

---

# Multi-layered Network Exploration via Random Walks: From Offline Optimization to Online Learning

---

Xutong Liu<sup>1</sup> Jinhang Zuo<sup>2</sup> Xiaowei Chen<sup>3</sup> Wei Chen<sup>4</sup> John C.S. Lui<sup>1</sup>

## Abstract

Multi-layered network exploration (MuLaNE) problem is an important problem abstracted from many applications. In MuLaNE, there are multiple network layers where each node has an importance weight and each layer is explored by a random walk. The MuLaNE task is to allocate total random walk budget  $B$  into each network layer so that the total weights of the unique nodes visited by random walks are maximized. We systematically study this problem from offline optimization to online learning. For the offline optimization setting where the network structure and node weights are known, we provide greedy based constant-ratio approximation algorithms for overlapping networks, and greedy or dynamic-programming based optimal solutions for non-overlapping networks. For the online learning setting, neither the network structure nor the node weights are known initially. We adapt the combinatorial multi-armed bandit framework and design algorithms to learn random walk related parameters and node weights while optimizing the budget allocation in multiple rounds, and prove that they achieve logarithmic regret bounds. Finally, we conduct experiments on a real-world social network dataset to validate our theoretical results.

## 1. Introduction

Network exploration is a fundamental paradigm of searching/exploring in order to discover information and resources available at nodes in a network, and random walk is often

used as an effective tool for network exploration (Lv et al., 2002; Gleich, 2015; Wilder et al., 2018). In this paper, we study the multi-layered network exploration via random walks problem, which can model many real-world applications, including resource searching in peer-to-peer (P2P) networks and web surfing in online social networks (OSNs).

In P2P networks, a user wants to find resources that are scattered on nodes (i.e. peers) in different P2P networks via some platform-specified strategies. A commonly used search strategy is based on multiple random walks (Lv et al., 2002). Since different resources have different importance, the search quality of different random walkers varies. Moreover, the resource-search process typically has a life span, i.e., a total time-limit or hop-limit for random walks. So the user’s goal is to decide how to allocate limited budgets to different random walkers to find as many important resources as possible. Another application is the web surfing, where users want to find information in different OSNs, by looking at posts via others’ home-pages (Lerman & Jones, 2006). Once the user arrives at one of his friends’ home-pages in a particular OSN, he could find some information and continues to browse the home-pages of his friends’ friends, which can be regarded as a random walk process. Since the user only has a finite duration in web surfing, the similar question is how to allocate his time in exploring different OSNs so to get the maximum amount of useful information.

We abstract the above application scenarios as the **Multi-Layered Network Exploration** (MuLaNE) problem. In MuLaNE, we model the overall network to be explored (e.g., combining different OSNs) as a multi-layered network  $\mathcal{G}$  that consists of  $m$  layers  $L_1, \dots, L_m$ . Each layer  $L_i$  (e.g., a single OSN) is represented by a weighted directed graph  $\mathcal{G}_i(\mathcal{V}_i, \mathcal{E}_i, w_i)$ , where  $\mathcal{V}_i$  is the set of nodes to be explored (e.g., users’ home-pages with importance weight  $\sigma_u$  for  $u \in \mathcal{V}_i$ ),  $\mathcal{E}_i \subseteq \mathcal{V}_i \times \mathcal{V}_i$  is the set of directed edges (e.g., social links), and  $w_i$  is the edge weight function on edges  $\mathcal{E}_i$ . Each layer  $L_i$  is associated with an explorer (or random walker)  $W_i$  and a fixed starting distribution  $\alpha_i$  on nodes  $\mathcal{V}_i$ . Explorer  $W_i$  starts on a node in  $\mathcal{G}_i$  following the distribution  $\alpha_i$ , and then walks on network  $\mathcal{G}_i$  following outgoing edges with probability proportional to edge weights. We assume that each random walk step will cost one unit of the bud-

---

<sup>1</sup>Department of Computer Science and Engineering, The Chinese University of Hong Kong, Hong Kong SAR, China <sup>2</sup>Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA, USA <sup>3</sup>Bytedance, Mountain View, CA, USA <sup>4</sup>Microsoft Research, Beijing, China. Correspondence to: Xutong Liu <liuxt@cse.cuhk.edu.hk>, Wei Chen <weic@microsoft.com>, John C.S. Lui <cslui@cse.cuhk.edu.hk>.

get<sup>1</sup>. Given a total budget constraint  $B$  and individual layer budget constraint  $c_i$ , the MuLaNE task is to find the optimal budget allocation  $\mathbf{k}=(k_1, \dots, k_m)$  with  $\sum_{i=1}^m k_i \leq B$ ,  $0 \leq k_i \leq c_i$  such that the expected total weights of *unique* nodes visited are maximized.

In real applications, the network structure  $\mathcal{G}$ , the node weights  $\sigma$  and the starting distributions  $\alpha_i$ 's may not be known in advance, so we consider both the *offline optimization* cases when  $\mathcal{G}$ ,  $\sigma$  and  $\alpha_i$ 's are known and the *online learning* case when  $\mathcal{G}$ ,  $\sigma$  and  $\alpha_i$ 's are unknown. Moreover, different layers may have overlapping vertices (e.g., homepage of the same user may appear in different OSNs), and the starting distributions  $\alpha_i$ 's may or may not be stationary distributions. In this paper, we provide a systematic study of all these case combinations.

For the offline optimization, we first consider the general overlapping setting and provide constant approximation algorithms based on the lattice submodularity property for non-stationary  $\alpha_i$ 's and the diminishing-return (DR) submodularity property for stationary  $\alpha_i$ 's. For the special non-overlapping setting, we design a dynamic programming algorithm that finds the exact optimal solution for MuLaNE.

In the online learning, we conduct multiple rounds of exploration, each of which has the same budget constraints  $B$  and  $c_i$ 's. After each round of exploration, the total weights of unique nodes visited in this round are the reward for this round, and the trajectory of every explorer in every layer and the importance weights of visited nodes are observed as the feedback, which could be used to learn information about the network for the benefit of future explorations.

We adapt the combinatorial multi-armed bandit (CMAB) framework and the CUCB algorithm (Chen et al., 2016) to our setting, and design online learning algorithms to minimize the regret, which is the difference between the cumulative reward achieved by the exact or approximate offline algorithm and that achieved by our learning algorithm, over  $T$  rounds. We show that directly learning the graph structure  $\mathcal{G}$  is inefficient. Instead, we define intermediate random variables in MuLaNE as the base arms in CMAB and learn these intermediate parameters, which can sufficiently determine the rewards to guide our budget allocation. Moreover, the node weight is not revealed until this node is first visited, which further complicates the design of online exploration algorithms. For the overlapping case, we adapt the CUCB algorithm to address the unrevealed node weights and the extra constraint of monotonicity for the intermediate parameters. We further improve the analysis by leveraging on special properties in our setting and show logarithmic regret bounds in this case. For the non-overlapping case,

<sup>1</sup>Our model can be extended to move multiple steps with one unit of budget.

we define more efficient intermediate parameters and use the exact offline algorithm, and thus achieve a better regret bound. Finally, we conduct experiments on a real-world multi-layered social network dataset to validate the effectiveness of both our offline and online algorithms.

Our contributions can be summarized as follows: (1) We are the first to model the multi-layered network exploration via random walks problem (MuLaNE) as an abstraction for many real-world applications (2) We provide a systematic study of the MuLaNE problem via theoretical analysis and empirical validation by considering offline and online settings for both overlapping and non-overlapping multi-layered networks. Due to the space limit, proofs are included in the supplementary material.

### 1.1. Related Work

Network exploration via random walks has been studied in various application contexts such as community detection (Pons & Latapy, 2005), centrality measuring (Gleich, 2015), large-scale network sampling (Li et al., 2019), and influence maximization (Wilder et al., 2018). Multiple random walks has also been used, e.g. in (Lv et al., 2002) for query resolution in peer-to-peer networks. However, none of these studies address the budgeted MuLaNE problem. MuLaNE is also related to the influence maximization (IM) problem (Kempe et al., 2003; Chen et al., 2009) and can be viewed as a special variant of IM. Different from the standard Independent Cascade (IC) model (Wang et al., 2012) in IM, which randomly broadcasts to all its neighbors, the propagation process of MuLaNE is a random walk that selects one neighbor. Moreover, each unit of budget in MuLaNE is used to propagate one random-walk step, not to select one seed node as in IC.

The offline budget allocation problem has been studied by Alon et al. (2012); Soma et al. (2014), the latter of which propose the lattice submodularity and a constant approximation algorithm based on this property. Our offline overlapping MuLaNE setting is based on a similar approach, but we provide a better approximation ratio compared to the original analysis in (Alon et al., 2012). Moreover, we further show that the stationary setting enjoys DR-submodularity, leading to an efficient algorithm with a better approximation ratio.

The multi-armed bandit (MAB) problem is first studied by Robbins (1952) and then extended by many studies (cf. (Bubeck & Cesa-Bianchi, 2012)). Our online MuLaNE setting fits into the general Combinatorial MAB (CMAB) framework of (Gai et al., 2012; Chen et al., 2016). CUCB is proposed as a general algorithm for CMAB (Chen et al., 2016). Our study includes several adaptations, such as handling unknown node weights, defining intermediate random variables as base arms and so on.

Chen et al. (2018) study the community exploration problem,

which is essentially a special case of MuLaNE with non-overlapping complete graphs. As a result, their technique is different and much simpler than ours.

The rest of the paper is organized as follows. Section 2 states the settings of MuLaNE; Section 3 states the equivalent bipartite coverage model to derive the explicit form for our reward function; Section 4 states offline algorithms for four offline settings with provable approximation guarantee and running time analysis; Section 5 states two online learning algorithms with regret analysis; Empirical results are shown in Section 6; and Section 7 concludes the paper.

## 2. Problem Settings

**Basic Notations.** In this paper, we use  $\mathbb{R}$  and  $\mathbb{Z}$  to denote the sets of real numbers and integers, respectively, and the associated subscript  $\geq 0$  and  $> 0$  denote their non-negative and positive subsets, respectively. Suppose we want to explore a multi-layered network  $\mathcal{G}(\mathcal{V}, \mathcal{E}, \boldsymbol{\sigma})$ , consisting of  $m$  layers  $L_1, \dots, L_m$  and  $N$  unique nodes  $\mathcal{V}$  with fixed node weights  $\boldsymbol{\sigma} \in [0, 1]^{\mathcal{V}}$ . Let  $[m]$  denote the set  $\{1, 2, \dots, m\}$ . Each layer  $L_i, i \in [m]$ , represents a subgraph of  $\mathcal{G}$  we could explore and is modeled as a weighted digraph  $\mathcal{G}_i(\mathcal{V}_i, \mathcal{E}_i; w_i)$ , where  $\mathcal{V}_i \subseteq \mathcal{V}$  is the set of vertices in layer  $L_i$ ,  $\mathcal{E}_i \subseteq \mathcal{V}_i \times \mathcal{V}_i$  is the set of edges in  $\mathcal{G}_i$ , and  $w_i : \mathcal{E}_i \rightarrow \mathbb{R}_{>0}$  is the edge weight function associating each edge of  $\mathcal{G}_i$ . For any two layers  $L_i$  and  $L_j$ , we say they are *non-overlapping* if  $\mathcal{V}_i \cap \mathcal{V}_j = \emptyset$ .  $\mathcal{G}$  is called a *non-overlapping* multi-layered network if any two layers are non-overlapping; otherwise  $\mathcal{G}$  is an *overlapping* multi-layered network. Let  $n_i = |\mathcal{V}_i|$  denote the size of layer  $L_i$ , and the adjacency matrix of  $\mathcal{G}_i$  is defined as  $\mathbf{A}_i \in \mathbb{R}_{>0}^{n_i \times n_i}$ , where  $\mathbf{A}_i[u, v] = w_i((u, v))$  if  $(u, v) \in \mathcal{E}_i$  and 0 otherwise.

**Exploration Rule.** Each layer  $L_i$  is associated with an explorer  $W_i$ , a budget  $k_i \in \mathbb{Z}_{>0}$  and an initial starting distribution  $\boldsymbol{\alpha}_i = (\alpha_{i,u})_{u \in \mathcal{V}_i} \in \mathbb{R}_{>0}^{n_i}$  with  $\sum_{u \in \mathcal{V}_i} \alpha_{i,u} = 1$ . The exploration process of  $W_i$  is identical to applying *weighted random walks* within layer  $L_i$ . Specifically, the explorer  $W_i$  starts the exploration from a random node  $v_1 \in \mathcal{V}_i$  with probability  $\alpha_{i,v_1}$ ; it consumes one unit of budget and continues the exploration by walking from  $v_1$  to one of its out-neighbors, say  $v_2$ , with probability  $\mathbf{P}_i[v_1, v_2]$ , where  $\mathbf{P}_i \in \mathbb{R}_{>0}^{n_i \times n_i}$  is the transition probability matrix, and  $\mathbf{P}_i[u, v] := \mathbf{A}_i[u, v] / (\sum_{w: (u,w) \in \mathcal{E}_i} \mathbf{A}_i[u, w])$ . Consider visiting the initial node  $v_1$  as the first step, the process is repeated until the explorer  $W_i$  walks  $k_i$  steps and visits  $k_i$  nodes (with possibly duplicated nodes). Note that one can easily generalize our results by moving  $\lambda_i \in \mathbb{Z}_{>0}$  steps with one unit of budget for  $W_i$ , and for simplicity, we set  $\lambda_i = 1$ .

**Reward Function.** We define the exploration trajectory for explorer  $W_i$  after exploring  $k_i$  nodes as  $\Phi(i, k_i) := (X_{i,1}, \dots, X_{i,k_i})$ , where  $X_{i,j} \in \mathcal{V}_i$  denotes the node visited

at  $j$ -th step, and  $\Pr(X_{i,1}=u) = \alpha_{i,u}, u \in \mathcal{V}_i$ . The reward for  $\Phi(i, k_i)$  is defined as the total weights of *unique* nodes visited by  $W_i$ , i.e.,  $\sum_{v \in \mathcal{V}_i} \sum_{j=1}^{k_i} \mathbb{1}_{\{X_{i,j}=v\}} \sigma_v$ . Considering trajectories of all  $W_i$ 's, the total reward is the total weights of *unique* nodes visited by all random walkers, i.e.,  $\sum_{i=1}^m \sum_{j=1}^{k_i} \mathbb{1}_{\{X_{i,j}=v\}} \sigma_v$ . For notational simplicity, let  $\mathcal{G} := (\mathcal{G}_1, \dots, \mathcal{G}_m)$ ,  $\boldsymbol{\alpha} := (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_m)$  and  $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_{|\mathcal{V}|})$  be the parameters of a problem instance and  $\mathbf{k} := (k_1, \dots, k_m)$  be the allocation vector. For overlapping multi-layered network  $\mathcal{G}$ , the total expected reward is

$$r_{\mathcal{G}, \boldsymbol{\alpha}, \boldsymbol{\sigma}}(\mathbf{k}) := \mathbb{E}_{\Phi(1,k_1), \dots, \Phi(m,k_m)} \sum_{i=1}^m \sum_{j=1}^{k_i} \mathbb{1}_{\{X_{i,j}=v\}} \sigma_v, \quad (1)$$

where its explicit formula will be discussed later in Sec. 3. For non-overlapping  $\mathcal{G}$ , the reward function can be simplified and written as the summation over separated layers,

$$r_{\mathcal{G}, \boldsymbol{\alpha}, \boldsymbol{\sigma}}(\mathbf{k}) := \sum_{i=1}^m \mathbb{E}_{\Phi(i,k_i)} \sum_{j=1}^{k_i} \mathbb{1}_{\{X_{i,j}=v\}} \sigma_v. \quad (2)$$

**Problem Formulation.** We are interested in the *budget allocation problem*: how to allocate the total budget  $B \in \mathbb{Z}_{>0}$  to the  $m$  explorers so as to maximize the total weights of unique nodes visited. Also, we assume there is a budget constraint  $\mathbf{c} := (c_1, \dots, c_m)$  such that the allocated budget  $k_i$  should not exceed  $c_i$ , i.e.,  $k_i \leq c_i$ , for  $i \in [m]$ .

**Definition 1.** Given graph structures  $\mathcal{G} := (\mathcal{G}_1, \dots, \mathcal{G}_m)$ , starting distributions  $\boldsymbol{\alpha} := (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_m)$ , total budget  $B$  and budget constraints  $\mathbf{c} := (c_1, \dots, c_m)$ , the **Multi-Layered Network Exploration problem**, denoted as **MuLaNE**, is formulated as the following optimization problem,

$$\max_{\mathbf{k}} r_{\mathcal{G}, \boldsymbol{\alpha}, \boldsymbol{\sigma}}(\mathbf{k}) \text{ s.t. } \mathbf{k} \in \mathbb{Z}_{>0}^m \leq \mathbf{c}, \sum_{i=1}^m k_i \leq B, \quad (3)$$

where  $r_{\mathcal{G}, \boldsymbol{\alpha}, \boldsymbol{\sigma}}(\mathbf{k})$  is given by Eq. (1) or Eq. (2). We also need to consider the following settings:

**Offline Setting.** For the offline setting, all problem instance parameters  $(\mathcal{G}, \boldsymbol{\alpha}, \boldsymbol{\sigma}, \mathbf{c}, B)$  are given, and we aim to find the optimal budget allocation  $\mathbf{k}^*$  determined by Eq. (3).

**Online Setting.** For the online setting, we consider  $T$ -round explorations. Before the exploration, we only know an upper bound of the total number of nodes in  $\mathcal{G}$  and the number of layers  $m$ ,<sup>2</sup> but we do not know about the network structure  $\mathcal{G}$ , the starting distributions  $\boldsymbol{\alpha}$  or node weights  $\boldsymbol{\sigma}$ . In round  $t \in [T]$ , we choose the budget allocation  $\mathbf{k}_t := (k_{t,1}, \dots, k_{t,m})$  only based on observations from previous rounds, where  $k_{t,i} \in \mathbb{Z}_{>0}$  is the budget allocated to the  $i$ -th layer and  $\sum_{i=1}^m k_{t,i} \leq B, \mathbf{k}_t \leq \mathbf{c}$ . Define  $\mathbf{k}_t$  as the *action* taken in round  $t$ . By taking the action  $\mathbf{k}_t$  we mean that we interact with the environment and the random explorer  $W_i$ , which is part of the environment, would

<sup>2</sup>More precisely we only need to know the number of independent explorers. If two explorers explore on the same layer, it is equivalent as two layers with identical graph structures.

explore  $k_{t,i}$  steps and generate the exploration trajectory  $(i; k_{t,i}) = (X_{i;1}, \dots, X_{i;k_{t,i}})$ . After we take the action  $(i; k_{t,i})$ , the exploration trajectory  $(i; k_{t,i})$  for each layer  $L_i$  as well as the fixed importance weight  $w_{t,i}$  of  $u \in L_i$  is revealed as the feedback, which we leverage on to learn parameters related to the graph structure: the starting distribution  $\mu$  and the node weights, so that we can select better actions in future rounds. The reward we gain in round  $t$  is the total weights of unique nodes visited by all random explorers in round  $t$ . Our goal is to design an efficient online learning algorithm  $\mathcal{A}$  to give us guidance on taking actions and gain as much cumulative reward as possible in  $T$  rounds.

In general, the online algorithm has to deal with the exploration-exploitation tradeoff. The cumulative regret is a commonly used metric to evaluate the performance of an online learning algorithm  $\mathcal{A}$ . Formally, the  $T$ -round  $(\epsilon, \delta)$ -approximation regret  $\mathcal{R}_{\mathcal{A}}$  is:

$$\text{Reg}_{\mathcal{A}}(T) = T \cdot r_{G^*}(k^*) - \mathbb{E} \sum_{t=1}^T r_{G^*}(k_t^{\mathcal{A}}); \quad (4)$$

where  $r_{G^*}(k)$  is the reward value for the optimal budget allocation  $k$ ,  $k_t^{\mathcal{A}}$  is budget allocation selected by the learning algorithm  $\mathcal{A}$  in round  $t$ , the expectation is taken over the randomness of the algorithm and the exploration trajectories in all  $T$  rounds, and  $(\epsilon, \delta)$  is the approximation guarantee of the offline oracle as explained below. Similar to other online learning frameworks (Chen et al., 2016; Wang & Chen, 2017), we assume that the online learning algorithm has access to an  $(\epsilon, \delta)$ -approximation oracle, which for problem instance  $(G; \mu; c; B)$  outputs an action  $k$  such that  $\Pr(r_{G^*}(k) \geq r_{G^*}(k^*) - \epsilon) \geq 1 - \delta$ . We also remark that the actual oracle we use takes certain intermediate parameters as inputs instead of  $(G; \mu; c; B)$ .

**Submodularity and DR-Submodularity Over Integer lattices.** To solve the MuLaNE problem, we leverage on the submodular and DR-submodular properties of the reward function. For any  $x, y \in Z_0^m$ , we denote  $x \vee y; x \wedge y \in Z_0^m$  as the coordinate-wise maximum and minimum of these two vectors, i.e.,  $(x \vee y)_i = \max\{x_i; y_i\}; (x \wedge y)_i = \min\{x_i; y_i\}$ . We define a function  $f : Z_0^m \rightarrow \mathbb{R}$  over the integer lattice  $Z_0^m$  as a submodular function if the following inequality holds for any  $x, y \in Z_0^m$ :

$$f(x \vee y) + f(x \wedge y) \geq f(x) + f(y); \quad (5)$$

Let  $\text{supp}(x \vee y)$  denote the set  $\{i \in [m] : x_i > y_i\}$ ,  $e_i = (0, \dots, 1, \dots, 0)$  be the one-hot vector whose  $i$ -th element is 1 and 0 otherwise, and  $x \vee y$  means  $x_i \vee y_i$  for all  $i \in [m]$ . We define a function  $f : Z_0^m \rightarrow \mathbb{R}$  as a DR-submodular (diminishing return submodular) function

<sup>3</sup>We further consider random node weights with unknown mean vector  $\mu$ , see the Appendix F.3 for details.

$$f(y \vee x) - f(y) \leq f(x \vee y) - f(x); \quad (6)$$

We say a function  $f$  is monotone if for any  $x \preceq y, f(x) \leq f(y)$ . Note that for a function  $f : Z_0^m \rightarrow \mathbb{R}$ , submodularity does not imply DR-submodularity over integer lattices. In fact, the former is weaker than the latter, that is, a DR-submodular function is always a submodular function, but not vice versa. However, for a typical submodular function  $f : \{0, 1\}^m \rightarrow \mathbb{R}$  defined on a set, they are equivalent.

### 3. Equivalent Bipartite Coverage Model

In order to derive the explicit formulation of the reward function  $r_{G^*}(k)$  in Eq. (1), we construct an undirected bipartite coverage graph  $\mathcal{G}(W; V; E^0)$ , where  $W = \{W_1, \dots, W_m\}$  denotes  $m$  random explorers  $V = \cup_{i \in [m]} V_i$  denotes all possible distinct nodes to be explored in  $G$ , and the edge set  $E^0 = \{(W_i; u) | u \in V_j; i \in [m]\}$  indicating whether node  $u$  could be visited by  $W_i$ . For each edge  $(W_i; u) \in E^0$ , we associate it with  $c_i + 1$  visiting probabilities denoted  $\mathcal{P}_{i,u}(k_i)$  for  $k_i \in [c_i]$ .  $\mathcal{P}_{i,u}(k_i)$  represents the probability that the nodes visited by the random walker  $W_i$  given the budget  $k_i$ , i.e.,  $\mathcal{P}_{i,u}(k_i) = \Pr(u \in (i; k_i))$ . Since  $W_i$  can never visit the node outside the  $i$ -th layer (i.e.,  $u \notin G_i$ ), we set  $\mathcal{P}_{i,u}(k_i) = 0$  if  $(W_i; u) \notin E^0$ . Then, given budget allocation  $k$ , the probability of a node  $u$  visited by at least one random explorer is  $\Pr(u; k) = 1 - \prod_{i \in [m]} (1 - \mathcal{P}_{i,u}(k_i))$  because each  $W_i$  has the independent probability  $\mathcal{P}_{i,u}(k_i)$  to visit  $u$ .

By summing over all possible nodes, the reward function is

$$r_{G^*}(k) = \sum_{u \in V} \mathbb{1}_u \prod_{i \in [m]} (1 - \mathcal{P}_{i,u}(k_i)) \quad (7)$$

According to Eq. (2), for non-overlapping multi-layered network, we can rewrite the reward function as:

$$r_{G^*}(k) = \prod_{i \in [m]} \sum_{u \in V_i} \mathcal{P}_{i,u}(k_i); \quad (8)$$

**Remark.** The bipartite coverage model is needed to integrate the graph structure and the random walk exploration mechanisms into  $\mathcal{P}_{i,u}(k_i)$  to determine the reward function. It also handles the scenario where each layer is explored by multiple random walkers (see Appendix A).

#### 3.1. Properties of the visiting probability $\mathcal{P}_{i,u}(k_i)$

The quantities  $\mathcal{P}_{i,u}(k_i)$ 's in Eq. (7) and (8) are crucial for both our offline and online algorithms, and thus we provide their analytical formulas and properties here. In order to analyze the property of the reward function, we apply the absorbing Markov Chain technique to calculate  $\mathcal{P}_{i,u}(k_i)$ . For  $u \notin G_i, \mathcal{P}_{i,u}(k_i) = 0$ ; For  $u \in G_i$ , we create an absorbing Markov Chain  $\mathcal{P}_i(u) \in \mathbb{R}^{n_i \times n_i}$  by setting the target

node  $u$  as the absorbing node. We derive the corresponding transition matrix  $P_i(u)$  as  $P_i(u)[v; ] = \begin{cases} 1 & \text{if } v = u \\ P_i[v; ] & \text{otherwise, where } P_i[v; ] \end{cases}$  denotes the row vector corresponding to the node  $v$ , and  $\mathbf{e}_u = (0; \dots; 0; 1; 0; \dots; 0)^T$  denotes the one-hot vector with 1 at the  $u$ -th entry and 0 elsewhere. Intuitively,  $P_i(u)$  corresponds to the transition probability matrix  $G$  after removing all out-edges of  $u$  and adding a self loop to itself in the original graph  $G$ . The random walk  $W_i$  will be trapped in  $u$  if it ever visits  $u$ . We observe that  $P_{i;u}(k_i)$  equals to the probability the random walker stays in the absorbing node  $u$  at step  $k_i \geq 2$  (trivially,  $P_{i;u}(0) = 0$ ),

$$P_{i;u}(k_i) = \sum_j P_i(u)^{k_i - 1} \mathbf{e}_u: \quad (9)$$

Then define the marginal gain  $g_{i;u}(k_i)$  at step  $k_i - 1$  as,

$$g_{i;u}(k_i) = P_{i;u}(k_i) - P_{i;u}(k_i - 1): \quad (10)$$

The physical meaning of  $g_{i;u}(k_i)$  is the probability that node  $u$  is visited exactly at the  $k_i$ -th step and not visited before  $k_i$ . Now, we can show  $g_{i;u}(k_i)$  is non-negative,

Lemma 1.  $g_{i;u}(k_i) \geq 0$  for any  $i \in [m], u \in V_i, k_i \geq 2$ .

This means  $P_{i;u}(k_i)$  is non-decreasing with respect to step  $k_i$ . In other words, the more budgets we allocate to the higher probability  $W_i$  can visit  $u$  in  $k_i$  steps.

Starting from arbitrary distributions: Although  $P_{i;u}(k_i)$  is monotone (non-decreasing) w.r.t  $k_i$ ,  $g_{i;u}(k_i)$  may not be monotonic non-increasing under arbitrary starting distributions. Intuitively, there may exist one critical step such that  $P_{i;u}(k_i)$  suddenly increases by a large value (see Appendix B.1). This means that  $P_{i;u}(k_i)$  lacks the diminishing return (or “discretely concave”) property, which many problems rely on to provide good optimization (Kapralov et al., 2013; Soma & Yoshida, 2015). In other words, we are dealing with a more challenging non-concave optimization problem for the discrete budget allocation.

Starting from the stationary distribution: If our walkers start with the stationary distributions, the following property holds:  $g_{i;u}(k_i)$  is non-increasing w.r.t  $k_i$ .

Lemma 2.  $g_{i;u}(k_i + 1) \leq g_{i;u}(k_i) \leq 0$  for any  $i \in [m], u \in V_i, k_i \geq 2$ , if  $\mathbf{e}_i = \mathbf{e}_i$ , where  $\mathbf{e}_i = \sum_j P_i = \mathbf{e}_i$ .

## 4. Offline Optimization for MuLaNE

In this section, we first consider the general case, where layers are overlapping and starting distributions are arbitrary. Next, we consider the starting distribution is the stationary distribution, and give solutions with better solution quality and time complexity. Then we analyze special cases where layers are non-overlapping and give optimal solutions for

<sup>4</sup>When related to matrix operations, we treat all vectors as column vectors by default.

arbitrary distributions. The summary for of ine models and algorithmic results are presented in Table 1.

### 4.1. Overlapping MuLaNE

Starting from arbitrary distributions. Based on the equivalent bipartite coverage model, one can observe that our problem formulation is a generalization of the Probabilistic Maximum Coverage (PMC) (Chen et al., 2016) problem, which is NP-hard and has a  $(1 - \epsilon)$  approximation based on submodular set function maximization. However, our problem is more general in that we want to select multi-sets from  $V$  with budget constraints, and the reward function in general does not have the DR-submodular property for an arbitrary starting distribution. Nevertheless, we have the following lemma to solve our problem.

Lemma 3. For any network  $G$ , distribution  $\mathbf{e}$  and weights  $\{r_G; \cdot\} : Z_{>0}^m \rightarrow \mathbb{R}$  is monotone and submodular.

Leveraging on the monotone submodular property, we design a Budget Effective Greedy algorithm (Alg. 1). The core of Alg. 1 is the BEG procedure. Let  $(i; b; k)$  be the per-unit marginal gain  $(r_G; \cdot; (k + b) - r_G; \cdot; (k)) = b$  for allocating  $b$  more budgets to layer  $i$ , which equals to  $P_{i;u} \sum_{u \in V_i} (1 - P_{i;u}(k_i)) (P_{i;u}(k_i + b) - P_{i;u}(k_i)) = b$ :  $(11)$

BEG procedure consists of two parts and maintains a queue  $Q$ , where any pair  $(i; b) \in Q$  represents a tentative plan of allocating  $b$  more budgets to layer  $i$ . The first part is built around the while loop (line 6-10), where each iteration greedily selects the  $(i; b)$  pair in  $Q$  such that the per-unit marginal gain  $(i; b; k)$  is maximized. The second part is a for loop (line 12-13), where in the  $i$ th round we attempt to allocate all  $b_i$  budgets to layer  $i$  and replace the current best budget allocation if we have a larger reward.

Theorem 1. Algorithm 1 obtains a  $(1 - \epsilon) \approx 0.357$ -approximate solution, where  $\mathbf{e}$  is the solution of equation  $\mathbf{e} = 2 \mathbf{e}$ , to the overlapping MuLaNE problem.

Line 1 uses  $O(m \sum_{k=1}^{n_{\max}^3})$  time to pre-calculate visiting probabilities based on Eq. (9), where  $n_{\max} = \max_i |V_i|$ . In the BEG procedure, the while loop contains  $B$  iterations, the size of queue  $Q$  is  $O(m \sum_{k=1}^{n_{\max}})$ , and line 7 uses  $O(n_{\max})$  to calculate  $(i; b; k)$  by proper pre-computation and update (see Appendix C.3), thus the time complexity of Algorithm 1 is  $O(B \sum_{k=1}^{n_{\max}} m n_{\max} + m \sum_{k=1}^{n_{\max}^3})$ .

Remark 1. The idea of combining the greedy algorithm with enumerating solutions on one layer is also adopted in previous works (Khuller et al., 1999; Alon et al., 2012), but they only give a  $\frac{1}{2}(1 - \epsilon) \approx 0.316$ -approximation analysis. In this paper, we provide a novel analysis with a better  $(1 - \epsilon) \approx 0.357$ -approximation, see Appendix C.1 for details.

Table 1. Summary of the of ine models and algorithms.

| Overlapping? | Starting distribution | Algorithm               | Aprx ratio       | Time complexity  |
|--------------|-----------------------|-------------------------|------------------|--|
| X            | Arbitrary             | Budget Effective Greedy | $(1 - e^{-1})^5$ | $O(B \sum_{k=1}^m n_{\max} + \sum_{k=1}^m n_{\max}^3)$ |
| X            | Stationary            | Myopic Greedy           | $(1 - e^{-1})$   | $O(B \sum_{k=1}^m n_{\max} + \sum_{k=1}^m n_{\max}^3)$ |
|              | Arbitrary             | Dynamic Programming     | 1                | $O(B \sum_{k=1}^m m + \sum_{k=1}^m n_{\max}^3)$        |
|              | Stationary            | Myopic Greedy           | 1                | $O(B \log m + \sum_{k=1}^m n_{\max}^3)$                |

**Algorithm 1 Budget Effective Greedy (BEG) Algorithm for the Overlapping MuLaNE**

Input: Network  $G$ , starting distributions  $\{p_i\}$ , node weights  $\{w_i\}$ , budget  $B$ , constraints  $\{c_i\}$ .  
 Output: Budget allocation  $k$ .

- 1: Compute visiting probabilities  $\{p_{i;u}(b)\}_{i \in [m]; u \in [2^V]; b \in [c_i]}$  according to Eq. (9).
- 2:  $k := \text{BEG}(\{p_{i;u}(b)\}_{i \in [m]; u \in [2^V]; b \in [c_i]}, B, c)$ .
- 3: Procedure  $\text{BEG}(\{p_{i;u}(b)\}_{i \in [m]; u \in [2^V]; b \in [c_i]}, B, c)$
- 4: Let  $k := (k_1; \dots; k_m) \in [0, K]^m$ .
- 5: Let  $Q = \{(i; b_i) \mid i \in [m]; 1 \leq b_i \leq c_i\}$ .
- 6: while  $K > 0$  and  $Q \neq \emptyset$ ; do
- 7:  $(i; b) = \arg \max_{(i; b) \in Q} (i; b; k) = b$ . Eq. (11)
- 8:  $k_i = k_i + b, K = K - b$ .
- 9: Modify all pairs  $(i; b) \in Q$  to  $(i; b - b)$ .
- 10: Remove all pairs  $(i; b) \in Q$  such that  $b = 0$ .
- 11: end while
- 12: for  $i \in [m]$  do
- 13: if  $r_{G; i}(k_i) > r_{G; i}(k_i - 1)$ , then  $k_i = k_i - 1$ .
- 14: end for
- 15: return  $k := (k_1; \dots; k_m)$ .
- 16: end Procedure

Remark 2. Another algorithm (Alon et al., 2012) with a better approximation ratio is to use partial enumeration techniques (i.e., BEGE), which can achieve  $(1 - e^{-1})$  approximation ratio. This is the best possible solution in polynomial time unless  $P=NP$ . But the time complexity is prohibitively high in  $O(B^4 m^4 \sum_{k=1}^m n_{\max} + m \sum_{k=1}^m n_{\max}^3)$ . The algorithm and the analysis are provided in the Appendix D.2.

Starting from the stationary distribution. We also consider the special case where each random exploration starts from the stationary distribution  $p_i$  with  $\sum_i p_i = 1$ . In this case, we have the following stronger DR-submodularity.

Lemma 4. For any network  $G$ , stationary distributions and node weights  $\{w_i\}$ , function  $r_{G; i}(\cdot) : Z_0^m \rightarrow \mathbb{R}$  is monotone and DR-submodular.

Since the reward function is DR-submodular, the BEG procedure can be replaced by the simple MG procedure in Alg. 2 with a better approximation ratio. The time complexity is also improved to  $O(B \sum_{k=1}^m n_{\max} + m \sum_{k=1}^m n_{\max}^3)$ .

Theorem 2. Algorithm 2 obtains a  $(1 - e^{-1})$ -approximate

**Algorithm 2 Myopic Greedy (MG) Algorithm for MuLaNE**

- 1: Same input, output and line 1-2 as in Alg. 1, except replacing BEG with MG procedure below.
- 2: Procedure  $\text{MG}(\{p_{i;u}(b)\}_{i \in [m]; u \in [2^V]; b \in [c_i]}, B, c)$
- 3: Let  $k := (k_1; \dots; k_m) \in [0, K]^m$ .
- 4: while  $K > 0$  do
- 5:  $i = \arg \max_{i \in [m]; k_i < c_i} (i; 1; k)$ . Eq. (11)
- 6:  $k_i = k_i + 1, K = K - 1$ .
- 7: end while
- 8: return  $k = (k_1; \dots; k_m)$ .
- 9: end Procedure

solution to the overlapping MuLaNE with the stationary starting distributions.

4.2. Non-overlapping MuLaNE

For non-overlapping MuLaNE, we are able to achieve the exact optimal solution: for the stationary starting distribution, a slight modification of the greedy algorithm Alg. 2 gives the optimal solution, while for an arbitrary starting distribution, we design a dynamic programming algorithm to compute the optimal solution. Due to the space constraint, the details are included in Appendix E.

5. Online Learning for MuLaNE

In the online setting, we continue to study both overlapping and non-overlapping MuLaNE problem. However, the network structure  $G$ , the distribution  $\{p_i\}$  and node weights are not known a priori. Instead, the only information revealed to the decision maker includes total budget  $B$ , the number of layers  $m$ , the number of the target nodes  $n_i$  (or an upper bound of it) and the budget constraint.

5.1. Online Algorithm for Overlapping MuLaNE

For the unknown network structure and starting distributions, we bypass the transition matrices  $P_{i;u}(b)$  and directly estimate the visiting probabilities  $\{p_{i;u}(b)\}_{i \in [m]; u \in [2^V]; b \in [c_i]}$ . This avoids the analysis for how the estimated  $\{p_{i;u}(b)\}$  affects the performance of online algorithms, which could be unbounded since we consider general graph structures and the reward function given by Eq. (7) is highly non-linear in  $\{p_{i;u}(b)\}$ . Moreover, we can save the matrix calculation by directly using  $\{p_{i;u}(b)\}$ ,

which is efficient even for large networks.

Specifically, we maintain a set of base arms  $A = \{(i; u; b) \mid i \in [m]; u \in V; b \in [c]\}$ , where the total number of arms is  $|A| = \sum_{i \in [m]} |V| \cdot |c|$ . For each base arm  $(i; u; b) \in A$ , we denote  $\mu_{i;u;b}$  as the true value of each base arm, i.e.,  $\mu_{i;u;b} = P_{i;u}(b)$ . For the unknown node weights, we maintain the optimistic weight  $\hat{\mu}_u = 1$  if  $u \in V$  has not been visited. After  $u$  is first visited and its true value  $\mu_u$  is revealed, we replace  $\hat{\mu}_u$  with  $\mu_u$ . With a little abuse of the notation, we use  $\mu$  to denote the unknown intermediate parameters and  $r; (k)$  to denote the reward  $r; (k)$ .

We present our algorithm in Alg. 3, which is an adaptation of the CUCB algorithm for the general Combinatorial Multi-arm Bandit (CMAB) framework (Chen et al., 2016) to our setting. Notice that we use  $A$  as defined above in the algorithm, and  $V$  is defined using  $V$ , which is the set of node ids and should not be known before the learning process starts. This is not an issue, because at the beginning we can create  $|V|$  (which is known) placeholders for the node ids, and once a new node is visited, we immediately replace one of the placeholders with the new node id. This appearing in the above definition of  $A$  is just for notational convenience.

In Alg. 3, we maintain an unbiased estimation of the visiting probability  $P_{i;u}(b)$ , denoted as  $\hat{T}_{i;u;b}$ . Let  $T_{i;u;b}$  record the number of times arm  $(i; u; b)$  is played so far and  $\hat{\mu}_{i;u;b}$  denote the optimistic importance weight. In each round  $t$ , we compute the confidence radius  $\beta_{i;u;b}$  in line 5, which controls the level of exploration. The confidence radius is larger when the arm  $(i; u; b)$  is not explored often (i.e.,  $T_{i;u;b}$  is small), and thus motivates more exploration. Due to the randomness of the exploration process, the upper confidence bound (UCB) value  $\hat{\mu}_{i;u;b} + \beta_{i;u;b}$  could be decreasing w.r.t. but our of line oracle BEG can only accept non-decreasing UCB values (otherwise the  $(1 - \epsilon; 1)$  approximation is not guaranteed). Therefore, in line 8, we increase the UCB value  $\hat{\mu}_{i;u;b}$  and set it to  $\max_{j \in [b]} \hat{\mu}_{i;u;j}$ . This is the adaption of the CUCB algorithm to our case, and thus we name our algorithm as CUCB-MAX. After we apply the  $(1 - \epsilon; 1)$  approximate solution  $k$  given by the BEG oracle (i.e., Alg. 1), we get  $m$  trajectories as feedbacks. In line 11, we can update the unknown weights of visited nodes. In line 13, for base arm  $(i; u; b)$  with  $b = k_i$ , we update corresponding statistics by the Bernoulli random variable  $Y_{i;u;b} \in \{0, 1\}$  indicating whether node  $u$  is visited by  $W_i$  in first  $b$  steps.

**Regret Analysis.** We define the reward gap  $\Delta_k = \max(0; r; (k) - \mu^*_{i;u;b})$  for all feasible action  $k$  satisfying  $\sum_{i=1}^m k_i = B$ ,  $0 \leq k_i \leq c_i$ , where  $k^*$  is the optimal solution for parameters,  $\mu^*_{i;u;b} = 1 - \epsilon$  is the approximation ratio of the  $(1 - \epsilon; 1)$ -approximate of line oracle. For each base arm  $(i; u; b)$ , we define  $\mu_{i;u;b}^{\min} = \min_{k > 0; k_i = b} \Delta_k$  and  $\mu_{i;u;b}^{\max} = \max_{k > 0; k_i = b} \Delta_k$ . As a convention, if there

**Algorithm 3 CUCB-MAX Algorithm for the MuLaNE**

```

Input: Budget  $B$ , number of layers  $m$ , number of nodes  $|V|$ , constraints  $c$ , of line oracle BEG.
1: For each arm  $(i; u; b) \in A$ ,  $T_{i;u;b} = 0, \hat{\mu}_{i;u;b} = 0$ .
2: For each node  $u \in V, \mu_u = 1$ .
3: for  $t = 1; 2; 3; \dots; T$  do
4:   for  $(i; u; b) \in A$  do
5:      $\beta_{i;u;b} = \sqrt{\frac{3 \ln t}{2 T_{i;u;b}}}$ .
6:      $\tilde{\mu}_{i;u;b} = \min\{\hat{\mu}_{i;u;b} + \beta_{i;u;b}; 1\}$ .
7:   end for
8:   For  $(i; u; b) \in A, \hat{\mu}_{i;u;b} = \max_{j \in [b]} \tilde{\mu}_{i;u;j}$ .
9:    $k = \text{BEG}(\{(\mu_{i;u;b})_{(i;u;b) \in A}, (\mu_u)_{u \in V}, B, c\})$ .
10:  Apply budget allocation  $k$ , which gives trajectories
11:   $X := (X_{i;1}; \dots; X_{i;k_i})_{i \in [m]}$  as feedbacks.
12:  For any visited node  $u \in V$ , receive its node weight  $\mu_u$  and set  $\mu_u = \mu_u$ .
13:  For any  $(i; u; b) \in A, Y_{i;u;b} = \begin{cases} 1 & \text{if } u \in \{X_{i;1}; \dots; X_{i;k_i}\} \\ 0 & \text{otherwise.} \end{cases}$ 
14:  For  $(i; u; b) \in A$ , update  $T_{i;u;b}$  and  $\hat{\mu}_{i;u;b}$ :
       $T_{i;u;b} = T_{i;u;b} + 1, \hat{\mu}_{i;u;b} = \hat{\mu}_{i;u;b} + (Y_{i;u;b} - \hat{\mu}_{i;u;b}) / T_{i;u;b}$ ;
       $\hat{\mu}_{i;u;b} = \max\{\hat{\mu}_{i;u;b}, \tilde{\mu}_{i;u;b}\}$ ;
end for
    
```

is no action  $k$  with  $k_i = b$  such that  $\Delta_k > 0$ , we define  $\mu_{i;u;b}^{\min} = 1$  and  $\mu_{i;u;b}^{\max} = 0$ . Let  $\mu_{i;u;b}^{\min} = \min_{(i;u;b) \in A} \mu_{i;u;b}^{\min}$  and  $\mu_{i;u;b}^{\max} = \max_{(i;u;b) \in A} \mu_{i;u;b}^{\max}$ . The following theorem summarizes the regret bound for Alg. 3.

**Theorem 3.** Algorithm 3 has the following distribution-dependent  $(1 - \epsilon; 1)$  approximation regret,

$$\text{Reg} : (T) \leq \sum_{(i;u;b) \in A} \left( \frac{108mj|V| \ln T}{\mu_{i;u;b}^{\min}} + 2|A| + \frac{2}{3}|A| \mu_{i;u;b}^{\max} \right)$$

**Remark 1.** Looking at the above distribution dependent bound, we have the  $\Theta(\log T)$  approximation regret, which is asymptotically tight. Coefficient  $m|V|$  in the leading term corresponds to the number of edges in the complete bipartite coverage graph. Notice that we cannot use the true edge number  $\sum_{i \in [m]} |V|$ , because the learning algorithm does not know which nodes are contained in each layer, and has to explore all visiting possibilities given by the default complete bipartite graph. The set of base arms has some redundancy due to the correlation between these base arms, and thus it is unclear if the summation over all base arms in the regret bound is tight. For the non-overlapping case, we further reduce the number of base arms to achieve better regret bounds, but for the overlapping case, how to further reduce base arms to achieve a tighter regret bound is a challenging open question left for the future work.

**Remark 2.** The  $(1 - \epsilon; 1)$  approximate regret is determined by the of line oracle BEG that we plug in Line 9 and can be replaced by  $(1 - \epsilon; 1)$  regret using BEGE or even by the exact regret if the oracle can obtain the optimal budget allocation. The usage of BEG is a trade-off we

make between computational efficiency and learning efficiency and empirically, it performs well as we shall see in Section 6.

Remark 3. The full proof of the above theorem is included in the Appendix F, where we rely on the following properties of  $r_{i;b}(k)$  to bound the regret.

Property 1. (Monotonicity). The reward  $r_{i;b}(k)$  is monotonically increasing, i.e., for any budget allocation  $\mathbf{a}$  and any two vectors  $\mathbf{v} = (v_{i;b})_{(i;b) \in \mathcal{A}}$ ,  $\mathbf{v}' = (v'_{i;b})_{(i;b) \in \mathcal{A}}$  and any node weights  $\mathbf{w}$ , we have  $r_{i;b}(k) \geq r'_{i;b}(k)$ , if  $v_{i;b} \geq v'_{i;b}$  and  $v_{i;b} \geq 0, v'_{i;b} \geq 0, \forall (i;b) \in \mathcal{A}; v \geq v'$ .

Property 2. (1-Norm Bounded Smoothness). The reward function  $r_{i;b}(k)$  satisfies the 1-norm bounded smoothness condition, i.e., for any budget allocation  $\mathbf{a}$  and any two vectors  $\mathbf{v} = (v_{i;b})_{(i;b) \in \mathcal{A}}$ ,  $\mathbf{v}' = (v'_{i;b})_{(i;b) \in \mathcal{A}}$  and any node weights  $\mathbf{w}$ , we have  $|r_{i;b}(k) - r'_{i;b}(k)| \leq \sum_{(i;b) \in \mathcal{A}} |v_{i;b} - v'_{i;b}|$ .

(a) Offline, overlapping. (b) Offline, non-overlapping.

(c) Online, overlapping. (d) Online, non-overlapping.

Figure 1. Above: total weights of unique nodes visited for offline algorithms. Below: regret for online algorithms when  $B = 3000$ .

We emphasize that our algorithm and analysis differ from the original CUCB algorithm (Chen et al., 2016) as follows. First, we have the additional regret caused by the over-estimated weights for unvisited nodes,  $\sum_{(i;b) \in \mathcal{A}} w_{i;b}$  term in property 2. We carefully bound this term based on the observation that  $w_{i;b}$  is small and decreasing quickly before  $(i;b)$  is first visited. Next, we have to take the max (line 8) to guarantee the UCB value is monotone  $\forall (i;b)$  since our BEG oracle can only output  $(1-\epsilon)$ -approximation with monotone inputs. Due to the above operation,  $(i;b)$ 's UCB value depends on the feedback from all arms  $(i; j)$  for  $j \in \mathcal{B}$  (set in line 12). So we should update all these arms (line 13) to guarantee that the estimates to all these arms are accurate enough. Finally, directly following the standard CMAB result would have a larger regret, because arms in  $\mathcal{I}$  are defined as triggered arms, but only arms in  $\mathcal{I}^0 = \{(i;b) \in \mathcal{A} \mid k_{i;b} = 0\}$  affect the rewards. So we conceptually view arms in  $\mathcal{I}^0$  as triggered arms and use a tighter 1-Norm Bounded Smoothness condition as given above to derive a tighter regret bound. This improves the coefficient of the leading  $\ln T$  term in the distribution dependent regret by a factor of  $\frac{1}{B} = O(1/B)$ , and the  $1/\epsilon$  term is smaller since the original definition would have  $\frac{1}{\epsilon} = \min_{k > 0; b \in \mathcal{B}} k \cdot b$ .

### 5.2. Online Algorithm for Non-overlapping Case

For the non-overlapping case, we have per-layer marginal gains as our base arms. Concretely, we maintain a set of base arms  $\mathcal{A} = \{(i;b) \mid (i;b) \in \mathcal{A}\}$ . For each base arm  $(i;b) \in \mathcal{A}$ , let  $v_{i;b} = \sum_{u \in \mathcal{V}} w_u (P_{i;u}(b) - P_{i;u}(b-1))$  be the true marginal gain of assigning budget to layer  $i$ . We apply the standard CUCB algorithm to this setting and call the resulting algorithm CUCB-MG (Algorithm 9 in Appendix G). Note that in the non-overlapping setting we can

solve the offline problem exactly, so we can use the exact offline oracle to solve the online problem and achieve an exact regret bound. This is the major advantage over the overlapping setting where we can only achieve an approximate bound. Define  $r_{i;b}(k) = r_{i;b}(k)$  for all feasible actions  $k$ , and  $r_{i;b}^{\min} = \min_{k > 0; k_i \leq b} r_{i;b}(k)$ , CUCB-MG has  $O(\sum_{(i;b) \in \mathcal{A}} \frac{1}{r_{i;b}^{\min}} \ln T)$  regret bound.

## 6. Experiments

**Dataset and settings.** We conduct experiments on a real-world multi-layered network FF-TW-YT, which contains  $m = 3$  layers representing users' social connections in FriendFeed (FF), Twitter (TW) and YouTube (YT) (Dickison et al., 2016). In total, FF-TW-YT has 407 distinct vertices representing users and 836 directed edges representing connections ("who follows whom") among users. The statistics for each layer is summarized in Table 2. We transform the FF-TW-YT network  $\mathcal{G}(V; E)$  into a symmetric directed network (by adding a new edge  $(v; u) \in E$  if  $(u; v) \in E$ ) because a user can be visited via her followers or followees. Each edge weight is set to be uniform and the node weights are set to be  $\text{unif}(0; 1)$  uniformly at random. Each random walker always starts from the smallest node-id in each layer and we set constraints equal to the total budget  $B$ . Note that in order to test the non-overlapping case, we use the same network but relabel node-ids so that

Table 2. Statistics for FF-TW-YT network

| Layer         | FriendFeed | Twitter | YouTube |
|---------------|------------|---------|---------|
| # of vertices | 5,540      | 5,702   | 663     |
| # of edges    | 31,921     | 42,327  | 614     |



Table 3. Running time (seconds) for of line and online algorithms.

|      | B=2.6k | B=2.8k | B=3.0k |     | B=2.6k | B=2.8k | B=3.0k | Overlapping? | X      |       |
|------|--------|--------|--------|-----|--------|--------|--------|--------------|--------|-------|
| BEG  | 0.274  | 0.316  | 0.363  | DP  | 0.038  | 0.044  | 0.050  | BEG          | 1.22   | NA    |
| BEGE | 34.37  | 45.51  | 59.20  | MG  | 0.008  | 0.009  | 0.010  | BEGE         | 60.03  | NA    |
| OPT  | 91.16  | 98.42  | 105.63 | OPT | 70.10  | 75.78  | 81.11  | DP           | NA     | 0.86  |
|      |        |        |        |     |        |        |        | OPT          | 107.33 | 82.01 |

(a) Running time of of line algorithms for the overlapping case with different budgets B.

(b) Running time of of line algorithms for the non-overlapping case with different budgets B.

(c) Per-round running time for CUCB-MAX (or CUCB-MG) with different oracles when B=3.0k.

they do not overlap between different layers. To handle overlapping case, the results are similar, but the difference the randomness, we repeat 1000 times and present the averaged total weights of unique nodes visited for of line of DP. For the online setting, all CUCB-MAX/CUCB-MG optimization. We calculate the regret by comparing with curves outperform the baselines. This demonstrates empirically that CUCB-MAX algorithm can effectively learn the unknown parameters while optimizing the objective. For average over 200 independent experiments to provide the computational efficiency of online learning algorithms, mean regret with 95% confidence interval. To evaluate the efficiency, we also present the running time for CUCB-MAX for both of line and online algorithms in Table 3.

Algorithms in comparison. For the of line setting, we present the results for Alg. 1 (denoted as BEG), Alg. 2 (denoted as MG), BEG with partial enumeration (denoted as BEGE) and Alg. 7 (denoted as DP). We provide two baselines PROP-S and PROP-W, which allocates the budget proportional to the layer size and proportional to the total weights if we allocate B budgets to that layer, respectively. The optimal solution (denoted as OPT) is also provided by enumerating all possible budget allocations. For online settings, we consider CUCB-MAX (Alg. 3) and CUCB-MG (Alg. 9) algorithms. We shrink the confidence interval by  $\epsilon$ , i.e.,  $\epsilon_i; u; b$   $\epsilon_i; u; b$ , to speed up the learning, though our theoretical regret bound requires  $\epsilon = 1$ . For baselines, we consider the EMP algorithm which always allocates according to the empirical mean, and the Greedy algorithm which allocates budgets according to empirical mean with probability 1 and allocates a B budgets to the  $i$ -th layer with probability  $\epsilon = m$ . We also compare with the Thompson sampling (TS) method (Wang & Chen, 2018), which uses Beta distribution  $\text{Beta}(\alpha; \beta)$  (where  $\alpha = \beta = 1$  initially) as prior distribution for each base arm.

Experimental results. We show the results for MuLaNE problems in Figure 1. For the of line overlapping case, both BEG and MG outperform two baselines PROP-W and PROP-S in receiving total weights. Although not guaranteed by the theory, BEG are empirically close to BEGE and the optimal solution (OPT). As for the computational efficiency, in Table 3a, the BEG is at least two orders of magnitude (e.g., 163 times when B=3.0k) faster than BEGE and OPT. Combining that the reward of BEG is empirically close to BEGE and the optimal solution, this shows that BEG is empirically better than BEGE and OPT. For the of line non-

overlapping case, the results are similar, but the difference is that we have the theoretical guarantee for the optimality of DP. For the online setting, all CUCB-MAX/CUCB-MG optimization. We calculate the regret by comparing with curves outperform the baselines. This demonstrates empirically that CUCB-MAX algorithm can effectively learn the unknown parameters while optimizing the objective. For average over 200 independent experiments to provide the computational efficiency of online learning algorithms, mean regret with 95% confidence interval. To evaluate the efficiency, we also present the running time for CUCB-MAX with different oracles in Table 3c. CUCB-MAX with BEG is 50 times faster than BEGE, which is consistent with our theoretical analysis. The results for different budgets are consistent with B = 3000, which are included in the Appendix H.2. Results and analysis for stationary starting distributions are also in the Appendix H.1.

## 7. Conclusions and Future Work

This paper formulates the multi-layered network exploration via random walks (MuLaNE) as a budget allocation problem, requiring that the total weights of distinct nodes visited on the multi-layered network is maximized. For the of line setting, we propose four algorithms for MuLaNE according to the specification of multi-layered network (overlapping or non-overlapping) and starting distributions (arbitrary or stationary), each of which has a provable guarantee on approximation factors and running time. We further study the online setting where network structure and the node weights are not known a priori. We propose the CUCB-MAX algorithm for overlapping MuLaNE and the CUCB-MG algorithm for the non-overlapping case, both of which are bounded by  $\mathcal{O}(\log T)$  (approximate) regret. Finally, we conduct experiments on a social network dataset to show the empirical performance of our algorithms.

There are many compelling directions for the future study. For example, it would be interesting to extend our problem where the decision maker can jointly optimize the starting distribution and the budget allocation. One could also study the adaptive MuLaNE by using the feedback from the exploration results of the previous steps to determine the exploration strategy for future steps.

## 8. Acknowledgement

The work of John C.S. Lui was supported in part by the GRF 14200420.

## References

- Alon, N., Gamzu, I., and Tennenholtz, M. Optimizing budget allocation among channels and in uencers. In Proceedings of the 21st international conference on World Wide Webpp. 381–388. ACM, 2012.
- Bubeck, S. and Cesa-Bianchi, N. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. Foundations and Trends in Machine Learning5(1):1–122, 2012.
- Chen, W., Wang, Y., and Yang, S. Efficient influence maximization in social networks. In Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data miningpp. 199–208, 2009.
- Chen, W., Wang, Y., Yuan, Y., and Wang, Q. Combinatorial multi-armed bandit and its extension to probabilistically triggered arms. The Journal of Machine Learning Research 17(1):1746–1778, 2016.
- Chen, X., Huang, W., Chen, W., and Lui, J. C. Community exploration: From offline optimization to online learning. In Advances in Neural Information Processing Systems pp. 5474–5483, 2018.
- Dickison, M. E., Magnani, M., and Rossi, L. Multilayer social networks Cambridge University Press, 2016.
- Dubhashi, D. P. and Panconesi, G. Concentration of measure for the analysis of randomized algorithms Cambridge University Press, 2009.
- Gai, Y., Krishnamachari, B., and Jain, R. Combinatorial network optimization with unknown variables: Multi-armed bandits with linear rewards and individual observations. IEEE/ACM Transactions on Networking (TON) 20(5): 1466–1478, 2012.
- Gleich, D. F. Pagerank beyond the web. SIAM Review 57(3):321–363, 2015.
- Kapralov, M., Post, I., and Vonáček, J. Online submodular welfare maximization: Greedy is optimal. Proceedings of the twenty-fourth annual ACM-SIAM symposium on Discrete algorithmspp. 1216–1225. SIAM, 2013.
- Kempe, D., Kleinberg, J., and Tardós, G. Maximizing the spread of influence through a social network. Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data miningpp. 137–146, 2003.
- Khuller, S., Moss, A., and Naor, J. S. The budgeted maximum coverage problem. Information processing letters 70(1):39–45, 1999.
- Kijima, M. Markov processes for stochastic modeling volume 6. CRC Press, 1997.
- Krause, A., Leskovec, J., Guestrin, C., VanBriesen, J., and Faloutsos, C. Efficient sensor placement optimization for securing large water distribution networks. Journal of Water Resources Planning and Management 134(6): 516–526, 2008.
- Lerman, K. and Jones, L. Social browsing on ickr. Xiv preprint cs/0612047, 2006.
- Li, Y., Wu, Z., Lin, S., Xie, H., Lv, M., Xu, Y., and Lui, J. C. Walking with perception: Efficient random walk sampling via common neighbor awareness. 2019 IEEE 35th International Conference on Data Engineering (ICDE) pp. 962–973. IEEE, 2019.
- Lv, Q., Cao, P., Cohen, E., Li, K., and Shenker, S. Search and replication in unstructured peer-to-peer networks. In Proceedings of the 16th international conference on Supercomputingpp. 84–95, 2002.
- Pons, P. and Latapy, M. Computing communities in large networks using random walks. International symposium on computer and information sciencespp. 284–293. Springer, 2005.
- Raginsky, M., Sason, I., et al. Concentration of measure inequalities in information theory, communications, and coding. Foundations and Trends in Communications and Information Theory10(1-2):1–246, 2013.
- Robbins, H. Some aspects of the sequential design of experiments. Bulletin of the American Mathematical Society 58(5):527–535, 1952.
- Serfozo, R. Basics of applied stochastic processes Springer Science & Business Media, 2009.
- Soma, T. and Yoshida, Y. A generalization of submodular cover via the diminishing return property on the integer lattice. In Advances in Neural Information Processing Systemspp. 847–855, 2015.
- Soma, T., Kakimura, N., Inaba, K., and Kawarabayashi, K.-i. Optimal budget allocation: Theoretical guarantee and efficient algorithm. In International Conference on Machine Learningpp. 351–359, 2014.
- Wang, C., Chen, W., and Wang, Y. Scalable influence maximization for independent cascade model in large-scale social networks. Data Mining and Knowledge Discovery 25(3):545–576, 2012.

Wang, Q. and Chen, W. Improving regret bounds for combinatorial semi-bandits with probabilistically triggered arms and its applications. *Advances in Neural Information Processing Systems* pp. 1161–1171, 2017.

Wang, S. and Chen, W. Thompson sampling for combinatorial semi-bandits. *International Conference on Machine Learning* pp. 5114–5122, 2018.

Wilder, B., Immorlica, N., Rice, E., and Tambe, M. Maximizing influence in an unknown social network. *Thirty-Second AAAI Conference on Artificial Intelligence* 2018.

## Supplementary Material

The supplementary material is organized as follows.

We first discuss how we handle the multiple random walkers in Section A. We provide proofs and examples for properties of the visiting probability in Section B. Proofs of the optimization for overlapping MuLaNE are provided in Section C. We provide the detailed budget effective greedy algorithm with partial enumeration (BEGE) and its analysis in Section D. Proofs of the optimization for non-overlapping MuLaNE are provided in Section E. We state the detailed analysis of online learning for overlapping MuLaNE in Section F. We state the detailed analysis of online learning for non-overlapping MuLaNE in Section G. Supplemental experiments are provided in Section H.

### A. Handling the Multiple Random Walkers

(a) Overlapping. (b) Non-overlapping.

Figure 2. Two types of multi-layered networks.

We can handle the scenario where each layer is explored by multiple random walkers using the bipartite coverage model. Without loss of generality, suppose we want to add a new random walker  $W_1^0$  to the layer  $L_1$ . A new node  $W_1^0$  can be easily added to representing the new random explorer, and new edges  $(W_1^0, u) | u \in V_1, g$  are added to  $E^0$ . Thus, we can use the same algorithms and analysis to solve the optimization problem on the newly constructed bipartite coverage graph.

### B. Proofs and Examples for Properties of the Visiting Probability $P_{i;u}(k_i)$

#### B.1. Starting From Arbitrary Distributions

Lemma 1.  $g_{i;u}(k_i) \geq 0$  for any  $i \in [m], u \in V_i, k_i \in \mathbb{Z}_{>0}$ .

Proof. For analysis, we use the following equivalent formulation for  $k_i \geq 2$ , (trivially  $P_{i;u}(0) = 0$  and  $P_{i;u}(1) = \mathbb{1}_{i;u}$ ),

$$P_{i;u}(k_i) = ((\mathbb{1}_{i;u})^{\mathbb{1}_{i;u}} (\mathbb{1}_{i;u})_{uk_i} + \mathbb{1}_{i;u}); \quad (12)$$

where  $(\mathbb{1}_{i;u})_{uk_i} = (I + T_i(u) + \dots + T_i(u)^{k_i-2})(p_i)_u$ ,  $(p_i)_u = (P_i[\cdot; u])_u$ ,  $T_i(u) = (P_i)_{u; u}$ , (we use Eq.

(2.67) in (Kijima, 1997) to derive  $\mathbb{1}_{i;u}$ ). Note that  $\mathbb{1}_{i;u} \in \mathbb{R}^{n-1}$  is the vector obtained by deleting the  $i$ -th element from  $\mathbb{1} \in \mathbb{R}^n$ , and  $P_{i;u} \in \mathbb{R}^{(n-1) \times (n-1)}$  is the matrix obtained by deleting the  $i$ -th row and the  $i$ -th column from  $P \in \mathbb{R}^{n \times n}$ .

We further derive the marginal gain  $g_{i;u}(k_i)$  at step  $k_i \geq 2$  (trivially  $g_{i;u}(1) = \mathbb{1}_{i;u}$ ) as,

$$g_{i;u}(k_i) = ((\mathbb{1}_{i;u})^{\mathbb{1}_{i;u}} T_i(u)^{k_i-2} (p_i)_u); \quad (13)$$

Now, we can show  $g_{i;u}(k_i)$  is non-negative because any element of  $(\mathbb{1}_{i;u})^{\mathbb{1}_{i;u}}$ ,  $T_i(u)$  and  $(p_i)_u$  are non-negative, which means  $P_{i;u}(k_i)$  is non-decreasing with respect to step  $k_i$ .

An example showing  $g_{i;u}(k_i)$  is not monotone: Consider a path  $P_3$  with three nodes as the  $i$ -th layer  $G(V_i; E_i)$ , where  $V_i = \{u, v, w\}$  and  $E_i = \{(u, v); (v, u); (v, w); (w, v)\}$ . If the  $W_i$  always starts from the left-most node and chooses the right-most node as our target node. Then  $g_{i;u}(1) = g_{i;u}(2) = 0$  but  $g_{i;u}(3) > 0$  since at least three steps are needed to visit the node  $w$ , which shows that  $g_{i;u}(k_i)$  is not always non-increasing.

#### B.2. Starting From the Stationary Distribution

Lemma 2.  $g_{i;u}(k_i + 1) \geq g_{i;u}(k_i) \geq 0$  for any  $i \in [m], u \in V_i, k_i \in \mathbb{Z}_{>0}$ , if  $\mathbb{1}_{i;u} = \mathbb{1}_i$ , where  $\mathbb{1}_i = \mathbb{1} - \mathbb{1}_{i;u}$ .

Proof. Consider any layer  $L_i$  with transition probability matrix  $P_i$ , if we start from the stationary distribution  $\pi_i$ , the probability that node  $u \in V_i$  is ever visited in the first  $k_i$  steps is  $P_{i;u}(k_i) = ((\mathbb{1}_{i;u})^{\mathbb{1}_{i;u}} (\mathbb{1}_{i;u})_{uk_i} + \mathbb{1}_{i;u})$ , where  $(\mathbb{1}_{i;u})_{uk_i} = (I + T_i(u) + \dots + T_i(u)^{k_i-2})(p_i)_u$ ,  $T_i(u) = (P_i)_{u; u}$ ,  $(p_i)_u = (P_i[\cdot; u])_u$ . Then the marginal gain for node  $u$  is  $g_{i;u}(k_i) = P_{i;u}(k_i) - P_{i;u}(k_i - 1) = ((\mathbb{1}_{i;u})^{\mathbb{1}_{i;u}} T_i(u)^{k_i-2} (p_i)_u$  for  $k_i \geq 2$  and  $g_{i;u}(k_i) = \mathbb{1}_{i;u}$  when  $k_i = 1$ .

Define the margin of the marginal gain as  $\mathbb{1}_{i;u}(k_i) = g_{i;u}(k_i + 1) - g_{i;u}(k_i)$ . When  $k_i = 1$ ,  $\mathbb{1}_{i;u}(1) = g_{i;u}(2) - g_{i;u}(1) = ((\mathbb{1}_{i;u})^{\mathbb{1}_{i;u}} (p_i)_u - \mathbb{1}_{i;u}) = \mathbb{1}_{i;u} P[u; u] - \mathbb{1}_{i;u} = \mathbb{1}_{i;u} P[u; u] - \mathbb{1}_{i;u} \leq 0$ . When  $k_i \geq 2$ ,  $\mathbb{1}_{i;u}(k_i) = g_{i;u}(k_i + 1) - g_{i;u}(k_i) = ((\mathbb{1}_{i;u})^{\mathbb{1}_{i;u}} (T_i(u)^{k_i-1} - T_i(u)^{k_i-2})(p_i)_u$ . Because  $\mathbb{1}_i = \mathbb{1} - \mathbb{1}_{i;u}$ , we have  $(\mathbb{1}_{i;u})^{\mathbb{1}_{i;u}} = (\mathbb{1}_i)^{\mathbb{1}_i} P_i(u) = \mathbb{1}_{i;u} (q_i)_u^{\mathbb{1}_{i;u}}$ , where  $(q_i)_u^{\mathbb{1}_{i;u}} = (P_i[u; \cdot])_u$ . Thus,  $\mathbb{1}_{i;u}(k_i) = ((\mathbb{1}_{i;u})^{\mathbb{1}_{i;u}} (q_i)_u^{\mathbb{1}_{i;u}} T_i(u)^{k_i-2} (p_i)_u) \geq 0$  because any element in  $(\mathbb{1}_{i;u})^{\mathbb{1}_{i;u}}$ ,  $(q_i)_u^{\mathbb{1}_{i;u}}$  and  $P_i(u)$  is non-negative.

More interestingly, the stationary distribution  $\pi_i$  is the only starting distribution for  $P_i$  such that any  $u \in V_i; k_i \in \mathbb{Z}_{>0}$ ,  $g_{i;u}(k_i + 1) \geq g_{i;u}(k_i) \geq 0$  when  $P_i$  is ergodic and there are no self loops in  $G$ .

Proof. With a little abuse of the notation, we use to denote the transition probability matrix  $P \in \mathbb{R}^{n \times n}$  and let  $P_{i,j}$  be the element in the  $i$ -th row and the  $j$ -th column. Since  $P$  is ergodic, according to Theorem 54 in (Serfozo, 2009), there exists a unique and positive stationary distribution  $\pi = (\pi_1, \dots, \pi_n)$ , i.e.,  $P\pi = \pi$  and  $\pi_j > 0$ . Any starting distribution  $\mu$  can be represented by  $\mu + \epsilon$ , where  $\epsilon = (\epsilon_1, \dots, \epsilon_n)$  is a perturbation vector, and  $\epsilon_j = \epsilon_j - 1$  for  $j \in [2, n]$ . We have the following equation for the margin of marginal gains for node  $u$  in the first two steps,

$$\begin{aligned} \mu_u &= g_{i;u}(1) - g_{i;u}(2) = (\mu_u + \epsilon_u) \sum_{j \in u} (j + u) P_{j;u} \\ &= \mu_u P_{uu} + \epsilon_u \sum_{j \in u} P_{j;u} \end{aligned}$$

Since  $G$  has no self loops, i.e.,  $P_{u;u} = 0$ , we have  $\mu_u = \sum_{j \in u} P_{j;u} \epsilon_j$ , and we can verify that  $\sum_{u \in [n]} \mu_u = \sum_{u \in [n]} \sum_{j \in u} P_{u;j} \epsilon_j = 0$ . Therefore, we have to guarantee  $\mu_u = 0$  for all  $u$ , otherwise there will exist a node  $u$  such that  $\mu_u < 0$ . To ensure  $\mu_u = 0$ , we need to ensure  $\epsilon_j = 0$  for all  $j$ . Again, since there exists a unique and positive stationary distribution for  $P$  and  $\sum_{u \in [n]} \mu_u = 0$ , we can derive  $\epsilon_j = 0$ , where  $\epsilon_j$  has to be 0. Therefore, combined with Lemma 2, the stationary distribution is the only starting distribution such that for any  $i$ , any  $u \in V_i$ ;  $k_i \in Z_{>0}$ ,  $g_{i;u}(k_i + 1) - g_{i;u}(k_i) = 0$ .

### C. Proofs of Of ine Optimization for Overlapping MuLaNE

#### C.1. Starting from the arbitrary distribution

Lemma 3. For any network  $G$ , distribution  $\pi$  and weights  $r_{G; \cdot}(\cdot) : Z^m_0 \rightarrow \mathbb{R}$  is monotone and submodular.

Proof. By definition, we need to show  $r_{G; \cdot}(x \wedge y) + r_{G; \cdot}(x \vee y) \geq r_{G; \cdot}(x) + r_{G; \cdot}(y)$  for any  $x, y \in Z^m_0$ , and  $r_{G; \cdot}(x) \geq r_{G; \cdot}(y)$  if  $x \geq y$ .

(Monotonicity.) By Eq. (7),  $r_{G; \cdot}(k) = \sum_{v \in V} \sum_{i \in [2, m]} (1 - P_{i;v}(k_i)) \prod_{j \in v} (P_{j;v}(x_j + 1) - P_{j;v}(k_j))$ . Since  $P_{j;v}(x_j) \geq P_{j;v}(k_j)$ , for any  $j \in [2, m]$ ;  $v \in V$ ,  $x_j \geq k_j$ , we have

$$\begin{aligned} r_{G; \cdot}(x \vee y) - r_{G; \cdot}(x) &= \sum_{v \in V} \sum_{i \in [2, m]} ((1 - P_{i;v}(x_i)) - (1 - P_{i;v}(y_i))) \prod_{j \in v} (P_{j;v}(x_j + 1) - P_{j;v}(k_j)) \\ &\geq 0; \end{aligned}$$

for any  $x \in Z^m_0$ ;  $j \in [2, m]$ . Then we can use the above inequality repeatedly to show that  $r_{G; \cdot}(x) \geq r_{G; \cdot}(y)$  when  $x \geq y$ .

(Submodularity.) For submodular property, it is sufficient to prove  $\sum_{i \in [2, m]} (1 - P_{i;v}(k_i))$  is submodular for any  $v \in V$ , because a positive weighted sum of submodular function is still submodular. We will rely on the following lemma to prove the submodularity of  $\sum_{i \in [2, m]} (k_i)$ .

Lemma 5. Function  $f : Z^m_0 \rightarrow \mathbb{R}$  is submodular if and only if

$$\begin{aligned} f(x + e_i) - f(x) &\leq f(x + e_j + e_i) - f(x + e_j); \\ \text{for any } x \in Z^m_0 \text{ and } i \neq j. \end{aligned} \tag{14}$$

Proof of Lemma 5. (If part.) We first prove, if inequality (14) holds, the following inequality holds,

$$f(x + e_i) - f(x) \leq f(y + e_i) - f(y) \tag{15}$$

for any  $x \leq y$  and  $i \in [2, m]$  such that  $x_i = y_i$ .

Let  $I_0 = \{i \in [2, m] : x_i = y_i\}$ ;  $I_1 = \{i \in [2, m] : x_i < y_i\}$ . For any  $x \leq y$ , we denote the elements in  $I_1$  by  $i_1, \dots, i_s$  and write  $y = x + \sum_{j=1}^s e_{i_j}$ , where  $i_j = y_{i_j} - x_{i_j}$ . For any  $i \in I_0$ , we have

$$\begin{aligned} f(x + e_i) - f(x) &\leq f(x + e_{i_1} + e_i) - f(x + e_{i_1}) \\ &\leq f(x + 2e_{i_1} + e_i) - f(x + 2e_{i_1}) \\ &\vdots \\ &\leq f(x + e_{i_1} + e_{i_1} + e_i) - f(x + e_{i_1} + e_{i_1}) \\ &\leq f(x + e_{i_1} + e_{i_1} + e_{i_2} + e_i) - f(x + e_{i_1} + e_{i_1} + e_{i_2}) \\ &\leq f(x + e_{i_1} + e_{i_1} + e_{i_2} + e_{i_2} + e_i) - f(x + e_{i_1} + e_{i_1} + e_{i_2} + e_{i_2}) \\ &\vdots \\ &\leq f(x + \sum_{j=1}^s e_{i_j} + e_i) - f(x + \sum_{j=1}^s e_{i_j}) \\ &= f(y + e_i) - f(y); \end{aligned} \tag{16}$$

Then for any  $i \in I_0$  and  $a \in Z^m_0$ , we have

$$\begin{aligned} f(x + a + e_i) - f(x) &\leq \sum_{j=1}^a (f(x + e_j + e_i) - f(x + (j-1)e_j + e_i)) \\ &\leq \sum_{j=1}^a (f(y + e_j + e_i) - f(y + (j-1)e_j + e_i)) \\ &= f(y + a + e_i) - f(y); \end{aligned} \tag{17}$$

The first inequality holds because of Inequality (16), the fact  $x + (j-1)e_j \leq y + (j-1)e_j$  and  $x_j + j - 1 = y_j + j - 1$ .

Then for any  $x, y \in Z^m_0$ , let  $I_2 = \{i \in [2, m] : x_i > y_i\}$  =

$f(x_1, \dots, x_s)$ . We have

$$\begin{aligned} f(x) &= f(x \wedge y) \\ &= \sum_{i=1}^s f(x \wedge y + (x_i, y_i)) \\ &= \sum_{i=1}^s f(x \wedge y + (x_i, y_i)) \\ &= \sum_{i=1}^s f(y + (x_i, y_i)) \\ &= f(x \vee y) = f(y) \end{aligned}$$

The inequality is derived from inequality (17) because  $\sum_{j=1}^{l-1} (x_{i_j}, y_{i_j}) \leq x \wedge y + \sum_{j=1}^{l-1} (x_{i_j}, y_{i_j})$  for any  $0 \leq l \leq s$  and  $(x \wedge y)_{i_l} = y_{i_l}$ , which concludes the if part.

(Only if part.) Assume  $f$  is submodular, let  $a = x + (i, b)$ ;  $b = x + (j, i) \in j$ , we have  $a \wedge b = x + (i, i)$ ,  $a \vee b = x + (i, j)$ .  $f(a \wedge b) = f(x)$ ,  $f(a \vee b) = f(x + (i, j))$ .  $f(a \wedge b) + f(a \vee b) = f(x) + f(x + (i, j))$ .

Then, by Lemma 5 and the explicit formula of the reward function given by Eq. (7), we can prove  $g(x; v) = g(x + (i, v)) + g(x + (j, v)) - g(x + (i, j, v))$  for any  $x \in Z^m_0, v \in V$  and  $i, j \in [m]$ , where  $g(x; v) = \sum_{i \in [m]} (1 - P_{i,v}(x_i))$ . This holds due to the fact that the left hand side equals to  $\sum_{i \in [m]} (1 - P_{i,v}(x_i)) - (1 - P_{i,v}(k_j)) - (1 - P_{i,v}(k_i + 1)) + (1 - P_{i,v}(k_i))$  and the right hand side equals to  $\sum_{i \in [m]} (1 - P_{i,v}(x_i)) - (1 - P_{i,v}(k_j + 1)) - (1 - P_{i,v}(k_i + 1)) + (1 - P_{i,v}(k_i))$ , and the left hand side is larger or equal to the right hand side because  $(1 - P_{i,v}(k_j)) \geq (1 - P_{i,v}(k_j + 1))$ . By summation over all nodes  $i \in [m]$  with node weights  $v \in [0, 1]$ , we can prove the reward function is submodular.

Theorem 1. Algorithm 1 obtains a  $(1 - e^{-\epsilon})$ -approximate solution, where  $\epsilon$  is the solution of equation  $e^{-\epsilon} = 2^{-\epsilon}$ , to the overlapping MuLaNE problem.

Proof. For theoretical analysis, we first give a modified version of Alg. 1 in Alg. 4. Both algorithms provide the same solution  $k$  given the same problem instance  $(G; \omega; B; c)$ . To see this fact, Alg. 4 considers invalid tentative allocations  $(i, b)$  (adding it will exceed the total budget constraint) and remove them in line 13 of Alg. 4, while Alg. 1 only

Algorithm 4 Equivalent Budget Effective Greedy Algorithm (BEG) for the Overlapping MuLaNE.

Input: Network  $G$ , starting distributions  $\omega$ , node weights  $\omega$ , budget  $B$ , constraints  $c$ .

Output: Budget allocation  $k$ .

- 1: Compute visiting probabilities  $(P_{i,u}(b))_{i \in [m]; u \in V; b \in [c_i]}$  according to Eq. (9).
- 2:  $k = \text{BEG}((P_{i,u}(b))_{i \in [m]; u \in V; b \in [c_i]}, \omega, B, c)$ .
- 3: Procedure  $\text{BEG}((P_{i,u}(b))_{i \in [m]; u \in V; b \in [c_i]}, \omega, B, c)$
- 4: Let  $k := (k_1, \dots, k_m) \leftarrow (0, \dots, 0)$ .
- 5: Let  $Q = \{(i, b_i) \mid i \in [m]; 1 \leq b_i \leq c_i\}$ .
- 6: while  $K > 0$  and  $Q \neq \emptyset$ ; do
- 7:  $(i, b) \leftarrow \arg \max_{(i,b) \in Q} (i, b; k) = b$ . Eq. (11)
- 8: if  $b \leq K$  then
- 9:  $k_i \leftarrow k_i + b, K \leftarrow K - b$ .
- 10: Modify pairs  $(i, b) \in Q$  to  $(i, b - b)$ .
- 11: Remove pairs  $(i, b) \in Q$  such that  $b = 0$ .
- 12: else
- 13: Remove  $(i, b)$  from  $Q$ .
- 14: end if
- 15: end while
- 16: for  $i \in [m]$  do
- 17: if  $r_{G_i}; (c_i - i) > r_{G_i}; (k_i)$ , then  $k_i \leftarrow c_i$ .
- 18: end for
- 19: return  $k := (k_1, \dots, k_m)$ .

considers valid allocations by directly removing invalid allocations in advance in line 10 of Alg. 1.

With a little abuse of the notation, we use  $k$  to represent the reward  $r_{G_i}; (k)$  of a given problem instance. Let  $k \in Z^m_0$  denote the optimal budget allocation and  $k^s \in Z^m_0$  denote the budget allocation before entering the  $s$ -th iteration of the while loop (line 6-13) in Alg. 4. After entering the  $s$ -th iteration, the algorithm tries to extend the current budget allocation  $k^s$  by choosing the pair  $(i, b)$  in line 7, which we denote as  $(i^s; b^s)$ . Let  $s$  be the first iteration we can not extend the current solution, i.e.,  $k^s = k^{s+1}$  and  $k^j < k^{j+1}$  for  $j = 1, \dots, s - 1$ . If we can always extend the current solution, we set  $s$  to be  $B + 1$ . For analysis, we temporarily add  $(i^s; b^s)$  (in the algorithm, this pair is removed by line 13 in the  $s$ -th iteration) to form a "virtual" budget allocation  $k^{s+1} = k^s + b^s \cdot e_{i^s}$ . Let  $k_g \in Z^m_0$  denote the solution returned by Alg. 4.

We first introduce lemmas describing two important properties given by the submodularity over the integer lattice.

Lemma 6. (Soma et al., 2014). Let  $\ell : Z^m_0 \rightarrow \mathbb{R}$  be a

submodular function. For any  $y \in Z^m_0$ , we have,

$$f(x \cup y) - f(x) \leq \sum_{i \in \text{supp}(y \setminus x)} (f(x \cup (y_i, x_i)) - f(x)) \quad (18)$$

Lemma 7. (Soma et al., 2014). Let  $r : Z^m_0 \rightarrow \mathbb{R}$  be a monotone submodular function. For any  $y \in Z^m_0$  with  $x \cup y$  and  $i \in [m]$  we have,

$$f(x \cup k_i) - f(x) \leq f(y \cup k_i) - f(y) \quad (19)$$

Then, we have the following lemma.

Lemma 8. For  $j = 1, \dots, s$ ,

$$r(k^{j+1}) \leq (1 - \frac{b_j}{B})r(k^j) + \frac{b_j}{B}r(k) \quad (20)$$

Proof of Lemma 8. This is because

$$\begin{aligned} r(k) &= r(k \cup k^j) - \sum_{i \in \text{supp}(k \setminus k^j)} (r(k^j \cup (k_i, k_i^j)) - r(k^j)) \\ &= r(k^j) + \sum_{i \in \text{supp}(k \setminus k^j)} (r(k^j \cup (k_i, k_i^j)) - r(k^j)) \\ &= r(k^j) + \sum_{i \in \text{supp}(k \setminus k^j)} (r(k^j \cup (k_i, k_i^j)) - r(k^j)) \quad (\text{Let } k_i = k_i, k_i^j = k_i^j) \\ &= r(k^j) + \sum_{i \in \text{supp}(k \setminus k^j)} \frac{r(k^{j+1}) - r(k^j)}{b_j} \\ &= r(k^j) + B \frac{r(k^{j+1}) - r(k^j)}{b_j} \quad (21) \end{aligned}$$

The second inequality comes from Lemma 6, the third inequality holds because of the greedy procedure in line 7 and the last inequality holds because  $\sum_{i \in \text{supp}(k \setminus k^j)} b_i \leq B$ . By rearranging terms, Inequality (20) holds.

Next, We can prove the following lemma.

Lemma 9. For  $l = 1, \dots, s$ ,

$$\begin{aligned} r(k^{l+1}) &\leq r(k^1) \prod_{j=1}^l (1 - \frac{b_j}{B}) \\ &\quad + r(k) \sum_{j=1}^l \prod_{i=1}^{j-1} (1 - \frac{b_i}{B}) \frac{b_j}{B} \quad (22) \end{aligned}$$

Proof of Lemma 9. We can prove this lemma by induction on  $l$ . When  $l = 1$ , the lemma holds due to Lemma 8. Assume that the lemma holds for  $l$ , we have the following inequality holds,

$$\begin{aligned} r(k^{l+1}) &\leq (1 - \frac{b_l}{B})r(k^l) + \frac{b_l}{B}r(k) \\ &\leq (1 - \frac{b_l}{B})r(k^1) \prod_{j=1}^{l-1} (1 - \frac{b_j}{B}) \\ &\quad + (1 - \frac{b_l}{B})r(k) \sum_{j=1}^{l-1} \prod_{i=1}^{j-1} (1 - \frac{b_i}{B}) \frac{b_j}{B} + \frac{b_l}{B}r(k) \\ &= r(k^1) \prod_{j=1}^l (1 - \frac{b_j}{B}) + r(k) \sum_{j=1}^l \prod_{i=1}^{j-1} (1 - \frac{b_i}{B}) \frac{b_j}{B} \quad (23) \end{aligned}$$

where the first inequality is due to Lemma 8 by setting  $l = l$ , and the second inequality is by the assumption for  $l$ . By induction, Lemma 9 holds.

We then consider the following cases.

Case 1. Suppose the total budget used for is larger or equal to  $B$ , i.e.,  $\sum_{j=1}^s b_j \geq B$ , where  $\alpha \in [0, 1]$ .

We have the following inequality.

$$r(k^s) \leq (1 - e^{-\alpha})r(k) \quad (24)$$

This is due to Lem. 9 by setting  $l = s - 1$ , combined with the fact that  $r(k^1) = 0$  and  $\sum_{j=1}^{s-1} (1 - \frac{b_j}{B}) \leq e^{-\alpha}$ . The later fact holds because,

$$\begin{aligned} \sum_{j=1}^{s-1} (1 - \frac{b_j}{B}) &= (s-1) \sum_{j=1}^{s-1} \frac{1}{s-1} \log(1 - \frac{b_j}{B}) \\ &\leq (s-1) \log \sum_{j=1}^{s-1} \frac{1}{s-1} \frac{b_j}{B} \\ &= (s-1) \log(1 - \frac{\alpha}{s-1}); \end{aligned}$$

where the first inequality holds because of the Jensen's Inequality (Raginsky et al., 2013) and the second inequality holds because  $\sum_{j=1}^{s-1} \frac{b_j}{B} = \alpha$ . Then we can easily check  $\sum_{j=1}^{s-1} (1 - \frac{b_j}{B}) \leq (1 - \frac{\alpha}{s-1})^{s-1} \leq e^{-\alpha}$ .

Case 2. Suppose the total budget  $\sum_{j=1}^s b_j \leq B$ . Then, we have  $\alpha > (1 - \frac{\alpha}{s-1})B$ . We can prove the following inequality holds,

$$r(k_g) \leq (1 - \frac{1}{2})r(k) \quad (25)$$

This is due to Inequality (21), we have

$$r(k) \leq r(k^s) + B \frac{r(k^{s+1}) - r(k^s)}{b^s}$$

$$r(k^s) + \frac{r(k^{s+1}) - r(k^s)}{1}$$

$$(1 + \frac{1}{1})r(k_g);$$

where the second inequality is due to  $b > (1 - \epsilon)B$  and the last equality is due to the fact that  $r(k^{s+1}) - r(k^s) \leq r(k_g) - r(k^s)$ . To see the above fact, we assume without loss of generality the pair  $(i^s, b^s)$  improves  $k^s$  towards the optimal budget allocation, i.e.  $k_{i^s}^{s+1} = b^s + k_{i^s}^s - k_{i^s}^s$ . Otherwise, if  $b^s + k_{i^s}^s > k_{i^s}^s$ , we can safely delete  $(i^s, b^s)$  in the queue  $Q$  and does not affect the greedy solution, the optimal solution and the analysis. Let  $r_{\max} = \max_{i \in [m]} r(G_i)$ , we have  $r(k^{s+1}) - r(k^s) \leq r(k^s - k_{i^s} - i^s) - r(k^s) = r(k_{i^s} - i^s) - r(0) \leq r_{\max} - r(k_g)$ . Also, we can obtain that  $r(k^s) \leq r(k_g)$  since  $k^s \leq k_g$ . By rearranging the terms, Inequality (25) holds.

Combining Inequality (24) and (25), we have

$$r(k_g) \leq \min_{\epsilon \in [0;1]} \max(1 - \epsilon; 1 - \frac{1}{2})r(k)$$

$$(1 - \epsilon)r(k);$$

where  $\epsilon$  is the solution for equation  $\epsilon = 2 - \frac{1}{\epsilon}$ .

### C.2. Starting From the Stationary Distribution

Lemma 4. For any network  $G$ , stationary distributions and node weights  $\mu$ , function  $r_{G; \mu; \cdot}(\cdot) : Z^{m_0} \rightarrow \mathbb{R}$  is monotone and DR-submodular.

Proof. By definition, we need to show  $r_{G; \mu; \cdot}(y + j) \geq r_{G; \mu; \cdot}(y) + r_{G; \mu; \cdot}(x + j) - r_{G; \mu; \cdot}(x)$ , and  $r_{G; \mu; \cdot}(x) \geq r_{G; \mu; \cdot}(y)$  for any  $x \geq y, j \in [m]$ .

We first introduce the following lemma to help us to prove the DR-submodularity.

Lemma 10. Function  $f$  is DR-submodular if and only if

$$f(x + i) - f(x) \geq f(x + j + i) - f(x + j);$$

for any  $x \in Z^{m_0}$ .

(26)

Proof of Lemma 10 (If part.) We can easily check this direction holds by using the similar argument for Inequality (16), where the only difference is we consider  $\mathcal{X} \subseteq [m]$  instead of  $\mathcal{I} \subseteq I_0$ .

(Only if part.) We can set  $y := x^0 + j, x := x^0 + i$ , for any  $i, j \in [m]; x^0 \in Z^{m_0}$ , and use Inequality (15) to

show the only if part holds.

Since  $r_{G; \mu; \cdot}(k)$  is submodular for arbitrary starting distributions  $\mu$ , Inequality (14) holds. Then consider any layer  $\mathcal{X} \subseteq [m]$  and budget allocation  $x \in Z^{m_0}$ , it is sufficient to show  $r_{G; \mu; \cdot}(x + 2j) - r_{G; \mu; \cdot}(x + j) \geq r_{G; \mu; \cdot}(x + j) - r_{G; \mu; \cdot}(x)$ . Since the left hand side minus the right hand side equals to  $\sum_{v \in V} \sum_{i \in [m]; i \in \mathcal{X}} (1 - P_{i;v}(x_i)) (P_{j;v}(x_j) - 2P_{j;u}(x_j + 1) + P_{j;v}(x_j + 2))$ , we only need to show  $\sum_{v \in V} \sum_{i \in [m]; i \in \mathcal{X}} (1 - P_{i;v}(x_i)) (g_{j;v}(x_j + 2) - g_{j;v}(x_j + 1)) \geq 0$ . The above inequality holds because of Lemma 2.

Theorem 2. Algorithm 2 obtains a  $(1 - \epsilon)$ -approximate solution to the overlapping MuLaNE with the stationary starting distributions.

Proof. We can observe that line 7 always select the pair  $(i; b)$  with  $b = 1$  because  $r_{G; \mu; \cdot}(k + b \cdot i) - r_{G; \mu; \cdot}(k + (b - 1) \cdot i) \geq r_{G; \mu; \cdot}(k + i) - r_{G; \mu; \cdot}(k)$  for arbitrary  $k, i$  and  $b$ . Thus, we have  $\epsilon = B + 1$  and by the similar argument for Inequality (24), we have  $r(k_g) = r(k^{B+1}) = (1 - \epsilon)r(k)$ , which completes the proof.

### C.3. Efficiently Evaluating the Reward Function

One key issue to derive the budget allocation is to efficiently evaluate the reward function  $r_{G; \mu; \cdot}(k)$  and its marginal gains  $(i; b; k)$ . Since we need to repetitively use the visiting probabilities  $P_{i;u}(k_i)$ , we pre-calculate  $P_{i;u}(k_i)$  in  $O(m \cdot k \cdot n_{\max}^3)$  time in our algorithms based on Eq. (9), where  $n_{\max} = \max_i |V_i|$ . For the overlapping case, given the current budget allocation  $k$ , we maintain a value  $\rho_u = \sum_{i \in [m]} (1 - P_{i;u}(k_i))$  for each node  $u \in V$ . The marginal gain  $(i; b; k) = \sum_{u \in V} \rho_u \cdot (1 - \frac{1 - P_{i;u}(k_i + b)}{1 - P_{i;u}(k_i)}) = b$  can be evaluated in  $O(n_{\max})$  time. Then, we update all  $\rho_u = \rho_u \cdot \frac{1 - P_{i;u}(k_i + b)}{1 - P_{i;u}(k_i)}$  in  $O(n_{\max})$  after we allocate more budgets to layer  $\mathcal{X}$ . Therefore, we can use  $O(n_{\max})$  in total to evaluate  $(i; b; k)$ . In practice, lazy evaluation (Krause et al., 2008) and parallel computing can be used to further accelerate our algorithm.

## D. Budget Effective Greedy Algorithm With Partial Enumeration and Its Analysis

### D.1. Algorithm

The algorithm is shown in Alg. 5.

### D.2. Analysis

Theorem 4. The Algorithm 5 obtains a  $(1 - \epsilon)$ -approximate solution to the overlapping MuLaNE with ar-



**Algorithm 5 Budget Effective Greedy Algorithm with Partial Enumeration (BEGE).**

Input: Graph  $G$ , starting distributions  $\{p_i\}$ , node weights  $\{w_i\}$ , budget  $B$ , constraints  $\mathcal{C}$ .

Output: Budget allocation  $k$ .

```

1: Compute visiting probabilities  $\{P_{i;u}(b)\}_{i \in [2:m]; u \in \mathcal{V}; b \in [c_i]}$  according to Eq. (7).
2:  $k \leftarrow \text{BEGE}(\{P_{i;u}(b)\}_{i \in [2:m]; u \in \mathcal{V}; b \in [c_i]}, \{w_i\}, B, \mathcal{C})$ .
3: Procedure BEGE( $\{P_{i;u}(b)\}_{i \in [2:m]; u \in \mathcal{V}; b \in [c_i]}, \{w_i\}, B, \mathcal{C}$ )
4: Let  $k_{\max} \leftarrow 0$ .
5:  $S \leftarrow \{k = (k_1; \dots; k_m) \mid 0 \leq k_i \leq c_i; \sum_{i \in [2:m]} k_i \leq B; \text{If } k_i > 0 \text{ g. } \dots\}$ . S contains all partial solutions which allocate partial budgets to at most three layers
6: for  $k \in S$  do
7:    $K \leftarrow B - \sum_{i \in [2:m]} k_i$ .
8:   Let  $Q \leftarrow \{(i; b_i) \mid i \in [2:m]; 1 \leq b_i \leq c_i\}$ .
9:   while  $K > 0$  and  $Q \neq \emptyset$  do
10:     $(i; b) \leftarrow \arg \max_{(i; b) \in Q} (i; b; k) = b$ . Eq. (11)
11:    if  $b \leq K$  then
12:       $k_i \leftarrow k_i + b, K \leftarrow K - b$ .
13:      Modify all pairs  $(i; b) \in Q$  to  $(i; b - b)$ .
14:      Remove all pairs  $(i; b) \in Q$  such that  $b = 0$ .
15:    else
16:      Remove  $(i; b)$  from  $Q$ .
17:    end if
18:  end while
19:  if  $r_{G; \cdot}(k) > r_{G; \cdot}(k_{\max})$ , then  $k_{\max} \leftarrow k$ .
20: end for
21: end Procedure

```

bitrary starting distributions.

Proof. Suppose that we start from a partial solution  $k^0 \in \mathbb{Z}^m \in \mathbb{0}$ . Let us first reorder the optimal solution according to a non-increasing marginal gain ordering. Namely, the marginal gain of the pair  $(s_1; b_1)$  with respect to the empty pair is the highest among all other pairs, the marginal gain of the pair  $(s_2; b_2)$  with respect to the solution consisting of pair  $(s_1; b_1)$  is the highest among all remaining pairs, and so on. To be concrete,  $k = (k_1; \dots; k_m)$  is selected to maximize the following equation:

$$\sum_{j=1}^m \sum_{i=1}^m \sum_{b_j=1}^{c_j} \sum_{b_i=1}^{c_i} (k_j + b_j; k_i + b_i) : (27)$$

Then, we try to bound the "virtual" marginal gain  $\Delta r(k^{s+1}) - \Delta r(k^s)$ , and recall that  $k^s$  is the first iteration we can not extend the current solution. Without loss of generality, we assume the virtual pair  $(i^s; b^s)$  improves  $k^s$  towards the

optimal budget allocation, i.e.,  $k_i^s + k_{i^s}^s > k_{i^s}^s$ . Otherwise, if  $k_i^s + k_{i^s}^s \leq k_{i^s}^s$ , we can safely delete  $(i^s; b^s)$  in the queue  $Q$  and does not affect the greedy solution, the optimal solution and the analysis.

If we start from the initial solution  $k^1$  such that  $k^1$  matches  $(s_1; b_1); (s_2; b_2); (s_3; b_3)$ , i.e.,  $k_{s_i}^1 = b_i$  for  $i = 1; 2; 3$ , we have  $\Delta r(k^{s+1}) - \Delta r(k^s) = r(k^{s+1}) - r(k^s) - r(k^s - k_{i^s}^s - b_{i^s}^s) + r(k^s) - r(k_{i^s}^s - b_{i^s}^s) + r(0) - r(b_{i^s}^s)$ . Moreover,  $\Delta r(k^s) = r(k^{s+1}) - r(k^s) - r(k^s - k_{i^s}^s - b_{i^s}^s) + r(k^s) - r(k_{i^s}^s - b_{i^s}^s) + r(b_{i^s}^s) - r(b_{i^s}^s + b_{i^s}^s)$ . Similarly, we have  $\Delta r(k^s) = r(k^{s+1}) - r(k^s) - r(k^s - k_{i^s}^s - b_{i^s}^s) + r(k^s) - r(k_{i^s}^s - b_{i^s}^s) + r(b_{i^s}^s) - r(b_{i^s}^s + b_{i^s}^s)$ . By adding above inequalities, we have  $\Delta r(k^1) = 3$ .

Now we can use Lemma 9 by setting  $\mu = \Delta r(k^s) - \Delta r(k^1) = 3$  and the fact  $\sum_{j=1}^m b_j \leq B$  to show  $\Delta r(k^{s+1}) - \Delta r(k^1) + (1 - \epsilon)(\Delta r(k^s) - \Delta r(k^1)) = 3$ . Combining  $\Delta r(k^s) - \Delta r(k^1) = 3$ , we have  $\Delta r(k^s) - \Delta r(k^1) \geq (1 - \epsilon)(\Delta r(k^s) - \Delta r(k^1))$ , which completes the proof.

The time complexity of Alg. 5 is  $O(B^4 m^4 k c k_1 n_{\max} + k c k_1 m n_{\max}^3)$ . This is because the number of all partial solutions is  $|S| = O(B^3 m^3)$ , and for any partial enumeration  $k \in S$ , the time complexity is the same order as the Alg. 1, i.e.,  $O(B k c k_1 m n_{\max})$ , where  $O(n_{\max})$  is the time to evaluate  $(i; b; k)$  as discussed before.

**E. Algorithms and Analysis of Offline Optimization for Non-overlapping MuLaNE**

**E.1. Starting From the Stationary Distribution**

According to Eq. (8) and (10), the reward function can be rewritten as

$$r_{G; \cdot}(k) = \sum_{i \in [2:m]} \sum_{b \in [c_i]} g_i(b); \quad (28)$$

where  $g_i(b)$  represents the layer-level marginal gain when we allocate one more budget (from  $b-1$  to  $b$ ) to layer  $i$ , which is given by

$$g_i(b) = \sum_{v \in \mathcal{V}_i} (P_{i;v}(b) - P_{i;v}(b-1)); \quad (29)$$

From the definition of the reward function, we have two observations. First, budgets allocated to a layer will not affect the reward of other layers because layers are non-overlapping. Second, the layer-level marginal gain  $g_i(b)$  for the  $i$ -th layer is non-increasing with respect to budget  $b$ . Based on these two observations, we propose the myopic greedy algorithm (Alg. 6), which is a slight modification of Alg. 2. Alg. 6 takes the multi-layered network  $G = (G_1; \dots; G_m)$ , starting distribution  $\{p_1; \dots; p_m\}$ , node weights  $\{w_i\} = (w_1; \dots; w_{|V_j|})$  and total budget  $B$  as inputs. It

consists of  $B$  rounds to compute the optimal budget allocation  $k^*$ . In each round, line 7 in Alg. 6 selects the layer with the largest layer-level marginal gain and allocate one budget to layer  $i$  until total  $B$  budgets are used up.

**Algorithm 6 Myopic Greedy Algorithm for the Non-overlapping MuLaNE**

Input: Network  $G$ , starting distributions  $\rho$ , budget  $B$ , constraints  $c$ .  
 Output: Budget allocation  $k$ .  
 1: Compute visiting probabilities  $P_{i;u}(b)$  according to Eq. (8).  
 2: Compute layer-level marginal gain  $g_i(b)$  according to Eq. (29).  
 3:  $k = \text{MG}((g_i(b))_{i \in [m]; b \in [c_i]}, B, c)$   
 4: Procedure  $\text{MG}((g_i(b))_{i \in [m]; b \in [c_i]}, B, c)$   
 5: Let  $k := (k_1, \dots, k_m)$ ,  $K = B$ .  
 6: while  $K > 0$  do  
 7:  $i = \arg \max_{j \in [m]; k_j < c_j} g_j(k_j + 1)$  using the priority queue  
 8:  $k_i = k_i + 1, K = K - 1$ .  
 9: end while return  $k = (k_1, \dots, k_m)$ .  
 10: end Procedure

**Theorem 5.** The Algorithm 6 obtains the optimal budget allocation to the non-overlapping MuLaNE with the stationary starting distributions.

**Proof.** Define a two dimensional array  $M \in \mathbb{R}^{c_1 + c_2 + \dots + c_m \times m}$ , where  $(i; j)$ -th entry is the  $j$ -th step layer-level marginal gain for layer  $i$ , i.e.,  $M[i; j] = \sum_{v \in V_i} g_{i;v}(j)$ , for  $i \in [m]; j \in [c_i]$ . Given the budget allocation  $k = (k_1, \dots, k_m)$ , the expected reward  $r_{G; \rho}(k)$  can be written as the sum of elements in  $M$ , i.e.,  $r_{G; \rho}(k) = \sum_{i=1}^m \sum_{j=1}^{k_i} M[i; j]$ . Because  $g_{i;v}(k_i)$  is non-increasing with respect to  $k_i$ , the element in each row or the layer-level marginal gain is non-increasing. Hence, the greedy method at step  $t$  chooses the  $t$ -th largest element in  $M$ . At steps  $t = B$ , the greedy policy selects all  $B$  largest elements and the corresponding budget allocation maximizes the expected reward  $r_{G; \rho}(k)$ .

Alg. 6 uses  $O(k c k_1 m n_{\max}^3)$  to compute visiting probabilities  $P_{i;u}(k_i)$  and  $O(k c k_1 m n_{\max})$  to compute layer-wise marginal gains  $g_i(b)$ . Then  $(i; 1; k)$  can be evaluated in  $O(\log m)$  using the priority queue, which is repeated for  $B$  iterations. Therefore, the time complexity of Alg. 6 is  $O(B \log m + k c k_1 m n_{\max}^3)$ .

**E.2. Starting From Arbitrary Distributions**

When random explorers start from any arbitrary distribution, Alg. 6 can not obtain the optimal solution. This is because the layer-level marginal gain  $g_i(b)$  is not non-increasing

**Algorithm 7 Dynamic Programming (DP) Algorithm for the Non-overlapping MuLaNE**

Input: Network  $G$ , starting distributions  $\rho$ , node weights  $w$ , budget  $B$ , constraints  $c$ .  
 Output: Budget allocation  $k$ .  
 1: Compute visiting probabilities  $P_{i;u}(b)$  according to Eq. (8).  
 2: Compute layer-level marginal gain  $g_i(b)$  according to Eq. (29).  
 3:  $k = \text{DP}((g_i(b))_{i \in [m]; b \in [c_i]}, B, c)$   
 4: Procedure  $\text{DP}((g_i(b))_{i \in [m]; b \in [c_i]}, B, c)$   
 5: for  $i \in [m]; b \in [c_i]$  do  
 6: Set  $A[i; b] = 0, V[i; b] = 0$ .  
 7: end for  
 8: for  $i = 0$  to  $m - 1$  do  
 9: for  $b \in [B]$  do  
 10:  $j = \arg \max_{j \in [c_{i+1}]} (V[i; b - j] + r_{G; \rho}(j, i+1))$   
 11:  $A[i + 1; b] = A[i; b - j] + j^{-i+1}$ .  
 12:  $V[i + 1; b] = V[i; b - j] + r_{G; \rho}(j, i+1)$ .  
 13: end for  
 14: end for return  $A[m; B]$ .  
 15: end Procedure

with respect to  $b$ . So we have to adopt a more general technique, dynamic programming (DP), to solve the budget allocation problem. The key idea is to keep a DP budget allocation table  $A \in \mathbb{R}^{(m+1) \times (B+1)}$ , where the  $(i; b)$ -th entry  $A[i; b]$  saves the optimal budget allocation by allocating  $b$  budgets to the first  $i$  layers. Another DP table  $V \in \mathbb{R}^{(m+1) \times (B+1)}$  saves the value of the optimal reward, where the  $(i; b)$ -th entry  $V[i; b]$  corresponds to the optimal reward when setting  $A[i; b]$  as the budget allocation. In the  $i^0$ -th outer loop and  $b^0$ -th inner loop in Alg 7, we have already obtained the optimal budget allocation  $A[i^0; b^0]$  with  $i = 0, \dots, i^0, b \in [B]$ , which help us to find the optimal amount of budget to  $(i^0 + 1)$ -th layer (in line 10), and thus we can obtain the  $A[i^0 + 1; b^0]$  and  $V[i^0 + 1; b^0]$  accordingly.

**Theorem 6.** Algorithm 7 obtains the optimal budget allocation to the non-overlapping MuLaNE with an arbitrary starting distribution.

Alg. 7 uses  $O(k c k_1 m n_{\max}^3)$  to compute visiting probabilities  $P_{i;u}(k_i)$  and  $O(k c k_1 m n_{\max})$  to compute layer-wise marginal gains  $g_i(b)$ . In the DP procedure, we can also pre-calculate all rewards  $r_{G; \rho}(j, i+1)$  in  $O(k c k_1 m)$ . Then DP procedure takes  $O(k c k_1 m)$  to find  $j$ , which is repeated for  $B m$  iterations. Therefore, the Alg. 7 has the time complexity of  $O(B k c k_1 m + k c k_1 m n_{\max}^3)$ .

## F. Analysis of Online Learning for Overlapping MuLaNE

### F.1. Proof of 1-Norm Bounded Smoothness.

Property 1. (Monotonicity). The reward  $r_k(\mathbf{k})$  is monotonically increasing, i.e., for any budget allocation any two vectors  $\mathbf{q} = (q_{i;u;b})_{i \in [m], u \in [U], b \in [B]}$ ,  $\mathbf{q}' = (q'_{i;u;b})_{i \in [m], u \in [U], b \in [B]}$  and any node weights  $\mathbf{k}$ , we have  $r_k(\mathbf{k}; \mathbf{q}) \leq r_k(\mathbf{k}; \mathbf{q}')$ , if  $q_{i;u;b} \leq q'_{i;u;b}$  and  $\forall i \in [m], u \in [U], b \in [B]$ .

Proof. According to Eq. (7) and since  $q_{i;u;b} \leq q'_{i;u;b}$ , for  $(i; u; b) \in \mathcal{A}$ , we have  $(1 - \sum_{i \in [m]} (1 - q_{i;u;k_i})) \leq (1 - \sum_{i \in [m]} (1 - q'_{i;u;k_i}))$ . Therefore for any  $\mathbf{u} \in \mathcal{U}$ , we have  $\sum_{i \in [m]} (1 - q_{i;u;k_i}) \leq \sum_{i \in [m]} (1 - q'_{i;u;k_i})$ .

By summing up both sides over  $\mathcal{U}$ , we have  $r_k(\mathbf{q}) \leq r_k(\mathbf{q}')$ .

Property 2. (1-Norm Bounded Smoothness). The reward function  $r_k(\mathbf{k})$  satisfies the 1-norm bounded smoothness condition, i.e., for any budget allocation any two vectors  $\mathbf{q} = (q_{i;u;b})_{i \in [m], u \in [U], b \in [B]}$ ,  $\mathbf{q}' = (q'_{i;u;b})_{i \in [m], u \in [U], b \in [B]}$  and any node weights  $\mathbf{k}$ , we have  $|r_k(\mathbf{k}; \mathbf{q}) - r_k(\mathbf{k}; \mathbf{q}')| \leq \sum_{i \in [m], u \in [U], b \in [B]} |q_{i;u;b} - q'_{i;u;b}|$ .

Proof. The left-hand side:

$$\begin{aligned} & |r_k(\mathbf{k}; \mathbf{q}) - r_k(\mathbf{k}; \mathbf{q}')| \\ &= \sum_{u \in \mathcal{U}} \sum_{i \in [m]} |q_{i;u;k_i} - q'_{i;u;k_i}| \\ &= \sum_{u \in \mathcal{U}} \sum_{i \in [m]} |q_{i;u;k_i} - q'_{i;u;k_i}| \\ &= \sum_{u \in \mathcal{U}} \sum_{i \in [m]} |q_{i;u;k_i} - q'_{i;u;k_i}| \\ &= \sum_{u \in \mathcal{U}} \sum_{i \in [m]} |q_{i;u;k_i} - q'_{i;u;k_i}| \\ &= \sum_{u \in \mathcal{U}} \sum_{i \in [m]} |q_{i;u;k_i} - q'_{i;u;k_i}| \end{aligned}$$

where  $q_{i;u} = (1 - \sum_{i \in [m]} q_{i;u;k_i})$ ;  $q'_{i;u} = (1 - \sum_{i \in [m]} q'_{i;u;k_i})$ .

For the first term, it can be derived that,

$$\begin{aligned} & \sum_{u \in \mathcal{U}} \sum_{i \in [m]} |q_{i;u;k_i} - q'_{i;u;k_i}| \\ &= \sum_{u \in \mathcal{U}} \sum_{i \in [m]} |q_{i;u;k_i} - q'_{i;u;k_i}| \\ &= \sum_{u \in \mathcal{U}} \sum_{i \in [m]} |q_{i;u;k_i} - q'_{i;u;k_i}| \end{aligned}$$

where last inequality holds because  $q_{i;u;k_i} \in [0, 1]$  and  $q'_{i;u;k_i} \in [0, 1]$ , for any  $j \in [m]; u \in \mathcal{U}$ .

For the second term, we have

$$\begin{aligned} & \sum_{j \in [m], u \in \mathcal{U}} |q_{j;u;k_j} - q'_{j;u;k_j}| \\ &= \sum_{j \in [m], u \in \mathcal{U}} |q_{j;u;k_j} - q'_{j;u;k_j}| \\ &= \sum_{j \in [m], u \in \mathcal{U}} |q_{j;u;k_j} - q'_{j;u;k_j}| \end{aligned}$$

where the inequality holds because of the Weierstrass product inequality. Plugging the above two terms into the previous inequality, 1-Norm Bounded Smoothness is satisfied.

### F.2. Regret Analysis for Overlapping MuLaNE

In this section, we give the regret analysis for overlapping MuLaNE.

#### F.2.1. FACTS

We utilize the following tail bound in our analysis.

Fact 1 (Hoeffding's Inequality (Dubhashi & Panconesi, 2009).) Let  $Y_1, \dots, Y_n$  be independent and identically distributed (i.e., i.i.d) random variables with common support  $[0, 1]$  and mean  $\mu$ . Let  $Z = Y_1 + \dots + Y_n$ . Then for all  $\epsilon$ ,

$$\Pr\{Z - n\mu \geq \epsilon\} \leq 2e^{-\frac{\epsilon^2}{2n}}$$

#### F.2.2. PROOF DETAILS

Let  $\mathcal{A}$  denote the set containing all base arms,  $\mathcal{A} = \{f(i; u; b) | i \in [m]; u \in [U]; b \in [B]\}$ . We add subscript  $t$  to denote the value of a variable at the end of round  $t$ , e.g.  $q_{i;u;b;t}$ , where  $T$  is the total number of rounds. For example,  $T_{i;u;b;t}$  denotes the total number of times that arm

$(i; u; b)$  is played at the end of round  $t$ . Let us first introduce a definition of an unlikely event, where  $\hat{\Delta}_{i;u;b;t-1}$  is not accurate as expected.

**Definition 2.** We say that the sampling is nice at the beginning of round  $t$ , if for every arm  $(i; u; b) \in \mathcal{A}$ ,  $j \wedge_{i;u;b;t-1}$

$i;u;b;j < i;u;b;t$ , where  $i;u;b;t = \frac{3 \ln t}{2T_{i;u;b;t-1}}$ . Let  $N_t$  be such event.

**Lemma 11.** For each round  $t \geq 1$ ,  $\Pr[N_t] \leq \frac{2}{3} \sum_{(i;u;b) \in \mathcal{A}} \frac{1}{T_{i;u;b;t-1}}$

**Proof.** For each round  $t \geq 1$ , we have

$$\begin{aligned} \Pr[N_t] &= \Pr\left[\bigvee_{(i;u;b) \in \mathcal{A}} \bigwedge_{j \in [2]^{[b]}} \hat{\Delta}_{i;u;b;t-1} > \frac{3 \ln t}{2T_{i;u;b;t-1}}\right] \\ &\leq \sum_{(i;u;b) \in \mathcal{A}} \Pr\left[\bigwedge_{j \in [2]^{[b]}} \hat{\Delta}_{i;u;b;t-1} > \frac{3 \ln t}{2T_{i;u;b;t-1}}\right] \\ &= \sum_{(i;u;b) \in \mathcal{A}} \Pr\left[\bigwedge_{j \in [2]^{[b]}} \left| \frac{1}{k} \sum_{s=1}^k Y_{i;u;b}^{[j]} - \frac{1}{k} \sum_{s=1}^k Y_{i;u;b}^{[j]} \right| > \frac{3 \ln t}{2T_{i;u;b;t-1}}\right] \\ &\leq \sum_{(i;u;b) \in \mathcal{A}} \sum_{j \in [2]^{[b]}} \frac{2}{t^3} \sum_{s=1}^k \Pr\left[\sum_{i=1}^s Y_{i;u;b}^{[j]} > \frac{3 \ln t}{2T_{i;u;b;t-1}}\right] \\ &\leq \sum_{(i;u;b) \in \mathcal{A}} \sum_{j \in [2]^{[b]}} \frac{2}{t^3} \sum_{s=1}^k \Pr\left[\sum_{i=1}^s Y_{i;u;b}^{[j]} > \frac{3 \ln t}{2T_{i;u;b;t-1}}\right] \\ &\leq \sum_{(i;u;b) \in \mathcal{A}} \sum_{j \in [2]^{[b]}} \frac{2}{t^3} \sum_{s=1}^k \Pr\left[\sum_{i=1}^s Y_{i;u;b}^{[j]} > \frac{3 \ln t}{2T_{i;u;b;t-1}}\right] \end{aligned}$$

Given  $T_{i;u;b;t-1} = k$ ,  $\hat{\Delta}_{i;u;b;t-1}$  is the average of  $k$  i.i.d. random variables  $Y_{i;u;b}^{[1]}, \dots, Y_{i;u;b}^{[k]}$ , where  $Y_{i;u;b}^{[j]}$  is the Bernoulli random variable (computed in Alg. 3 line 12) when the arm  $(i; u; b)$  is played for the  $j$ -th time during the execution. That is,  $\hat{\Delta}_{i;u;b;t-1} = \frac{1}{k} \sum_{j=1}^k Y_{i;u;b}^{[j]}$ . Note that the independence of  $Y_{i;u;b}^{[1]}, \dots, Y_{i;u;b}^{[k]}$  comes from the fact we select a new starting position following the starting distribution  $\eta$  at the beginning of each round. And the last inequality uses the Hoeffding Inequality (Fact 1).

Then we can use monotonicity and 1-norm bounded smoothness properties (Property 1 and Property 2) to bound the reward gap  $k_t = r_{i;u;b}(k_t) - r_{i;u;b}(k)$  between the optimal action  $k$  and the action  $k_t = (k_{t,1}, \dots, k_{t,m})$  selected by our algorithm. To achieve this, we introduce a positive real number  $M_{i;u;b}$  for each arm  $(i; u; b)$  and define  $M_{k_t} = \max_{(i;u;b) \in \mathcal{A}} M_{i;u;b}$ , where  $S_t = f(i; u; b) \sum_{j \in [2]^{[b]}} \frac{1}{T_{i;u;b;t-1}}$ . Define  $E_t = \sum_{j \in [2]^{[b]}} \frac{1}{T_{i;u;b;t-1}} |V_j|$ , and let

$$\tau(M; s) = \begin{cases} 8 & \text{if } s = 0; \\ 3 \frac{6 \ln T}{s} & \text{if } 1 \leq s \leq \tau(M); \\ 0 & \text{if } s > \tau(M) + 1; \end{cases}$$

$$\tau(M) = \frac{54 E^2 \ln T}{M^2} c$$

**Lemma 12.** If  $f_{k_t} M_{k_t} g$  and  $N_t$  holds, we have

$$\sum_{(i;u;b) \in \mathcal{A}} \sum_{j \in [2]^{[b]}} \tau(M_{i;u;b}; T_{i;u;b;t-1}) \leq \frac{1}{2} S_t \quad (30)$$

**Proof.** First, we can observe,  $f(i; u; b) \sum_{j \in [2]^{[b]}} \frac{1}{T_{i;u;b;t-1}}$  is upper bounded by  $\frac{1}{T_{i;u;b;t-1}}$ , i.e.,  $\frac{1}{T_{i;u;b;t-1}} \geq \frac{1}{T_{i;u;b;t}}$ , when  $N_t$  holds. This is due to

$$\begin{aligned} \frac{1}{T_{i;u;b;t-1}} &= \max_{j \in [2]^{[b]}} \frac{1}{T_{i;u;b;t}} \\ &\geq \min_{j \in [2]^{[b]}} \frac{1}{T_{i;u;b;t}} + \frac{1}{T_{i;u;b;t}} \\ &= \frac{1}{T_{i;u;b;t}} + \frac{1}{T_{i;u;b;t}} \end{aligned}$$

where the last inequality comes from the fact that  $j \wedge_{i;u;b;t-1} i;u;b;j < i;u;b;t$  when  $N_t$  holds.

Then, notice that the right hand of Inequality (30) is non-negative, it is trivially satisfied when  $k_t = 0$ . We only need to consider  $k_t > 0$ . By  $N_t$  and Property 1, We have the following inequalities,

$$r_{i;u;b}(k_t) - r_{i;u;b}(k) = \sum_{j \in [2]^{[b]}} \frac{1}{T_{i;u;b;t-1}} (k_{t,j} - k_j)$$

The first inequality is due to the oracle outputs  $\theta_{i;u;b;t}$  (approximate solution given parameters  $\theta_{i;u;b;t}$  where  $\theta_{i;u;b;t} = 1$  by definition). The second inequality is due to Condition 1 (monotonicity).

Then, by Condition 2, we have

$$\begin{aligned} \sum_{(i;u;b) \in \mathcal{A}} \sum_{j \in [2]^{[b]}} \frac{1}{T_{i;u;b;t-1}} (k_{t,j} - k_j) &\leq \sum_{(i;u;b) \in \mathcal{A}} \sum_{j \in [2]^{[b]}} \frac{1}{T_{i;u;b;t-1}} (k_{t,j} - k_j) \\ &\leq \sum_{(i;u;b) \in \mathcal{A}} \sum_{j \in [2]^{[b]}} \frac{1}{T_{i;u;b;t-1}} (k_{t,j} - k_j) \end{aligned}$$

where the last inequality is due to the following observation:  $u \in \mathcal{U}$  if and only if  $u$  has not been visited so far, so  $\hat{\Delta}_{i;u;b;t-1} = 0$  for any  $i \in [m]; b \in [B]$  and  $i;u;b;t = \max_{j \in [2]^{[b]}} \frac{1}{T_{i;u;b;t-1}} + \frac{1}{T_{i;u;b;t-1}} g = \frac{1}{T_{i;u;b;t-1}}$ .

Next, we can bound  $k_t$  by bounding  $\frac{1}{T_{i;u;b;t-1}}$  by applying a transformation. As  $k_t = M_{k_t} g$  holds by assumption, so  $\sum_{(i;u;b) \in \mathcal{A}} \sum_{j \in [2]^{[b]}} \frac{1}{T_{i;u;b;t-1}} \leq \frac{1}{2} S_t$

$$\begin{aligned}
 & \sum_{(i;u;b) \in S_t} \sum_{j \in [b]} M_{k_t} \cdot \text{We have} \\
 & \sum_{(i;u;b) \in S_t} \sum_{j \in [b]} \left( \sum_{i;u;b} j + \sum_{i;u;b} j + \sum_{i;u;b} j \right) \\
 & \sum_{(i;u;b) \in S_t} \sum_{j \in [b]} \left( \sum_{i;u;b} j + \sum_{i;u;b} j + \sum_{i;u;b} j \right) \\
 & = 2 \sum_{(i;u;b) \in S_t} \sum_{j \in [b]} \left( \sum_{i;u;b} j + \sum_{i;u;b} j + \sum_{i;u;b} j \right) \frac{M_{k_t}}{2jEg} \\
 & = 2 \sum_{(i;u;b) \in S_t} \sum_{j \in [b]} \left( \sum_{i;u;b} j + \sum_{i;u;b} j + \sum_{i;u;b} j \right) \frac{M_{i;u;b}}{2jEg} : \\
 & \tag{31}
 \end{aligned}$$

By  $N_t$ , for all  $(i; u; b) \in A$ , we have:

$$\begin{aligned}
 & \sum_{j \in [b]} \sum_{i;u;b} j + \sum_{i;u;b} j + \sum_{i;u;b} j \frac{M_{i;u;b}}{2jEg} \\
 & = \sum_{i;u;b} \sum_{j \in [b]} j + \sum_{i;u;b} \sum_{j \in [b]} j + \sum_{i;u;b} \sum_{j \in [b]} j \frac{M_{i;u;b}}{2jEg} \\
 & \quad (l = \arg \max_{j \in [b]} \sum_{i;u;b} j) \\
 & \sum_{i;u;b} \sum_{j \in [b]} j + \sum_{i;u;b} \sum_{j \in [b]} j \frac{M_{i;u;b}}{2jEg} \\
 & \min_{i;u;b} \sum_{j \in [b]} j + \sum_{i;u;b} \sum_{j \in [b]} j \frac{M_{i;u;b}}{2jEg} \\
 & \min_{i;u;b} \sum_{j \in [b]} j + \sum_{i;u;b} \sum_{j \in [b]} j \frac{M_{i;u;b}}{2jEg} \\
 & \quad \left( s = \frac{3 \ln T}{2T_{i;u;b;t} - 1} \right) \\
 & \min_{i;u;b} \sum_{j \in [b]} j + \sum_{i;u;b} \sum_{j \in [b]} j \frac{M_{i;u;b}}{2jEg} :
 \end{aligned}$$

where the first inequality is due to  $\sum_{i;u;b} j$  for  $j \in [b]$ , the second inequality is due to the event  $N_t$  and the third inequality holds because  $\sum_{i;u;b} j \geq T_{i;u;b;t} - 1$  for  $j \in [b]$ .

Now consider the following cases:

Case 1. If  $T_{i;u;b;t} - 1 \geq T(M_{i;u;b})$ , we have  $\sum_{i;u;b} j \geq \sum_{i;u;b} j + \sum_{i;u;b} j \frac{M_{i;u;b}}{2jEg} \geq \min_{i;u;b} \sum_{j \in [b]} j + \sum_{i;u;b} \sum_{j \in [b]} j \frac{M_{i;u;b}}{2jEg} \geq 1$ .

Case 2. If  $T_{i;u;b;t} - 1 < T(M_{i;u;b;t}) + 1$ , then  $\sum_{i;u;b} j \geq \sum_{i;u;b} j + \sum_{i;u;b} j \frac{M_{i;u;b}}{2jEg}$ , so  $\sum_{i;u;b} j \geq \sum_{i;u;b} j + \sum_{i;u;b} j \frac{M_{i;u;b}}{2jEg} = 0 = \frac{1}{2} T(M_{i;u;b;t} - 1)$ .

Combine these two cases and inequality (31), we have the following inequality, which concludes the lemma.

$$\sum_{(i;u;b) \in S_t} \sum_{j \in [b]} T(M_{i;u;b}; T_{i;u;b;t} - 1) :$$

dependent (1 - e^{-ε}) approximation regret,

$$\text{Reg} ; (T) \leq \sum_{(i;u;b) \in A} \frac{108mjVj \ln T}{\min_{i;u;b}} + 2jAj + \frac{2}{3}jAj \max :$$

Proof of Theorem 3 By above lemmas, we have

$$\begin{aligned}
 & \sum_{t=1}^T \mathbb{E} \left[ \sum_{(i;u;b) \in S_t} \sum_{j \in [b]} T(M_{i;u;b}; T_{i;u;b;t} - 1) \right] \\
 & = \sum_{(i;u;b) \in A} \sum_{t=1}^T \sum_{j \in [b]} T(M_{i;u;b}; T_{i;u;b;t} - 1) \\
 & \quad \sum_{(i;u;b) \in A} \sum_{s=0}^{T-1} \sum_{j \in [b]} T(M_{i;u;b}; s) \\
 & \quad \sum_{(i;u;b) \in A} \sum_{s=0}^{T-1} \sum_{j \in [b]} \int_0^s \frac{1}{3} \frac{\ln T}{s} ds \\
 & \quad \sum_{(i;u;b) \in A} \sum_{s=0}^{T-1} \sum_{j \in [b]} \int_0^s \frac{1}{3} \frac{\ln T}{s} ds \\
 & \quad \sum_{(i;u;b) \in A} \sum_{s=0}^{T-1} \sum_{j \in [b]} \frac{108Eg \ln T}{M_{i;u;b}} :
 \end{aligned}$$

For distribution dependent bound, let us set  $\sum_{i;u;b} j = \sum_{i;u;b} j$  and thus the event  $N_t$  always holds, i.e.,  $\sum_{i;u;b} j < M_{k_t} g = 0$ . We have

$$\sum_{t=1}^T \mathbb{E} \left[ \sum_{(i;u;b) \in S_t} \sum_{j \in [b]} \frac{108Eg \ln T}{\min_{i;u;b}} \right] :$$

By Lem. 11, Prf:  $N_t g \geq 2jAj t^2$ , then

$$\mathbb{E} \left[ \sum_{t=1}^T \sum_{(i;u;b) \in S_t} \sum_{j \in [b]} 2jAj t^2 \max_{i;u;b} \frac{2}{3}jAj \max : \right]$$

By our choice of  $M_{i;u;b}$ ,

$$\sum_{t=1}^T \mathbb{E} \left[ \sum_{(i;u;b) \in S_t} \sum_{j \in [b]} T(M_{i;u;b}; T_{i;u;b;t} - 1) \right] = 0 :$$

Theorem 3. Algorithm 3 has the following distribution- By the definition of the regret,

$$\begin{aligned} \text{Reg}(T) &= \sum_{t=1}^T \sum_{(i;u;b) \in \mathcal{A}} \left( \frac{108E^0 j A_j \ln T}{M_{i;u;b}} + 2jA_j + \frac{2}{3}jA_j \right)_{\max} \\ &= \sum_{(i;u;b) \in \mathcal{A}} \left( \frac{108mjV_j \ln T}{M_{i;u;b}} + 2jA_j + \frac{2}{3}jA_j \right)_{\max} \end{aligned}$$

For distribution independent bound, we take  $M_{i;u;b} = M$ , then we have  $\text{Reg}(T) \leq TM$ .

Then

$$\begin{aligned} \text{Reg}(T) &\leq \sum_{(i;u;b) \in \mathcal{A}} \left( \frac{108E^0 j A_j \ln T}{M_{i;u;b}} + 2jA_j + \frac{2}{3}jA_j \right)_{\max} + TM \\ &= 2^p \frac{108E^0 j A_j \ln T}{M} + 2jA_j + \frac{2}{3}jA_j \max + TM \\ &= 2^p \frac{108mjV_j \ln T}{M} + 2jA_j + \frac{2}{3}jA_j \max + 2jA_j \end{aligned}$$

### F.3. Extension for Random Node Weights

In this section, we extend the deterministic node weight for node  $u \in V$  to a random weight. Specifically, we denote  $\tilde{w}_u \in [0, 1]$  as the random weight when we first visit (after the first visit we will not get any reward so on so forth), and denote  $w_u := E[\tilde{w}_u]$ .

For the offline setting, the reward function in Eq.(7) remains the same since  $w_u$  is independent of all other random variables. For the online setting, since the reward function remains unchanged, property 1 and property 2 still hold. However, we need to change the way we maintain the optimistic node weight  $w_u$ .

We present our CUCB-MAX-R algorithm in Alg. 8. Compared with Alg. 3, the major difference is that now represent the upper confidence bound (UCB) value of  $\ln$  in line 2, we denote  $T_v^0$  as the number of times nodes played so far and  $\hat{w}_v$  as the empirical mean of  $w_v$ . In line 9,  $\hat{w}_v^0$  is the confidence radius and  $u_i$  is the UCB value of node

### Algorithm 8 CUCB-MAX-R: Combinatorial Upper Confidence Bound algorithm for the overlapping MuLaNE, with Random node weights

Input: Budget  $B$ , number of layers  $m$ , number of nodes  $|V|$ , constraints  $c$ , of the oracle BEG.

- 1: For each arm  $(i; u; b) \in \mathcal{A}$ ,  $T_{i;u;b}^0 = 0, \hat{w}_{i;u;b} = 0$ .
- 2: For each node  $v \in V$ ,  $T_v^0 = 0, \hat{w}_v = 0$ .
- 3: for  $t = 1; 2; 3; \dots; T$  do
- 4: for  $(i; u; b) \in \mathcal{A}$  do
- 5:  $u_{i;u;b} = \frac{3 \ln t = (2T_{i;u;b}^0)}$ .
- 6:  $\tilde{w}_{i;u;b} = \min \hat{w}_{i;u;b} + u_{i;u;b}; 1g$ .
- 7: end for
- 8: for  $v \in V$  do
- 9:  $\hat{w}_v^0 = \frac{3 \ln t = (2T_v^0)}$ .
- 10:  $\hat{w}_v = \min \hat{w}_v + \hat{w}_v^0; 1g$ .
- 11: end for
- 12: For  $(i; u; b) \in \mathcal{A}$ ,  $u_{i;u;b} = \max_{j \in [b]} \tilde{w}_{i;u;j}$ .
- 13:  $k = \text{BEG}((u_{i;u;b})_{(i;u;b) \in \mathcal{A}}, (\hat{w}_v)_{v \in V}, B, c)$ .
- 14: Apply budget allocation  $k$ , which gives trajectories  $X := (X_{i;1}; \dots; X_{i;k_i})_{i \in [m]}$  as feedbacks.
- 15: For any visited node  $v \in V$ ,  $X_{i;1}; \dots; X_{i;k_i}g$ , receive its random node weight  $w_v$ , update  $T_v^0 = T_v^0 + 1, \hat{w}_v = \hat{w}_v + (\tilde{w}_v - \hat{w}_v) = T_v^0$ .
- 16: For any  $(i; u; b) \in \mathcal{A}$ ,  $u_{i;u;b} = f(i; u; b) \in \mathcal{A}_j \in \mathcal{A}_k, g, Y_{i;u;b} = 1$  if  $u \in f(X_{i;1}; \dots; X_{i;k_i}g)$  and 0 otherwise.
- 17: For  $(i; u; b) \in \mathcal{A}$ , update  $T_{i;u;b} = T_{i;u;b} + 1, \hat{w}_{i;u;b} = \hat{w}_{i;u;b} + (Y_{i;u;b} - \hat{w}_{i;u;b}) = T_{i;u;b}$ .
- 18: end for

weight  $w_u$ . In line 15, we update the empirical mean for any visited  $v$ .

Regret analysis Let  $\mathcal{A}$  denote the set containing all base arms for visiting probabilities, i.e.  $\mathcal{A} = \{(i; u; b) \in [m]; u \in V; b \in [c]g$ . With a bit of abuse of notation, we also use  $\mathcal{A}$  to denote the set of base arms corresponding to the node weights. Therefore,  $\mathcal{A}[V]$  is the set of all base arms, which is different from the Section F.2 that only has  $\mathcal{A}$  as base arms and does not have base arms for random weights. Following the notations of (Wang & Chen, 2017), we denote  $X_a \in [0, 1]$  be the random outcome of base arm  $a \in \mathcal{A}[V]$  and denote  $\mathbb{D}$  as the unknown joint distribution of random outcomes  $\mathbb{D} = (X_a)_{a \in \mathcal{A}[V]}$  with unknown mean  $\mu_a = E_{\mathbb{D}}[X_a]$ . Let  $k$  be any feasible budget allocation, we denote  $p_a^{\mathbb{D};k}$  the probability that action  $k$  triggers arm  $a$  (i.e. outcome  $X_a$  is observed) when the unknown environment distribution is  $\mathbb{D}$ . Specifically,  $p_a^{\mathbb{D};k} = 1$  for  $a \in f(i; u; b) \in [m]; u \in V; b = k; g, p_a^{\mathbb{D};k} = 1_{i \in [m]}(1_{i;a;k_i})$  for  $a \in V$  and  $p_a^{\mathbb{D};k} = 0$  for the rest of base arms. Therefore, we can rewrite the

inequality in property 2 as

$$\begin{aligned}
 & j r ; (k) X r o ; o(k) j \\
 & ( u j i ; u ; b \quad 0 \quad i ; u ; b j ) \\
 & i 2 [ m ] ; u 2 V ; b = k_i \quad 0 \quad 1 \\
 & X \quad j u \quad 0 j @ 1 \quad Y \quad ( 1 \quad i ; u ; b ) A \\
 & u 2 V \quad X \quad i 2 [ m ] \\
 & X \quad p_a^{D ; k} j a \quad 0 j ; \quad ( 32 ) \\
 & a 2 A
 \end{aligned}$$

where we abuse the notation to denote  $a$  and  $0$  as  $0$  if  $a \in \mathcal{A}_2$  in the last inequality. We can observe that the above rewritten property 2 can be viewed as a special case of the 1-norm TPM bounded smoothness condition in (Wang & Chen, 2017).

Now we can apply the similar proof in Section B.3 in (Wang & Chen, 2017), together with how we deal with the  $j_{i;u;b;t}$  in Section F.2.2 to get a distribution-dependent regret bound.

Let  $K = m|V| + |V|$  be the maximum number of triggered arms,  $J_A[V] = Bm|V| + |V|$  be the number of base arms and  $\epsilon_0 = \max_{k} \epsilon_k$ . For any feasible budget allocation  $k$  and any base arm  $a \in \mathcal{A}[V]$ , the gap  $\text{gap}_k = \max_{a \in \mathcal{A}[V]} (r_a(k) - r^*(k))$  and we define  $\frac{a}{\min} = \inf_{k: p_a^{D;k} > 0, k > 0} \text{gap}_k$ ,  $\frac{a}{\max} = \sup_{k: p_a^{D;k} > 0, k > 0} \text{gap}_k$ . As a convention, if there is no action such that  $p_a^{D;k} > 0$  and  $k > 0$ , we define  $\frac{a}{\min} = +1$ ;  $\frac{a}{\max} = 0$ . We define  $\min_{a \in \mathcal{A}_2} \frac{a}{\min}$  and  $\max_{a \in \mathcal{A}_2} \frac{a}{\max}$ .

Theorem 7. Algorithm 8 has the following distribution-dependent  $(1 - \epsilon)$  approximation regret,

$$\begin{aligned}
 \text{Reg}_\epsilon(T) & \leq \sum_{a \in \mathcal{A}[V]} \frac{576K \ln T}{\frac{a}{\min}} \\
 & + 4J_A[V] + \sum_{a \in \mathcal{A}[V]} \frac{2K}{\frac{a}{\min}} \epsilon_0 + 2 \frac{2}{6} \max;
 \end{aligned}$$

We can also have a distribution-independent bound,

$$\begin{aligned}
 \text{Reg}_\epsilon(T) & \leq 12^p \overline{J_A[V]} K T \ln T \\
 & + 2J_A[V] + \frac{T}{18 \ln T} \epsilon_0 + 2 \frac{2}{6} \max;
 \end{aligned}$$

Remark. Compared with the regret bound in Thm. 3 in the main text, we can see the base arms now become  $|V|$  instead of  $A$ , which is worse than the regret bound of the fixed unknown weight setting. This is a trade-off since we need more estimation to handle the uncertain random node weight, which incurs additional regrets. It also shows our non-trivial analysis of the fixed unknown weights carefully merge the error term into base arms and incurs only constant-factor of extra regrets.

## G. Online Learning for Non-overlapping MuLaNE

### G.1. Online Learning Algorithm for Non-overlapping MuLaNE

For non-overlapping MuLaNE, we estimate layer-level marginal gains as our parameters, rather than probabilities  $P_{i;u}(b)$  and node weights  $\theta$ . By doing so, we can largely simplify our analysis since we no longer need to handle the unknown weights and the extra constraint of  $\theta(b)$  as in the overlapping setting. Concretely, we maintain a set of base arms  $\mathcal{A} = \{f(i; b) | i \in \mathcal{I}_2[m]; b \in \mathcal{B}[c_i]\}$ , and let  $|\mathcal{A}| = \sum_{i \in \mathcal{I}_2[m]} c_i$  be the total number of these arms. For each base arm  $(i; b) \in \mathcal{A}$ , denote  $\theta_{i;b}$  as the truth value of each base arm, i.e.  $\theta_{i;b} = \sum_{u \in \mathcal{U}} \theta_u (P_{i;u}(b) - P_{i;u}(b^*))$ . Based on the definition of base arms, we can see that we do not need to explicitly estimate the node weights or probabilities  $P_{i;u}(b)$ , so we will use  $\theta(k)$  to denote the reward function  $r_{\mathcal{G}}(\cdot; (k))$ .

We present our algorithm in Alg. 9, which has three differences compared with the Alg. 3 for overlapping MuLaNE. First, the expected reward of action is  $r(k) = \sum_{i=1}^m \sum_{b=1}^{k_i} \theta_{i;b}$ , which is a linear reward function and satisfies two new Monotonicity and 1-Norm Bounded Smoothness conditions (see Condition 1 and 2). Second, we use the Dynamic Programming method (Alg. 7) as our Oracle, which does not require the parameters  $\theta$  to be increasing. Moreover, the ofine oracle always outputs optimal solutions. Third, we use a new random variable  $Y_{i;b}$ , which indicates the gain of weights which depends on whether a new node is visited by random walker from the layer exactly at the  $b$ -th step, i.e.  $Y_{i;b} = \sum_{j=1}^S \mathbb{1}_{f(X_{ij}) = u} X_{ij}$  for base arm  $(i; b)$  such that  $b \leq k_i$ .

Condition 1 (Monotonicity). The reward  $r(k)$  is monotonically increasing, i.e., for any budget allocation and any two vectors  $\theta = (\theta_{i;b})_{(i;b) \in \mathcal{A}}$ ,  $\theta^0 = (\theta^0_{i;b})_{(i;b) \in \mathcal{A}}$ , we have  $r(k) \geq r^0(k)$ , if  $\theta_{i;b} \geq \theta^0_{i;b}$ , for  $(i; b) \in \mathcal{A}$ .

Condition 2 (1-Norm Bounded Smoothness). The reward function  $r(k)$  satisfies the 1-norm bounded smoothness, i.e., for any budget allocation  $k$  and any two vectors  $\theta = (\theta_{i;b})_{(i;b) \in \mathcal{A}}$ ,  $\theta^0 = (\theta^0_{i;b})_{(i;b) \in \mathcal{A}}$ , we have  $|r(k) - r^0(k)| \leq \sum_{(i;b) \in \mathcal{A}} |\theta_{i;b} - \theta^0_{i;b}|$ .

We define the reward gap  $\text{gap}_k = r(k) - r^*(k)$  for all feasible action  $k$  satisfying  $\sum_{i=1}^m k_i = B$ ,  $0 \leq k_i \leq c_i$ . For each base arm  $(i; b)$ , we define  $\frac{i;b}{\min} = \min_{k > 0; k_i \leq c_i} \text{gap}_k$  and  $\frac{i;b}{\max} = \max_{k > 0; k_i \leq c_i} \text{gap}_k$ . As a convention, if there is no action with  $k_i \leq c_i$  such that  $k > 0$ , we define  $\frac{i;b}{\min} = 1$  and  $\frac{i;b}{\max} = 0$ . Also, we define  $\frac{\min}{\min} = \min_{(i;b) \in \mathcal{A}_2} \frac{i;b}{\min}$  and  $\frac{\max}{\max} = \max_{(i;b) \in \mathcal{A}_2} \frac{i;b}{\max}$ .

Theorem 8. CUCB-MG has the following distribution-

Algorithm 9 CUCB-MG: Combinatorial Upper Confidence Bound (CUCB) algorithm for non-overlapping MuLaNE, using Marginal Gains as the base arms

Input: Budget  $B$ , number of layers  $s$ , of the oracle DP.

- 1: For each  $(i; b) \in \mathcal{A}$ ,  $T_{i;b} = 0$ ,  $\hat{\mu}_{i;b} = 1$ .
- 2: for  $t = 1; 2; 3; \dots; T$  do
- 3: for  $(i; b) \in \mathcal{A}$  do
- 4:  $\hat{\mu}_{i;b} = \frac{3 \ln t}{2T_{i;b}}$ .
- 5:  $\hat{\mu}_{i;b} = \min\{\hat{\mu}_{i;b} + \frac{1}{t}, 1\}$ .
- 6: end for
- 7:  $k_t = \text{DP}(\hat{\mu}_{i;b}, B, c)$
- 8: Apply budget allocation  $k_t$ , which gives trajectories  $X_t := (X_{i;1}, \dots, X_{i;k_t})_{i \in [m]}$  as feedbacks.
- 9: For  $(i; b) \in \mathcal{A}$ ,  $\mu_{i;b} = f(i; b) \mathbb{1}_{k_t \geq b}$ ,  $Y_{i;b} = \sum_{j=1}^{S_{i;b}} f(X_{i;j}) \mathbb{1}_{k_t \geq b}$ .
- 10: For  $(i; b) \in \mathcal{A}$ , update  $T_{i;b}$  and  $\hat{\mu}_{i;b}$ :
- 11:  $T_{i;b} = T_{i;b} + 1$ ,  $\hat{\mu}_{i;b} = \hat{\mu}_{i;b} + (Y_{i;b} - \hat{\mu}_{i;b}) \mathbb{1}_{k_t \geq b}$ .
- 12: end for

Proof. For each round  $t$ , we have

$$\begin{aligned} \text{Prf: } N_t &= \sum_{(i;b) \in \mathcal{A}} \sum_{j=1}^{S_{i;b}} \frac{3 \ln t}{2T_{i;b}} \mathbb{1}_{k_t \geq b} \\ &= \sum_{(i;b) \in \mathcal{A}} \sum_{j=1}^{S_{i;b}} \frac{3 \ln t}{2T_{i;b}} \mathbb{1}_{k_t \geq b} \\ &= \sum_{(i;b) \in \mathcal{A}} \sum_{k=1}^{k_t} \frac{3 \ln t}{2T_{i;b}} \mathbb{1}_{k_t \geq k} \\ &= \sum_{(i;b) \in \mathcal{A}} \sum_{k=1}^{k_t} \frac{2}{t^3} \mathbb{1}_{k_t \geq k} \end{aligned}$$

When  $T_{i;b} = k$ ,  $\hat{\mu}_{i;b}$  is the average of  $k$  i.i.d. random variables  $Y_{i;b}^{[1]}, \dots, Y_{i;b}^{[k]}$ , where  $Y_{i;b}^{[j]}$  is the Bernoulli random variable of arm  $(i; b)$  when it is played for the  $j$ -th time during the execution. That is,  $\hat{\mu}_{i;b} = \frac{1}{k} \sum_{j=1}^k Y_{i;b}^{[j]}$ .

Note that the independence  $Y_{i;b}^{[1]}, \dots, Y_{i;b}^{[k]}$  comes from the fact we select a new starting position following the starting distribution  $\mu_i$  at the beginning of each round. The last inequality uses the Hoeffding Inequality (Fact 1).

dependent regret bound,

$$\text{Reg}(T) \leq \sum_{(i;b) \in \mathcal{A}} \frac{48B \ln T}{\min_{i;b} c} + 2jA_j + \frac{2}{3}jA_j \max_{i;b} c$$

## G.2. Regret Analysis for Non-overlapping MuLaNE

In this section, we give the regret analysis for CUCB-MG, which applies the standard CUCB algorithm (Wang & Chen, 2017) to this setting.

### G.2.1. PROOF DETAILS

Let  $\mathcal{A}$  denote the set containing all base arms,  $\mathcal{A} = \{(i; b) \mid i \in [m]; b \in [c]\}$ . We add subscript to denote the value of a variable at the end of round  $t$ . For example,  $T_{i;b,t}$  denotes the total times of arm  $(i; b) \in \mathcal{A}$  is played at the end of round  $t$ . Let us first introduce a definition of an unlikely event that  $\hat{\mu}_{i;b,t}$  is not accurate as expected.

Definition 3. We say that the sampling is nice at the beginning of round  $t$ , if for every arm  $(i; b) \in \mathcal{A}$ ,  $j \hat{\mu}_{i;b,t} < \mu_{i;b} < j \hat{\mu}_{i;b,t} + \frac{1}{t}$ , where  $\mu_{i;b} = \frac{3 \ln t}{2T_{i;b}}$ . Let  $N_t$  be such event.

Lemma 13. For each round  $t$ ,  $\text{Prf: } N_t \leq 2jA_j t^{-2}$

Then we can use monotonicity and 1-norm bounded smoothness property (Condition 1 and Condition 2) to bound the reward gap  $k_t = r(k_t) - r(k_t)$  between the optimal action  $k_t$  and the action  $k_t = (k_{t;1}, \dots, k_{t;m})$  selected by our algorithm  $A$ . To achieve this, we introduce a positive real number  $M_{i;b}$  for each arm  $(i; b)$  and define  $M_{k_t} = \max_{(i;b) \in \mathcal{S}_t} M_{i;b}$ , where  $\mathcal{S}_t = \{(i; b) \in \mathcal{A} \mid k_{t,i} \geq b\}$  and  $jS_j = B$ . Define,

$$\tau(M; s) = \begin{cases} 8 & \text{if } s = 0; \\ \frac{2}{s} & \text{if } 1 \leq s \leq \tau(M); \\ 0 & \text{if } s > \tau(M) + 1; \end{cases}$$

where

$$\tau(M) = \frac{24B^2 \ln T}{M^2} c$$

Lemma 14. If  $k_t = M_{k_t} g$  and  $N_t$  holds, we have

$$\sum_{(i;b) \in \mathcal{S}_t} \tau(M_{i;b}; T_{i;b,t}) \leq 2jA_j t^{-2} \quad (33)$$

Proof. First, we can observe  $\mu_{i;b} < \hat{\mu}_{i;b}$ , for  $(i; b) \in \mathcal{A}$ , when  $N_t$  holds. This comes from the fact  $\hat{\mu}_{i;b} = \min\{\hat{\mu}_{i;b} + \frac{1}{t}, 1\}$  and  $j \hat{\mu}_{i;b} < \mu_{i;b} < j \hat{\mu}_{i;b} + \frac{1}{t}$ .



Notice that the right hand of Inequality (33) is non-negative, it is trivially satisfied when  $\Delta_{k_t} = 0$ . We only need to consider  $\Delta_{k_t} > 0$ . By  $\mathcal{N}_t$  and Condition 1, We have the following inequalities,

$$r_{-t}(\mathbf{k}_t) \geq r_{-t}(\mathbf{k}) \geq r(\mathbf{k}) = \Delta_{k_t} + r(\mathbf{k}_t).$$

The first inequality is due to the oracle outputs the optimal solution given parameters  $\bar{\mu}_t$ , the second inequality is due to Condition 1 (monotonicity).

Then, by Condition 2, we have

$$\Delta_{k_t} \leq r_{-t}(\mathbf{k}_t) - r(\mathbf{k}_t) \leq \sum_{(i,b) \in \mathcal{S}_t} |\bar{\mu}_{i,b,t} - \mu_{i,b}|.$$

Next, we can bound  $\Delta_{k_t}$  by bounding  $\bar{\mu}_{i,b,t} - \mu_{i,b}$  by applying a transformation. As  $\{\Delta_{k_t} \geq M_{k_t}\}$  holds by assumption, so  $\sum_{(i,b) \in \mathcal{S}_t} |\bar{\mu}_{i,b,t} - \mu_{i,b}| \geq \Delta_{k_t} \geq M_{k_t}$ . We have

$$\begin{aligned} \Delta_{k_t} &\leq \sum_{(i,b) \in \mathcal{S}_t} |\bar{\mu}_{i,b,t} - \mu_{i,b}| \\ &\leq -M_{k_t} + 2 \sum_{(i,b) \in \mathcal{S}_t} |\bar{\mu}_{i,b,t} - \mu_{i,b}| \\ &= 2 \sum_{(i,b) \in \mathcal{S}_t} (|\bar{\mu}_{i,b,t} - \mu_{i,b}| - \frac{M_{k_t}}{2B}) \\ &\leq 2 \sum_{(i,b) \in \mathcal{S}_t} (|\bar{\mu}_{i,b,t} - \mu_{i,b}| - \frac{M_{i,b}}{2B}). \end{aligned} \quad (34)$$

By  $\mathcal{N}_t$ , we have:

$$\begin{aligned} |\bar{\mu}_{i,b,t} - \mu_{i,b}| - \frac{M_{i,b}}{2B} &\leq \min\{2\rho_{i,b,t}, 1\} - \frac{M_{i,b}}{2B} \\ &\leq \min\{2 \frac{3 \ln T}{2T_{i,b,t} - 1}, 1\} - \frac{M_{i,b}}{2B}. \end{aligned}$$

Now consider the following cases:

**Case 1.** If  $T_{i,b,t} - 1 \leq \ell_T(M_{i,b})$ , we have  $\bar{\mu}_{i,b,t} - \mu_{i,b} - \frac{M_{i,b}}{2B} \leq \min\{2 \frac{3 \ln T}{2T_{i,b,t} - 1}, 1\} - \frac{M_{i,b}}{2B} \leq \frac{1}{2} \kappa_T(M_{i,b}, T_{i,b,t} - 1)$ .

**Case 2.** If  $T_{i,b,t} - 1 \geq \ell_T(M_{i,b}) + 1$ , then  $2 \frac{3 \ln T}{2T_{i,b,t} - 1} \leq \frac{M_{i,b}}{2B}$ , so  $\bar{\mu}_{i,b,t} - \mu_{i,b} - \frac{M_{i,b}}{2B} \leq 0 = \frac{1}{2} \kappa_T(M_{i,b}, T_{i,b,t} - 1)$ .

Combine these two cases and inequality (34), we have the following inequality, which concludes the lemma.

$$\Delta_{k_t} \leq \sum_{(i,b) \in \mathcal{S}_t} \kappa_T(M_{i,b}, T_{i,b,t} - 1).$$

**Theorem 8.** *CUCB-MG has the following distribution-*

*dependent regret bound,*

$$\text{Reg}_\mu(T) \leq \sum_{(i,b) \in \mathcal{A}} \frac{48B \ln T}{\Delta_{\min}^{i,b}} + 2|\mathcal{A}| + \frac{\pi^2}{3} |\mathcal{A}| \Delta_{\max}.$$

*Proof of Theorem 8.* By above lemmas, we have

$$\begin{aligned} &\sum_{t=1}^T \mathbb{1}(\{\Delta_{k_t} \geq M_{k_t}\} \wedge \mathcal{N}_t) \Delta_{k_t} \\ &\leq \sum_{t=1}^T \sum_{(i,b) \in \mathcal{S}_t} \kappa_T(M_{i,b}, T_{i,b,t} - 1) \\ &= \sum_{(i,b) \in \mathcal{A}} \sum_{t=1}^T \sum_{(i,b) \in \mathcal{S}_t} \kappa_T(M_{i,b}, T_{i,b,t} - 1) \\ &\leq \sum_{(i,b) \in \mathcal{A}} \sum_{s=0}^{\ell_T(M_{i,b})} \kappa_T(M_{i,b}, s) \\ &\leq 2|\mathcal{A}| + \sum_{(i,b) \in \mathcal{A}} \sum_{s=1}^{\ell_T(M_{i,b})} \frac{\ell_T(M_{i,b})}{2} \frac{\ln T}{6 \frac{\ln T}{s}} \\ &\leq 2|\mathcal{A}| + \sum_{(i,b) \in \mathcal{A}} \int_0^{\ell_T(M_{i,b})} \frac{\ell_T(M_{i,b})}{2} \frac{\ln T}{6 \frac{\ln T}{s}} ds \\ &\leq 2|\mathcal{A}| + \sum_{(i,b) \in \mathcal{A}} \frac{4}{6 \ln T} \ell_T(M_{i,b}) \\ &\leq 2|\mathcal{A}| + \sum_{(i,b) \in \mathcal{A}} \frac{48B \ln T}{M_{i,b}}. \end{aligned}$$

Let us set  $M_{i,b} = \Delta_{\min}^{i,b}$  and thus the event  $\{\Delta_{k_t} \geq M_{k_t}\}$  always holds. We have

$$\begin{aligned} &\sum_{t=1}^T \mathbb{1}(\mathcal{N}_t \wedge \{\Delta_{k_t} \geq M_{k_t}\}) \Delta_{k_t} \\ &\leq 2|\mathcal{A}| + \sum_{(i,b) \in \mathcal{A}} \frac{48B \ln T}{\Delta_{\min}^{i,b}}. \end{aligned}$$

By Lem. 13,  $\Pr\{\neg \mathcal{N}_t\} \leq 2|\mathcal{A}|t^{-2}$ , then we have

$$\mathbb{E} \sum_{t=1}^T \mathbb{1}(\neg \mathcal{N}_t) \Delta_{k_t} \leq \sum_{t=1}^T 2|\mathcal{A}|t^{-2} \Delta_{\max} \leq \frac{\pi^2}{3} |\mathcal{A}| \Delta_{\max}.$$

By our choice of  $M_{i,b}$ ,

$$\sum_{t=1}^T \mathbb{1}(\Delta_{k_t} < M_{k_t}) \Delta_{k_t} = 0.$$

By the definition of the regret,

$$\begin{aligned}
 \text{Reg}(T) &= \mathbb{E} \sum_{t=1}^T \Delta_{k_t} \\
 &\leq \mathbb{E} \sum_{t=1}^T \mathbb{1}_{\{\Delta_{k_t} < M_{k_t}\}} \Delta_{k_t} \\
 &\quad + \sum_{t=1}^T \mathbb{1}_{\{\Delta_{k_t} \geq M_{k_t}\}} \Delta_{k_t} \\
 &\leq \sum_{(i,b) \in \mathcal{A}} \frac{48B \ln T}{\Delta_{min}^{i,b}} + 2|\mathcal{A}| + \frac{\pi^2}{3} |\mathcal{A}| \Delta_{max}.
 \end{aligned}$$

For distribution independent bound, we take  $M_{i,b} = M = \frac{48B|\mathcal{A}| \ln T}{T}$ , then we have  $\sum_{t=1}^T \mathbb{1}_{\{\Delta_{k_t} \leq M_{k_t}\}} \leq TM$ .

Then

$$\begin{aligned}
 \text{Reg}_\mu(T) &\leq \sum_{(i,b) \in \mathcal{A}} \frac{48B \ln T}{M_{i,b}} + 2|\mathcal{A}| + \frac{\pi^2}{3} |\mathcal{A}| \Delta_{max} + TM \\
 &\leq \frac{48B|\mathcal{A}| \ln T}{M} + 2|\mathcal{A}| + \frac{\pi^2}{3} |\mathcal{A}| \Delta_{max} + TM \\
 &= 2 \frac{48B|\mathcal{A}| \ln T}{M} + \frac{\pi^2}{3} |\mathcal{A}| \Delta_{max} + 2|\mathcal{A}| \\
 &\leq 14 \frac{48B|\mathcal{A}| \ln T}{M} + \frac{\pi^2}{3} |\mathcal{A}| \Delta_{max} + 2|\mathcal{A}|.
 \end{aligned}$$

## H. Supplemental Experiments

### H.1. Experimental Results for Stationary Distributions

In this section, we show experimental results with explorers starting from stationary distributions in Fig. 3. Specifically, each explorer  $W_i$  starts from the node  $u$  in layer  $L_i$  with probability proportional to the out-degree of  $u$ . For offline overlapping case, the total weights of unique nodes visited for BEG and MG are the same, and outperforms two baselines, which accords with our theoretical analysis. Similar to the Sec. 6, BEG and MG are empirically close to the optimal solution. For offline non-overlapping case, DP and MG give us the same *optimal* solution. For the online learning algorithms, note that we still use BEG rather than MG as the offline oracle, this is due to the UCB values do not have the diminishing return properties. The same applies for the non-overlapping case, where DP (instead of MG) is used as the offline oracle. The trend of the curves and the analysis are similar to that in the Sec. 6.

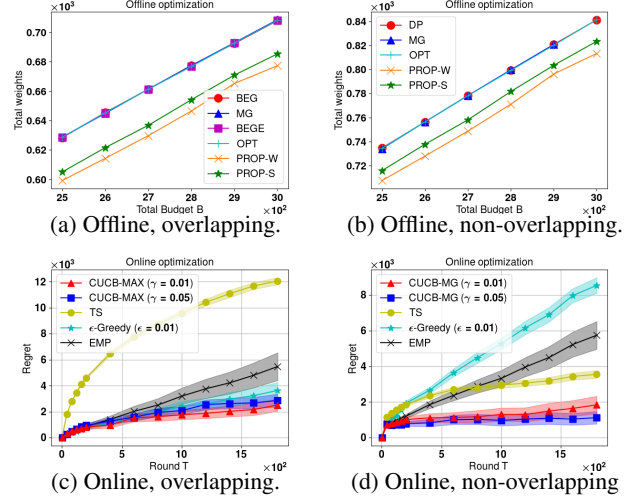


Figure 3. Explorers start with stationary distributions. Above: total weights of unique nodes visited given by different offline algorithms. Below: regret for different online algorithms when  $B = 3000$ .

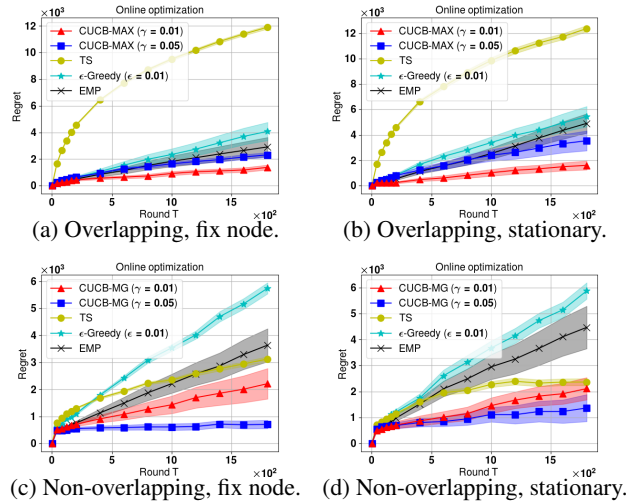


Figure 4. Regret for different online algorithms, when total budget  $B = 2000$ .

### H.2. Experimental Results for Online Learning Algorithms with Different Budgets $B$

As shown in Fig. 4 and Fig. 5, although there are some fluctuations, the trend of the curves and the analysis are almost the same as that when  $B = 3000$ , which shows our experimental results are consistent for different budgets.

(a) Overlapping, fix node.      (b) Overlapping, stationary.

(c) Non-overlapping, fix node.    (d) Non-overlapping, stationary.

*Figure 5.* Regret for different online algorithms, when total budget  $B = 4000$ .