

Training a Single Bandit Arm: Supplementary Material

A Proof of Proposition 1

Proof of Proposition 1. For any policy π , we have that

$$\begin{aligned}
 \mathcal{R}_T(\pi, \boldsymbol{\nu}) &= \mathbb{E}(\max_{i \in [K]} \bar{U}_T^i) \\
 &= \mathbb{E}(\max_{i \in [K]} (\sum_{t=1}^T U_{n_{t-1}+1}^i \mathbb{1}_{\{I_t=i\}})) \\
 &\stackrel{(a)}{\leq} \mathbb{E} \left(\sum_{t=1}^T \max_{i \in [K]} (U_{n_{t-1}+1}^i \mathbb{1}_{\{I_t=i\}}) \right) \\
 &= \sum_{t=1}^T \mathbb{E} \left(\max_{i \in [K]} (U_{n_{t-1}+1}^i \mathbb{1}_{\{I_t=i\}}) \right) \\
 &\stackrel{(b)}{=} \sum_{t=1}^T \mathbb{E} \left(U_{n_{t-1}+1}^{I_t} \max_{i \in [K]} (\mathbb{1}_{\{I_t=i\}}) \right) \\
 &= \sum_{t=1}^T \mathbb{E} \left(U_{n_{t-1}+1}^{I_t} \right) \\
 &= \sum_{t=1}^T \mathbb{E} \left(\mathbb{E} \left(U_{n_{t-1}+1}^{I_t} \mid \mathcal{H}_t \right) \right) \\
 &\stackrel{(c)}{=} \sum_{t=1}^T \mathbb{E}(\mu_{I_t}) \leq \mu_1 T.
 \end{aligned}$$

Here, (a) is obtained due to pushing the max inside the sum; (b) is obtained because $U_{n_{t-1}+1}^i \geq 0$ for all i ; and (c) holds because the reward for an arm in a period is independent of the past history of play and observations. Thus, the reward of $\mu_1 T$ is the highest that one can obtain under any policy. And this reward can, in fact, be obtained by the policy of always picking arm 1. This shows that

$$\sup_{\pi \in \Pi} \mathcal{R}_T(\pi, \boldsymbol{\nu}) = \mathcal{R}_T^*(\boldsymbol{\nu}).$$

□

B Proof of Theorem 1

The proof of Theorem 1 relies on the following key result.

Proposition 2. *Consider a class $\mathcal{V} = \mathcal{M}^K$ of K -armed stochastic bandits and let $(\pi_T)_{T \in \mathbb{N}}$ be a consistent sequence of policies for \mathcal{V} . Then, for all $\alpha \in (0, 1]$ and $\nu \in \mathcal{V}$ such that the optimal arm k^* is unique,*

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}_{\nu} \left[n_{\lceil T^\alpha \rceil}^i \right]}{\log(T)} \geq \frac{\alpha}{d_{\inf}(\nu_i, \mu^*, \mathcal{M})}$$

holds for each suboptimal arm $i \neq k^*$ in ν , where μ^* is the highest mean.

Proof of Proposition 2. In what follows, we denote \mathbb{P}_{ν} to be the probability distribution induced by the policy π on events until time T under bandit ν , and we let \mathbb{E}_{ν} denote the corresponding expectation.

Let $\text{Reg}_{\text{SUM}, T}(\pi, \nu)$ denote the expected regret of the sum objective after T pulls of policy π under the bandit instance ν , which can be defined as

$$\text{Reg}_{\text{SUM}, T}(\pi, \nu) = \mu^* T - \mathbb{E}_{\nu} \left(\sum_{t=1}^T X_t \right) \quad (5)$$

$$= \mu^* T - \mathbb{E}_{\nu} \left(\sum_{i=1}^K \bar{U}_T^i \right), \quad (6)$$

where $X_t = U_{n_t}^{I_t}$, which is the reward due to the arm pulled at time t , and $\bar{U}_t^i = \sum_{n=1}^{n_t} U_n^i$, which is the cumulative reward obtained from arm i until time t . We need the following two lemmas for our proof.

Lemma 1. *Fix $\alpha \in (0, 1]$ and a policy π . Consider a K -armed bandit instance ν with $\mu^* \triangleq \mu_1 \geq \mu_2 \geq \dots \geq \mu_K$. Fix a suboptimal arm i and let $A_i = \left\{ n_{\lceil T^\alpha \rceil}^i > \frac{T^\alpha}{2} \right\}$. Then,*

$$\text{Reg}_{\text{SUM}, T}(\pi, \nu) > \mathbb{P}_{\nu}(A_i) \frac{T^\alpha \Delta_i}{2}.$$

Lemma 2. *Fix $\alpha \in (0, 1]$ and a policy π . Consider a K -armed bandit instance ν with $\mu^* \triangleq \mu_1 \geq \mu_2 \geq \dots \geq \mu_K$. Fix a suboptimal arm i and construct another K -armed bandit instance ν' satisfying $\mu'_i > \mu^* = \mu_1 \geq \mu_2 \geq \dots \geq \mu_{i-1} \geq \mu_{i+1} \geq \dots \geq \mu_K$. Let $A_i^c = \left\{ n_{\lceil T^\alpha \rceil}^i \leq \frac{T^\alpha}{2} \right\}$. Then,*

$$\text{Reg}_{\text{SUM}, T}(\pi, \nu') \geq \mathbb{P}_{\nu'}(A_i^c) \frac{T^\alpha (\mu'_i - \mu^*)}{2}.$$

The proof of Lemma 1 is presented below at the end of this section. The proof of Lemma 2 is similar and hence is omitted.

Fix $\alpha \in (0, 1]$. We proceed by constructing a second bandit ν' . Fix a suboptimal arm i , i.e., $\Delta_i > 0$, and let $\nu'_j = \nu_j$ for $j \neq i$ and pick a $\nu'_i \in \mathcal{M}$ such that $D(\nu_i, \nu'_i) \leq d_i + \epsilon$ and $\mu'_i > \mu^*$ for some arbitrary $\epsilon > 0$.

Let μ_i (μ'_i) be the mean of arm i in ν (ν') and $d_i \triangleq d_{\inf}(\nu_i, \mu^*, \mathcal{M})$. Recall that $d_{\inf}(\nu, \mu^*, \mathcal{M}) = \inf_{\nu' \in \mathcal{M}} \{D(\nu, \nu') : \mu(\nu') > \mu^*\}$ where $\mu(\nu)$ denotes the mean of distribution ν .

Since any lower bound on the regret for the sum objective implies the same lower bound on the max objective,

using Lemma 1 and Lemma 2, we have the following:

$$\begin{aligned}
 \text{Reg}_T(\pi, \boldsymbol{\nu}) + \text{Reg}_T(\pi, \boldsymbol{\nu}') &\geq \text{Reg}_{\text{SUM}, T}(\pi, \boldsymbol{\nu}) + \text{Reg}_{\text{SUM}, T}(\pi, \boldsymbol{\nu}') \\
 &> \frac{T^\alpha}{2} (\mathbb{P}_{\boldsymbol{\nu}}(A_i) \Delta_i + \mathbb{P}_{\boldsymbol{\nu}'}(A_i^c) (\mu'_i - \mu^*)) \\
 &\geq \frac{T^\alpha}{2} \min\{\Delta_i, (\mu'_i - \mu^*)\} (\mathbb{P}_{\boldsymbol{\nu}}(A_i) + \mathbb{P}_{\boldsymbol{\nu}'}(A_i^c)) \\
 &= \frac{T^\alpha}{2} \min\{\Delta_i, (\mu'_i - \mu^*)\} (\bar{\mathbb{P}}_{\boldsymbol{\nu}}(A_i) + \bar{\mathbb{P}}_{\boldsymbol{\nu}'}(A_i^c)) \\
 &\geq \frac{T^\alpha}{4} \min\{\Delta_i, (\mu'_i - \mu^*)\} \exp\left(-\mathbb{E}_{\boldsymbol{\nu}} \left[n_{\lceil T^\alpha \rceil}^i\right] (d_i + \epsilon)\right). \tag{7}
 \end{aligned}$$

Here, $\bar{\mathbb{P}}_{\boldsymbol{\nu}}$ ($\bar{\mathbb{P}}_{\boldsymbol{\nu}'}$) is the probability distribution induced by the policy π on events until time $\lceil T^\alpha \rceil$ under bandit $\boldsymbol{\nu}$ ($\boldsymbol{\nu}'$). The equality then results from the fact that the two events $\{n_{\lceil T^\alpha \rceil}^i > \frac{T^\alpha}{2}\}$ and $\{n_{\lceil T^\alpha \rceil}^i \leq \frac{T^\alpha}{2}\}$ depend only on the play until time $\lceil T^\alpha \rceil$. The last inequality follows from using the Bretagnolle-Huber inequality and divergence decomposition (see Theorem 14.2 and Lemma 15.1 in Lattimore and Szepesvári (2018), respectively) combined with the fact that $D(\nu_i, \nu'_i) \leq d_i + \epsilon$:

$$\bar{\mathbb{P}}_{\boldsymbol{\nu}}(A_i) + \bar{\mathbb{P}}_{\boldsymbol{\nu}'}(A_i^c) \geq \frac{1}{2} \exp(-D(\bar{\mathbb{P}}_{\boldsymbol{\nu}}, \bar{\mathbb{P}}_{\boldsymbol{\nu}'})) \geq \frac{1}{2} \exp\left(-\mathbb{E}_{\boldsymbol{\nu}} \left[n_{\lceil T^\alpha \rceil}^i\right] (d_i + \epsilon)\right), \tag{8}$$

where the events A_i and A_i^c are defined as they have been in Lemmas 1 and 2 for the fixed arm i .

Rearranging Equation 7, we obtain

$$\frac{\mathbb{E}_{\boldsymbol{\nu}} \left[n_{\lceil T^\alpha \rceil}^i\right]}{\log(T)} > \frac{1}{d_i + \epsilon} \frac{\log\left(\frac{T^\alpha \min\{\Delta_i, \mu'_i - \mu^*\}}{4(\text{Reg}_T(\pi, \boldsymbol{\nu}) + \text{Reg}_T(\pi, \boldsymbol{\nu}'))}\right)}{\log(T)}, \tag{9}$$

and taking the limit inferior yields

$$\begin{aligned}
 \liminf_{T \rightarrow \infty} \frac{\mathbb{E}_{\boldsymbol{\nu}} \left[n_{\lceil T^\alpha \rceil}^i\right]}{\log(T)} &> \frac{1}{d_i + \epsilon} \liminf_{T \rightarrow \infty} \frac{\log\left(\frac{T^\alpha \min\{\Delta_i, \mu'_i - \mu^*\}}{4(\text{Reg}_T(\pi, \boldsymbol{\nu}) + \text{Reg}_T(\pi, \boldsymbol{\nu}'))}\right)}{\log(T)}, \\
 &= \frac{1}{d_i + \epsilon} \liminf_{T \rightarrow \infty} \frac{\alpha \log(T) + \log(\beta_i) - \log(4) - \log(\text{Reg}_T(\pi, \boldsymbol{\nu}) + \text{Reg}_T(\pi, \boldsymbol{\nu}'))}{\log(T)} \\
 &= \frac{1}{d_i + \epsilon} \left(\alpha - \limsup_{T \rightarrow \infty} \frac{\log(\text{Reg}_T(\pi, \boldsymbol{\nu}) + \text{Reg}_T(\pi, \boldsymbol{\nu}'))}{\log(T)}\right) \\
 &\geq \frac{\alpha}{d_i + \epsilon}, \tag{10}
 \end{aligned}$$

where $\beta_i = \min\{\Delta_i, \mu'_i - \mu^*\}$.

Since π is a consistent policy over the class \mathcal{V} , we can find a constant c_p for any $p > 0$ such that $\text{Reg}_T(\pi, \boldsymbol{\nu}) + \text{Reg}_T(\pi, \boldsymbol{\nu}') \leq c_p T^p$, which implies

$$\limsup_{T \rightarrow \infty} \frac{\log(\text{Reg}_T(\pi, \boldsymbol{\nu}) + \text{Reg}_T(\pi, \boldsymbol{\nu}'))}{\log(T)} \leq \limsup_{T \rightarrow \infty} \frac{p \log(T) + \log(c_p)}{\log(T)} = p. \tag{11}$$

Then, Equation 10 follows from Equation 11 and the fact that $p > 0$ is arbitrary. Since $\epsilon > 0$ is arbitrary as well, we have

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}_{\boldsymbol{\nu}} \left[n_{\lceil T^\alpha \rceil}^i\right]}{\log(T)} \geq \frac{\alpha}{d_i} \tag{12}$$

for each $i \neq k^*$, i.e., each suboptimal arm i in $\boldsymbol{\nu}$. \square

We present the proof of Lemma 1 before proceeding with the proof of Theorem 1.

Proof of Lemma 1. Recall that I_t is the arm pulled at time t and X_t is the reward due to arm pulled at time t , i.e., $X_t \sim \nu_{I_t}$. Then, due to, e.g., Lemma 4.5 in Lattimore and Szepesvári (2018), we can decompose the expected regret as

$$\text{Reg}_{\text{SUM},T}(\pi, \nu) = \sum_{\substack{j=1 \\ j \neq i}}^K \Delta_j \mathbb{E}_\nu(n_T^j) + \Delta_i \mathbb{E}_\nu(n_T^i). \quad (13)$$

Due to the non-negativity of expected number of pulls and the suboptimality gaps, we have

$$\text{Reg}_{\text{SUM},T}(\pi, \nu) \geq \Delta_i \mathbb{E}_\nu(n_T^i).$$

Now, we look at $\mathbb{E}_\nu(n_T^i)$:

$$\begin{aligned} \mathbb{E}_\nu(n_T^i) &= \mathbb{E}_\nu(n_T^i | A_i) \mathbb{P}_\nu(A_i) + \mathbb{E}_\nu(n_T^i | A_i^c) \mathbb{P}_\nu(A_i^c) \\ &\geq \mathbb{E}_\nu(n_T^i | A_i) \mathbb{P}_\nu(A_i) \\ &\stackrel{(a)}{>} \frac{T^\alpha}{2} \mathbb{P}_\nu(A_i), \end{aligned} \quad (14)$$

where (a) is due to event $A_i = \left\{ n_{\lceil T^\alpha \rceil}^i > \frac{T^\alpha}{2} \right\}$. Finally, we have

$$\text{Reg}_{\text{SUM},T}(\pi, \nu) > \mathbb{P}_\nu(A_i) \frac{T^\alpha \Delta_i}{2}. \quad (15)$$

□

Proof of Theorem 1. Let k^* denote the unique optimal arm in ν and, without loss of generality, let $k^* = 1$, i.e., $\mu^* = \mu_1$. Let I^* denote the arm with the highest cumulative reward after T pulls and recall that n_T^i denotes the number of pulls spent on arm i until time T . Since all of the following expectations are over ν , we drop the subscript of ν hereafter. We first look at the expected regret:

$$\text{Reg}_T(\pi, \nu) = \mu^* T - \mathbb{E} \left[\max \left(\bar{U}_T^1, \bar{U}_T^2, \dots, \bar{U}_T^K \right) \right] \quad (16)$$

$$\stackrel{(a)}{=} \mathbb{E} \left[\sum_{t=1}^T U_t^1 \right] - \mathbb{E} \left[\max \left(\bar{U}_T^1, \bar{U}_T^2, \dots, \bar{U}_T^K \right) \mathbb{1}_{\{I^*=1\}} \right] - \mathbb{E} \left[\max \left(\bar{U}_T^1, \bar{U}_T^2, \dots, \bar{U}_T^K \right) \mathbb{1}_{\{I^* \neq 1\}} \right] \quad (17)$$

$$\stackrel{(b)}{=} \mathbb{E} \left[\sum_{t=1}^T U_t^1 \mathbb{1}_{\{I^*=1\}} \right] + \mathbb{E} \left[\sum_{t=1}^T U_t^1 \mathbb{1}_{\{I^* \neq 1\}} \right] - \mathbb{E} \left[\sum_{t=1}^{n_T^1} U_t^1 \mathbb{1}_{\{I^*=1\}} \right] - \sum_{i \neq 1} \mathbb{E} \left[\sum_{t=1}^{n_T^i} U_t^i \mathbb{1}_{\{I^*=i\}} \right] \quad (18)$$

$$= \mathbb{E} \left[\left(\sum_{t=1}^T U_t^1 - \sum_{t=1}^{n_T^1} U_t^1 \right) \mathbb{1}_{\{I^*=1\}} \right] + \sum_{i \neq 1} \mathbb{E} \left[\left(\sum_{t=1}^T U_t^1 - \sum_{t=1}^{n_T^i} U_t^i \right) \mathbb{1}_{\{I^*=i\}} \right] \quad (19)$$

$$\stackrel{(c)}{\geq} \mathbb{E} \left[\left(\sum_{t=1}^T U_t^1 - \sum_{t=1}^{n_T^1} U_t^1 \right) \mathbb{1}_{\{I^*=1\}} \right] + \sum_{i \neq 1} \mathbb{E} \left[\left(\sum_{t=1}^T U_t^1 - \sum_{t=1}^T U_t^i \right) \mathbb{1}_{\{I^*=i\}} \right] \quad (20)$$

$$= \mathbb{E} \left[\left(\sum_{t=n_T^1+1}^T U_t^1 \right) \mathbb{1}_{\{I^*=1\}} \right] + \sum_{i \neq 1} \mathbb{E} \left[\left(\sum_{t=1}^T U_t^1 - \sum_{t=1}^T U_t^i \right) \mathbb{1}_{\{I^*=i\}} \right] \quad (21)$$

$$\stackrel{(d)}{=} \mu^* \mathbb{E} \left[(T - n_T^1) \mathbb{1}_{\{I^*=1\}} \right] + \sum_{i \neq 1} \mathbb{E} \left[\left(\sum_{t=1}^T U_t^1 - \sum_{t=1}^T U_t^i \right) \mathbb{1}_{\{I^*=i\}} \right] \quad (22)$$

$$\stackrel{(e)}{=} \mu^* \sum_{i \neq 1} \mathbb{E} \left[n_T^i \mathbb{1}_{\{I^*=1\}} \right] + \sum_{i \neq 1} \mathbb{E} \left[\left(\sum_{t=1}^T U_t^1 - \sum_{t=1}^T U_t^i \right) \mathbb{1}_{\{I^*=i\}} \right]. \quad (23)$$

Here, (a) is due to the fact that $\mathbb{E}[U_t^1] = \mu^*$ for $t \in [T]$. (b) follows from the definition of I^* . (c) results from $n_T^i \leq T$ for $i \in [K]$. (d) is due to the fact that the future rewards from the first arm is independent of the past history of play and observations of policy π . Finally, (e) follows from the identity $T = \sum_{i=1}^K n_T^i$.

We first focus on bounding the second term in the Expression 23. In order to do that, for each suboptimal arm $i, i \neq 1$, define a “good” event

$$G_i = \left\{ \bar{U}_T^1 > \bar{U}_T^i + \frac{T\Delta_i}{2} \right\}.$$

Notice that, for $i \neq 1$, $\Delta_i > 0$.

We proceed by showing that event G_i occurs with high probability. To that end, consider the complement event

$$\begin{aligned} P(G_i^c) &= P\left(\bar{U}_T^1 \leq \bar{U}_T^i + \frac{T\Delta_i}{2}\right) \\ &= P\left(\frac{\bar{U}_T^1 - \bar{U}_T^i}{T} - (\mu_1 - \mu_i) \leq \frac{\Delta_i}{2} - (\mu_1 - \mu_i)\right). \end{aligned} \quad (24)$$

By Hoeffding’s inequality,

$$P(G_i^c) \leq \exp\left(-\frac{2T^2\left(\frac{\Delta_i}{2}\right)^2}{4T}\right) = \exp\left(-\frac{T\Delta_i^2}{8}\right),$$

since $-1 \leq U_{t_1}^1 - U_{t_2}^i \leq 1$ for any pair $t_1, t_2 \in [T]$. We thus also have that

$$P(I^* = i, G_i^c) \leq \exp\left(-\frac{T\Delta_i^2}{8}\right). \quad (25)$$

We then have

$$\begin{aligned} &\mathbb{E}\left[\left(\sum_{t=1}^T U_t^1 - \sum_{t=1}^T U_t^i\right) \mathbb{1}_{\{I^*=i\}}\right] \\ &= \mathbb{E}\left[\left(\sum_{t=1}^T U_t^1 - \sum_{t=1}^T U_t^i\right) \mid I^* = i, G_i\right] P(I^* = i, G_i) + \mathbb{E}\left[\left(\sum_{t=1}^T U_t^1 - \sum_{t=1}^T U_t^i\right) \mid I^* = i, G_i^c\right] P(I^* = i, G_i^c) \\ &\geq \frac{T\Delta_i}{2} P(I^* = i, G_i) - TP(I^* = i, G_i^c) \\ &\geq \frac{T\Delta_i}{2} P(I^* = i, G_i) - O(1) \\ &\geq 0 - O(1). \end{aligned} \quad (26)$$

Thus the second term in (23) is lower bounded by a (instance-dependent) constant.

Next, we bound the first term in (23). To do so, we first need an upper bound on $P(I^* = i)$ for any $i \neq 1$. By consistency of policy π , we have that $\text{Reg}_T(\pi, \nu) \leq o(T^p)$ for every $p > 0$. Thus from (23) and (26), for any $i \neq 1$, we have that

$$o(T^p) \geq \mathbb{E}\left[\left(\sum_{t=1}^T U_t^1 - \sum_{t=1}^T U_t^i\right) \mathbb{1}_{\{I^*=i\}}\right] \geq \frac{T\Delta_i}{2} P(I^* = i, G_i) - O(1). \quad (28)$$

This implies that for any $i \neq 1$,

$$P(I^* = i, G_i) \leq o(T^{p-1}), \quad (29)$$

for every $p > 0$. Finally, (25) and (29) together imply that, for any $i \neq 1$, $P(I^* = i) = P(I^* = i, G_i) + P(I^* = i, G_i^c) \leq o(T^{p-1})$ for every $p > 0$.

Finally, we are ready to derive a lower bound on the first term in the expression (23). For any $\alpha \in (0, 1)$, we have

$$\begin{aligned}
 \mathbb{E}[n_{\lceil T^\alpha \rceil}^i] &= \mathbb{E} \left[n_{\lceil T^\alpha \rceil}^i \mathbb{1}_{\{I^*=1\}} \right] + \mathbb{E} \left[n_{\lceil T^\alpha \rceil}^i \mathbb{1}_{\{I^* \neq 1\}} \right] \\
 &\leq \mathbb{E} \left[n_{\lceil T^\alpha \rceil}^i \mathbb{1}_{\{I^*=1\}} \right] + \lceil T^\alpha \rceil P(I^* \neq 1) \\
 &\leq \mathbb{E} \left[n_T^i \mathbb{1}_{\{I^*=1\}} \right] + \lceil T^\alpha \rceil P(I^* \neq 1) \\
 &\leq \mathbb{E} \left[n_T^i \mathbb{1}_{\{I^*=1\}} \right] + o(T^{\alpha+p-1}),
 \end{aligned} \tag{30}$$

for every $p > 0$. But then from Proposition 2, we have

$$\frac{\alpha}{d_i} \leq \liminf_{T \rightarrow \infty} \frac{\mathbb{E}[n_{\lceil T^\alpha \rceil}^i]}{\log T} \tag{31}$$

$$\leq \liminf_{T \rightarrow \infty} \frac{\mathbb{E} \left[n_T^i \mathbb{1}_{\{I^*=1\}} \right]}{\log T} + \liminf_{T \rightarrow \infty} \frac{o(T^{\alpha+p-1})}{\log T}. \tag{32}$$

By choosing a p such that $0 < p < 1 - \alpha$, we have that $\liminf_{T \rightarrow \infty} \frac{o(T^{\alpha+p-1})}{\log T} = 0$. And thus, for every $\alpha \in (0, 1)$, we have

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E} \left[n_T^i \mathbb{1}_{\{I^*=1\}} \right]}{\log T} \geq \frac{\alpha}{d_i}, \tag{33}$$

which implies that

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E} \left[n_T^i \mathbb{1}_{\{I^*=1\}} \right]}{\log T} \geq \frac{1}{d_i}. \tag{34}$$

Finally, putting everything together, from (23), (27), and (34), we have

$$\begin{aligned}
 \liminf_{T \rightarrow \infty} \frac{\text{Reg}_T(\pi, \nu)}{\log T} &\geq \liminf_{T \rightarrow \infty} \mu^* \sum_{i \neq 1} \frac{\mathbb{E} \left[n_T^i \mathbb{1}_{\{I^*=1\}} \right]}{\log T} - \liminf_{T \rightarrow \infty} \sum_{i \neq 1} \frac{O(1)}{\log T} \\
 &\geq \sum_{i \neq 1} \frac{\mu^*}{d_i}.
 \end{aligned} \tag{35}$$

Plugging in the definition of d_i and substituting k^* back in place give the desired result. \square

C Proof of Theorem 2

Proof of Theorem 2. First we fix a policy $\pi \in \Pi$. Let $\Delta \triangleq (K-1)^{1/3}/(2T^{1/3})$. We construct two bandit environments with different reward distributions for each of the arms and show that π cannot perform well in both environments simultaneously.

We first specify the reward distribution for the arms in the base environment, denoted as the bandit $\nu = \{\nu_1, \dots, \nu_K\}$. Assume that the reward for all of the arms have the Bernoulli distribution, i.e., $\nu_i \sim \text{Bernoulli}(\mu_i)$. We let $\mu_1 = \frac{1}{2} + \Delta$, and $\mu_i = \frac{1}{2}$ for $2 \leq i \leq K$. We let P_ν denote the probability distribution induced over events until time T under policy π in this first environment, i.e., in bandit ν . Let E_ν denote the expectation under P_ν .

Define $n_{\lceil \Delta T \rceil}^i$ as the (random) number of pulls spent on arm $i \in \{1, \dots, K\}$ until time $\lceil \Delta T \rceil$ (note that $\sum_{i=1}^K n_{\lceil \Delta T \rceil}^i = \lceil \Delta T \rceil$) under policy π . Specifically, $n_{\lceil \Delta T \rceil}^1$ is the total (random) number of pulls spent on the first arm under policy π until time $\lceil \Delta T \rceil$. Under policy π , let l^* denote the arm in the set $[K] \setminus \{1\}$ that is pulled the least in expectation until time $\lceil \Delta T \rceil$, i.e., $l^* \in \arg \min_{2 \leq i \leq K} E_\nu(n_{\lceil \Delta T \rceil}^i)$. Then clearly, we have that $E_\nu(n_{\lceil \Delta T \rceil}^{l^*}) \leq \frac{\lceil \Delta T \rceil}{K-1}$.

Having defined l^* , we can now define the second environment, denoted as the bandit $\nu' = \{\nu'_1, \dots, \nu'_K\}$. Again, assume that the reward for all of the arms have the Bernoulli distribution, i.e., $\nu'_i \sim \text{Bernoulli}(\mu'_i)$. We let $\mu'_1 = \frac{1}{2} + \Delta$, $\mu'_i = \frac{1}{2}$ for $[2 \leq i \leq K] \setminus \{l^*\}$, and $\mu'_{l^*} = \frac{1}{2} + 2\Delta$. We let $P_{\nu'}$ denote the probability distribution induced over events until time T under policy π in this second environment, i.e., in bandit ν' . Let $E_{\nu'}$ denote the expectation under $P_{\nu'}$.

With some abuse of notation, for any event B , we define:

$$\text{Reg}_T(\pi, \nu, B) = \mu^* TP_\nu(B) - E_\nu(\max(\bar{U}_T^1, \bar{U}_T^2, \dots, \bar{U}_T^K) \mathbb{1}_B). \quad (36)$$

It is then clear that $\text{Reg}_T(\pi, \nu) = \text{Reg}_T(\pi, \nu, B) + \text{Reg}_T(\pi, \nu, B^c)$. We need the following two results for our proof.

Lemma 3. *Fix a policy π . Consider the K -armed bandit instance ν with Bernoulli rewards and mean vector $\mu = (\frac{1}{2} + \Delta, \frac{1}{2}, \frac{1}{2}, \dots, \frac{1}{2})$, where $\Delta < \frac{1}{2}$. Consider the event $A = \{n_{\lceil \Delta T \rceil}^1 \leq \frac{\Delta T}{2}\}$. Then we have,*

$$\text{Reg}_T(\pi, \nu, A) \geq \frac{\Delta T}{4} P_\nu(A) - 2\sqrt{T \log(KT)} - 2.$$

The proof of Lemma 3 is presented below in this section. A similar argument shows the following.

Lemma 4. *Fix a policy π . Consider the K -armed bandit instance ν' with Bernoulli rewards and mean vector $\mu' = (\frac{1}{2} + \Delta, \frac{1}{2}, \frac{1}{2}, \dots, \frac{1}{2}, \frac{1}{2} + 2\Delta)$, where $\Delta < \frac{1}{4}$. Consider the event $A^c = \{n_{\lceil \Delta T \rceil}^1 > \frac{\Delta T}{2}\}$. Then we have,*

$$\text{Reg}_T(\pi, \nu', A^c) \geq \frac{\Delta T}{4} P_{\nu'}(A^c) - 2\sqrt{T \log(KT)} - 2.$$

The proof of Lemma 4 is omitted since it is almost identical to that of Lemma 3. These two facts result in the following two inequalities:

$$\text{Reg}_T(\pi, \nu, A) \geq P_\nu \left(n_{\lceil \Delta T \rceil}^1 \leq \frac{\Delta T}{2} \right) \Omega(\Delta T), \text{ and} \quad (37)$$

$$\text{Reg}_T(\pi, \nu', A^c) \geq P_{\nu'} \left(n_{\lceil \Delta T \rceil}^1 > \frac{\Delta T}{2} \right) \Omega(\Delta T). \quad (38)$$

Note that here we have ignored the lower order $\sqrt{T \log(KT)}$ terms since $\Delta T = \Theta(T^{2/3} K^{1/3})$. Now, using the Bretagnolle-Huber inequality (see Theorem 14.2 in Lattimore and Szepesvári (2018)), we have,

$$\text{Reg}_T(\pi, \nu, A) + \text{Reg}_T(\pi, \nu', A^c) \geq \Omega(\Delta T) \left(P_\nu \left(n_{\lceil \Delta T \rceil}^1 \leq \frac{\Delta T}{2} \right) + P_{\nu'} \left(n_{\lceil \Delta T \rceil}^1 > \frac{\Delta T}{2} \right) \right) \quad (39)$$

$$= \Omega(\Delta T) \left(\bar{P}_\nu \left(n_{\lceil \Delta T \rceil}^1 \leq \frac{\Delta T}{2} \right) + \bar{P}_{\nu'} \left(n_{\lceil \Delta T \rceil}^1 > \frac{\Delta T}{2} \right) \right) \quad (40)$$

$$\geq \Omega(\Delta T) \exp(-D(\bar{P}_\nu, \bar{P}_{\nu'})). \quad (41)$$

Here, \bar{P}_ν ($\bar{P}_{\nu'}$) is the probability distribution induced by the policy π on events until time $\lceil \Delta T \rceil$ under bandit ν (ν'). The first equality then results from the fact that the two events $\{n_{\lceil \Delta T \rceil}^1 \leq \frac{\Delta T}{2}\}$ and $\{n_{\lceil \Delta T \rceil}^1 > \frac{\Delta T}{2}\}$ depend only on the play until time $\lceil \Delta T \rceil$. In the second inequality, which results from the Bretagnolle-Huber inequality, $D(\bar{P}_\nu, \bar{P}_{\nu'})$ is the relative entropy, or the Kullback-Leibler (KL) divergence between the distributions \bar{P}_ν and $\bar{P}_{\nu'}$ respectively. We can upper bound $D(\bar{P}_\nu, \bar{P}_{\nu'})$ as,

$$D(\bar{P}_\nu, \bar{P}_{\nu'}) = \mathbb{E}_\nu(n_{\lceil \Delta T \rceil}^{l^*}) D(\nu_{l^*}, \nu'_{l^*}) \leq \frac{\lceil \Delta T \rceil}{K-1} D(\nu_{l^*}, \nu'_{l^*}) \lesssim \frac{8\Delta^3 T}{K-1}, \quad (42)$$

where ν_{l^*} (ν'_{l^*}) denotes the reward distribution of arm l^* in the first (second) environment. The first equality results from divergence decomposition (see Lemma 15.1 in Lattimore and Szepesvári (2018)) and the fact no arm other than l^* offers any distinguishability between ν and ν' . The next inequality follows from the fact that $\mathbb{E}_\nu[n_{\lceil \Delta T \rceil}^{l^*}] \leq (\lceil \Delta T \rceil)/(K-1)$, since by definition, l^* is the arm that is pulled the least in expectation until time $\lceil \Delta T \rceil$ in bandit ν under π . Now $D(\nu_{l^*}, \nu'_{l^*})$ is simply the relative entropy between the distributions Bernoulli(1/2) and Bernoulli(1/2 + 2Δ), which, by elementary calculations, can be shown to be at most $8\Delta^2$, resulting in the final inequality. Thus, we finally have,

$$\text{Reg}_T(\pi, \nu, A) + \text{Reg}_T(\pi, \nu', A^c) \geq \Omega(\Delta T) \exp\left(-\frac{8\Delta^3 T}{K-1}\right).$$

Substituting $\Delta = (K-1)^{1/3}/(2T^{1/3})$ gives

$$\text{Reg}_T(\pi, \nu, A) + \text{Reg}_T(\pi, \nu', A^c) \geq \Omega\left((K-1)^{1/3} T^{2/3}\right). \quad (43)$$

Equation 43 along with

$$\text{Reg}_T(\pi, \nu, A^c) \geq -O(\sqrt{T \log(KT)}) \quad \text{and} \quad (44)$$

$$\text{Reg}_T(\pi, \nu', A) \geq -O(\sqrt{T \log(KT)}), \quad (45)$$

imply that

$$\text{Reg}_T(\pi, \nu) + \text{Reg}_T(\pi, \nu') \geq \Omega\left((K-1)^{1/3} T^{2/3}\right). \quad (46)$$

Finally, using $2 \max\{a, b\} \geq a + b$ gives the desired lower bound on the regret.

Showing Equations 44 and 45 is an easy exercise:

$$\begin{aligned} \text{Reg}_T(\pi, \nu, A^c) &= \mu^* TP_\nu(A^c) - \mathbb{E}_\nu(\max(\bar{U}_T^1, \bar{U}_T^2, \dots, \bar{U}_T^K) \mathbb{1}_{A^c}) \\ &\geq \mu^* TP_\nu(A^c) - \mathbb{E}_\nu(\max(\sum_{t=1}^T U_t^1, \sum_{t=1}^T U_t^2, \dots, \sum_{t=1}^T U_t^K) \mathbb{1}_{A^c}) \\ &\stackrel{(a)}{\geq} \mu^* TP_\nu(A^c) - \mu^* TP_\nu(A^c) - 2\sqrt{T \log(KT)} - 2 \\ &= -2\sqrt{T \log(KT)} - 2. \end{aligned} \quad (47)$$

Here, (a) follows from an argument essentially identical to the one in the proof of Lemma 3 below and we do not repeat it here for brevity. Similarly, we can show that

$$\text{Reg}_T(\pi, \nu', A) \geq -2\sqrt{T \log(KT)} - 2. \quad (48)$$

□

Proof of Lemma 3. We first have that

$$\mathbb{E}_\nu(\max(\bar{U}_T^1, \bar{U}_T^2, \dots, \bar{U}_T^K) \mathbb{1}_A) = \mathbb{E}_\nu(\max(\sum_{t=1}^{n_T^1} U_t^1, \sum_{t=1}^{n_T^2} U_t^2, \dots, \sum_{t=1}^{n_T^K} U_t^K) \mathbb{1}_A) \quad (49)$$

$$\leq \mathbb{E}_\nu(\max(\sum_{t=1}^{T-\lceil \frac{T\Delta}{2} \rceil} U_t^1, \sum_{t=1}^T U_t^2, \dots, \sum_{t=1}^T U_t^K) \mathbb{1}_A). \quad (50)$$

Defining $T_1 = T - \lceil \frac{T\Delta}{2} \rceil$, and $T_i = T$ for all $i > 1$, consider the “good” event

$$G = \left\{ \left| \sum_{t=1}^{T_j} U_t^j - \mu_j T_j \right| \leq \sqrt{T \log(KT)} \text{ for all } j \right\}.$$

Since $U_t^j \in [0, 1]$, by Hoeffding’s inequality, we have that for any $T' \leq T$,

$$\begin{aligned} \mathbb{P}_\nu \left(\left| \sum_{t=1}^{T'} U_t^j - \mu_j T' \right| \leq \sqrt{T \log(KT)} \right) &\geq 1 - 2 \exp\left(-\frac{2(\sqrt{T \log(KT)})^2}{T'}\right) \\ &\geq 1 - 2 \exp\left(-\frac{2(\sqrt{T \log(KT)})^2}{T}\right) \\ &= 1 - \frac{2}{K^2 T^2} \geq 1 - \frac{2}{KT}. \end{aligned}$$

Hence, by the union bound we have that $\mathbb{P}(G) \geq 1 - \frac{2}{T}$. Thus we finally have,

$$\begin{aligned} &\mathbb{E}_\nu \left(\max \left(\sum_{t=1}^{T - \lceil \frac{T\Delta}{2} \rceil} U_t^1, \sum_{t=1}^T U_t^2, \dots, \sum_{t=1}^T U_t^K \right) \mathbb{1}_A \right) \\ &= \mathbb{E}_\nu \left(\max \left(\sum_{t=1}^{T - \lceil \frac{T\Delta}{2} \rceil} U_t^1, \sum_{t=1}^T U_t^2, \dots, \sum_{t=1}^T U_t^K \right) \mathbb{1}_{A, G} \right) \\ &\quad + \mathbb{E}_\nu \left(\max \left(\sum_{t=1}^{T - \lceil \frac{T\Delta}{2} \rceil} U_t^1, \sum_{t=1}^T U_t^2, \dots, \sum_{t=1}^T U_t^K \right) \mathbb{1}_A \mid G^c \right) \mathbb{P}_\nu(G^c) \\ &\leq \mathbb{E}_\nu \left(\left(\max_{i \in [K]} \mu_i T_i + 2\sqrt{T \log(KT)} \right) \mathbb{1}_{A, G} \right) + \frac{2}{T} \times T \\ &\leq \max \left(\left(\frac{1}{2} + \Delta \right) \left(T - \frac{\Delta T}{2} \right), \frac{T}{2} \right) \mathbb{P}_\nu(A) + 2\sqrt{T \log(KT)} + 2 \\ &\stackrel{(a)}{=} \left(\frac{1}{2} + \Delta \right) \left(T - \frac{\Delta T}{2} \right) \mathbb{P}_\nu(A) + 2\sqrt{T \log(KT)} + 2. \end{aligned} \tag{51}$$

Here, (a) follows from the fact that $\Delta < \frac{1}{2}$. Thus, from Equations 50 and 51, we finally have,

$$\begin{aligned} \text{Reg}_T(\pi, \nu, A) &= \left(\frac{1}{2} + \Delta \right) T \mathbb{P}_\nu(A) - \mathbb{E}_\nu \left(\max \left(\bar{U}_T^1, \bar{U}_T^2, \dots, \bar{U}_T^K \right) \mathbb{1}_A \right) \\ &\geq \left(\frac{1}{2} + \Delta \right) T \mathbb{P}_\nu(A) - \left(\frac{1}{2} + \Delta \right) \left(T - \frac{\Delta T}{2} \right) \mathbb{P}_\nu(A) - 2\sqrt{T \log(KT)} - 2 \\ &= \left(\frac{1}{2} + \Delta \right) \frac{\Delta T}{2} \mathbb{P}_\nu(A) - 2\sqrt{T \log(KT)} - 2 \\ &\geq \frac{\Delta T}{4} \mathbb{P}_\nu(A) - 2\sqrt{T \log(KT)} - 2. \end{aligned} \tag{52}$$

□

D Proof of Theorem 3

The proof of Theorem 3 utilizes two technical lemmas. The first one is the following.

Lemma 5. *Let $\delta \in (0, 1)$, and X_1, X_2, \dots , be a sequence of independent 0-mean 1-Sub-Gaussian random variables. Let $\bar{\mu}_t = \frac{1}{t} \sum_{s=1}^t X_s$. Then for any $x > 0$,*

$$P\left(\exists t > 0 : \bar{\mu}_t + \sqrt{\frac{4}{t} \log^+ \left(\frac{1}{\delta t^{3/2}}\right)} + x < 0\right) \leq \frac{39\delta}{x^3}.$$

Its proof is similar to the proof of Lemma 9.3 in Lattimore and Szepesvári (2018), which we present below for completeness.

Proof of Lemma 5. We have,

$$\begin{aligned} & P\left(\exists t > 0 : \bar{\mu}_t + \sqrt{\frac{4}{t} \log^+ \left(\frac{1}{\delta t^{3/2}}\right)} + x < 0\right) \\ &= P\left(\exists t > 0 : t\bar{\mu}_t + \sqrt{4t \log^+ \left(\frac{1}{\delta t^{3/2}}\right)} + tx < 0\right) \\ &\leq \sum_{i=0}^{\infty} P\left(\exists t \in [2^i, 2^{i+1}] : t\bar{\mu}_t + \sqrt{4t \log^+ \left(\frac{1}{\delta t^{3/2}}\right)} + tx < 0\right) \\ &\leq \sum_{i=0}^{\infty} P\left(\exists t \in [0, 2^{i+1}] : t\bar{\mu}_t + \sqrt{2^{i+2} \log^+ \left(\frac{1}{\delta 2^{(i+1) \cdot 3/2}}\right)} + 2^i x < 0\right) \\ &\leq \sum_{i=0}^{\infty} \exp\left(-\frac{\left(\sqrt{2^{i+2} \log^+ \left(\frac{1}{\delta 2^{(i+1) \cdot 3/2}}\right)} + 2^i x\right)^2}{2^{i+2}}\right) \\ &\leq \delta \sum_{i=0}^{\infty} 2^{(i+1) \cdot 3/2} \exp(-2^{i-2} x^2), \end{aligned} \tag{53}$$

where the first inequality follows from a union bound on a geometric grid. The second inequality is used to set up the argument to apply Theorem 9.2 in Lattimore and Szepesvári (2018) and the third inequality is due to its application. The fourth inequality follows from $(a+b)^2 \geq a^2 + b^2$ for $a, b \geq 0$. Then, using a property of unimodal functions ($\sum_{j=c}^d f(j) \leq \max_{i \in [c, d]} f(i) + \int_c^d f(i) di$ for a unimodal function f), the term $2^{(i+1) \cdot 3/2} \exp(-2^{i-2} x^2)$ can be upper bounded by $\frac{42\delta}{e^{3/2} x^3} + \delta \int_0^{\infty} (2^{3/2})^{i+1} \exp(-x^2 2^{i-2}) di$. Evaluating the integral to $\frac{8\sqrt{2\pi}}{\log(2)} \frac{1}{x^3}$, we get

$$P\left(\exists t > 0 : \bar{\mu}_t + \sqrt{\frac{4}{t} \log \frac{1}{\delta t^{3/2}}} + x < 0\right) \leq \frac{39\delta}{x^3}. \tag{54}$$

□

The second result we need is Lemma 8.2 from Lattimore and Szepesvári (2018), which we present below for completeness.

Lemma 6. *Lattimore and Szepesvári (2018) Let X_1, X_2, \dots , be a sequence of independent 0-mean 1-Sub-Gaussian random variables. Let $\bar{\mu}_t = \frac{1}{t} \sum_{s=1}^t X_s$. Let $\epsilon > 0$, and $a > 0$, and define*

$$\kappa = \sum_{t=1}^T \mathbb{1}\{\bar{\mu}_t + \sqrt{\frac{2a}{t}} > \epsilon\}.$$

Then $E[\kappa] \leq 1 + \frac{2}{\epsilon^2} (a + \sqrt{a\pi} + 1)$.

Proof of Theorem 3. Let 1 denote the first arm and i^* denote the arm used in the Commit phase of ADA-ETC. We first define a random variable that quantifies the lowest value of the index of arm 1 can take with respect to its true mean across τ pulls.

$$\Delta \triangleq \left(\mu_1 - \min_{n \leq \tau} \left(\bar{\mu}_n^1 + \sqrt{\frac{4}{n} \log \left(\frac{T}{Kn^{3/2}} \right)} \mathbb{1}_{\{n < \tau\}} \right) \right)^+.$$

The following bound is instrumental for our analysis. For any $x \geq 0$,

$$\begin{aligned} P(\Delta > x) &= P \left(\exists n \leq \tau : \bar{\mu}_n^1 + \sqrt{\frac{4}{n} \log \left(\frac{T}{Kn^{3/2}} \right)} \mathbb{1}_{\{n < \tau\}} < \mu_1 - x \right) \\ &\leq P \left(\exists n < \tau : \bar{\mu}_n^1 + \sqrt{\frac{4}{n} \log \left(\frac{T}{Kn^{3/2}} \right)} < \mu_1 - x \right) + P(\bar{\mu}_\tau^1 < \mu_1 - x) \\ &\stackrel{(a)}{\leq} \min \left(1, \frac{39K}{Tx^3} + \exp(-2\tau x^2) \right) \end{aligned} \quad (55)$$

$$\stackrel{(b)}{\leq} \min \left(1, \frac{40K}{Tx^3} \right). \quad (56)$$

Here, (a) follows from Lemma 5 and Hoeffding's inequality, and (b) follows by the definition of τ and since $\exp(-2\alpha^2/3) \leq 1/\alpha$ for all $\alpha \geq 0$.

We next decompose the regret into the regret from wasted pulls in the Explore phase and the regret from committing to a suboptimal arm in the Commit phase. Let ω be the random time when the Explore phase ends. Let r_ω^i be the reward earned from arm i until time ω . Then the expected regret in the event that $\{i^* = i\}$ is bounded by:

$$\mathbb{E} \left(\left(T\mu_1 - \left(T - \sum_{j \neq i} n_\omega^j - n_\omega^i \right) \mu_i - r_\omega^i \right) \mathbb{1}_{\{i^* = i\}} \right). \quad (57)$$

Note that this expression assumes that the cumulative reward of arm i will be chosen to compete against $T\mu_1$ at the end of time T ; however, if there is an arm with a higher cumulative reward, then the resulting regret can only be lower. Thus the total expected regret is bounded by:

$$\begin{aligned} &\sum_{i=1}^K \mathbb{E} \left(\left(T\mu_1 - \left(T - \sum_{j \neq i} n_\omega^j - n_\omega^i \right) \mu_i - r_\omega^i \right) \mathbb{1}_{\{i^* = i\}} \right) \\ &\stackrel{(a)}{\leq} \sum_{i=1}^K \mathbb{E}(T\Delta_i \mathbb{1}_{\{i^* = i\}}) + \mu_1 \sum_{i=1}^K \mathbb{E}(n_\omega^i \mathbb{1}_{\{i^* \neq i\}}) + \sum_{i=1}^K \mathbb{E} \left((n_\omega^i \mu_i - r_\omega^i) \mathbb{1}_{\{i^* = i\}} \right) \\ &\stackrel{(b)}{=} \sum_{i=1}^K \mathbb{E}(T\Delta_i \mathbb{1}_{\{i^* = i\}}) + \mu_1 \sum_{i=1}^K \mathbb{E}(n_\omega^i \mathbb{1}_{\{i^* \neq i\}}) + \sum_{i=1}^K \mathbb{E} \left((\tau \mu_i - r_\omega^i) \mathbb{1}_{\{i^* = i\}} \right) \\ &= \sum_{i=1}^K \mathbb{E}(T\Delta_i \mathbb{1}_{\{i^* = i\}}) + \mu_1 \sum_{i=1}^K \mathbb{E}(n_\omega^i \mathbb{1}_{\{i^* \neq i\}}) + \sum_{i=1}^K P(i^* = i) (\tau \mu_i - \mathbb{E}(r_\omega^i | i^* = i)) \\ &\stackrel{(c)}{=} \sum_{i=1}^K \mathbb{E}(T\Delta_i \mathbb{1}_{\{i^* = i\}}) + \mu_1 \sum_{i=1}^K \mathbb{E}(n_\omega^i \mathbb{1}_{\{i^* \neq i\}}) + \sum_{i=1}^K P(i^* = i) \left(\tau \mu_i - \sum_{n=1}^{\tau} \mathbb{E}(U_n^i | i^* = i) \right) \\ &\stackrel{(d)}{\leq} \underbrace{\sum_{i=1}^K \mathbb{E}(T\Delta_i \mathbb{1}_{\{i^* = i\}})}_{\text{Regret from misidentifications in Commit phase}} + \underbrace{\mu_1 \sum_{i=1}^K \mathbb{E}(n_\omega^i \mathbb{1}_{\{i^* \neq i\}})}_{\text{Regret from wasted pulls in the Explore phase}}. \end{aligned} \quad (58)$$

Here, (a) results from rearranging terms, and from the fact that $\mu_i \leq \mu_1$. Both (b) and (c) result from the fact that in the event that $\{i^* = i\}$, $n_\omega^i = \tau$. (d) holds since, by a standard stochastic dominance argument, $\tau \mu_i \leq \sum_{n=1}^{\tau} \mathbb{E}(U_n^i | i^* = i)$.

We bound these two terms one by one.

Regret from Explore. First, note that an instance-independent bound on the regret from Explore is simply $K\tau = K\lceil \frac{T^{2/3}}{K^{2/3}} \rceil = O(K^{1/3}T^{2/3})$, which is the maximum number of pulls possible before ADA-ETC enters the Commit phase. Hence, we now focus on deriving an instance-dependent bound. We have that

$$\begin{aligned} \mathbb{E}\left(\sum_{i=1}^K n_{\omega}^i \mathbb{1}_{\{i^* \neq i\}}\right) &\leq \mathbb{E}\left(\sum_{i=2}^K n_{\omega}^i\right) + \tau P(i^* \neq 1) \\ &= \mathbb{E}\left(\sum_{i \geq 2: \Delta \leq \frac{\Delta_i}{2}} n_{\omega}^i\right) + \mathbb{E}\left(\sum_{i \geq 2: \Delta > \frac{\Delta_i}{2}} n_{\omega}^i\right) + \tau P(i^* \neq 1). \end{aligned} \quad (59)$$

We first bound the first term. Define the random variable

$$\eta_i = \sum_{n=1}^{\tau} \mathbb{1}\left\{\bar{\mu}_n^i + \sqrt{\frac{4}{n} \log\left(\frac{T}{Kn^{3/2}}\right)} \mathbb{1}_{\{n < \tau\}} \geq \mu_i + \frac{\Delta_i}{2}\right\}.$$

Then in the event that $\Delta \leq \frac{\Delta_i}{2}$, we have that $n_{\omega}^i \leq \eta_i$. We also have that $n_{\omega}^i \leq \tau$. And thus in the event that $\Delta \leq \frac{\Delta_i}{2}$, we have $n_{\omega}^i \leq \min(\eta_i, \tau)$. Hence the first term above is bounded as:

$$\sum_{i=2}^K P(\Delta \leq \frac{\Delta_i}{2}) \mathbb{E}(\min(\eta_i, \tau)) \leq \sum_{i=2}^K P(\Delta \leq \frac{\Delta_i}{2}) \min(\mathbb{E}(\eta_i), \tau) \leq \sum_{i=2}^K \min(\mathbb{E}(\eta_i), \tau).$$

We can now bound $\mathbb{E}(\eta_i)$ as follows:

$$\begin{aligned} \mathbb{E}(\eta_i) &\leq 1 + \mathbb{E}\left(\sum_{n=1}^{\tau-1} \mathbb{1}\left\{\bar{\mu}_n^i + \sqrt{\frac{4}{n} \log\left(\frac{T}{Kn^{3/2}}\right)} \geq \mu_i + \frac{\Delta_i}{2}\right\}\right) \\ &= 1 + \mathbb{E}\left(\sum_{n=1}^{\tau-1} \mathbb{1}\left\{\bar{\mu}_n^i + \sqrt{\frac{4}{n} \log^+\left(\frac{T}{Kn^{3/2}}\right)} \geq \mu_i + \frac{\Delta_i}{2}\right\}\right) \\ &\stackrel{(a)}{\leq} 1 + \frac{1}{\Delta_i^2} + \mathbb{E}\left(\sum_{n=1}^{\tau-1} \mathbb{1}\left\{\bar{\mu}_n^i + \sqrt{\frac{4}{n} \log^+\left(\frac{T\Delta_i^3}{K}\right)} \geq \mu_i + \frac{\Delta_i}{2}\right\}\right) \\ &\stackrel{(b)}{\leq} 1 + \frac{1}{\Delta_i^2} + \frac{8}{\Delta_i^2} \left(2 \log^+\left(\frac{T\Delta_i^3}{K}\right) + \sqrt{2\pi \log^+\left(\frac{T\Delta_i^3}{K}\right)} + 1\right) \\ &\leq \frac{10}{\Delta_i^2} + \frac{16}{\Delta_i^2} \log^+\left(\frac{T\Delta_i^3}{K}\right) + \frac{24}{\Delta_i^2} \sqrt{\log^+\left(\frac{T\Delta_i^3}{K}\right)}. \end{aligned} \quad (60)$$

Here, (a) is due to lower bounding $1/n^{3/2}$ by Δ_i^3 , and adding $1/\Delta^2$ for the first $1/\Delta^2$ time periods where this lower bound doesn't hold. (b) is due to Lemma 6. The final inequality results from the fact that $\Delta_i \leq 1$ and from trivially bounding $2\pi \leq 9$. Thus, we finally have,

$$\mathbb{E}\left(\sum_{i \geq 2: \Delta \leq \frac{\Delta_i}{2}} n_{\omega}^i\right) \leq \sum_{i=2}^K \min\left(\frac{10}{\Delta_i^2} + \frac{16}{\Delta_i^2} \log^+\left(\frac{T\Delta_i^3}{K}\right) + \frac{24}{\Delta_i^2} \sqrt{\log^+\left(\frac{T\Delta_i^3}{K}\right)}, \tau\right). \quad (61)$$

We now focus on the second term in Equation 59. Note that we have $n_{\omega}^i \leq \tau$, and hence,

$$\mathbb{E}\left(\sum_{i \geq 2: \Delta > \frac{\Delta_i}{2}} n_{\omega}^i\right) \leq \tau \sum_{i=2}^K P(\Delta > \frac{\Delta_i}{2}) \leq \tau \sum_{i=2}^K \min\left(1, \frac{320K}{T\Delta_i^3}\right). \quad (62)$$

Here the second inequality follows from Equation 56. Next, we focus on the third term in Equation 59. We have:

$$\begin{aligned}
 P(i^* \neq 1) &= P(i^* \neq 1 \text{ and } \Delta \leq \frac{\Delta_2}{2}) + P(i^* \neq 1 \text{ and } \Delta > \frac{\Delta_2}{2}) \\
 &\leq \min \left(1, \sum_{i=2}^K P(i^* = i \text{ and } \Delta \leq \frac{\Delta_2}{2}) + P(\Delta > \frac{\Delta_2}{2}) \right) \\
 &\leq \min \left(1, \sum_{i=2}^K P(i^* = i \text{ and } \Delta \leq \frac{\Delta_2}{2}) + \frac{320K}{T\Delta_2^3} \right). \tag{63}
 \end{aligned}$$

Here the final inequality again follows from Equation 56. Now in the event that $\Delta \leq \Delta_2/2$, $i^* = i$ implies that there is some $n \leq \tau$ such that $\text{LCB}_n^i = \bar{\mu}_n^i - \bar{\mu}_n^i \mathbb{1}_{\{n < \tau\}} > \mu_i + \Delta_i/2$. Thus, we have,

$$\begin{aligned}
 \sum_{i=2}^K P(i^* = i \text{ and } \Delta \leq \frac{\Delta_2}{2}) &\leq \sum_{i=2}^K P(\exists n \leq \tau : \bar{\mu}_n^i - \bar{\mu}_n^i \mathbb{1}_{\{n < \tau\}} > \mu_i + \Delta_i/2) \\
 &= \sum_{i=2}^K P(\bar{\mu}_\tau^i > \mu_i + \Delta_i/2) \\
 &\stackrel{(a)}{\leq} \sum_{i=2}^K \exp(-\frac{\tau\Delta_i^2}{2}) \stackrel{(b)}{\leq} \sum_{i=2}^K \frac{8K}{T\Delta_i^3}. \tag{64}
 \end{aligned}$$

Here, (a) follows from Hoeffding's inequality, and (b) follows from the definition of τ and the fact that $\exp(-\alpha^{2/3}/2) \leq 8/\alpha$ for $\alpha \geq 0$. Thus we finally have

$$\begin{aligned}
 \tau P(i^* \neq 1) &\leq \tau \min(1, \sum_{i=2}^K \frac{8K}{T\Delta_i^3} + \frac{320K}{T\Delta_2^3}) \\
 &\leq \tau \min(1, \sum_{i=2}^K \frac{328K}{T\Delta_i^3}). \tag{65}
 \end{aligned}$$

Thus, combining Equations 61, 62, and 65, we have that the regret from the Explore phase is bounded by

$$\begin{aligned}
 &\mu_1 \sum_{i=2}^K \min \left(\frac{10}{\Delta_i^2} + \frac{16}{\Delta_i^2} \log^+ \left(\frac{T\Delta_i^3}{K} \right) + \frac{24}{\Delta_i^2} \sqrt{\log^+ \left(\frac{T\Delta_i^3}{K} \right)}, \tau \right) \\
 &+ \mu_1 \tau \sum_{i=2}^K \min(1, \frac{320K}{T\Delta_i^3}) + \mu_1 \tau \min(1, \sum_{i=2}^K \frac{328K}{T\Delta_i^3}) \\
 &\leq \mu_1 \sum_{i=2}^K \min \left(\frac{10}{\Delta_i^2} + \frac{16}{\Delta_i^2} \log^+ \left(\frac{T\Delta_i^3}{K} \right) + \frac{24}{\Delta_i^2} \sqrt{\log^+ \left(\frac{T\Delta_i^3}{K} \right)}, \tau \right) \\
 &+ \mu_1 \tau \sum_{i=2}^K \min(2, \frac{628K}{T\Delta_i^3}). \tag{66}
 \end{aligned}$$

Here the inequality results from the fact that $\min(1, a) + \min(1, b) \leq \min(2, a + b)$ for $a, b > 0$. This finishes our derivation of a distribution dependent bound on the regret from the Explore phase. We next focus on the regret arising from misidentification in the Commit phase.

Regret from Commit. This regret is upper bounded by

$$\mathbb{E} \left(\sum_{i: \Delta \leq \frac{\Delta_i}{2}} \mathbb{1}_{\{i^* = i\}} T\Delta_i \right) + \mathbb{E} \left(\sum_{i: \Delta > \frac{\Delta_i}{2}} \mathbb{1}_{\{i^* = i\}} T\Delta_i \right). \tag{67}$$

We now get instance dependent and independent bounds on each of the first two terms.

An instance dependent bound on $\mathbf{E}(\sum_{i:\Delta \leq \frac{\Delta_i}{2}} \mathbb{1}_{\{i^*=i\}} T\Delta_i)$. In the event that $\Delta \leq \Delta_i/2$, $i^* = i$ implies that there is some $n \leq \tau$ such that $\text{LCB}_n^i = \bar{\mu}_n^i - \bar{\mu}_n^i \mathbb{1}_{\{n < \tau\}} > \mu_i + \Delta_i/2$. Thus, we have,

$$\mathbf{E}\left(\sum_{i:\Delta \leq \frac{\Delta_i}{2}} \mathbb{1}_{\{i^*=i\}} T\Delta_i\right) \leq \sum_{i=2}^K P(\exists n \leq \tau : \bar{\mu}_n^i - \bar{\mu}_n^i \mathbb{1}_{\{n < \tau\}} > \mu_i + \Delta_i/2) T\Delta_i. \quad (68)$$

Now, we have,

$$P(\exists n \leq \tau : \bar{\mu}_n^i - \bar{\mu}_n^i \mathbb{1}_{\{n < \tau\}} > \mu_i + \Delta_i/2) = P(\bar{\mu}_\tau^i > \mu_i + \Delta_i/2) \leq \exp(-\frac{\tau \Delta_i^2}{2}). \quad (69)$$

Here the final inequality follows from Hoeffding's inequality. Thus we finally have,

$$\mathbf{E}\left(\sum_{i:\Delta \leq \frac{\Delta_i}{2}} \mathbb{1}_{\{i^*=i\}} T\Delta_i\right) \leq \sum_{i=2}^K \exp(-\frac{\tau \Delta_i^2}{2}) T\Delta_i. \quad (70)$$

An instance independent bound on $\mathbf{E}(\sum_{i:\Delta \leq \frac{\Delta_i}{2}} \mathbb{1}_{\{i^*=i\}} T\Delta_i)$. We have

$$\begin{aligned} \mathbf{E}\left(\sum_{i:\Delta < \frac{\Delta_i}{2}} \mathbb{1}_{\{i^*=i\}} T\Delta_i\right) &= T^{2/3} K^{1/3} \sqrt{2 \log K} + \mathbf{E}\left(\sum_{i:\Delta < \frac{\Delta_i}{2}; \Delta_i \geq \frac{K^{1/3} \sqrt{2 \log K}}{T^{1/3}}} \mathbb{1}_{\{i^*=i\}} T\Delta_i\right) \\ &\stackrel{(a)}{\leq} T^{2/3} K^{1/3} \sqrt{2 \log K} + \mathbf{E}\left(\sum_{i:\Delta_i \geq \frac{K^{1/3} \sqrt{2 \log K}}{T^{1/3}}} \exp(-\frac{\tau \Delta_i^2}{2}) T\Delta_i\right) \\ &\stackrel{(b)}{\leq} T^{2/3} K^{1/3} \sqrt{2 \log K} + T^{2/3} K^{1/3} \sqrt{2 \log K}. \end{aligned} \quad (71)$$

Here, (a) follows for the same reason as the derivation of the bound in Equation 70. Next, observe that the function $\exp(-\frac{\tau x^2}{2})x$ is maximized at $x^* = \sqrt{2/\tau} = \sqrt{2}K^{1/3}/T^{1/3}$. But since $\Delta_i \geq \sqrt{2 \log K} K^{1/3}/T^{1/3} \geq \sqrt{2}K^{1/3}/T^{1/3}$, by the unimodality of $\exp(-\frac{\tau x^2}{2})x$, we have

$$\exp(-\frac{\tau \Delta_i^2}{2}) T\Delta_i \leq \exp(-\log K) T^{2/3} K^{1/3} \sqrt{2 \log K} = \frac{1}{K} T^{2/3} K^{1/3} \sqrt{2 \log K}.$$

Hence (b) follows.

An instance dependent bound on $\mathbf{E}(\sum_{i:\Delta > \frac{\Delta_i}{2}} \mathbb{1}_{\{i^*=i\}} T\Delta_i)$.

$$\begin{aligned} \mathbf{E}\left(\sum_{i:\Delta > \frac{\Delta_i}{2}} \mathbb{1}_{\{i^*=i\}} T\Delta_i\right) &\leq \mathbf{E}\left(\max_{i \in [K]} T\Delta_i \mathbb{1}_{\{\Delta > \frac{\Delta_i}{2}\}}\right) \\ &= P(\Delta > \frac{\Delta_K}{2}) T\Delta_K + \sum_{i=1}^{K-1} P(\frac{\Delta_{i+1}}{2} \geq \Delta > \frac{\Delta_i}{2}) T\Delta_i \\ &= P(\Delta > \frac{\Delta_K}{2}) T\Delta_K + \sum_{i=1}^{K-1} \left(P(\Delta > \frac{\Delta_i}{2}) - P(\Delta > \frac{\Delta_{i+1}}{2})\right) T\Delta_i \\ &= \sum_{i=2}^K P(\Delta > \frac{\Delta_i}{2}) T(\Delta_i - \Delta_{i-1}) \\ &\leq \sum_{i=2}^K \min(1, \frac{320K}{T\Delta_i^3}) T(\Delta_i - \Delta_{i-1}). \end{aligned} \quad (72)$$

Here the final inequality again follows from Equation 56.

An instance independent bound on $\mathbf{E}(\sum_{i:\Delta>\frac{\Delta_i}{2}} \mathbb{1}_{\{i^*=i\}} T\Delta_i)$. We have,

$$\mathbf{E}\left(\sum_{i:\Delta>\frac{\Delta_i}{2}} \mathbb{1}_{\{i^*=i\}} T\Delta_i\right) \leq \mathbf{E}(2T\Delta \sum_{i=1}^K \mathbb{1}_{\{i^*=i\}}) = \mathbf{E}(2T\Delta) = 2T\mathbf{E}(\Delta). \quad (73)$$

We then look at $\mathbf{E}(\Delta)$. We have,

$$\mathbf{E}(\Delta) = \int_0^\infty \mathbf{P}(\Delta > x) dx \leq \int_0^\infty \min\left(1, \frac{40K}{Tx^3}\right) dx.$$

This integral evaluates to

$$\int_0^{\frac{(40K)^{1/3}}{T^{1/3}}} dx + \int_{\frac{(40K)^{1/3}}{T^{1/3}}}^\infty \frac{40K}{Tx^3} dx \leq 2\frac{(40K)^{1/3}}{T^{1/3}}.$$

Combining these results, we have

$$\mathbf{E}(\Delta) \leq 2\frac{(40K)^{1/3}}{T^{1/3}}. \quad (74)$$

Thus we finally have,

$$\mathbf{E}\left(\sum_{i:\Delta>\frac{\Delta_i}{2}} \mathbb{1}_{\{i^*=i\}} T\Delta_i\right) \leq 4(40K)^{1/3}T^{2/3}. \quad (75)$$

The final instance-dependent bound follows from Equations 66, 70, and 72. The instance-independent bound follows from the fact that the regret from the Explore phase is at most $K\tau = O(T^{2/3}K^{1/3})$ and from Equations 71 and 75. \square