

# Adversarial Combinatorial Bandits with Switching Costs

Yanyan Dong and Vincent Y. F. Tan, *Senior Member, IEEE*

**Abstract**—We study the problem of adversarial combinatorial bandit with a switching cost  $\lambda$  for a switch of each selected arm in each round, considering both the bandit feedback and semi-bandit feedback settings. In the oblivious adversarial case with  $K$  base arms and time horizon  $T$ , we derive lower bounds for the minimax regret and design algorithms to approach them. To prove these lower bounds, we design stochastic loss sequences for both feedback settings, building on an idea from previous work in Dekel et al. (2014). The lower bound for bandit feedback is  $\tilde{\Omega}((\lambda K)^{\frac{1}{3}}(TI)^{\frac{2}{3}})$  while that for semi-bandit feedback is  $\tilde{\Omega}((\lambda KI)^{\frac{1}{3}}T^{\frac{2}{3}})$  where  $I$  is the number of base arms in the combinatorial arm played in each round. To approach these lower bounds, we design algorithms that operate in batches by dividing the time horizon into batches to restrict the number of switches between actions. For the bandit feedback setting, where only the total loss of the combinatorial arm is observed, we introduce the BATCHED-EXP2 algorithm which achieves a regret upper bound of  $\tilde{O}((\lambda K)^{\frac{1}{3}}T^{\frac{2}{3}}I^{\frac{1}{3}})$  as  $T$  tends to infinity. In the semi-bandit feedback setting, where all losses for the combinatorial arm are observed, we propose the BATCHED-BROAD algorithm which achieves a regret upper bound of  $\tilde{O}((\lambda K)^{\frac{1}{3}}(TI)^{\frac{2}{3}})$ .

**Index Terms**—multi-armed bandits, adversarial bandits, combinatorial bandits, switching costs, online optimization

## I. INTRODUCTION

The classical multi-armed bandit (MAB) problem is a sequential decision making game between an agent and an environment [1], where the agent plays the arms sequentially to minimize the total loss over time. After each arm is played, the agent receives some feedback in the form of a loss (or gain) associated with the chosen arm. In many applications such as the financial trading or reconfiguration in industrial environments, there is a cost  $\lambda > 0$  for a switch of each selected arm in each round, which must be considered to assess the overall performance of the algorithms designed for them. For example, Guha and Munagala [2] constructed a sensor network problem and refined probabilistic models of sensed data at various nodes, which costs energy in transferring from the current node to a new node. Shi et al. [3] introduced

an application of the switching cost in edge computing with artificial intelligence, where the edge server can only utilize a small number of machine learning models in each round to learn the best subset of models based on the feedback. Downloading a model that is not in the current edge server from the cloud incurs a switching cost.

The problem of MAB with switching costs has been studied for both the oblivious adversarial case and the stochastic case. We will focus on the oblivious adversarial setting, where losses are generated by an arbitrary deterministic source before the game. In the oblivious adversarial case with  $K$  arms and time horizon  $T$ , Arora et al. [4] refined the Exp3 algorithm to achieve a regret upper bound<sup>1</sup> of  $\tilde{O}(K^{\frac{1}{3}}T^{\frac{2}{3}})$  when the switching cost  $\lambda = 1$ . Later, Dekel et al. [5] proved that the upper bound is tight by showing that the minimax regret lower bound is  $\tilde{\Omega}((\lambda K)^{\frac{1}{3}}T^{\frac{2}{3}})$ . Rouyer et al. [6] proposed an algorithm which is a modification of the Tsallis-Switch algorithm to achieve an upper bound of  $O((\lambda K)^{\frac{1}{3}}T^{\frac{2}{3}})$ . Without switching costs, the minimax regret of the adversarial MAB problem is  $\Theta(\sqrt{TK})$  [7], [8].

Combinatorial bandits form a general extension of the standard framework, which is a linear bandit with the special combinatorial action set  $\mathcal{A} \subseteq \{0, 1\}^K$ . We generalize the problem of MAB with switching costs by considering the combinatorial problem with  $I$  arms played in each round, where the set of the played arms is called a *combinatorial arm*. There are many practical applications of combinatorial problems with switching costs. For example, a hospital may plan to experiment on a drug that is known to be a combination of  $I$  components of  $K$  (raw material) components. In our parlance, there are  $K$  base arms and the combinatorial arm contains  $I$  out of the  $K$  base arms. The quality of each chosen component depends on certain unknown complex effects of the environment and patients, and this may be modeled by an adversarial setting. The overall effect of the drug is the sum of the qualities of all the individual chosen components. There is, however, a non-negligible purchasing (or switching) cost when the agent decides to swap one component for another from one time step to the next. In [3], a special combinatorial problem was considered when the costly full feedback was available and a switching cost was added in each round. Our setting does not allow full feedback and only considers bandit feedback and semi-bandit feedback. Under bandit feedback, the player can only observe the total loss of the played combinatorial arm while under semi-bandit feedback, all the losses for the

This work is supported by funding from a Ministry of Education Academic Research Fund (AcRF) Tier 2 grant under grant number A-8000423-00-00 and AcRF Tier 1 grants under grant numbers A-8000189-01-00 and A-8000980-00-00.

Yanyan Dong is with the School of Science and Engineering, The Chinese University of Hong Kong, Shenzhen, Shenzhen 518172, China (email: yanyan-dong@link.cuhk.edu.cn). This work was carried out when Yanyan Dong was a Research Fellow at the Department of Electrical and Computer Engineering, National University of Singapore.

Vincent Y. F. Tan is with the Department of Mathematics and the Department of Electrical and Computer Engineering, National University of Singapore (e-mail: vtan@nus.edu.sg), Singapore, 119077.

The work was accepted in IEEE Transactions on Information Theory.

<sup>1</sup>In this paper, we use  $\tilde{O}$  and  $\tilde{\Omega}$  to denote the big-O and big omega notations ignoring any logarithmic factors in  $T$ .

combinatorial arm are observed.

In this paper, we will focus on analyzing the regret of adversarial combinatorial bandit when there are switching costs. We derive lower bounds for the minimax regret and propose the algorithms that approximately meet the lower bounds under both feedback scenarios. To prove the lower bounds, we design stochastic loss sequences for both bandit feedback and semi-bandit feedback, which generalize the idea in [5] for the combinatorial scenarios. Under different types of feedback, the loss sequences designed are different. For a fixed time, the loss sequence under bandit feedback uses the same Gaussian noise for different base arms while the loss sequence under semi-bandit feedback uses i.i.d. Gaussian noises for different base arms. We show that the lower bound for bandit feedback is  $\tilde{\Omega}((\lambda K)^{\frac{1}{3}}(TI)^{\frac{2}{3}})$  while that for semi-bandit feedback is  $\tilde{\Omega}((\lambda KI)^{\frac{1}{3}}T^{\frac{2}{3}})$  where  $I$  is the number of base arms in the combinatorial arm played in each round. Dekel et al. [4], Rouyer et al. [6] and Shi et al. [3] all utilize the batch-based algorithms to restrict the number of switches between actions by dividing the whole time horizon into batches and forcing the algorithm to play the same action for all the rounds within a batch. We also utilize this technique in our algorithms under our combinatorial setting. In the bandit feedback setting, we introduce the BATCHED-EXP2 algorithm with John's exploration, which is a batched version of the Exp2 algorithm with John's exploration [9] and achieves a regret bound of  $\tilde{O}((\lambda K)^{\frac{1}{3}}T^{\frac{2}{3}}I^{\frac{1}{3}})$  when the time horizon  $T$  tends to infinity. In the semi-bandit feedback setting, we introduce the BATCHED-BROAD algorithm, which is a batched version of the Online Mirror Descent algorithm with log-barrier regularizer (BROAD) in [10] and achieves a regret upper bound of  $\tilde{O}((\lambda K)^{\frac{1}{3}}(TI)^{\frac{2}{3}})$ .

In the remainder of this paper, we first formulate the problem and introduce our main results in §II. Then §III presents the main ideas for proving the lower bound for two different types of feedback. In §IV, we introduce two algorithms designed for the bandit feedback and semi-bandit feedback respectively. §V is dedicated to the compare our algorithms with some baselines by numerical experiments. In Appendix, we provide complete proofs for the lower bounds.

## II. PROBLEM FORMULATION AND MAIN RESULTS

We use  $[n]$  to denote the set  $\{1, \dots, n\}$  for any positive integer  $n$  and  $x_{1:T}$  to denote the sequence  $\{x_i\}_{i=1}^T$ . We consider an adversarial combinatorial bandit problem with  $K$  base arms  $[K]$  and a switching cost  $\lambda$  for each switched arm with  $\lambda > 0$ . The loss vectors  $l_t \in [0, 1]^K$  for  $t \geq 1$  are arbitrarily generated by the adversary and do not depend on the actions taken by the learner, where the  $i$ -th component of  $l_t$  is the loss incurred if arm  $i$  is pulled at time  $t$ . Let  $\mathcal{A} = \{A \in \{0, 1\}^K : \|A\|_1 = I\}$  be the set of all combinatorial arms with  $I$  base arms where the  $i$ -th components of  $A$  is one if arm  $i$  is pulled and zero otherwise. For the sake of simplicity, we define the set of base arms within a combinatorial arm  $\mathcal{I} \in \mathcal{A}$  as  $\{i \in [K] : \mathcal{I}_i = 1\}$ , where  $\mathcal{I}_i$  is the  $i$ -th component in  $\mathcal{I}$ . In each time  $t \geq 1$ , the player pulls a combinatorial arm  $A_t \in \mathcal{A}$  and incurs a loss of  $\langle A_t, l_t \rangle = \sum_{i=1}^K A_{t,i} l_{t,i}$

and  $l_{t,i}$  are respectively the  $i$ -th component of the vectors  $A_t \in \{0, 1\}^K$  and  $l_t \in [0, 1]^K$ , where  $\langle \cdot, \cdot \rangle$  denotes the inner product operation. We consider both bandit and semi-bandit feedback [11]. Under bandit feedback and after the combinatorial arm  $A_t$  is pulled at round  $t$ , the player can only observe the feedback  $X_t = \langle A_t, l_t \rangle \in [0, I]$  while under semi-bandit feedback, all the losses for the combinatorial arm  $X_t = A_t \circ l_t \in [0, 1]^K$  are observed where  $\circ$  stands for the element-wise multiplication. From time  $t-1$  to  $t$ , there is a switching cost of  $\lambda \cdot d(A_t, A_{t-1})$  for the player where  $d(A_t, A_{t-1}) \triangleq \frac{1}{2} \|A_t \oplus A_{t-1}\|_1$ , i.e.,  $d(A_t, A_{t-1})$  measures the number of arms switched from the combinatorial arm  $A_{t-1}$  pulled at time  $t-1$  to  $A_t$  pulled at time  $t$ . We set  $A_0 = 0$  and then the first action  $A_1$  will always incur a switching cost of  $\lambda I$ . The cumulative loss over  $T$  rounds for the player equals

$$\sum_{t=1}^T \langle A_t, l_t \rangle + \lambda \sum_{t=1}^T d(A_t, A_{t-1}).$$

Given the loss sequence  $l_{1:T} \in [0, 1]^{K \times T}$  and action sequence  $A_{1:T} \in \mathcal{A}^T$ , define the  $\lambda$ -switching regret as

$$\begin{aligned} R_\lambda(A_{0:T}, l_{1:T}) \\ \triangleq \sum_{t=1}^T \langle A_t, l_t \rangle + \lambda \sum_{t=1}^T d(A_t, A_{t-1}) - \min_{A \in \mathcal{A}} \sum_{t=1}^T \langle A, l_t \rangle. \end{aligned}$$

A policy  $\pi = \pi_{1:T}$  is composed of all the conditional distributions over actions at each time  $t$  given past actions and feedback, i.e.,  $\pi_t(\cdot | A_1, X_1, \dots, A_{t-1}, X_{t-1})$ . The set of all policies in  $\mathcal{A}^T$  is denoted as  $\Pi$  and the set of all deterministic loss sequences in  $[0, 1]^{K \times T}$  is denoted as  $\mathcal{L}$ . We define the expected  $\lambda$ -switching regret when the loss functions  $l_{1:T} \in \mathcal{L}$  are specified and a policy  $\pi \in \Pi$  is employed as follows,

$$R_\lambda(\pi, l_{1:T}) = \mathbb{E} [R_\lambda(A_{0:T}, l_{1:T})],$$

where the expectation is taken over the player's randomized choice of actions under the policy  $\pi$ . In this paper, we use regret to represent for  $\lambda$ -switching regret for simplicity. When  $\lambda = 0$ , i.e., the switching cost is not considered, we write  $R(\pi, l_{1:T}) = R_0(\pi, l_{1:T})$  for simplicity and we call  $R(\pi, l_{1:T})$  the *pseudo-regret*. We use the *minimax expected  $\lambda$ -switching regret* to measure the difficulty of the game as in [5], which is defined as follows,

$$R_\lambda^* = \inf_{\pi \in \Pi} \sup_{l_{1:T} \in \mathcal{L}} R_\lambda(\pi, l_{1:T}).$$

Under bandit feedback, the following theorem states that there exists a loss sequence  $l_{1:T}$  such that  $R_\lambda(\pi, l_{1:T})$  is  $\Omega\left(\frac{(\lambda K)^{\frac{1}{3}}(TI)^{\frac{2}{3}}}{\log_2 T}\right)$  for any player policy  $\pi \in \Pi$ , which implies that  $R_\lambda^* = \Omega\left(\frac{(\lambda K)^{\frac{1}{3}}(TI)^{\frac{2}{3}}}{\log_2 T}\right)$ .

**Theorem 1.** *Consider the combinatorial bandit with switching costs under the bandit feedback. For any player policy  $\pi \in \Pi$ , there exists a loss sequence  $l_{1:T} \in \mathcal{L}$  that incurs an expected  $\lambda$ -switching regret of*

$$R_\lambda(\pi, l_{1:T}) = \Omega\left(\frac{(\lambda K)^{\frac{1}{3}}(TI)^{\frac{2}{3}}}{\log_2 T}\right),$$

provided that  $K \geq 3I$  and  $T \geq \max\{\frac{\lambda K}{T}, 8\}$ .

Under semi-bandit feedback, the following theorem states that there always exists a loss sequence  $l_{1:T}$  such that  $R_\lambda(\pi, l_{1:T})$  is  $\Omega(\frac{(\lambda KI)^{\frac{1}{3}} T^{\frac{2}{3}}}{\log_2 T})$  for any player policy  $\pi \in \Pi$ , which implies that  $R_\lambda^* = \Omega(\frac{(\lambda KI)^{\frac{1}{3}} T^{\frac{2}{3}}}{\log_2 T})$ .

**Theorem 2.** *Consider the combinatorial bandit with switching costs under the semi-bandit feedback. For any player policy  $\pi \in \Pi$ , there exists a loss sequence  $l_{1:T} \in \mathcal{L}$  that incurs an expected  $\lambda$ -switching regret of*

$$R_\lambda(\pi, l_{1:T}) = \Omega\left(\frac{(\lambda KI)^{\frac{1}{3}} T^{\frac{2}{3}}}{\log_2 T}\right),$$

provided that  $K \geq 3I$  and  $T \geq \max\{\frac{\lambda K}{T^2}, 6\}$ .

We observe that the orders of the lower bound of regret in Theorems 1 and 2 differ only in the combinatorial size  $I$ . Also, the lower bound in Theorem 1 is greater than that in Theorem 2 in terms of the order of  $I$  due to the fact that semi-bandit feedback results in more observations for the player compared with bandit feedback.

To prove Theorems 1 and 2, we design two stochastic loss sequences for bandit feedback and semi-bandit feedback respectively, which are presented in Algorithms 1 and 2 of Sections III-A and III-B, respectively. The design of two loss sequences generalizes the idea in [5] for the combinatorial scenarios. For both loss sequences, we apply Yao's minimax principle [12], which asserts that the regret of a randomized player against the worst-case adversary is at least the minimax regret of the optimal deterministic player against a stochastic loss sequence. By employing this principle, we are able to prove the two theorems using the constructed adversaries. Under different types of feedback, the adversaries designed are different. For a fixed time, the adversary under bandit feedback uses the same Gaussian noise for different base arms while the adversary under semi-bandit feedback uses i.i.d. Gaussian noises for different base arms.

We also design two algorithms for two types of feedback that can be shown to be almost optimal when compared to the lower bounds. Similar to [3], [4], [6], we utilize the batched algorithm to limit the number of switches between actions by dividing the whole time horizon  $T$  into batches and forcing the algorithm to play the same combinatorial arm for all the rounds within a batch. For bandit feedback, the EXP2 with John's exploration algorithm [9] is the most efficient among the existing algorithms when switching cost is not considered. We introduce a refined version of this algorithm called the BATCHED-EXP2 algorithm with John's exploration when the switching cost is involved, which achieves a regret bound as shown in below.

**Theorem 3 (Informal).** *Consider the combinatorial bandit problem under bandit feedback. For any adversary  $l_{1:T} \in \mathcal{L}$ , the policy  $\pi$  of BATCHED-EXP2 with John's exploration distribution as detailed in Algorithm 3 achieves a  $\lambda$ -switching regret of*

$$R_\lambda(\pi, l_{1:T}) = O((\lambda K)^{\frac{1}{3}} T^{\frac{2}{3}} I^{\frac{4}{3}}).$$

TABLE I: Comparison Between the Lower Bounds and Upper Bounds under Two Types of Feedback

	bandit feedback	semi-bandit feedback
lower bound	$\tilde{\Omega}((\lambda K)^{\frac{1}{3}} (TI)^{\frac{2}{3}})$	$\tilde{\Omega}((\lambda KI)^{\frac{1}{3}} T^{\frac{2}{3}})$
upper bound	$O((\lambda K)^{\frac{1}{3}} T^{\frac{2}{3}} I^{\frac{4}{3}})$	$\tilde{O}((\lambda K)^{\frac{1}{3}} T^{\frac{2}{3}} I^{\frac{2}{3}} + KI)$

For semi-bandit feedback, the BROAD algorithm [10] achieves the minimax regret when switching cost is not considered. We introduce a refined version of this algorithm called the BATCHED-BROAD algorithm when the switching costs are involved, which achieves a regret bound as shown in below.

**Theorem 4 (Informal).** *Consider the combinatorial bandit problem under semi-bandit feedback. For any adversary  $l_{1:T} \in \mathcal{L}$ , the policy  $\pi$  of BATCHED-BROAD as detailed in Algorithm 4 achieves a  $\lambda$ -switching regret of*

$$R_\lambda(\pi, l_{1:T}) = \tilde{O}((\lambda K)^{\frac{1}{3}} (TI)^{\frac{2}{3}} + KI).$$

We compare the lower bounds and upper bounds in Table I, which shows that the regret gap between BATCHED-EXP2 with John's exploration and the lower bound scales at most as  $I^{\frac{2}{3}}$  and that between BATCHED-BROAD and the lower bound scales at most  $I^{\frac{1}{3}}$ . Closing these gaps appears to be challenging and is left for future work.

### III. LOWER BOUND ANALYSIS

Following the method in [5], we apply Yao's minimax principle [12] to prove Theorems 1 and 2. The principle states that the regret of a randomized player against the worst-case loss sequence is at least the minimax regret of the optimal deterministic player against a stochastic loss sequence. Thus Theorem 1 (resp. Theorem 2) holds if we can construct a stochastic sequence of loss vectors  $L_{1:T}$  (each  $L_t = (L_{t,1}, \dots, L_{t,K}) \in [0, 1]^K$  is a random vector) such that

$$\begin{aligned} R_\lambda(\pi, L_{1:T}) &\triangleq \mathbb{E} \left[ \sum_{t=1}^T \langle A_t, L_t \rangle + \lambda \sum_{t=1}^T d(A_t, A_{t-1}) \right] - \min_{A \in \mathcal{A}} \sum_{t=1}^T \langle A, L_t \rangle \\ &= \Omega\left(\frac{(\lambda K)^{\frac{1}{3}} (TI)^{\frac{2}{3}}}{\log_2 T}\right) \quad \left( \text{resp. } \Omega\left(\frac{(\lambda KI)^{\frac{1}{3}} T^{\frac{2}{3}}}{\log_2 T}\right) \right), \end{aligned} \quad (1)$$

for any deterministic player policy  $\pi$ , where the expectation is taken over the adversary's randomized choice of loss vectors. In the following two subsections, we will judiciously construct specific loss sequences for the two feedback scenarios, which are key to proving Theorems 1 and 2.

#### A. Proof of Theorem 1: Bandit Feedback

In this section, we provide the proof of Theorem 1. First we obtain Lemma 5 for the stochastic sequence of loss vectors  $L_{1:T}$  in Algorithm 1 which is constructed by generalizing the loss sequence in [5]. Let  $\text{clip}(\alpha) := \min\{\max\{\alpha, 0\}, 1\}$  and

---

**Algorithm 1:** The Combinatorial Identical-Noise (CIN) loss sequence

---

**Input :** Time horizon  $T$ , switching cost  $\lambda$ , number of base arms  $K$  and combinatorial arm size  $I$ .

*Step 1:* Set

$$\begin{aligned} \epsilon &= \frac{(\lambda K)^{\frac{1}{3}}(IT)^{-\frac{1}{3}}}{9 \log_2 T}, \\ \sigma &= \frac{1}{6 \sqrt{\log_2 T \log_2 \frac{4T(\lambda+\epsilon)}{\epsilon}}}. \end{aligned} \quad (2)$$

Choose  $\chi \in \mathcal{A}$  uniformly at random and then generate  $W_t$ , for  $t = [T]$  according to (7).

*Step 2:* For all  $t \in [T]$  and  $x \in [K]$ , set  $x$ -th components of  $\tilde{L}_t = (\tilde{L}_{t,1}, \dots, \tilde{L}_{t,K})$  and  $L_t = (L_{t,1}, \dots, L_{t,K})$  as

$$\tilde{L}_{t,x} = W_t + \frac{1}{2} - \epsilon \chi_x, \quad L_{t,x} = \text{clip}(\tilde{L}_{t,x}). \quad (3)$$

**Output:** Loss sequence  $L_{1:T}$ .

---

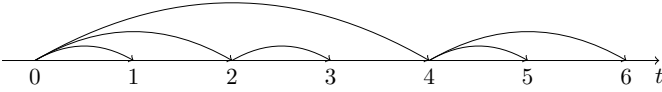


Fig. 1: An illustration of the definition of parent time  $\rho(t)$ . There is an arrow from each time  $t$ -th parent time  $\rho(t)$  to  $t$ . For example,  $\rho(3) = 2$  and  $\rho(4) = 0$ .

let  $\chi_x$  denote the  $x$ -th coordinate of a vector  $\chi \in \mathcal{A}$ . Let the *parent time* of  $t$  be

$$\rho(t) = t - 2^{\delta(t)}, \quad \delta(t) = \max\{i \geq 0 : 2^i \text{ divides } t\}. \quad (4)$$

See Figure 1 for an illustration of  $\rho(t)$ . We say a time slot  $t'$  is an *ancestor* for  $t$  if  $t' = \rho^{(c)}(t)$  for some positive integer  $c$  and  $\rho^{(c)}$  stands for the  $c$  iterated composition of  $\rho$ . Given the function  $\rho$ , recursively define as in [5] *the set of all ancestors* of  $t$  as  $\mathcal{S}(t)$  by  $\mathcal{S}(0) = \emptyset$  and

$$\mathcal{S}(t) = \mathcal{S}(\rho(t)) \cup \{\rho(t)\}, \quad t \in [T]. \quad (5)$$

The *depth* of  $\rho$  is then defined as  $d(\rho) = \max_{t \in [T]} |\mathcal{S}(t)|$ . As in [5], we define  $\text{cut}(t) = \{s \in [T] : \rho(s) < t \leq s\}$ , which is the set of time slots that are separated from their parent by  $t$ . The *width* of  $\rho$  is defined as

$$w(\rho) = \max_{t \in [T]} |\text{cut}(t)|. \quad (6)$$

The design of  $\rho$  in (4) guarantees that the depth  $d(\rho)$  and width  $w(\rho)$  are both upper bounded by  $\log_2 T + 1$  [5]. Define  $W_t$  for  $t = 1, \dots, T$  recursively by setting  $W_0 = 0$  and

$$W_t = W_{\rho(t)} + \xi_t, \quad \forall t \in [T], \quad (7)$$

where  $\xi_t$ ,  $t = 1, \dots, T$  are independent zero-mean, variance  $\sigma^2$  Gaussian variables. The design of the parent time function  $\rho(t)$  and  $W_t$  guarantees that  $\tilde{L}_{t,x}$  lies in  $[0, 1]$  with high probability, which allows us to control the difference between the  $\lambda$ -switching regrets under  $L_{1:T}$  and  $\tilde{L}_{1:T}$  when the same deterministic strategy is applied. Then we can first analyze the regret bound under the loss sequence  $\tilde{L}_{1:T}$ , which is easier

due to the fact that the random variables in it are sums of Gaussian random variables. This allows us to bound the regret under  $L_{1:T}$ . The detailed proof for this guarantee is provided in Lemma 12 of Appendix A.

The differences between the CIN loss sequence in Algorithm 1 and the loss sequence in [5] include the presence of  $I$  best arms in each round instead of one and the variations in the values of the parameters  $\epsilon$  and  $\sigma$ . As in [5], the definition of  $W_t$  induces the common uncertainty of all base arms. At each time  $t$ , the losses of the base arms in  $\chi$  are all the same and greater than other base arms by a constant  $\epsilon$ . As the player can only observe the total loss of all the base arms in the chosen combinatorial arm, it is difficult to figure out whether the loss observed is induced by the randomness of  $W_t$  or due to the chosen arms being better than other ones if the algorithm does not switch the chosen combinatorial arm for some time. Therefore, by the constructed loss sequence  $L_t$ , we can prove the lower bound in Theorem 1.

By Yao's minimax principle, Theorem 1 can be proved if the following lemma is verified.

**Lemma 5.** *Consider the combinatorial bandit problem with switching costs under the bandit feedback. Let  $L_{1:T}$  be the stochastic sequence of loss vectors defined in Algorithm 1. When  $K \geq 3I$  and  $T \geq \max\{\frac{\lambda K}{T}, 8\}$ , for any deterministic player's policy  $\pi$ , we have*

$$R_\lambda(\pi, L_{1:T}) \geq \frac{(\lambda K)^{\frac{1}{3}}(TI)^{\frac{2}{3}}}{260 \log_2 T}.$$

To prove Lemma 5, we need to analyze the expected regret under an arbitrary deterministic policy  $\pi : \mathbb{R}^{I \times T} \rightarrow \mathcal{A}^T$  when the loss sequence is  $L_{1:T}$ . Note that under  $L_{1:T}$ , a deterministic policy  $\pi$  yields an action sequence  $A_{1:T} \in \mathcal{A}^T$  so that  $A_t$  is a function of the player's past observations  $X_{1:t-1}$  with  $X_t = \langle A_t, L_t \rangle$ . We first analyze the expected regret under the same deterministic policy  $\pi$  and the loss sequence  $\tilde{L}_{1:T}$ , which is defined in (3) of Algorithm 1. Similar to  $L_{1:T}$ , the deterministic policy  $\pi$  yields an action sequence  $\tilde{A}_{1:T} \in \mathcal{A}^T$  so that  $\tilde{A}_t$  is a function of past observations  $Y_{0:t-1}$  under  $\tilde{L}_{1:T}$ , where

$$Y_0 = 1/2, \quad Y_t = \langle \tilde{A}_t, \tilde{L}_t \rangle. \quad (8)$$

Define the conditional probability measures

$$\mathcal{Q}_{\mathcal{I}}(\cdot) = P(\cdot | \chi = \mathcal{I}), \quad \mathcal{I} \in \mathcal{A}, \quad (9)$$

$$\mathcal{Q}_0(\cdot) = P(\cdot | \chi = \emptyset), \quad (10)$$

where  $\mathcal{Q}_{\mathcal{I}}$  and  $\mathcal{Q}_0$  are the probability distributions under the adversaries with  $\chi = \mathcal{I}$  and  $\chi = \emptyset$ , respectively. Thus  $\mathcal{Q}_0(\cdot)$  is the probability distribution when all arms incur the same loss. Let  $\tilde{\mathcal{F}}$  be the  $\sigma$ -algebra generated by the player's observations  $Y_{1:T}$ . The total variation distance between  $\mathcal{Q}_0$  and  $\mathcal{Q}_{\mathcal{I}}$  over  $\tilde{\mathcal{F}}$  is defined as

$$d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) = \sup_{A \in \tilde{\mathcal{F}}} |\mathcal{Q}_0(A) - \mathcal{Q}_{\mathcal{I}}(A)|. \quad (11)$$

This distance captures the player's ability to identify whether combinatorial arm  $\mathcal{I}$  is better than or equivalent to the other combinatorial arms based on the loss values he observes [5]. In the following, we first give a key lemma that relates the

player's ability to identify the best action to the number of switches he performs to or from base arms in  $\mathcal{I}$ .

Define  $Y_S = \{Y_t\}_{t \in S}$  and let  $\Delta(Y_S|Y_{S'})$  be the KL divergence between the distribution of  $Y_S$  conditioned on  $Y_{S'}$  under  $\mathcal{Q}_0$  and  $\mathcal{Q}_{\mathcal{I}}$ , i.e.,

$$\Delta(Y_S|Y_{S'}) \triangleq \mathbb{E}_{\mathcal{Q}_0} \left[ \ln \frac{\mathcal{Q}_0(Y_S|Y_{S'})}{\mathcal{Q}_{\mathcal{I}}(Y_S|Y_{S'})} \right]. \quad (12)$$

Using the chain rule,

$$\Delta(Y_{0:T}) \triangleq \Delta(Y_{0:T}|\emptyset) = \sum_{t=1}^T \Delta(Y_t|Y_{S(t)}).$$

The analysis in [5] focused on the evaluation of  $\Delta(Y_{0:T})$ , i.e. the Kullback-Leibler (KL) divergence between  $\mathcal{Q}_0(Y_{0:T})$  and  $\mathcal{Q}_{\mathcal{I}}(Y_{0:T})$ ,  $d_{\text{KL}}(\mathcal{Q}_0(Y_{0:T})\|\mathcal{Q}_{\mathcal{I}}(Y_{0:T}))$ , which has a close connection to the number of switches of arms in their setting and therefore leads to the lower bound of the regret. As [5], we define

$$M_{\mathcal{I}} \triangleq 2 \sum_{t=1}^T \langle A_t \oplus A_{t-1}, \mathcal{I} \rangle \text{ and } \tilde{M}_{\mathcal{I}} \triangleq 2 \sum_{t=1}^T \langle \tilde{A}_t \oplus \tilde{A}_{t-1}, \mathcal{I} \rangle, \quad (13)$$

as the total numbers of switches the player performs to or from base arms in  $\mathcal{I}$  during the whole time horizon under  $L_{1:T}$  and  $\tilde{L}_{1:T}$ , respectively. Define the total numbers of switches in the whole time horizon under  $L_{1:T}$  and  $\tilde{L}_{1:T}$  as

$$M \triangleq \sum_{t=1}^T d(A_t, A_{t-1}) \text{ and } \tilde{M} \triangleq \sum_{t=1}^T d(\tilde{A}_t, \tilde{A}_{t-1}). \quad (14)$$

Notice that for  $\forall i \in [K]$ , we have  $|\{\mathcal{I} \in \mathcal{A} : \mathcal{I}_i = 1\}| = \binom{K-1}{I-1}$ , where  $\mathcal{I}_i$  is the  $i$ -th component of  $\mathcal{I}$ , and then the sum over all  $\mathcal{I} \in \mathcal{A}$  of  $M_{\mathcal{I}}$  (i.e.,  $\sum_{\mathcal{I} \in \mathcal{A}} M_{\mathcal{I}}$ ) is  $\binom{K-1}{I-1}$  times the total number of switches the player performs to or from base arms in  $[K]$ . Since each switch is counted twice in the ‘‘to’’ and ‘‘from’’ directions, respectively, we conclude that  $\sum_{\mathcal{I} \in \mathcal{A}} M_{\mathcal{I}} = 2 \binom{K-1}{I-1} M$  and  $\sum_{\mathcal{I} \in \mathcal{A}} \tilde{M}_{\mathcal{I}} = 2 \binom{K-1}{I-1} \tilde{M}$ .

The upper bound for  $\Delta(Y_{0:T})$  in terms of  $\tilde{M}_{\mathcal{I}}$  and  $w(\rho)$  can be derived as follows, which directly leads to the upper bound for  $d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}})$  in terms of  $\tilde{M}_{\mathcal{I}}$  and  $w(\rho)$ . Based on Lemma 6, we can prove Lemma 5, and the whole proof is provided in Appendix A.

**Lemma 6.** *Under the loss sequence  $\tilde{L}_{1:T}$ , which is defined in (3) of Algorithm 1 with  $\chi = \mathcal{I} \in \mathcal{A}$ , it holds that  $\Delta(Y_{0:T}) \leq \frac{\epsilon^2}{2I\sigma^2} w(\rho) \mathbb{E}_{\mathcal{Q}_0}[\tilde{M}_{\mathcal{I}}]$ , which implies that*

$$d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \leq \frac{\epsilon}{2\sigma\sqrt{I}} \sqrt{w(\rho) \mathbb{E}_{\mathcal{Q}_0}[\tilde{M}_{\mathcal{I}}]}.$$

*Proof:* For any combinatorial arm  $A \in \mathcal{A}$ ,  $\langle A, \mathcal{I} \rangle$  is the number of optimal arms in the combinatorial arm. Under  $\mathcal{Q}_{\mathcal{I}}$ , by (3) and (7), we have

$$Y_{\rho(t)} = \langle \tilde{A}_{\rho(t)}, \tilde{L}_{\rho(t)} \rangle = (W_{\rho(t)} + \frac{1}{2})I - \epsilon \cdot \langle \tilde{A}_{\rho(t)}, \mathcal{I} \rangle,$$

and

$$Y_t = \langle \tilde{A}_t, \tilde{L}_t \rangle = (W_{\rho(t)} + \frac{1}{2} + \xi_t)I - \epsilon \cdot \langle \tilde{A}_t, \mathcal{I} \rangle.$$

Under  $\mathcal{Q}_0$ , similarly we have

$$Y_{\rho(t)} = (W_{\rho(t)} + \frac{1}{2})I, \quad Y_t = (W_{\rho(t)} + \frac{1}{2} + \xi_t)I.$$

Then the distribution of  $Y_t$  conditioned on  $Y_{S(t)}$  is  $N(Y_{\rho(t)}, I^2\sigma^2)$  under  $\mathcal{Q}_0$  and  $N(Y_{\rho(t)} + N_t\epsilon, I^2\sigma^2)$  under  $\mathcal{Q}_{\mathcal{I}}$ , where  $N_t = \langle \tilde{A}_{\rho(t)}, \mathcal{I} \rangle - \langle \tilde{A}_t, \mathcal{I} \rangle$  is the difference between the numbers of arms in  $\mathcal{I}$  played at time  $\rho(t)$  and  $t$ , and  $\epsilon$  is defined in (2) of Algorithm 1. Therefore,

$$\begin{aligned} \Delta(Y_t|Y_{\rho(t)}) &= \sum_{i=-I}^I \mathcal{Q}_0(N_t = i) d_{\text{KL}}(N(0, I^2\sigma^2) \| N(i\epsilon, I^2\sigma^2)) \\ &= \sum_{i=1}^I \mathcal{Q}_0(|N_t| = i) \frac{i^2\epsilon^2}{2I^2\sigma^2}, \end{aligned}$$

where  $\{|N_t| = i\}$  is the event that the player switched at least  $i$  arms from or to base arms in  $\mathcal{I}$  between rounds  $\rho(t)$  and  $t$ . We observe that  $N'_t \triangleq \langle \tilde{A}_{\rho(t)} \oplus \tilde{A}_t, \mathcal{I} \rangle$  is the total number of switches the player performs to or from base arms in  $\mathcal{I}$  between rounds  $\rho(t)$  and  $t$ . Then  $|N_t| \leq N'_t$  and

$$\begin{aligned} \Delta(Y_{0:T}) &= \frac{\epsilon^2}{2I^2\sigma^2} \sum_{t=1}^T \sum_{i=1}^I \mathcal{Q}_0(|N_t| = i) i^2 \\ &\leq \frac{\epsilon^2}{2I\sigma^2} \sum_{t=1}^T \sum_{i=1}^I \mathcal{Q}_0(|N_t| = i) i \\ &= \frac{\epsilon^2}{2I\sigma^2} \sum_{t=1}^T \mathbb{E}_{\mathcal{Q}_0}[|N_t|] \\ &\leq \frac{\epsilon^2}{2I\sigma^2} \sum_{t=1}^T \mathbb{E}_{\mathcal{Q}_0}[N'_t]. \end{aligned} \quad (15)$$

We have  $N'_t = \sum_{i \in [K]} \mathcal{I}_i \cdot \mathbb{1}_{Z_{t,i}}$ , where  $Z_{t,i} = \{\tilde{A}_{\rho(t),i} \neq \tilde{A}_{t,i}\}$  ( $\tilde{A}_{x,i}$  is used to denote the  $i$ -th component of  $\tilde{A}_x$ ). Let

$$M_i = |\{t \in [T] : \tilde{A}_{t-1,i} \neq \tilde{A}_{t,i} \text{ or } \tilde{A}_{t,i} \neq \tilde{A}_{t+1,i}\}|,$$

denote the total number of switches the player performs to or from action  $i$  during the whole time horizon. We have  $\tilde{M}_{\mathcal{I}} = \sum_{i \in [K]} \mathcal{I}_i \cdot M_i$ . Let  $s_{1:M_i,i}$  denote the time slots of switches from or to arm  $i$ , i.e.  $\tilde{A}_{s_{j,i-1,i}} \neq \tilde{A}_{s_{j,i,i}}$  or  $\tilde{A}_{s_{j,i,i}} \neq \tilde{A}_{s_{j,i+1,i}}$  for any  $j \in \{1, \dots, M_i\}$ . Since the event  $Z_{t,i}$  implies that there exists at least one time  $s$  of switch from or to action  $i$ , such that  $t \in \text{cut}(s)$ , we have

$$\begin{aligned} \sum_{t=1}^T N'_t &= \sum_{t=1}^T \sum_{i \in [K]} \mathcal{I}_i \mathbb{1}_{Z_{t,i}} \leq \sum_{i \in [K]} \mathcal{I}_i \sum_{t \in \text{cut}(s_{r,i})} \mathbb{1}_{Z_{t,i}} \\ &\leq \sum_{i \in [K]} \mathcal{I}_i \sum_{r=1}^{M_i} |\text{cut}(s_{r,i})| \leq \sum_{i \in [K]} \mathcal{I}_i M_i w(\rho) = \tilde{M}_{\mathcal{I}} w(\rho), \end{aligned} \quad (16)$$

where  $w(\rho)$  is the width of  $\rho$  defined in (6). Therefore,

$$\Delta(Y_{0:T}) \leq \frac{\epsilon^2}{2I\sigma^2} w(\rho) \mathbb{E}_{\mathcal{Q}_0}[\tilde{M}_{\mathcal{I}}].$$

---

**Algorithm 2:** The Combinatorial Diverse-Noise (CDN) loss sequence

---

**Input :** Time horizon  $T$ , number of actions  $K$  and combinatorial arm size  $I$

*Step 1:* Set  $\sigma = \frac{1}{(9 \log_2 T)}$  and  $\epsilon = \frac{(\lambda K)^{\frac{1}{3}} I^{-\frac{2}{3}} T^{-\frac{1}{3}}}{9 \log_2 T}$ . Choose  $\chi \in \mathcal{A}$  uniformly at random and then generate  $W_t^i$  for  $t \in [T]$  and  $i \in [K]$  according to (18).

*Step 2:* For  $\forall t \in [T]$  and  $\forall x \in [K]$ , set  $x$ -th components of  $\tilde{L}_t$  and  $L_t$  as

$$\tilde{L}_{tx} = W_t^x + \frac{1}{2} - \epsilon \chi_x, \quad L_{tx} = \text{clip}(\tilde{L}_{tx}). \quad (17)$$

**Output:** Loss sequence  $L_{1:T}$ .

---

By Pinsker's inequality [13, Lemma 11.6.1], we have

$$\begin{aligned} d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) &= \sup_{A \in \tilde{\mathcal{F}}} |\mathcal{Q}_0(A) - \mathcal{Q}_{\mathcal{I}}(A)| \\ &\leq \frac{\epsilon}{2\sigma\sqrt{I}} \sqrt{w(\rho) \mathbb{E}_{\mathcal{Q}_0}[\tilde{M}_{\mathcal{I}}]}. \end{aligned}$$

It may be possible to improve Lemma 6 to obtain a tighter lower bound. In the proof of Lemma 6, there is an inequality  $|N_t| \leq N'_t$  used in (15), where  $N_t = \langle \tilde{A}_{\rho(t)}, \mathcal{I} \rangle - \langle \tilde{A}_t, \mathcal{I} \rangle$  is the difference between numbers of arms in  $\mathcal{I}$  played at time  $\rho(t)$  and  $t$  and  $N'_t \triangleq \langle \tilde{A}_{\rho(t)} \oplus \tilde{A}_t, \mathcal{I} \rangle$  is the total number of switches the player performs to or from base arms in  $\mathcal{I}$  between rounds  $\rho(t)$  and  $t$ . It is easy to observe that in many cases of  $\tilde{A}_{\rho(t)}$  and  $\tilde{A}_t$ , this inequality is not tight and could even be quite loose. Therefore, the design of a loss sequence that tightens this inequality is a good future research direction to possibly obtain a tighter lower bound. ■

### B. Proof of Theorem 2: Semi-bandit Feedback

Similar to the analysis of lower bound under bandit feedback in §III-A, we will prove Theorem 2 by constructing a stochastic loss sequence  $L_{1:T}$  in Algorithm 2. Let  $\rho(t)$  be defined according to (4). For any  $i \in [K]$ , define  $W_t^i$  for  $t \in [T]$  and recursively by  $W_0^i = 0$  and

$$W_t^i = W_{\rho(t)}^i + \xi_t^i, \quad \forall t \in [T], \quad (18)$$

where  $\xi_t^i$ ,  $t \in [T]$ ,  $i \in [K]$  are independent zero-mean, variance  $\sigma^2$  Gaussian variables. The difference between the loss sequences in Algorithm 2 and Algorithm 1 lies in the independent and identically distributed (i.i.d.) nature of the added Gaussian noises  $\xi_t^i$  for each arm  $i \in [K]$ . This modification effectively tackles the challenge presented by the semi-bandit feedback scenario, where all losses for each base arm in the selected combinatorial arm are observed. Under the loss sequence  $\tilde{L}_{1:T}$  in Algorithm 1 and semi-bandit feedback, the observed losses for each base arm in the chosen combinatorial arm at time  $t \in [T]$  when  $\chi = \emptyset$  are all the same while the losses for each base arm when  $\chi \in \mathcal{A}$  may not be the same and differ by  $\epsilon$  (see definition of  $\epsilon$  in Algorithm 1). Then the supports for observed losses under the adversaries in Algorithm 1 when  $\chi = \emptyset$  and  $\chi \in \mathcal{A}$  are

different and thus the KL divergence on the observed losses under two adversaries is infinite. To overcome this issue, the loss sequence in Algorithm 2 is introduced, which induces the same support for the observed losses under both  $\chi = \emptyset$  and  $\chi \in \mathcal{A}$ . Further details will be provided in the proof of the following Lemma 7.

By Yao's minimax principle, Theorem 2 can be proved if the following lemma is verified.

**Lemma 7.** Consider the combinatorial bandit problem with switching costs under the semi-bandit feedback. Let  $L_{1:T}$  be the stochastic sequence of loss functions defined in Algorithm 2. When  $K \geq 3I$  and  $T \geq \max\{\frac{\lambda K}{I^2}, 6\}$ , for any deterministic player  $\pi$ , we have

$$R_{\lambda}(\pi, L_{1:T}) \geq \frac{(\lambda KI)^{\frac{1}{3}} T^{\frac{2}{3}}}{60 \log_2 T}.$$

To prove Lemma 7, we need to analyze the expected regret under an arbitrary deterministic policy  $\pi : [0, 1]^{I \times T} \rightarrow \mathcal{A}^T$  when the loss sequence is  $L_{1:T}$ . Note that under  $L_{1:T}$ , a deterministic policy  $\pi$  yields an action sequence  $A_{1:T} \in \mathcal{A}^T$  so that  $A_t$  is a function of the player's past observations  $Z_{1:t-1}$  with  $Z_t = A_t \circ L_t$ . Similar to §III-A, we first analyze the regret under the loss sequence  $\tilde{L}_{1:T}$  defined in Algorithm 2 that the player would suffer on the deterministic policy  $\pi \circ \text{clip}$ . The policy  $\pi \circ \text{clip}$  yields the same action sequence  $A_{1:T}$  with that of  $L_{1:T}$  under  $\tilde{L}_{1:T}$ . Thus we only need to analyze the regret under the loss sequence  $\tilde{L}_{1:T}$  and the action sequence  $A_{1:T}$ . Let

$$Y_0 = \frac{1}{2}, \quad Y_t = A_t \circ \tilde{L}_t. \quad (19)$$

and  $Y_{t,j} = \tilde{L}_{t,i} A_{t,i}$ . Define  $Y_S = \{Y_t\}_{t \in S}$  and let  $\Delta(Y_S | Y_{S'})$  be the KL divergence between the distribution of  $Y_S$  conditioned on  $Y_{S'}$  under  $\mathcal{Q}_0$  and  $\mathcal{Q}_{\mathcal{I}}$  as defined in (12), where  $\mathcal{Q}_0$  and  $\mathcal{Q}_{\mathcal{I}}$  are defined as in (9) and (10) under the loss sequence  $\tilde{L}_{1:T}$  in Algorithm 2.

In the following lemma, we derive a relation between  $\Delta(Y_{0:T})$  and  $M_{\mathcal{I}}$ , where  $M_{\mathcal{I}}$  is the total number of switches the player performs to or from base arms in  $\mathcal{I}$  during the whole time horizon and defined as (13) in §III-A. Based on the inequality in the following Lemma 8, we can prove Lemma 7 by the similar verification with [5] and the whole proof is detailed in Appendix B.

**Lemma 8.** Under the loss sequence  $\tilde{L}_{1:T}$ , which is defined in (17) of Algorithm 2 with  $\chi = \mathcal{I} \in \mathcal{A}$ , it holds that  $\Delta(Y_{0:T}) \leq \frac{\epsilon^2}{2\sigma^2} w(\rho) \mathbb{E}_{\mathcal{Q}_0}[M_{\mathcal{I}}]$ .

*Proof:* Using the chain rule,

$$\Delta(Y_{0:T}) \triangleq \Delta(Y_{0:T} | \emptyset) = \sum_{t=1}^T \Delta(Y_t | Y_{\mathcal{S}(t)}),$$

where  $\mathcal{S}(t)$  is defined as in (5). Since  $Y_{t,j}$  are independent for different  $j \in [K]$  given  $Y_{\mathcal{S}(t)}$ ,

$$\Delta(Y_t | Y_{\mathcal{S}(t)}) = \sum_{i: A_{t,i}=1} \Delta(Y_{t,i} | Y_{\mathcal{S}(t)}).$$

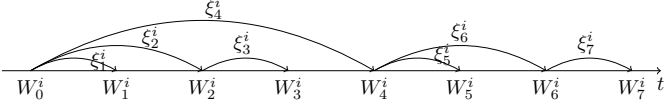


Fig. 2: An illustration of the definition of  $W_t^i$  for  $i \in [K]$ . The value of  $W_t^i$  is obtained by summing the i.i.d. Gaussian variables  $\xi_{t'}^i$ 's on the edges along the path from  $W_0^i$ , i.e. summing over all  $t' \in \mathcal{S}(t) \cup \{t\} \setminus \{0\}$ . For example,  $W_7^i = \xi_4^i + \xi_6^i + \xi_7^i$ .

Let  $\mathcal{I}_i$  denote the  $i$ -th component of  $\mathcal{I}$ . In the following, we analyze the value of  $\Delta(Y_{t,i}|Y_{\mathcal{S}(t)})$  for  $t \in [T]$  and  $i \in [K]$  such that  $A_{t,i} = 1$ .

- 1) When  $\mathcal{I}_i = 0$ , the distributions of  $Y_{t,i}$  conditioned on  $Y_{\mathcal{S}(t)}$  are the same under both  $\mathcal{Q}_0$  and  $\mathcal{Q}_{\mathcal{I}}$ . Specifically, when  $i$  is not the optimal arm, the probability distributions of  $Y_{t,i} = \tilde{L}_{t,i}$  conditioned on  $Y_{\mathcal{S}(t)}$  under  $\mathcal{Q}_0$  and  $\mathcal{Q}_{\mathcal{I}}$  are the same since the definitions of  $\tilde{L}_{t,i}$  in Algorithm 2 are the same under  $\mathcal{Q}_0$  and  $\mathcal{Q}_{\mathcal{I}}$  when  $i \notin \mathcal{I}$ . Thus  $\Delta(Y_{t,i}|Y_{\mathcal{S}(t)}) = 0$ .
- 2) When  $\mathcal{I}_i = 1$  and  $A_{t',i} = 0$  for all  $t' \in \mathcal{S}(t)$ , the distributions of  $Y_{t,i}$  conditioned on  $Y_{\mathcal{S}(t)}$  are  $N(\frac{1}{2}, |\mathcal{S}(t)|\sigma^2)$  under  $\mathcal{Q}_0$  and  $N(\frac{1}{2} - \epsilon, |\mathcal{S}(t)|\sigma^2)$  under  $\mathcal{Q}_{\mathcal{I}}$ . Then

$$\Delta(Y_{t,i}|Y_{\mathcal{S}(t)}) = d_{\text{KL}}(N(0, |\mathcal{S}(t)|\sigma^2) \| N(\epsilon, |\mathcal{S}(t)|\sigma^2)).$$

For example, suppose  $t = 7$  and  $A_{4,i} = A_{6,i} = 0$ , then we have  $W_7^i = \frac{1}{2} + \xi_4^i + \xi_6^i + \xi_7^i$  under  $\mathcal{Q}_0$  and  $W_7^i = \frac{1}{2} - \epsilon + \xi_4^i + \xi_6^i + \xi_7^i$  under  $\mathcal{Q}_{\mathcal{I}}$  (see the illustration in Figure 2). Thus  $W_7^i \sim N(\frac{1}{2}, 3\sigma^2)$  under  $\mathcal{Q}_0$  and  $W_7^i \sim N(\frac{1}{2} - \epsilon, 3\sigma^2)$  under  $\mathcal{Q}_{\mathcal{I}}$ .

- 3) When  $A_{t',i} = 1$  for some  $t' \in \mathcal{S}(t)$  and  $A_{r,i} = 0$  for  $r \in \mathcal{S}(t)$  with  $r > t'$ , the distributions of  $Y_{t,i}$  conditioned on  $Y_{\mathcal{S}(t)}$  are both  $N(Y_{t',i}, c\sigma^2)$  under  $\mathcal{Q}_0$  and  $\mathcal{Q}_{\mathcal{I}}$ , where  $c = |\mathcal{S}(t) \cup \{t\} \setminus \{0, \dots, t'\}|$ . Then  $\Delta(Y_{t,i}|Y_{\mathcal{S}(t)}) = 0$ . For example, suppose  $t = 7$ ,  $A_{4,i} = 1$  and  $A_{6,i} = 0$ , then we have  $W_7^i = Y_{4,i} + \xi_6^i + \xi_7^i$  (see the illustration in Figure 2) and thus  $W_7^i \sim N(Y_{4,i}, 2\sigma^2)$  when  $Y_{4,i}$  has been observed.

Therefore, by setting

$$N_t^* = |\{j \in [K] : A_{t,j} = \mathcal{I}_j = 1, A_{t',j} = 0, \forall t' \in \mathcal{S}(t)\}|,$$

we have

$$\begin{aligned} \Delta(Y_t|Y_{\mathcal{S}(t)}) &= \sum_{i=1}^I \mathcal{Q}_0(N_t^* = i) i d_{\text{KL}}(N(0, |\mathcal{S}(t)|\sigma^2) \| N(\epsilon, |\mathcal{S}(t)|\sigma^2)) \\ &= \sum_{i=1}^I \mathcal{Q}_0(N_t^* = i) \frac{i\epsilon^2}{2|\mathcal{S}(t)|\sigma^2}. \end{aligned}$$

We observe that  $\{N_t^* = i\}$  is the event that the player switched

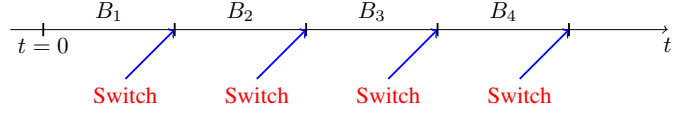


Fig. 3: An illustration of the batches in algorithm. The whole time horizon are divided into batches and during each batch, the player does not change the choice of the combinatorial arm.

at least  $i$  arms to actions in  $\mathcal{I}$  between rounds  $\rho(t)$  and  $t$ . Then

$$\begin{aligned} \Delta(Y_{0:T}) &= \frac{\epsilon^2}{2\sigma^2} \sum_{t=1}^T \frac{1}{|\mathcal{S}(t)|} \sum_{i=1}^I \mathcal{Q}_0(N_t^* = i) i \\ &\leq \frac{\epsilon^2}{2\sigma^2} \sum_{t=1}^T \mathbb{E}_{\mathcal{Q}_0}[N_t^*] \\ &\leq \frac{\epsilon^2}{2\sigma^2} \sum_{t=1}^T \mathbb{E}_{\mathcal{Q}_0}[N_t'] \\ &\leq \frac{\epsilon^2}{2\sigma^2} w(\rho) \mathbb{E}_{\mathcal{Q}_0}[M_{\mathcal{I}}], \end{aligned} \quad (20)$$

where  $N_t' \triangleq \langle \tilde{A}_{\rho(t)} \oplus \tilde{A}_t, \mathcal{I} \rangle$  is the total number of switches the player performs to or from base arms in  $\mathcal{I}$  between rounds  $\rho(t)$  and  $t$ , and the last inequality holds due to (16). ■

It may also be possible to improve Lemma 8 to obtain a tighter lower bound. In the proof of Lemma 8, there is an inequality  $N_t^* \leq N_t'$  used in (20), where

$$N_t^* = |\{j \in [K] : A_{t,j} = \mathcal{I}_j = 1, A_{t',j} = 0, \forall t' \in \mathcal{S}(t)\}|,$$

and  $N_t' \triangleq \langle A_{\rho(t)} \oplus A_t, \mathcal{I} \rangle$ . It is clear that for many cases of  $A_{\rho(t)}$  and  $A_t$ , this inequality is not tight and could even be quite loose. Therefore, the design of a loss sequence that tightens this inequality constitutes a good future research direction to possibly obtain a tighter lower bound.

#### IV. ALGORITHM FOR BANDIT FEEDBACK AND SEMI-BANDIT FEEDBACK

In this section, we will introduce our algorithms for the two types of feedback. We will use the batched algorithm to restrict the number of switches between actions by dividing the whole time horizon into batches and forcing the algorithm to play the same action for all the rounds within a batch as shown in Figure 3.

##### A. Algorithm for Bandit Feedback

The Exp2 with John's exploration algorithm [9] is an efficient algorithm for the combinatorial bandit problem under bandit feedback. In Algorithm 3, we introduce a refinement of this algorithm, called BATCHED-EXP2 with John's exploration to take into account switching costs, where John's exploration distribution can be obtained in [9, § 7.3.2]. In this section, we prove the following theorem, which is a formal version of Theorem 3, to obtain a bound for the regret of the proposed algorithm.

**Algorithm 3:** BATCHED-EXP2 with John's exploration

**Input :** John's exploration distribution  $\mu$  over  $\mathcal{A}$ ;  
 batch lengths  $B_1, \dots, B_N$  s.t.  $\sum_{i=1}^N B_n = T$ ;  
 mixing coeff.  $\gamma \in (0, 1)$  and learning rate  $\eta$ ;  
 $q_1 = (\frac{1}{|\mathcal{A}|}, \dots, \frac{1}{|\mathcal{A}|}) \in \mathbb{R}^{|\mathcal{A}|}$ .

**for**  $1 \leq n \leq N$  **do**

(a) Let  $p_n = (1 - \gamma)q_n + \gamma\mu$ , and select a combinatorial arm  $A(n)$  with respect to  $p_n$ . Pull the selected combinatorial arm for  $B_n$  times, which then incurs a loss  $X(n) = \langle A(n), l(n) \rangle$ , where  $l(n) = \sum_{b \in [B_n]} l(n, b)$  and  $l(n, b)$  is the loss vector at the  $b$ -th time due to the pulling of  $A(n)$  in this batch.

(b) Estimate the loss vector  $l(n)$  by  $\tilde{l}(n) = X(n)\Sigma_{n-1}^+ A(n)$ , with  $\Sigma_{n-1} = \mathbb{E}_{A \sim p_n}[AA^T]$  where  $\Sigma_{n-1}^+$  is the pseudo-inverse of  $\Sigma_{n-1}$ .

(c) Update the exponential weights. That is, for all  $A \in \mathcal{A}$ ,

$$q_{n+1}(A) = \frac{q_n(A) \exp(-\eta \langle A, \tilde{l}(n) \rangle)}{\sum_{A' \in \mathcal{A}} q_n(A') \exp(-\eta \langle A', \tilde{l}(n) \rangle)}.$$

**end**

**Theorem 9.** Let  $\pi$  be the policy of BATCHED-EXP2 with John's exploration distribution. The time horizon  $T$  is divided into  $N$  batches with lengths satisfying  $B_n = B = \left\lceil \lambda^{\frac{2}{3}} K^{-\frac{1}{3}} T^{\frac{1}{3}} I^{-\frac{1}{3}} \right\rceil$  for  $1 \leq n \leq \lfloor \frac{T}{B} \rfloor$  and  $B_N = T - (N-1)B$  with  $N = \lfloor \frac{T}{B} \rfloor + 1$ . Let  $\gamma = \eta BIK$  and  $\eta = \sqrt{\frac{\ln(K)}{3NK(BI)^2}}$ . Then for any adversary  $l_{1:T} \in \mathcal{L}$ , the  $\lambda$ -switching regret satisfies

$$R_\lambda(\pi, l_{1:T}) = O((\lambda K)^{\frac{1}{3}} T^{\frac{2}{3}} I^{\frac{4}{3}}).$$

*Proof of Theorem 9:* By definition, we have  $l(n) = \sum_{b=1}^{B_n} l_{(n-1)B+b} \in [0, B]$ . Since  $A(n) \in \mathcal{A}$ , the accumulated loss in each batch satisfies

$$X(n) = \langle A(n), l(n) \rangle \leq BI.$$

By [9, Theorem 7.6], when  $\gamma = \eta BIK$  and  $\eta = \sqrt{\frac{\ln|\mathcal{A}|}{3NK(BI)^2}}$ , the pseudo-regret satisfies

$$\begin{aligned} R(\pi, l_{1:T}) &\leq 2BI\sqrt{3NK \ln|\mathcal{A}|} \\ &\leq 4\lambda^{\frac{2}{3}} K^{-\frac{1}{3}} T^{\frac{1}{3}} I^{-\frac{1}{3}} \sqrt{6\lambda^{-\frac{2}{3}} K^{\frac{1}{3}} T^{\frac{2}{3}} I^{\frac{1}{3}} KI \ln \frac{eK}{I}} \\ &= 4\sqrt{6 \ln \frac{eK}{I}} \lambda^{\frac{1}{3}} K^{\frac{1}{3}} T^{\frac{2}{3}} I^{\frac{4}{3}}. \end{aligned}$$

Thus

$$\begin{aligned} R_\lambda(\pi, l_{1:T}) &\leq R_T(\pi, l_{1:T}) + \lambda IN \\ &\leq 4\sqrt{6 \ln \frac{eK}{I}} \lambda^{\frac{1}{3}} K^{\frac{1}{3}} T^{\frac{2}{3}} I^{\frac{4}{3}} + 2\lambda^{\frac{1}{3}} K^{\frac{1}{3}} T^{\frac{2}{3}} I^{\frac{4}{3}} \\ &= \left(4\sqrt{6 \ln \frac{eK}{I}} + 2\right) \lambda^{\frac{1}{3}} K^{\frac{1}{3}} T^{\frac{2}{3}} I^{\frac{4}{3}}. \end{aligned}$$

**Algorithm 4:** BATCHED-BROAD

**Define:**  $F_n(a) = \frac{1}{\eta_n} \sum_{i=1}^K \ln \frac{1}{a_i}$ .

**Input :**  $\eta_1 = \eta$ ;  
 batch lengths  $B_1, \dots, B_N$  s.t.  $\sum_{i=1}^N B_n = T$ ;  
 $n = 1, N_0 = 0$ .

**for**  $\beta = 0, 1, \dots$  **do**

$a'_n = \arg \min_{a \in \text{Co}(\mathcal{A})} F_1(a)$ .

**while**  $n \leq N$  **do**

1)  $a_n = \arg \min_{a \in \text{Co}(\mathcal{A})} \{D_{F_n}(a, a'_n)\}$ , where  $D_{F_n}$  is defined as in (21).

2) Sample  $A(n)$  such that  $\mathbb{E}[A(n)] = a_n$  and then pull it for  $B_n$  times. Observe  $A(n) \circ l(n, b)$  for  $b \in [B_n]$  and incur a loss

$$X(n) = \langle A(n), l(n) \rangle,$$

where  $l(n) = \sum_{b \in [B_n]} l(n, b)$  with  $l(n, b)$  being the loss vector at  $b$ -th time of pulling  $A(n)$  in this batch.

3) Compute the estimator  $\hat{l}(n)$  with

$$(\hat{l}(n))_i = \frac{(A(n))_i (l(n))_i}{(a_n)_i},$$

for  $i \in [K]$ , where we use  $(v)_i$  to denote the  $i$ th component in  $v$ .

4) Update

$$a'_{n+1} = \arg \min_{a \in \text{co}(\mathcal{A})} \langle a, \hat{l}(n) \rangle + D_{F_n}(a, a'_n).$$

**if**  $\sum_{s=N_{\beta+1}}^n \|A(s) \circ l(s)\|_2^2 \geq \frac{K \ln T}{3\eta^2}$  **then**

$\eta_{n+1} \leftarrow \eta_n/2$ ,  $N_{\beta+1} \leftarrow n$ ,  $n \leftarrow n+1$ ;  
     **break**;

**end**

$\eta_{n+1} \leftarrow \eta_n$ ,  $n \leftarrow n+1$ .

**end**

**B. Algorithm for Semi-Bandit Feedback**

We propose the BATCHED-BROAD algorithm as stated in Algorithm 4, based on the BROAD algorithm in [10], which is an Online Mirror Descent algorithm with log-barrier regularizer. For a regularizer  $F: \mathbb{R}^K \rightarrow \mathbb{R}$ , define

$$D_F(p, q) \triangleq F(p) - F(q) - \langle \nabla F(q), p - q \rangle. \quad (21)$$

In the following theorem, we first prove BATCHED-BROAD can achieve a  $\lambda$ -switching regret as shown in Theorem 4.

**Theorem 10.** Let  $\pi$  be the policy of BATCHED-BROAD. The time horizon  $T$  is divided into  $N$  batches with lengths satisfying  $B_n = B = \left\lceil (TI)^{\frac{1}{3}} \lambda^{\frac{2}{3}} K^{-\frac{1}{3}} + 1 \right\rceil$ , for  $n = 1, \dots, N-1$  and  $B_N = T - \sum_{n=1}^{N-1} B_n$  with  $N = \lfloor \frac{T}{B} \rfloor + 1$ . Let  $\eta = \min\{\frac{1}{18IB^2}, \frac{1}{81}\}$ . Then for any adversary  $l_{1:T} \in \mathcal{L}$ , when  $T \geq \frac{K}{I\lambda^2}$  the  $\lambda$ -switching regret satisfies

$$R_\lambda(\pi, l_{1:T}) = \tilde{O}((\lambda K)^{\frac{1}{3}} (TI)^{\frac{2}{3}} + KI).$$



In the following lemma, we first give the generalized analysis of BROAD algorithm for losses in varying intervals over time. Based on the lemma, we can then prove Theorem 10.

**Lemma 11.** *Consider a general combinatorial multi-armed bandit problem with semi-bandit feedback where  $l_t \in [0, b_t]^K$  and  $\mathcal{A} = \{A \in \{0, 1\}^K : \|A\|_1 = I\}$ . Let the agent's policy  $\pi$  be obtained by Algorithm 4 with batch lengths  $B_n = 1$  for all  $n \leq N$ . If for all  $t \leq T$ ,  $\eta_t \leq \min\{\frac{1}{18Ib_t^2}, \frac{1}{8I}\}$ , then the pseudo-regret satisfies*

$$R(\pi, l_{1:T}) = O\left(\sqrt{(KI \ln T) \sum_{t=1}^T b_t^2} + KI \ln T\right)$$

*Proof:* For the combinatorial arm  $A_t$  pulled at time  $t$  and the loss vector  $l_t \in [0, b_t]^K$  at time  $t$ , we have

$$\|A_t \circ l_t\|_2^2 = \sum_{i:A_{t,i}=1} l_{t,i}^2 \leq I b_t^2.$$

Then we have  $\eta_t \|A_t \circ l_t\|_2^2 \leq \frac{1}{18}$  and thus by [10, Theorem 8], we have

$$R(\pi, l_{1:T}) = O\left(\mathbb{E}\left[\sqrt{(K \ln T) \sum_{t=1}^T \|A_t \circ l_t\|_2^2} + KI \ln T\right]\right),$$

where the expectation is taken over the randomized choice of  $A_t$ . Since  $\|A_t \circ l_t\|_2^2 \leq I b_t^2$ , the proof is completed. ■

*Proof of Theorem 10:* In  $n$ -th batch, the accumulated loss vector  $l(n) \in [0, B_n]^K$ . By Lemma 11, when  $T \geq \frac{K}{I\lambda^2}$  the pseudo-regret satisfies

$$\begin{aligned} R(\pi, l_{1:T}) &= O\left(\sqrt{(KI \ln T) \sum_{n=1}^N B_n^2} + KI \ln T\right) \\ &\leq O\left(\sqrt{(2KI \ln T)T(TI)^{\frac{1}{3}}\lambda^{\frac{2}{3}}K^{-\frac{1}{3}}} + KI \ln T\right) \\ &= O\left(\sqrt{2 \ln T}(\lambda K)^{\frac{1}{3}}(TI)^{\frac{2}{3}} + KI \ln T\right). \end{aligned}$$

Since  $N \leq \frac{T}{(TI)^{\frac{1}{3}}\lambda^{\frac{2}{3}}K^{-\frac{1}{3}}} + 1$ , the  $\lambda$ -switching regret satisfies

$$\begin{aligned} R_\lambda(\pi, l_{1:T}) &\leq R(\pi, l_{1:T}) + \lambda IN \\ &= O(\sqrt{\ln T}(\lambda K)^{\frac{1}{3}}(TI)^{\frac{2}{3}} + KI \ln T). \end{aligned}$$

## V. NUMERICAL RESULTS

In this section, we present the numerical results to compare our algorithms BATCHED-EXP2 with John's exploration in Algorithm 3 and BATCHD-BROAD in Algorithm 4 with some baselines from the literature after adding batches in which the played combinatorial arm does not change within each batch. For bandit feedback, our baseline algorithm is the EXP3 algorithm [14] and we modified it to be a batched algorithm called BATCHED-EXP3. For the semi-bandit feedback, we choose the Follow-the-Regularized-Leader algorithm with hybrid regularizer [15]  $F(a) = \sum_{i=1}^K -\sqrt{a_i} + \gamma(1 - a_i) \log(1 - a_i)$  and the unnormalized negentropy potential [1]  $F(a) = \sum_{i=1}^K (a_i \ln a_i - a_i)$ . We call the modification of

these two algorithms the BATCHED-HYBRID and BATCHED-NEGENTROPY, respectively.

Since the optimal adversarial adversary is difficult to design, we use the CIN loss sequence and CDN loss sequence given by Algorithms 1 and 2 in §III for bandit and semi-bandit feedback, respectively. Besides the lower-bound traces, we also design a stochastically constrained (SC) adversary which is very similar to that used in [15]. Specifically, the time horizon of length  $T$  is split into phases:

$$\underbrace{1, 2, \dots, t_1}_{T_1}, \underbrace{t_1 + 1, \dots, t_2}_{T_2}, \dots, \underbrace{t_{n-1}, \dots, T}_{T_n},$$

where the length of phase  $i$  is  $T_i = \lfloor 1.6^i \rfloor$ . The loss for each arm  $i$  at time  $t$  is set to be an independent Bernoulli distribution with mean

$$\mu_{ti} = \begin{cases} 1 - \check{\alpha}\lambda & \text{if } i \leq I \\ 1 & \text{else} \end{cases}$$

if  $t$  belongs to an odd phase and

$$\mu_{ti} = \begin{cases} 0 & \text{if } i \leq I \\ \check{\alpha}\lambda & \text{else} \end{cases}$$

otherwise. In the above setting,  $\lambda$  is the switching cost. We denote the SC adversary with parameter  $\check{\alpha}$  by SC( $\check{\alpha}$ ) adversary. Note that the mean of the optimal arm oscillates between being close to 1 and close to 0 to create a challenging environment for our bandit algorithms.

### A. Bandit Feedback

We use two types of adversaries to compare the performance of BATCHED-EXP2 with John's exploration and BATCHED-EXP3 algorithm. The batch length in the algorithms is fixed to be  $B = \lfloor 3\lambda^{\frac{2}{3}}K^{-\frac{1}{3}}(TI)^{\frac{1}{3}} \rfloor$ .

First, we use the lower-bound trace CIN adversary that we designed in Algorithm 1 with  $\sigma = \frac{10}{9 \log_2 T}$  and  $\epsilon = \frac{10(\lambda K)^{\frac{1}{3}}(IT)^{-\frac{1}{3}}}{9 \log_2 T}$ . From Figure 4a, we observe that the  $\lambda$ -switching regret of BATCHED-EXP2 with John's exploration is much smaller than that of BATCHED-EXP3 when  $K = 10$ ,  $I = 3$ ,  $\lambda = 1$ . From Figure 4b, we observe that the  $\lambda$ -switching regret of BATCHED-EXP2 with John's exploration is much smaller than that of BATCHED-EXP3 for a smaller value of  $\lambda = 0.1$ , showing that even when the switching cost is small, our algorithm outperforms the benchmark significantly. In Figure 4e, we compare the  $\lambda$ -switching regret for different values of  $I$  when  $K = 20$ ,  $T = 10000$  and  $\lambda = 1$ . It is observed that the regret grows as  $I^{0.304}$ ; the dependence appears to be loose with respect to the upper bound of  $I^{\frac{1}{3}}$  in Theorem 9 and suggests that the BATCHED-EXP2 algorithm works well under the CIN loss sequence. Also, we observe that  $I^{0.304}$  is also loose in terms of lower bound  $I^{\frac{2}{3}}$  in Theorem 1, which can be explained by the following statement. Under the CIN loss sequence  $L_{1:T}$  and the policy  $\pi^{\text{Exp2}}$  of the BATCHED-EXP2 algorithm, we delineate two reasons that explain why the  $\lambda$ -switching regret  $R_\lambda(\pi^{\text{Exp2}}, L_{1:T})$  is not lower bounded by  $\Omega(I^{\frac{2}{3}})$ .

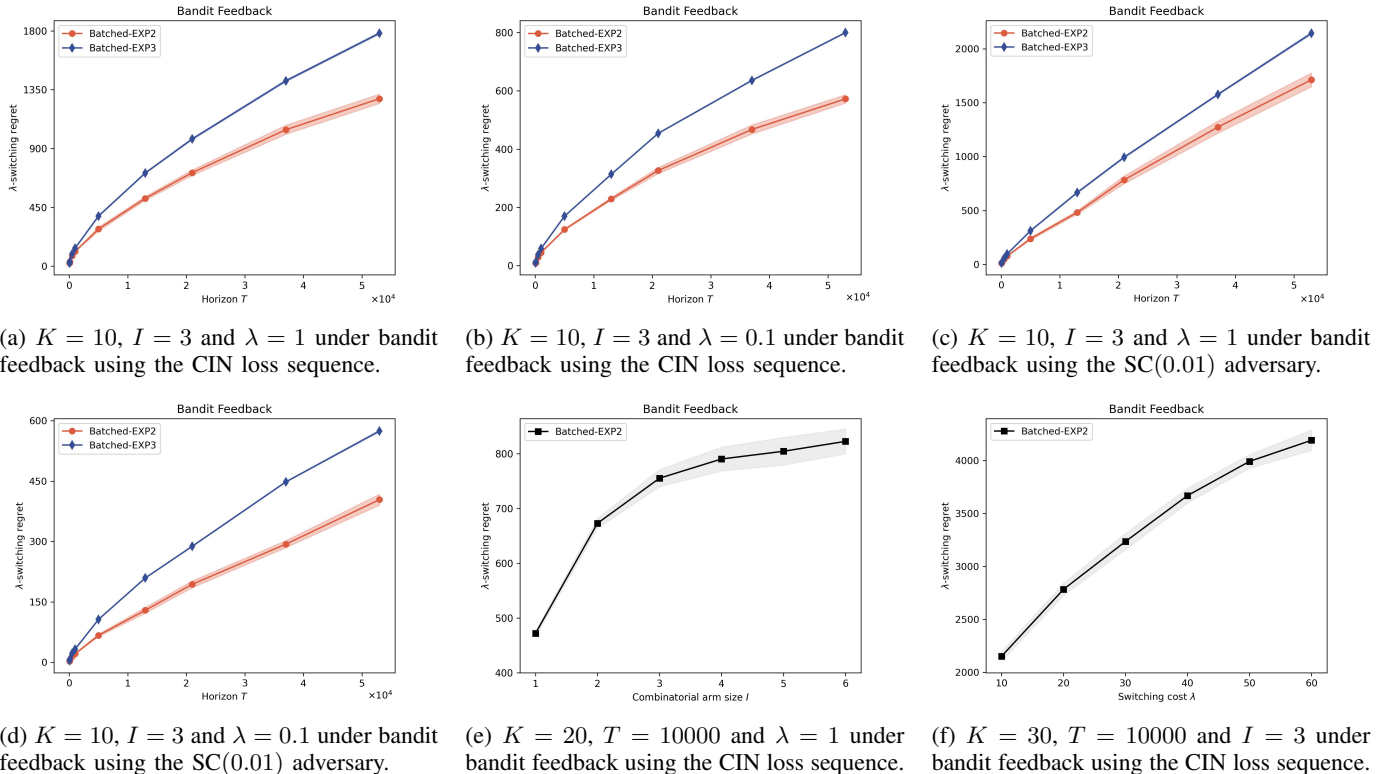


Fig. 4: Comparison of the performance of different algorithms under Bandit Feedback

- Let  $\mathcal{V}$  be the set of all stochastic loss sequences  $V_{1:T}$  ( $V_t \in [0, 1]^K$  is a random vector). Given the combinatorial arm size  $K$ , the switching cost  $\lambda$  and the time horizon  $T$ , we have

$$\begin{aligned} & \inf_{\pi \in \Pi} \sup_{V_{1:T} \in \mathcal{V}} R_\lambda(\pi, V_{1:T}) \\ & \geq \inf_{\pi \in \Pi} \sup_{l_{1:T} \in \mathcal{L}} R_\lambda(\pi, l_{1:T}) \geq \Omega(I^{\frac{2}{3}}), \end{aligned}$$

where the last inequality holds due to Theorem 1. It is not comparable between  $R_\lambda(\pi^{\text{Exp2}}, L_{1:T})$  and  $\inf_{\pi \in \Pi} \sup_{V_{1:T} \in \mathcal{V}} R_\lambda(\pi, V_{1:T})$ . Thus  $R_\lambda(\pi^{\text{Exp2}}, L_{1:T})$  may not be  $\Omega(I^{\frac{2}{3}})$ .

- Given the combinatorial arm size  $K$ , the switching cost  $\lambda$  and the time horizon  $T$ , by Lemma 5, for any deterministic player  $\pi$ , we have

$$R_\lambda(\pi, L_{1:T}) \geq \Omega(I^{\frac{2}{3}}),$$

under the CIN loss sequence  $L_{1:T}$ . Since the policy  $\pi^{\text{Exp2}}$  of the BATCHED-EXP2 algorithm is stochastic,  $R_\lambda(\pi^{\text{Exp2}}, L_{1:T})$  may not be  $\Omega(I^{\frac{2}{3}})$ .

In Figure 4f, we compare the  $\lambda$ -switching regret for different values of  $\lambda$  when  $K = 30, T = 10000$ , and  $I = 3$ . It is observed that the regret grows as  $\lambda^{0.379}$ . Note that our theoretical results in Theorems 1 and 9 say that the expected  $\lambda$ -switching regret scales as  $\Theta(\lambda^{1/3})$  when  $T, K$  and  $I$  are fixed. Even though the empirical observation of the regret scaling as  $\lambda^{0.379}$  cannot be directly compared to the theoretical result of  $\lambda^{1/3}$  because, among other reasons, the loss sequence constructed here is, in fact, stochastically constrained, the fact

that the exponents of  $\lambda$  are not too far from each other is reassuring.

The second trace we used is the SC(0.01) adversary. From Figure 4c and Figure 4d, we observe that the  $\lambda$ -switching regrets of BATCHED-EXP2 with John's exploration are both smaller than those of BATCHED-EXP3 when  $K = 10, I = 3$  and  $\lambda = 1$ , and  $K = 10, I = 3$  and  $\lambda = 0.1$ , respectively. This again corroborates the efficacy of our proposed methods.

## B. Semi-bandit Feedback

We compare BATCHED-BROAD, BATCHED-HYBRID and BATCHED-NEGENTROPY algorithm under the lower-bound trace CDN adversary that we designed in Algorithm 2 with  $\sigma = \frac{10}{9 \log_2 T}$  and  $\epsilon = \frac{10(\lambda k)^{\frac{1}{3}} I^{-\frac{2}{3}} T^{-\frac{1}{3}}}{9 \log_2 T}$  and the SC(0.005) adversary. The batch length of the algorithms is fixed to be  $B = \left\lceil 3\lambda^{\frac{2}{3}} K^{-\frac{1}{3}} T^{\frac{1}{3}} I^{\frac{2}{3}} \right\rceil$ .

Under the CDN adversary, we have the following results. From Figure 5a, we observe that the  $\lambda$ -switching regret of BATCHED-BROAD is much smaller than that of BATCHED-HYBRID and BATCHED-NEGENTROPY when  $K = 10, I = 3, \lambda = 1$ . From Figure 5b, we observe that the  $\lambda$ -switching regret of BATCHED-BROAD is much smaller than that of BATCHED-HYBRID and BATCHED-NEGENTROPY for a smaller value  $\lambda = 0.1$ , showing that even when the switching cost is small, our algorithm outperforms the benchmark significantly. In Figure 5e, we compare the  $\lambda$ -switching regret for different values of  $I$  when  $K = 40$ ,

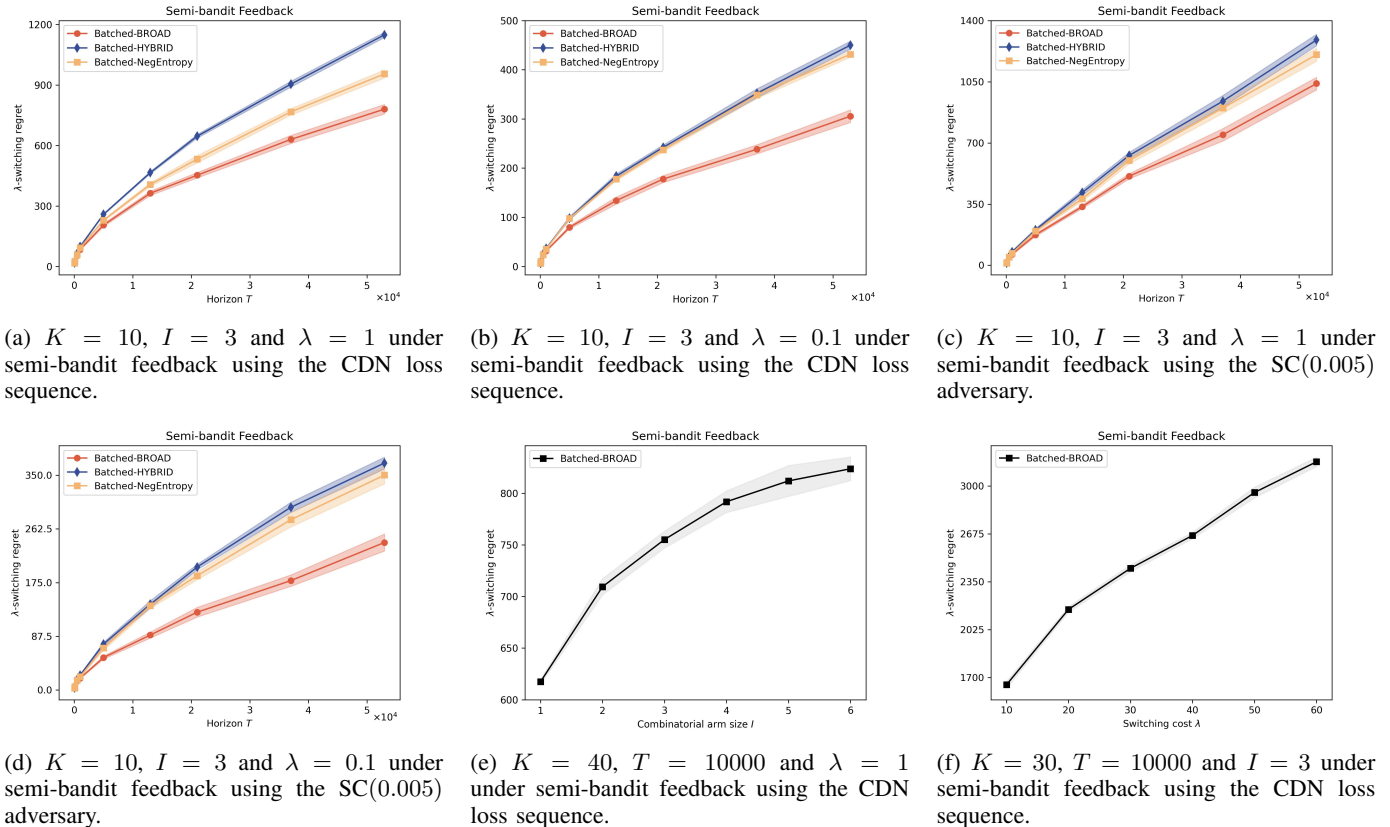


Fig. 5: Comparison of the performances of different algorithms under Semi-bandit Feedback

$T = 10000$  and  $\lambda = 1$ . It is observed that the regret grows as  $I^{0.163}$ ; the dependence appears to be loose with respect to the upper bound of  $I^{\frac{2}{3}}$  in Theorem 10 and suggests that the BATCHED-BROAD algorithm works well under the CDN loss sequence. Also, we observe that  $I^{0.163}$  is also loose in terms of lower bound  $I^{\frac{1}{3}}$  in Theorem 2, which is possible by a similar reasoning as that for bandit feedback in §V-A. In Figure 5f, we compare the  $\lambda$ -switching regret for different values of  $\lambda$  when  $K = 30, T = 10000$ , and  $I = 3$ . It is observed that the regret grows as  $\lambda^{0.356}$ , which is again close to  $\lambda^{1/3}$  as given by our theoretical result in Theorems 2 and 10.

Under the SC(0.005) adversary, from Figure 5c and Figure 5d, we observe that the  $\lambda$ -switching regrets of BATCHED-BROAD are both smaller than that of BATCHED-HYBRID and BATCHED-NEGENTROPY when  $K = 10, I = 3$  and  $\lambda = 1$ , and  $K = 10, I = 3$  and  $\lambda = 0.1$ , respectively.

## VI. CONCLUSION

We derived lower bounds for the minimax regret for the problem of adversarial combinatorial bandit with a switching cost  $\lambda$  for each changed arm in each round. We also designed algorithms that operate in batches to approach the lower bounds. Our findings provide insights into the inherent difficulty of the problem and suggest efficient approaches to minimize switching costs and optimize the overall performance in terms of regret. Further research involves deriving tighter bounds in both directions for both bandit and semi-bandit feedback. Also, other sets of combinatorial arms such as

those involved in the shortest path problem, ranking problems, and multitask problems can also be considered when switching costs are involved.

## APPENDIX

### A. Proof of Lemma 5

Given the constructed stochastic loss sequence  $L_{1:T}$  defined in Algorithm 1, we now want to analyze the player's expected regret under an arbitrary deterministic policy  $\pi$  which yields an action sequence  $A_{1:T} \in \mathcal{A}^T$  so that  $A_t$  is a function of the player's past observations  $X_{1:t-1}$  with  $X_t = \langle A_t, L_t \rangle$ . First we define

$$\tilde{R} \triangleq \sum_{t=1}^T \langle A_t, L_t \rangle + \lambda \sum_{t=1}^T d(A_t, A_{t-1}) - \min_{A \in \mathcal{A}} \sum_{t=1}^T \langle A, L_t \rangle. \quad (22)$$

We also define the regret with respect to the unclipped stochastic loss functions  $\tilde{L}_{1:T}$  defined in Algorithm 1 under the same deterministic policy  $\pi$  which yields an action sequence  $\tilde{A}_{1:T} \in \mathcal{A}^T$  so that  $\tilde{A}_t$  is a function of the player's past observations  $Y_{1:t-1}$  with  $Y_t = \langle \tilde{A}_t, \tilde{L}_t \rangle$ . Let

$$\tilde{R} \triangleq \sum_{t=1}^T \langle \tilde{A}_t, \tilde{L}_t \rangle + \lambda \sum_{t=1}^T d(\tilde{A}_t, \tilde{A}_{t-1}) - \min_{A \in \mathcal{A}} \sum_{t=1}^T \langle A, \tilde{L}_t \rangle.$$

Then  $\mathbb{E}[R] = R_\lambda(\pi, L_{1:T})$  where the expectation in  $\mathbb{E}[R]$  is taken over the adversary's randomized choice of the loss

sequences  $L_{1:T}$ , and  $R_\lambda(\pi, L_{1:T})$  is defined in (1). The next lemma shows that in expectation, the regret  $\mathbb{E}[R]$  can be lower bounded in terms of  $\mathbb{E}[\tilde{R}]$  (the expectation in  $\mathbb{E}[\tilde{R}]$  is taken over the adversary's randomized choice of the loss sequences  $\tilde{L}_{1:T}$ ).

**Lemma 12.** *Assume that  $T \geq \max\{\frac{\lambda K}{T}, 6\}$ . Then  $\mathbb{E}[R] \geq \mathbb{E}[\tilde{R}] - \frac{\epsilon TI}{4}$ .*

*Proof:* We consider the event  $B = \{\forall t : L_t = \tilde{L}_t\}$ , and first show that  $P(B) \geq 1 - \frac{\epsilon}{4(\lambda + \epsilon)}$ . For  $\delta = \frac{\epsilon}{4(\lambda + \epsilon)}$ , by [5, Lemma 1] we have that with probability at least  $1 - \delta$ ,

$$|W_t| \leq \sigma \sqrt{2d(\rho) \ln \frac{T}{\delta}} \leq 2\sigma \sqrt{\log_2 T \log_2 \frac{4T(\lambda + \epsilon)}{\epsilon}},$$

for all  $t \in [T]$ , where the last inequality holds due to  $d(\rho) \leq \log_2 T + 1$  by [5, Lemma 2]. Thus, setting  $\sigma = \frac{1}{6\sqrt{\log_2 T \log_2 \frac{4T(\lambda + \epsilon)}{\epsilon}}}$ , we obtain that

$$P\left(\forall t \in [T], \frac{1}{2} + W_t \in \left[\frac{1}{6}, \frac{5}{6}\right]\right) \geq 1 - \delta.$$

For  $T \geq \max\{\frac{\lambda K}{T}, 6\}$ , we have  $\epsilon < \frac{1}{6}$  and thus  $\tilde{L}_t(x) \in [0, 1]$  for all  $x \in [K]$  whenever  $1/2 + W_t \in [\frac{1}{6}, \frac{5}{6}]$ . This implies that  $P(B) \geq 1 - \delta$ .

If  $B$  occurs then  $R = \tilde{R}$ ; otherwise,  $\tilde{R} - R \leq (\lambda + \epsilon)TI$  since  $R, \tilde{R} \in [0, (\lambda + \epsilon)TI]$ . Therefore,

$$\mathbb{E}[\tilde{R}] - \mathbb{E}[R] = \mathbb{E}[\tilde{R} - R | \neg B] \cdot P(\neg B) \leq \frac{\epsilon TI}{4}.$$

Let  $\mathcal{Q}_0$  and  $\mathcal{Q}_I$  follow previous definition in (9) and (10),  $\tilde{\mathcal{F}}$  be the  $\sigma$ -algebra generated by  $X_{1:T}$  defined in (8) and  $d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_I)$  follow the definition in (11).  $\tilde{M}_I$  and  $\tilde{M}$  are defined in (13) and (14), respectively. Then we have the following lemma that bounds total variation from above.

*Remark 1.* Note that  $\tilde{A}_t$  is a deterministic function of its past observations  $Y_{1:t-1}$ ; thus the  $\sigma$ -algebra generated by  $\tilde{A}_{1:T}$  is a subset of  $\tilde{\mathcal{F}}$ .

**Lemma 13.** *It holds that*

$$\frac{1}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_I) \leq \frac{\epsilon}{\sigma \sqrt{2K}} \sqrt{(\log_2 T + 1) \mathbb{E}_{\mathcal{Q}_0}[\tilde{M}]}.$$

*Proof:* By [5, Lemma 2], the width  $w(\rho) \leq \log_2 T + 1$ . Then by Lemma 6, we have

$$d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_I) \leq \frac{\epsilon}{2\sigma\sqrt{I}} \sqrt{(\log_2 T + 1) \mathbb{E}_{\mathcal{Q}_0}[\tilde{M}_I]}.$$

Then using the concavity of the squared root function and by  $\sum_{\mathcal{I} \in \mathcal{A}} \tilde{M}_I = 2 \binom{K-1}{I-1} \tilde{M}$ , it holds that

$$\begin{aligned} & \frac{1}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_I) \\ & \leq \frac{\epsilon}{2\sigma\sqrt{I}} \sqrt{\log_2 T + 1} \sum_{\mathcal{I} \in \mathcal{A}} \frac{1}{\binom{K}{I}} \sqrt{\mathbb{E}_{\mathcal{Q}_0}[\tilde{M}_I]} \\ & \leq \frac{\epsilon}{2\sigma\sqrt{I}} \sqrt{\log_2 T + 1} \sqrt{\mathbb{E}_{\mathcal{Q}_0} \left[ \frac{1}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} \tilde{M}_I \right]} \\ & = \frac{\epsilon}{\sqrt{2}\sigma} \sqrt{\frac{(\log_2 T + 1) \mathbb{E}_{\mathcal{Q}_0}[\tilde{M}]}{K}}. \end{aligned}$$

**Lemma 14.** *It holds that*

$$\mathbb{E}[\tilde{R}] \geq \epsilon TI \left(1 - \frac{I}{K}\right) - \frac{\epsilon TI}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_I) + \lambda \mathbb{E}[\tilde{M}].$$

*Proof:* For any  $i \in [K]$ , let  $T_i$  denote the number of rounds the player picks arm  $i$  in the action sequence  $\tilde{A}_{1:T}$ . So we can write  $\tilde{R} = \epsilon \left( TI - \sum_{i \in [K]} \chi_i T_i \right) + \lambda \tilde{M}$ , where we use  $\chi_i$  to denote the  $i$ -th component of  $\chi$ . Also, we use  $\mathcal{I}_i$  to denote the  $i$ -th component of  $\mathcal{I} \in \mathcal{A}$  in the following. Since  $\chi \in \mathcal{A}$  is selected uniformly at random in Algorithm 1, we have

$$\begin{aligned} \mathbb{E}[\tilde{R}] & = \frac{1}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} \mathbb{E} \left[ \epsilon \left( TI - \sum_{i \in [K]} \mathcal{I}_i T_i \right) + \lambda \tilde{M} \mid \chi = \mathcal{I} \right] \\ & = \epsilon TI - \frac{\epsilon}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} \sum_{i \in [K]} \mathcal{I}_i \mathbb{E}_{\mathcal{Q}_I} [T_i] + \lambda \mathbb{E}[\tilde{M}]. \end{aligned}$$

For all  $i \in [K]$  and  $t \in [T]$ , the event  $\{\tilde{A}_{t,i} = 1\}$  belongs to the  $\sigma$ -field  $\tilde{\mathcal{F}}$  by Remark 1 ( $\tilde{A}_{x,i}$  is used to denote the  $i$ -th component of  $\tilde{A}_x$ ), so we have

$$\mathcal{Q}_I(\tilde{A}_{t,i} = 1) - \mathcal{Q}_0(\tilde{A}_{t,i} = 1) \leq d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_I).$$

Summing over  $t \in [T]$  yields

$$\mathbb{E}_{\mathcal{Q}_I} [T_i] - \mathbb{E}_{\mathcal{Q}_0} [T_i] \leq T d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_I).$$

Summing over  $i \in [K]$  such that  $\mathcal{I}_i = 1$  yields

$$\sum_{i \in [K]} \mathcal{I}_i \left( \mathbb{E}_{\mathcal{Q}_I} [T_i] - \mathbb{E}_{\mathcal{Q}_0} [T_i] \right) \leq TI d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_I).$$

Summing over  $\mathcal{I} \in \mathcal{A}$  yields

$$\sum_{\mathcal{I} \in \mathcal{A}} \sum_{i \in [K]} \mathcal{I}_i \left( \mathbb{E}_{\mathcal{Q}_I} [T_i] - \mathbb{E}_{\mathcal{Q}_0} [T_i] \right) \leq TI \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_I).$$

Thus

$$\begin{aligned} & \sum_{\mathcal{I} \in \mathcal{A}} \sum_{i \in [K]} \mathcal{I}_i \mathbb{E}_{\mathcal{Q}_I} [T_i] \\ & \leq \sum_{\mathcal{I} \in \mathcal{A}} \sum_{i \in [K]} \mathcal{I}_i \mathbb{E}_{\mathcal{Q}_0} [T_i] + TI \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_I) \\ & = \binom{K-1}{I-1} TI + TI \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_I). \end{aligned}$$

Therefore,

$$\begin{aligned}\mathbb{E}[\tilde{R}] &\geq \epsilon TI - \frac{\epsilon}{\binom{K}{I}} \left( \binom{K-1}{I-1} TI + TI \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \right) \\ &\quad + \lambda \mathbb{E}[\tilde{M}] \\ &= \epsilon TI \left( 1 - \frac{I}{K} \right) - \frac{\epsilon TI}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) + \lambda \mathbb{E}[\tilde{M}].\end{aligned}$$

*Proof of Lemma 5:* We first prove Lemma 5 for deterministic policies that make no more than  $S_0 = \frac{\epsilon TI}{\lambda}$  switches. For algorithms with this property, we have

$$\mathcal{Q}_0(\tilde{M} > \epsilon TI) = \mathcal{Q}_{\mathcal{I}}(\tilde{M} > \epsilon TI) = 0.$$

As the event  $\{\tilde{M} \geq m\}$  is in  $\tilde{\mathcal{F}}$  by Remark 1, then

$$\begin{aligned}\mathbb{E}_{\mathcal{Q}_0}[\tilde{M}] - \mathbb{E}_{\mathcal{Q}_{\mathcal{I}}}[\tilde{M}] &= \sum_{m=1}^{S_0} \left( \mathcal{Q}_0(\tilde{M} \geq m) - \mathcal{Q}_{\mathcal{I}}(\tilde{M} \geq m) \right) \\ &\leq \frac{\epsilon TI}{\lambda} d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}).\end{aligned}$$

Then

$$\begin{aligned}\mathbb{E}_{\mathcal{Q}_0}[\tilde{M}] - \mathbb{E}[\tilde{M}] &= \frac{1}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} \left( \mathbb{E}_{\mathcal{Q}_0}[\tilde{M}] - \mathbb{E}_{\mathcal{Q}_{\mathcal{I}}}[\tilde{M}] \right) \\ &\leq \frac{\epsilon TI}{\lambda \binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}).\end{aligned}\quad (23)$$

Combining (23) with Lemma 12 and Lemma 14, we obtain

$$\begin{aligned}\mathbb{E}[R] &\geq \epsilon TI \left( 1 - \frac{I}{K} - \frac{1}{4} \right) - \frac{2\epsilon TI}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \\ &\quad + \lambda \mathbb{E}_{\mathcal{Q}_0}[\tilde{M}].\end{aligned}$$

By Lemma 13, and  $\log_2 T + 1 \leq 2 \log_2 T$ , we have

$$\frac{1}{\binom{K}{I}} \sum_{\mathcal{I} \subset [K]} d_{\text{TV}}^{\tilde{\mathcal{F}}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \leq \frac{\epsilon}{\sigma \sqrt{K}} \sqrt{(\log_2 T) \mathbb{E}_{\mathcal{Q}_0}[\tilde{M}]}.$$

Using the notation  $m = \sqrt{\mathbb{E}_{\mathcal{Q}_0}[\tilde{M}]}$  and when  $K \geq 3I$ ,

$$\mathbb{E}[R] \geq \frac{5}{12} \epsilon TI - \frac{2\epsilon^2 TI}{\sigma \sqrt{K}} \sqrt{\log_2 T} m + \lambda m^2,$$

where the right hand side is minimized when  $m = \frac{\epsilon^2 TI \sqrt{\log_2 T}}{\lambda \sigma \sqrt{K}}$ . Thus the right-hand side is lower bounded by  $\frac{5\epsilon TI}{12} - \frac{\epsilon^4 T^2 I^2 \log_2 T}{\lambda \sigma^2 K}$ . Using our choice of  $\sigma = \frac{1}{6\sqrt{\log_2 T \log_2 \frac{4T(\lambda+\epsilon)}{\epsilon}}}$  and  $\epsilon = \frac{(\lambda K)^{\frac{1}{3}} (IT)^{-\frac{1}{3}}}{9 \log_2 T}$ , we derive

$$\mathbb{E}[R] \geq \frac{5(\lambda K)^{\frac{1}{3}} (TI)^{\frac{2}{3}}}{108 \log_2 T} - \frac{4(\lambda K)^{\frac{1}{3}} (TI)^{\frac{2}{3}} \log_2 \left( \frac{4T(\lambda+\epsilon)}{\epsilon} \right)}{9^3 (\log_2 T)^2}.\quad (24)$$

If  $\lambda < \epsilon$ , the right hand side of (24) is lower bounded by  $\frac{(\lambda K)^{\frac{1}{3}} (TI)^{\frac{2}{3}}}{30 \log_2 T}$  under the assumption that  $T \geq 8$ . Otherwise, if  $\epsilon \leq \lambda \leq T$ , we have

$$\epsilon = \frac{(\lambda K)^{\frac{1}{3}} (IT)^{-\frac{1}{3}}}{9 \log_2 T} \geq \frac{(3\lambda)^{1/3}}{9(\log_2 T) T^{1/3}} \geq \frac{(\lambda)^{1/3}}{7T^{4/3}},\quad (25)$$

where the first inequality is due to  $K \geq 3I$ . Combining  $\lambda \geq \epsilon$  and (25), we get  $\lambda \geq \frac{T^{-2}}{20}$  and then  $\epsilon \geq \frac{2T^{-2}}{39}$ . Thus when  $\epsilon \leq \lambda \leq T$  and  $T \geq 8$ , the right hand side of (24) is lower bounded by  $\frac{(\lambda K)^{\frac{1}{3}} (TI)^{\frac{2}{3}}}{130 \log_2 T}$ . Therefore, for any  $K \geq 3I$  and  $T \geq \max\{\frac{\lambda K}{I}, 8\}$ , it holds that

$$\mathbb{E}[R] \geq \frac{(\lambda K)^{\frac{1}{3}} (TI)^{\frac{2}{3}}}{130 \log_2 T}\quad (26)$$

For any general algorithm that has an arbitrary number of switches, we can turn it to a new algorithm that makes at most  $S_0$  switches by halting the algorithm once it makes  $S_0$  switches and repeating the last action in the remaining rounds. The regret  $R^*$  of the new algorithm (as defined in (22) under new algorithm) equals  $R$  when  $M \leq S_0$  and when  $M > S_0$ ,

$$R^* \leq R + \epsilon TI \leq 2R,$$

since  $R \geq \lambda S_0$ . Thus  $\mathbb{E}[R^*] \leq 2\mathbb{E}[R]$ . Since  $\mathbb{E}[R^*]$  is lower bounded by the right-hand side of (26), this implies the claimed lower bound on the expected regret of any deterministic player. ■

## B. Proof of Lemma 7

Given the constructed stochastic loss sequence  $L_{1:T}$  defined in Algorithm 2, we now want to analyze the player's expected regret under an arbitrary deterministic policy  $\pi$  which yields an action sequence  $A_{1:T} \in \mathcal{A}^T$  so that  $A_t$  is a function of the player's past observations  $Z_{1:t-1}$  with  $Z_t = A_t \circ L_t$ . Following the definition in (22) for  $R$ , we analyze the expected regret  $\mathbb{E}[R]$  in the new semi-bandit feedback setting and CDN loss sequence  $L_{1:T}$  given in Algorithm 2. Let  $\mathcal{Q}_0$  and  $\mathcal{Q}_{\mathcal{I}}$  follow previous definition in (9) and (10). Let  $\mathcal{F}$  be the  $\sigma$ -algebra generated by the observations  $Z_{1:T}$ , where  $Z_t = L_t \circ A_t$ . The total variation  $d_{\text{TV}}^{\mathcal{F}}$  is defined as in (11) with respect to the  $\sigma$ -algebra  $\mathcal{F}$ . For  $Y_{0:T}$  defined in (19) of §III-B and the action sequence  $A_{1:T}$ , we define

$$R' \triangleq \sum_{t=1}^T \langle A_t, \tilde{L}_t \rangle + \lambda \sum_{t=1}^T d(A_t, \tilde{A}_{t-1}) - \min_{A \in \mathcal{A}} \sum_{t=1}^T \langle A, \tilde{L}_t \rangle.$$

*Remark 2.* Note that  $A_t$  is a deterministic function of its past observations  $Z_{1:t-1}$ , thus the  $\sigma$ -algebra generated by  $A_{1:T}$  is a subset of  $\mathcal{F}$ .

**Lemma 15.** *It holds that*

$$\frac{1}{\binom{K}{I}} \sum_{\mathcal{I} \subset [K]} d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \leq \frac{\epsilon}{\sigma \sqrt{2K}} \sqrt{I(\log_2 T + 1) \mathbb{E}_{\mathcal{Q}_0}[M]}.$$

*Proof:* By Lemma 8, Pinsker's inequality [13, Lemma 11.6.1] and  $w(\rho) \leq \log_2 T + 1$ , we have

$$d_{\text{TV}}^{\mathcal{F}'}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \leq \frac{\epsilon}{2\sigma} \sqrt{(\log_2 T + 1) \mathbb{E}_{\mathcal{Q}_0}[M_{\mathcal{I}}]},$$

where  $\mathcal{F}'$  is the  $\sigma$ -algebra generated by  $Y_{0:T}$  (defined in (19) of §III-B). Since  $Z_{1:T}$  is a function of  $Y_{0:T}$ , we have  $\mathcal{F} \subset \mathcal{F}'$  which implies  $d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \leq \frac{\epsilon}{2\sigma} \sqrt{(\log_2 T + 1) \mathbb{E}_{\mathcal{Q}_0}[M_{\mathcal{I}}]}$ .

Then using the concavity of the squared root function and by  $\sum_{\mathcal{I} \subset [K]} M_{\mathcal{I}} = 2 \binom{K-1}{I-1} M$ , it holds that

$$\begin{aligned} & \frac{1}{\binom{K}{I}} \sum_{\mathcal{I} \subset [K]} d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \\ & \leq \frac{\epsilon}{2\sigma} \sqrt{\log_2 T + 1} \sum_{\mathcal{I} \subset [K]} \frac{1}{\binom{K}{I}} \sqrt{\mathbb{E}_{\mathcal{Q}_0}[M_{\mathcal{I}}]} \\ & \leq \frac{\epsilon}{2\sigma} \sqrt{\log_2 T + 1} \sqrt{\mathbb{E}_{\mathcal{Q}_0} \left[ \frac{1}{\binom{K}{I}} \sum_{\mathcal{I} \subset [K]} M_{\mathcal{I}} \right]} \\ & = \frac{\epsilon}{\sigma \sqrt{2K}} \sqrt{I(\log_2 T + 1) \mathbb{E}_{\mathcal{Q}_0}[M]}. \end{aligned}$$

**Lemma 16.** Assume that  $T \geq \max\{\frac{\lambda K}{T^2}, \frac{I}{6}\}$ . Then  $\mathbb{E}[R] \geq \mathbb{E}[R'] - \frac{\epsilon TI}{6}$ .

*Proof:* We consider the event  $B = \{\forall t: L_t = \tilde{L}_t\}$ , and first show that  $P(B) \geq 5/6$ . For  $\delta = \frac{I}{T} \leq \frac{1}{6}$ , by [5, Lemma 1] we have that with probability at least  $\frac{5}{6}$ , for all  $t \in [T]$  and  $i \in [K]$ ,

$$|W_t^i| \leq \sigma \sqrt{2d(\rho) \log_2 \frac{TI}{\delta}} \leq 3\sigma \log_2 T,$$

where the last inequality is due to  $d(\rho) \leq \log_2 T + 1$  by [5, Lemma 2]. Thus, setting  $\sigma = \frac{1}{9 \log T}$  we obtain that

$$P\left(\forall t \in [T], i \in [K], \frac{1}{2} + W_t^i \in \left[\frac{1}{6}, \frac{5}{6}\right]\right) \geq \frac{5}{6}.$$

For  $T \geq \max(\frac{\lambda K}{T^2}, 6)$ , we have  $\epsilon < \frac{1}{6}$  and thus  $\tilde{L}_t(x) \in [0, 1]$  for all  $x \in [K]$  whenever  $\frac{1}{2} + W_t^i \in [\frac{1}{6}, \frac{5}{6}]$ . This implies that  $P(B) \geq \frac{5}{6}$ .

If  $B$  occurs,  $R = R'$ ; otherwise,  $\lambda M \leq R \leq R' \leq \lambda M + \epsilon TI$ , so that  $R' - R \leq \epsilon TI$ . Therefore,

$$\mathbb{E}[R'] - \mathbb{E}[R] = \mathbb{E}[R' - R | \neg B] \cdot P(\neg B) \leq \frac{\epsilon TI}{6}.$$

**Lemma 17.** It holds that

$$\mathbb{E}[R'] \geq \epsilon TI \left(1 - \frac{I}{K}\right) - \frac{\epsilon TI}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) + \lambda \mathbb{E}[M]$$

*Proof:* For any  $i \in [K]$ , let  $T_i$  denote the number of times the player picks arm  $i$  when the time horizon is  $T$ . So we can write  $R' = \epsilon \left(TI - \sum_{i \in [K]} \chi_i T_i\right) + \lambda M$ , where  $\chi_i$  is the  $i$ -th component of  $\chi$ . We also use  $\mathcal{I}_i$  is the  $i$ -th component of  $\mathcal{I} \in \mathcal{A}$ . Thus

$$\begin{aligned} \mathbb{E}[R'] &= \frac{1}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} \mathbb{E} \left[ \epsilon \left(TI - \sum_{i \in [K]} \mathcal{I}_i T_i\right) + \lambda M \mid \chi = \mathcal{I} \right] \\ &= \epsilon TI - \frac{\epsilon}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} \sum_{i \in [K]} \mathcal{I}_i \mathbb{E}_{\mathcal{Q}_{\mathcal{I}}} [T_i] + \lambda \mathbb{E}[M]. \end{aligned}$$

For all  $i \in [K]$  and  $t \in [T]$ , the event  $\{A_{t,i} = 1\}$  belongs to the  $\sigma$ -field  $\mathcal{F}$  by Remark 2, so

$$\mathcal{Q}_{\mathcal{I}}(A_{t,i} = 1) - \mathcal{Q}_0(A_{t,i} = 1) \leq d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}).$$

Summing over  $t \in [T]$  yields

$$\mathbb{E}_{\mathcal{Q}_{\mathcal{I}}} [T_i] - \mathbb{E}_{\mathcal{Q}_0} [T_i] \leq T d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}).$$

Summing over  $i \in [K]$  such that  $\mathcal{I}_i = 1$  yields

$$\sum_{i \in [K]} \mathcal{I}_i \left( \mathbb{E}_{\mathcal{Q}_{\mathcal{I}}} [T_i] - \mathbb{E}_{\mathcal{Q}_0} [T_i] \right) \leq TI d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}).$$

Summing over  $\mathcal{I} \in \mathcal{A}$  yields

$$\sum_{\mathcal{I} \in \mathcal{A}} \sum_{i \in [K]} \mathcal{I}_i \left( \mathbb{E}_{\mathcal{Q}_{\mathcal{I}}} [T_i] - \mathbb{E}_{\mathcal{Q}_0} [T_i] \right) \leq TI \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}).$$

Thus

$$\begin{aligned} & \sum_{\mathcal{I} \in \mathcal{A}} \sum_{i \in [K]} \mathcal{I}_i \mathbb{E}_{\mathcal{Q}_{\mathcal{I}}} [T_i] \\ & \leq \sum_{\mathcal{I} \in \mathcal{A}} \sum_{i \in [K]} \mathcal{I}_i \mathbb{E}_{\mathcal{Q}_0} [T_i] + TI \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \\ & = \binom{K-1}{I-1} TI + TI \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbb{E}[R'] & \geq \epsilon TI - \frac{\epsilon}{\binom{K}{I}} \left( \binom{K-1}{I-1} TI + TI \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \right) \\ & \quad + \lambda \mathbb{E}[M] \\ & = \epsilon TI \left(1 - \frac{I}{K}\right) - \frac{\epsilon TI}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) + \lambda \mathbb{E}[M] \end{aligned}$$

*Proof of Lemma 7:* We first prove the theorem for deterministic players that make no more than  $S_0 = \epsilon TI / \lambda$  switches. For algorithms with this property, we have

$$\mathcal{Q}_0 \left( M > \frac{\epsilon TI}{\lambda} \right) = \mathcal{Q}_{\mathcal{I}} \left( M > \frac{\epsilon TI}{\lambda} \right) = 0.$$

By Remark 2, the event  $\{M \geq m\} \in \mathcal{F}$  which implies

$$\begin{aligned} \mathbb{E}_{\mathcal{Q}_0}[M] - \mathbb{E}_{\mathcal{Q}_{\mathcal{I}}}[M] &= \sum_{m=1}^{\epsilon TI / \lambda} (\mathcal{Q}_0(M \geq m) - \mathcal{Q}_{\mathcal{I}}(M \geq m)) \\ &\leq \frac{\epsilon TI}{\lambda} d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}). \end{aligned}$$

Then

$$\begin{aligned} \mathbb{E}_{\mathcal{Q}_0}[M] - \mathbb{E}[M] &= \frac{1}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} \left( \mathbb{E}_{\mathcal{Q}_0}[M] - \mathbb{E}_{\mathcal{Q}_{\mathcal{I}}}[M] \right) \\ &\leq \frac{\epsilon TI}{\lambda \binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \end{aligned}$$

Combining this with Lemma 16 and Lemma 17, we obtain

$$\begin{aligned} \mathbb{E}[R] & \geq \epsilon \left(1 - \frac{I}{K} - \frac{1}{6}\right) - \frac{2\epsilon TI}{\binom{K}{I}} \sum_{\mathcal{I} \in \mathcal{A}} d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \\ & \quad + \lambda \mathbb{E}_{\mathcal{Q}_0}[M]. \end{aligned}$$

By Corollary 15, and  $\log_2 T + 1 \leq 2 \log_2 T$ ,

$$\frac{1}{\binom{K}{I}} \sum_{\mathcal{I} \subset [K]} d_{\text{TV}}^{\mathcal{F}}(\mathcal{Q}_0, \mathcal{Q}_{\mathcal{I}}) \leq \frac{\epsilon}{\sigma \sqrt{K}} \sqrt{I(\log_2 T) \mathbb{E}_{\mathcal{Q}_0}[M]}.$$

Using the notation  $m = \sqrt{\mathbb{E}_{Q_0}[M]}$  and when  $K \geq 3I$ ,

$$\mathbb{E}[R] \geq \frac{\epsilon TI}{2} - \frac{2\epsilon^2 T I^{3/2}}{\sigma \sqrt{K}} \sqrt{\log_2 T} m + \lambda m^2,$$

where the right hand side is minimized when  $m = \frac{\epsilon^2 T I^{3/2} \sqrt{\log_2 T}}{\lambda \sigma \sqrt{K}}$ . Thus the right hand side is lower bounded by  $\frac{\epsilon TI}{2} - \frac{\epsilon^4 T^2 I^3 \log_2 T}{\lambda \sigma^2 K}$ . Using our choice of  $\sigma = \frac{1}{9 \log_2 T}$  and  $\epsilon = \frac{(\lambda K)^{\frac{1}{3}} I^{-\frac{2}{3}} T^{-\frac{1}{3}}}{9 \log_2 T}$ , gives

$$\mathbb{E}[R] \geq \frac{(\lambda KI)^{\frac{1}{3}} T^{\frac{2}{3}}}{30 \log_2 T}. \quad (27)$$

For any general algorithm that has an arbitrary number of switches, we can turn it to a new algorithm that makes at most  $S_0$  switches by halting the algorithm once it makes  $S_0$  switches and repeating the last action in the remaining rounds. The regret  $R^*$  of the new algorithm equals  $R$  when  $M \leq S_0$  and when  $M > S_0$ ,

$$R^* \leq R + \epsilon TI \leq 2R,$$

since  $R \geq \lambda S_0$ . Thus  $\mathbb{E}[R^*] \leq 2\mathbb{E}[R]$ . Since  $\mathbb{E}[R^*]$  is lower bounded by the right-hand side of (27), this implies the claimed lower bound on the expected regret of any deterministic player. ■

## REFERENCES

- [1] T. Lattimore and C. Szepesvári, *Bandit Algorithms*. Cambridge University Press, 2020.
- [2] S. Guha and K. Munagala, “Multi-armed bandits with metric switching costs,” in *Automata, Languages and Programming: 36th International Colloquium, ICALP 2009, Rhodes, Greece, July 5-12, 2009, Proceedings, Part II 36*. Springer Berlin Heidelberg, 2009, pp. 496–507.
- [3] M. Shi, X. Lin, and L. Jiao, “Power-of-2-arms for bandit learning with switching costs,” in *Proceedings of the Twenty-Third International Symposium on Theory, Algorithmic Foundations, and Protocol Design for Mobile Networks and Mobile Computing*, 2022, pp. 131–140.
- [4] R. Arora, O. Dekel, and A. Tewari, “Online bandit learning against an adaptive adversary: From regret to policy regret,” in *Proceedings of the 29th International Conference on Machine Learning*. PMLR, 2012, pp. 1747–1754.
- [5] O. Dekel, J. Ding, T. Koren, and Y. Peres, “Bandits with switching costs:  $T^{2/3}$  regret,” in *Proceedings of the 46th Annual ACM Symposium on Theory of Computing*, 2014, pp. 459–467.
- [6] C. Rouyer, Y. Seldin, and N. Cesa-Bianchi, “An algorithm for stochastic and adversarial bandits with switching costs,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 9127–9135.
- [7] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire, “The non-stochastic multiarmed bandit problem,” *SIAM Journal on Computing*, vol. 32, no. 1, pp. 48–77, 2002.
- [8] N. Cesa-Bianchi and G. Lugosi, *Prediction, Learning, and Games*. Cambridge university press, 2006.
- [9] S. Bubeck, “Introduction to online optimization,” *Lecture notes*, vol. 2, pp. 1–86, 2011.
- [10] C.-Y. Wei and H. Luo, “More adaptive algorithms for adversarial bandits,” in *Conference On Learning Theory*. PMLR, 2018, pp. 1263–1291.
- [11] R. Combes, M. Sadeh Talebi, A. Proutiere, and M. Lelarge, “Combinatorial bandits revisited,” *Advances in Neural Information Processing Systems*, vol. 28, 2015.
- [12] A. C.-C. Yao, “Probabilistic computations: Toward a unified measure of complexity,” in *18th Annual Symposium on Foundations of Computer Science (sfcs 1977)*. IEEE Computer Society, 1977, pp. 222–227.
- [13] T. M. Cover, *Elements of Information Theory*. John Wiley & Sons, 1999.
- [14] J.-Y. Audibert, S. Bubeck, and G. Lugosi, “Regret in online combinatorial optimization,” *Mathematics of Operations Research*, vol. 39, no. 1, pp. 31–45, 2014.

- [15] J. Zimmert, H. Luo, and C.-Y. Wei, “Beating stochastic and adversarial semi-bandits optimally and simultaneously,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 7683–7692.

**Yanyan Dong** received her B.S. degree from Jilin University in 2017, and Ph.D. degree from The Chinese University of Hong Kong, Shenzhen in 2022. She was a Research Fellow at the National University of Singapore from Dec. 2022 to Dec. 2023. Her research interests include information theory, coding theory, network coding, and machine learning.

**Vincent Y. F. Tan** (Senior Member, IEEE) was born in Singapore in 1981. He received the B.A. and M.Eng. degrees in electrical and information science from Cambridge University in 2005, and the Ph.D. degree in electrical engineering and computer science (EECS) from the Massachusetts Institute of Technology (MIT) in 2011. He is currently a Professor with the Department of Mathematics and the Department of Electrical and Computer Engineering (ECE), National University of Singapore (NUS). His research interests include information theory, machine learning, and statistical signal processing.

Dr. Tan is an elected member of the IEEE Information Theory Society Board of Governors. He was an IEEE Information Theory Society Distinguished Lecturer from 2018 to 2019. He received the MIT EECS Jin-Au Kong Outstanding Doctoral Thesis Prize in 2011, the NUS Young Investigator Award in 2014, the Singapore National Research Foundation (NRF) Fellowship (Class of 2018), and the NUS Young Researcher Award in 2019. He is currently serving as a Senior Area Editor for the *IEEE Transactions on Signal Processing* and as an Associate Editor in Machine Learning and Statistics for the *IEEE Transactions on Information Theory*. He also regularly serves as an Area Chair of prominent machine learning conferences such as the *International Conference on Learning Representations (ICLR)* and the *Conference on Neural Information Processing Systems (NeurIPS)*.