

Power-of-2-Arms for Bandit Learning With Switching Costs

Ming Shi
Purdue University
West Lafayette, IN, USA
sming@purdue.edu

Xiaojun Lin
Purdue University
West Lafayette, IN, USA
linx@ecn.purdue.edu

Lei Jiao
University of Oregon
Eugene, OR, USA
jiao@cs.uoregon.edu

ABSTRACT

Motivated by edge computing with artificial intelligence, in this paper we study a bandit-learning problem with switching costs. Existing results in the literature either incur $\Theta(T^{\frac{2}{3}})$ regret with bandit feedback, or rely on free full-feedback in order to reduce the regret to $O(\sqrt{T})$. In contrast, we expand our study to incorporate two new factors. First, full feedback could incur a cost. Second, the player may choose 2 (or more) arms at a time, in which case she is free to use any one of the chosen arms to calculate loss, and switching costs are incurred only when she changes the set of chosen arms. For the setting where the player pulls only one arm at a time, our new regret lower-bound shows that, even when costly full-feedback is added, the $\Theta(T^{\frac{2}{3}})$ regret still cannot be improved. However, the dependence on the number of arms may be improved when the full-feedback cost is small. In contrast, for the setting where the player can choose 2 (or more) arms at a time, we provide a novel online learning algorithm that achieves a lower $O(\sqrt{T})$ regret. Further, our new algorithm does not need any full feedback at all. This sharp difference therefore reveals the surprising power of choosing 2 (or more) arms for this type of bandit-learning problems with switching costs. Both our new algorithm and regret analysis involve several new ideas in choosing the primary and secondary arms, tuning the weight decay parameter within and across episodes, and using the loss differences in the weight updates, which may be of independent interest.

KEYWORDS

Bandit learning, switching costs, regret analysis, edge computing with artificial intelligence

1 INTRODUCTION

In this paper, we are interested in bandit learning with switching costs, which can be used to model many practical decision-making problems that not only face significant uncertainty, but also incur costs for changing decisions. Consider edge computing with artificial intelligence (Edge AI) [7] as an example, where an edge server close to the end users downloads machine learning (ML) models from the cloud to process incoming inference requests. As the underlying ground-truth model of the data changes in uncertain ways (which is often referred to as concept drift [13]), the best ML model also changes in time. However, because of the limited capability of the edge server, it can often only accommodate a small number of ML models. Thus, the edge server needs to learn which subset of ML models should be used, based on the feedback (e.g., inference losses) observed. Further, downloading an ML model (which is not currently on the edge server) from the cloud incurs communication overhead, which can be modelled by a switching cost β_1 . Hence, the edge server has to carefully select the ML models to

reduce the total inference losses and switching costs in the long run, which thus corresponds to a bandit-learning problem with switching costs. Other examples of such problems can be found in transportation networks [2], data-center networks [19], wireless communication [4] and cyber-physical systems [15], etc.

In the online learning literature, it is well-known that the existence of switching costs significantly changes the nature of the regret. Specifically, in an adversarial setting (which we will focus on in this paper), for bandit learning *without* switching costs, the Exp3 algorithm can attain $O(\sqrt{T})$ regret over a time-horizon T [3]. However, once the switching cost is added, the regret (for the setting where only one arm can be pulled at each time) increases substantially to $\Theta(T^{\frac{2}{3}})$ [1]. A matching lower bound in [10] suggests that such an increased regret is unavoidable. While this result may be somewhat discouraging, it leaves many important open questions, as we explained below. Note that since ML models in Edge AI corresponds to arms in bandit learning, we use the word “model” and “arm” interchangeably in the rest of the paper.

First, in practice, in addition to pulling one arm, there are often other ways to obtain feedback. For example, in Edge AI, the edge server could send the data to the cloud for analysis. In this case, the feedback from all ML models can be obtained, beyond the model already deployed on the edge server. This is somewhat analogous to the full-feedback setting studied in [14]. Reference [14] shows that, if the full feedback can be obtained with zero costs, the regret for bandit learning with switching costs will remain at $O(\sqrt{T})$, which would have been much lower than that of [1] where only bandit feedback is available. However, in practice, feedback from the cloud also incurs non-negative costs due to multiple reasons, e.g., communication costs, latency and privacy issues [7]. *Thus, the regret for bandit learning with both switching costs and full-feedback costs remains an open problem.*

Second, instead of holding only one ML model at each time, in Edge AI, the edge server can usually accommodate $M \geq 2$ ML models at each time. In this setting, using any of these M models for inference does not incur a switching/downloading cost, and at each time the feedback from all these M models (currently on the edge server) can be observed. This setting is thus most similar to a bandit-learning problem with limited advice [18], where $M \geq 2$ arms can be chosen at each time. However, [18] only studied the case without switching costs, where the regret is $O(\sqrt{T})$ regardless of whether one ($M = 1$) or more ($M \geq 2$) arms are chosen at each time. Our setting is also related to bandit-learning problems with semi-bandit feedback [8] and side information [2]. The studies for semi-bandit feedback [8] typically do not consider switching costs either. Although the side-information setting [2] has been studied with switching costs, it is somewhat different from ours because switching within the $M \geq 2$ arms also incur switching costs there. Partly due to this difference, the regret [2] remains at $\Theta(T^{\frac{2}{3}})$. In

summary, it remains an open problem whether in our setting, choosing $M \geq 2$ arms can improve the regret.

In this paper, we provide new answers to the above-mentioned two open problems. First, we study the case when $M = 1$, i.e., only one arm can be pulled at each time, and there is a switching cost β_1 to change the arm and a full-feedback cost β_2 to obtain feedback from all arms. As we discussed earlier, the latter action corresponds to the edge server sending data to the cloud for analysis. We provide a lower bound of the regret, which grows as $\Theta(T^{\frac{2}{3}})$. In other words, when only one arm can be pulled ($M = 1$), adding costly full-feedback will not fundamentally change how regret depends on T . However, our lower bound does suggest that utilizing costly full-feedback may change the multiplication factor in front of $T^{\frac{2}{3}}$. In some settings, this factor can be reduced from $O(K^{\frac{1}{3}})$ to $O((\ln K)^{\frac{1}{3}})$, where K is the total number of arms. This lower bound is obtained by constructing two new type of adversaries (please see Sec. 3.2) that forces any online learning algorithm to either switch arms or use costly full-feedback for at least $\Omega(T^{\frac{2}{3}})$ number of times, in order to obtain a loss no greater than the optimal static loss plus $O(T^{\frac{2}{3}})$. The proof of the lower bound involves an analysis of the Kullback-Leibler (KL) divergence (i.e., relative entropy) on a hidden Markov model, which is of independent interest. Moreover, we provide an algorithm called Randomized Online Learning With Costly Full-Feedback (ROCF) that achieves a regret that matches the lower bound (please see Sec. 3.3).

Second, we study the setting when $M \geq 2$, i.e., more than one arm can be chosen at each time and one of them is used to incur the loss, while there are still switching costs and full-feedback costs. Surprisingly, here we provide a new online learning algorithm, called Randomized Online Learning With Working Groups (ROW), that can achieve a regret of $O(\sqrt{T})$ without even using full feedback (see Theorem 4.1), which significantly improves the $\Theta(T^{\frac{2}{3}})$ regret for $M = 1$. In other words, having the flexibility to accommodate one additional model (i.e., $M = 2$) almost brings comparable benefit as having free full-feedback [14]. *To the best of our knowledge, this sharp transition from $M = 1$ to $M \geq 2$ has never been reported in the literature for bandit learning with switching costs*¹. This may be seen as somewhat analogous to the “power-of-2” routing in load balancing [16] (where sampling two queues can attain comparable reduction to delay as sampling all queues), which is why we refer to it as the “power-of-2-arms”. As M increases, the regret of ROW further decreases. Using a trivial lower bound for bandit learning with free full-feedback [5, 14], we conclude that the dependence of the regret of ROW on T must be optimal.

To achieve the improved $O(\sqrt{T})$ regret, ROW employs several new ideas. First, since $M \geq 2$ arms can be chosen at each time, in addition to choosing the “best” arm that has been observed so far (which we refer to as the primary arm), ROW also has the flexibility to choose $M - 1$ other arms (which we refer to as the secondary arms). We refer to the union of primary and secondary arms as the working group. In order to fully utilize such flexibility at minimal switching costs, the first idea of ROW is to fix a primary arm over $O(\sqrt{T})$ time-slots (which we refer to as an episode), and switch

the secondary arms $\lceil \frac{K-1}{M-1} \rceil$ times during an episode, each time to a new subset of secondary arms that have not yet been chosen in this episode. In other words, each episode is divided into $\lceil \frac{K-1}{M-1} \rceil$ number of sub-episodes, and switching only occurs at the end of each sub-episode. In this way, ROW only makes a constant number of switches within each episode (and $\Theta(\sqrt{T})$ switches for all the time), but it can obtain not only the feedback of the primary arm for the entire episode, but also the feedback of every other arms for $\frac{1}{\lceil \frac{K-1}{M-1} \rceil}$ fraction of the episode. Intuitively, this way of obtaining feedback incurs much lower costs than using costly full-feedback to obtain the same amount of feedback (for any K and $\beta_2 > 0$ independent of T), which is also the reason that ROW does not use costly full-feedback at all. Note that such saving is only possible when $M \geq 2$. As we have discussed earlier, for $M = 1$, either the switching cost or the full-feedback cost has to be $\Omega(T^{\frac{2}{3}})$ to attain low losses.

However, just using the above idea alone is insufficient to produce the $O(\sqrt{T})$ regret. The reason is that the feedback obtained is highly correlated in time. This is because each subset of secondary arms is retained for the whole sub-episode (whose length is also $O(\sqrt{T})$). It is known that such correlation tends to increase the regret. Indeed, we can construct two counter-examples (please see Sec. 4.1) to show that, if we merely use episodic versions of existing bandit-learning algorithms, e.g., Exp3 [3], the regret will still be very high. ROW utilizes a second crucial idea to overcome this difficulty. Our key observation is that, whenever such a sub-episode with highly-correlated feedback occurs, one of arms in the current working group (either the primary arm or a secondary arm) will likely be consistently better than other arms. Then, ROW will try to switch to the better arm more quickly within the sub-episode, and thus improve the regret. Specifically, recall that in Exp3 [3], each new feedback $\tilde{l}(t)$ reduces the weight of an arm by a factor $e^{-\eta \tilde{l}(t)}$. The parameter η thus determines how fast Exp3 responds to new feedback information, and it must be set to a particular value to achieve the minimum $O(\sqrt{T})$ regret (for bandit learning without switching costs). To accomplish this faster switching within a sub-episode, our proposed ROW algorithm will use a larger weight-decay parameter η_2 within each sub-episode, while using a smaller parameter η_1 across episodes. However, η_2 cannot be too large either. Otherwise, the regret will be poor for sub-episodes where the feedback is not correlated in time. In Sec. 4.2.2, we give a sufficient condition on how much η_2 should be larger than η_1 to strike the right balance. We note that this idea of using two different weight-decay parameters is new and may be of independent interest.

Finally, since in each episode the primary arm will receive much more feedback than the secondary arms, this creates a bias in the overall quality of feedback at the end of each episode. This bias issue is resolved by using instead the loss differences between the primary and secondary arms (please see our Idea 3 in Sec. 4.1). Our proof for the $O(\sqrt{T})$ regret carefully captures the effect of the above ideas by utilizing several new techniques (please see Sec. 4.2 for details).

¹Note that for bandit learning *without* switching costs, choosing $M \geq 2$ arms will improve the regret, but it cannot alter the dependence on T [18].

2 PROBLEM FORMULATION

In this section, we provide the problem formulation for our bandit-learning problem with switching costs and full-feedback costs. Moreover, we present a motivating example based on edge computing with artificial intelligence (Edge AI), which has received extensive attention recently [7, 23]. Finally, we introduce the performance metric.

2.1 Bandit Learning With Switching Costs and Full-Feedback Costs

A player interacts with the adversary/environment sequentially in time. Let $\mathcal{K} \triangleq \{1, 2, \dots, K\}$ denote the set of all arms and let M be an integer, $1 \leq M < K$. In each time-slot $t = 1, \dots, T$, first the player chooses M arms among all K arms. Let $\hat{\mathbb{K}}(t)$ denote the set of the M arms chosen at time t . The player uses one of the arms in $\hat{\mathbb{K}}(t)$ as the active arm, which is denoted by $k(t)$. The loss of this arm, $l_{k(t)}(t)$, will be used to calculate the loss and regret later. In addition, the losses $l_k(t)$ of all arms $k \in \hat{\mathbb{K}}(t)$ are observed by the player. The loss $l_k(t)$ can be any arbitrary value in $[0, 1]$. In this paper, we study both the cases when $M = 1$ and $2 \leq M < K$. When $M = 1$, $\hat{\mathbb{K}}(t)$ only contains the active arm $k(t)$ and only the loss of this arm is observed. In this case, we simply say that the player ‘‘pulls’’ the single arm $k(t)$ at time t . On the other hand, when $2 \leq M < K$, in addition to the loss of the active arm $k(t)$, the losses of other $M - 1$ arms in $\hat{\mathbb{K}}(t)$ are also observed.

Next, for every arm that is newly added to the set $\hat{\mathbb{K}}(t)$, a switching cost $\beta_1 > 0$ will be incurred. Thus, the switching cost at time t is $\beta_1 \sum_{k \in \hat{\mathbb{K}}(t)} \mathbf{1}_{\{k \notin \hat{\mathbb{K}}(t-1)\}}$, where $\mathbf{1}_E$ is an indicator function (i.e., $\mathbf{1}_E = 1$ if the event E is true, and $\mathbf{1}_E = 0$ otherwise). As typically assumed in bandit-learning problems [2, 3, 10, 14, 21], we assume that $\hat{\mathbb{K}}(0) = \Phi$ is empty. In addition to the feedback from the M arms in $\hat{\mathbb{K}}(t)$, at each time t , the player can choose to obtain full feedback of time t for all the arms (including those not in $\hat{\mathbb{K}}(t)$) at a cost β_2 . Let $z(t) = 1$ if the player chooses to obtain the full feedback at time t , and $z(t) = 0$ otherwise. Therefore, the total cost is

$$\text{Cost}(1 : T) \triangleq \sum_{t=1}^T \left\{ l_{k(t)}(t) + \beta_1 \sum_{k \in \hat{\mathbb{K}}(t)} \mathbf{1}_{\{k \notin \hat{\mathbb{K}}(t-1)\}} + \beta_2 z(t) \right\}. \quad (1)$$

2.2 An Example Motivated by Edge AI

We consider an Edge AI setting where an edge server collaborates with a remote cloud. The edge server runs machine learning (ML) models on an online stream of input data to predict their labels. (For example, in an E-commerce recommendation system, the input data at each time contains the customer data, item data and web shop transactions, etc. The input data will be used by the edge server to return the recommendations, i.e., the predicted labels of what the customer is interested in.) We assume that K ML models are already trained and available in the remote cloud. However, due to the limited capability of the edge server, only M models can be deployed at the edge server at each time. Since it is unknown which ML model works best, the edge server needs to use the feedback (e.g., the actual product picked by the customer) to learn which subset of ML models it should deploy. (In practice, both the

underlying distribution of the input data and the mapping from data to labels change in time due to the so-called concept drift [13]. Therefore, the best model(s) also changes in time. As a result, this learning process may be performed again after a concept drift.)

This learning process can be modelled as the bandit-learning problem described above. Each arm corresponds to one of the K ML models. At each time t , the edge server chooses the subset $\hat{\mathbb{K}}(t)$ of M models, which correspond to the M arms chosen in bandit learning. This subset $\hat{\mathbb{K}}(t)$ may be the same as the subset $\hat{\mathbb{K}}(t - 1)$ chosen at last time $t - 1$, or it may differ, in which case a switching cost $\beta_1 \sum_{k \in \hat{\mathbb{K}}(t)} \mathbf{1}_{\{k \notin \hat{\mathbb{K}}(t-1)\}}$ for downloading the ML models that are not currently on the edge server will be incurred. Note that this switching cost is assumed to be proportional to the number of ML models (which are not currently on the edge server) downloaded at time t . Then, the input data $\vec{X}(t)$ is revealed. The edge server will use the models in $\hat{\mathbb{K}}(t)$ to infer the label of $\vec{X}(t)$. Further, it will use the result $\vec{Y}'_{k(t)}(t)$ of one of the models $k(t) \in \hat{\mathbb{K}}(t)$, to return to the end user. This model $k(t)$ then corresponds to the active arm in bandit learning. Next, the true label $\vec{Y}(t)$ of $\vec{X}(t)$ is revealed. The edge server can then calculate the inference loss $l_k(t)$ for each ML model $k \in \hat{\mathbb{K}}(t)$, based on the difference between the inferred label $\vec{Y}'_k(t)$ and the true label $\vec{Y}(t)$, e.g., using the squared loss (i.e., $l_k(t) = \|\vec{Y}(t) - \vec{Y}'_k(t)\|^2$) [11].

At the end of time t , the edge server may also choose to consult the cloud for the quality of all ML models. In that case, it sends the data $\vec{X}(t)$ to the cloud. After the cloud processes this data with all ML models $k \in \mathcal{K}$, the edge server can retrieve the inference-loss $l_k(t)$ of all the ML models. Clearly, it incurs additional computation/communication overhead to obtain such feedback from the cloud, which we model by the full-feedback cost β_2 .

2.3 Performance Metric

We use the regret [1, 3, 10, 14, 20] as the performance metric. For an online learning algorithm π , its total cost $\text{Cost}^\pi(1 : T)$ is given by (1), which includes both switching costs and full-feedback costs. For the optimal static solution OPT, it knows the future losses in advance, and hence can choose only one arm/model throughout the time-horizon. The cost of OPT is then given by $\text{Cost}^{\text{OPT}}(1 : T) = \min_{k \in \mathcal{K}} \sum_{t=1}^T l_k(t) + \beta_1$, where there is only one switching cost β_1 at the beginning of the time-horizon, and there is no full-feedback cost. The regret of algorithm π is defined to be the worst-case difference between the expected total cost of algorithm π and the total cost of OPT, i.e.,

$$R^\pi(T) \triangleq \sup_{l_{1:K}(1:T)} \left\{ \mathbb{E}_\pi [\text{Cost}^\pi(1 : T)] - \text{Cost}^{\text{OPT}}(1 : T) \right\}, \quad (2)$$

where the expectation is taken over the possible randomness of the algorithm π , and $l_{1:K}(1 : T)$ denotes the losses $l_k(t)$ of all arms $k \in [1, K]$ for all time $t \in [1, T]$. Our goal is to design an online learning algorithm with a regret as low as possible.

3 THE CASE OF $M = 1$

In this section, we focus on the case when $M = 1$, i.e., the player (e.g., edge-server) can pull only one arm (e.g., model) at each time. We are interested in studying whether adding full feedback with a

cost β_2 can alter the regret of bandit learning with switching costs. Recall that in this case, the active arm $k(t)$ is the only arm in $\hat{\mathcal{K}}(t)$. As we mentioned in the introduction, when full feedback is free, it has been shown in [14] that using full feedback will improve the regret from $\Theta(T^{\frac{2}{3}})$ to $O(\sqrt{T})$. However, since in our model the full feedback incurs a cost, it is no longer clear whether the regret can still be improved. In this section, we first give a lower bound on the regret when the cost of full feedback is considered. Second, we provide an algorithm called Randomized Online Learning With Costly Full-Feedback (ROCF), which attains a regret that matches the lower bound. Our main conclusion is that adding costly full-feedback will not change the dependence of the regret on T , but may change the multiplication factor as a function of K .

3.1 A Lower Bound on the Regret ($M=1$)

We first present a lower bound on the regret of any online algorithm.

THEOREM 3.1. *Consider bandit learning with switching costs and full-feedback costs introduced in Sec. 2.1. When $M = 1$, the regret of any online algorithm π must be lower-bounded as follows,*

$$R^\pi(T) \geq \underline{R}^\pi(T) \triangleq \max \left\{ C_1 \beta_a^{\frac{1}{3}} (\log_2 K)^{\frac{1}{3}} T^{\frac{2}{3}}, C_2 \beta_b^{\frac{1}{3}} T^{\frac{2}{3}} \right\}, \quad (3)$$

where

$$\begin{aligned} \beta_a &= \min \left\{ \frac{3}{2} \beta_1, \beta_2 \right\}, \quad \beta_b = \min \left\{ \frac{3}{4} K \beta_1, \beta_2 \right\}, \\ C_1 &= \sqrt[3]{\frac{4}{9 \ln 2}} \cdot \frac{1}{72 (\log_2 T - \log_2 \log_2 K)}, \quad \text{and} \\ C_2 &= \sqrt[3]{\frac{4}{9 \ln 2}} \cdot \frac{1}{72 \log_2 T}. \end{aligned}$$

We can see from Theorem 3.1 that, even when the costly full-feedback is added, as long as $M = 1$, $\Theta(T^{\frac{2}{3}})$ is still the optimal regret for bandit learning with switching costs. This is in sharp contrast to the case of free full-feedback [14], where the regret can be improved to $O(\sqrt{T})$. While this result may be somewhat discouraging, the costly full-feedback does play some role in the multiplication factor in front of $T^{\frac{2}{3}}$, which depends on the relative magnitude of β_1 and β_2 . Intuitively, when the full-feedback cost β_2 is large, the online learning algorithm would rather switch to obtain feedback than using costly full-feedback. On the other hand, when β_2 is small, the online learning algorithm should avoid switching and obtain feedback from costly full-feedback. Thus, we expect that costly full-feedback will be more useful in the latter case than in the former case. The conclusion of Theorem 3.1 shows this difference precisely. Specifically, we can make the following observations.

(i) When $\beta_2 \geq \frac{3}{4} K \beta_1$, the lower bound $\underline{R}^\pi(T)$ in (3) is equal to

$$\max \left\{ C_1 \left(\frac{3}{2} \beta_1 \right)^{\frac{1}{3}} (\log_2 K)^{\frac{1}{3}} T^{\frac{2}{3}}, C_2 \left(\frac{3}{4} \beta_1 \right)^{\frac{1}{3}} K^{\frac{1}{3}} T^{\frac{2}{3}} \right\}. \quad (4)$$

As K increases, the second term in (4) quickly dominates. This means that, when the full-feedback cost β_2 is high, the regret of any online learning algorithm π will at least increase as $\beta_1^{\frac{1}{3}} K^{\frac{1}{3}} T^{\frac{2}{3}}$. Note that this expression is the same as the regret (for bandit learning with switching costs) when there is no full feedback at all [10]. This

Algorithm 1 The Multivariate Hidden Markov (MHM) adversary

Parameters: Choose ϵ and σ according to (12).

Initialization: Choose k^* uniformly from \mathcal{K} .

for $t = 1 : T$ **do**

Step 1: Generate the value of $G(t)$ according to (7).

Step 2: Generate the losses of each arm $k \in \mathcal{K}$ as follows,

$$l_k(t) = G(t) + \frac{1}{2} - \epsilon \cdot \mathbf{1}_{\{k=k^*\}} + \gamma_k(t), \quad (6)$$

where $\gamma_k(t) \sim \mathcal{N}(0, \sigma^2)$ are *i.i.d.* Gaussian random variables with zero-mean and σ^2 -variance.

end for

observation is consistent with our intuition that, when β_2 is large, the online algorithm cannot benefit from costly full-feedback.

(ii) When $\beta_2 < \frac{3}{4} K \beta_1$, the lower bound $\underline{R}^\pi(T)$ in (3) is equal to

$$\max \left\{ C_1 \beta_a^{\frac{1}{3}} (\log_2 K)^{\frac{1}{3}} T^{\frac{2}{3}}, C_2 \beta_2^{\frac{1}{3}} T^{\frac{2}{3}} \right\}. \quad (5)$$

As K increases, the first term in (5) quickly dominates. This means that, when the full-feedback cost β_2 is not high, the regret of any online algorithm π will at least increase as $\beta_a^{\frac{1}{3}} (\ln K)^{\frac{1}{3}} T^{\frac{2}{3}}$. If in addition $\beta_2 \leq \frac{3}{2} \beta_1$, we have $\beta_a^{\frac{1}{3}} (\ln K)^{\frac{1}{3}} T^{\frac{2}{3}} = \beta_2^{\frac{1}{3}} (\ln K)^{\frac{1}{3}} T^{\frac{2}{3}}$, which is smaller than $\beta_1^{\frac{1}{3}} K^{\frac{1}{3}} T^{\frac{2}{3}}$. Compared with the earlier case (with large β_2), our regret expression here has the same dependence on T , but now increases more slowly as a function of the total number K of arms. This observation is also consistent with our intuition that, when β_2 is small, the online algorithm can benefit from costly full-feedback more.

Finally, we note that the division of the two cases depends on the value of $K\beta_1$ and β_2 . The intuition is that, with K switches, an online algorithm may also attain the feedback from all K arms. Thus, it makes sense to compare $K\beta_1$ with β_2 to determine which type of feedback is more effective.

3.2 Lower Bound Analysis

To prove Theorem 3.1, we design two important adversaries, which are shown in Sec. 3.2.1 and Sec. 3.2.2. The first adversary captures the dependence of the regret on T . The second adversary uses the first adversary as a building block, which allows us to refine the dependence of the regret on K . For both adversaries, we make use of Yao's principle [22] that the worst-case expected regret $R^\pi(T)$ of a randomized online algorithm π is lower-bounded by the expected regret of the best deterministic online algorithm against a randomized adversary. Thus, in the following we focus on designing randomized adversaries, and studying the regret of deterministic online algorithms. Recall that $\mathcal{K} = \{1, \dots, K\}$.

3.2.1 Multivariate Hidden Markov (MHM) Adversary. In this section, we provide the first randomized adversary, called Multivariate Hidden Markov (MHM) adversary, which generalizes the idea in [10]. Please see Algorithm 1.

Specifically, Step 1 in Algorithm 1 is the same as that used by the adversary introduced in [10]. That is, for each time t , define the parent time of t as $\rho(t) \triangleq t - 2^{\delta(t)}$, where $\delta(t) \triangleq \max\{\delta \mid t \equiv 0$

(mod 2^δ). The main reason that the parent time $\rho(t)$ is $2^{\delta(t)}$ time-slot ahead of time t is to guarantee that with high probability, the generated losses $l_k(t)$ are in $[0, 1]$. Please see Appendix E for the concrete proof of this guarantee. Then, Step 1 of MHM generates a Gaussian process $G(t)$ in the following way,

$$G(t) = G(\rho(t)) + \xi(t), \text{ for all time } t \in [1, T], \quad (7)$$

where $G(0) = 0$, and $\xi(t) \sim \mathcal{N}(0, \sigma^2)$ are *i.i.d.* Gaussian random variables with zero-mean and σ^2 -variance. As in [10], this process $G(t)$ creates a common uncertainty across all arms. As a result, if an online algorithm does not switch arms, it will have a difficult time figuring out whether the losses experienced on the chosen arms are due to this common process $G(t)$, or due to the chosen arms being inferior to other arms. In Step 2, the first three terms² in (6) are also the same as that used in [10]. However, (6) differs from the adversary of [10] in the fourth term. This additional term adds a Gaussian noise $\gamma_k(t)$ to the loss $l_k(t)$ of each arm at each time. This additional noise is critical because our online algorithm π can use costly full-feedback, which is not considered in [10]. Intuitively, without this noise $\gamma_k(t)$, by using one round of costly full-feedback, the online algorithm can know the losses of all arms in the same time-slot. Then, the online algorithm will immediately know which arm is the optimal one (i.e., the arm with a loss that is ϵ lower). In contrast, the additional noise in (6) eliminates the possibility for such a trivial solution. We refer to this adversary as Multivariate Hidden Markov (MHM) because the hidden loss $l^{\text{hi}}(t) \triangleq l_{k(t)}(t) - \gamma_{k(t)}(t)$ satisfies the Markov property.

As we explain below, this additional noise $\gamma_k(t)$ causes new difficulties in the proof of the lower bound. We follow the approach in [10] to derive the regret lower-bound of any deterministic online algorithm π against the MHM adversary. Specifically, let $\mathcal{P}_{k^*}(\cdot)$ denote the probability measure under the setting where one optimal arm k^* incurs ϵ lower cost than other arms, as in (6). Let $\mathcal{P}_0(\cdot)$ denotes the probability measure when $\epsilon = 0$, i.e., the arm k^* is statistically the same as other arms. In addition, let $l^{\text{ob}}(\cdot)$ denote the observed losses of the online learning algorithm. Then, the analysis in [10] focuses on estimating the Kullback-Leibler (KL) divergence $D_{\text{KL}}(\mathcal{P}_{k^*}(l^{\text{ob}}(1:T)) \| \mathcal{P}_0(l^{\text{ob}}(1:T)))$, which then leads to the lower bound on the regret. However, for our MHM adversary, the additional noise $\gamma_k(t)$ incurs a new difficulty. Recall that $\rho(t)$ is the parent (time) of t , and thus t is the child (time) of $\rho(t)$. Let $\bar{\rho}(t)$ denote the set of the predecessors of time t , i.e. its parent, parent's parent, etc. Similarly, let $\underline{\rho}(t)$ denote the set of the descendants of time t . Note that without $\gamma_k(t)$, the observed loss $l^{\text{ob}}(t)$ would have been a Gaussian process $G(t)$ plus a fixed constant $\frac{1}{2}$ or $\frac{1}{2} - \epsilon$. Thus, $l^{\text{ob}}(t)$ would have satisfied a form of Markov property [12, p. 235], i.e., conditioned on current observed losses, the conditional probability distribution of future losses at a descendant time in $\underline{\rho}(t)$ is independent of past losses at any predecessor time in $\bar{\rho}(t)$. Then, the proof could use the chain rule of KL divergence [9, p. 23]. In contrast, with the additional noise $\gamma_k(t)$, the observed loss $l^{\text{ob}}(t)$ does not satisfy the Markov property any more. This is because, conditioned on the observed losses at time t , past observed losses still provide information for the statistics of the future losses. For

²The first three terms in (6) guarantees that the expected values of the losses are $\frac{1}{2}$ and $\frac{1}{2} - \epsilon$ for the sub-optimal arms $k \neq k^*$ and the optimal arm k^* , respectively.

example, by taking the average of the losses observed at all predecessors in $\bar{\rho}(t)$, we can average out $\gamma_k(t)$ across time, and thus estimate the mean value of the loss at a descendant time in $\underline{\rho}(t)$ with a higher accuracy. Therefore, we cannot use the chain rule directly, and must find a new way to bound the KL divergence.

To overcome this new difficulty, we develop a result on the KL divergence of hidden Markov models [9, p. 69]. Specifically, notice that the hidden loss $l^{\text{hi}}(t) \triangleq l_{k(t)}(t) - \gamma_{k(t)}(t)$, i.e., the loss in (6) but with $\gamma_{k(t)}(t)$ removed, satisfies the Markov property. Then, using the chain rule of probability, we can show that

$$\begin{aligned} D_{\text{KL}}(\mathcal{P}_{k^*}(l^{\text{ob}}(1:T), l^{\text{hi}}(1:T)) \| \mathcal{P}_0(l^{\text{ob}}(1:T), l^{\text{hi}}(1:T))) \\ = D_{\text{KL}}(\mathcal{P}_{k^*}(l^{\text{ob}}(1:T) | l^{\text{hi}}(1:T)) \| \mathcal{P}_0(l^{\text{ob}}(1:T) | l^{\text{hi}}(1:T))) \\ + D_{\text{KL}}(\mathcal{P}_{k^*}(l^{\text{hi}}(1:T)) \| \mathcal{P}_0(l^{\text{hi}}(1:T))), \end{aligned} \quad (8)$$

where the conditional KL divergence is defined to be [9, p. 22]

$$\begin{aligned} D_{\text{KL}}(\mathcal{P}_{k^*}(l^{\text{ob}}(1:T) | l^{\text{hi}}(1:T)) \| \mathcal{P}_0(l^{\text{ob}}(1:T) | l^{\text{hi}}(1:T))) \\ \triangleq \mathbb{E}_{\mathcal{P}_{k^*}(l^{\text{hi}}(1:T))} \left[\mathbb{E}_{\mathcal{P}_{k^*}(l^{\text{ob}}(1:T) | l^{\text{hi}}(1:T))} \left[\ln \frac{\mathcal{P}_{k^*}(l^{\text{ob}}(1:T) | l^{\text{hi}}(1:T))}{\mathcal{P}_0(l^{\text{ob}}(1:T) | l^{\text{hi}}(1:T))} \middle| l^{\text{hi}}(1:T) \right] \right]. \end{aligned} \quad (9)$$

Similarly, we can show that

$$\begin{aligned} D_{\text{KL}}(\mathcal{P}_{k^*}(l^{\text{ob}}(1:T), l^{\text{hi}}(1:T)) \| \mathcal{P}_0(l^{\text{ob}}(1:T), l^{\text{hi}}(1:T))) \\ = D_{\text{KL}}(\mathcal{P}_{k^*}(l^{\text{hi}}(1:T) | l^{\text{ob}}(1:T)) \| \mathcal{P}_0(l^{\text{hi}}(1:T) | l^{\text{ob}}(1:T))) \\ + D_{\text{KL}}(\mathcal{P}_{k^*}(l^{\text{ob}}(1:T)) \| \mathcal{P}_0(l^{\text{ob}}(1:T))) \\ \geq D_{\text{KL}}(\mathcal{P}_{k^*}(l^{\text{ob}}(1:T)) \| \mathcal{P}_0(l^{\text{ob}}(1:T))) \end{aligned} \quad (10)$$

where the inequality is because the KL divergence is always non-negative [9, p. 26]. Combining (8) and (10), we have that

$$\begin{aligned} D_{\text{KL}}(\mathcal{P}_{k^*}(l^{\text{ob}}(1:T)) \| \mathcal{P}_0(l^{\text{ob}}(1:T))) \\ \leq D_{\text{KL}}(\mathcal{P}_{k^*}(l^{\text{ob}}(1:T) | l^{\text{hi}}(1:T)) \| \mathcal{P}_0(l^{\text{ob}}(1:T) | l^{\text{hi}}(1:T))) \\ + D_{\text{KL}}(\mathcal{P}_{k^*}(l^{\text{hi}}(1:T)) \| \mathcal{P}_0(l^{\text{hi}}(1:T))). \end{aligned} \quad (11)$$

The first term on the right-hand-side of (11) can be easily calculated at each time, since conditioned on the hidden loss $l^{\text{hi}}(t)$, the observed loss $l^{\text{ob}}(t)$ is only due to *i.i.d.* Gaussian variables $\gamma_k(t)$. The second term on the right-hand-side of (11) can be calculated by using the chain rule of the KL divergence, since the hidden loss $l^{\text{hi}}(t)$ satisfies the Markov property. We can then obtain Lemma 3.2 below for the regret lower-bound against the MHM adversary.

LEMMA 3.2. *Consider bandit learning with switching costs and full-feedback costs introduced in Sec. 2.1. When $M = 1$, by choosing*

$$\epsilon = \sqrt[3]{\frac{4}{9 \ln 2}} \cdot \frac{1}{9 \log_2 T} \cdot \beta_b^{\frac{1}{3}} T^{-\frac{1}{3}} \text{ and } \sigma = \frac{1}{9 \log_2 T}, \quad (12)$$

Algorithm 2 The Dividing Set (DS) adversary

Initialization: A strictly positive integer $n \in \mathbb{Z}^{++}$, the total number of arms $K = 2^n$, and the set $\hat{\mathbb{k}}^*$ of optimal arms that begins with all arms, i.e., $\hat{\mathbb{k}}^*(0) = \mathcal{K}$.

for $j = 1 : n$ **do**

Step 1: Shrink the set of optimal arms randomly by half. Specifically, form the universe of the optimal arm set (with half size) as follows,

$$\Lambda(j) \triangleq \left\{ \hat{\mathbb{k}}^* \mid \hat{\mathbb{k}}^* \subseteq \hat{\mathbb{k}}^*(j-1), |\hat{\mathbb{k}}^*| = \frac{|\hat{\mathbb{k}}^*(j-1)|}{2} \right\}, \quad (14)$$

where $|\hat{\mathbb{k}}^*|$ denotes the cardinality of the set $\hat{\mathbb{k}}^*$. Then, choose $\hat{\mathbb{k}}^*(j)$ uniformly from $\Lambda(j)$.

Step 2: Restart and run a subroutine $\Psi(K, \hat{\mathbb{k}}^*(j), \frac{T}{n})$ to generate the losses in the j -th episode..

end for

the regret of any online learning algorithm π against the MHM adversary is lower-bounded as follows: for $T \geq \max\{\beta_b, 6K\}$,

$$R^\pi(T) \geq \sqrt[3]{\frac{4}{9 \ln 2}} \cdot \frac{1}{72 \log_2 T} \cdot \beta_b^{\frac{1}{3}} T^{\frac{2}{3}}, \quad (13)$$

where $\beta_b = \min\{\frac{3}{4}K\beta_1, \beta_2\}$.

Please see Appendix A for the complete proof of Lemma 3.2. From Lemma 3.2, we can see that the regret lower-bound produced by MHM corresponds to the second term in (3). Note that it correctly captures the dependence of the regret on T , but the dependence on K still needs to be refined.

3.2.2 The Dividing Set (DS) Adversary. To further refine the dependence on K , in this section we provide the second randomized adversary, called Dividing Set (DS). Please see Algorithm 2. We note that this “dividing set” idea could be used to produce sharper regret lower-bounds (in terms of their dependence on K) for other bandit-learning problems, and thus may also be of independent interest.

Specifically, the DS adversary starts with $K = 2^n$ arms, where n is a strictly positive integer. DS initializes the set of optimal arms to be $\hat{\mathbb{k}}^*(0) = \mathcal{K}$. Next, DS divides the entire time-horizon into n episodes, each with $\frac{T}{n}$ time-slots³. At the beginning of the j -th ($j = 1, \dots, n$) episode, the DS adversary uniformly chooses half of the arms from the last optimal-arm set $\hat{\mathbb{k}}^*(j-1)$ to form the new optimal-arm set $\hat{\mathbb{k}}^*(j)$ (i.e., Step 1 in Algorithm 2). Then, DS restarts and runs a subroutine $\Psi(K, \hat{\mathbb{k}}^*(j), \frac{T}{n})$ (i.e., Step 2 in Algorithm 2) with $K = 2$ (treating the two halves of the arms as $K = 2$ arms). This subroutine generates adversarial losses for $K = 2$ arms over $\frac{T}{n}$ time-slots, and it uses the arms in $\hat{\mathbb{k}}^*(j)$ as the optimal arms. In particular, we will use the MHM adversary for this sub-routine later. Suppose that the subroutine adversary $\Psi(K, \hat{\mathbb{k}}^*(j), \frac{T}{n})$ incurs a regret lower-bound of $O(T^\zeta)$ for each episode. Lemma 3.3 below provides a regret lower-bound for the entire time.

³For ease of exposition, we assume that T is divisible by n . All results can trivially be extended to the case when T is not divisible by n .

LEMMA 3.3. Suppose that a subroutine $\Psi(K, \hat{\mathbb{k}}^*, T)$ can generate adversarial losses over T time-slots for a bandit problem with K arms and the set of optimal arms given by $\hat{\mathbb{k}}^*$. Further, suppose that for any online algorithm π , $\Psi(K, \hat{\mathbb{k}}^*, T)$ can produce a regret lower-bound of $f(\log_2 T) \cdot T^\zeta$ ($0 < \zeta < 1$). Then, the DS adversary guarantees that the regret of any online algorithm π is lower-bounded as follows,

$$R^\pi(T) \geq f\left(\log_2 \frac{T}{\log_2 K}\right) (\log_2 K)^{1-\zeta} T^\zeta. \quad (15)$$

Please see Appendix F for the complete proof of Lemma 3.3. As we discussed earlier, for our setting we will simply use the MHM adversary as the subroutine for DS. In particular, MHM places all arms into two categories. All arms in the optimal-arm set $\hat{\mathbb{k}}^*$ are viewed as a single optimal arm, i.e., $l_k(t) = G(t) + \frac{1}{2} - \epsilon + \gamma_1(t)$ for all $k \in \hat{\mathbb{k}}^*$, and all arms outside the optimal-arm set $\hat{\mathbb{k}}^*$ are viewed as a single sub-optimal arm, i.e., $l_k(t) = G(t) + \frac{1}{2} + \gamma_2(t)$ for all $k \notin \hat{\mathbb{k}}^*$. In other words, effectively there are only two arms, $K = 2$, for MHM in each episode of DS. In this case, β_b in (13) is equal to $\beta_a \triangleq \min\{\frac{3}{2}\beta_1, \beta_2\}$. Thus, combining Lemma 3.2 and Lemma 3.3, we then get the regret lower-bound as

$$R^\pi(T) \geq \sqrt[3]{\frac{4}{9 \ln 2}} \cdot \frac{1}{72 \log_2 \left(\frac{T}{\log_2 K}\right)} \cdot \beta_a^{\frac{1}{3}} (\log_2 K)^{\frac{1}{3}} T^{\frac{2}{3}}, \quad (16)$$

The result of Theorem 3.1 then follows by combining (13) and (16).

3.3 Randomized Online Learning With Costly Full-Feedback (ROCF)

According to our discussion in Sec. 3.1, to match the lower bound $\underline{R}^\pi(T)$ in Theorem 3.1, we should achieve a regret $R^\pi(T)$, such that

$$R^\pi(T) \leq \begin{cases} a_1 \beta_1^{\frac{1}{3}} K^{\frac{1}{3}} T^{\frac{2}{3}}, & \text{if } \beta_2 \geq \frac{3}{4}K\beta_1, \\ a_2 \beta_2^{\frac{1}{3}} (\ln K)^{\frac{1}{3}} T^{\frac{2}{3}}, & \text{if } \beta_2 < \frac{3}{4}K\beta_1, \end{cases} \quad (17)$$

where a_1 and a_2 are positive constants. In this section, we provide an algorithm called Randomized Online Learning With Costly Full-Feedback (ROCF). Please see Algorithm 3. ROCF divides time into $\lceil \frac{T}{\tau} \rceil$ episodes of length τ . One arm is chosen at the beginning of an episode, and is kept throughout the episode. Across episodes, depending on the values of β_1 and β_2 , ROCF either apply the decision of Exp3 [3] without asking for costly full-feedback, or use the costly full-feedback in one random time-slot within each episode and apply the decision of the “shrinking dartboard” algorithm of [14].

Specifically, (i) when the full-feedback cost β_2 is large, i.e., $\beta_2 \geq \frac{3}{4}K\beta_1$, ROCF applies the Exp3 algorithm [3] across episodes. More specifically, at the end of the last time-slot of the u -th episode, ROCF computes the losses as follows,

$$\tilde{L}_k^{\text{ROCF}}[u] = \begin{cases} \frac{L_k[u]}{\hat{p}_k^{\text{ROCF}}[u]}, & \text{if } k = k^{\text{ROCF}}[u], \\ 0, & \text{if } k \neq k^{\text{ROCF}}[u], \end{cases} \quad (18)$$

where $L_k[u] \triangleq \sum_{t=t_u}^{t_u+\tau-1} l_k(t)$, and $k^{\text{ROCF}}[u]$ is the active arm used in the u -th episode. Next, using the computed losses, ROCF updates

Algorithm 3 Randomized Online Learning With Costly Full-Feedback (ROCF)

Parameters: Choose η and τ according to (22).

Initialization: $w_k^{\text{ROCF}}[1] = 1$ and $p_k^{\text{ROCF}}[1] = \frac{1}{K}$, for all $k \in \mathcal{K}$.

for $u = 1 : \lceil \frac{T}{\tau} \rceil$ (The u -th episode starts from $t_u = (u-1)\tau + 1$ to $t_u + \tau - 1$.) **do**

Step 1: At the beginning of the first time-slot t_u , pick an arm for the entire episode as follows.

if $u == 1$ (i.e., the first episode) **then**

Pick an arm $k^{\text{ROCF}}[1]$ from all arms $k \in \mathcal{K}$ according to the probability $p_k^{\text{ROCF}}[1]$.

else

if $\beta_2 \geq \frac{3}{4}K\beta_1$ **then**

Pick an arm $k^{\text{ROCF}}[u]$ from all arms $k \in \mathcal{K}$ according to the probability $p_k^{\text{ROCF}}[u]$.

else

With probability $p^{\text{ns}}[u] = \frac{w_{k^{\text{ROCF}}[u-1]}^{\text{ROCF}}[u]}{w_{k^{\text{ROCF}}[u-1]}^{\text{ROCF}}[u-1]}$, keep the previous arm, i.e., $k^{\text{ROCF}}[u] = k^{\text{ROCF}}[u-1]$. With probability $1 - p^{\text{ns}}[u]$, pick an arm $k^{\text{ROCF}}[u]$ from all arms $k \in \mathcal{K}$ according to the probability $p_k^{\text{ROCF}}[u]$.

end if

end if

Step 2: Uniformly choose a time $\tilde{t}[u]$ from $[t_u, t_u + \tau - 1]$.

Step 3: (Inside each episode.)

for $t = t_u : t_u + \tau - 1$ **do**

Pull the arm $k^{\text{ROCF}}[u]$ and use it as the active arm.

if $\beta_2 < \frac{3}{4}K\beta_1$ and $t == \tilde{t}[u]$ **then**

Ask for full feedback.

end if

end for

Step 4: At the end of the last time-slot of the u -th episode, compute the losses for all arms $k \in \mathcal{K}$ according to (18) (if $\beta_2 \geq \frac{3}{4}K\beta_1$) or (21) (if $\beta_2 < \frac{3}{4}K\beta_1$). Then, update the weights $w_k^{\text{ROCF}}[u+1]$ and probabilities $p_k^{\text{ROCF}}[u+1]$ according to (19) and (20), respectively.

end for

the weights and probabilities for all arms $k \in \mathcal{K}$ as follows,

$$w_k^{\text{ROCF}}[u+1] = w_k^{\text{ROCF}}[u] \cdot e^{-\eta \tilde{L}_k^{\text{ROCF}}[u]}, \quad (19)$$

$$p_k^{\text{ROCF}}[u+1] = \frac{w_k^{\text{ROCF}}[u+1]}{\sum_{k=1}^K w_k^{\text{ROCF}}[u+1]}, \quad (20)$$

where η is a tunable parameter (i.e., Step 4 in Algorithm 3). Then, at the beginning of the first time-slot of the $(u+1)$ -th episode, according to the updated probabilities $p_k^{\text{ROCF}}[u+1]$, ROCF picks an arm $k^{\text{ROCF}}[u+1]$ from all arms $k \in \mathcal{K}$ (i.e., Step 1 in Algorithm 3).

(ii) When the full-feedback cost β_2 is small, i.e., $\beta_2 < \frac{3}{4}K\beta_1$, ROCF asks for the costly full-feedback in one random time-slot within each episode and applies the decision of the “shrinking dartboard” algorithm of [14]. More specifically, in the u -th episode, ROCF asks for full feedback at time $\tilde{t}[u]$, which is uniformly chosen in the episode (i.e., Step 2 and Step 3 in Algorithm 3). Then, at the end of the last time-slot of the u -th episode, ROCF computes the

losses using only the full feedback, i.e.,

$$\tilde{L}_k^{\text{ROCF}}[u] = l_k(\tilde{t}[u]), \text{ for all arm } k \in \mathcal{K}. \quad (21)$$

Next, ROCF updates the weights and probabilities according to (19) and (20), while using (21) for the loss $\tilde{L}_k^{\text{ROCF}}[u]$ (i.e., Step 4 in Algorithm 3). Further, at the beginning of the first time-slot of the $(u+1)$ -th episode, with probability $p^{\text{ns}}[u+1] = \frac{w_{k^{\text{ROCF}}[u]}^{\text{ROCF}}[u+1]}{w_{k^{\text{ROCF}}[u]}^{\text{ROCF}}[u]}$,

ROCF keeps the arm $k^{\text{ROCF}}[u]$, i.e., $k^{\text{ROCF}}[u+1] = k^{\text{ROCF}}[u]$. With probability $1 - p^{\text{ns}}[u+1]$, ROCF picks an arm $k^{\text{ROCF}}[u+1]$ from all arms $k \in \mathcal{K}$ according to the updated probabilities $p_k^{\text{ROCF}}[u+1]$.

For both cases, ROCF keeps using the arm $k^{\text{ROCF}}[u+1]$ as the active arm for all time-slots in the $(u+1)$ -th episode (i.e., Step 3 in Algorithm 3).

3.4 Regret Analysis

In Theorem 3.4 below, we show the upper bound of the regret attained by ROCF.

THEOREM 3.4. Consider bandit learning with switching costs and costly full-feedback introduced in Sec. 2.1. Choose

$$\eta = \left(\min \left\{ \frac{3}{4}K\beta_1, \beta_2 \right\} \right)^{-\frac{1}{3}} (\ln K)^{\frac{2}{3}} T^{-\frac{2}{3}}, \text{ and} \quad (22a)$$

$$\tau = \left[(f(\beta_1, \beta_2))^{\frac{1}{3}} (\ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} \right], \quad (22b)$$

where $f(\beta_1, \beta_2) = \frac{9\beta_2^2}{16K}$, if $\beta_2 \geq \frac{3}{4}K\beta_1$, and $f(\beta_1, \beta_2) = \beta_2^2$, otherwise. For $T \geq \max \left\{ \frac{128K \ln K}{9\beta_1^2}, \frac{8 \ln K}{\beta_2^2} \right\}$, the regret of ROCF is upper-bounded as follows,

$$R^{\text{ROCF}}(T) \leq \begin{cases} b_1 \beta_1^{\frac{1}{3}} (K \ln K)^{\frac{1}{3}} T^{\frac{2}{3}}, & \text{if } \beta_2 \geq \frac{3}{4}K\beta_1, \\ \frac{7}{2} \beta_2^{\frac{1}{3}} (\ln K)^{\frac{1}{3}} T^{\frac{2}{3}} + b_2, & \text{if } \beta_2 < \frac{3}{4}K\beta_1, \end{cases} \quad (23)$$

where $b_1 = \left(\frac{3}{2} \sqrt[3]{\frac{3}{4}} + 2 \sqrt[3]{\frac{16}{9}} \right)$ and $b_2 = \beta_1 \ln K (1 + 2/\beta_2)$.

By comparing (17) and Theorem 3.4, we can see that the regret of ROCF matches the lower bound $R^{\pi}(T)$ up to a $(\ln K)^{\frac{1}{3}}$ factor. In particular, when the full-feedback cost β_2 is small, i.e., $\beta_2 < \frac{3}{4}K\beta_1$, the regret is indeed improved from $O(K^{\frac{1}{3}} T^{\frac{2}{3}})$ to $O((\ln K)^{\frac{1}{3}} T^{\frac{2}{3}})$. Please see Appendix G for the complete proof of Theorem 3.4.

4 THE POWER-OF-2-ARMS ($M \geq 2$)

In this section, we proceed to the case when $M \geq 2$. In contrast to the previous section where we show that adding costly full-feedback does not change the $\Theta(T^{\frac{2}{3}})$ regret, here we provide a new algorithm that utilizes the flexibility of having 2 (or more) arms and successfully improves the regret to $O(\sqrt{T})$.

4.1 Randomized Online Learning With Working Groups (ROW)

We call our new algorithm Randomized Online Learning With Working Groups (ROW). Please see Algorithm 4. We start with describing the high-level skeleton of ROW. Recall that $\mathcal{K} = \{1, \dots, K\}$.

Algorithm 4 Randomized Online Learning With Working Groups (ROW)

Parameters: Choose η_2, τ_2, η_1 and τ_1 according to (46).

Initialization: $w_k^{\text{ROW}}[1] = 1$ and $p_k^{\text{ROW}}[1] = \frac{1}{K}$, for all $k \in \mathcal{K}$.

for $u = 1 : \lceil \frac{T}{\tau_1} \rceil$ (The u -th episode starts from $t_u = (u - 1)\tau_1 + 1$ to $t_u + \tau_1 - 1$.) **do**

Step 1: At the beginning of the first time-slot t_u , according to probability $p_k^{\text{ROW}}[u]$ calculated in (24), choose a primary arm $k_0^{\text{ROW}}[u]$ from all arms $k \in \mathcal{K}$ for the entire episode.

for $v = 1 : \frac{\tau_1}{\tau_2}$ (The v -th sub-episode starts from $t_{u,v} = (u - 1)\tau_1 + (v - 1)\tau_2 + 1$ to $t_{u,v} + \tau_2 - 1$.) **do**

Step 2: At the beginning of the first time-slot $t_{u,v}$, uniformly choose the set $\hat{\mathcal{K}}_{M-1}^{\text{ROW}}[u, v]$ of $M - 1$ secondary arms from the not-yet-been-chosen arms in $\mathcal{K} - \left(\bigcup_{v'=1}^{v-1} \hat{\mathcal{K}}_{M-1}^{\text{ROW}}[u, v'] \cup \{k_0^{\text{ROW}}[u]\} \right)$. Then, form the working group by the primary arm and secondary arms, i.e., $\hat{\mathcal{K}}^{\text{ROW}}[u, v] = \{k_0^{\text{ROW}}[u]\} \cup \hat{\mathcal{K}}_{M-1}^{\text{ROW}}[u, v]$.

Step 3: Initialize the weights $\hat{w}_k^{\text{ROW}}(t_{u,v})$ and probabilities $\hat{p}_k^{\text{ROW}}(t_{u,v})$ of all arms $k \in \hat{\mathcal{K}}^{\text{ROW}}[u, v]$ according to (25) and (26), respectively.

for $t = t_{u,v} : t_{u,v} + \tau_2 - 1$ **do**

Step 4: Use an arm $k \in \hat{\mathcal{K}}^{\text{ROW}}[u, v]$ as the active arm according to the updated probability $\hat{p}_k^{\text{ROW}}(t)$.

Step 5: Update the weights $\hat{w}_k^{\text{ROW}}(t)$ and probabilities $\hat{p}_k^{\text{ROW}}(t)$ of all arms $k \in \hat{\mathcal{K}}^{\text{ROW}}[u, v]$ according to (27) and (26), respectively.

end for

end for

Step 6: At the end of the last time-slot of the u -th episode, update the weights $w_k^{\text{ROW}}[u+1]$ and probabilities $p_k^{\text{ROW}}[u+1]$ of all arms $k \in \mathcal{K}$ according to (29) and (24), respectively.

end for

Idea 1: Note that in order to obtain the $O(\sqrt{T})$ regret, we can switch or use costly full-feedback at most $O(\sqrt{T})$ number of times. The first idea of ROW is thus to design an effective way to rotate a working group (of M arms) through all K arms, so that plenty of feedback can be obtained for all the arms, while incurring $O(\sqrt{T})$ switching costs and zero full-feedback costs. Specifically, ROW divides the entire time-horizon into $U = \lceil \frac{T}{\tau_1} \rceil$ episodes, each with $\tau_1 = \Theta(\sqrt{T})$ time-slots. In the first time-slot $t_u = (u - 1)\tau_1 + 1$ of the u -th ($u = 1, \dots, U$) episode, ROW chooses a primary arm $k_0^{\text{ROW}}[u]$. This primary arm $k_0^{\text{ROW}}[u]$ will be fixed for all τ_1 time-slots in the u -th episode. In addition, ROW divides each episode into $V = \lceil \frac{K-1}{M-1} \rceil$ sub-episodes, each with $\tau_2 = \frac{\tau_1}{V}$ time-slots. In the rest of this paper, we refer to the v -th sub-episode in the u -th episode as sub-episode (u, v) . At the beginning of the first time-slot $t_{u,v} = (u - 1)\tau_1 + (v - 1)\tau_2 + 1$ of sub-episode (u, v) , ROW uniformly chooses $M - 1$ secondary arms from the arms that have not yet been chosen in the u -th episode⁴ (i.e., Step 2 in Algorithm 4). We let

⁴When $K - 1$ is not divisible by $M - 1$, the number of the remaining unchosen arms in the last (i.e., V -th) sub-episode may be less than $M - 1$. In this case, after choosing

$\hat{\mathcal{K}}_{M-1}^{\text{ROW}}[u, v]$ denote the set of the $M - 1$ secondary arms chosen in sub-episode (u, v) . Let $\hat{\mathcal{K}}^{\text{ROW}}[u, v] = \{k_0^{\text{ROW}}[u]\} \cup \hat{\mathcal{K}}_{M-1}^{\text{ROW}}[u, v]$ denote the working group formed by the primary arm and secondary arms. The working group $\hat{\mathcal{K}}^{\text{ROW}}[u, v]$ will be fixed for the whole sub-episode (u, v) .

Notice that by using this idea, ROW only switches at the boundaries of sub-episodes and never uses full feedback. Therefore, by tuning τ_2 to be $\Theta(\sqrt{T})$, the total switching cost is guaranteed to be $\Theta(\sqrt{T})$, and the total full-feedback cost is 0. More importantly, with this idea, we not only have the feedback for the primary arm for the entire episode, but also have the feedback for each secondary arm for $\frac{1}{V}$ fraction of each episode. Intuitively, this way of obtaining feedback incurs much lower costs than using costly full-feedback. For example, if we want to obtain the same amount of feedback by using costly full-feedback alone, we would have to incur a full-feedback cost equal to $\Theta(\sqrt{T})$ in every episode! This is also the reason that ROW does not use full feedback at all.

We now describe the rest of the details of ROW. At the beginning of the first time-slot of the u -th ($u = 1, \dots, U$) episode, each arm $k \in \mathcal{K}$ is associated with a weight $w_k^{\text{ROW}}[u]$, which is initialized to be $w_k^{\text{ROW}}[1] = 1$ (we will describe how to update $w_k^{\text{ROW}}[u]$ from $w_k^{\text{ROW}}[u - 1]$ shortly). Then, from all arms $k \in \mathcal{K}$, ROW chooses a primary arm $k_0^{\text{ROW}}[u]$ with probability (i.e., Step 1 in Algorithm 4)

$$p_k^{\text{ROW}}[u] = \frac{w_k^{\text{ROW}}[u]}{\sum_{k=1}^K w_k^{\text{ROW}}[u]}. \quad (24)$$

Then, at the beginning of the first time-slot of each sub-episode (u, v) , the $M - 1$ secondary arms $\hat{\mathcal{K}}_{M-1}^{\text{ROW}}[u, v]$ are chosen uniformly and rotated through all of the rest $K - 1$ arms as we described earlier (i.e., Step 2 in Algorithm 4).

Further, within each sub-episode (u, v) we solve a bandit-learning problem with the set of arms restricted to the chosen working group. Note that this restricted version of the bandit-learning problem has no switching cost (since any arm $k \in \hat{\mathcal{K}}^{\text{ROW}}[u, v]$ can be used as the active arm without incurring switching costs), and also has full feedback (from all the arms $k \in \hat{\mathcal{K}}^{\text{ROW}}[u, v]$). Thus, we can directly use the full-feedback version of the Exp3 algorithm inside each sub-episode (u, v) . Specifically, in the first time-slot $t_{u,v}$ of sub-episode (u, v) , ROW initializes the weights of all the arms $k \in \mathcal{K}$ as follows (i.e., Step 3 in Algorithm 4),

$$\hat{w}_k^{\text{ROW}}(t_{u,v}) = w_k^{\text{ROW}}[u], \quad (25)$$

i.e., to be the values of the weights at the beginning of the entire episode u . Then, for each time $t = t_{u,v}, \dots, t_{u,v} + \tau_2 - 1$, each arm $k \in \hat{\mathcal{K}}^{\text{ROW}}[u, v]$ is used as the active arm $k^{\text{ROW}}(t)$ with probability (i.e., Step 4 and Step 5 in Algorithm 4)

$$\hat{p}_k^{\text{ROW}}(t) = \frac{\hat{w}_k^{\text{ROW}}(t)}{\sum_{k \in \hat{\mathcal{K}}^{\text{ROW}}[u, v]} \hat{w}_k^{\text{ROW}}(t)}. \quad (26)$$

After the losses $l_k(t)$ of all the arms $k \in \hat{\mathcal{K}}^{\text{ROW}}[u, v]$ are obtained for time t , ROW updates their weights with a tunable parameter

all those unchosen arms, ROW uniformly chooses the secondary arms from the arms that have not yet been chosen for the V -th sub-episode.

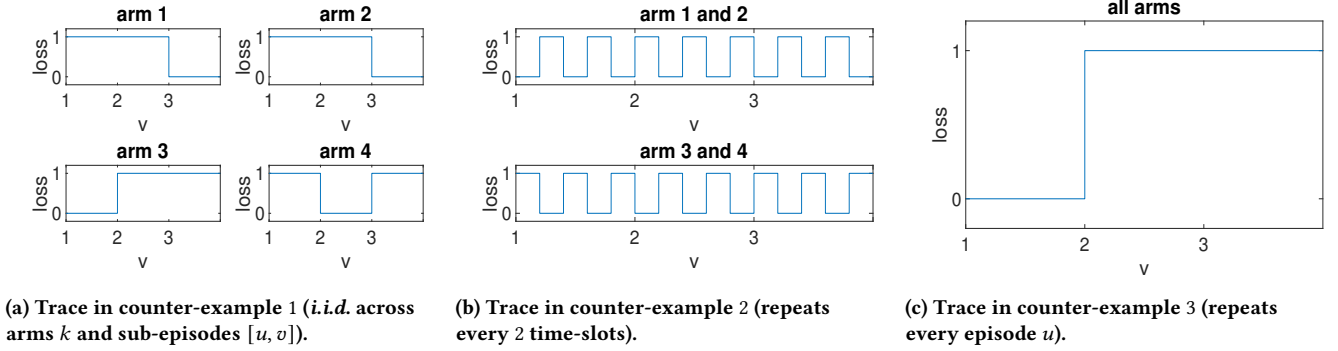


Figure 1: One realization of the counter-example traces in one episode.

η_2 as follows (i.e., Step 5 in Algorithm 4),

$$\hat{w}_k^{\text{ROW}}(t+1) = \hat{w}_k^{\text{ROW}}(t) \cdot e^{-\eta_2 l_k(t)}, \quad (27)$$

and then proceeds to the next time-slot $t+1$. Note that the weights $\hat{w}_k^{\text{ROW}}(t)$ are reset by (25) in the first time-slot $t = t_{u,v}$ of each sub-episode (u, v) .

Finally, at the end of the last time-slot of the entire episode u , ROW collects all the feedback received during the episode. Next, during the sub-episodes that arm k was chosen for the working group, ROW subtracts the loss of the primary arm from the corresponding loss of this arm k . Then, the resulting value is divided by the conditional probability that k is chosen as a secondary arm (conditioned on k not being the primary arm), i.e., $\frac{M-1}{K-1}$. Precisely, we let $v_u(k) \triangleq \{v \mid v = 1, \dots, V, k \in \hat{\mathbb{K}}^{\text{ROW}}[u, v]\}$ denote the sub-episodes (u, v) when the arm k was chosen in the working group. Let $L_k[u, v_u(k)] \triangleq \sum_{v \in v_u(k)} \sum_{t=t_{u,v}}^{t_{u,v}+\tau_2-2} l_k(t)$ denote the sum of the losses of arm k in sub-episodes (u, v) (except the last time-slot $t = t_{u,v} + \tau_2 - 1$) for all $v \in v_u(k)$. Then, ROW computes the loss difference of each arm $k \in \mathcal{K}$ as follows,

$$\tilde{L}_k^{\text{ROW}}[u] = \frac{L_k[u, v_u(k)] - L_{k_0^{\text{ROW}}}[u][u, v_u(k)]}{\frac{M-1}{K-1}}. \quad (28)$$

Note that for the primary arm $k_0^{\text{ROW}}[u]$, the loss difference is $\tilde{L}_{k_0^{\text{ROW}}}[u] = 0$, which is also consistent with (28). Then, ROW updates the weights for all the arms $k \in \mathcal{K}$ with a tunable parameter η_1 as follows (i.e., Step 6 in Algorithm 4),

$$w_k^{\text{ROW}}[u+1] = w_k^{\text{ROW}}[u] \cdot e^{-\eta_1 \tilde{L}_k[u]}, \quad (29)$$

which becomes the initial weights for the next episode $u+1$. In (46), we give the values of all parameters of ROW, i.e., η_1 , η_2 , τ_1 and τ_2 .

Readers familiar with bandit-learning algorithms may have already noticed two other crucial differences in ROW. First, a different weight-decay parameter η_2 is used to update weights in (27) within the episode, compared with the parameter η_1 that is used in (29) across episodes. Second, when updating the weights across episodes in (29), we use the difference between the loss of an arm and that of the primary arm, instead of using the absolute loss of the arm directly. In the following, we explain why these two differences (i.e., our idea 2 and idea 3) are crucial for achieving the $O(\sqrt{T})$ regret.

Idea 2: Use different weight-decay parameters η_2 and η_1 . Recall that in every episode, ROW can obtain at least $\frac{1}{V}$ fraction of feedback from every arm. We would have hoped that this amount of feedback is sufficient for attaining a low $O(\sqrt{T})$ regret. Indeed, consider an alternate bandit-learning problem where the feedback of each arm is obtained independently with probability $\frac{1}{V}$ in every time-slot. It is not difficult to show that Exp3 [3] using this amount of feedback will attain the $O(\sqrt{T})$ regret.

However, compared with the above alternate problem, the difficulty we are facing here is that in ROW the feedback becomes highly correlated in time. Indeed, the secondary arms are fixed during the whole sub-episode. Thus, we either have all feedback of an arm, or have none for the whole sub-episode. Below, we construct two counter-examples to illustrate the difficulties in dealing with such correlation. For ease of exposition, we use $l(t_1 : t_2) \triangleq [l(t), \text{ for all } t = t_1, t_1 + 1, \dots, t_2]$ to collect $l(t)$ from $t = t_1$ to $t = t_2$.

Counter-example 1: Consider $K = 4$ arms and $M = 2$. For each arm k , in each sub-episode (u, v) , $l_k(t_{u,v} : t_{u,v} + \tau_2 - 1) = 0$ with probability $\frac{1}{2}$, and $l_k(t_{u,v} : t_{u,v} + \tau_2 - 1) = 1$ with probability $\frac{1}{2}$. The losses are independent across arms k and across sub-episodes $[u, v]$. Please see Fig. 1a for this loss trace in one episode. Using this counter-example, we show why existing bandit-learning method, Exp3 [1], could lead to a poor regret. Let us consider the optimal static loss. First, the expected total loss of each arm is trivially $\mathbb{E}[L] = \frac{T}{2}$. Second, let us estimate the variance of the total loss of each arm. Since the loss is a constant within a sub-episode, the higher correlation in time leads to a higher variance in the total loss of each arm. Specifically, for each arm, the variance of its total loss in a sub-episode⁵ is $\Theta(\tau_2^2)$. Thus the variance of its total loss across T time-slots is $\text{Var}(L) = \frac{T}{\tau_2} \cdot \Theta(\tau_2^2) = \Theta(T^{\frac{3}{2}})$. Thus, one of the K arms may incur a total loss that is smaller than the average by $\Theta(\sqrt{\text{Var}(L)})$. As a result, the total loss of the optimal static decision OPT is $\mathbb{E}[L] - \Theta(\sqrt{\text{Var}(L)}) = \frac{T}{2} - \Theta(T^{\frac{3}{4}})$. (This estimate can also be obtained by applying the random walk analysis [21, p. 111].) Next, we consider the total loss of the episodic version of Exp3 [1]. Such version of Exp3 picks an arm k_0 at the beginning of an episode, and use it as the active arm for the entire episode. Since the loss in each episode is independent, the total loss of such Exp3 will be the

⁵In contrast, if the losses were *i.i.d.* in time, the variance should have been $\Theta(\tau_2)$.

average loss of each arm in this counter-example, i.e., $\frac{T}{2}$. Therefore, the regret would be $\Theta(T^{\frac{3}{4}})$.

Counter-example 1 clearly illustrates why the higher correlation in time leads to a higher regret for the episodic version of Exp3. To overcome this difficulty, we make an important observation. In this setting with highly correlated losses, we observe that one arm (with losses 0) will be consistently better than the other arms (with losses 1) in each sub-episode. We may then beat the average loss by switching to the better arm within a sub-episode. Indeed, with $M = 2$, the chance that one of the two arms incurs zero loss is $\frac{3}{4}$. Thus, if we can switch to the better arm (with losses 0) quickly within a sub-episode, we may attain a total loss approximately equals to $\frac{T}{4}$, which would have beaten the optimal static decision OPT. This counter-example thus suggests why it is important to use Exp3 [3] inside each sub-episode (in addition to across episodes).

However, it is still highly non-trivial to choose the parameter η of Exp3 within each sub-episode. One possible thought is that, we can think of each sub-episode as a bandit-learning problem with $\tau_2 = \Theta(\sqrt{T})$ time-slots. Then, if we view the better arm within the sub-episode as the static optimal arm, we would have to use $\eta = \Theta(T^{-\frac{1}{4}})$ in order to attain the minimal regret against the better arm. However, this choice of η would have been too large, as can be seen in the counter-example below.

Counter-example 2: Consider $K = 4$ arms and $M = 2$. For arms $k = 1, 2$, $l_k(t) = 0$ for all odd time-slots t , and $l_k(t) = 1$ for all even time-slots t . For arms $k = 3, 4$, $l_k(t) = 1$ for all odd time-slots t , and $l_k(t) = 0$ for all even time-slots t . Please see Fig. 1b for this loss trace in one episode. Using this counter-example, we can see why using Exp3 [3] with a parameter $\eta = \Theta(T^{-\frac{1}{4}})$ could lead to a poor regret. Let us consider the optimal static loss. Since the total loss of every arm is $\frac{T}{2}$, the optimal static loss is $\frac{T}{2}$. Next, we consider the total loss of Exp3. Notice that the probability of each arm is initialized to be the same, i.e., $\frac{1}{K}$, at time $t = 1$. Then, at each time, suppose that all arms have been observed almost the same number of times. Thus, the probabilities of all arms would be about the same. However, whenever an arm with loss $l_{k_1}(t) = 0$ and an arm with loss $l_{k_2}(t) = 1$ are observed simultaneously, at the next time $t + 1$ Exp3 will use the arm k_1 as the active arm with a probability higher by approximately $\Theta(\eta)$. According to counter-example 2, $l_{k_1}(t + 1) = 1$. Thus, Exp3 will suffer an additional loss $\Theta(\eta)$ approximately at each time. Hence, the total loss of Exp3 will be $\frac{T}{2} + \Theta(\eta T) = \frac{T}{2} + \Theta(T^{\frac{3}{4}})$. Therefore, the regret would be $\Theta(T^{\frac{3}{4}})$.

Counter-example 2 clearly indicates that, in order to attain the $O(\sqrt{T})$ regret, the parameter η_2 should be no larger than $O(T^{-\frac{1}{2}})$. However, since a sub-episode is of length much smaller than T , we conjecture that η_2 still needs to be larger than η_1 (the latter is used across episodes), so that ROW converges fast to the better arm inside the chosen working group. Lemma 4.4 in Sec. 4.2.2 will provide the exact condition on how η_2 and η_1 should be tuned to obtain the $O(\sqrt{T})$ regret.

Idea 3: Use the loss difference from the primary arm to update weights across episodes. We next describe why it is also crucial to use the loss difference in (28) instead of the absolute loss of each arm. Recall that at the end of each episode, we receive τ_1 feedback from the primary arm, but only $\tau_2 = \frac{\tau_1}{V}$ feedback from each secondary

arm. Intuitively, this bias will also increase the variance of the total losses accumulated in the past, which again leads to a higher regret. The following counter-example illustrates this difficulty.

Counter-example 3: Consider $K = 4$ arms and $M = 2$. In the first sub-episode of each episode, the loss of each arm at each time is 0. For all subsequent sub-episodes of each episode, the loss of each arm at each time is 1. Please see Fig. 1c for this loss trace in one episode. In the literature, the standard way to deal with this bias in the amount of feedback is to divide the observed loss by the probability that the arm is observed [1, 3, 5]. For each arm, this probability is $p_k[u] + (1 - p_k[u])\frac{M-1}{K-1}$, where $p_k[u]$ is the probability that arm k is chosen as the primary arm, and $(1 - p_k[u])\frac{M-1}{K-1}$ is the probability that arm k is chosen as the secondary arm in a sub-episode. With this mechanism, the estimated losses will be $\tilde{L}_k[u] = \frac{2\tau_2}{p_k[u] + (1 - p_k[u])\frac{M-1}{K-1}}$ when k is the primary arm, $\tilde{L}_k[u] = 0$ when k is a secondary arm that is chosen in the first ($v = 1$) sub-episode, and $\tilde{L}_k[u] = \frac{\tau_2}{p_k[u] + (1 - p_k[u])\frac{M-1}{K-1}}$ when k is a secondary arm that is chosen in the subsequent ($v = 2, 3$) sub-episodes. Suppose that $p_k[u] = \frac{1}{K}$ is the same across all arms. Then, the denominator is actually the same across all arms, but the numerator will still lead to a significant variance. Indeed, since the primary arm is chosen randomly with probability $p_k[u] = \frac{1}{K}$, it is not hard to verify that the total estimated loss of each arm over an episode will have a variance of $\Theta(\tau_2^2)$. In contrast, if full feedback was available, all arms would have a total loss equal to $2\tau_2$ in an episode, and the variance would have been zero. It is easy to show that, with this additional $\Theta(\tau_2^2)$ gap in the variance, the regret of Exp3 [1] is still $O(T^{\frac{3}{2}})$, which is much larger than $O(\sqrt{T})$.

Counter-example 3 thus suggests that, instead of dividing the loss by the probability of observing an arm, we need some new ways to deal with the above bias issue. Precisely, in (28), ROW updates the estimated loss by the difference of the loss of each secondary arm and that of the primary arm. In addition, the loss difference of the primary arm is simply 0. Returning to counter-example 3, the new estimated loss will then be $\tilde{L}_k[u] = 0$ for all the arms $k \in \mathcal{K}$. Thus, the additional variance $\Theta(\tau_2^2)$ of the estimated losses has been eliminated, which is also crucial for attaining the $O(\sqrt{T})$ regret.

4.2 Regret Analysis

In Theorem 4.1 below, we show the upper bound of the regret attained by ROW. For ease of exposition, we focus on the case when $K - 1$ is divisible by $M - 1$. (It is not difficult to extend to the case when $K - 1$ is not divisible by $M - 1$. Please see Appendix L for details.)

THEOREM 4.1. *Consider bandit learning with switching costs and full-feedback costs introduced in Sec. 2.1. When $M \geq 2$, the regret of ROW can be upper-bounded as follows, for $T \geq \frac{448(K-1)^2 \ln K}{\frac{5}{2} + 2\beta_1}$,*

$$R^{\text{ROW}}(T) \leq 8b_1 \frac{K-1}{M-1} \sqrt{\ln K} \sqrt{T} + b_2, \quad (30)$$

where $b_1 = \sqrt{\frac{5}{2} + 2b_3\beta_1}$, $b_2 = b_3\beta_1 + 1$ and $b_3 = \min\{M, K - M\}$.

In Sec. 3 when $M = 1$, the optimal regret is $\Theta(T^{\frac{2}{3}})$ for bandit learning with switching costs and full-feedback costs. In sharp contrast, now with $M \geq 2$, ROW achieves a significantly lower regret equals to $O(\sqrt{T})$. Moreover, ROW never uses full feedback. Further, as M increases, the regret of ROW can be further reduced. *To the best of our knowledge, this is the first result in the literature to utilize the flexibility of choosing $M \geq 2$ arms to improved the regret to $O(\sqrt{T})$ for bandit learning with switching costs.* Furthermore, using a trivial lower bound for bandit learning with free full-feedback [5, 14], we can conclude that the $O(\sqrt{T})$ regret cannot be further improved.

The rest of this section is devoted to the proof of Theorem 4.1. Due to the three new ideas in ROW, new analytical techniques are needed to capture the evolution of the weights, which are also of independent interest. In order to relate the loss of ROW to that of the optimal static loss, our analysis below is carried out in three steps. First, inside each sub-episode, we relate the total loss of ROW in each sub-episode to a log-sum-exp function $g_2[u, v]$ of the parameter η_2 and the feedback from the chosen working group. Second, at the end of each episode, we relate $g_2[u, v]$ of all sub-episodes to another log-sum-exp function $g_1[u]$ of the parameter η_1 and the loss difference $\tilde{L}_k^{\text{ROW}}[u]$. Third, across all episodes, we relate $g_1[u]$ to the optimal static loss. Combining these three steps, the total loss of ROW will then be related to the optimal static loss. In the following, we let $\mathcal{H}[u-1]$ denotes the σ -algebra generated by the observation of ROW from time $t = 1$ to $t = (u-1)\tau_1$. Let $L_k[u, v] \triangleq \sum_{t=t_{u,v}}^{t_{u,v}+\tau_2-2} l_k(t)$.

4.2.1 Inside each sub-episode. We start by relating the expected loss of ROW inside each sub-episode (u, v) to a log-sum-exp function $g_2[u, v]$ (see Lemma 4.2). This function $g_2[u, v]$ will then be further related to the variance of the feedback from the chosen working group $\hat{\mathbb{K}}^{\text{ROW}}[u, v]$ in the sub-episode (see Lemma 4.3). Recall that in (25), the weights $\hat{w}_k^{\text{ROW}}(t_{u,v})$ in the first time-slots of all sub-episodes are initialized to be the weights $w_k^{\text{ROW}}[u]$ at the beginning of the episode u . Thus, given a same working group, the probabilities $\hat{p}_k^{\text{ROW}}(t_{u,v})$ are also the same at the beginning of all sub-episode v in an episode u . We let

$$\hat{p}_k^{\text{ROW}}[u] \triangleq \frac{w_k^{\text{ROW}}[u]}{\sum_{k \in \hat{\mathbb{K}}^{\text{ROW}}[u,v]} w_k^{\text{ROW}}[u]} \quad (31)$$

denote this common probability.

LEMMA 4.2. *For each sub-episode (u, v) , given the history $\mathcal{H}[u-1]$ and the chosen working group $\hat{\mathbb{K}}[u, v]$, we have*

$$\sum_{t=t_{u,v}}^{t_{u,v}+\tau_2-1} \sum_{k \in \hat{\mathbb{K}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}(t) l_k(t) \leq g_2[u, v] + \frac{1}{2} \eta_2 \tau_2 + 1, \quad (32)$$

where

$$g_2[u, v] \triangleq -\frac{1}{\eta_2} \ln \left(\sum_{k \in \hat{\mathbb{K}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}[u] e^{-\eta_2 L_k[u, v]} \right). \quad (33)$$

On the left-hand-side of (32), the probability $\hat{p}_k^{\text{ROW}}(t)$ is the probability of using arm k as the active arm. Thus, the left-hand-side of (32) represents the conditional (conditioned on the working

group $\hat{\mathbb{K}}^{\text{ROW}}[u, v]$ and history $\mathcal{H}[u-1]$) expected loss of ROW in sub-episode (u, v) . Hence, (32) upper-bounds the conditional expected loss of ROW by a log-sum-exp function $g_2[u, v]$ and the term $\frac{1}{2} \eta_2 \tau_2 + 1$. We make two important comments. First, the value of $g_2[u, v]$ is approximated dominated by the arm with the smallest loss $L_k[u, v]$ (whenever the corresponding probability $\hat{p}_k^{\text{ROW}}[u]$ is non-zero). (32) thus confirms that ROW is trying to switch to the “better” arm in the working group. Second, the gap $\frac{1}{2} \eta_2 \tau_2$ is much smaller than the gap $\frac{1}{2} \eta_2 \tau_2^2$ incurred by the episodic version of Exp3 [1]. Note that the above-mentioned two conclusions precisely capture our ideas 1 and 2, which together allow ROW to converge quickly to the better arm in the working group. Please see Appendix H for the complete proof of Lemma 4.2.

The following lemma then relates $g_2[u, v]$ to the expectation and variance of the feedback from the chosen working group in the sub-episode, which will be useful when we move to the second-step of studying the weight updates at the end of each episode.

LEMMA 4.3. *For each sub-episode (u, v) , given the history $\mathcal{H}[u-1]$ and the chosen working group $\hat{\mathbb{K}}^{\text{ROW}}[u, v]$, if $\eta_2 \tau_2 \leq \ln 2$, we have*

$$g_2[u, v] \leq \mathbb{E} \left[L[u, v] | \mathcal{H}[u-1], \hat{\mathbb{K}}^{\text{ROW}}[u, v] \right] - \frac{\eta_2}{8} \cdot \text{Var} \left(L[u, v] | \mathcal{H}[u-1], \hat{\mathbb{K}}^{\text{ROW}}[u, v] \right), \quad (34)$$

where the expectation is taken with regard to the randomness in $\hat{p}_k^{\text{ROW}}[u]$, i.e.,

$$\begin{aligned} \mathbb{E} \left[L[u, v] | \mathcal{H}[u-1], \hat{\mathbb{K}}^{\text{ROW}}[u, v] \right] &\triangleq \sum_{k \in \hat{\mathbb{K}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}[u] L_k[u, v], \\ \text{Var} \left(L[u, v] | \mathcal{H}[u-1], \hat{\mathbb{K}}^{\text{ROW}}[u, v] \right) &\triangleq \sum_{k \in \hat{\mathbb{K}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}[u] \\ &\cdot \left(L_k[u, v] - \mathbb{E} \left[L[u, v] | \mathcal{H}[u-1], \hat{\mathbb{K}}^{\text{ROW}}[u, v] \right] \right)^2. \end{aligned}$$

Notice that the expectation and variance on the right-hand-side of (34) are for the feedback from the working group $\hat{\mathbb{K}}^{\text{ROW}}[u, v]$. Thus, Lemma 4.3 shows that the log-sum-exp function $g_2[u, v]$ can be related to the expectation and variance of the feedback from the chosen working group. Given the working group $\hat{\mathbb{K}}^{\text{ROW}}[u, v]$, Lemma 4.3 is proved by applying the Taylor expansion of the e^{-x} function to $g_2[u, v]$. Please see Appendix I for the complete proof of Lemma 4.3.

4.2.2 Relating the loss upper-bound at the end of a sub-episode to the weights across episodes. Lemma 4.2 provides an upper bound on the loss of ROW at the end of each sub-episode (u, v) . Note that this upper bound depends on η_2 . On the other hand, at the end of each episode u , we calculate the weights according to (29). Notice that not only is $\tilde{L}_k^{\text{ROW}}[u]$ in (29) different from $L_k[u, v]$ in (33), the parameter η_2 is also different from η_1 . Thus, we need a way to convert the loss upper-bound in Lemma 4.2 for each sub-episode to a form that depends on the weights calculated by (29). This is accomplished by Lemma 4.4 below. Further, this lemma gives a sufficient condition on how to tune the parameters η_2 and η_1 .

Specifically, notice that the loss difference $\tilde{L}_k^{\text{ROW}}[u]$ calculated in (28) is a difference from the loss of the primary arm $k_0^{\text{ROW}}[u]$. We let $g_2[u]$ denote the sum of $g_2[u, v]$ for all sub-episodes v , minus

a term that corresponds to the loss of the primary arm, i.e.,

$$\begin{aligned} g_2[u] &\triangleq \sum_{v=1}^V g_2[u, v] - \sum_{v=1}^V L_{k_0^{\text{ROW}}[u]}[u, v] \\ &= -\frac{1}{\eta_2} \sum_{v=1}^V \ln \left(\sum_{k \in \hat{\mathbb{K}}^{\text{ROW}}[u, v]} \hat{p}_k^{\text{ROW}}[u] e^{-\eta_2 \mathcal{L}_k^{\text{ROW}}[u, v]} \right), \end{aligned} \quad (35)$$

where $\mathcal{L}_k^{\text{ROW}}[u, v] = L_k[u, v] - L_{k_0^{\text{ROW}}[u]}[u, v]$.

LEMMA 4.4. *If the parameters η_2 , τ_2 , η_1 and τ_1 satisfy that*

$$\eta_2 \geq 16 \left(\frac{K-1}{M-1} \right)^2 \cdot \eta_1, \quad \eta_2 \tau_2 \leq \ln 2 \quad \text{and} \quad \eta_1 \tau_1 \leq \ln 2, \quad (36)$$

we have

$$\begin{aligned} &\mathbb{E}_{\hat{\mathbb{K}}^{\text{ROW}}[u, 1:V]} \left[g_2[u] \middle| \mathcal{H}[u-1] \right] \\ &\leq \mathbb{E}_{\hat{\mathbb{K}}^{\text{ROW}}[u, 1:V]} \left[g_1[u] \middle| \mathcal{H}[u-1] \right], \end{aligned} \quad (37)$$

where the expectation is taken with respect to the randomness in the working groups, and

$$g_1[u] \triangleq -\frac{1}{\eta_1} \ln \left(\sum_{k=1}^K p_k^{\text{ROW}}[u] e^{-\eta_1 \tilde{L}_k^{\text{ROW}}[u]} \right). \quad (38)$$

The log-sum-exp function $g_2[u]$ on the left-hand-side of (37) is related to $g_2[u, v]$ through (35), which is then related to the loss of ROW in each sub-episode through (32). The log-sum-exp function $g_1[u]$ on the right-hand-side of (37) is related to the weights calculated at the end of the episode. Thus, Lemma 4.4 relates the loss upper-bound at the end of each sub-episode to the weights across episodes, and (36) confirms our conjecture that η_2 should be larger than η_1 .

Please see Appendix M for the complete proof of Lemma 4.4. In the following, we sketch the key steps (Step 1 - Step 3 below) for proving Lemma 4.4, which may also be of independent interest.

Sketch of proof of Lemma 4.4:

Step-1: Similar to Lemma 4.3, we can derive a lower bound of $g_1[u]$ by relating it to the expectation and variance of the loss differences.

LEMMA 4.5. *For each episode u , given the history $\mathcal{H}[u-1]$ and the chosen working groups $\hat{\mathbb{K}}^{\text{ROW}}[u, 1:V]$, if $\eta_1 \tau_1 \leq \ln 2$, we have*

$$\begin{aligned} g_1[u] &\geq \mathbb{E} \left[\tilde{L}^{\text{ROW}}[u] \middle| \mathcal{H}[u-1], \hat{\mathbb{K}}^{\text{ROW}}[u, 1:V] \right] \\ &\quad - \eta_1 \cdot \text{Var} \left(\tilde{L}^{\text{ROW}}[u] \middle| \mathcal{H}[u-1], \hat{\mathbb{K}}^{\text{ROW}}[u, 1:V] \right), \end{aligned} \quad (39)$$

where the expectation is taken with regard to the randomness in $p_k^{\text{ROW}}[u]$, i.e.,

$$\begin{aligned} &\mathbb{E} \left[\tilde{L}^{\text{ROW}}[u] \middle| \mathcal{H}[u-1], \hat{\mathbb{K}}^{\text{ROW}}[u, 1:V] \right] \triangleq \sum_{k=1}^K p_k^{\text{ROW}}[u] \tilde{L}^{\text{ROW}}[u], \\ &\text{Var} \left(\tilde{L}^{\text{ROW}}[u] \middle| \mathcal{H}[u-1], \hat{\mathbb{K}}^{\text{ROW}}[u, 1:V] \right) \triangleq \sum_{k=1}^K p_k^{\text{ROW}}[u] \\ &\quad \cdot \left(\tilde{L}^{\text{ROW}}[u] - \mathbb{E} \left[\tilde{L}^{\text{ROW}}[u] \middle| \mathcal{H}[u-1], \hat{\mathbb{K}}^{\text{ROW}}[u, 1:V] \right] \right)^2. \end{aligned}$$

Please see Appendix J for the complete proof of Lemma 4.5.

Step-2: Lemma 4.4 is then proved by mainly comparing the expectations of (34) and (39) with regard to the randomness in the working groups. Here, we use the help of a fictitious “full feedback” system, where we assume that there is an oracle who knows the losses from all arms in all time-slots during the episode. Further, this oracle assigns the probability distribution $p_k^{\text{ROW}}[u]$ on the arms.

It is easy to show that the expectations of both working-group feedback and the loss differences are related to the expectation of the fictitious “full feedback”. Further, Lemma 4.6 and Lemma 4.7 below show that the variances of both the working-group feedback and loss differences can also be related to the variance of full feedback, given by $\text{Var}(L[u, v] | \mathcal{H}[u-1])$ in the lemma below.

LEMMA 4.6. *For each sub-episode (u, v) , given the history $\mathcal{H}[u-1]$, we have*

$$\begin{aligned} &\mathbb{E}_{\hat{\mathbb{K}}^{\text{ROW}}[u, v]} \left[\text{Var} \left(L[u, v] \middle| \hat{\mathbb{K}}^{\text{ROW}}[u, v] \right) \middle| \mathcal{H}[u-1] \right] \\ &\geq \frac{M-1}{K-1} \cdot \text{Var} \left(L[u, v] \middle| \mathcal{H}[u-1] \right), \end{aligned} \quad (40)$$

where

$$\begin{aligned} &\text{Var} \left(L[u, v] \middle| \mathcal{H}[u-1] \right) \\ &\triangleq \sum_{k=1}^K p_k^{\text{ROW}}[u] \left(L[u, v] - \sum_{k=1}^K p_k^{\text{ROW}}[u] L[u, v] \right)^2. \end{aligned}$$

The variance on the left-hand-side of (40) is for the losses from the feedback in the working group $\hat{\mathbb{K}}^{\text{ROW}}[u, v]$. The outside expectation is taken over all possible working groups. The variance on the right-hand-side of (40) is for the fictitious “full feedback”. Intuitively, if the right-hand-side of (40) is strictly positive, there must be some difference among the losses of the arms. Then, even when a random subset of arms is chosen into the working group, we should still see some variance. That is the intuition why the left-hand-side of (40) must also be strictly positive, which is the conclusion in Lemma 4.6. Moreover, as M increases, the constant factor $\frac{M-1}{K-1}$ increases to be closer to 1. This is one of the reasons that the regret of ROW decreases with M . In sharp contrast, when $M=1$, we have $\frac{M-1}{K-1} = 0$. Indeed, in this case, no matter how large the variance of full feedback is, the variance on the left-hand-side of (40) will always be equal to 0. This is one of the reasons for the sharp transition from the $O(T^{\frac{2}{3}})$ regret when $M=1$ to the $O(\sqrt{T})$ regret when $M \geq 2$. Please see Appendix K for the complete proof of Lemma 4.6.

Step-3: However, the fictitious “full feedback” is not available to the online learning algorithm. Hence, Lemma 4.6 is not very useful unless we can related the full feedback to the loss difference that we design in (28). This is exactly the purpose of Lemma 4.7 below.

LEMMA 4.7. *For each episode u , given the history $\mathcal{H}[u-1]$, we have*

$$\begin{aligned} &\sum_{v=1}^V \text{Var} \left(L[u, v] \middle| \mathcal{H}[u-1] \right) \geq \frac{M-1}{2(K-1)} \\ &\quad \cdot \mathbb{E}_{\hat{\mathbb{K}}^{\text{ROW}}[u, 1:V]} \left[\text{Var} \left(\tilde{L}^{\text{ROW}}[u] \middle| \hat{\mathbb{K}}^{\text{ROW}}[u, 1:V] \right) \middle| \mathcal{H}[u-1] \right]. \end{aligned} \quad (41)$$

Different from Lemma 4.6, Lemma 4.7 focuses on the variance of the loss differences $\tilde{L}^{\text{ROW}}[u]$, as in the right-hand-side of (41). Moreover, the expectation is taken over all possible sequences of the working groups for the whole episode. Thus, the variance of full feedback on the left-hand-side of (41) is also summed over all sub-episodes v . Intuitively, if the right-hand-side of (41) is strictly positive, there must exist some difference across the secondary arms when comparing with the common primary arm. Then, the differences among the secondary arms cannot all be 0. This means there must be some variance of the full feedback. This is the intuition why the left-hand-side of (41) must also be strictly positive, which is the conclusion in Lemma 4.7. Similar to that in (40), as M increases, the constant factor $\frac{M-1}{2(K-1)}$ increases to be closer to 1. This is another reason that the regret of ROW decreases with M . In sharp contrast, when $M = 1$, we have $\frac{M-1}{2(K-1)} = 0$, which again implies a sharp transition from $M = 1$ to $M \geq 2$. By comparing the constant factors in (34) and (40) with that in (39) and (41), we can see that to obtain (37), η_2 needs to be larger than $16 \left(\frac{K-1}{M-1}\right)^2 \cdot \eta_1$. Please see Appendix L for the complete proof of Lemma 4.7.

REMARK 1. *Lemma 4.7 is the result of using our idea 3. In other words, without our idea 3 for constructing the loss differences $\tilde{L}_k^{\text{ROW}}[u]$ in (28), Lemma 4.7 may not hold. For example, in the counter-example 3 that we introduced in Sec. 4.1, the variance of full feedback is 0. Without our idea 3, the variance of the absolute loss from the feedback in all sub-episodes will be $\Theta(\tau_2^2)$, which would have made Lemma 4.7 invalid. In contrast, with our idea 3, the loss difference is the difference from the loss of the primary arm, which will be 0 for all arms. Thus, the variance of the loss differences of ROW in each episode is 0, which is the same as the variance of full feedback.*

Combining Lemma 4.3, Lemma 4.5, Lemma 4.6, and Lemma 4.7, we can then prove Lemma 4.4. The detailed proof is available in Appendix M.

Up to now, by combining (32), (35) and (37) for all sub-episode v and episode u , we can relate the total loss of ROW to $g_1[u]$ as follows,

$$\begin{aligned} & \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[\sum_{v=1}^V \sum_{t=t_{u,v}}^{t_{u,v}+\tau_2-1} \sum_{k \in \hat{k}^{\text{ROW}}[u,v]} \hat{P}_k^{\text{ROW}}(t) \right. \right. \\ & \quad \left. \left. \cdot l_k(t) - \sum_{v=1}^V L_{k_0^{\text{ROW}}[u]}[u,v] \mathcal{H}[u-1] \right] \right\} \\ & \leq \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[g_1[u] \mathcal{H}[u-1] \right] \right\} + \frac{1}{2} \eta_2 T + VU. \end{aligned} \quad (42)$$

In the next section, we show how to relate the first term on the right-hand-side of (42) to the optimal static loss.

4.2.3 Relating the upper-bound of the total loss of ROW to the optimal static loss. Lemma 4.8 below relates the sum of $g_1[u]$ on the right-hand-side of (42) to the optimal static loss of OPT.

LEMMA 4.8. *We have the following inequality,*

$$\begin{aligned} & \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[g_1[u] \mathcal{H}[u-1] \right] \right\} \\ & \leq \text{Cost}^{\text{OPT}}(1:T) + \frac{\ln K}{\eta_1} \\ & \quad - \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E} \left[\sum_{v=1}^V L_{k_0^{\text{ROW}}[u]}[u,v] \mathcal{H}[u-1] \right] \right\}. \end{aligned} \quad (43)$$

In (43), the term on the left-hand-side is one of the terms in the upper bound of the total loss of ROW, i.e., the first term on the right-hand-side of (42). The first term on the right-hand-side is the optimal static loss. The second term on the right-hand-side of (43) can be obtained by following the Exp3 analysis [3]. The third term on the right-hand-side of (43) is because the loss of the primary arm is subtracted in $g_2[u]$ (see (35)). This term also appears on the left-hand-side of (42), which will eventually be cancelled. Please see Appendix N for the complete proof of Lemma 4.8.

4.2.4 The final regret. Since ROW only switches at the boundaries of the sub-episodes, the total switching cost of ROW can be upper-bounded as follows,

$$\sum_{t=1}^T \sum_{k \in \hat{k}^{\text{ROW}}(t)} \beta_1 \mathbf{1}_{\{k \notin \hat{k}^{\text{ROW}}(t-1)\}} \leq \min\{M, K-M\} \cdot \beta_1 \left\lceil \frac{T}{\tau_2} \right\rceil. \quad (44)$$

Next, since ROW never asks for full feedback, the total full-feedback cost of ROW is 0. Combining (42), (43) and (44), we can see that the regret of ROW is upper-bounded as follows,

$$R^{\text{ROW}}(T) \leq \frac{\ln K}{\eta_1} + \frac{1}{2} \eta_2 T + \min\{M, K-M\} \cdot \beta_1 \left\lceil \frac{T}{\tau_2} \right\rceil + \left\lceil \frac{T}{\tau_2} \right\rceil. \quad (45)$$

Then, by choosing

$$\begin{cases} \eta_2 = \frac{c_1 c_2}{\sqrt{T}}, & \tau_2 = \left\lceil \frac{\ln 2}{c_1 c_2} \sqrt{T} \right\rceil, \\ \eta_1 = \frac{c_1}{c_2 \sqrt{T}}, & \tau_1 = \left\lceil \frac{K-1}{M-1} \left\lceil \frac{\ln 2}{c_1 c_2} \sqrt{T} \right\rceil \right\rceil, \end{cases} \quad (46)$$

where $c_1 = \sqrt{\frac{\ln K}{\frac{5}{2} + \min\{M, K-M\} \cdot 2\beta_1}}$ and $c_2 = \frac{4(K-1)}{M-1}$, we have

$$\begin{aligned} R^{\text{ROW}} & \leq \frac{8(K-1)}{M-1} \sqrt{\frac{5}{2} + \min\{M, K-M\} \cdot 2\beta_1} \sqrt{\ln K} \sqrt{T} \\ & \quad + \min\{M, K-M\} \cdot \beta_1 + 1, \end{aligned} \quad (47)$$

for $T \geq \frac{448(K-1)^2 \ln K}{\frac{5}{2} + 2\beta_1}$. The result of Theorem 4.1 then follows. Please see Appendix O for the complete proof of Theorem 4.1.

5 NUMERICAL RESULTS

In this section, we present numerical results comparing our new algorithms ROCF introduced in Sec. 3.3 (for $M = 1$) and ROW introduced in Sec. 4.1 (for $M \geq 2$) with the episodic version of Exp3 proposed in [1]. According to [1], the theoretical regret of the episodic version of Exp3 is $\Theta(K^{\frac{1}{3}} T^{\frac{2}{3}})$.

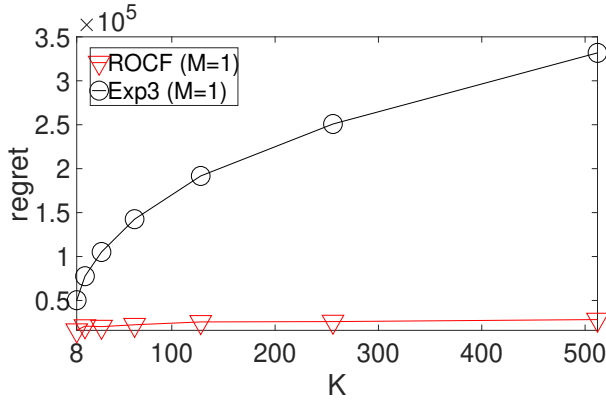


Figure 2: Compare the regrets of ROCF and the episodic version of Exp3.

5.1 The Case of $M = 1$

In the case with $M = 1$, we compare the regret of ROCF (that we proposed in Sec. 3.3) with that of the episodic version of Exp3 (proposed in [1]). As we discussed in Sec. 3.1, when $\beta_2 < \frac{3}{4}K\beta_1$, ROCF improves the dependence of the regret on the number K of arms from $K^{\frac{1}{3}}$ to $(\ln K)^{\frac{1}{3}}$. Thus, here we focus on the case when the full-feedback cost β_2 is smaller than the switching cost β_1 . Specifically, we let the switching cost and full-feedback cost be $\beta_1 = 10$ and $\beta_2 = 1$, respectively. We use the lower-bound trace that we designed in Sec. 3.2, where the DS adversary runs the MHM adversary as the subroutine. We consider $T = 10^6$ time-slots. We compare how the regret increases with the number of arms K . Please see Fig. 2. From Fig. 2, we can see that the regret of ROCF is much smaller than that of Exp3, especially when K is large. For example, when $K = 512$, the regret of Exp3 is around 3.32×10^5 . In contrast, the regret of ROCF is only about 2.83×10^4 .

5.2 The case of $M \geq 2$

In the case with $M \geq 2$, we compare the regret of ROW (that we proposed in Sec. 4.1) with that of ROCF (that we proposed in Sec. 3.3) and the episodic version of Exp3 (proposed in [1]). We consider ROW with $M = 2$ and ROW with $M = 3$. In Fig. 3, we use both the lower-bound trace that we designed in Sec. 3.2 and the three counter-example traces that we designed in Sec. 4.1. We consider $K = 4$ arms, the full-feedback cost $\beta_2 = 1$ and the switching cost $\beta_1 = 1$. We compare how the regret increases with the time length T . From Fig. 3, we can see that for all 4 traces, the regret of ROW (with $M = 2$ and with $M = 3$) is much smaller than that of Exp3 (and ROCF). For example, when using counter-example 3 and $T = \sqrt{10} \times 10^6$, the regret of Exp3 is around 2.61×10^4 . In contrast, the regret of ROW with $M = 2$ is only about 3.22×10^3 , confirming the power of using 2 arms. Moreover, we can see that when M increases, the gap between the regret of ROW and that of Exp3 further increases. Specifically, take the case when using counter-example 3 and $T = \sqrt{10} \times 10^6$ as an example. When M increases from 2 to 3, the regret of ROW further decreases from about 3.22×10^3 to 945.85.

In Fig. 3, the regret of ROCF (for $M = 1$) is also smaller than that of Exp3. This is because the choice of $\beta_1 = 1$ and $\beta_2 = 1$ for Fig. 3 satisfies $\beta_2 \leq \frac{3}{4}K\beta_1$. As we show in (3) and (23), this is the range where costly full-feedback is helpful for reducing the regret when $M = 1$. In Fig. 4, we present a different set of results when $\beta_1 = 0.1$. (The other parameters are the same as in Fig. 3.) When β_1 decreases to 0.1, i.e., $\beta_2 > \frac{3}{4}K\beta_1$, using costly full-feedback is no longer that helpful for $M = 1$. Thus, we can observe that the gap between the regret of ROCF and that of Exp3 diminishes. In contrast, since ROW does not use full feedback, it still shows significant reduction in regret compared with other algorithms.

6 CONCLUSION

In this paper, we investigate bandit-learning problems with switching costs and full-feedback costs. First, when only $M = 1$ arm is pulled at each time, we provide a lower bound (and a matching upper bound) of the regret. Our new bounds show that adding costly full-feedback will not alter the $\Theta(T^{\frac{2}{3}})$ regret for $M = 1$, while the dependence on K could be improved when the full-feedback cost β_2 is small. Second, when $M \geq 2$ arms can be chosen at each time, we provide a novel online learning algorithm ROW that improves the regret to $O(\sqrt{T})$ without even using full feedback. Our result thus reveals that having 2 (or more) arms is surprisingly as powerful as having free full-feedback, for obtaining a low regret in bandit-learning problems with switching costs. Our algorithm ROW and regret analysis involve several new ideas, e.g., using different weight-decay parameters inside and across episodes. Our numerical results confirm that the regret of our algorithm ROW is much smaller than that of the episodic version of Exp3.

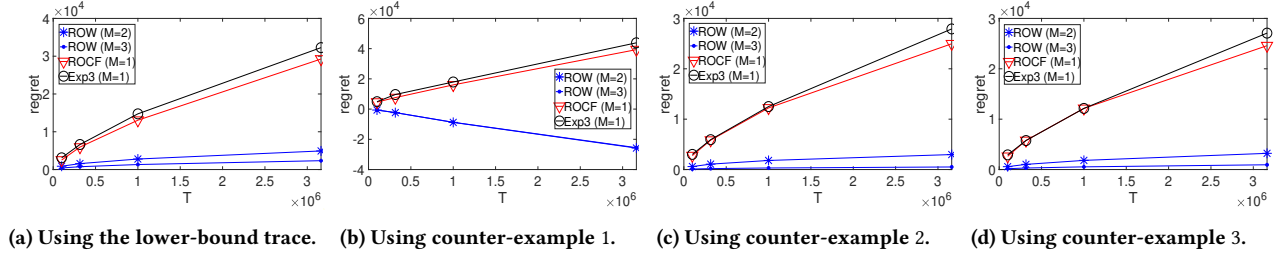
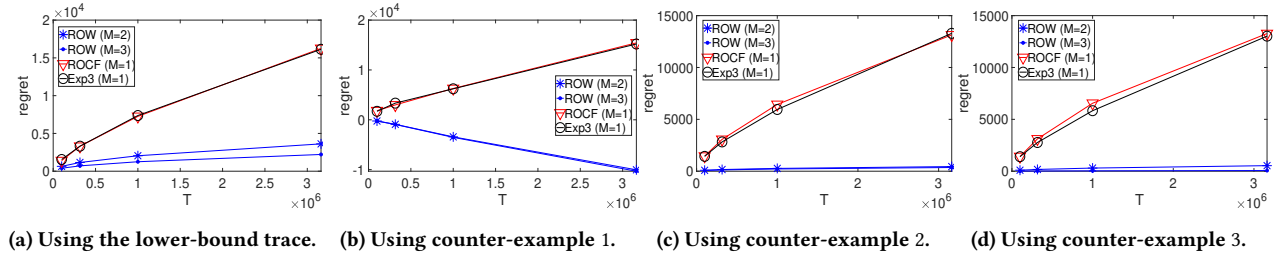
There are several interesting directions of future work. First, notice that we study the static regret. It would be interesting to extend our study to the dynamic regret, where the optimal arm changes in time. Second, ROW assumes the knowledge of the time length T . It would be useful to extend ROW to the setting where T is not known in advance.

ACKNOWLEDGMENTS

This work has been partially supported by NSF grants CNS-2113893 and CNS-2047719.

REFERENCES

- [1] Raman Arora, Ofer Dekel, and Ambuj Tewari. 2012. Online bandit learning against an adaptive adversary: from regret to policy regret. In *Proceedings of 29th International Conference on Machine Learning*. 1747–1754.
- [2] Raman Arora, Teodor Vanislavov Marinov, and Mehryar Mohri. 2019. Bandits with feedback graphs and switching costs. In *Advances in Neural Information Processing Systems*. 10397–10407.
- [3] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. 2002. The nonstochastic multiarmed bandit problem. *SIAM journal on computing* 32, 1 (2002), 48–77.
- [4] Dirk Bergemann and Juuso Välimäki. 2006. Bandit problems. *Cowles Foundation discussion paper* (2006).
- [5] Avrim Blum and Yishay Mansour. 2007. Learning, regret minimization, and equilibria. *Algorithmic Game Theory* (2007).
- [6] A. Borodin and R. El-Yaniv. 1997. On randomization in online computation. In *Proceedings of Computational Complexity, Twelfth Annual IEEE Conference*. 226–238. <https://doi.org/10.1109/CCC.1997.612318>
- [7] Jiayi Chen and Xukan Ran. 2019. Deep learning with edge computing: A review. *Proc. IEEE* 107, 8 (2019), 1655–1674.

Figure 3: Compare the regrets of ROW, ROCF and the episodic version of Exp3. ($\beta_1 = 1, \beta_2 = 1$.)Figure 4: Compare the regrets of ROW, ROCF and the episodic version of Exp3. ($\beta_1 = 0.1, \beta_2 = 1$.)

- [8] Richard Combes, Mohammad Sadegh Talebi Mazraeh Shahi, and Alexandre Proutiere. 2015. Combinatorial bandits revisited. In *Advances in Neural Information Processing Systems*. 2116–2124.
- [9] Thomas M Cover. 1999. *Elements of information theory*. John Wiley & Sons.
- [10] Ofer Dekel, Jian Ding, Tomer Koren, and Yuval Peres. 2014. Bandits with switching costs: $T^{2/3}$ regret. In *Proceedings of 46th annual ACM symposium on Theory of computing*. 459–467.
- [11] Pedro Domingos. 2000. A unified bias-variance decomposition. In *Proceedings of 17th International Conference on Machine Learning*. 231–238.
- [12] Rick Durrett. 2019. *Probability: theory and examples*. Cambridge university press.
- [13] João Gama, Indrė Žliobaitė, Albert Bifet, Mykola Pechenizkiy, and Abdelhamid Bouchachia. 2014. A survey on concept drift adaptation. *ACM computing surveys* 46, 4 (2014), 1–37.
- [14] Sascha Guelen, Berthold Vöcking, and Melanie Winkler. 2010. Regret minimization for online buffering problems using the weighted majority algorithm. In *Conference on Learning Theory*. Citeseer, 132–143.
- [15] Sudipto Guha and Kamesh Munagala. 2009. Multi-armed bandits with metric switching costs. In *International Colloquium on Automata, Languages, and Programming*. Springer, 496–507.
- [16] Michael Mitzenmacher. 2001. The power of two choices in randomized load balancing. *IEEE Transactions on Parallel and Distributed Systems* 12, 10 (2001), 1094–1104.
- [17] J v Neumann. 1928. Zur theorie der gesellschaftsspiele. *Mathematische annalen* 100, 1 (1928), 295–320.
- [18] Yevgeny Seldin, Peter Bartlett, Koby Crammer, and Yasin Abbasi-Yadkori. 2014. Prediction with limited advice and multiarmed bandits with paid observations. In *Proceedings of 31st International Conference on Machine Learning*. 280–287.
- [19] Shai Shalev-Shwartz. 2011. Online learning and online convex optimization. *Foundations and trends in Machine Learning* 4, 2 (2011), 107–194.
- [20] Ming Shi, Xiaojun Lin, and Sonia Fahmy. 2021. Competitive Online Convex Optimization With Switching Costs and Ramp Constraints. *IEEE/ACM Transactions on Networking* 29, 2 (2021), 876–889. <https://doi.org/10.1109/TNET.2021.3053910>
- [21] Aleksandrs Slivkins. 2019. Introduction to multi-armed bandits. *Foundations and Trends® in Machine Learning* 12, 1-2 (2019), 1–286.
- [22] Andrew Chi-Chin Yao. 1977. Probabilistic computations: Toward a unified measure of complexity. In *18th Annual Symposium on Foundations of Computer Science*. IEEE Computer Society, 222–227.
- [23] Zhi Zhou, Xu Chen, En Li, Liekang Zeng, Ke Luo, and Junshan Zhang. 2019. Edge intelligence: Paving the last mile of artificial intelligence with edge computing. *Proc. IEEE* 107, 8 (2019), 1738–1762.

A PROOF OF LEMMA 3.2

First of all, according to Yao’s principle [6, 17, 22], we have

$$\begin{aligned} R^\pi(T) &\geq \mathbb{E}_{I_{1:K}(1:T)} \left[\mathbb{E}_\pi \left[\text{Cost}^\pi(1:T) \right] - \text{Cost}^{\text{OPT}}(1:T) \right] \\ &= \mathbb{E}_\pi \left[\mathbb{E}_{I_{1:K}(1:T)} \left[\text{Cost}^\pi(1:T) - \text{Cost}^{\text{OPT}}(1:T) \right] \right] \\ &\geq \min_{\pi'} \left\{ \mathbb{E}_{I_{1:K}(1:T)} \left[\text{Cost}^{\pi'}(1:T) - \text{Cost}^{\text{OPT}}(1:T) \right] \right\}, \end{aligned} \quad (48)$$

i.e., the worst-case expected regret $R^\pi(T)$ of a randomized online algorithm π against an oblivious adversary is lower-bounded by the expected regret of the best deterministic online algorithm π' against a randomized adversary. Thus, the regret lower-bound that the MHM adversary provides is a lower bound of $R^\pi(T)$.

In the following, we prove Lemma 3.2, the lower bound provided by the MHM adversary. We first provide some notations in Sec. A.1, followed by some supporting lemmas in Sec. A.2 that will be used in the final proof in Sec A.3.

A.1 Notations

In the following, we use $\mathcal{P}_{k^*}(\cdot)$ to denote the conditional probability-measure given the optimal arm k^* , i.e.,

$$\mathcal{P}_{k^*}(\cdot) \triangleq \Pr \{ \cdot | k^* \}, k^* = 1, \dots, K. \quad (49)$$

We use $\mathcal{P}_0(\cdot)$ to denote an auxiliary conditional probability-measure given that there is no optimal arm (i.e., the expected losses of all arms are the same), i.e.,

$$\mathcal{P}_0(\cdot) \triangleq \Pr \{ \cdot | k^* = 0 \}. \quad (50)$$

We use $l^{\text{ob}}(t)$ to denote the observed loss by the online algorithm π , i.e.,

$$l^{\text{ob}}(t) \triangleq \begin{cases} l_{k(t)}(t), & \text{if } z(t) = 0, \\ l_{1:K}(t), & \text{if } z(t) = 1. \end{cases} \quad (51)$$

We use $l^{\text{hi}}(t)$ to denote the hidden loss of the pulled arm $k(t)$, i.e.,

$$l^{\text{hi}}(t) \triangleq l_{k(t)}(t) - Y_{k(t)}(t). \quad (52)$$

A.2 Intermediate Steps

In this subsection, we relate the expected cost-difference (between π and OPT) to the total variation distance (between $\mathcal{P}_0(l^{\text{ob}}(1:T))$ and $\mathcal{P}_{k^*}(l^{\text{ob}}(1:T))$).

First, we characterize the KL divergence between $\mathcal{P}_0(l^{\text{ob}}(1:T))$ and $\mathcal{P}_{k^*}(l^{\text{ob}}(1:T))$ in Lemma A.1 below. We use $D_{\text{KL}}(Q_1(\cdot)\|Q_2(\cdot))$ to denote the Kullback-Leibler (KL) divergence (i.e., relative entropy) between two probability measures $Q_1(\cdot)$ and $Q_2(\cdot)$, i.e.,

$$D_{\text{KL}}(Q_1(\cdot)\|Q_2(\cdot)) \triangleq \mathbb{E}_{Q_1} \left[\ln \left(\frac{Q_1(\cdot)}{Q_2(\cdot)} \right) \right].$$

LEMMA A.1. *The KL divergence between the probability measure $\mathcal{P}_0(l^{\text{ob}}(1:T))$ and $\mathcal{P}_{k^*}(l^{\text{ob}}(1:T))$ of the entire observed loss-sequence $l^{\text{ob}}(1:T)$ is upper-bounded as follows: for $T \geq 2$,*

$$D_{\text{KL}} \left(\mathcal{P}_0(l^{\text{ob}}(1:T)) \middle\| \mathcal{P}_{k^*}(l^{\text{ob}}(1:T)) \right) \leq \frac{\log_2 T \cdot \epsilon^2}{2\sigma^2} \cdot \left\{ 2\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}_{k^*}^s] + 2\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}} \neq k^*] + K\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}} = k^*] \right\}, \quad (53)$$

where $\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}_{k^*}^s]$ denotes the expected number (under the probability measure \mathcal{P}_0) of times that the algorithm switches from or to the optimal arm k^* , $\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}} \neq k^*]$ denotes the expected number (under the probability measure \mathcal{P}_0) of times the algorithm asks for costly full-feedback when the optimal arm k^* is not pulled, and $\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}} = k^*]$ denotes the expected number (under the probability measure \mathcal{P}_0) of times the algorithm asks for costly full-feedback when the optimal arm k^* is pulled.

Please see Appendix B for the complete proof of Lemma A.1. From Lemma A.1, we can then characterize the total variation distance between $\mathcal{P}_0(l^{\text{ob}}(1:T))$ and $\mathcal{P}_{k^*}(l^{\text{ob}}(1:T))$, averaged over all k^* , in Lemma A.2 below. We use $D_{\text{TV}}(Q_1(\cdot)\|Q_2(\cdot))$ to denote the total variation distance between two probability measures $Q_1(\cdot)$ and $Q_2(\cdot)$,

$$D_{\text{TV}}(Q_1(\cdot)\|Q_2(\cdot)) \triangleq \sup_{A \in \mathcal{F}} |Q_1(A) - Q_2(A)|,$$

where \mathcal{F} denotes the σ -algebra of the sample space.

LEMMA A.2. *The average total-variation-distance between the probability measure $\mathcal{P}_0(l^{\text{ob}}(1:T))$ and $\mathcal{P}_{k^*}(l^{\text{ob}}(1:T))$ of the entire observed loss-sequence $l^{\text{ob}}(1:T)$ is upper-bounded as follows: for $T \geq 2$,*

$$\begin{aligned} & \frac{1}{K} \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0(l^{\text{ob}}(1:T)) \middle\| \mathcal{P}_{k^*}(l^{\text{ob}}(1:T)) \right) \\ & \leq \frac{\sqrt{\ln 2} \cdot \epsilon}{2\sigma} \cdot \sqrt{\log_2 T} \cdot \sqrt{\frac{4}{K} \cdot \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^s] + 3 \cdot \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}}]}, \quad (54) \end{aligned}$$

where $\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^s]$ denotes the expected number (under the probability measure \mathcal{P}_0) of times that the algorithm switches, and $\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}}]$ denotes the expected number (under the probability measure \mathcal{P}_0) of times that the algorithm asks for costly full-feedback.

Please see Appendix C for the complete proof of Lemma A.2. The above bound on total variation distance then allows us to lower-bound the regret in Lemma A.3 below.

LEMMA A.3. *The expected regret of any deterministic online algorithm π is lower-bounded as follows,*

$$\begin{aligned} & \mathbb{E} \left[\text{Cost}^\pi(1:T) - \text{Cost}^{\text{OPT}}(1:T) \right] \\ & \geq \frac{\epsilon T}{2} - \frac{\epsilon T}{K} \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0(l^{\text{ob}}(1:T)) \middle\| \mathcal{P}_{k^*}(l^{\text{ob}}(1:T)) \right) \\ & \quad + \beta_1 \mathbb{E} [\mathbf{N}^s] + \beta_2 \mathbb{E} [\mathbf{N}^{\text{ck}}], \quad (55) \end{aligned}$$

where the expectation \mathbb{E} is with respect to both $\mathcal{P}_{k^*}(\cdot)$ and the randomness of choosing the optimal arm k^* .

Please see Appendix D for the complete proof of Lemma A.3. In Lemma A.4 below, we take care of the possibility that adversary inputs $l_k(t)$ generated by MHM may exceed the admitted range $[0, 1]$.

LEMMA A.4. *Let $l'_k(t)$ denote the clipped loss of $l_k(t)$, i.e.,*

$$l'_k(t) = \min\{\max\{l_k(t), 0\}, 1\}.$$

Next, we use Reg' to denote the regret of the decision sequence $k(1:T)$ made by the online algorithm under the clipped loss $l'_k(t)$, i.e.,

$$\text{Reg}' \triangleq \sum_{t=1}^T l'_{k(t)}(t) + \beta_1 \mathbf{N}^s - \sum_{k^*=1}^K l'_{k^*}(t) - \beta_1.$$

Similarly, we use Reg to denote the regret of the same decision sequence $k(1:T)$ but under the unclipped loss $l_k(t)$, i.e.,

$$\text{Reg} \triangleq \sum_{t=1}^T l_{k(t)}(t) + \beta_1 \mathbf{N}^s - \sum_{k^*=1}^K l_{k^*}(t) - \beta_1.$$

Then, we have

$$\mathbb{E}[\text{Reg}'] \geq \mathbb{E}[\text{Reg}] - \frac{\epsilon T}{6}, \quad (56)$$

where the expectation \mathbb{E} is with respect to both $\mathcal{P}_{k^*}(\cdot)$ and the randomness of choosing the optimal arm k^* .

Please see Appendix E for the complete proof of Lemma A.4.

A.3 Final Steps

PROOF. In this subsection, by using Lemma A.1-Lemma A.4, we derive the lower-bound of the regret. In the following, we first focus on analyzing the regret of any deterministic online algorithm satisfying the following two conditions for any loss sequence: (Ci) the total switching-cost is less than or equal to ϵT , (Cii) the total full-feedback-cost is less than or equal to ϵT . We will relax this assumption at the end of the proof (please see the end of this section).

First, for any deterministic online algorithm satisfying conditions (Ci) and (Cii), we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^s] - \mathbb{E} [\mathbf{N}^s] \\ &= \frac{1}{K} \sum_{k^*=1}^K \{ \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^s] - \mathbb{E}_{\mathcal{P}_{k^*}} [\mathbf{N}^s] \} \\ &\leq \frac{\epsilon T}{\beta_1 K} \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0 \left(I^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(I^{\text{ob}}(1:T) \right) \right), \quad (57) \end{aligned}$$

where the inequality holds because, under condition (Ci) above, $\mathbf{N}^s \leq \frac{\epsilon T}{\beta_1}$. In addition, we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}}] - \mathbb{E} [\mathbf{N}^{\text{ck}}] \\ &= \frac{1}{K} \sum_{k^*=1}^K \{ \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}}] - \mathbb{E}_{\mathcal{P}_{k^*}} [\mathbf{N}^{\text{ck}}] \} \\ &\leq \frac{\epsilon T}{\beta_2 K} \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0 \left(I^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(I^{\text{ob}}(1:T) \right) \right), \quad (58) \end{aligned}$$

where the inequality holds because, under condition (Cii) above, $\mathbf{N}^{\text{ck}} \leq \frac{\epsilon T}{\beta_2}$.

Next, combining (55)-(58), we have

$$\begin{aligned} & \mathbb{E} \left[\text{Cost}^\pi(1:T) - \text{Cost}^{\text{OPT}}(1:T) \mid l_k(t) \in [0, 1], \right. \\ & \quad \left. \text{for all } k \in [1, K], t \in [1, T] \right] \\ &\geq \frac{\epsilon T}{2} - \frac{\epsilon T}{K} \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0 \left(I^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(I^{\text{ob}}(1:T) \right) \right) \\ & \quad + \beta_1 \mathbb{E} [\mathbf{N}^s] + \beta_2 \mathbb{E} [\mathbf{N}^{\text{ck}}] - \frac{\epsilon T}{6} \\ &\geq \frac{\epsilon T}{3} - \frac{3\epsilon T}{K} \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0 \left(I^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(I^{\text{ob}}(1:T) \right) \right) \\ & \quad + \beta_1 \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^s] + \beta_2 \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}}], \quad (59) \end{aligned}$$

where the first inequality is because of (55) and (56), the second inequality is because of (57) and (58). Then, according to (54), we have

$$\begin{aligned} & \mathbb{E} \left[\text{Cost}^\pi(1:T) - \text{Cost}^{\text{OPT}}(1:T) \mid l_k(t) \in [0, 1], \right. \\ & \quad \left. \text{for all } k \in [1, K], t \in [1, T] \right] \\ &\geq \frac{\epsilon T}{3} - \frac{3\sqrt{\ln 2} \sqrt{\log_2 T} \cdot T \epsilon^2}{2\sigma} \sqrt{\frac{4}{K} \cdot \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^s] + 3 \cdot \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}}]} \\ & \quad + \beta_1 \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^s] + \beta_2 \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}}]. \quad (60) \end{aligned}$$

Finally, according to (48), to get the lower bound of $R^\pi(T)$, we only need to derive the minimal value of the right-hand-side of (60) over all possible values of $\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^s]$ and $\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}}]$, for a carefully chosen set of ϵ and σ . We first focus on the second to fourth terms on the right-hand-side of (60). Consider a function

$$f(x, y) = -a\sqrt{\frac{4}{K}x + 3y} + \beta_1 x + \beta_2 y,$$

where $a = \frac{3\sqrt{\ln 2} \sqrt{\log_2 T} \cdot T \epsilon^2}{2\sigma}$ and $x, y \geq 0$.

(I) If $\beta_2 \geq \frac{3}{4}K\beta_1$, we have

$$\begin{aligned} f(x, y) &= -a\sqrt{\frac{4}{K}x + 3y} + \beta_1 \frac{K}{4} \left(\frac{4}{K}x + 3y \right) + \left(\beta_2 - \frac{3K}{4}\beta_1 \right) y \\ &\stackrel{z=\frac{4}{K}x+3y}{=} -a\sqrt{z} + \beta_1 \frac{K}{4}z + \left(\beta_2 - \frac{3K}{4}\beta_1 \right) y. \end{aligned}$$

Then, $f(x, y)$ is minimized at $y = 0$. Thus, we have, if $\beta_2 \geq \frac{3}{4}K\beta_1$,

$$\begin{aligned} f(x, y) &\geq -a\sqrt{z} + \beta_1 \frac{K}{4}z \\ &\geq -\frac{a^2}{\beta_1 K}, \quad (61) \end{aligned}$$

where the last inequality becomes an equality when $z = \frac{4a^2}{\beta_1 K^2}$. Then, combining (60) and (61), we have

$$\begin{aligned} & \mathbb{E} \left[\text{Cost}^\pi(1:T) - \text{Cost}^{\text{OPT}}(1:T) \mid l_k(t) \in [0, 1], \right. \\ & \quad \left. \text{for all } k \in [1, K], t \in [1, T] \right] \\ &\geq \frac{\epsilon T}{3} - \frac{a^2}{\beta_1 K} \\ &= \frac{\epsilon T}{3} - \frac{9 \ln 2 \cdot \log_2 T \cdot T^2}{4\beta_1 K \sigma^2} \epsilon^4. \end{aligned}$$

Using our choice of ϵ and σ in (12), we have

$$\begin{aligned} & \mathbb{E} \left[\text{Cost}^\pi(1:T) - \text{Cost}^{\text{OPT}}(1:T) \mid l_k(t) \in [0, 1], \right. \\ & \quad \left. \text{for all } k \in [1, K], t \in [1, T] \right] \\ &\geq \sqrt[3]{\frac{1}{3 \ln 2}} \cdot \frac{1}{36 \log_2 T} \cdot \beta_1^{\frac{1}{3}} K^{\frac{1}{3}} T^{\frac{2}{3}}. \end{aligned}$$

(II) If $\beta_2 < \frac{3}{4}K\beta_1$, similarly, we have

$$\begin{aligned} f(x, y) &\stackrel{z=\frac{4}{K}x+3y}{=} -a\sqrt{z} + \left(\beta_1 - \frac{4}{3K}\beta_2 \right) x + \frac{\beta_2}{3}z \\ &\geq -a\sqrt{z} + \frac{\beta_2}{3}z \\ &\geq -\frac{3a^2}{4\beta_2}, \quad (62) \end{aligned}$$

where the last inequality becomes an equality when $z = \frac{9a^2}{4\beta_2^2}$. Then, combining (60) and (62), we have

$$\begin{aligned} & \mathbb{E} \left[\text{Cost}^\pi(1:T) - \text{Cost}^{\text{OPT}}(1:T) \mid l_k(t) \in [0, 1], \right. \\ & \quad \left. \text{for all } k \in [1, K], t \in [1, T] \right] \\ &\geq \frac{\epsilon T}{3} - \frac{3a^2}{4\beta_2} \\ &= \frac{\epsilon T}{3} - \frac{27 \ln 2 \cdot \log_2 T \cdot T^2}{16\beta_2 \sigma^2} \epsilon^4. \end{aligned}$$

Using our choice of ϵ and σ in (12), we have

$$\begin{aligned} & \mathbb{E} \left[\text{Cost}^\pi(1:T) - \text{Cost}^{\text{OPT}}(1:T) \middle| l_k(t) \in [0, 1], \right. \\ & \qquad \qquad \qquad \left. \text{for all } k \in [1, K], t \in [1, T] \right] \\ & \geq \sqrt[3]{\frac{4}{9 \ln 2}} \cdot \frac{1}{36 \log_2 T} \cdot \beta_2^{\frac{1}{3}} T^{\frac{2}{3}}. \end{aligned}$$

Up to this point, we have proved the lower bound of the regret for any deterministic online algorithm satisfying the two conditions: (Ci) the total switching-cost is less than or equal to ϵT , (Cii) the total full-feedback-cost is less than or equal to ϵT . For any algorithm not satisfying these two conditions, similar to the conclusion in [10], its regret can be lower-bounded by the regret of a modified version of this algorithm that satisfies these two conditions. Specifically, for any online algorithm π that violates conditions (Ci) and/or (Cii) for some loss sequence, we can construct the modified version π' of π as follows: π' follows π until the time when π has already incurred a total switching-cost or full-feedback-cost equal to ϵT . For all subsequent time-slots, π' uses a fixed decision. Let us now demonstrate the relation between the expected regret of π and that of π' . For any loss sequence, if π satisfies condition (Ci) and (Cii), the regret of π' is equal to that of π . Otherwise, notice that the expected difference between the loss of the optimal arm and that of any other arm for each time-slot is at most ϵ , conditioned on all decisions that occur before time t . Thus, the regret of π' is upper-bounded by the regret of π plus ϵT , which is further upper-bounded by twice of the regret of π (since the regret of π in this case must be no smaller than either the switching cost or feedback cost, which is at least ϵT). Combining these two cases together, we can draw the conclusion that the expected regret of π must be lower-bound by half of that of π' (the latter satisfies both condition (Ci) and condition (Cii) for any loss sequence). \square

B PROOF OF LEMMA A.1

For the convenience of the reader, we re-state Lemma A.1 below.

LEMMA A.1. *The KL divergence between the probability measure $\mathcal{P}_0(l^{ob}(1:T))$ and $\mathcal{P}_{k^*}(l^{ob}(1:T))$ of the entire observed loss-sequence $l^{ob}(1:T)$ is upper-bounded as follows: for $T \geq 2$,*

$$\begin{aligned} D_{\text{KL}} \left(\mathcal{P}_0(l^{ob}(1:T)) \middle| \mathcal{P}_{k^*}(l^{ob}(1:T)) \right) & \leq \frac{\log_2 T \cdot \epsilon^2}{2\sigma^2} \\ & \cdot \left\{ 2\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}_{k^*}^s] + 2\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{ck} \neq k^*] + K\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{ck} = k^*] \right\}, \end{aligned} \quad (63)$$

where $\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}_{k^*}^s]$ denotes the expected number (under the probability measure \mathcal{P}_0) of times that the algorithm switches from or to the optimal arm k^* , $\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{ck} \neq k^*]$ denotes the expected number (under the probability measure \mathcal{P}_0) of times the algorithm asks for costly full-feedback when the optimal arm k^* is not pulled, and $\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{ck} = k^*]$ denotes the expected number (under the probability measure \mathcal{P}_0) of times the algorithm asks for costly full-feedback when the optimal arm k^* is pulled.

We now present the proof in steps.

B.1 Initial Steps

First, we prove that (see discussion in Sec. 3 why related to the hidden loss sequence l^{hi} is essential)

$$\begin{aligned} D_{\text{KL}} \left(\mathcal{P}_0(l^{ob}(1:T)) \middle| \mathcal{P}_{k^*}(l^{ob}(1:T)) \right) \\ \leq D_{\text{KL}} \left(\mathcal{P}_0(l^{ob}(1:T) | l^{\text{hi}}(1:T)) \middle| \mathcal{P}_{k^*}(l^{ob}(1:T) | l^{\text{hi}}(1:T)) \right) \\ + D_{\text{KL}} \left(\mathcal{P}_0(l^{\text{hi}}(1:T)) \middle| \mathcal{P}_{k^*}(l^{\text{hi}}(1:T)) \right). \end{aligned} \quad (64)$$

(Please see (9) for the definition of the conditional KL divergence.) This is because (i)

$$\begin{aligned} & D_{\text{KL}} \left(\mathcal{P}_0(l^{ob}(1:T), l^{\text{hi}}(1:T)) \middle| \mathcal{P}_{k^*}(l^{ob}(1:T), l^{\text{hi}}(1:T)) \right) \\ & = \mathbb{E}_{\mathcal{P}_0} \left[\log \left(\frac{\mathcal{P}_0(l^{ob}(1:T), l^{\text{hi}}(1:T))}{\mathcal{P}_{k^*}(l^{ob}(1:T), l^{\text{hi}}(1:T))} \right) \right] \\ & = \mathbb{E}_{\mathcal{P}_0} \left[\log \left(\frac{\mathcal{P}_0(l^{ob}(1:T) | l^{\text{hi}}(1:T)) \cdot \mathcal{P}_0(l^{\text{hi}}(1:T))}{\mathcal{P}_{k^*}(l^{ob}(1:T) | l^{\text{hi}}(1:T)) \cdot \mathcal{P}_{k^*}(l^{\text{hi}}(1:T))} \right) \right] \\ & = \mathbb{E}_{\mathcal{P}_0} \left[\log \left(\frac{\mathcal{P}_0(l^{ob}(1:T) | l^{\text{hi}}(1:T))}{\mathcal{P}_{k^*}(l^{ob}(1:T) | l^{\text{hi}}(1:T))} \right) \right] \\ & \quad + \mathbb{E}_{\mathcal{P}_0} \left[\log \left(\frac{\mathcal{P}_0(l^{\text{hi}}(1:T))}{\mathcal{P}_{k^*}(l^{\text{hi}}(1:T))} \right) \right] \\ & = D_{\text{KL}} \left(\mathcal{P}_0(l^{ob}(1:T) | l^{\text{hi}}(1:T)) \middle| \mathcal{P}_{k^*}(l^{ob}(1:T) | l^{\text{hi}}(1:T)) \right) \\ & \quad + D_{\text{KL}} \left(\mathcal{P}_0(l^{\text{hi}}(1:T)) \middle| \mathcal{P}_{k^*}(l^{\text{hi}}(1:T)) \right), \end{aligned} \quad (65)$$

where the second equality is because of the linearity of the expectation, and (ii)

$$\begin{aligned} & D_{\text{KL}} \left(\mathcal{P}_0(l^{ob}(1:T), l^{\text{hi}}(1:T)) \middle| \mathcal{P}_{k^*}(l^{ob}(1:T), l^{\text{hi}}(1:T)) \right) \\ & = \mathbb{E}_{\mathcal{P}_0} \left[\log \left(\frac{\mathcal{P}_0(l^{ob}(1:T), l^{\text{hi}}(1:T))}{\mathcal{P}_{k^*}(l^{ob}(1:T), l^{\text{hi}}(1:T))} \right) \right] \\ & = \mathbb{E}_{\mathcal{P}_0} \left[\log \left(\frac{\mathcal{P}_0(l^{\text{hi}}(1:T) | l^{ob}(1:T)) \cdot \mathcal{P}_0(l^{ob}(1:T))}{\mathcal{P}_{k^*}(l^{\text{hi}}(1:T) | l^{ob}(1:T)) \cdot \mathcal{P}_{k^*}(l^{ob}(1:T))} \right) \right] \\ & = \mathbb{E}_{\mathcal{P}_0} \left[\log \left(\frac{\mathcal{P}_0(l^{\text{hi}}(1:T) | l^{ob}(1:T))}{\mathcal{P}_{k^*}(l^{\text{hi}}(1:T) | l^{ob}(1:T))} \right) \right] \\ & \quad + \mathbb{E}_{\mathcal{P}_0} \left[\log \left(\frac{\mathcal{P}_0(l^{ob}(1:T))}{\mathcal{P}_{k^*}(l^{ob}(1:T))} \right) \right] \\ & = D_{\text{KL}} \left(\mathcal{P}_0(l^{\text{hi}}(1:T) | l^{ob}(1:T)) \middle| \mathcal{P}_{k^*}(l^{\text{hi}}(1:T) | l^{ob}(1:T)) \right) \\ & \quad + D_{\text{KL}} \left(\mathcal{P}_0(l^{ob}(1:T)) \middle| \mathcal{P}_{k^*}(l^{ob}(1:T)) \right) \\ & \geq D_{\text{KL}} \left(\mathcal{P}_0(l^{ob}(1:T)) \middle| \mathcal{P}_{k^*}(l^{ob}(1:T)) \right), \end{aligned} \quad (66)$$

where the last inequality is because

$$D_{\text{KL}}\left(\mathcal{P}_0\left(l^{\text{hi}}(1:T)|l^{\text{ob}}(1:T)\right)\|\mathcal{P}_{k^*}\left(l^{\text{hi}}(1:T)|l^{\text{ob}}(1:T)\right)\right) \geq 0.$$

(This is because the KL divergence is always non-negative, i.e., $D_{\text{KL}} \geq 0$ [9, p. 26].) Combining (65) and (66), we get (64).

According to (64) we have

$$\begin{aligned} & D_{\text{KL}}\left(\mathcal{P}_0\left(l^{\text{ob}}(1:T)\right)\|\mathcal{P}_{k^*}\left(l^{\text{ob}}(1:T)\right)\right) \\ & \leq \mathbb{E}_{\mathcal{P}_0}\left[\log\left(\frac{\mathcal{P}_0\left(l^{\text{ob}}(1:T)|l^{\text{hi}}(1:T)\right)}{\mathcal{P}_{k^*}\left(l^{\text{ob}}(1:T)|l^{\text{hi}}(1:T)\right)}\right)\right] \\ & \quad + \mathbb{E}_{\mathcal{P}_0}\left[\log\left(\frac{\mathcal{P}_0\left(l^{\text{hi}}(1:T)\right)}{\mathcal{P}_{k^*}\left(l^{\text{hi}}(1:T)\right)}\right)\right] \\ & = \mathbb{E}_{\mathcal{P}_0}\left[\log\left(\prod_{t=1}^T \frac{\mathcal{P}_0\left(l^{\text{ob}}(t)|l^{\text{hi}}(t)\right)}{\mathcal{P}_{k^*}\left(l^{\text{ob}}(t)|l^{\text{hi}}(t)\right)}\right)\right] \\ & \quad + \mathbb{E}_{\mathcal{P}_0}\left[\log\left(\prod_{t=1}^T \frac{\mathcal{P}_0\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)}{\mathcal{P}_{k^*}\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)}\right)\right] \\ & = \sum_{t=1}^T \left\{ \mathbb{E}_{\mathcal{P}_0}\left[\log\left(\frac{\mathcal{P}_0\left(l^{\text{ob}}(t)|l^{\text{hi}}(t)\right)}{\mathcal{P}_{k^*}\left(l^{\text{ob}}(t)|l^{\text{hi}}(t)\right)}\right)\right] \right. \\ & \quad \left. + \mathbb{E}_{\mathcal{P}_0}\left[\log\left(\frac{\mathcal{P}_0\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)}{\mathcal{P}_{k^*}\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)}\right)\right] \right\} \\ & = \sum_{t=1}^T \left\{ D_{\text{KL}}\left(\mathcal{P}_0\left(l^{\text{ob}}(t)|l^{\text{hi}}(t)\right)\|\mathcal{P}_{k^*}\left(l^{\text{ob}}(t)|l^{\text{hi}}(t)\right)\right) \right. \\ & \quad \left. + D_{\text{KL}}\left(\mathcal{P}_0\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)\|\mathcal{P}_{k^*}\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)\right) \right\}, \end{aligned} \tag{67}$$

where the first equality is because (i) given $l^{\text{hi}}(t)$, $l^{\text{ob}}(t)$ is conditionally independent of $l^{\text{ob}}(t')$ for all $t' \neq t$, (ii) of the chain rule, and given $l^{\text{hi}}(\rho(t))$, $l^{\text{hi}}(t)$ is conditionally independent of $l^{\text{hi}}(t')$ for all $t' \neq \rho(t)$ and $t' < t$. The second equality is because of the linearity of the expectation.

B.2 Intermediate Steps

In this subsection, we calculate the conditional KL divergence

$$D_{\text{KL}}\left(\mathcal{P}_0\left(l^{\text{ob}}(t)|l^{\text{hi}}(t)\right)\|\mathcal{P}_{k^*}\left(l^{\text{ob}}(t)|l^{\text{hi}}(t)\right)\right)$$

and

$$D_{\text{KL}}\left(\mathcal{P}_0\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)\|\mathcal{P}_{k^*}\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)\right)$$

for each time t . We use $\mathcal{N}(\epsilon, \sigma^2)$ to denote the Gaussian distribution with mean ϵ and variance σ^2 . We use $\mathcal{N}_K(\vec{\mu}, \Sigma)$ to denote the multi-variate Gaussian distribution with K dimensions, mean vector equal to $\vec{\mu}$, and covariance matrix equal to Σ . We use $\vec{\mu}_a(b)$ to denote the mean vector with all entries equal to b but the k^* -th entry equal to

a . We use $\Sigma(\sigma^2)$ to denote the covariance matrix with all entries not on the diagonal equal to 0, and all entries on the diagonal equal to σ^2 . We consider two cases one-by-one: the case 1 when the online algorithm does not ask for costly full-feedback at time t and the case 2 when the online algorithm asks for costly full-feedback at time t .

Before we elaborate on case 1, we state a standard result that, for two Gaussian distributions $\mathcal{N}(\mu_1, \sigma_1^2)$ and $\mathcal{N}(\mu_2, \sigma_2^2)$, we have

$$D_{\text{KL}}\left(\mathcal{N}(\mu_1, \sigma_1^2)\|\mathcal{N}(\mu_2, \sigma_2^2)\right) = D_{\text{KL}}\left(\mathcal{N}(0, \sigma_1^2)\|\mathcal{N}(\mu_2 - \mu_1, \sigma_2^2)\right). \tag{68}$$

This is because the KL divergence between two Gaussian distributions $\mathcal{N}(\mu_1, \sigma_1^2)$ and $\mathcal{N}(\mu_2, \sigma_2^2)$ is

$$D_{\text{KL}}\left(\mathcal{N}(\mu_1, \sigma_1^2)\|\mathcal{N}(\mu_2, \sigma_2^2)\right) = \log\left(\frac{\sigma_2}{\sigma_1}\right) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}, \tag{69}$$

which depends on the relative difference between the means, but not the absolute values of the means.

(i) **Case 1:** If the online algorithm does not ask for costly full-feedback, i.e., $z(t) = 0$, the observed loss $l^{\text{ob}}(t) = l_{k(t)}(t)$ is a scalar. Then, we have

$$\begin{aligned} & D_{\text{KL}}\left(\mathcal{P}_0\left(l^{\text{ob}}(t)|l^{\text{hi}}(t)\right)\|\mathcal{P}_{k^*}\left(l^{\text{ob}}(t)|l^{\text{hi}}(t)\right)\right) \\ & = D_{\text{KL}}\left(\mathcal{N}(0, \sigma^2)\|\mathcal{N}(0, \sigma^2)\right) = 0, \end{aligned}$$

where the first equality is because of (68) and the fact that, conditioned on $l^{\text{hi}}(t)$, $l^{\text{ob}}(t)$ follows a Gaussian distribution with mean $l^{\text{hi}}(t)$ and variance σ^2 (due to the noise $y_{k(t)}(t)$) under both probability measures $\mathcal{P}_0(\cdot)$ and $\mathcal{P}_{k^*}(\cdot)$.

(i.a) If $k(t) = k(\rho(t))$, i.e., the arm pulled at time t is the same as the arm pulled at time $\rho(t)$, we have

$$\begin{aligned} & D_{\text{KL}}\left(\mathcal{P}_0\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)\|\mathcal{P}_{k^*}\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)\right) \\ & = D_{\text{KL}}\left(\mathcal{N}(0, \sigma^2)\|\mathcal{N}(0, \sigma^2)\right) = 0, \end{aligned}$$

where the first equality is because of (68) and the fact that, conditioned on $l^{\text{hi}}(\rho(t))$, $l^{\text{hi}}(t)$ follows a Gaussian distribution with mean $l^{\text{hi}}(\rho(t))$ and variance σ^2 (due to the noise $\xi(t)$) under both probability measures $\mathcal{P}_0(\cdot)$ and $\mathcal{P}_{k^*}(\cdot)$.

(i.b) If $k(t) \neq k^*$ and $k(\rho(t)) = k^*$, i.e., the arm pulled at time t is not the optimal arm k^* but the arm pulled at time $\rho(t)$ is k^* , we have,

$$\begin{aligned} & D_{\text{KL}}\left(\mathcal{P}_0\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)\|\mathcal{P}_{k^*}\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)\right) \\ & = D_{\text{KL}}\left(\mathcal{N}(0, \sigma^2)\|\mathcal{N}(\epsilon, \sigma^2)\right) = \frac{\epsilon^2}{2\sigma^2}. \end{aligned}$$

Compared with case (i.a), the difference here is that, under $\mathcal{P}_{k^*}(\cdot)$, there is an additional gap ϵ (due to the additional $-\epsilon$ when generating $l_{k^*}(\rho(t))$ in (6)) in the mean of $l^{\text{hi}}(t)$. As a result, the KL divergence is not 0 any more.

(i.c) If $k(t) = k^*$ and $k(\rho(t)) \neq k^*$, similar to case (i.b), we have

$$\begin{aligned} & D_{\text{KL}}\left(\mathcal{P}_0\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)\|\mathcal{P}_{k^*}\left(l^{\text{hi}}(t)|l^{\text{hi}}(\rho(t))\right)\right) \\ & = D_{\text{KL}}\left(\mathcal{N}(0, \sigma^2)\|\mathcal{N}(-\epsilon, \sigma^2)\right) = \frac{\epsilon^2}{2\sigma^2}. \end{aligned}$$

(i.d) If $k(t) \neq k(\rho(t))$, $k(t) \neq k^*$ and $k(\rho(t)) \neq k^*$, i.e., the arms pulled at time t and $\rho(t)$ are not the optimal arm and are different, we have

$$\begin{aligned} D_{\text{KL}} \left(\mathcal{P}_0 \left(l^{\text{hi}}(t) | l^{\text{hi}}(\rho(t)) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{hi}}(t) | l^{\text{hi}}(\rho(t)) \right) \right) \\ = D_{\text{KL}} \left(\mathcal{N}(0, \sigma^2) \parallel \mathcal{N}(0, \sigma^2) \right) = 0, \end{aligned}$$

where the first equality is because of (68) and the fact that, conditioned on $l^{\text{hi}}(\rho(t))$, $l^{\text{hi}}(t)$ follows a Gaussian distribution with mean $l^{\text{hi}}(\rho(t))$ and variance σ^2 (due to the noise $\xi(t)$) under both probability measures $\mathcal{P}_0(\cdot)$ and $\mathcal{P}_{k^*}(\cdot)$.

Having checked case 1, we now move on to case 2. Before we start, we state a standard result that, for two multi-variate Gaussian distributions $\mathcal{N}_K(\bar{\mu}_1, \Sigma_1)$ and $\mathcal{N}_K(\bar{\mu}_2, \Sigma_2)$, we have

$$D_{\text{KL}} \left(\mathcal{N}(\bar{\mu}_1, \Sigma_1) \parallel \mathcal{N}(\bar{\mu}_2, \Sigma_2) \right) = D_{\text{KL}} \left(\mathcal{N}(0, \Sigma_1) \parallel \mathcal{N}(\bar{\mu}_2 - \bar{\mu}_1, \Sigma_2) \right). \quad (70)$$

This is because the KL divergence between two multi-variate Gaussian distributions $\mathcal{N}_K(\bar{\mu}_1, \Sigma_1)$ and $\mathcal{N}_K(\bar{\mu}_2, \Sigma_2)$ is

$$\begin{aligned} D_{\text{KL}} \left(\mathcal{N}_K(\bar{\mu}_1, \Sigma_1) \parallel \mathcal{N}_K(\bar{\mu}_2, \Sigma_2) \right) \\ = \frac{1}{2} \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} - K + \text{tr} \left(\Sigma_2^{-1} \Sigma_1 \right) - (\bar{\mu}_2 - \bar{\mu}_1)^T \Sigma_2^{-1} (\bar{\mu}_2 - \bar{\mu}_1) \right], \end{aligned} \quad (71)$$

which again depends on the relative difference between the means, but not the absolute values of the means.

(ii) **Case 2:** If the online algorithm asks for costly full-feedback, i.e., $z(t) = 1$, the observed loss $l^{\text{ob}}(t) = l_{1:K}(t)$ is a K -dimension vector. Further, note that even though $l^{\text{hi}}(t)$ is defined based only on the chosen arm $k(t)$, we can immediately infer the hidden loss of all other arms. Indeed, under \mathcal{P}_0 , the hidden losses of all arms are the same. Under \mathcal{P}_{k^*} , there is only a $-\epsilon$ difference between the hidden losses of the optimal arm and that of all other arms.

(ii.a) If $k(t) \neq k^*$ and $k(\rho(t)) \neq k^*$, i.e., the arms pulled at both time t and time $\rho(t)$ are not the optimal arm, we have

$$\begin{aligned} D_{\text{KL}} \left(\mathcal{P}_0 \left(l^{\text{ob}}(t) | l^{\text{hi}}(t) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{ob}}(t) | l^{\text{hi}}(t) \right) \right) \\ = D_{\text{KL}} \left(\mathcal{N}_K \left(0, \Sigma(\sigma^2) \right) \parallel \mathcal{N}_K \left(\bar{\mu}_{-\epsilon}(0), \Sigma(\sigma^2) \right) \right) = \frac{\epsilon^2}{2\sigma^2}, \end{aligned}$$

where the first equality is because of (70) and the fact that, (I) conditioned on $l^{\text{hi}}(t)$, under the probability measure $\mathcal{P}_0(\cdot)$, $l^{\text{ob}}(t)$ follows a multi-variate Gaussian distribution with mean $\bar{\mu}_{l^{\text{hi}}(t)}(l^{\text{hi}}(t))$ and covariance matrix $\Sigma(\sigma^2)$ (due to the noise $\gamma_k(t)$), and (II) conditioned on $l^{\text{hi}}(t)$, under the probability measure $\mathcal{P}_{k^*}(\cdot)$, $l^{\text{ob}}(t)$ follows a multi-variate Gaussian distribution with mean $\bar{\mu}_{l^{\text{hi}}(t)-\epsilon}(l^{\text{hi}}(t))$ and covariance matrix $\Sigma(\sigma^2)$ (due to the noise $\gamma_k(t)$ and the additional term $-\epsilon$ in (6) for arm k^*). In addition, similar to case (i.d), we have

$$\begin{aligned} D_{\text{KL}} \left(\mathcal{P}_0 \left(l^{\text{hi}}(t) | l^{\text{hi}}(\rho(t)) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{hi}}(t) | l^{\text{hi}}(\rho(t)) \right) \right) \\ = D_{\text{KL}} \left(\mathcal{N}(0, \sigma^2) \parallel \mathcal{N}(0, \sigma^2) \right) = 0, \end{aligned}$$

(ii.b) If $k(t) = k(\rho(t)) = k^*$, i.e., the arms pulled at both time t and time $\rho(t)$ are the optimal arm, we have

$$\begin{aligned} D_{\text{KL}} \left(\mathcal{P}_0 \left(l^{\text{ob}}(t) | l^{\text{hi}}(t) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{ob}}(t) | l^{\text{hi}}(t) \right) \right) \\ = D_{\text{KL}} \left(\mathcal{N}_K \left(0, \Sigma(\sigma^2) \right) \parallel \mathcal{N}_K \left(\bar{\mu}_0(\epsilon), \Sigma(\sigma^2) \right) \right) = \frac{(K-1)\epsilon^2}{2\sigma^2}, \end{aligned}$$

where the first equality is because of (70) and the fact that, (I) conditioned on $l^{\text{hi}}(t)$, under the probability measure $\mathcal{P}_0(\cdot)$, $l^{\text{ob}}(t)$ follows a multi-variate Gaussian distribution with mean $\bar{\mu}_{l^{\text{hi}}(t)}(l^{\text{hi}}(t))$ and covariance matrix $\Sigma(\sigma^2)$ (due to the noise $\gamma_k(t)$), and (II) conditioned on $l^{\text{hi}}(t)$, under the probability measure $\mathcal{P}_{k^*}(\cdot)$, $l^{\text{ob}}(t)$ follows a multi-variate Gaussian distribution with mean $\bar{\mu}_{l^{\text{hi}}(t)}(l^{\text{hi}}(t) + \epsilon)$ and covariance matrix $\Sigma(\sigma^2)$ (due to the noise $\gamma_k(t)$ and the additional term $-\epsilon$ in (6) for arm k^*). In addition, similar to case (i.a), we have

$$\begin{aligned} D_{\text{KL}} \left(\mathcal{P}_0 \left(l^{\text{hi}}(t) | l^{\text{hi}}(\rho(t)) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{hi}}(t) | l^{\text{hi}}(\rho(t)) \right) \right) \\ = D_{\text{KL}} \left(\mathcal{N}(0, \sigma^2) \parallel \mathcal{N}(0, \sigma^2) \right) = 0. \end{aligned}$$

(ii.c) If $k(t) \neq k^*$ and $k(\rho(t)) = k^*$, similar to case (ii.a), we have

$$\begin{aligned} D_{\text{KL}} \left(\mathcal{P}_0 \left(l^{\text{ob}}(t) | l^{\text{hi}}(t) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{ob}}(t) | l^{\text{hi}}(t) \right) \right) \\ = D_{\text{KL}} \left(\mathcal{N}_K \left(0, \Sigma(\sigma^2) \right) \parallel \mathcal{N}_K \left(\bar{\mu}_{-\epsilon}(0), \Sigma(\sigma^2) \right) \right) = \frac{\epsilon^2}{2\sigma^2}. \end{aligned}$$

In addition, similar to case (i.b), we have

$$\begin{aligned} D_{\text{KL}} \left(\mathcal{P}_0 \left(l^{\text{hi}}(t) | l^{\text{hi}}(\rho(t)) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{hi}}(t) | l^{\text{hi}}(\rho(t)) \right) \right) \\ = D_{\text{KL}} \left(\mathcal{N}(0, \sigma^2) \parallel \mathcal{N}(\epsilon, \sigma^2) \right) = \frac{\epsilon^2}{2\sigma^2}. \end{aligned}$$

(ii.d) If $k(t) = k^*$ and $k(\rho(t)) \neq k^*$, similar to case (ii.b), we have

$$\begin{aligned} D_{\text{KL}} \left(\mathcal{P}_0 \left(l^{\text{ob}}(t) | l^{\text{hi}}(t) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{ob}}(t) | l^{\text{hi}}(t) \right) \right) \\ = D_{\text{KL}} \left(\mathcal{N}_K \left(0, \Sigma(\sigma^2) \right) \parallel \mathcal{N}_K \left(\bar{\mu}_0(\epsilon), \Sigma(\sigma^2) \right) \right) = \frac{(K-1)\epsilon^2}{2\sigma^2}, \end{aligned}$$

In addition, similar to case (i.c), we have

$$\begin{aligned} D_{\text{KL}} \left(\mathcal{P}_0 \left(l^{\text{hi}}(t) | l^{\text{hi}}(\rho(t)) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{hi}}(t) | l^{\text{hi}}(\rho(t)) \right) \right) \\ = D_{\text{KL}} \left(\mathcal{N}(0, \sigma^2) \parallel \mathcal{N}(-\epsilon, \sigma^2) \right) = \frac{\epsilon^2}{2\sigma^2}. \end{aligned}$$

B.3 Final Steps

Now, based on the results that we obtained in Sec. B.1 and Sec B.2, we provide the final proof of Lemma A.1.

PROOF. First, we use $I^s(t)$ to indicate whether the online algorithm switches from or to the optimal arm k^* from time $\rho(t)$ to t . That is, $I^s(t) = 1$ if the online algorithm switches from or to the optimal arm k^* from time $\rho(t)$ to t , and $I^s(t) = 0$ otherwise. Moreover, we use $I^{\text{ck}}(t)$ to indicate whether the online algorithm asks for costly full-feedback at time t . Next, combining all above

cases and (67), we have

$$\begin{aligned}
& D_{\text{KL}} \left(\mathcal{P}_0 \left(l^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{ob}}(1:T) \right) \right) \\
& \leq \sum_{t=1}^T \left[\mathcal{P}_0 \left(\mathcal{I}^{\text{ck}}(t) = 0, \mathcal{I}^{\text{s}}(t) = 1 \right) \cdot \frac{\epsilon^2}{2\sigma^2} \right. \\
& \quad \text{(due to cases (i.b) and (i.c))} \\
& \quad + \mathcal{P}_0 \left(\mathcal{I}^{\text{ck}}(t) = 1, k(t) \neq k^* \right) \cdot \frac{\epsilon^2}{\sigma^2} \\
& \quad \text{(due to cases (ii.a) and (ii.c))} \\
& \quad \left. + \mathcal{P}_0 \left(\mathcal{I}^{\text{ck}}(t) = 1, k(t) = k^* \right) \cdot \frac{K\epsilon^2}{2\sigma^2} \right] \\
& \quad \text{(due to cases (ii.b) and (ii.d))} \\
& \leq \frac{\epsilon^2}{2\sigma^2} \cdot \mathbb{E}_{\mathcal{P}_0} \left[\sum_{t=1}^T \mathbf{1}_{\{\mathcal{I}^{\text{s}}(t)=1\}} \right] \\
& \quad + \frac{\epsilon^2}{\sigma^2} \cdot \mathbb{E}_{\mathcal{P}_0} \left[\sum_{t=1}^T \mathbf{1}_{\{\mathcal{I}^{\text{ck}}(t)=1, k(t) \neq k^*\}} \right] \\
& \quad + \frac{K\epsilon^2}{2\sigma^2} \cdot \mathbb{E}_{\mathcal{P}_0} \left[\sum_{t=1}^T \mathbf{1}_{\{\mathcal{I}^{\text{ck}}(t)=1, k(t)=k^*\}} \right], \quad (72)
\end{aligned}$$

where $\mathbf{1}_E$ is an indicator function (i.e., $\mathbf{1}_E = 1$ if the event E is true, and $\mathbf{1}_E = 0$ otherwise). Recall that the event $\{\mathcal{I}^{\text{s}}(t) = 1\}$ means that there is a switch from or to k^* between time $\rho(t)$ and t . In contrast, $\mathbf{N}_{k^*}^{\text{s}}$ counts the number of switches from or to k^* between adjacent times. To relate $\sum_{t=1}^T \mathbf{1}_{\{\mathcal{I}^{\text{s}}(t)=1\}}$ to $\mathbf{N}_{k^*}^{\text{s}}$, we follow the proof in [10]. First, the event $\{\mathcal{I}^{\text{s}}(t) = 1\}$ implies that there exists at least one switch from or to k^* in some adjacent time-slots between time $\rho(t)$ and t . Next, we use $\{S_i\}_{i=1, \mathbf{N}_{k^*}^{\text{s}}}$ to denote the time-slots from 1 to T that each switching between adjacent time-slots occurs. Then, we have

$$\begin{aligned}
\sum_{t=1}^T \mathbf{1}_{\{\mathcal{I}^{\text{s}}(t)=1\}} & \leq \sum_{i=1}^{\mathbf{N}_{k^*}^{\text{s}}} \sum_{t \in [1, T]: \rho(t) < S_i \leq t} \mathbf{1}_{\{\mathcal{I}^{\text{s}}(t)=1\}} \\
& \leq \sum_{i=1}^{\mathbf{N}_{k^*}^{\text{s}}} |\{t \in [1, T] : \rho(t) < S_i \leq t\}| \\
& \leq (\log_2 T + 1) \cdot \mathbf{N}_{k^*}^{\text{s}},
\end{aligned}$$

where $|\cdot|$ denotes the cardinality of a set, and the last inequality is because of Lemma 2 in [10]. Since $T \geq 2$, we have

$$\sum_{t=1}^T \mathbf{1}_{\{\mathcal{I}^{\text{s}}(t)=1\}} \leq 2 \log_2 T \cdot \mathbf{N}_{k^*}^{\text{s}}.$$

Moreover, we have (I) $\sum_{t=1}^T \mathbf{1}_{\{\mathcal{I}^{\text{ck}}(t)=1, k(t) \neq k^*\}}$ is equal to the number of times using full-feedback \mathbf{N}^{ck} when k^* is not pulled, and (II)

$\sum_{t=1}^T \mathbf{1}_{\{\mathcal{I}^{\text{ck}}(t)=1, k(t)=k^*\}}$ is equal to the number of times using full-feedback \mathbf{N}^{ck} when k^* is pulled. Hence, we have

$$\begin{aligned}
& D_{\text{KL}} \left(\mathcal{P}_0 \left(l^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{ob}}(1:T) \right) \right) \\
& \leq \log_2 T \cdot \left\{ \frac{\epsilon^2}{\sigma^2} \cdot \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}_{k^*}^{\text{s}}] + \frac{\epsilon^2}{\sigma^2} \cdot \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}} | \neq k^*] \right. \\
& \quad \left. + \frac{K\epsilon^2}{2\sigma^2} \cdot \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}} | = k^*] \right\},
\end{aligned}$$

which concludes the proof. \square

C PROOF OF LEMMA A.2

For the convenience of the reader, we re-state Lemma A.2 below.

LEMMA A.2. *The average total-variation-distance between the probability measure $\mathcal{P}_0 \left(l^{\text{ob}}(1:T) \right)$ and $\mathcal{P}_{k^*} \left(l^{\text{ob}}(1:T) \right)$ of the entire observed loss-sequence $l^{\text{ob}}(1:T)$ is upper-bounded as follows: for $T \geq 2$,*

$$\begin{aligned}
& \frac{1}{K} \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0 \left(l^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{ob}}(1:T) \right) \right) \\
& \leq \frac{\sqrt{\ln 2} \cdot \epsilon}{2\sigma} \cdot \sqrt{\log_2 T} \cdot \sqrt{\frac{4}{K} \cdot \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{s}}] + 3 \cdot \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}}]}, \quad (73)
\end{aligned}$$

where $\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{s}}]$ denotes the expected number (under the probability measure \mathcal{P}_0) of times that the algorithm switches, and $\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}}]$ denotes the expected number (under the probability measure \mathcal{P}_0) of times that the algorithm asks for costly full-feedback.

PROOF. According to Pinsker's inequality (Theorem 12.6.1 in [9]), we have

$$\begin{aligned}
& D_{\text{TV}} \left(\mathcal{P}_0 \left(l^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{ob}}(1:T) \right) \right) \\
& \leq \sqrt{\frac{\ln 2}{2}} \cdot \sqrt{D_{\text{KL}} \left(\mathcal{P}_0 \left(l^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{ob}}(1:T) \right) \right)}.
\end{aligned}$$

Then, according to Lemma A.1, we have

$$\begin{aligned}
& D_{\text{TV}} \left(\mathcal{P}_0 \left(l^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(l^{\text{ob}}(1:T) \right) \right) \\
& \leq \frac{\sqrt{\ln 2}}{2} \cdot \sqrt{\log_2 T} \cdot \frac{\epsilon}{\sigma} \\
& \quad \cdot \sqrt{2\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}_{k^*}^{\text{s}}] + 2\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}} | \neq k^*] + K\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}} | = k^*]}.
\end{aligned}$$

Thus, according to Jensen's inequality, we have

$$\begin{aligned}
 & \frac{1}{K} \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0 \left(I^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(I^{\text{ob}}(1:T) \right) \right) \\
 & \leq \frac{\sqrt{\ln 2}}{2} \cdot \sqrt{\log_2 T} \cdot \frac{\epsilon}{\sigma} \cdot \frac{1}{K} \\
 & \quad \cdot \sum_{k^*=1}^K \sqrt{2\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}_{k^*}^s] + 2\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}} \neq k^*] + K\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}} = k^*]} \\
 & \leq \frac{\sqrt{\ln 2}}{2} \cdot \sqrt{\log_2 T} \cdot \frac{\epsilon}{\sigma} \cdot \left\{ \frac{2}{K} \sum_{k^*=1}^K \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}_{k^*}^s] \right. \\
 & \quad \left. + \frac{2}{K} \sum_{k^*=1}^K \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}} \neq k^*] + \sum_{k^*=1}^K \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}} = k^*] \right\}^{\frac{1}{2}}.
 \end{aligned}$$

Finally, we have

$$\begin{aligned}
 & \frac{1}{K} \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0 \left(I^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(I^{\text{ob}}(1:T) \right) \right) \\
 & \leq \frac{\sqrt{\ln 2}}{2} \cdot \sqrt{\log_2 T} \cdot \frac{\epsilon}{\sigma} \cdot \sqrt{\frac{4}{K} \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^s] + 3\mathbb{E}_{\mathcal{P}_0} [\mathbf{N}^{\text{ck}}]},
 \end{aligned}$$

where the inequality is because, when we sum over all k^* , (I) the event for switching from or to each k^* at time t are counted twice (i.e., when $k^* = k(t)$ and $k^* = k(t-1)$), and (II) the event for asking for costly full-feedback when $k(t) \neq k^*$ are counted $K-1$ times (i.e., when $k^* \in [1, K] - \{k(t)\}$).

□

D PROOF OF LEMMA A.3

For the convenience of the reader, we re-state Lemma A.3 below.

LEMMA A.3. *The expected regret of any deterministic online algorithm π is lower-bounded as follows,*

$$\begin{aligned}
 & \mathbb{E} \left[\text{Cost}^\pi(1:T) - \text{Cost}^{\text{OPT}}(1:T) \right] \\
 & \geq \frac{\epsilon T}{2} - \frac{\epsilon T}{K} \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0 \left(I^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(I^{\text{ob}}(1:T) \right) \right) \\
 & \quad + \beta_1 \mathbb{E} [\mathbf{N}^s] + \beta_2 \mathbb{E} [\mathbf{N}^{\text{ck}}], \tag{74}
 \end{aligned}$$

where the expectation \mathbb{E} is with respect to both $\mathcal{P}_{k^*}(\cdot)$ and the randomness of choosing the optimal arm k^* .

PROOF. First, we let \mathbf{N}_{k^*} denote the number of times that the algorithm pulls the optimal arm k^* . Then, we have

$$\begin{aligned}
 & \mathbb{E} \left[\text{Cost}^\pi(1:T) - \text{Cost}^{\text{OPT}}(1:T) \right] \\
 & = \frac{1}{K} \sum_{k=1}^K \mathbb{E} \left[\epsilon (T - \mathbf{N}_{k^*}) + \beta_1 \mathbb{E} [\mathbf{N}^s] + \beta_2 \mathbb{E} [\mathbf{N}^{\text{ck}}] \mid k^* = k \right] \\
 & = \epsilon T - \frac{\epsilon}{K} \sum_{k^*=1}^K \mathbb{E}_{\mathcal{P}_{k^*}} [\mathbf{N}_{k^*}] + \beta_1 \mathbb{E} [\mathbf{N}^s] + \beta_2 \mathbb{E} [\mathbf{N}^{\text{ck}}].
 \end{aligned}$$

Next, since $\mathbf{N}_{k^*} \leq T$, we have

$$\begin{aligned}
 & \mathbb{E}_{\mathcal{P}_{k^*}} [\mathbf{N}_{k^*}] - \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}_{k^*}] \\
 & \leq T \cdot D_{\text{TV}} \left(\mathcal{P}_0 \left(I^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(I^{\text{ob}}(1:T) \right) \right).
 \end{aligned}$$

Thus, we have

$$\begin{aligned}
 & \sum_{k^*=1}^K \mathbb{E}_{\mathcal{P}_{k^*}} [\mathbf{N}_{k^*}] \leq T \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0 \left(I^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(I^{\text{ob}}(1:T) \right) \right) \\
 & \quad + \sum_{k^*=1}^K \mathbb{E}_{\mathcal{P}_0} [\mathbf{N}_{k^*}] \\
 & = T \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0 \left(I^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(I^{\text{ob}}(1:T) \right) \right) + T.
 \end{aligned}$$

Hence, we have

$$\begin{aligned}
 & \mathbb{E} \left[\text{Cost}^\pi(1:T) - \text{Cost}^{\text{OPT}}(1:T) \right] \\
 & \geq \epsilon T - \frac{\epsilon T}{K} \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0 \left(I^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(I^{\text{ob}}(1:T) \right) \right) \\
 & \quad - \frac{\epsilon T}{K} + \beta_1 \mathbb{E} [\mathbf{N}^s] + \beta_2 \mathbb{E} [\mathbf{N}^{\text{ck}}] \\
 & \geq \frac{\epsilon T}{2} - \frac{\epsilon T}{K} \sum_{k^*=1}^K D_{\text{TV}} \left(\mathcal{P}_0 \left(I^{\text{ob}}(1:T) \right) \parallel \mathcal{P}_{k^*} \left(I^{\text{ob}}(1:T) \right) \right) \\
 & \quad + \beta_1 \mathbb{E} [\mathbf{N}^s] + \beta_2 \mathbb{E} [\mathbf{N}^{\text{ck}}],
 \end{aligned}$$

where the last inequality is because $K \geq 2$.

□

E PROOF OF LEMMA A.4

For the convenience of the reader, we re-state Lemma A.4 below.

LEMMA A.4. *Let $l'_k(t)$ denote the clipped loss of $l_k(t)$, i.e.,*

$$l'_k(t) = \min\{\max\{l_k(t), 0\}, 1\}.$$

Next, we use Reg' to denote the regret of the decision sequence $\mathbf{k}(1:T)$ made by the online algorithm under the clipped loss $l'_k(t)$, i.e.,

$$\text{Reg}' \triangleq \sum_{t=1}^T l'_{\mathbf{k}(t)}(t) + \beta_1 \mathbf{N}^s - \sum_{t=1}^T l'_{k^*}(t) - \beta_1.$$

Similarly, we use Reg to denote the regret of the same decision sequence $\mathbf{k}(1:T)$ but under the unclipped loss $l_k(t)$, i.e.,

$$\text{Reg} \triangleq \sum_{t=1}^T l_{\mathbf{k}(t)}(t) + \beta_1 \mathbf{N}^s - \sum_{t=1}^T l_{k^*}(t) - \beta_1.$$

Then, we have

$$\mathbb{E}[\text{Reg}'] \geq \mathbb{E}[\text{Reg}] - \frac{\epsilon T}{6}, \tag{75}$$

where the expectation \mathbb{E} is with respect to both $\mathcal{P}_{k^*}(\cdot)$ and the randomness of choosing the optimal arm k^* .

To prove Lemma A.4, we first prove Lemma E.1 below.

LEMMA E.1. *With a probability larger than $\frac{5}{6}$, the loss $l_{1:K}(1 : T)$ generated by the MHM adversary is in the feasible range $[0, 1]$. Specifically, for $T \geq \max\{\beta_b, 6K\}$,*

$$\Pr \{l_k(t) \in [0, 1], \text{ for all } k \in [1, K], t \in [1, T]\} \geq \frac{5}{6}. \quad (76)$$

PROOF. (Proof of Lemma E.1.)

First, we upper-bound the variance of the generated loss $l_k(t)$. Notice that the parent time in (7) is defined in a same way as that in [10]. Moreover, similar to definition 1 in [10], we define the depth of the Gaussian process $G(1 : T)$ to be

$$d_\rho(G) \triangleq \max_{t \in [1, T]} \left\{ |\bar{\rho}(t)| + |\underline{\rho}(t)| + 1 \right\},$$

i.e., the maximum number of the precedents and the descendants plus 1 (for time t itself). According to Lemma 2 in [10], the depth $d_\rho(G)$ is upper-bounded by $\lfloor \log_2 T \rfloor + 1$. Thus, the variance of $G(t)$ is upper-bound by $(\lfloor \log_2 T \rfloor + 1) \cdot \sigma^2$. Remember that MHM adds a new Gaussian noise $\gamma_k(t)$ with σ^2 variance. Therefore, the variance of $G(t) + \gamma_k(t)$ is upper-bounded by $(\lfloor \log_2 T \rfloor + 2) \cdot \sigma^2$, which is less than or equal to $2 \log_2 T \cdot \sigma^2$ when $T \geq 6K$.

Next, we can lower-bound the probability in (76). Since a standard Gaussian variable x satisfies that $\Pr \{|x| \geq y\} \leq e^{-\frac{y^2}{2}}$, we infer that

$$\begin{aligned} & \Pr \left\{ |G(t) + \gamma_k(t)| \geq \sqrt{4 \cdot 2 \log_2 T \cdot \sigma^2 \cdot \ln T} \right\} \\ &= \Pr \left\{ \left| \frac{G(t) + \gamma_k(t)}{\sqrt{2 \log_2 T \cdot \sigma^2}} \right| \geq \sqrt{4 \ln T} \right\} \\ &\leq e^{-2 \ln T} \\ &= \frac{1}{T^2}. \end{aligned}$$

According to the union bound, we have that for all $T \geq 6K$,

$$\begin{aligned} & \Pr \left\{ \max_{k \in [1, K]} \max_{t \in [1, T]} |G(t) + \gamma_k(t)| \leq \sqrt{4 \cdot 2 \log_2 T \cdot \sigma^2 \cdot \ln T} \right\} \\ &\geq 1 - \frac{K}{T} \geq \frac{5}{6}. \end{aligned}$$

Moreover, according to (12) that $\sigma = \frac{1}{9 \log_2 T}$, we have

$$\sqrt{4 \cdot 2 \log_2 T \cdot \sigma^2 \cdot \ln T} \leq 3 \sigma \log_2 T = \frac{1}{3}.$$

Therefore, we have that for all $T \geq 6K$,

$$\Pr \left\{ \max_{k \in [1, K]} \max_{t \in [1, T]} |G(t) + \gamma_k(t)| \leq \frac{1}{3} \right\} \geq \frac{5}{6}.$$

Hence, we have

$$\begin{aligned} & \Pr \left\{ G(t) + \frac{1}{2} + \gamma_k(t) \in \left[\frac{1}{6}, \frac{5}{6} \right], \text{ for all } k \in [1, K], t \in [1, T] \right\} \\ &\geq \frac{5}{6}. \end{aligned} \quad (77)$$

Finally, according to (12), we have $\epsilon \leq \frac{1}{6}$ for $T \geq \max\{\beta_b, 6K\}$. Lemma E.1 then follows. \square

PROOF. (Proof of Lemma A.4.)

We follow the arguments in [10]. We use Ω to denote the event $\{l_k(t) \in [0, 1], \text{ for all } k \in [1, K], t \in [1, T]\}$. If Ω occurs, we have $\text{Reg} = \text{Reg}'$. If Ω does not occur, note that the expected difference between the loss of the optimal arm and that of any other arm at any time t is at most ϵ , conditioned on all decisions that occurs before time t . Thus, we have

$$\mathbb{E} [\text{Reg} - \text{Reg}' | \neg \Omega] \leq \epsilon T.$$

Hence, we have

$$\mathbb{E} [\text{Reg}] - \mathbb{E} [\text{Reg}'] = \mathbb{E} [\text{Reg} - \text{Reg}' | \neg \Omega] \cdot \Pr(\neg \Omega) \leq \frac{\epsilon T}{6}, \quad (78)$$

which concludes the proof. \square

F PROOF OF LEMMA 3.3

PROOF. To prove Lemma 3.3, we consider a static solution OPT' using the single arm in $\hat{k}^*(n)$ for all time $t = 1, \dots, T$. Since the length of each episode is $\frac{T}{n}$ and $\hat{k}^*(n)$ is the optimal arm all the time, the subroutine $\Psi(K, \hat{k}^*, \frac{T}{n})$ produces a regret lower-bound equal to $\Omega \left(f \left(\log_2 \frac{T}{n} \right) \left(\frac{T}{n} \right)^\zeta \right)$ for each episode. Moreover, since the total cost of the optimal static solution OPT must be smaller than or equal to that of OPT' , we have,

$$\begin{aligned} R^\pi(T) &\geq n \cdot \Omega \left(f \left(\log_2 \frac{T}{n} \right) \left(\frac{T}{n} \right)^\zeta \right) \\ &= \Omega \left(f \left(\log_2 \frac{T}{n} \right) n^{1-\zeta} T^\zeta \right) \\ &= \Omega \left(f \left(\log_2 \frac{T}{\log_2 K} \right) (\log_2 K)^{1-\zeta} T^\zeta \right), \end{aligned} \quad (79)$$

where the last step is because $n = \log_2 K$. \square

G PROOF OF THEOREM 3.4

PROOF. First, we prove that in the first time-slot of the u -th ($u = 1, \dots, U$) episode, for both cases (i.e., $\beta_2 \geq \frac{3}{4}K\beta_1$ and $\beta_2 < \frac{3}{4}K\beta_1$), the probability of picking each arm k is

$$\Pr \{k^{\text{ROCF}}[u] = k\} = p_k^{\text{ROCF}}[u]. \quad (80)$$

(i) When $\beta_2 \geq \frac{3}{4}K\beta_1$, (80) is trivially true.

(ii) When $\beta_2 < \frac{3}{4}K\beta_1$, we prove (80) by mathematical induction (similar to the argument in [14]).

Base case: (80) is obviously true for $u = 1$. Next, we assume the

induction hypothesis that (80) is true for $u = u_0$. Then, we have

$$\begin{aligned}
 & Pr \left\{ k^{\text{ROCF}}[u_0 + 1] = k \right\} \\
 &= Pr \left\{ k^{\text{ROCF}}[u_0] = k \right\} \cdot p^{\text{ns}}[u_0 + 1] \\
 &\quad + \sum_{k'=1}^K Pr \left\{ k^{\text{ROCF}}[u_0] = k' \right\} (1 - p^{\text{ns}}[u_0 + 1]) \cdot p_k^{\text{ROCF}}[u_0 + 1] \\
 &= p_k^{\text{ROCF}}[u_0] \cdot \frac{w_k^{\text{ROCF}}[u_0 + 1]}{w_k^{\text{ROCF}}[u_0]} \\
 &\quad + \sum_{k'=1}^K p_{k'}^{\text{ROCF}}[u_0] \left(1 - \frac{w_{k'}^{\text{ROCF}}[u_0 + 1]}{w_{k'}^{\text{ROCF}}[u_0]} \right) \cdot p_k^{\text{ROCF}}[u_0 + 1].
 \end{aligned}$$

Then, according to the definition of the probability $p_k^{\text{ROCF}}[u]$ in (20), we have

$$\begin{aligned}
 & Pr \left\{ k^{\text{ROCF}}[u_0 + 1] = k \right\} \\
 &= \frac{w_k^{\text{ROCF}}[u_0]}{\sum_{k=1}^K w_k^{\text{ROCF}}[u_0]} \cdot \frac{w_k^{\text{ROCF}}[u_0 + 1]}{w_k^{\text{ROCF}}[u_0]} + \sum_{k'=1}^K \frac{w_{k'}^{\text{ROCF}}[u_0]}{\sum_{k=1}^K w_k^{\text{ROCF}}[u_0]} \\
 &\quad \cdot \left(1 - \frac{w_{k'}^{\text{ROCF}}[u_0 + 1]}{w_{k'}^{\text{ROCF}}[u_0]} \right) \cdot \frac{w_k^{\text{ROCF}}[u_0 + 1]}{\sum_{k=1}^K w_k^{\text{ROCF}}[u_0 + 1]} \\
 &= \frac{w_k^{\text{ROCF}}[u_0 + 1]}{\sum_{k=1}^K w_k^{\text{ROCF}}[u_0]} \\
 &\quad + \sum_{k'=1}^K \frac{w_{k'}^{\text{ROCF}}[u_0] - w_{k'}^{\text{ROCF}}[u_0 + 1]}{\sum_{k=1}^K w_k^{\text{ROCF}}[u_0]} \cdot \frac{w_k^{\text{ROCF}}[u_0 + 1]}{\sum_{k=1}^K w_k^{\text{ROCF}}[u_0 + 1]}.
 \end{aligned}$$

Noting that

$$\frac{w_k^{\text{ROCF}}[u_0 + 1]}{\sum_{k=1}^K w_k^{\text{ROCF}}[u_0]} = \sum_{k'=1}^K \frac{w_{k'}^{\text{ROCF}}[u_0 + 1]}{\sum_{k=1}^K w_k^{\text{ROCF}}[u_0]} \cdot \frac{w_k^{\text{ROCF}}[u_0 + 1]}{\sum_{k=1}^K w_k^{\text{ROCF}}[u_0 + 1]},$$

we thus have

$$\begin{aligned}
 Pr \left\{ k^{\text{ROCF}}[u_0 + 1] = k \right\} &= \frac{w_k^{\text{ROCF}}[u_0 + 1]}{\sum_{k=1}^K w_k^{\text{ROCF}}[u_0 + 1]} \\
 &= p_k^{\text{ROCF}}[u_0 + 1],
 \end{aligned}$$

where the last equality is because of (20).

Now, we can calculate the regret attained by ROCF for both cases.

(i) If $\beta_2 \geq \frac{3}{4}K\beta_1$, according to Exp3 analysis [1], we have that,

$$R^{\text{ROCF}}(T) \leq \frac{\ln K}{\eta} + \frac{1}{2}\eta K \tau T + \beta_1 \frac{T}{\tau}.$$

According to (22), we have $\tau = \left\lceil \left(\frac{3}{4}\beta_1 \right)^{\frac{2}{3}} (K \ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} \right\rceil$. Thus, we have

$$\begin{aligned}
 R^{\text{ROCF}}(T) &\leq \frac{\ln K}{\eta} + \frac{1}{2}\eta K \left[\left(\frac{3}{4}\beta_1 \right)^{\frac{2}{3}} (K \ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} \right] T \\
 &\quad + \beta_1 \frac{T}{\left[\left(\frac{3}{4}\beta_1 \right)^{\frac{2}{3}} (K \ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} \right]} \\
 &\leq \frac{\ln K}{\eta} + \frac{1}{2}\eta K \left(\frac{3}{4}\beta_1 \right)^{\frac{2}{3}} (K \ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} T \\
 &\quad + \beta_1 \frac{T}{\left(\frac{3}{4}\beta_1 \right)^{\frac{2}{3}} (K \ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} - 1}. \tag{81}
 \end{aligned}$$

We consider the first two terms and the last term on the right-hand-side of (81) one-by-one. For the first two terms, according to (22)

that $\eta = \left(\frac{3K}{4}\beta_1 \right)^{-\frac{1}{3}} (\ln K)^{\frac{2}{3}} T^{-\frac{2}{3}}$, we have

$$\frac{\ln K}{\eta} + \frac{1}{2}\eta K \left(\frac{3}{4}\beta_1 \right)^{\frac{2}{3}} (K \ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} T = \frac{3}{2} \sqrt[3]{\frac{3}{4}} \beta_1^{\frac{1}{3}} (K \ln K)^{\frac{1}{3}} T^{\frac{2}{3}}. \tag{82}$$

For the last term, since $\left(\frac{3}{4}\beta_1 \right)^{\frac{2}{3}} (K \ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} \geq 2$ when $T \geq \frac{128K \ln K}{9\beta_1^2}$, we have

$$\left(\frac{3}{4}\beta_1 \right)^{\frac{2}{3}} (K \ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} - 1 \geq \frac{1}{2} \left(\frac{3}{4}\beta_1 \right)^{\frac{2}{3}} (K \ln K)^{-\frac{1}{3}} T^{\frac{1}{3}}. \tag{83}$$

Combining (81)-(83), we have

$$R^{\text{ROCF}}(T) \leq \left(\frac{3}{2} \sqrt[3]{\frac{3}{4}} + 2\sqrt[3]{\frac{16}{9}} \right) \beta_1^{\frac{1}{3}} (K \ln K)^{\frac{1}{3}} T^{\frac{2}{3}}.$$

(ii) If $\beta_2 < \frac{3}{4}K\beta_1$, according to the shrinking-dartboard analysis [14], we have that

$$R^{\text{ROCF}}(T) \leq \frac{\ln K}{\eta} + \frac{1}{2}\eta \tau T + \beta_1 \left(\ln K + \eta \frac{T}{\tau} \right) + \beta_2 \frac{T}{\tau}.$$

According to (22), we have $\tau = \left\lceil \beta_2^{\frac{2}{3}} (\ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} \right\rceil$. Thus, we have

$$\begin{aligned}
 R^{\text{ROCF}}(T) &\leq \frac{\ln K}{\eta} + \frac{1}{2}\eta \left[\beta_2^{\frac{2}{3}} (\ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} \right] T \\
 &\quad + \beta_2 \frac{T}{\left[\beta_2^{\frac{2}{3}} (\ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} \right]} + \beta_1 \left(\ln K + \eta \frac{T}{\left[\beta_2^{\frac{2}{3}} (\ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} \right]} \right) \\
 &\leq \frac{\ln K}{\eta} + \frac{1}{2}\eta \beta_2^{\frac{2}{3}} (\ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} T \\
 &\quad + \beta_2 \frac{T}{\beta_2^{\frac{2}{3}} (\ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} - 1} + \beta_1 \left(\ln K + \eta \frac{T}{\beta_2^{\frac{2}{3}} (\ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} - 1} \right). \tag{84}
 \end{aligned}$$

We consider the first two terms and the last two terms on the right-hand-side of (84) one-by-one. For the first two terms, according to (22) that $\eta = \beta_2^{-\frac{1}{3}} (\ln K)^{\frac{2}{3}} T^{-\frac{2}{3}}$, we have

$$\frac{\ln K}{\eta} + \frac{1}{2}\eta\beta_2^{\frac{2}{3}} (\ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} T = \frac{3}{2}\beta_2^{\frac{1}{3}} (\ln K)^{\frac{1}{3}} T^{\frac{2}{3}}. \quad (85)$$

For the last two terms, since $\beta_2^{\frac{2}{3}} (\ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} \geq 2$ when $T \geq \frac{8 \ln K}{\beta_2^2}$, we have

$$\beta_2^{\frac{2}{3}} (\ln K)^{-\frac{1}{3}} T^{\frac{1}{3}} - 1 \geq \frac{1}{2}\beta_2^{\frac{2}{3}} (\ln K)^{-\frac{1}{3}} T^{\frac{1}{3}}. \quad (86)$$

Combining (84)-(86), we have

$$\begin{aligned} R^{\text{ROCF}}(T) &\leq \frac{3}{2}\beta_2^{\frac{1}{3}} (\ln K)^{\frac{1}{3}} T^{\frac{2}{3}} + 2\beta_2^{\frac{1}{3}} (\ln K)^{\frac{1}{3}} T^{\frac{2}{3}} \\ &\quad + \beta_1 (\ln K + 2 \ln K / \beta_2) \\ &\leq \frac{7}{2}\beta_2^{\frac{1}{3}} (\ln K)^{\frac{1}{3}} T^{\frac{2}{3}} + \beta_1 \ln K (1 + 2/\beta_2). \end{aligned}$$

□

H PROOF OF LEMMA 4.2

PROOF. To prove Lemma 4.2, we start from focus on the term on the left-hand-side of (32). Specifically, for each sub-episode (u, v) , given the history $\mathcal{H}[u-1]$ and the chosen working group $\hat{\mathbb{k}}^{\text{ROW}}[u, v]$, we have that for each time $t \in [t_{u,v}, t_{u,v} + \tau_2 - 2]$,

$$\begin{aligned} &-\frac{1}{\eta_2} \ln \left(\frac{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{w}_k^{\text{ROW}}(t+1)}{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{w}_k^{\text{ROW}}(t)} \right) \\ &= -\frac{1}{\eta_2} \ln \left(\frac{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{w}_k^{\text{ROW}}(t) e^{-\eta_2 l_k(t)}}{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{w}_k^{\text{ROW}}(t)} \right) \\ &= -\frac{1}{\eta_2} \ln \left(\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}(t) e^{-\eta_2 l_k(t)} \right), \end{aligned} \quad (87)$$

where the first equality is because of the updates of the weight $\hat{w}_k^{\text{ROW}}(t)$ in (27), and the second equality is because of the updates of the probability $\hat{p}_k^{\text{ROW}}(t)$ in (26). Next, from (87), we have

$$\begin{aligned} &-\frac{1}{\eta_2} \ln \left(\frac{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{w}_k^{\text{ROW}}(t+1)}{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{w}_k^{\text{ROW}}(t)} \right) \\ &\geq -\frac{1}{\eta_2} \ln \left(\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}(t) \left(1 - \eta_2 l_k(t) + \frac{1}{2}\eta_2^2 l_k^2(t) \right) \right) \\ &= -\frac{1}{\eta_2} \ln \left(1 - \eta_2 \cdot \sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}(t) l_k(t) \right. \\ &\quad \left. + \frac{1}{2}\eta_2^2 \cdot \sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}(t) l_k^2(t) \right), \end{aligned} \quad (88)$$

where the inequality is because $e^{-x} \leq 1 - x + \frac{1}{2}x^2$ for all $x \in [0, 1]$ and $\eta_2 l_k(t) \in [0, 1]$, and the equality is because of the re-arranging of the terms. Then, from (88), we have

$$\begin{aligned} &-\frac{1}{\eta_2} \ln \left(\frac{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{w}_k^{\text{ROW}}(t+1)}{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{w}_k^{\text{ROW}}(t)} \right) \\ &\geq -\frac{1}{\eta_2} \left(-\eta_2 \cdot \sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}(t) l_k(t) \right. \\ &\quad \left. + \frac{1}{2}\eta_2^2 \cdot \sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}(t) l_k^2(t) \right) \\ &\geq \sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}(t) l_k(t) - \frac{1}{2}\eta_2, \end{aligned} \quad (89)$$

where the first inequality is because $\ln(1-x) \leq -x$ for all x , and the second inequality is because $l_k^2(t) \leq 1$ for all k and t .

From now on, by utilizing the relation in (89), we relate the expected total loss of ROW inside each sub-episode to the log-sum-exp function $g_2[u, v]$. Recall that for $g_2[u, v]$, we define $L_k[u, v] \triangleq \sum_{t=t_{u,v}}^{t_{u,v}+\tau_2-2} l_k(t)$, and $\hat{p}_k^{\text{ROW}}[u] \triangleq \frac{w_k^{\text{ROW}}[u]}{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} w_k^{\text{ROW}}[u]}$. By moving

the term $\frac{1}{2}\eta_2$ from the right-hand-side of (89) to the left-hand-side, and then taking the sum of both sides of (89) for all time-slots $t \in [t_{u,v}, t_{u,v} + \tau_2 - 2]$, we have

$$\begin{aligned} &\sum_{t=t_{u,v}}^{t_{u,v}+\tau_2-2} \sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}(t) l_k(t) \\ &\leq -\frac{1}{\eta_2} \sum_{t=t_{u,v}}^{t_{u,v}+\tau_2-2} \ln \left(\frac{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{w}_k^{\text{ROW}}(t+1)}{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{w}_k^{\text{ROW}}(t)} \right) + \frac{1}{2}\eta_2(\tau_2 - 1) \\ &\leq -\frac{1}{\eta_2} \ln \left(\frac{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{w}_k^{\text{ROW}}(t_{u,v} + \tau_2 - 1)}{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{w}_k^{\text{ROW}}(t_{u,v})} \right) + \frac{1}{2}\eta_2\tau_2 \\ &= -\frac{1}{\eta_2} \ln \left(\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}[u] e^{-\eta_2 L_k[u,v]} \right) + \frac{1}{2}\eta_2\tau_2, \end{aligned} \quad (90)$$

where the first equality is because of the telescoping sum. The second equality of (90) is because of the update of the weight $\hat{w}_k^{\text{ROW}}(t)$ in (27) and (25), and the definition of $\hat{p}_k^{\text{ROW}}[u]$ in (31). Finally, notice that for any working group $\hat{\mathbb{k}}^{\text{ROW}}[u, v]$, we have that, at the last time-slot $t = t_{u,v} + \tau_2 - 1$ of the sub-episode,

$$\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}(t_{u,v} + \tau_2 - 1) l_k(t_{u,v} + \tau_2 - 1) \leq 1. \quad (91)$$

By combining (90) and (91), we can get (32). □

REMARK 2. Notice that when T is not divisible by τ_2 , the number of time-slots in the last sub-episode of the last episode may be smaller

than τ_2 . In this case, we only need to take the sum of both sides of (89) up to time T . As a result, the loss $L_k[u, v]$ only contains the loss $l_k(t)$ up to time T . Moreover, the term $\frac{1}{2}\eta_2\tau_2$ in (90) will be $\frac{1}{2}\eta_2 \cdot \text{mod}(T, \tau_2)$, where $\text{mod}(T, \tau_2)$ denote the remainder when T is divided by τ_2 . Then, for the last sub-episode of the last episode in this case, the term $\frac{1}{2}\eta_2\tau_2$ in (32) will also be $\frac{1}{2}\eta_2 \cdot \text{mod}(T, \tau_2)$. (This is why the sum of $\frac{1}{2}\eta_2\tau_2$ over all sub-episodes in (133) is equal to $\frac{1}{2}\eta_2T$.) However, for the convenience of elaboration, we simply use $\frac{1}{2}\eta_2\tau_2$ for both cases that T being and not being divisible by τ_2 .

I PROOF OF LEMMA 4.3

To prove Lemma 4.3, we first prove Proposition I.1 below. Lemma 4.3 then directly follows Proposition I.1.

PROPOSITION I.1. *Consider the log-sum-exp function*

$$-\frac{1}{\eta} \ln \left(\sum_{k=1}^K p_k e^{-\eta l_k} \right),$$

where $\sum_{k=1}^K p_k = 1$, and $0 \leq p_k \leq 1$ for all k . If $\eta \cdot \max_{k=1, \dots, K} l_k \leq \ln 2$, and $l_k \geq 0$ for all k , we have

$$-\frac{1}{\eta} \ln \left(\sum_{k=1}^K p_k e^{-\eta l_k} \right) \leq \mathbb{E}[l] - \frac{\eta}{8} \cdot \text{Var}(l), \quad (92)$$

where $\mathbb{E}[l] = \sum_{k=1}^K p_k l_k$ and $\text{Var}(l) = \sum_{k=1}^K p_k (l_k - \mathbb{E}[l])^2$.

PROOF. (Proof of Proposition I.1.)

First, we have

$$\begin{aligned} & -\frac{1}{\eta} \ln \left(\sum_{k=1}^K p_k e^{-\eta l_k} \right) \\ &= -\frac{1}{\eta} \ln \left[\sum_{k=1}^K p_k e^{-\eta(l_k - \mathbb{E}[l] + \mathbb{E}[l])} \right] \\ &= -\frac{1}{\eta} \ln \left[\sum_{k=1}^K p_k e^{-\eta(l_k - \mathbb{E}[l])} \cdot e^{-\eta \mathbb{E}[l]} \right] \\ &= -\frac{1}{\eta} \ln \left[\sum_{k=1}^K p_k e^{-\eta(l_k - \mathbb{E}[l])} \right] + \mathbb{E}[l]. \end{aligned} \quad (93)$$

Notice that we assume $\eta \cdot \max_{k=1, \dots, K} l_k \leq \ln 2$, and $l_k \geq 0$ for all k .

Thus, $\eta(l_k - \mathbb{E}[l]) \leq \ln 2$ and $\eta^2(l_k - \mathbb{E}[l])^2 \leq \ln 2$ for all k . Next, from (93), we have

$$\begin{aligned} & -\frac{1}{\eta} \ln \left(\sum_{k=1}^K p_k e^{-\eta l_k} \right) \\ & \leq -\frac{1}{\eta} \ln \left[\sum_{k=1}^K p_k \left(1 - \eta(l_k - \mathbb{E}[l]) + \frac{1}{4}\eta^2(l_k - \mathbb{E}[l])^2 \right) \right] + \mathbb{E}[l] \\ & = -\frac{1}{\eta} \ln \left[1 + \frac{1}{4}\eta^2 \sum_{k=1}^K p_k (l_k - \mathbb{E}[l])^2 \right] + \mathbb{E}[l], \end{aligned} \quad (94)$$

where the inequality is because (i) $e^{-x} \geq 1 - x + \frac{1}{4}x^2$ for all $x \leq \ln 2$ and (ii) $\eta(l_k - \mathbb{E}[l]) \leq \ln 2$ for all k , and the equality is because

$\mathbb{E}[l] = \sum_{k=1}^K p_k l_k$. Finally, from (94), we have

$$\begin{aligned} & -\frac{1}{\eta} \ln \left(\sum_{k=1}^K p_k e^{-\eta l_k} \right) \\ & \leq -\frac{1}{\eta} \ln \left(e^{\frac{1}{8}\eta^2 \sum_{k=1}^K p_k (l_k - \mathbb{E}[l])^2} \right) + \mathbb{E}[l] \\ & = \mathbb{E}[l] - \frac{\eta}{8} \cdot \text{Var}(l). \end{aligned}$$

where the inequality is because (i) $1 + 2x \geq e^x$ for all $x \in [0, \ln 2]$ and (ii) $\eta^2(l_k - \mathbb{E}[l])^2 \leq \ln 2$ for all k . \square

PROOF. (Proof of Lemma 4.3.)

Notice that we have $\eta_2\tau_2 \leq \ln 2$, and $L_k[u, v] \in [0, \tau_2]$ for all k . Hence, Proposition I.1 implies that Lemma 4.3 is true. \square

J PROOF OF LEMMA 4.5

To prove Lemma 4.5, we first prove Proposition J.1 below. Lemma 4.5 then directly follows Proposition J.1.

PROPOSITION J.1. *Consider the log-sum-exp function*

$$-\frac{1}{\eta} \ln \left(\sum_{k=1}^K p_k e^{-\eta l_k} \right),$$

where $\sum_{k=1}^K p_k = 1$, and $0 \leq p_k \leq 1$ for all k . If $\eta \cdot |l_k| \leq \ln 2$ for all k , we have

$$-\frac{1}{\eta} \ln \left(\sum_{k=1}^K p_k e^{-\eta l_k} \right) \geq \mathbb{E}[l] - \eta \cdot \text{Var}(l), \quad (95)$$

where $\mathbb{E}[l] = \sum_{k=1}^K p_k l_k$ and $\text{Var}(l) = \sum_{k=1}^K p_k (l_k - \mathbb{E}[l])^2$.

PROOF. (Proof of Proposition J.1.)

First, notice that Proposition I.1 and Proposition J.1 consider the same log-sum-exp function $-\frac{1}{\eta} \ln \left(\sum_{k=1}^K p_k e^{-\eta l_k} \right)$. Thus, we still have (93). However, different from Proposition I.1 that gives an upper bound of the log-sum-exp function, Proposition J.1 gives a lower bound.

Since $\eta \cdot |l_k| \leq \ln 2$ for all k , we have $\eta(l_k - \mathbb{E}[l]) \geq -2 \ln 2$. Next, from (93), we have

$$\begin{aligned} & -\frac{1}{\eta} \ln \left(\sum_{k=1}^K p_k e^{-\eta l_k} \right) \\ & \geq -\frac{1}{\eta} \ln \left[\sum_{k=1}^K p_k \left(1 - \eta(l_k - \mathbb{E}[l]) + \eta^2(l_k - \mathbb{E}[l])^2 \right) \right] + \mathbb{E}[l] \\ & = -\frac{1}{\eta} \ln \left[1 + \eta^2 \sum_{k=1}^K p_k (l_k - \mathbb{E}[l])^2 \right] + \mathbb{E}[l], \end{aligned} \quad (96)$$

where the inequality is because (i) $e^{-x} \leq 1 - x + x^2$ for all $x \geq -2 \ln 2$ and (ii) $\eta(l_k - \mathbb{E}[l]) \geq -2 \ln 2$, and the equality is because

$\mathbb{E}[L] = \sum_{k=1}^K p_k l_k$. Finally, from (96), we have

$$\begin{aligned} & -\frac{1}{\eta} \ln \left(\sum_{k=1}^K p_k e^{-\eta l_k} \right) \\ & \geq -\frac{1}{\eta} \ln \left(e^{\eta^2 \sum_{k=1}^K p_k (l_k - \mathbb{E}[L])^2} \right) + \mathbb{E}[L] \\ & = \mathbb{E}[L] - \eta \cdot \text{Var}(L), \end{aligned}$$

where the inequality is because $1 + x \leq e^x$ for all x . \square

PROOF. (Proof of Lemma 4.5.)

Notice that we have $\eta_1 \tau_1 \leq \ln 2$, and $\tilde{L}_k^{\text{ROW}}[u, v] \in [-\tau_1, \tau_1]$ for all k . Hence, Proposition J.1 implies that Lemma 4.5 is true. \square

K PROOF OF LEMMA 4.6

In Proposition K.1 below, we develop a new expression for the variance

$$\text{Var}(L) \triangleq \sum_{k=1}^K p_k \left(l_k - \sum_{k=1}^K p_k l_k \right)^2. \quad (97)$$

Proposition K.1 will be used to prove Lemma 4.6 and Lemma 4.7.

PROPOSITION K.1. *For the variance $\text{Var}(L)$ in (97), we have*

$$\text{Var}(L) = \frac{1}{2} \cdot \sum_{\substack{k_1, k_2=1, \\ k_1 \neq k_2}}^K p_{k_1} p_{k_2} (l_{k_1} - l_{k_2})^2. \quad (98)$$

PROOF. (Proof of Proposition K.1.)

Lemma K.1 is true because

$$\begin{aligned} & \sum_{\substack{k_1, k_2=1, \\ k_1 \neq k_2}}^K p_{k_1} p_{k_2} (l_{k_1} - l_{k_2})^2 \\ & = \sum_{\substack{k_1, k_2=1 \\ k_1 \neq k_2}}^K p_{k_1} p_{k_2} (l_{k_1}^2 + l_{k_2}^2) - \sum_{\substack{k_1, k_2=1 \\ k_1 \neq k_2}}^K 2p_{k_1} p_{k_2} l_{k_1} l_{k_2} \\ & = 2 \sum_{k=1}^K p_k \sum_{\substack{k'=1 \\ k' \neq k}}^K p_{k'} l_k^2 - 2 \sum_{k_1=1}^K p_{k_1} l_{k_1} \left(\sum_{\substack{k=1 \\ k \neq k_1}}^K p_k l_k \right) \\ & = 2 \sum_{k=1}^K p_k (1 - p_k) l_k^2 - 2 \sum_{k_1=1}^K p_{k_1} l_{k_1} \left(\sum_{k=1}^K p_k l_k - p_{k_1} l_{k_1} \right) \\ & = 2 \sum_{k=1}^K p_k l_k^2 - 2 \left(\sum_{k=1}^K p_k l_k \right)^2 \\ & = 2 \cdot \text{Var}(L). \end{aligned}$$

\square

We can now proceed with the proof of Lemma 4.6.

PROOF. (Proof of Lemma 4.6.)

First, according to Proposition K.1, we have

$$\begin{aligned} & \text{Var} \left(L[u, v] | \mathcal{H}[u-1], \hat{\mathbb{K}}^{\text{ROW}}[u, v] \right) \\ & = \frac{1}{2} \cdot \sum_{\substack{k_1, k_2 \in \hat{\mathbb{K}}^{\text{ROW}}[u, v], \\ k_1 \neq k_2}} \hat{p}_{k_1}^{\text{ROW}}[u] \hat{p}_{k_2}^{\text{ROW}}[u] (L_{k_1}[u, v] - L_{k_2}[u, v])^2. \end{aligned} \quad (99)$$

Next, we derive (i) the relation between $\hat{p}_k^{\text{ROW}}[u]$ and $p_k^{\text{ROW}}[u]$, and (ii) the probability of choosing each working group $\hat{\mathbb{K}}^{\text{ROW}}[u, v]$. For (i), recall from (24) that $p_k^{\text{ROW}}[u] = \frac{w_k^{\text{ROW}}[u]}{\sum_{k=1}^K w_k^{\text{ROW}}[u]}$. Moreover, recall from (31) that $\hat{p}_k^{\text{ROW}}[u] \triangleq \frac{w_k^{\text{ROW}}[u]}{\sum_{k \in \hat{\mathbb{K}}^{\text{ROW}}[u, v]} w_k^{\text{ROW}}[u]}$, which is calculated based on the chosen working-group $\hat{\mathbb{K}}^{\text{ROW}}[u, v]$. Then, we have

$$\begin{aligned} \hat{p}_k^{\text{ROW}}[u] & = \frac{w_k^{\text{ROW}}[u]}{\sum_{k \in \hat{\mathbb{K}}^{\text{ROW}}[u, v]} w_k^{\text{ROW}}[u]} \\ & = \frac{w_k^{\text{ROW}}[u]}{\sum_{k=1}^K w_k^{\text{ROW}}[u]} \cdot \frac{\sum_{k=1}^K w_k^{\text{ROW}}[u]}{\sum_{k \in \hat{\mathbb{K}}^{\text{ROW}}[u, v]} w_k^{\text{ROW}}[u]} \\ & = p_k^{\text{ROW}}[u] \cdot \frac{1}{\sum_{k \in \hat{\mathbb{K}}^{\text{ROW}}[u, v]} \frac{w_k^{\text{ROW}}[u]}{\sum_{k=1}^K w_k^{\text{ROW}}[u]}} \\ & = p_k^{\text{ROW}}[u] \cdot \frac{1}{\sum_{k \in \hat{\mathbb{K}}^{\text{ROW}}[u, v]} p_k^{\text{ROW}}[u]}, \end{aligned} \quad (100)$$

where the third equality and fourth equality are because of the update of the probability $p_k^{\text{ROW}}[u]$ in (24). In other words, the conditional probability $\hat{p}_k^{\text{ROW}}[u]$ is simply the probability $p_k^{\text{ROW}}[u]$ divided by the sum of the $p_k^{\text{ROW}}[u]$ inside the chosen working-group $\hat{\mathbb{K}}^{\text{ROW}}[u, v]$. For (ii), the probability of choosing each working group $\hat{\mathbb{K}}^{\text{ROW}}[u, v]$ is

$$\begin{aligned} & \Pr \left\{ \hat{\mathbb{K}}^{\text{ROW}}[u, v] | \mathcal{H}[u-1] \right\} \\ & = \sum_{k \in \hat{\mathbb{K}}^{\text{ROW}}[u, v]} \left[\Pr \left\{ k_0^{\text{ROW}}[u] = k | \mathcal{H}[u-1] \right\} \right. \\ & \quad \cdot \left. \Pr \left\{ \hat{\mathbb{K}}_{M-1}^{\text{ROW}}[u, v] = \hat{\mathbb{K}}^{\text{ROW}}[u, v] - \{k\} | k_0^{\text{ROW}}[u] = k, \mathcal{H}[u-1] \right\} \right] \\ & = \sum_{k \in \hat{\mathbb{K}}^{\text{ROW}}[u, v]} p_k^{\text{ROW}}[u] \cdot \frac{1}{\binom{K-1}{M-1}}, \end{aligned} \quad (101)$$

where the last equality is because of (24) and because ROW chooses the secondary arms uniformly (see Step 2 in Algorithm 4). Then,

from (99)-(101), we have

$$\begin{aligned}
 & \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,v]} \left[\text{Var} \left(L[u,v] \middle| \hat{\mathbb{k}}^{\text{ROW}}[u,v] \right) \middle| \mathcal{H}[u-1] \right] \\
 &= \sum_{\hat{\mathbb{k}}^{\text{ROW}}[u,v]} Pr \left\{ \hat{\mathbb{k}}^{\text{ROW}}[u,v] \middle| \mathcal{H}[u-1] \right\} \\
 & \quad \cdot \frac{1}{2} \cdot \sum_{\substack{k_1, k_2 \in \hat{\mathbb{k}}^{\text{ROW}}[u,v], \\ k_1 \neq k_2}} \hat{p}_{k_1}^{\text{ROW}}[u] \hat{p}_{k_2}^{\text{ROW}}[u] (L_{k_1}[u,v] - L_{k_2}[u,v])^2 \\
 &= \frac{1}{2} \sum_{\hat{\mathbb{k}}^{\text{ROW}}[u,v]} \sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} p_k^{\text{ROW}}[u] \cdot \frac{1}{\binom{K-1}{M-1}} \\
 & \quad \cdot \sum_{\substack{k_1, k_2 \\ \in \hat{\mathbb{k}}^{\text{ROW}}[u,v], \\ k_1 \neq k_2}} \frac{p_{k_1}^{\text{ROW}}[u] p_{k_2}^{\text{ROW}}[u]}{\left(\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} p_k^{\text{ROW}}[u] \right)^2} (L_{k_1}[u,v] - L_{k_2}[u,v])^2, \tag{102}
 \end{aligned}$$

where the first equality is because of (99), and the second equality is because of (100) and (101). By re-arranging the terms on the right-hand-side of (102), we have

$$\begin{aligned}
 & \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,v]} \left[\text{Var} \left(L[u,v] \middle| \hat{\mathbb{k}}^{\text{ROW}}[u,v] \right) \middle| \mathcal{H}[u-1] \right] \\
 &= \frac{1}{2 \cdot \binom{K-1}{M-1}} \sum_{\hat{\mathbb{k}}^{\text{ROW}}[u,v]} \sum_{\substack{k_1, k_2 \in \hat{\mathbb{k}}^{\text{ROW}}[u,v], \\ k_1 \neq k_2}} \frac{p_{k_1}^{\text{ROW}}[u] p_{k_2}^{\text{ROW}}[u]}{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} p_k^{\text{ROW}}[u]} \\
 & \quad \cdot (L_{k_1}[u,v] - L_{k_2}[u,v])^2. \tag{103}
 \end{aligned}$$

Finally, from (103), we have

$$\begin{aligned}
 & \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,v]} \left[\text{Var} \left(L[u,v] \middle| \hat{\mathbb{k}}^{\text{ROW}}[u,v] \right) \middle| \mathcal{H}[u-1] \right] \\
 & \geq \frac{1}{2 \cdot \binom{K-1}{M-1}} \sum_{\hat{\mathbb{k}}^{\text{ROW}}[u,v]} \sum_{\substack{k_1, k_2 \in \hat{\mathbb{k}}^{\text{ROW}}[u,v], \\ k_1 \neq k_2}} p_{k_1}^{\text{ROW}}[u] p_{k_2}^{\text{ROW}}[u] \\
 & \quad \cdot (L_{k_1}[u,v] - L_{k_2}[u,v])^2 \\
 &= \frac{\binom{K-2}{M-2}}{2 \cdot \binom{K-1}{M-1}} \sum_{\substack{k_1, k_2=1, \\ k_1 \neq k_2}}^K p_{k_1}^{\text{ROW}}[u] p_{k_2}^{\text{ROW}}[u] (L_{k_1}[u,v] - L_{k_2}[u,v])^2 \\
 &= \frac{M-1}{K-1} \cdot \text{Var} \left(L[u,v] \middle| \mathcal{H}[u-1] \right),
 \end{aligned}$$

where the inequality is because $\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} p_k^{\text{ROW}}[u] \leq 1$, the first equality is because there are $\binom{K-2}{M-2}$ working groups containing both arm k_1 and arm k_2 , and the second equality is because of Proposition K.1. \square

L PROOF OF LEMMA 4.7

PROOF. First, based on the definition of the variance, we have

$$\begin{aligned}
 & \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[\text{Var} \left(\tilde{L}^{\text{ROW}}[u] \middle| \hat{\mathbb{k}}^{\text{ROW}}[u,1:V] \right) \middle| \mathcal{H}[u-1] \right] \\
 &= \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[\sum_{k=1}^K p_k^{\text{ROW}}[u] \left(\tilde{L}_k^{\text{ROW}}[u] \right)^2 \right. \\
 & \quad \left. - \left(\sum_{k=1}^K p_k^{\text{ROW}}[u] \tilde{L}_k^{\text{ROW}}[u] \right)^2 \middle| \mathcal{H}[u-1] \right] \\
 & \leq \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[\sum_{k=1}^K p_k^{\text{ROW}}[u] \left(\tilde{L}_k^{\text{ROW}}[u] \right)^2 \middle| \mathcal{H}[u-1] \right]. \tag{104}
 \end{aligned}$$

Next, from (104), we have

$$\begin{aligned}
 & \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[\text{Var} \left(\tilde{L}^{\text{ROW}}[u] \middle| \hat{\mathbb{k}}^{\text{ROW}}[u,1:V] \right) \middle| \mathcal{H}[u-1] \right] \\
 & \leq \left(\frac{K-1}{M-1} \right)^2 \cdot \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[\sum_{k=1}^K p_k^{\text{ROW}}[u] \right. \\
 & \quad \left. \cdot \left(L_k[u, v_u(k)] - L_{k_0^{\text{ROW}}[u]}[u, v_u(k)] \right)^2 \middle| \mathcal{H}[u-1] \right] \\
 &= \left(\frac{K-1}{M-1} \right)^2 \cdot \sum_{k=1}^K \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[p_k^{\text{ROW}}[u] \right. \\
 & \quad \left. \cdot \left(L_k[u, v_u(k)] - L_{k_0^{\text{ROW}}[u]}[u, v_u(k)] \right)^2 \middle| \mathcal{H}[u-1] \right], \tag{105}
 \end{aligned}$$

where the inequality is because of the calculation of the loss difference $\tilde{L}_k^{\text{ROW}}[u]$ in (28), and the equality is because of the linearity of the expectation. Notice that the expectation $\mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]}$ is with respect to the randomness of $\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]$. Thus, it can be expanded into a sum over the randomness of the primary arm $k_0^{\text{ROW}}[u]$ and the randomness of the sub-episodes where each secondary arm k is chosen. Then, from (105), we have

$$\begin{aligned}
 & \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[\text{Var} \left(\tilde{L}^{\text{ROW}}[u] \middle| \hat{\mathbb{k}}^{\text{ROW}}[u,1:V] \right) \middle| \mathcal{H}[u-1] \right] \\
 & \leq \left(\frac{K-1}{M-1} \right)^2 \cdot \sum_{k=1}^K \sum_{k_0^{\text{ROW}}[u]=1}^K \left\{ p_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u] \right. \\
 & \quad \left. \cdot \sum_{v=1}^V \left[Pr \{v_u(k) = \{v\}\} \cdot p_k^{\text{ROW}}[u] \left(L_k[u, v] - L_{k_0^{\text{ROW}}[u]}[u, v] \right)^2 \right] \right\}.
 \end{aligned}$$

Moreover, when $k = k_0^{\text{ROW}}[u]$, we have $L_k[u, v] - L_{k_0^{\text{ROW}}[u]}[u, v] = 0$. Thus, we have

$$\begin{aligned}
& \mathbb{E}_{\tilde{k}^{\text{ROW}}[u, 1:V]} \left[\text{Var} \left(\tilde{L}^{\text{ROW}}[u] \Big|_{\tilde{k}^{\text{ROW}}[u, 1:V]} \right) \Big| \mathcal{H}[u-1] \right] \\
& \leq \left(\frac{K-1}{M-1} \right)^2 \cdot \sum_{k=1}^K \sum_{\substack{k_0^{\text{ROW}}[u]=1, \\ k_0^{\text{ROW}}[u] \neq k}}^K \left\{ p_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u] \right. \\
& \quad \cdot \sum_{v=1}^V \left[\text{Pr} \{v_u(k) = \{v\}\} \cdot p_k^{\text{ROW}}[u] \left(L_k[u, v] - L_{k_0^{\text{ROW}}[u]}[u, v] \right)^2 \right] \Big\} \\
& \leq \left(\frac{K-1}{M-1} \right) \cdot \sum_{k=1}^K \sum_{\substack{k_0^{\text{ROW}}[u]=1, \\ k_0^{\text{ROW}}[u] \neq k}}^K p_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u] \\
& \quad \cdot \sum_{v=1}^V \left[p_k^{\text{ROW}}[u] \left(L_k[u, v] - L_{k_0^{\text{ROW}}[u]}[u, v] \right)^2 \right] \\
& = \frac{2(K-1)}{M-1} \cdot \sum_{v=1}^V \text{Var} (L[u, v] | \mathcal{H}[u-1]), \tag{106}
\end{aligned}$$

where the first inequality is because the probability

$$\text{Pr} \{v_u(k) = \{v\}\} = \frac{M-1}{K-1}, \tag{107}$$

and the last equality is because of Proposition K.1. \square

REMARK 3. Notice that when $K-1$ is not divisible by $M-1$, some arm could be chosen as a secondary arm again in the V -th sub-episode. In this case, besides $v_u(k) = \{v\}$, we need to consider $v_u(k) = \{v, V\}$. Then, the first inequality in (106) becomes

$$\begin{aligned}
& \mathbb{E}_{\tilde{k}^{\text{ROW}}[u, 1:V]} \left[\text{Var} \left(\tilde{L}^{\text{ROW}}[u] \Big|_{\tilde{k}^{\text{ROW}}[u, 1:V]} \right) \Big| \mathcal{H}[u-1] \right] \\
& \leq \left(\frac{K-1}{M-1} \right)^2 \cdot \sum_{k=1}^K \sum_{\substack{k_0^{\text{ROW}}[u]=1, \\ k_0^{\text{ROW}}[u] \neq k}}^K \left\{ p_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u] \right. \\
& \quad \cdot \left[\sum_{v=1}^V \text{Pr} \{v_u(k) = \{v\}\} \cdot p_k^{\text{ROW}}[u] \left(L_k[u, v] - L_{k_0^{\text{ROW}}[u]}[u, v] \right)^2 \right. \\
& \quad \left. \left. + \sum_{v=1}^{V-1} \text{Pr} \{v_u(k) = \{v, V\}\} \cdot p_k^{\text{ROW}}[u] \left(L_k[u, v] + L_k[u, V] \right. \right. \right. \\
& \quad \left. \left. \left. - L_{k_0^{\text{ROW}}[u]}[u, v] - L_{k_0^{\text{ROW}}[u]}[u, V] \right)^2 \right] \right\}. \tag{108}
\end{aligned}$$

To get an upper bound for the right-hand-side of (108), let us first focus on the terms inside the bracket "[·]". The first term inside the

big bracket "[·]" is trivially upper-bounded by twice of itself,

$$\begin{aligned}
& \sum_{v=1}^V \text{Pr} \{v_u(k) = \{v\}\} \cdot p_k^{\text{ROW}}[u] \left(L_k[u, v] - L_{k_0^{\text{ROW}}[u]}[u, v] \right)^2 \\
& \leq 2 \sum_{v=1}^{V-1} \text{Pr} \{v_u(k) = \{v\}\} \cdot p_k^{\text{ROW}}[u] \left(L_k[u, v] - L_{k_0^{\text{ROW}}[u]}[u, v] \right)^2 \\
& \quad + 2 \cdot \text{Pr} \{v_u(k) = \{V\}\} \cdot p_k^{\text{ROW}}[u] \left(L_k[u, V] - L_{k_0^{\text{ROW}}[u]}[u, V] \right)^2. \tag{109}
\end{aligned}$$

Moreover, note that

$$\begin{aligned}
& \left(L_k[u, v] + L_k[u, V] - L_{k_0^{\text{ROW}}[u]}[u, v] - L_{k_0^{\text{ROW}}[u]}[u, V] \right)^2 \\
& \leq 2 \left(L_k[u, v] - L_{k_0^{\text{ROW}}[u]}[u, v] \right)^2 \\
& \quad + 2 \left(L_k[u, V] - L_{k_0^{\text{ROW}}[u]}[u, V] \right)^2.
\end{aligned}$$

Thus, the second term inside the big bracket "[·]" can be upper-bounded as follows,

$$\begin{aligned}
& \sum_{v=1}^{V-1} \text{Pr} \{v_u(k) = \{v, V\}\} \cdot p_k^{\text{ROW}}[u] \left(L_k[u, v] + L_k[u, V] \right. \\
& \quad \left. - L_{k_0^{\text{ROW}}[u]}[u, v] - L_{k_0^{\text{ROW}}[u]}[u, V] \right)^2 \\
& \leq 2 \sum_{v=1}^{V-1} \text{Pr} \{v_u(k) = \{v, V\}\} \cdot p_k^{\text{ROW}}[u] \left(L_k[u, v] - L_{k_0^{\text{ROW}}[u]}[u, v] \right)^2 \\
& \quad + 2 \sum_{v=1}^{V-1} \text{Pr} \{v_u(k) = \{v, V\}\} \cdot p_k^{\text{ROW}}[u] \left(L_k[u, V] - L_{k_0^{\text{ROW}}[u]}[u, V] \right)^2. \tag{110}
\end{aligned}$$

Note that the probability that the secondary arm k is chosen in the sub-episode $v = 1, \dots, V-1$ is simply $\frac{M-1}{K-1}$. Therefore, the sum of the first term on the right-hand-side of (109) and the first term on the right-hand-side of (110) is equal to

$$\begin{aligned}
& 2 \sum_{v=1}^{V-1} \left(\text{Pr} \{v_u(k) = \{v\}\} + \text{Pr} \{v_u(k) = \{v, V\}\} \right) \\
& \quad \cdot p_k^{\text{ROW}}[u] \left(L_k[u, v] - L_{k_0^{\text{ROW}}[u]}[u, v] \right)^2 \\
& = 2 \sum_{v=1}^{V-1} \frac{M-1}{K-1} \cdot p_k^{\text{ROW}}[u] \left(L_k[u, v] - L_{k_0^{\text{ROW}}[u]}[u, v] \right)^2. \tag{111}
\end{aligned}$$

Similarly, the probability that the secondary arm k is chosen in the sub-episode $v = V$ is also $\frac{M-1}{K-1}$. Therefore, the sum of the second term on the right-hand-side of (109) and the second term on the right-hand-side of (110) is equal to

$$\begin{aligned}
& 2 \left(\text{Pr} \{v_u(k) = \{V\}\} + \sum_{v=1}^{V-1} \text{Pr} \{v_u(k) = \{v, V\}\} \right) \\
& \quad \cdot p_k^{\text{ROW}}[u] \left(L_k[u, V] - L_{k_0^{\text{ROW}}[u]}[u, V] \right)^2 \\
& = 2 \cdot \frac{M-1}{K-1} \cdot p_k^{\text{ROW}}[u] \left(L_k[u, V] - L_{k_0^{\text{ROW}}[u]}[u, V] \right)^2. \tag{112}
\end{aligned}$$

Finally, combining (108)-(112), we have

$$\begin{aligned}
 & \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[\text{Var} \left(\hat{L}^{\text{ROW}}[u] \Big|_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \right) \Big| \mathcal{H}[u-1] \right] \\
 & \leq 2 \left(\frac{K-1}{M-1} \right) \cdot \sum_{k=1}^K \sum_{\substack{k_0^{\text{ROW}}[u]=1, \\ k_0^{\text{ROW}}[u] \neq k}}^K p_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u] \\
 & \quad \cdot \sum_{v=1}^V \left[p_k^{\text{ROW}}[u] \left(L_k[u,v] - L_{k_0^{\text{ROW}}[u]}[u,v] \right)^2 \right] \\
 & = \frac{4(K-1)}{M-1} \cdot \sum_{v=1}^V \text{Var}(L[u,v] | \mathcal{H}[u-1]). \tag{113}
 \end{aligned}$$

Therefore, when $K-1$ is not divisible by $M-1$, there will be an additional factor 2 in (41), i.e., (41) will become

$$\begin{aligned}
 & \sum_{v=1}^V \text{Var}(L[u,v] | \mathcal{H}[u-1]) \geq \frac{M-1}{4(K-1)} \\
 & \quad \cdot \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[\text{Var} \left(\hat{L}^{\text{ROW}}[u] \Big|_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \right) \Big| \mathcal{H}[u-1] \right]. \tag{114}
 \end{aligned}$$

This will affect the choice of the parameters, i.e., the relation between η_1 and η_2 in (36) will become $\eta_2 \geq 32 \left(\frac{K-1}{M-1} \right)^2 \cdot \eta_1$, and the constant c_2 in (46) will become $c_2 = \frac{4\sqrt{2}(K-1)}{M-1}$.

M PROOF OF LEMMA 4.4

PROOF. Our proof of Lemma 4.4 follows three steps. First, we upper-bound $\mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[g_2[u] | \mathcal{H}[u-1] \right]$. Second, we lower-bound $\mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[g_1[u] | \mathcal{H}[u-1] \right]$. Third, using the relation between η_1 and η_2 in (36), we relate these two bounds.

Step 1: Since $\eta_2 \tau_2 \leq \ln 2$, we have

$$\begin{aligned}
 & \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[g_2[u] | \mathcal{H}[u-1] \right] \\
 & = \sum_{v=1}^V \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,v]} \left[-\frac{1}{\eta_2} \right. \\
 & \quad \cdot \ln \left(\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}[u] e^{-\eta_2 L_k[u,v]} \right) \Big| \mathcal{H}[u-1] \Big] \\
 & \quad - \sum_{v=1}^V \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[L_{k_0^{\text{ROW}}[u]}[u,v] | \mathcal{H}[u-1] \right] \\
 & \leq \sum_{v=1}^V \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,v]} \left[\mathbb{E} \left[L[u,v] | \hat{\mathbb{k}}^{\text{ROW}}[u,v] \right] | \mathcal{H}[u-1] \right] \\
 & \quad - \frac{\eta_2}{8} \cdot \frac{M-1}{K-1} \cdot \sum_{v=1}^V \text{Var}(L[u,v] | \mathcal{H}[u-1]) \\
 & \quad - \sum_{v=1}^V \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[L_{k_0^{\text{ROW}}[u]}[u,v] | \mathcal{H}[u-1] \right] \tag{115}
 \end{aligned}$$

where the equality is because of the definition of $g_2[u]$ in (35) and the linearity of the expectation, and the inequality is because of

Lemma 4.3 and Lemma 4.6. Let us focus on the first term on the right-hand-side of (115). According to (31), we have

$$\begin{aligned}
 \hat{p}_k^{\text{ROW}}[u] & = \frac{w_k^{\text{ROW}}[u]}{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} w_k^{\text{ROW}}[u]} \\
 & = \frac{\frac{w_k^{\text{ROW}}[u]}{\sum_{k=1}^K w_k^{\text{ROW}}[u]}}{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \frac{w_k^{\text{ROW}}[u]}{\sum_{k=1}^K w_k^{\text{ROW}}[u]}} \\
 & = \frac{p_k^{\text{ROW}}[u]}{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} p_k^{\text{ROW}}[u]}. \tag{116}
 \end{aligned}$$

We let

$$\mathbb{E}[L[u,v] | \mathcal{H}[u-1]] \triangleq \sum_{k=1}^K p_k^{\text{ROW}}[u] L_k[u,v], \tag{117}$$

denote the expected loss of full feedback with regard to the randomness in $p_k^{\text{ROW}}[u]$. Thus, according to (116) and (117), we have

$$\begin{aligned}
 & \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,v]} \left[\mathbb{E} \left[L[u,v] | \hat{\mathbb{k}}^{\text{ROW}}[u,v] \right] | \mathcal{H}[u-1] \right] \\
 & = \sum_{\hat{\mathbb{k}}^{\text{ROW}}[u,v]} \text{Pr} \left\{ \hat{\mathbb{k}}^{\text{ROW}}[u,v] | \mathcal{H}[u-1] \right\} \\
 & \quad \cdot \sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}[u] L_k[u,v] \\
 & = \sum_{\hat{\mathbb{k}}^{\text{ROW}}[u,v]} \sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} p_k^{\text{ROW}}[u] \cdot \frac{1}{\binom{K-1}{M-1}} \\
 & \quad \cdot \sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} \frac{p_k^{\text{ROW}}[u]}{\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} p_k^{\text{ROW}}[u]} L_k[u,v] \\
 & = \frac{1}{\binom{K-1}{M-1}} \cdot \sum_{\hat{\mathbb{k}}^{\text{ROW}}[u,v]} \sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]} p_k^{\text{ROW}}[u] L_k[u,v] \\
 & = \mathbb{E}[L[u,v] | \mathcal{H}[u-1]], \tag{118}
 \end{aligned}$$

where the second equality is because of (101) and (116), and the fourth equality is because there are $\binom{K-1}{M-1}$ working groups containing each arm k . Then, from (115) and (118), we have

$$\begin{aligned}
 & \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[g_2[u] | \mathcal{H}[u-1] \right] \\
 & \leq \sum_{v=1}^V \mathbb{E}[L[u,v] | \mathcal{H}[u-1]] \\
 & \quad - \frac{\eta_2}{8} \cdot \frac{M-1}{K-1} \cdot \sum_{v=1}^V \text{Var}(L[u,v] | \mathcal{H}[u-1]) \\
 & \quad - \sum_{v=1}^V \mathbb{E}_{\hat{\mathbb{k}}^{\text{ROW}}[u,1:V]} \left[L_{k_0^{\text{ROW}}[u]}[u,v] | \mathcal{H}[u-1] \right], \tag{119}
 \end{aligned}$$

Step 2: Since $\eta_1 \tau_1 \leq \ln 2$, according to Lemma 4.5 and Lemma 4.7, we have

$$\begin{aligned} & \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[g_1[u] \middle| \mathcal{H}[u-1] \right] \\ & \geq \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[\mathbb{E} \left[\tilde{L}^{\text{ROW}}[u] \middle| \hat{k}^{\text{ROW}}[u,1:V] \right] \middle| \mathcal{H}[u-1] \right] \\ & \quad - \eta_1 \cdot \frac{2(K-1)}{M-1} \cdot \sum_{v=1}^V \text{Var} (L[u,v] | \mathcal{H}[u-1]). \end{aligned} \quad (120)$$

Let us focus on the first term on the right-hand-side of (120). We have

$$\begin{aligned} & \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[\mathbb{E} \left[\tilde{L}^{\text{ROW}}[u] \middle| \hat{k}^{\text{ROW}}[u,1:V] \right] \middle| \mathcal{H}[u-1] \right] \\ & = \sum_{k_0^{\text{ROW}}[u]=1}^K p_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u] \\ & \quad \cdot \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:k_0^{\text{ROW}}[u]]} \left[\sum_{k=1}^K p_k^{\text{ROW}}[u] \tilde{L}^{\text{ROW}}[u] \right] \\ & = \sum_{k_0^{\text{ROW}}[u]=1}^K p_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u] \\ & \quad \cdot \sum_{k=1}^K p_k^{\text{ROW}}[u] \cdot \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:k_0^{\text{ROW}}[u]]} \left[\tilde{L}^{\text{ROW}}[u] \right], \end{aligned} \quad (121)$$

where the second equality is because of the linearity of the expectation, and because the probability $p_k^{\text{ROW}}[u]$ is independent of the choice of the working group $\hat{k}^{\text{ROW}}[u,1:V]$. Notice that the conditional probability that arm k is chosen as a secondary arm in each sub-episode (conditioned on k not being the primary arm, but unconditioned on the events in other sub-episodes) is $\frac{M-1}{K-1}$. Thus, from (121), we have

$$\begin{aligned} & \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[\mathbb{E} \left[\tilde{L}^{\text{ROW}}[u] \middle| \hat{k}^{\text{ROW}}[u,1:V] \right] \middle| \mathcal{H}[u-1] \right] \\ & = \sum_{k_0^{\text{ROW}}[u]=1}^K p_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u] \sum_{k=1}^K p_k^{\text{ROW}}[u] \\ & \quad \cdot \sum_{v=1}^V \frac{M-1}{K-1} \frac{L_k[u, v_u(k)] - L_{k_0^{\text{ROW}}[u]}[u, v_u(k)]}{\frac{M-1}{K-1}} \\ & = \sum_{k_0^{\text{ROW}}[u]=1}^K p_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u] \sum_{k=1}^K p_k^{\text{ROW}}[u] \\ & \quad \cdot \sum_{v=1}^V \left[L_k[u, v_u(k)] - L_{k_0^{\text{ROW}}[u]}[u, v_u(k)] \right] \\ & = \sum_{v=1}^V \mathbb{E} [L[u, v] | \mathcal{H}[u-1]] \\ & \quad - \sum_{v=1}^V \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[L_{k_0^{\text{ROW}}[u]}[u, v] \middle| \mathcal{H}[u-1] \right]. \end{aligned} \quad (122)$$

When $k = k_0^{\text{ROW}}[u]$, $\tilde{L}_k^{\text{ROW}}[u] = \frac{L_k[u, v_u(k)] - L_{k_0^{\text{ROW}}[u]}[u, v_u(k)]}{\frac{M-1}{K-1}} = 0$. Thus, it does not affect the first equality of (122). Then, combining

(120) and (122), we have

$$\begin{aligned} & \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[g_1[u] \middle| \mathcal{H}[u-1] \right] \\ & \geq \sum_{v=1}^V \mathbb{E} [L[u, v] | \mathcal{H}[u-1]] \\ & \quad - \eta_1 \cdot \frac{2(K-1)}{M-1} \cdot \sum_{v=1}^V \text{Var} (L[u, v] | \mathcal{H}[u-1]) \\ & \quad - \sum_{v=1}^V \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[L_{k_0^{\text{ROW}}[u]}[u, v] \middle| \mathcal{H}[u-1] \right], \end{aligned} \quad (123)$$

Step 3: Finally, let us compare (119) and (123). The only difference is the second term on the right-hand-side. According to (36), $\eta_2 \geq 16 \left(\frac{K-1}{M-1} \right)^2 \cdot \eta_1$. Hence, we have (37). \square

N PROOF OF LEMMA 4.8

PROOF. We start from considering the first term on the left-hand-side of (43). First, we have that for all episodes $u = 1, \dots, U$,

$$\begin{aligned} & -\frac{1}{\eta_1} \ln \left(\sum_{k=1}^K p_k^{\text{ROW}}[u] e^{-\eta_1 \tilde{L}_k^{\text{ROW}}[u]} \right) \\ & = -\frac{1}{\eta_1} \ln \left(\sum_{k=1}^K \frac{w_k^{\text{ROW}}[u]}{\sum_{k=1}^K w_k^{\text{ROW}}[u]} e^{-\eta_1 \tilde{L}_k^{\text{ROW}}[u]} \right) \\ & = -\frac{1}{\eta_1} \ln \left(\frac{\sum_{k=1}^K w_k^{\text{ROW}}[u+1]}{\sum_{k=1}^K w_k^{\text{ROW}}[u]} \right), \end{aligned} \quad (124)$$

where the first equality is because of the update of the probability $p_k^{\text{ROW}}[u]$ in (24), and the second equality is because of the update of the weight $w_k^{\text{ROW}}[u]$ in (29). Then, according to (124), we have

$$\begin{aligned} & \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[-\frac{1}{\eta_1} \right. \right. \\ & \quad \left. \left. \cdot \ln \left(\sum_{k=1}^K p_k^{\text{ROW}}[u] e^{-\eta_1 \tilde{L}_k^{\text{ROW}}[u]} \right) \middle| \mathcal{H}[u-1] \right] \right\} \\ & = \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[-\frac{1}{\eta_1} \right. \right. \\ & \quad \left. \left. \cdot \ln \left(\frac{\sum_{k=1}^K w_k^{\text{ROW}}[u+1]}{\sum_{k=1}^K w_k^{\text{ROW}}[u]} \right) \middle| \mathcal{H}[u-1] \right] \right\} \\ & = \sum_{u=1}^U \mathbb{E}_{\text{ROW}} \left[-\frac{1}{\eta_1} \ln \left(\frac{\sum_{k=1}^K w_k^{\text{ROW}}[u+1]}{\sum_{k=1}^K w_k^{\text{ROW}}[u]} \right) \right]. \end{aligned}$$

Thus, we have

$$\begin{aligned}
 & \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\tilde{k}^{\text{ROW}}[u,1:V]} \left[-\frac{1}{\eta_1} \right. \right. \\
 & \quad \left. \left. \cdot \ln \left(\sum_{k=1}^K p_k^{\text{ROW}}[u] e^{-\eta_1 \tilde{L}_k^{\text{ROW}}[u]} \right) \middle| \mathcal{H}[u-1] \right] \right\} \\
 &= \mathbb{E}_{\text{ROW}} \left[\sum_{u=1}^U -\frac{1}{\eta_1} \ln \left(\frac{\sum_{k=1}^K w_k^{\text{ROW}}[u+1]}{\sum_{k=1}^K w_k^{\text{ROW}}[u]} \right) \right] \\
 &= \mathbb{E}_{\text{ROW}} \left[-\frac{1}{\eta_1} \ln \left(\frac{\sum_{k=1}^K w_k^{\text{ROW}}[U+1]}{\sum_{k=1}^K w_k^{\text{ROW}}[1]} \right) \right], \tag{125}
 \end{aligned}$$

where the first equality is because of the linearity of the expectation, and the fourth equality is because of the telescoping sum. Next, since $w_k^{\text{ROW}}[1] = 1$ for all k , from (125), we have

$$\begin{aligned}
 & \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\tilde{k}^{\text{ROW}}[u,1:V]} \left[-\frac{1}{\eta_1} \right. \right. \\
 & \quad \left. \left. \cdot \ln \left(\sum_{k=1}^K p_k^{\text{ROW}}[u] e^{-\eta_1 \tilde{L}_k^{\text{ROW}}[u]} \right) \middle| \mathcal{H}[u-1] \right] \right\} \\
 &= \mathbb{E}_{\text{ROW}} \left[-\frac{1}{\eta_1} \ln \sum_{k=1}^K w_k^{\text{ROW}}[U+1] \right] + \frac{\ln K}{\eta_1}. \tag{126}
 \end{aligned}$$

Notice that the second term on the right-hand-side of (126) is the first term on the right-hand-side of (43).

Then, let us focus on the first term on the right-hand-side of (126). We have

$$\begin{aligned}
 & \mathbb{E}_{\text{ROW}} \left[-\frac{1}{\eta_1} \ln \sum_{k=1}^K w_k^{\text{ROW}}[U+1] \right] \\
 & \leq \mathbb{E}_{\text{ROW}} \left[-\frac{1}{\eta_1} \ln w_{k^{\text{OPT}}}^{\text{ROW}}[U+1] \right] \\
 & = \mathbb{E}_{\text{ROW}} \left[-\frac{1}{\eta_1} \ln \left(e^{-\eta_1 \sum_{u=1}^U \tilde{L}_{k^{\text{OPT}}}^{\text{ROW}}[u]} \right) \right] \\
 & = \mathbb{E}_{\text{ROW}} \left[\sum_{u=1}^U \tilde{L}_{k^{\text{OPT}}}^{\text{ROW}}[u] \right], \tag{127}
 \end{aligned}$$

where the first inequality is because $\sum_{k=1}^K w_k^{\text{ROW}}[U+1] \geq w_{k^{\text{OPT}}}^{\text{ROW}}[U+1]$, the first equality is because of the update of the weight $w_k^{\text{ROW}}[u]$ in (29). Notice that if k^{OPT} is chosen as the primary arm, i.e., $k_0^{\text{ROW}}[u] = k^{\text{OPT}}$, we have $\tilde{L}_{k^{\text{OPT}}}^{\text{ROW}}[u] = 0$. This will not affect the last equality. On the other hand, the conditional probability that arm k^{OPT} is chosen as a secondary arm in each sub-episode (conditioned on k not being the primary arm, but unconditional on the events in other sub-episodes) is $\frac{M-1}{K-1}$. Thus, from (127), we

have

$$\begin{aligned}
 & \mathbb{E}_{\text{ROW}} \left[-\frac{1}{\eta_1} \ln \sum_{k=1}^K w_k^{\text{ROW}}[U+1] \right] \\
 & = \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left[\sum_{k_0^{\text{ROW}}[u]=1}^K p_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u] \right. \\
 & \quad \left. \cdot \sum_{v=1}^V \frac{M-1}{K-1} \frac{L_{k^{\text{OPT}}}[u,v] - L_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u,v]}{\frac{M-1}{K-1}} \right] \\
 & = \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left[\sum_{k_0^{\text{ROW}}[u]=1}^K p_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u] \right. \\
 & \quad \left. \cdot \sum_{v=1}^V \{L_{k^{\text{OPT}}}[u,v] - L_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u,v]\} \right]. \tag{128}
 \end{aligned}$$

Then, from (128), we have

$$\begin{aligned}
 & \mathbb{E}_{\text{ROW}} \left[-\frac{1}{\eta_1} \ln \sum_{k=1}^K w_k^{\text{ROW}}[U+1] \right] \\
 & \leq \sum_{t=1}^T l_{k^{\text{OPT}}}(t) \\
 & \quad - \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left[\sum_{k_0^{\text{ROW}}[u]=1}^K \left(p_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u] \sum_{v=1}^V L_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u,v] \right) \right]. \tag{129}
 \end{aligned}$$

where the inequality is because of the linearity of the expectation.

Finally, combining (126) and (129), we have (43). \square

O PROOF OF THEOREM 4.1

PROOF. We use

$$\text{Loss}^{\text{ROW}}(1:T) \triangleq \sum_{t=1}^T l_{k^{\text{ROW}}(t)}(t) \tag{130}$$

and

$$\text{Loss}^{\text{OPT}}(1:T) \triangleq \min_{k \in [1,K]} \sum_{t=1}^T l_k(t) \tag{131}$$

to denote the total loss of ROW and OPT, respectively. According to the upper bound of the switching costs of ROW in (44), we have

$$\begin{aligned}
 R^{\text{ROW}}(T) & \leq \max_{l_{i,K}(1:T)} \left\{ \mathbb{E}_{\text{ROW}} \left[\text{Loss}^{\text{ROW}}(1:T) \right] - \text{Loss}^{\text{OPT}}(1:T) \right\} \\
 & \quad + \min \{M, K-M\} \cdot \beta_1 \left[\frac{T}{\tau_2} \right]. \tag{132}
 \end{aligned}$$

In the following, we focus on calculating the worst-case difference between the expected total loss of ROW and the total loss of OPT, i.e., the first term on the right-hand-side of (132). First,

according to Lemma 4.2, we have

$$\begin{aligned}
& \mathbb{E}_{\text{ROW}} \left[\text{Loss}^{\text{ROW}}(1 : T) \right] \\
&= \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \sum_{v=1}^V \mathbb{E}_{\hat{k}^{\text{ROW}}[u,v]} \left[\sum_{t=t_{u,v}}^{t_{u,v}+\tau_2-1} \sum_{k \in \hat{k}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}(t) \right. \right. \\
&\quad \left. \left. \cdot l_k(t) \mathcal{H}[u-1] \right] \right\} \\
&\leq \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \sum_{v=1}^V \mathbb{E}_{\hat{k}^{\text{ROW}}[u,v]} \left[g_2[u, v] \right. \right. \\
&\quad \left. \left. + \frac{1}{2} \eta_2 \tau_2 + 1 \mathcal{H}[u-1] \right] \right\} \\
&= \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \sum_{v=1}^V \mathbb{E}_{\hat{k}^{\text{ROW}}[u,v]} \left[g_2[u, v] \mathcal{H}[u-1] \right] \right\} \\
&\quad + \frac{1}{2} \eta_2 T + \left\lceil \frac{T}{\tau_2} \right\rceil, \tag{133}
\end{aligned}$$

where the inequality is because of Lemma 4.2. Next, according to the relation between $g_2[u, v]$ and $g_2[u]$ in (35), from (133), we have

$$\begin{aligned}
& \mathbb{E}_{\text{ROW}} \left[\text{Loss}^{\text{ROW}}(1 : T) \right] \\
&\leq \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[g_2[u] \mathcal{H}[u-1] \right] \right\} \\
&\quad + \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[\sum_{v=1}^V L_{k_0^{\text{ROW}}[u]}[u, v] \mathcal{H}[u-1] \right] \right\} \\
&\quad + \frac{1}{2} \eta_2 T + \left\lceil \frac{T}{\tau_2} \right\rceil. \tag{134}
\end{aligned}$$

Notice that the values of the parameters in (46) satisfy the conditions in (36). Thus, Lemma 4.4 holds. Then, applying Lemma 4.4 to (134), we have

$$\begin{aligned}
& \mathbb{E}_{\text{ROW}} \left[\text{Loss}^{\text{ROW}}(1 : T) \right] \\
&\leq \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[g_1[u] \mathcal{H}[u-1] \right] \right\} \\
&\quad + \sum_{u=1}^U \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\hat{k}^{\text{ROW}}[u,1:V]} \left[\sum_{v=1}^V L_{k_0^{\text{ROW}}[u]}[u, v] \mathcal{H}[u-1] \right] \right\} \\
&\quad + \frac{1}{2} \eta_2 T + \left\lceil \frac{T}{\tau_2} \right\rceil. \tag{135}
\end{aligned}$$

Then, according to Lemma 4.8, from (135), we have

$$\mathbb{E}_{\text{ROW}} \left[\text{Loss}^{\text{ROW}}(1 : T) \right] \leq \sum_{t=1}^T l_{k^{\text{OPT}}}(t) + \frac{\ln K}{\eta_1} + \frac{1}{2} \eta_2 T + \left\lceil \frac{T}{\tau_2} \right\rceil.$$

Since (132)-(135) hold for all loss sequences $l_{1:K}(1 : T)$, we have

$$\begin{aligned}
& \mathbb{E}_{\text{ROW}} \left[\text{Loss}^{\text{ROW}}(1 : T) \right] - \sum_{t=1}^T l_{k^{\text{OPT}}}(t) \\
&\leq \frac{1}{2} \eta_2 T + \frac{\ln K}{\eta_1} + \left\lceil \frac{T}{\tau_2} \right\rceil, \text{ for all loss sequences } l_{1:K}(1 : T). \tag{136}
\end{aligned}$$

Finally, combining (132) and (136), we have

$$R^{\text{ROW}}(T) \leq \frac{1}{2} \eta_2 T + \frac{\ln K}{\eta_1} + \min \{M, K - M\} \cdot \beta_1 \left\lceil \frac{T}{\tau_2} \right\rceil + \left\lceil \frac{T}{\tau_2} \right\rceil. \tag{137}$$

(137) is the same as (45), which can be used as a reference for tuning the parameters.

In the following, we elaborate how we get the exact form of the final regret in Theorem 4.1. First, according to (46), we have $\eta_1 = \frac{\eta_2}{16 \left(\frac{K-1}{M-1} \right)^2}$ and $\tau_2 = \left\lfloor \frac{\ln 2}{\eta_2} \right\rfloor$. Thus, according to (137), the regret of ROW is upper-bounded as follows,

$$\begin{aligned}
R^{\text{ROW}}(T) &\leq \frac{1}{2} \eta_2 T + 16 \left(\frac{K-1}{M-1} \right)^2 \cdot \frac{\ln K}{\eta_2} \\
&\quad + (\min \{M, K - M\} \cdot \beta_1 + 1) \left\lceil \frac{T}{\eta_2} \right\rceil + \min \{M, K - M\} \cdot \beta_1 + 1 \\
&\leq \left(\frac{1}{2} + \frac{\min \{M, K - M\} \cdot \beta_1 + 1}{\ln 2 - \eta_2} \right) \eta_2 T + 16 \left(\frac{K-1}{M-1} \right)^2 \cdot \frac{\ln K}{\eta_2} \\
&\quad + \min \{M, K - M\} \cdot \beta_1 + 1. \tag{138}
\end{aligned}$$

The value of η_2 is then chosen to approximately minimize the right hand side. Specifically, according to (46), we have

$$\eta_2 = \sqrt{\frac{\ln K}{\frac{5}{2} + \min \{M, K - M\} \cdot 2\beta_1}} \cdot \frac{4(K-1)}{M-1} \cdot T^{-\frac{1}{2}}.$$

Thus, for $T \geq \frac{448(K-1)^2 \ln K}{\frac{5}{2} + 2\beta_1}$, we have

$$\begin{aligned}
\eta_2 &\leq \sqrt{\frac{16 \left(\frac{5}{2} + 2\beta_1 \right)}{448 \left(\frac{5}{2} + \min \{M, K - M\} \cdot 2\beta_1 \right) (M-1)^2}} \\
&\leq \sqrt{\frac{1}{28}} \leq \ln 2 - \frac{1}{2}. \tag{139}
\end{aligned}$$

Finally, combining (138) and (139), the regret of ROW is upper-bounded as follows,

$$\begin{aligned}
R^{\text{ROW}}(T) &\leq \left(\frac{1}{2} + 2 (\min \{M, K - M\} \cdot \beta_1 + 1) \right) \eta_2 T \\
&\quad + 16 \left(\frac{K-1}{M-1} \right)^2 \cdot \frac{\ln K}{\eta_2} + \min \{M, K - M\} \cdot \beta_1 + 1 \\
&\leq \frac{8(K-1)}{M-1} \sqrt{\frac{5}{2} + \min \{M, K - M\} \cdot 2\beta_1} \sqrt{\ln K} \sqrt{T} \\
&\quad + \min \{M, K - M\} \cdot \beta_1 + 1.
\end{aligned}$$

□