algorithm outputs NO. Hence, we conclude under SETH, there is no $2/t$ -approximation algorithm unless the amortized running time is $O(n^{1-\epsilon})$. \square

The lower bound can be extended to the LT model, under the additional constraints that the algorithm only selects seeds from a prescribed set $V' \subseteq V$. This is a natural assumption that has been made/discussed in previous work [18, 30, 42] for the LT model. The construction is similar to Theorem 4.4, with the exception that (1) the weight on edges between V_1 and V_2 are $1/n$ and all other edges have weight 1, (2) the node set V is partitioned into four parts $V_1 \cup V_2 \cup V_3 \cup V_4$. $|V_4| = nm^2t$ and each node in V_3 is connected to n nodes in V_4 . Detailed proof can be found in the full version of this paper.

Theorem 4.4. *Assuming SETH, for the fully dynamic influence maximization problem under LT model, if the algorithm is only allowed to select seed from a prescribed set, then no algorithm can achieve $2^{-(\log n)^{1-o(1)}}$ approximation unless the amortized running time is $n^{1-o(1)}$.*

5 Discussion

We study the dynamic influence maximization problem and provide sharp computational results on the incremental update model and the fully dynamic model. In the incremental model, we provide an algorithm that maintains a seed set with $(1 - 1/e - \epsilon)$ -approximation and has amortized running time $k \cdot \text{poly}(\log n, \epsilon^{-1})$, which matches the state of art offline IM algorithm up to poly-logarithmic factor. For the fully dynamic model, we prove that under SETH, no algorithm can achieve $2^{-(\log n)^{1-o(1)}}$ approximation unless the amortized running time is $n^{1-o(1)}$. There are a few interesting questions for future investigation: (1) Further improve the amortized running time in the incremental model. In particular, is it possible to reduce the amortized running time of the dynamic MAX-k coverage procedure to $O(1)$? (2) Investigate fully dynamic influence maximization problem under mild assumptions, e.g. what if the graph is bipartite?

Acknowledgement

Binghui Peng wishes to thank Xi Chen for useful discussions on dynamic submodular maximization. Binghui Peng is supported in part by Christos Papadimitriou’s NSF grants CCF-1763970 AF, CCF-1910700 AF and a softbank grant, and by Xi Chen’s NSF grants NSF CCF-1703925.

Broader Impact

Our work is mainly theoretical with no foreseeable ethical issues.

References

- [1] Amir Abboud, Raghavendra Addanki, Fabrizio Grandoni, Debmalya Panigrahi, and Barna Saha. Dynamic set cover: improved algorithms and lower bounds. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pages 114–125, 2019.
- [2] Amir Abboud, Aviad Rubinfeld, and Ryan Williams. Distributed pcp theorems for hardness of approximation in p. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 25–36. IEEE, 2017.
- [3] Charu C Aggarwal, Shuyang Lin, and Philip S Yu. On influential node discovery in dynamic social networks. In *Proceedings of the 2012 SIAM International Conference on Data Mining*, pages 636–647. SIAM, 2012.
- [4] Noga Alon, Iftah Gamzu, and Moshe Tennenholtz. Optimizing budget allocation among channels and influencers. In *Proceedings of the 21st international conference on World Wide Web*, pages 381–388, 2012.
- [5] Ashwinkumar Badanidiyuru, Baharan Mirzasoleiman, Amin Karbasi, and Andreas Krause. Streaming submodular maximization: Massive data summarization on the fly. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 671–680, 2014.
- [6] Ashwinkumar Badanidiyuru, Christos Papadimitriou, Aviad Rubinfeld, Lior Seeman, and Yaron Singer. Locally adaptive optimization: Adaptive seeding for monotone submodular functions. In *Proceedings of the twenty-seventh annual ACM-SIAM symposium on Discrete algorithms*, pages 414–429. SIAM, 2016.
- [7] Eric Balkanski, Adam Breuer, and Yaron Singer. Non-monotone submodular maximization in exponentially fewer iterations. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 2359–2370, 2018.
- [8] Eric Balkanski, Nicole Immorlica, and Yaron Singer. The importance of communities for learning to influence. *Advances in Neural Information Processing Systems*, 30:5862–5871, 2017.
- [9] Eric Balkanski, Aviad Rubinfeld, and Yaron Singer. An optimal approximation for submodular maximization under a matroid constraint in the adaptive complexity model. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pages 66–77, 2019.
- [10] Eric Balkanski and Yaron Singer. The adaptive complexity of maximizing a submodular function. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1138–1151, 2018.
- [11] Albert-Laszlo Barabási, Hawoong Jeong, Zoltan Néda, Erzsebet Ravasz, Andras Schubert, and Tamas Vicsek. Evolution of the social network of scientific collaborations. *Physica A: Statistical mechanics and its applications*, 311(3-4):590–614, 2002.
- [12] Christian Borgs, Michael Brautbar, Jennifer Chayes, and Brendan Lucier. Maximizing social influence in nearly optimal time. In *Proceedings of the twenty-fifth annual ACM-SIAM symposium on Discrete algorithms*, pages 946–957. SIAM, 2014.

- [13] Adam Breuer, Eric Balkanski, and Yaron Singer. The fast algorithm for submodular maximization. In *International Conference on Machine Learning*, pages 1134–1143. PMLR, 2020.
- [14] Niv Buchbinder and Moran Feldman. Submodular functions maximization problems., 2018.
- [15] Lijie Chen. On the hardness of approximate and exact (bichromatic) maximum inner product. In *33rd Computational Complexity Conference (CCC 2018)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2018.
- [16] Wei Chen, Laks VS Lakshmanan, and Carlos Castillo. Information and influence propagation in social networks. *Synthesis Lectures on Data Management*, 5(4):1–177, 2013.
- [17] Wei Chen and Binghui Peng. On adaptivity gaps of influence maximization under the independent cascade model with full-adoption feedback. In *30th International Symposium on Algorithms and Computation (ISAAC 2019)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2019.
- [18] Wei Chen, Binghui Peng, Grant Schoenebeck, and Biaoshuai Tao. Adaptive greedy versus non-adaptive greedy for influence maximization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 590–597, 2020.
- [19] Wei Chen, Yajun Wang, and Siyu Yang. Efficient influence maximization in social networks. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 199–208, 2009.
- [20] Wei Chen, Yajun Wang, Yang Yuan, and Qinshi Wang. Combinatorial multi-armed bandit and its extension to probabilistically triggered arms. *The Journal of Machine Learning Research*, 17(1):1746–1778, 2016.
- [21] Xiaodong Chen, Guojie Song, Xinran He, and Kunqing Xie. On influential nodes tracking in dynamic social networks. In *Proceedings of the 2015 SIAM International Conference on Data Mining*, pages 613–621. SIAM, 2015.
- [22] Kaito Fujii and Shinsaku Sakaue. Beyond adaptive submodularity: Approximation guarantees of greedy policy with adaptive submodularity ratio. In *International Conference on Machine Learning*, pages 2042–2051. PMLR, 2019.
- [23] Daniel Golovin and Andreas Krause. Adaptive submodularity: Theory and applications in active learning and stochastic optimization. *Journal of Artificial Intelligence Research*, 42:427–486, 2011.
- [24] Monika Henzinger, Sebastian Krinninger, and Danupon Nanongkai. Dynamic approximate all-pairs shortest paths: Breaking the $o(mn)$ barrier and derandomization. *SIAM Journal on Computing*, 45(3):947–1006, 2016.
- [25] Russell Impagliazzo and Ramamohan Paturi. On the complexity of k-sat. *Journal of Computer and System Sciences*, 62(2):367–375, 2001.
- [26] Dimitris Kalimeris, Gal Kaplun, and Yaron Singer. Robust influence maximization for hyperparametric models. In *International Conference on Machine Learning*, pages 3192–3200. PMLR, 2019.
- [27] Ehsan Kazemi, Marko Mitrovic, Morteza Zadimoghaddam, Silvio Lattanzi, and Amin Karbasi. Submodular streaming in all its glory: Tight approximation, minimum memory and low adaptive complexity. In *International Conference on Machine Learning*, pages 3311–3320. PMLR, 2019.
- [28] Ehsan Kazemi, Morteza Zadimoghaddam, and Amin Karbasi. Scalable deletion-robust submodular maximization: Data summarization with privacy and fairness constraints. In *International conference on machine learning*, pages 2544–2553. PMLR, 2018.
- [29] David Kempe, Jon Kleinberg, and Éva 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*, pages 137–146, 2003.

- [30] Sanjeev Khanna and Brendan Lucier. Influence maximization in undirected networks. In *Proceedings of the Twenty-fifth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1482–1496. SIAM, 2014.
- [31] Silvio Lattanzi, Slobodan Mitrović, Ashkan NorouziFard, Jakub Tarnawski, and Morteza Zadimoghaddam. Fully dynamic algorithm for constrained submodular optimization. *Advances in Neural Information Processing Systems*, 2020.
- [32] Jure Leskovec, Andreas Krause, Carlos Guestrin, Christos Faloutsos, Jeanne VanBriesen, and Natalie Glance. Cost-effective outbreak detection in networks. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 420–429, 2007.
- [33] Shuai Li, Fang Kong, Kejie Tang, Qizhi Li, and Wei Chen. Online influence maximization under linear threshold model. *Advances in Neural Information Processing Systems*, 33, 2020.
- [34] Yuchen Li, Ju Fan, Yanhao Wang, and Kian-Lee Tan. Influence maximization on social graphs: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 30(10):1852–1872, 2018.
- [35] David Liben-Nowell and Jon Kleinberg. The link-prediction problem for social networks. *Journal of the American society for information science and technology*, 58(7):1019–1031, 2007.
- [36] Xiaodong Liu, Xiangke Liao, Shanshan Li, Si Zheng, Bin Lin, Jingying Zhang, Lisong Shao, Chenlin Huang, and Liquan Xiao. On the shoulders of giants: incremental influence maximization in evolving social networks. *Complexity*, 2017, 2017.
- [37] Morteza Monemizadeh. Dynamic submodular maximization. *Advances in Neural Information Processing Systems*, 33, 2020.
- [38] Naoto Ohsaka, Takuya Akiba, Yuichi Yoshida, and Ken-ichi Kawarabayashi. Dynamic influence analysis in evolving networks. *Proceedings of the VLDB Endowment*, 9(12):1077–1088, 2016.
- [39] Binghui Peng and Wei Chen. adaptive influence maximization with myopic feedback. *Advances in Neural Information Processing Systems 32 pre-proceedings (NeurIPS 2019)*, 2019.
- [40] Aviad Rubinfeld. Hardness of approximate nearest neighbor search. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1260–1268, 2018.
- [41] Aviad Rubinfeld and Virginia Vassilevska Williams. Seth vs approximation. *ACM SIGACT News*, 50(4):57–76, 2019.
- [42] Grant Schoenebeck and Biaoshuai Tao. Influence maximization on undirected graphs: Toward closing the $(1-1/e)$ gap. *ACM Transactions on Economics and Computation (TEAC)*, 8(4):1–36, 2020.
- [43] Lior Seeman and Yaron Singer. Adaptive seeding in social networks. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pages 459–468. IEEE, 2013.
- [44] Youze Tang, Yanchen Shi, and Xiaokui Xiao. Influence maximization in near-linear time: A martingale approach. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, pages 1539–1554, 2015.
- [45] Youze Tang, Xiaokui Xiao, and Yanchen Shi. Influence maximization: Near-optimal time complexity meets practical efficiency. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*, pages 75–86, 2014.
- [46] Yanhao Wang, Qi Fan, and Yuchen Li. Tan, kian-lee. real-time influence maximization on dynamic social streams.(2017). In *Proceedings of the VLDB Endowment: 43rd International Conference on Very Large Data Bases, Munich, Germany, 2017 August 28-September*, volume 1, pages 805–816, 2017.
- [47] Virginia Vassilevska Williams. On some fine-grained questions in algorithms and complexity. In *Proceedings of the ICM*, volume 3, pages 3431–3472. World Scientific, 2018.

A Probabilistic tools

Lemma A.1 (Chernoff bound, the multiplicative form). *Let $X = \sum_{i=1}^n X_i$, where $X_i \in [0, 1]$ are independent random variables. Let $\mu = \mathbb{E}[X] = \sum_{i=1}^n \mathbb{E}[X_i]$. Then*

1. $\Pr[X \geq (1 + \delta)\mu] \leq \exp(-\delta^2\mu/(2 + \delta)), \forall \delta > 0;$
2. $\Pr[X \leq (1 - \delta)\mu] \leq \exp(-\delta^2\mu/2), \forall 0 < \delta < 1.$

Lemma A.2 (Hoeffding bound). *Let X_1, \dots, X_n denote n independent bounded variables in $[a_i, b_i]$. Let $X = \sum_{i=1}^n X_i$, then we have*

$$\Pr[|X - \mathbb{E}[X]| \geq t] \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

Lemma A.3 (Azuma bound, the multiplicative form). *Let $X_1, \dots, X_n \in [0, 1]$ be real valued random variable. Suppose*

$$\mathbb{E}[X_i | X_1, \dots, X_n] = \mu_i$$

holds for all $i \in [n]$ and let $\mu = \sum_{i=1}^n \mu_i$. Then, we have

$$\Pr\left[\sum_{i=1}^n X_i \geq (1 + \delta)\mu\right] \leq \exp(-\delta^2\mu/(2 + \delta))$$

$$\Pr\left[\sum_{i=1}^n X_i \leq (1 - \delta)\mu\right] \leq \exp(-\delta^2\mu/2).$$

Lemma A.4 (Azuma-Hoeffding bound). *Let X_0, \dots, X_n be a martingale sequence with respect to the filter $F_0 \subseteq F_2 \dots \subseteq F_n$ such that for $Y_i = X_i - X_{i-1}$, $i \in [n]$, we have that $|Y_i| = |X_i - X_{i-1}| \leq c_i$. Then*

$$\Pr[|X_t - Y_0| \geq t] \leq 2 \exp\left(-\frac{t^2}{2 \sum_{i=1}^n c_i^2}\right).$$

B Missing proof from Section 3.1

We first provide a proof of Lemma 3.3, which states the ESTIMATE procedure gives a good estimate on sampling probability.

Proof of Lemma 3.3. Let X_t be the fraction of edges that are checked in the t -th iteration of ESTIMATE (i.e., Line 3 of ESTIMATE). By definition, we have that $\mathbb{E}[X_t] = \text{AVG}$ and $X_t \in [0, 1]$ for all t . Note that K is precisely the minimum value K' such that $\sum_{t=1}^{K'} X_t \geq R$.

Let $K_1 = (1 - \epsilon) \cdot \frac{R}{\text{AVG}}$ and $K_2 = (1 + \epsilon) \cdot \frac{R}{\text{AVG}}$. We want to show that with high probability $K \in [K_1, K_2]$. Note that this event is exactly the intersection of the event $\sum_{t=1}^{K_1} X_t < R$ and the $\sum_{t=1}^{K_2} X_t > R$. For the first one, notice that

$$\mathbb{E}\left[\sum_{t=1}^{K_1} X_t\right] = K_1 \cdot \text{AVG} = (1 - \epsilon)R.$$

By the multiplicative form Chernoff bound, we have

$$\Pr\left[\sum_{t=1}^{K_1} X_t \geq R\right] \leq \Pr\left[\sum_{t=1}^{K_1} X_t \geq (1 - \epsilon) \sum_{t=1}^{K_1} \mathbb{E}[X_t]\right] \leq \exp\left(-\epsilon^2 \sum_{t=1}^{K_1} \mathbb{E}[X_t]/3\right)$$

$$\leq \exp(-\epsilon^2(1 - \epsilon)ck\epsilon^{-2} \log(n/\delta)/3) \leq \delta/n^{ck/6}.$$

For the second event, notice that

$$\mathbb{E}\left[\sum_{t=1}^{K_2} X_t\right] = K_2 \cdot \text{AVG} = (1 + \epsilon)R,$$

and by the multiplicative form Chernoff bound, we have

$$\begin{aligned} \Pr \left[\sum_{t=1}^{K_2} X_t \leq R \right] &< \Pr \left[\sum_{t=1}^{K_2} X_t \leq (1 - \epsilon/2) \sum_{t=1}^{K_2} \mathbb{E}[X_t] \right] \leq \exp \left(-\epsilon^2 \sum_{t=1}^{K_2} \mathbb{E}[X_t]/8 \right) \\ &\leq \exp(-\epsilon^2(1 + \epsilon)ck\epsilon^{-2} \log(n/\delta)/8) \leq \delta/n^{ck/8}. \end{aligned}$$

Taking an union bound, we conclude that with probability at least $1 - \frac{2\delta}{n^{ck/8}}$, $K \in [K_1, K_2]$. \square

We next prove Lemma 3.4, which states the initial steps for building \mathcal{H}_{est} and \mathcal{H}_{cv} is not large.

Proof of Lemma 3.4. We prove that with $1 - \frac{2\delta}{n^{ck/8}}$, the total number of steps for building \mathcal{H}_{est} are within $[(1 - \epsilon)^2 R, (1 + \epsilon)^2 R]$, the Lemma then follows by an union bound. Let Y_i ($i \in [n_0]$) be the fraction of edges checked when sampling the RR set R_i of the i -th node. Notice that if with probability $1 - p$, the algorithm won't sample R_i , and we take $Y_i = 0$ for this case. We have that $Y_i \in [0, 1]$ and

$$\mathbb{E} \left[\sum_{i=1}^{n_0} Y_i \right] = \frac{K}{n_0} \cdot n_0 \text{AVG} = K \cdot \text{AVG} \in [(1 - \epsilon)R, (1 + \epsilon)R]. \quad (3)$$

The last step holds due to Lemma 3.3.

Consequently, by the multiplicative form Chernoff bound, we have

$$\begin{aligned} \Pr \left[\sum_{i=1}^{n_0} Y_i \geq (1 + \epsilon)^2 R \right] &\leq \Pr \left[\sum_{i=1}^{n_0} Y_i \geq (1 + \epsilon) \sum_{i=1}^{n_0} \mathbb{E}[Y_i] \right] \leq \exp \left(-\epsilon^2 \sum_{i=1}^{n_0} \mathbb{E}[Y_i]/3 \right) \\ &\leq e^{-\epsilon^2(1-\epsilon)c\epsilon^{-2}k \log(n/\delta)/3} \leq \delta/n^{ck/6}. \end{aligned}$$

The first step and third follows from Eq. (3). We assume $\epsilon < 1/2$ in the last step.

Similarly, one has

$$\begin{aligned} \Pr \left[\sum_{i=1}^{n_0} Y_i \leq (1 - \epsilon)^2 R \right] &\leq \Pr \left[\sum_{i=1}^{n_0} Y_i \leq (1 - \epsilon) \sum_{i=1}^{n_0} \mathbb{E}[Y_i] \right] \leq \exp \left(-\epsilon^2 \sum_{i=1}^{n_0} \mathbb{E}[Y_i]/2 \right) \\ &\leq e^{-\epsilon^2(1-\epsilon)c\epsilon^{-2}k \log(n/\delta)/2} \leq \delta/n^{ck/8}. \end{aligned}$$

We conclude the proof here. \square

We next provide a proof for Lemma 3.5, which states that the total number of steps of constructing \mathcal{H}_{est} is around its expectation.

Proof of Lemma 3.5. We know that the random variable $Z_j \in [0, 1]$ since the total number edge checked for one RR set is at most $m \leq 2m_0$. Moreover, $\{Z_j\}_{j \geq 1}$ forms a martingale. Since J is the first number such that $\sum_{j=1}^J 2m_0 Z_j \geq 16Rm_0$, i.e., $\sum_{j=1}^J Z_j \geq 8R$, the event $J \in [J_1, J_2]$ is equivalent to the intersection of event $\sum_{j=1}^{J_1} Z_j > 8R$ and $\sum_{j=1}^{J_2} Z_j < 8R$. We bound the probability separately.

For the first event, we have

$$(1 - \epsilon)8R \leq \sum_{j=1}^{J_1} \mathbb{E}[Z_j] \leq (1 - \epsilon)8R + 1.$$

By the multiplicative form Azuma bound, we have

$$\begin{aligned} \Pr \left[\sum_{j=1}^{J_1} Z_j > 8R \right] &\leq \Pr \left[\sum_{j=1}^{J_1} Z_j \geq (1 + \epsilon/2) \sum_{j=1}^{J_1} \mathbb{E}[Z_j] \right] \leq \exp \left(-\epsilon^2 \sum_{j=1}^{J_1} \mathbb{E}[Z_j]/12 \right) \\ &\leq \exp(-\epsilon^2(1 - \epsilon)8c\epsilon^{-2}k \log(n/\delta)/12) \leq \frac{\delta}{n^{ck/3}}. \end{aligned} \quad (4)$$

Similarly, for the second event, we have that

$$(1 + \epsilon)8R - 1 \leq \sum_{j=1}^{J_2} \mathbb{E}[Z_j] \leq (1 + \epsilon)8R.$$

By the multiplicative form Azuma bound, we have

$$\begin{aligned} \Pr \left[\sum_{j=1}^{J_2} Z_j < 8R \right] &\leq \Pr \left[\sum_{j=1}^{J_2} Z_j \leq (1 - \epsilon/2) \sum_{j=1}^{J_2} \mathbb{E}[Z_j] \right] \leq \exp \left(-\epsilon^2 \sum_{j=1}^{J_2} \mathbb{E}[Z_j]/8 \right) \\ &\leq \exp(-\epsilon^2 8c\epsilon^{-2}k \log(n/\delta)/8) \leq \frac{\delta}{n^{ck}} \end{aligned} \quad (5)$$

Combine the Eq. (4) and Eq. (5) with union bound, we conclude that probability at least $1 - \frac{2\delta}{n^{ck/3}}$, the total number step for \mathcal{H}_{est} satisfies $J_1 \leq J \leq J_2$. \square

The total number of edge checked for \mathcal{H}_{cv} is bounded by $16(1 + \epsilon)^2 Rm_0$ with high probability, as stated in Lemma 3.6

Proof of Lemma 3.6. Similar as Lemma 3.5, consider the t -th edge to arrive, let $W_{t,i} \cdot 2m_0$ be the number of edges checked by INSERT-EDGE when it augments the RR set of the i -th node on \mathcal{H}_{cv} . Again, if INSERT-NODE and BUILDGRAPH did not sample the i -th node and its RR set, we simply have $W_{t,i} = 0$. We also write W_j to denote the j -th random variable in the sequence $W_{0,1}, \dots, W_{0,n_0}, W_{1,1}, \dots$, and it also forms a martingale. Recall J is the total number of steps executed by INSERT EDGE on constructing \mathcal{H}_{est} , and it is also the time where we stop augmenting \mathcal{H}_{cv} .

We know $W_j \in [0, 1]$ since $m \leq 2m_0$ and $\{W_j\}_{j \geq 1}$ forms a martingale. Since we sample \mathcal{H}_{est} and \mathcal{H}_{cv} separately, the stopping time is irrelevant to the realization of W_j . Condition on the event of Lemma 3.5, i.e. $J \in [J_1, J_2]$, we know that

$$(1 - \epsilon)8R \leq \sum_{j=1}^{J_1} \mathbb{E}[W_j] \leq \sum_{j=1}^J \mathbb{E}[W_j] \leq \sum_{j=1}^{J_2} \mathbb{E}[W_j] \leq (1 + \epsilon)8R.$$

Hence, by Azuma bound, we have

$$\begin{aligned} \Pr \left[\sum_{j=1}^J W_j \cdot 2m_0 \geq 16(1 + \epsilon)^2 Rm_0 \right] &= \Pr \left[\sum_{j=1}^J W_j \geq (1 + \epsilon)^2 8R \right] \\ &= \Pr \left[\sum_{i=1}^J W_j \geq (1 + \epsilon) \sum_{i=1}^J \mathbb{E}[Z_j] \right] \\ &\leq \exp \left(-\epsilon^2 \sum_{j=1}^J \mathbb{E}[W_j]/3 \right) \\ &\leq \exp(-\epsilon^2 8(1 - \epsilon)c\epsilon^{-2}k \log(\delta/n)/3) \\ &\leq \frac{\delta}{n^{ck}}. \end{aligned}$$

This concludes the proof. \square

We provide the proof of Lemma 3.8, which asserts that the normalized coverage function $f_{\text{cv},t}$ is an unbiased estimator on the influence spread function.

Proof of Lemma 3.8. The expected influence spread of a set S at time step t satisfies

$$\sigma_t(S) = \sum_{v \in V_t} \mathbb{E}[|S \cap R_{v,t}|] = \sum_{v \in V_t} \frac{n_0}{K} \mathbb{E}[x_{t,v,S}] = \mathbb{E}[f_{\text{cv},t}(S)].$$

The first step follows from the definition of influence spread function σ_t , the second step follows from our sampling process, the definition of $x_{t,v,S}$ and the fact that we sample the RR set of node v with probability $p = \frac{K}{n_0}$. \square

Next, we prove Lemma 3.9.

Proof of Lemma 3.9. Fix a fixed time step t . For any node set $S \subseteq V_t$, $|S| \leq k$, by Lemma 3.8, we have $\mathbb{E}[f_{cv,t}(S)] = \sigma_t(S) \leq \text{OPT}_t$. For convenience, we write $\mathbb{E}[f_{cv,t}(S)] = \lambda \text{OPT}_t$, and $\lambda \in [0, 1]$. Hence, we have

$$\begin{aligned} \Pr [f_{cv,t}(S) > \mathbb{E}[f_{cv,t}(S)] + \epsilon \text{OPT}_t] &= \Pr \left[f_{cv,t}(S) > \left(1 + \frac{\epsilon}{\lambda}\right) \mathbb{E}[f_{cv,t}(S)] \right] \\ &= \Pr \left[\sum_{v \in V_t} x_{t,v,S} > \left(1 + \frac{\epsilon}{\lambda}\right) \sum_{v \in V_t} \mathbb{E}[x_{t,v,S}] \right]. \end{aligned} \quad (6)$$

We divide into two cases. First, suppose $\lambda > \epsilon$. Since $x_{t,v,S} \in \{0, 1\}$ and they are independent, by the multiplicative form Chernoff bound, we have

$$\Pr \left[\sum_{v \in V_t} x_{t,v,S} > \left(1 + \frac{\epsilon}{\lambda}\right) \sum_{v \in V_t} \mathbb{E}[x_{t,v,S}] \right] \leq \exp \left(-\epsilon^2 \mathbb{E} \left[\sum_{v \in V_t} x_{t,v,S} \right] / 3\lambda^2 \right).$$

The exponent obeys

$$\begin{aligned} \epsilon^2 \sum_{v \in V_t} \mathbb{E}[x_{t,v,S}] / 3\lambda^2 &= \epsilon^2 \frac{K}{n_0} \lambda \text{OPT}_t / 3\lambda^2 \geq \epsilon^2 \frac{K}{3n_0} \text{OPT}_t \geq \epsilon^2 (1 - \epsilon) \frac{R}{3n_0 \text{AVG}} \text{OPT}_t \\ &\geq \epsilon^2 (1 - \epsilon) R / 3 \geq ck \log(n/\delta) / 6. \end{aligned} \quad (7)$$

where the first step follows from

$$\frac{n_0}{K} \sum_{v \in V_t} \mathbb{E}[x_{t,v,S}] = \mathbb{E}[f_{cv,t}(S)] = \lambda \text{OPT}_t, \quad (8)$$

the second step follows from $\lambda < 1$, the third comes from the condition of Lemma 3.3, i.e., $K \in [(1 - \epsilon) \cdot \frac{R}{\text{AVG}}, (1 + \epsilon) \cdot \frac{R}{\text{AVG}}]$. The fourth step follows from Lemma B.1 and the monotonicity of OPT, i.e.,

$$\text{AVG} \leq \frac{\text{OPT}_0}{n_0} \leq \frac{\text{OPT}_t}{n_0}. \quad (9)$$

Hence, when $\epsilon < \lambda$, one has

$$\Pr [f_{cv,t}(S) > \mathbb{E}[f_{cv,t}(S)] + \epsilon \text{OPT}_t] \leq \frac{\delta}{n^{ck/6}}.$$

Next, suppose $\lambda \leq \epsilon$, then $\epsilon/\lambda \geq 1$. The multiplicative form Chernoff bound

$$\Pr \left[\sum_{v \in V_t} x_{t,v,S} > \left(1 + \frac{\epsilon}{\lambda}\right) \sum_{v \in V_t} \mathbb{E}[x_{t,v,S}] \right] \leq \exp \left(-\epsilon \sum_{v \in V_t} \mathbb{E}[x_{t,v,S}] / 3\lambda \right)$$

The exponent satisfies

$$\epsilon \sum_{v \in V_t} \mathbb{E}[x_{t,v,S}] / 3\lambda = \epsilon \cdot \frac{K}{3n_0} \text{OPT}_t \geq \epsilon (1 - \epsilon) \frac{R}{3n_0 \text{AVG}} \text{OPT}_t \geq \epsilon (1 - \epsilon) R / 3 \geq ck \log(n/\delta) / 6.$$

where the first step follows from Eq. (8), the second step follows from the condition of Lemma 3.3, i.e., $K \in [(1 - \epsilon) \cdot \frac{R}{\text{AVG}}, (1 + \epsilon) \cdot \frac{R}{\text{AVG}}]$. The third step follows Eq. (9). Hence, when $\lambda \geq \epsilon$, one also has

$$\Pr [f_{cv,t}(S) > \mathbb{E}[f_{cv,t}(S)] + \epsilon \text{OPT}_t] \leq \frac{\delta}{n^{ck/6}}.$$

This proves the Eq. (1). The proof of Eq. (2) is similar. In particular, we have

$$\begin{aligned} \Pr [f_{cv,t}(S) < \mathbb{E}[f_{cv,t}(S)] - \epsilon \text{OPT}_t] &= \Pr \left[f_{cv,t}(S) < \left(1 - \frac{\epsilon}{\lambda}\right) \mathbb{E}[f_{cv,t}(S)] \right] \\ &= \Pr \left[\sum_{v \in V_t} x_{t,v,S} < \left(1 - \frac{\epsilon}{\lambda}\right) \sum_{v \in V_t} \mathbb{E}[x_{t,v,S}] \right]. \end{aligned}$$

It suffices to consider the case $\epsilon < \lambda$. The multiplicative form Chernoff bound implies

$$\begin{aligned} \Pr \left[\sum_{v \in V_t} x_{t,v,S} < \left(1 - \frac{\epsilon}{\lambda}\right) \sum_{v \in V_t} \mathbb{E}[x_{t,v,S}] \right] &\leq \exp \left(-\epsilon^2 \sum_{v \in V_t} \mathbb{E}[x_{t,v,S}] / 2\lambda^2 \right) \\ &\leq \exp(\epsilon^2(1 - \epsilon)R/2) \leq \frac{\delta}{n^{ck/4}}. \end{aligned}$$

The second step follows from Eq. (7). We conclude the proof here. \square

Lemma 3.10 indicates pointwise approximation is sufficient to carry over the approximation ratio between two problems. We provide a proof here.

Proof of Lemma 3.10. By an union bound over all sets of size at most k , we know that

$$\Pr_{g \sim D} [\exists S, |S| \leq k, |f(S) - g(S)| - \gamma] \leq n^k \frac{\delta}{n^{ck}} \leq \frac{\delta}{n^{ck/2}}.$$

Under the above event, we know that

$$f(S_g) \geq g(S_g) - \delta \geq \alpha \max_{S \subseteq V, |S| \leq k} g(S) - \gamma \geq \alpha \max_{S \subseteq V, |S| \leq k} f(S) - 2\gamma.$$

This concludes the proof. \square

Given an influence graph, let $\text{AVG} \cdot m$ be the expected number of steps taken by random sampling a RR set, and let OPT be the maximum (expected) influence spread of a seed set of size at most k , then one has

Lemma B.1 (Claim 3.3 in [12]). $\text{AVG} \leq \frac{\text{OPT}}{n}$.

We provide a proof for completeness.

Proof. Given a node v and an edge (u, w) , let $x_{v,e} = 1$ if the edge e is checked when one samples the RR set of node v . Then we have

$$\text{AVG} \cdot m = \frac{1}{n} \sum_{v \in V} \sum_{e \in E} \mathbb{E}[x_{v,e}] = \frac{1}{n} \sum_{e=(u,w) \in E} \mathbb{E}[|v : w \in R_v|].$$

Hence, we have

$$\text{AVG} \cdot m = \frac{1}{n} \sum_{e=(u,w) \in E} \mathbb{E}[|v : w \in R_v|] = \frac{1}{n} \sum_{e=(u,v) \in E} \sigma(v) \leq \frac{1}{n} \sum_{e=(u,v) \in E} \text{OPT} = \frac{m}{n} \text{OPT}.$$

We conclude the proof. \square

We provide proof for Lemma 3.11, which asserts with high probability, there are at most $O(\log n)$ iterations within a phase.

Proof of Lemma 3.11. For any $t \geq 0$, let n_t, m_t be the number of nodes and edges at the beginning the t -th iteration, and let $\text{AVG}_t \cdot m_t$ be the average steps of sampling a random RR set. We can assume there is at least one edge in the graph, and therefore, $\text{AVG}_0 \geq \frac{1}{n_0 m_0}$. Inside the phase, we must have $n_0 \leq n_t \leq 2n_0$ and $m_0 \leq m_t \leq 2m_0$ holds for any t . We prove that $\text{AVG}_{t+1} \geq 2 \text{AVG}_t$

holds with high probability. Notice that $\text{AVG}_t \leq 1$, this means the algorithm restarts for at most $O(\log n)$ times. By Lemma 3.3, the sample size K_t at the beginning of t -th iteration obeys

$$(1 - \epsilon) \cdot \frac{R}{\text{AVG}_t} \leq K_t \leq (1 + \epsilon) \cdot \frac{R}{\text{AVG}_t}. \quad (10)$$

with probability at least $1 - \frac{2\delta}{n^{ck/8}}$.

On the other side, for the t -th iteration, define $J_t, J_{t,1}, J_{t,2}$ similarly as Lemma 3.5. With probability at least $1 - \frac{2\delta}{n^{ck/8}}$, we have $J_t \in [J_{t,1}, J_{t,2}]$, and this indicates

$$8(1 - \epsilon)R \leq \sum_{j=1}^{J_{t,1}} \mathbb{E}[Z_j] \leq \sum_{j=1}^{J_t} \mathbb{E}[Z_j] \leq \sum_{j=1}^{J_{t,2}} \mathbb{E}[Z_j] \leq 8(1 + \epsilon)R. \quad (11)$$

The first and last step follow from the definition of $J_{t,1}$ and $J_{t,2}$, the second and third step follow from $J_t \in [J_{t,1}, J_{t,2}]$. Moreover, we also know that

$$\sum_{j=1}^{J_t} \mathbb{E}[Z_j] = \frac{K_t}{n_t} n_{t+1} \text{AVG}_{t+1}. \quad (12)$$

as we include the RR set of each n_t node with probability $\frac{K_t}{n_t}$. Therefore, we have

$$2K_t \text{AVG}_{t+1} \geq \frac{K_t}{n_t} n_{t+1} \text{AVG}_{t+1} \geq 8(1 - \epsilon)R \geq \frac{8(1 - \epsilon)}{1 + \epsilon} K_t \text{AVG}_t \geq 4K_t \text{AVG}_t$$

The first step follows from $n_{t+1} \leq 2n_0 \leq 2n_t$. The second step follows from Eq. (11) and Eq. (12). The third step comes from Eq. (10), and we assume $\epsilon < 1/3$ in the last step. Hence, we have proved with probability at least $1 - \frac{4\delta}{n^{ck/8}}$, $\text{AVG}_{t+1} \geq 2 \text{AVG}_t$. Taking an union bound over t and combining the fact that $\text{AVG}_0 \geq \frac{1}{n_0 m_0}$ and $\text{AVG}_t \leq 1$, we conclude with probability $1 - \frac{4\delta}{n^{ck/16}}$, the algorithm restarts at most $O(\log n)$ times within each phase. \square

We wrap up the proof of Lemma 3.2

Proof of Lemma 3.2. We first prove the correctness of the algorithm. By Lemma 3.8 and Lemma 3.9, we know that at any time step t , the *normalized* coverage function $f_{\text{cv},t}$ defined on \mathcal{H}_{cv} gives a good approximation on the influence spread function. In particular, we have that with probability at least $1 - \frac{4\delta}{n^{ck/8}}$, one has

$$|f_{\text{cv},t}(S) - \mathbb{E}[f_{\text{cv},t}(S)]| = |f_{\text{cv},t}(S) - \sigma_t(S)| \leq \epsilon \text{OPT}_t$$

Combining with Lemma 3.10, suppose we can solve the dynamic MAX-k coverage with approximation α , then our algorithm maintains a solution set S_t that satisfies

$$f(S) \geq (\alpha - 2\epsilon) \text{OPT}_t.$$

We next focus on the amortized running time. Since the number of edges and nodes can only doubled for most $O(\log n)$ times, there are $O(\log n)$ phases. While within one phase, by Lemma 3.11, with probability $1 - \frac{4\delta}{n^{ck/16}}$, our algorithm restarts for at most $O(\log n)$ times. Each time our algorithm restarts, it invokes the ESTIMATE procedure for once. This step takes $Rm_0 = m_0 ck \epsilon^{-2} \log n$ steps in total and has $O(k \epsilon^{-2} \log n)$ amortized time per update. The algorithm constructs \mathcal{H}_{est} and \mathcal{H}_{cv} , we calculate their cost separately. For constructing \mathcal{H}_{est} , our algorithm takes at most $16R_0 = 16m_0 ck \epsilon^{-2} \log n$ steps in total and has $O(k \epsilon^{-2} \log n)$ amortized time per update. For the construction of \mathcal{H}_{cv} , by Lemma 3.4 and Lemma 3.6, with probability at least $1 - \frac{9\delta}{n^{ck/8}}$, it takes less than $16(1 + \epsilon)^2 Rm_0 \leq 64m_0 ck \epsilon^{-2} \log n$ steps in total and $O(k \epsilon^{-2} \log n)$ amortized time per updates. Note that our algorithm not only needs to construct \mathcal{H}_{cv} , but also needs to maintains a set that has the (approximately) maximum coverage on \mathcal{H}_{cv} . This reduces to a dynamic MAX-k coverage problem, which by our assumption, can be solved in amortized running time of \mathcal{T}_{mk} . Taking an union bound over all steps and fix the constant c to be greater than 24, we conclude with probability at least $1 - \delta$, the overall amortized running time per update is bounded by

$$\log n \cdot \log n \cdot (k \epsilon^{-2} \log n + k \epsilon^{-2} \log n + k \epsilon^{-2} \log n \cdot \mathcal{T}_{\text{mk}}) \leq O(\mathcal{T}_{\text{mk}} k \epsilon^{-2} (\log n)^3).$$

This concludes the proof. \square

C Missing proof from Section 3.2

Lemma 3.13 ensures the approximation guarantee of the algorithm, we provide a detailed proof here.

Proof of Lemma 3.13. Fix a time step t , let OPT denote the value of the optimal solution, i.e. $\text{OPT} = \max_{S, |S| \leq k} f_t(S)$. For ease of notation, we drop the subscript t in the rest of the proof. There exists an index $i \in I$ such that

$$(1 + \epsilon)^i = \text{OPT}_i \leq \text{OPT} < \text{OPT}_{i+1} = (1 + \epsilon)^{i+1}.$$

We prove the i -th thread outputs a good solution set S_i .

First, suppose $|S_i| = k$. Let $s_{i,j}$ be the j -th element added to the set S_i , and denote $S_{i,j} = \{s_{i,1}, \dots, s_{i,j}\}$, $j \in [k]$. Our algorithm guarantees that

$$f(S_{i,j+1}) - f(S_{i,j}) \geq \frac{1}{k}(\text{OPT}_i - f(S_{i,j}))$$

Then, we have that

$$\begin{aligned} \text{OPT}_i - f(S_{i,k}) &= \text{OPT}_i - f(S_{i,k-1}) + f(S_{i,k-1}) - f(S_{i,k}) \\ &\leq \text{OPT}_i - f(S_{i,k-1}) - \frac{1}{k}(\text{OPT}_i - f(S_{i,k-1})) \\ &\leq \left(1 - \frac{1}{k}\right) (\text{OPT}_i - f(S_{i,k-1})) \\ &\quad \vdots \\ &\leq \left(1 - \frac{1}{k}\right)^k (\text{OPT}_i - f(\emptyset)) \\ &= \left(1 - \frac{1}{k}\right)^k \text{OPT}_i, \end{aligned}$$

and therefore,

$$\begin{aligned} f(S_i) &\geq \left(1 - \left(1 - \frac{1}{k}\right)^k\right) \text{OPT}_i \geq \left(1 - \frac{1}{e}\right) \text{OPT}_i \\ &\geq \left(1 - \frac{1}{e}\right) (1 + \epsilon)^{-1} \text{OPT} \geq \left(1 - \frac{1}{e} - \epsilon\right) \text{OPT}. \end{aligned}$$

On the otherside, if $|S_i| < k$, then we prove $f(S_i) \geq \text{OPT}_i \geq (1 + \epsilon)^{-1} \text{OPT}$. We prove by contradiction and assume $f(S_i) < \text{OPT}_i$ for now. Let the optimal solution be $O = \{o_1, \dots, o_k\}$. Then we claim that

$$f_{S_i}(o) < \frac{1}{k}(\text{OPT}_i - f(S_i)) \tag{13}$$

holds for all $o \in O$. The reason is that (i) if $o \in S_i$, then $f_{S_i}(o) = 0 < \frac{1}{k}(\text{OPT}_i - f(S_i))$. If $o \notin O$, since $|S_i| < k$, the above is guaranteed by our algorithm. Hence, we have

$$\text{OPT} = f(O) \leq f(S_i) + f_{S_i}(O) \leq f(S_i) + \sum_{j=1}^k f_{S_i}(o_j) < f(S_i) + k \cdot \frac{1}{k}(\text{OPT}_i - f(S_i)) = \text{OPT}_i.$$

The second step holds by monotonicity, the third step holds by submodularity and the fourth step holds by Eq. (13). This comes to a contradiction. Hence, we proved $f(S_i) \geq (1 - 1/e - \epsilon) \text{OPT}$ in both cases. \square

We next prove Lemma 3.14, which analyses the amortized running time

Proof of Lemma 3.14. It suffices to prove that for each thread $i \in I$, the amortized running time is $O(1)$. We specify some implementation details. For any set S , let $N(S)$ denote the all neighbors of S . We maintain a set V_R^i that includes all nodes covered by the current set S_i , i.e. $V_R^i = N(S_i)$. We also maintain a set V_u^i for each node u , which contains all element covered by node u in $V_R \setminus V_R^i$, i.e., $V_u^i = |N(u) \setminus V_R^i|$. Finally, we also maintain an order (on cardinality) over the set V_u^i . This is used in the REVOKE procedure, where we retrieve the node u with the maximum $|V_u^i|$ and compare it with $\frac{\text{OPT}_i - |V_R^i|}{k}$. We are going to prove that we can maintain these data structures and perform all necessary operations in $O(1)$ amortized time.

Let $P(t)$ denote the number of operations performed on the i -th thread up to time t . For each edge $e = (u, v)$, let X_e denote whether the edge e is covered by V_R^i , i.e. $X_e = 1$ if $v \in V_R^i$ and $X_e = 0$ otherwise. Similarly, for each node $v \in V_R$, let Y_v denote whether node v is included in V_R^i , i.e., $Y_v = 1$ if $v \in V_R^i$ and $Y_v = 0$ otherwise. For each node $u \in V_L$, let Z_u denote whether node u is included in S_i , i.e., $Z_u = 1$ if $u \in S_i$ and $Z_u = 0$ otherwise. Define the potential function $\Phi : t \rightarrow \mathbb{R}^+$:

$$\Phi(t) = 2|V_t| + 2|E_t| + 2 \sum_{e \in E_t} X_e + \sum_{v \in V_R} Y_v + 2 \sum_{u \in V_L} Z_u.$$

Our goal is to show $P(t) \leq \Phi(t)$. This is sufficient for our purpose as one can easily show $\Phi(t) \leq 5t$. We prove the claim by induction. The claim holds trivially for the base case $t = 0$. We next assume $t > 0$ and consider the time step t . If a new node arrives, then we have that $P(t) = P(t-1) + 1$. Since $|V_t| = |V_{t-1}| + 1$ and other terms of Φ won't decrease, we have $\Phi(t) \leq \Phi(t-1) + 1$. Suppose a new edge $e = (u, v)$ arrives. (1) If $v \in V_R^i$, that is, the node v has already been covered. Then $P(t) = P(t-1) + 2$ since we don't perform any additional operations. We also have $\Phi(t) = \Phi(t-1) + 2$ as $|E_t| = |E_{t-1}| + 2$ and the other term remains unchanged. (2.1) If $v \notin V_R^i$ and $u \in S_i$, then we need to expand the set $V_t \leftarrow V_{t-1} \cup \{v\}$ (one unit operation), delete node v from V_u^i if $v \in V_u^i$ ($\sum_{v \in V_R} |V_u^i \cap \{v\}|$ operations) and maintains the order of $\{V_{v'}^i\}_{v' \in V_R}$ (at most $\sum_{v \in V_R} |V_v^i \cap \{u\}|$ operations). Meanwhile, we have $|E_t| = |E_{t-1}| + 1$ and the term $2 \sum_{e \in E_t} (1 - X_e)$ would increase for $2 \sum_{v \in V_R} |V_v^i \cap \{u\}|$, as these edges change from uncovered to covered. Hence, we still have $P(t) - P(t-1) \leq \Phi(t) - \Phi(t-1)$. (2.2) If $v \notin V_R^i$ and $u \notin S_i$. This may only cause two unit operations if there is no node u satisfies $|V_u^i| \geq \frac{\text{OPT}_i - |V_R^i|}{k}$. This time, we have $P_t = P_t + 2$ $\Phi(t) = \Phi(t-1) + 2$ as $|E_t| = |E_{t-1}| + 1$. On the other side, if there exists some node u with large marginal. We need to add u to V_R^i (1 unit operation), add nodes in V_u^i to V_R^i ($|V_u^i|$ operations) and removes nodes in V_u^i from all other set $V_{u'}^i$ ($\sum_{v \in V_R} |V_u^i \cap V_{v'}^i|$ operations in total). We also want to maintain an order on $\{V_{v'}^i\}_{v' \in V_R}$, and this takes less than $\sum_{v \in V} |V_u^i \cap V_v^i|$ operations in total. Meanwhile, for the potential function, the term $2 \sum_{e \in E_t} X_e$ increases for $2 \sum_{v \in V} |V_u^i \cap V_v^i|$, as this the number of edges change from uncovered to covered. The term $\sum_{v \in V_R} Y_v$ increases for $(|V_u^i|)$ and term $2 \sum_{u \in V_L} Z_u$ will also increase by 2, as we augment the set S_i by 1. Hence, we still have $P(t) - P(t-1) \leq \Phi(t) - \Phi(t-1)$ in this case. Finally, we note the that algorithm may call REVOKE multiple times upon the arrival of a new edge, and for each call, we call do perform similary analysis as (2.2). Hence, we conclude that $P(t+1) - P(t) \leq \Phi(t) - \Phi(t-1)$ holds for all t . We conclude the proof here. \square

D Missing proof from Section 4

Proof of Theorem 4.3. We assume $k = 1$ in our reduction. Let $m = n^{o(1)}$, $t = 2^{(\log n)^{1-o(1)}}$. Given an instance \mathcal{A}, \mathcal{B} of the problem in Theorem 4.2, we reduce it to the dynamic influence maximization problem. We assume $|B_\tau| \geq t$ for all $\tau \in [n]$ as we can duplicate the ground element for t times. Consider the following influence graph $G = (V, E, p)$, where the node set V are partitioned into $V = V_1 \cup V_2 \cup V_3$. There are n nodes in V_1 , denoted as $v_{1,1}, \dots, v_{1,n}$. Intuitively, the i -th node corresponds to the set $A_i \in \mathcal{A}$. The set V_2 contains m nodes, denoted as $v_{2,1}, \dots, v_{2,m}$. Intuitively, they correspond to the ground set $[m]$. For any node $v_{1,i} \in V_1$ and node $v_{2,j} \in V_2$, there is a directed edge from $v_{1,i}$ to $v_{2,j}$, iff the j -th element is contained in the set A_i . We associate the influence probability 1 to every edge between V_1 and V_2 . The set V_3 contains $m^2 t$ nodes, denoted as $\{v_{3,j,\ell}\}_{j \in [m], \ell \in [mt]}$. There is a directed edge with influence probability 1 from node $v_{2,j}$ to $v_{3,j,\ell}$, for each $j \in [m], \ell \in [mt]$.

Consider the following update sequence of the DIM problem. The graph G is loaded first and then followed by n consecutive epochs. In the τ -th epoch, all edges between V_2 and V_3 are deleted, and for each $j \in B_\tau$, we add back the edge between $v_{2,j}$ and $v_{3,j,\ell}$ for all $\ell \in [mt]$.

We first calculate the total number of updates. It takes $n + m + mt = n^{1+o(1)}$ steps to insert all nodes in V and takes at most $mn + mt = n^{1+o(1)}$ to insert all edges in E . We delete/insert at most m^2t edges in each epoch, and since there are n epochs, the total number operations are bounded by $nm^2t = n^{1+o(1)}$. Hence, the total number of updates is at most $n^{1+o(1)}$.

Suppose on the contrary, there exists an algorithm for DIM problem that achieves $2/t$ -approximation in $n^{1-\epsilon}$ time, we then derive a contradiction to SETH. Under the above reduction, we output YES, if for some epoch $\tau \in [n]$, the DIM algorithm outputs a solution with influence spread greater than $2m|B_\tau|$. We output NO otherwise. Note the influence of a node can be computed in $m = n^{o(1)}$ times.

Completeness. Suppose there exists $A_i \in \mathcal{A}$, $B_\tau \in \mathcal{B}$ such that $B_\tau \subseteq A_i$. Then in the τ -th epoch, by taking node $v_{1,i}$ in the seed set, the influence spread at least $(mt + 1)|B_\tau| + 1$. Since the DIM algorithm gives $2/t$ -approximation, the influence is greater than $2m|B_\tau|$ in this case. Hence, we indeed output YES.

Soundness. Suppose $|A_i \cap B_\tau| < |B_\tau|/t$ for any $i, \tau \in [n]$, then we prove the influence spread is no more than $2m|B_\tau|$ for any epoch. This is clearly true for nodes in V_2 and V_3 , as their influence is no more $mt + 1 < 2|B_\tau|m$. Here, we use the fact that $|B_\tau| \geq t$. For nodes in V_1 , since the intersection of A_i and B_τ is less than $|B_\tau|/t$, and a node $v_{2,j} \in V_2$ has influence $1 + mt$ if $j \in B_\tau$ and it has influence 1 otherwise. We conclude for any node $v_{1,i}$, its influence is at most

$$1 + m + \frac{1}{t}|B_\tau|mt = 1 + m + |B_\tau|m \leq 2|B_\tau|m.$$

Hence, we output NO in this case.

In summary, the reduced DIM requires $n^{1+o(1)}$ updates and queries and it gives an answer for the problem in Theorem 4.2. Hence, we conclude under SETH, there is no $2/t$ -approximation algorithm unless the amortized running time is $n^{1-\epsilon}$. \square

Proof of Theorem 4.4. The reduction is similar to the one in Theorem 4.4. Let $m = n^{o(1)}$, $t = 2^{(\log n)^{1-o(1)}}$, $k = 1$. Given an instance \mathcal{A}, \mathcal{B} of the problem in Theorem 4.2, we assume $|B_\tau| \geq t$ for all $\tau \in [n]$. Consider the following influence graph $G = (V, E, w)$, where the node set V are partitioned into $V = V_1 \cup V_2 \cup V_3 \cup V_4$. There are n nodes in V_1 , denoted as $v_{1,1}, \dots, v_{1,n}$ and there are m nodes in V_2 , denoted as $v_{2,1}, \dots, v_{2,m}$. For any node $v_{1,i} \in V_1$ and node $v_{2,j} \in V_2$, there is a directed edge from $v_{1,i}$ to $v_{2,j}$, iff the j -th element is contained in the set A_1 . We associate the weight to be $1/n$ to every edge between V_1 and V_2 . The set V_3 contains m^2t nodes, denoted as $\{v_{3,j,\ell}\}_{j \in [m], \ell \in [mt]}$. There is a directed edge with weight 1 from node $v_{2,j}$ to $v_{3,j,\ell}$, for each $j \in [m], \ell \in [mt]$. The set V_4 contains nm^2t nodes, denoted as $\{v_{4,j,\ell,b}\}_{j \in [m], \ell \in [mt], b \in [n]}$. There is a directed edge with weight 1 from node $v_{3,j,\ell}$ to $v_{4,j,\ell,b}$, for each $j \in [m], \ell \in [mt], b \in [n]$. We assume the prescribed set is V_1 , that is, we are only allowed to select seeds from V_1 .

We use the same update sequence of the DIM problem. The graph G is loaded first and then followed by n consecutive epochs. In the τ -th epoch, all edges between V_2 and V_3 are deleted, and for each $j \in B_\tau$, we add back the edge between $v_{2,j}$ and $v_{3,j,\ell}$ for all $\ell \in [mt]$.

The total number of updates is still at most $n^{1+o(1)}$, as the total number of edges between V_3 and V_4 is at most $nmt^2 = n^{1+o(1)}$ and we only insert them once. Suppose on the contrary, there exists an algorithm for DIM problem that achieves $2/t$ -approximation in $n^{1-\epsilon}$ time, we then derive a contradiction to SETH. Under the above reduction, we output YES, if for some epoch $\tau \in [n]$, the DIM algorithm outputs a solution with influence spread greater than $2m|B_\tau|$. We output NO otherwise. Again, the influence of a node can be computed in $m = n^{o(1)}$ times.

Completeness. Suppose there exists $A_i \in \mathcal{A}$, $B_\tau \in \mathcal{B}$ such that $B_\tau \subseteq A_i$. Then in the τ -th epoch, by taking node $v_{1,i}$ in the seed set, the influence spread at least $\frac{1}{n}|B_\tau| \cdot nmt + 1 = |B_\tau|mt + 1$. Since the DIM algorithm gives $2/t$ -approximation, the influence is greater than $2m|B_\tau|$ in this case. Hence, we indeed output YES.

Soundness. Suppose $|A_i \cap B_\tau| < |B_\tau|/t$ for any $i, \tau \in [n]$, then we prove that no node in V_1 has influence spread more than $2m|B_\tau|$, in any epoch. Since the intersection of A_i and B_τ is less than $|B_\tau|/t$, and a node $v_{2,j} \in V_2$ has influence $1 + mt + mtn$ if $j \in B_\tau$ and it has influence 1 otherwise. We conclude for any node $v_{1,i}$, its influence is at most

$$1 + m + \frac{1}{t}|B_\tau| \cdot \frac{1}{n}(1 + mt + nmt) < 2|B_\tau|m.$$

Hence, we output NO in this case.

In summary, the reduced DIM requires $n^{1+o(1)}$ updates and queries, and it gives an answer for the problem in Theorem 4.2. Hence, we conclude under SETH, there is no $2/t$ -approximation algorithm unless the amortized running time is $n^{1-\epsilon}$. \square