

## Motivation

**Goal:** Identify the item having the highest averaged return with a given confidence.

**Typical guaranty:** Asymptotic optimality of the expected sample complexity.

⚠ Not informative for moderate confidence level !

📖 This paper: sample complexity upper bounds for **any confidence level** !

## Best-arm identification (BAI)

$K$  arms: arm  $i \in [K]$  is associated with a Gaussian distribution  $\mathcal{N}(\mu_i, 1)$ .

**Goal:** identify  $i^* = \arg \max_{i \in [K]} \mu_i$  with confidence  $1 - \delta \in (0, 1)$ .

**Algorithm:** at time  $n$ ,

- **Sequential test:** if the stopping time  $\tau_\delta$  is reached, then return the candidate answer  $\hat{i}_n$ , else

- **Sampling rule:** pull arm  $I_n$  and observe  $X_n \sim \mathcal{N}(\mu_{I_n}, 1)$ .

**Fixed-confidence:** given an confidence  $\delta \in (0, 1)$ , define a stopping time  $\tau_\delta$  which is  $\delta$ -correct, i.e.  $\mathbb{P}_\mu(\tau_\delta < +\infty, \hat{i}_{\tau_\delta} \neq i^*) \leq \delta$ , and

📖 Minimize the **expected sample complexity**  $\mathbb{E}_\mu[\tau_\delta]$ .

## Lower bound on the expected sample complexity

? What is the best one could achieve ?

📖 Garivier and Kaufmann (2016): For all  $\delta$ -correct algorithms and all Gaussian instances with  $\mu \in \mathbb{R}^K$ ,  $\liminf_{\delta \rightarrow 0} \mathbb{E}_\mu[\tau_\delta] / \log(1/\delta) \geq T^*(\mu)$  where

$$T^*(\mu) = \min_{\beta \in (0,1)} T_\beta^*(\mu) \quad \text{and} \quad T_\beta^*(\mu)^{-1} = \max_{w \in \Delta_K, w_{i^*} = \beta} \min_{j \neq i^*} \frac{1}{2} \frac{(\mu_{i^*} - \mu_j)^2}{1/\beta + 1/w_j}.$$

## TTUCB: UCB-based Top Two sampling rule

**Input:** fixed proportion  $\beta \in (0, 1)$  and function  $g : \mathbb{N} \rightarrow \mathbb{R}^+$ .

Get the **UCB leader**  $B_n = \arg \max_{i \in [K]} \{\mu_{n,i} + \sqrt{g(n)/N_{n,i}}\}$ ;

Get the TC challenger  $C_n \in \arg \min_{i \neq B_n} \frac{(\mu_{n,B_n} - \mu_{n,i})_+}{\sqrt{1/N_{n,B_n} + 1/N_{n,i}}}$ ;

Use **tracking** to get  $I_n = B_n$  if  $N_{n,B_n} \leq \beta L_{n+1,B_n}$ , otherwise  $I_n = C_n$ ;

**Output:** next arm to sample  $I_n$ .

$(N_{n,i}, \mu_{n,i})$ : number of pulls and empirical mean of arm  $i$  before time  $n$ .

$L_{n,i}$ : number of selection of arm  $i$  as leader before time  $n$ .

$N_{n,j}^i$ : number of pulls of arm  $j$  when arm  $i$  is leader before time  $n$ .

- Take  $\beta = 1/2$  since  $w^*(\mu)_{i^*} \leq 1/2$  and  $T_{1/2}^*(\mu)/T^*(\mu) \ll 2$  for most instances.
- Choose small  $g$  s.t.  $\mathbb{P}_\mu(\mathcal{E}_n) \geq 1 - Kn^{-s}$  with

$$\mathcal{E}_n = \{\forall (t, i) \in [n^{1/\alpha}] \times [K], \mu_i \in [\mu_{t,i} \pm \sqrt{g(t)/N_{t,i}}]\}$$

where  $\alpha, s > 1$ , e.g.  $g_u(n) = 2\alpha(1+s) \log n$ .

## $\delta$ -correct sequential test

? How to obtain a  $\delta$ -correct sequential test for Gaussian distributions ?

📖 **GLR stopping rule:** recommend  $\hat{i}_n \in \arg \max_{i \in [K]} \mu_{n,i}$  and stop at time

$$\tau_\delta = \inf\{n > K \mid \min_{i \neq \hat{i}_n} \frac{\mu_{n,\hat{i}_n} - \mu_{n,i}}{\sqrt{1/N_{n,\hat{i}_n} + 1/N_{n,i}}} \geq \sqrt{2c(n-1, \delta)}\}, \quad (1)$$

with  $c(n, \delta) \simeq \log(1/\delta) + 2 \log \log(1/\delta) + 4 \log(4 + \log(n/2))$ .

## Asymptotic confidence guarantees

**Theorem 1.** Let  $(\delta, \beta) \in (0, 1)^2$ . Combined with GLR stopping (1), the TTUCB algorithm is  $\delta$ -correct and asymptotically  $\beta$ -optimal for all  $\mu \in \mathbb{R}^K$  having distinct means, i.e. it satisfies  $\limsup_{\delta \rightarrow 0} \mathbb{E}_\mu[\tau_\delta] / \log(1/\delta) \leq T_\beta^*(\mu)$ .

**Limitations:** no guarantees (1) for **moderate regime** of  $\delta$  and (2) when sub-optimal arms share the **same mean**.

## Finite confidence guarantees

**Theorem 2.** Let  $\delta \in (0, 1)$ . Combined with GLR stopping (1), the TTUCB algorithm using  $\beta = 1/2$  and  $g_u$  with  $\alpha = s = 1.2$  satisfies that, for all  $\mu \in \mathbb{R}^K$  such that  $|i^*(\mu)| = 1$ ,

$$\mathbb{E}_\mu[\tau_\delta] \leq \inf_{x \in [0, (K-1)^{-1}]} \max\{T_0(\delta, x), C_\mu^{1,2}, C_0(x)^6, (2/\varepsilon)^{1.2}\} + 12K,$$

where  $\varepsilon \in (0, 1]$  and

$$C_\mu = \mathcal{O}(H(\mu) \log H(\mu)) \quad \text{with} \quad H(\mu) = 2\Delta_{\min}^{-2} + \sum_{i \neq i^*} 2(\mu_{i^*} - \mu_i)^{-2},$$

$$\limsup_{\delta \rightarrow 0} T_0(\delta, 0) / \log(1/\delta) \leq 2T_{1/2}^*(\mu),$$

$$C_0(x) = 2/(\varepsilon a_\mu(x)) + 1 \quad \text{with} \quad a_\mu(x) = (1-x)^{d_\mu(x)} \max_{i \neq i^*} \{w_{1/2}^*(\mu)_i, x/2\}$$

and  $d_\mu(x) = |\{i \neq i^* \mid w_{1/2}^*(\mu)_i < x/2\}|$ .

**Refined analysis:** Clipping  $\min_{i \neq i^*} w_{1/2}^*(\mu)_i$  by  $x/2$  yields  $C_0(x) = \mathcal{O}(K/\varepsilon)$ .

📖 Generic method that improves the analysis of APT (Locatelli et al, 2016).

**Table 1:** Upper bound on the sample complexity  $\tau_\delta$  in probability (§) or in expectation (†). The notation  $\tilde{\mathcal{O}}$  hides polylogarithmic factors. (\*) Upper bound on  $\mathbb{E}_\mu[\tau_\delta \mathbb{1}(\mathcal{E})]$  where  $\mathbb{P}[\mathcal{E}^c] \leq \gamma$ . (\*\*) Asymptotic bound holds for instances with distinct means. Ordered references: Kalyanakrishnan et al. (2012), Karnin et al. (2013), Jamieson et al. (2014), Degenne et al. (2019), Katz-Samuels et al. (2020), Wang et al. (2021), Barrier et al. (2022).

Algorithm	Asymptotic $\delta \rightarrow 0$	Finite $\delta$ when $H(\mu) \rightarrow +\infty$
LUCB1†	$\mathcal{O}(H(\mu) \log(1/\delta))$	$\mathcal{O}(H(\mu) \log H(\mu))$
Exp-Gap§	$\mathcal{O}(H(\mu) \log(1/\delta))$	$\mathcal{O}(\sum_{i \neq i^*} \Delta_i^{-2} \log \log \Delta_i^{-1})$
lil' UCB§	$\mathcal{O}(H(\mu) \log(1/\delta))$	$\mathcal{O}(\sum_{i \neq i^*} \Delta_i^{-2} \log \log \Delta_i^{-1})$
DKM†	$T^*(\mu) \log(1/\delta) + \tilde{\mathcal{O}}(\sqrt{\log(1/\delta)})$	$\tilde{\mathcal{O}}(KT^*(\mu)^2)$
Peace§	$\mathcal{O}(T^*(\mu) \log(1/\delta))$	$\mathcal{O}(H(\mu) \log(K/\Delta_{\min}))$
FWS†	$T^*(\mu) \log(1/\delta) + \mathcal{O}(\log \log(1/\delta))$	$\mathcal{O}(e^K H(\mu)^{19/2})$
EBS†*	$T^*(\mu) \log(1/\delta) + o(1)$	$\mathcal{O}(KH(\mu)^4/w_{\min}^2)$
<b>TTUCB†**</b>	$T_\beta^*(\mu) \log(1/\delta) + \mathcal{O}(\log \log(1/\delta))$	$\mathcal{O}((H(\mu) \log H(\mu))^\alpha)$

## Tracking instead of randomization

- **Fully deterministic** algorithm.
- Deterministic counts simplifies the non-asymptotic analysis.
- Faster convergence of  $N_{n,i^*}/n$  to  $\beta$ , at least in  $\mathcal{O}(1/n)$  instead of  $\mathcal{O}(1/\sqrt{n})$ .

**Lemma 1.** For all  $n > K$  and all  $i \in [K]$ , we have  $-1/2 \leq N_{n,i}^i - \beta L_{n,i} \leq 1$ .

## Generic regret minimizing leader

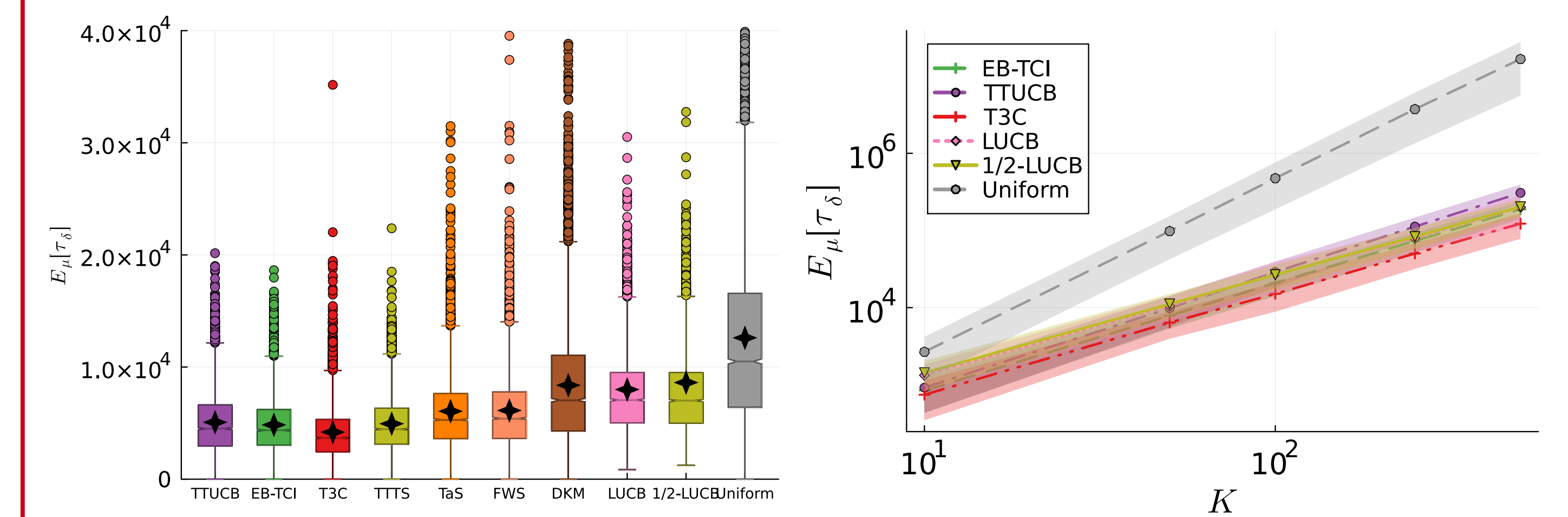
The Top Two method is a generic wrapper to convert any regret minimization algorithm into a best arm identification strategy.

**Sufficient condition:** Arm  $i^*$  is leader except for a sublinear number of times.

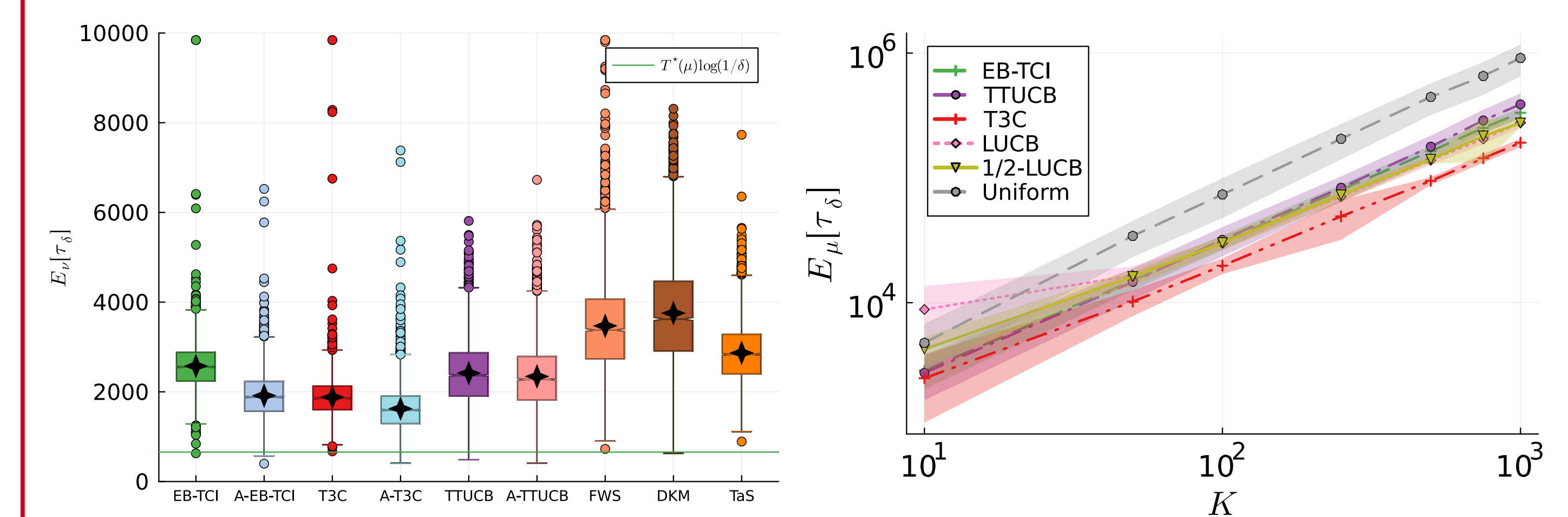
📖 Upper bound  $(N_{n,i})_{i \neq i^*}$  or  $\sum_{i \neq i^*} \Delta_i N_{n,i}$  under a concentration event.

**Lemma 2 (UCB).** Under  $\mathcal{E}_n$ , we have  $L_{n,i^*} \geq n - 24H(\mu) \log n - 2K - 1$ .

## Experiments



**Figure 1:** Empirical stopping time for  $\delta = 0.1$  on (a) random instances with  $\mu_1 = 0.6$  and  $\mu_i \sim \mathcal{U}([0.2, 0.5])$  for  $i \neq 1$  ( $K = 10$ ) and (b) instances  $\mu_i = 1 - \left(\frac{i-1}{K-1}\right)^{0.6}$  with  $H(\mu) = \Theta(K^{1.2})$ .



**Figure 2:** Empirical stopping time for  $\delta = 0.1$  on "1-sparse" instances: (a)  $(K, \mu_{i^*}, \Delta) = (35, 0, 0.5)$  with  $T_{1/2}^*(\mu)/T^*(\mu) \approx 3/2$  and (b)  $(\mu_{i^*}, \Delta) = (0, 0.25)$  with  $H(\mu) = \Theta(K)$ . Constant  $\beta = 1/2$  and adaptive proportions (A-), IDS (You et al., 2023) sets  $\beta_n = N_{n,C_n}/(N_{n,C_n} + N_{n,B_n})$ .

## Conclusion

1. First non-asymptotic analysis of Top Two algorithms, which holds for instances having a unique best arm.
2. Deterministic asymptotically  $\beta$ -optimal Top Two algorithm using UCB leader and tracking instead of randomization.