

# Fast and Regret-Optimal Best Arm Identification: Fundamental Limits and Low-Complexity Algorithms

Qining Zhang, Ph.D. Candidate  
EECS, University of Michigan, Ann Arbor

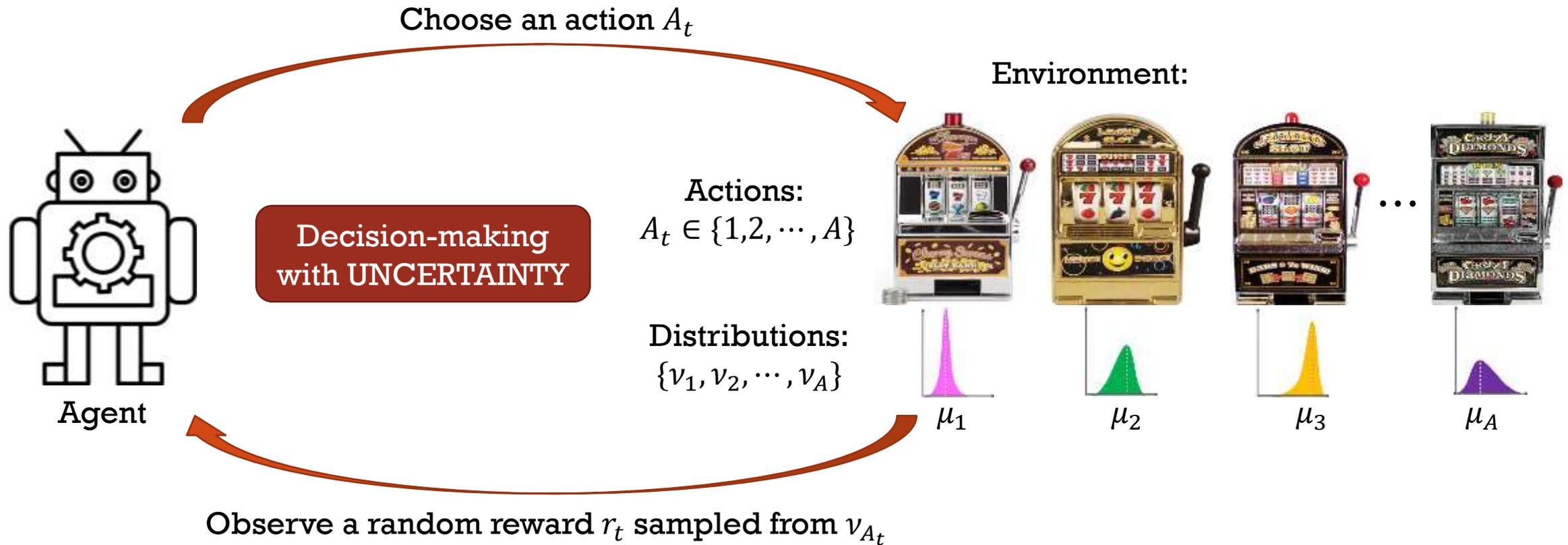
Joint work with my advisor Lei Ying (Michigan)

NeurIPS 2023

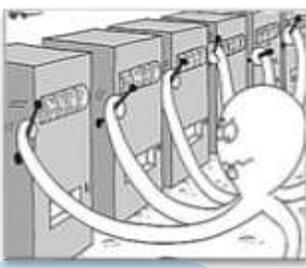


# Multi-armed Bandits

- Online decision making for  $T$  slots.



# Classic Views of MAB



- **Regret Minimization** (Lai and Robbins. 1985)

- Exploration v.s. Exploitation

$$\text{Reg}_\mu(T) = T\mu_{a^*} - \mathbb{E}_\mu \left[ \sum_{t=1}^T \mu_{A_t} \right]$$

Does not commit to any arm

- **Best Arm Identification** (Garivier, et al. 2016)

- Sample Complexity

$$\min \tau \text{ s.t. } \Pr(\hat{a}_\tau = a^*) \geq 1 - \delta$$

Over-exploration of suboptimal arms

- **Online MAB** (Auer et al. 2002)

- **Optimism**

UCB bonus

$$A_t = \operatorname{argmax}_a \bar{r}_{t-1}(a) + \sqrt{\frac{\text{clog}(T)}{N_{t-1}(a)}}$$

- **Offline MAB** (Rashidinejad et al. 2021)

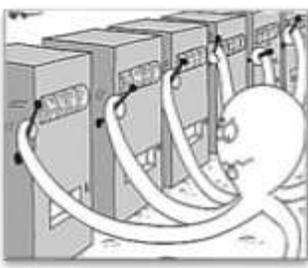
- **Pessimism**

LCB penalty

$$\hat{a} = \operatorname{argmax}_a \bar{r}_N(a) - \sqrt{\frac{\text{clog}(N)}{N(a)}}$$

What is the fundamental difference between online and offline data?

# What Happens in Real-World Applications



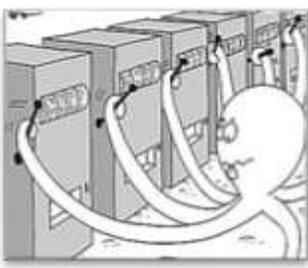
Clinic trials



Career choices



# Regret Optimal Best Arm Identification



- Two Goals

- Optimal cumulative regret.
- Commit to optimal action quickly.

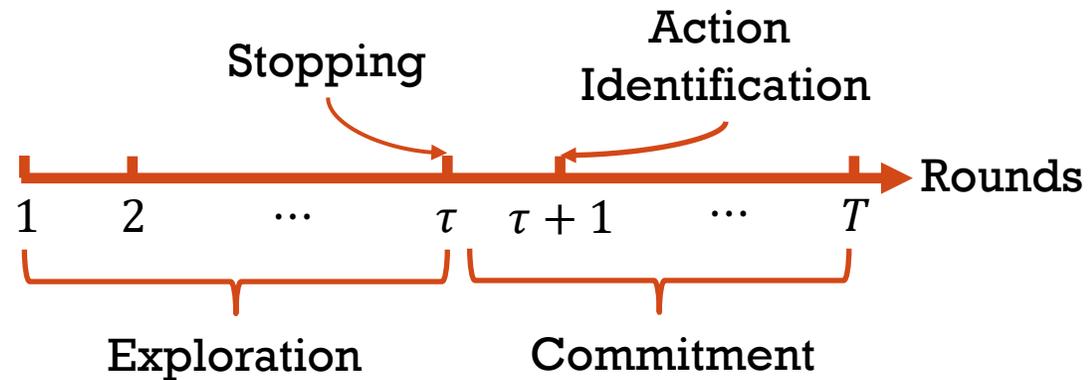
- Three Components

- Exploration
- Stopping
- Action Identification

→  $\min_{\pi \in \Pi_{RO}} \mathbb{E}[\tau]$  such that  $\Pr(\hat{a} \neq a^*) = \mathcal{O}(T^{-1})$ ,

where

→  $\Pi_{RO} = \left\{ \pi: \limsup_{T \rightarrow \infty} \frac{\text{Reg}_{\mu}^{\pi}(T)}{\log T} = \sum_{a \neq a^*} \frac{\Delta_a}{\text{KL}(\mu_a, \mu_{a^*}^*)} \right\}$ .



Can we design an algorithm for ROBAI?  
What are the fundamental limits?

# EOCP: Explore Optimistically then Commit Pessimistically

- Exploration
  - Modified-UCB.

$$A_t = \operatorname{argmax}_a \bar{r}_{t-1}(a) + \sqrt{\frac{2l}{N_{t-1}(a)}}$$

- Action Identification
  - Modified-LCB algorithm.

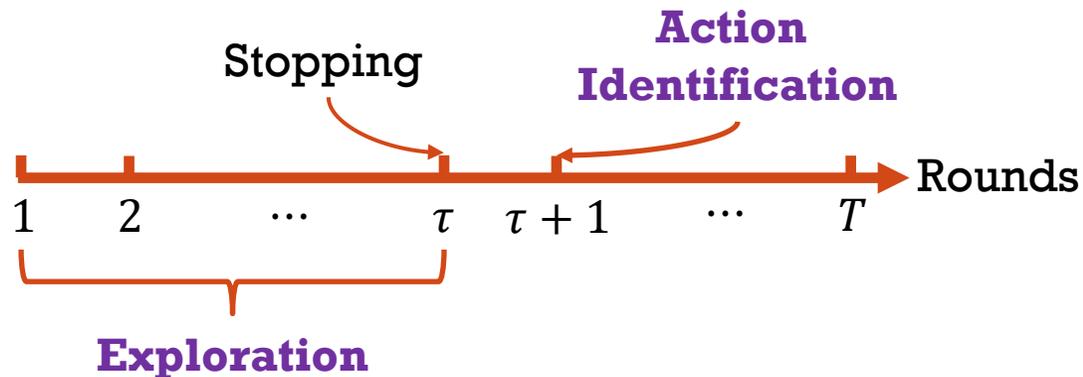
$$\hat{a} = \operatorname{argmax}_a \bar{r}_{t-1}(a) - \sqrt{\frac{2l}{N_{t-1}(a)}}$$

- Stopping
  - Pre-determined Stopping
  - Adaptive Stopping

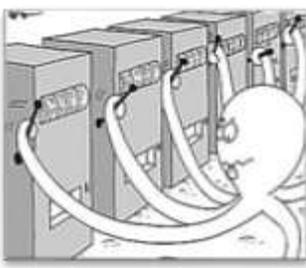
$$\text{Stop when } t = \frac{8Al}{\Delta_{\min}^2}$$

Lower bound on the minimum reward gap

$$\text{Stop when: } \max_a \min_{a'} N_t(a) - lN_t(a') > 1$$



# Main Results



- **Pre-determined Stopping Time:**

Theorem 1. (Pre-determined)

If we choose  $l = \log(T) + c_0\sqrt{\log(T)}$ , the regret of EOCP is bounded:

$$\text{Reg}_\mu^{\text{EOCP}}(T) \leq \sum_{a:\Delta_a>0} \frac{2}{\Delta_a} \log T + o(\log T).$$

And  $\text{SCC}_\mu^{\text{EOCP}}(T) = \mathcal{O}(\Delta_{\min}^{-2} \log T)$ .

- **Adaptive Stopping Time:**

Constant Optimal Regret

Theorem 2. (Adaptive)

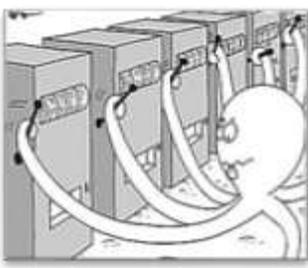
If we choose  $l = \log(T) + c_0\sqrt{\log(T)}$ , the regret of EOCP-UG is bounded:

$$\text{Reg}_\mu^{\text{EOCP}}(T) \leq \sum_{a:\Delta_a>0} \frac{2}{\Delta_a} \log T + o(\log T).$$

And  $\text{SCC}_\mu^{\text{EOCP}}(T) = \mathcal{O}(\Delta_{\min}^{-2} \log^2 T)$ .

Sample Complexity Loss

# Fundamentality



- **Commitment Time Limits for Regret Optimal Algorithms**

Theorem 3 (Informal).

For 2-armed Gaussian bandit, for any algorithm  $\pi$  with regret is  $\mathcal{O}(\log^c T)$  away from optimal, in pre-determined setting:

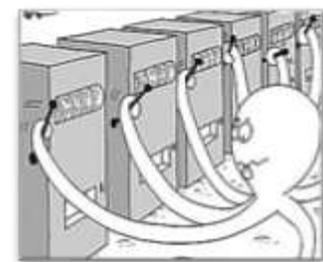
$$SCC_{\mu}^{\pi}(T) = \Omega\left(\frac{\log(T)}{\Delta^2}\right),$$

in adaptive setting:

$$SCC_{\mu}^{\pi}(T) = \Omega\left(\frac{\log^{2-c}(T)}{\Delta^2}\right).$$

- **EOCP matches the LB when  $\Delta$  is known a priori.**

# Comparison to Literature



Bandit Algorithm	Regret	Sample Complexity	Confidence
UCB <sub>(Auer, et al. 2002)</sub>	$\frac{2}{\Delta} \log(T)$	$T$	N/A
TS <sub>(Thompson. 1933)</sub>	$\frac{2}{\Delta} \log(T)$	$T$	N/A
BAI-ETC <sub>(Garivier, et al. 2016)</sub>	$\frac{4}{\Delta} \log(T)$	$\mathcal{O}\left(\frac{\log(T)}{\Delta^2}\right)$	$\tilde{\mathcal{O}}(T^{-1})$
EOCP <sub>(Ours)</sub>	$\frac{2}{\Delta} \log(T)$	$\mathcal{O}\left(\frac{\log(T)}{\Delta^2}\right)$	$\mathcal{O}(T^{-1})$
EOCP-UG <sub>(Ours)</sub>	$\frac{2}{\Delta} \log(T)$	$\mathcal{O}\left(\frac{\log^2(T)}{\Delta^2}\right)$	$\mathcal{O}(T^{-1})$
KL-EOCP <sub>(Ours)</sub>	$\frac{\Delta}{\text{KL}(\mu_2, \mu_1)} \log T$	$\mathcal{O}\left(\frac{\log(T)}{\text{KL}(\mu_2, \mu_1)}\right)$	$\mathcal{O}(T^{-1})$
Lower Bound (Gaussian)	$\frac{2}{\Delta} \log(T)$	$\mathcal{O}\left(\frac{\log(T)}{\Delta^2}\right)$	$\mathcal{O}(T^{-1})$
Lower Bound (General)	$\frac{\Delta}{\text{KL}(\mu_2, \mu_1)} \log T$	$\mathcal{O}\left(\frac{\log(T)}{\text{KL}(\mu_2, \mu_1)^2}\right)$	$\mathcal{O}(T^{-1})$

# Thanks!

# Questions?

<https://arxiv.org/abs/2309.00591>