

Le site qu'il te faut pour réviser. **Gratuitement.**

6ème

5ème

4ème

3ème

Seconde

1ère L

1ère ES

1ère S

Bac L

Bac ES

Bac S



SEQUENTIAL MACHINE LEARNING FOR ADAPTIVE EDUCATIONAL SYSTEMS

PhD defense of Julien SEZNEC, the 15th of December 2020

Supervised by:

Michal VALKO (Inria Lille / Deepmind)

Alessandro LAZARIC (Inria Lille / FAIR)

Jonathan BANON (Lelivrescolaire.fr)



lelivrescolaire.fr

OUTLINE

- 1. Bandits for Adaptive Educational Systems**
- 2. Rested Rotting Bandits**
- 3. Restless Rotting Bandits**

AFTERCLASSE



DIFFICULTÉ

PASSER

SUivant >

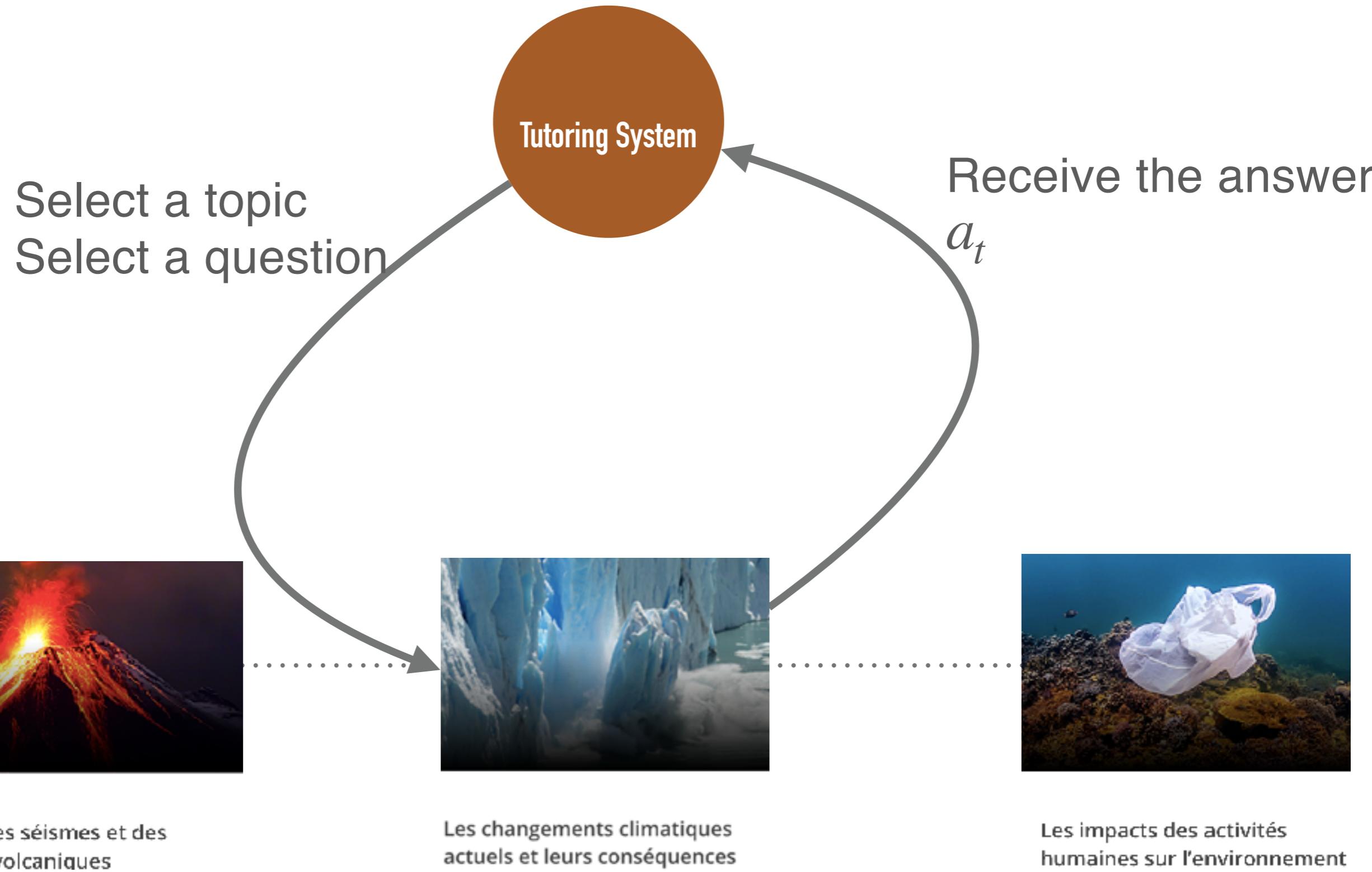
Dans une classe de 15 garçons et 11 filles, on choisit un élève au hasard. Classe ces probabilités selon qu'elles sont correctes ou incorrectes.

$$\text{Probabilité de choisir une fille} = \frac{15}{26}$$

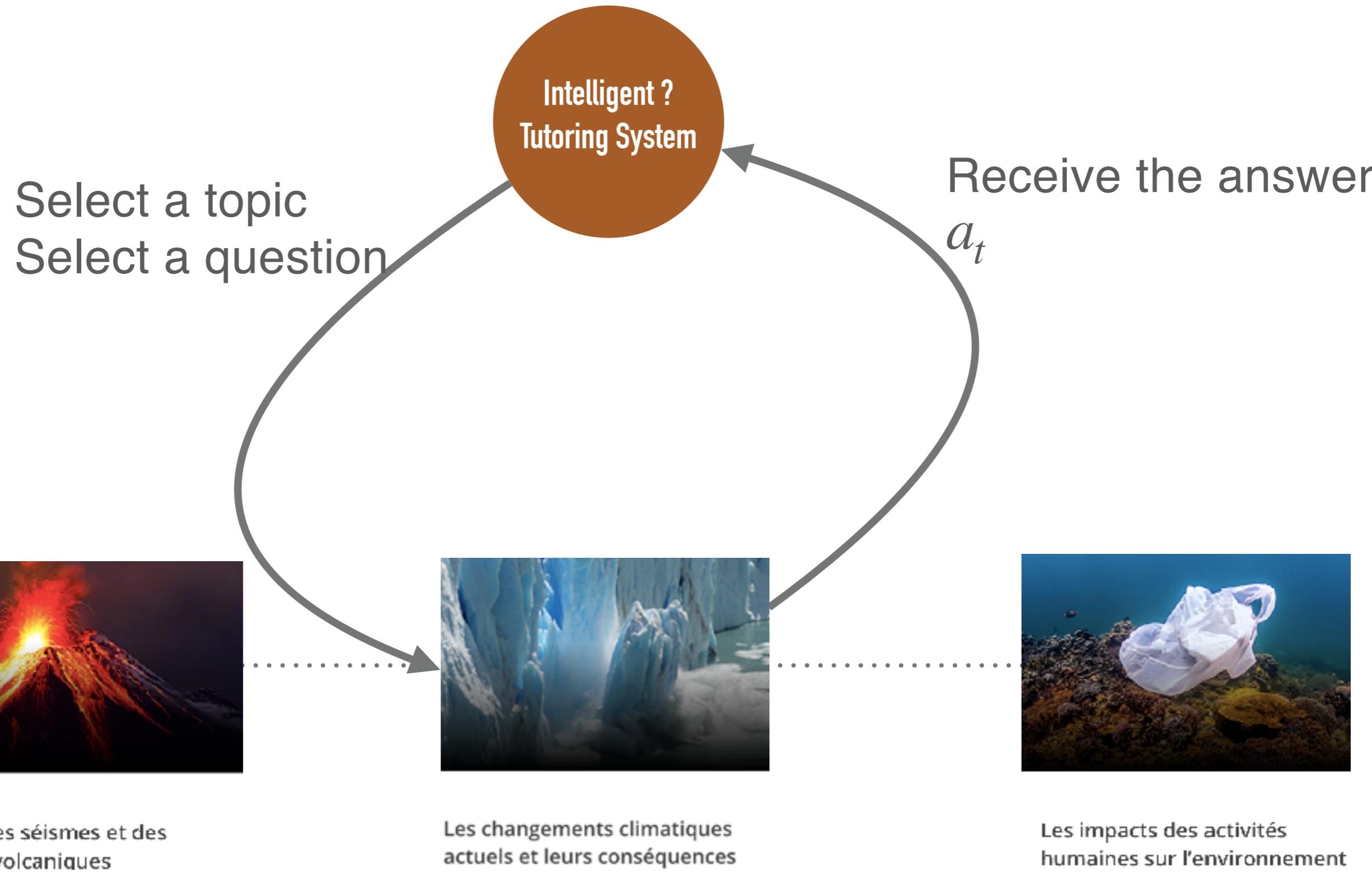
1 Probabilités incorrectes

2 Probabilités correctes

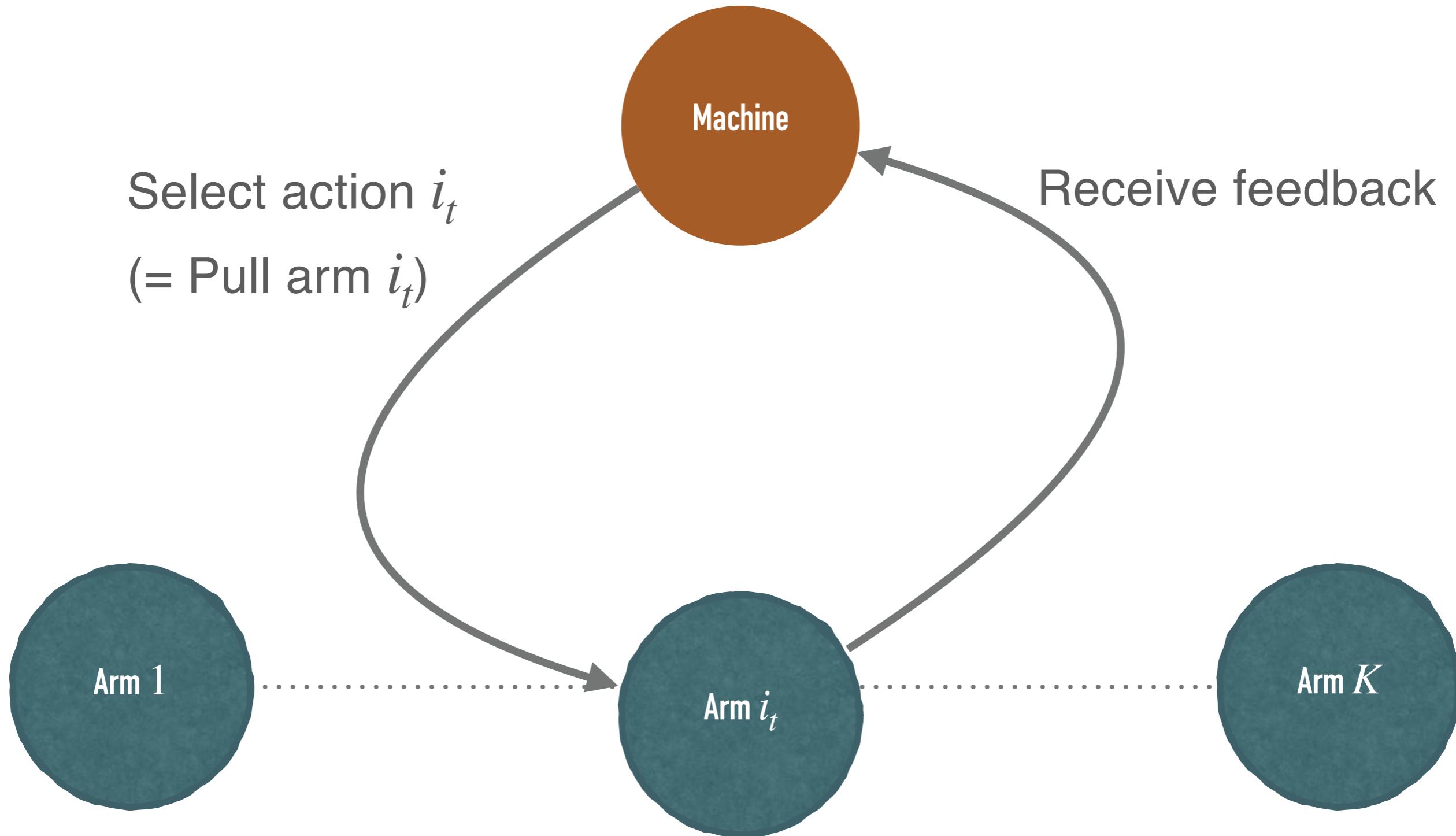
SEQUENTIAL ADAPTIVE EDUCATIONAL SYSTEMS



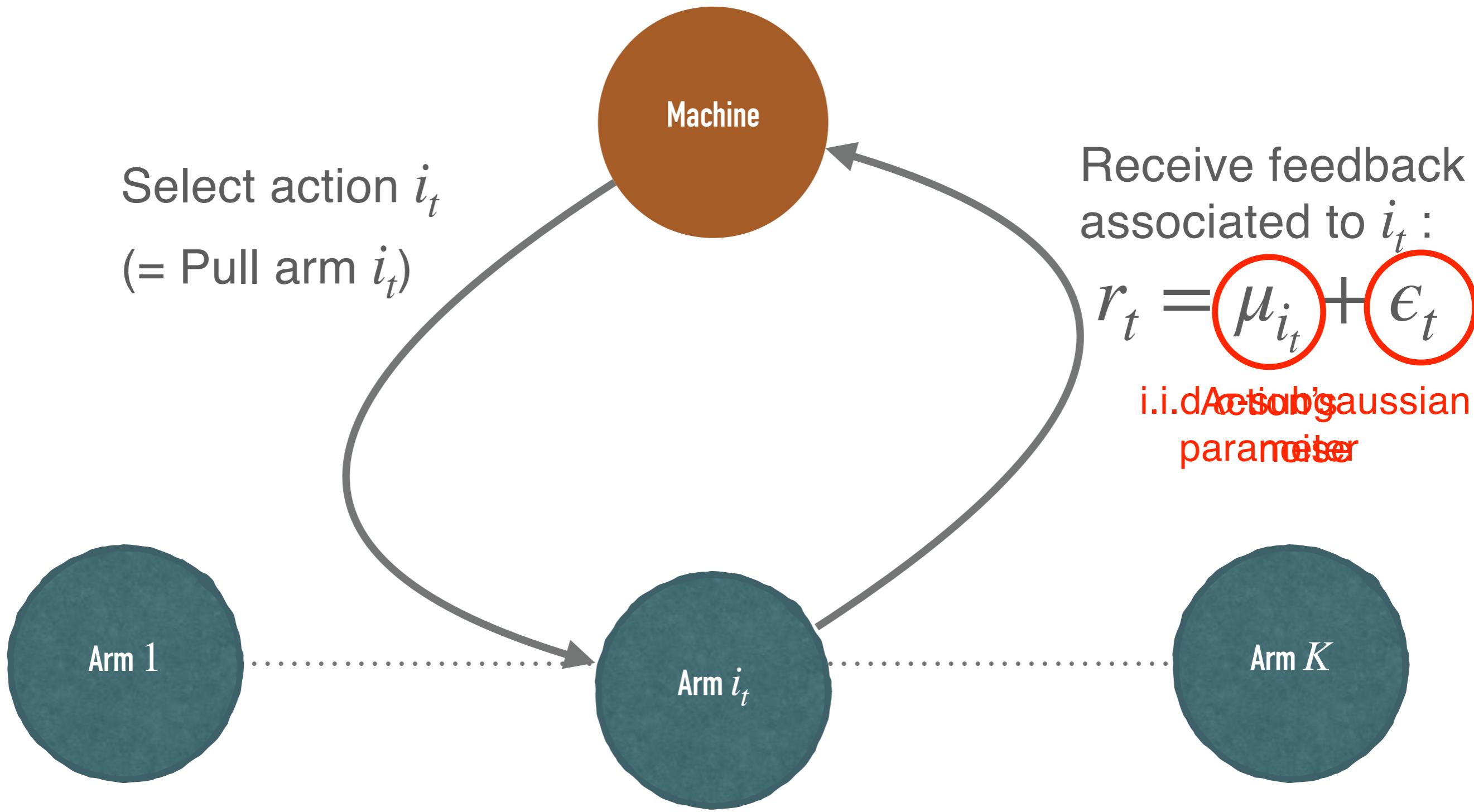
SEQUENTIAL ADAPTIVE EDUCATIONAL SYSTEMS



SEQUENTIAL EXPLORATION: BANDITS FEEDBACK



STATIONARY BANDITS: FEEDBACK



STATIONARY BANDITS: OBJECTIVE

- ▶ Standard view : Feedback = reward
- ▶ Objective : Maximize cumulative reward in expectation at the horizon T

$$J_T(\pi) = \mathbb{E} \left[\sum_{t \leq T} r_t \right] = \mathbb{E} \left[\sum_{t \leq T} \mu_{i_t} \right]$$

- ▶ Optimal oracle policy

$$\pi^\star(t) \in \operatorname{argmax}_{i \leq K} \mu_i$$

- ▶ Regret

$$R_T(\pi) = J_T(\pi^\star) - J_T(\pi) = \sum_{i \in \mathcal{K}} N_{i,T} \Delta_i$$

$$\left\{ \begin{array}{l} \Delta_i = \max_j \mu_j - \mu_i \\ N_{i,T} \text{ Number of pulls of arm } i \text{ at round } T \end{array} \right.$$

UPPER CONFIDENCE BOUND POLICY

Select the arm with the largest index:

$$\hat{\mu}_i + \sqrt{\frac{2\sigma^2 \log 1/\delta}{N_{i,t}}}$$

Expected reward + cost of information

Optimal problem-independent rate:

$$\Theta(\sigma\sqrt{KT})$$

Optimal problem-dependent rate:

$$\sum_{i, \Delta_i > 0} \frac{2\sigma \log(T)}{\Delta_i}$$

$$\delta = \max\left(\frac{KN_{i,t}}{t}, 1\right)$$

BANDITS IN EDUCATIONAL SYSTEMS

1. What is the objective ? What is the reward ?
2. What is the « best » action ?
3. How do actions impact observations ?
4. What can we learn with ~100 samples ?

BANDITS IN EDUCATIONAL SYSTEMS

Clément et al.
(2016)

EXP3

Target the topic
with the largest
improvement

$$r_t = \frac{1}{h} \sum_{i=0}^{h-1} (a_{t-i} - a_{t-h-i})$$

Melesko et al.
(2019)

UCB

Target the most
difficult topic

$$r_t = -a_t$$

Lin et al. (2020)

Thomson sampling

Target the most
difficult topic

$$r_t = -a_t$$

BANDITS IN EDUCATIONAL SYSTEMS

Clément et al.
(2016)

EXP3

Target the topic
with the largest
improvement

$$r_t = \frac{1}{h} \sum_{i=0}^{h-1} (a_{t-i} - a_{t-h-i})$$

Melesko et al.
(2019)

UCB

Target the most
difficult topic

$$r_t = -a_t$$

Lin et al. (2020)

Thomson sampling

Target the most
difficult topic

$$r_t = -a_t$$

OUTLINE

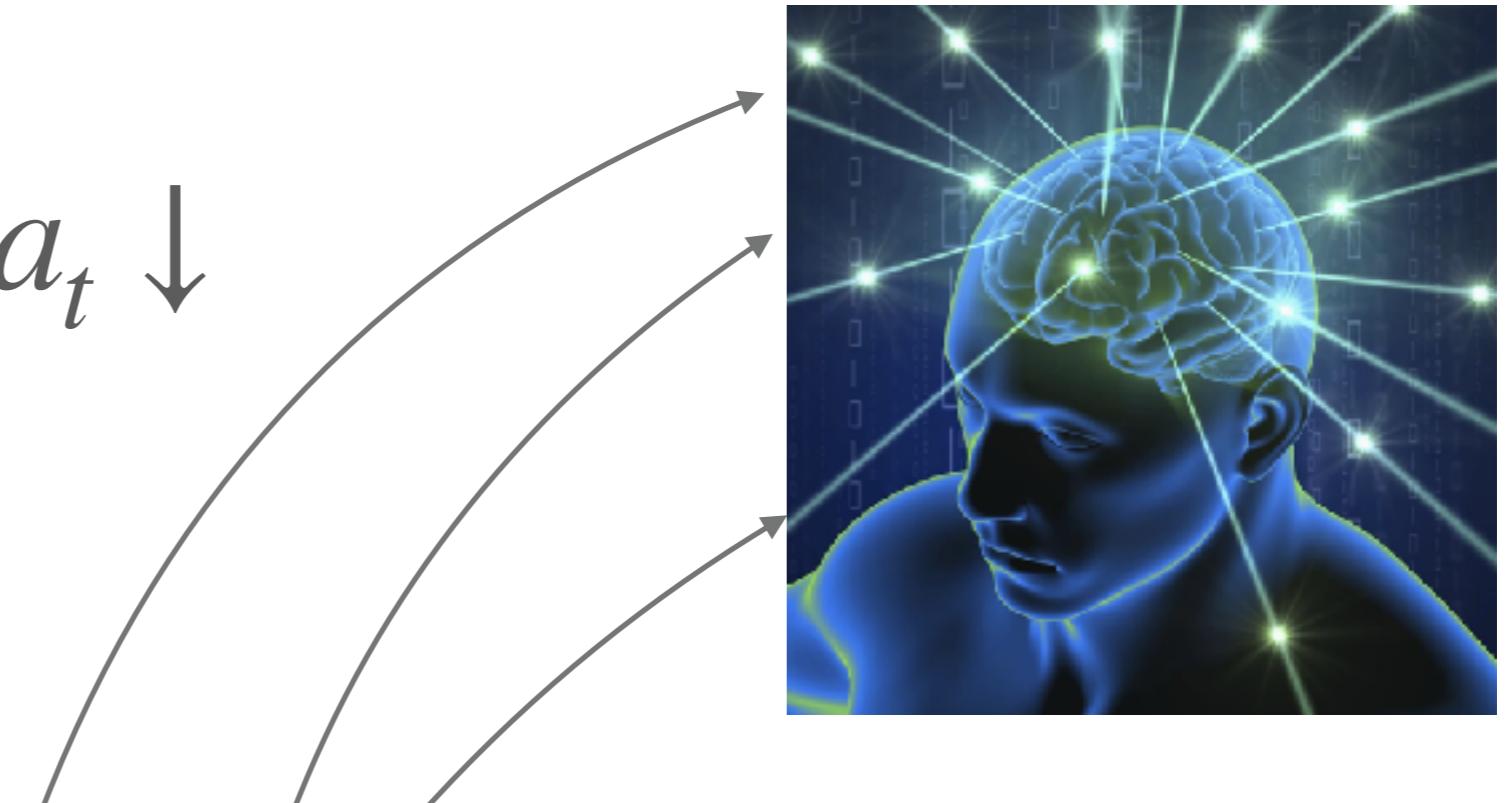
1. Bandits for Adaptive Educational Systems

2. Rested Rotting Bandits

3. Restless Rotting Bandits

WHEN BANDITS GO ROTTING ...

$$a_t \uparrow \implies r_t = -a_t \downarrow$$



L'origine des séismes et des éruptions volcaniques



Les changements climatiques actuels et leurs conséquences



Les impacts des activités humaines sur l'environnement

RESTED vs. RESTLESS BANDITS

Rested :

- Actions trigger non-stationarity
- No action = no change = rested
- $\mu_i(N_{i,t})$

Restless :

- Changes happen independently of the actions
- Action can change when they are not selected
- $\mu_i(t)$

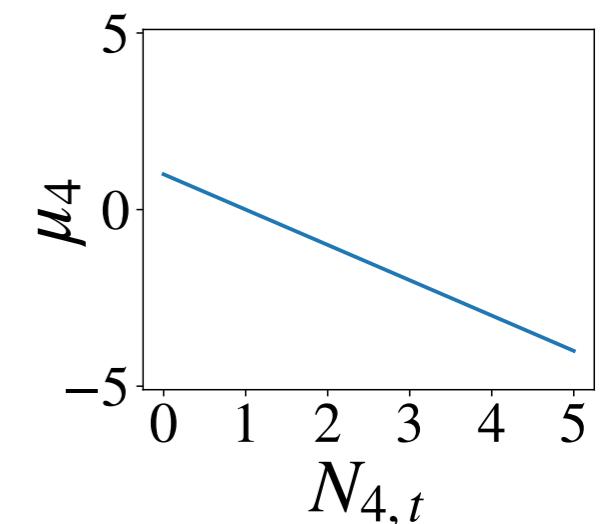
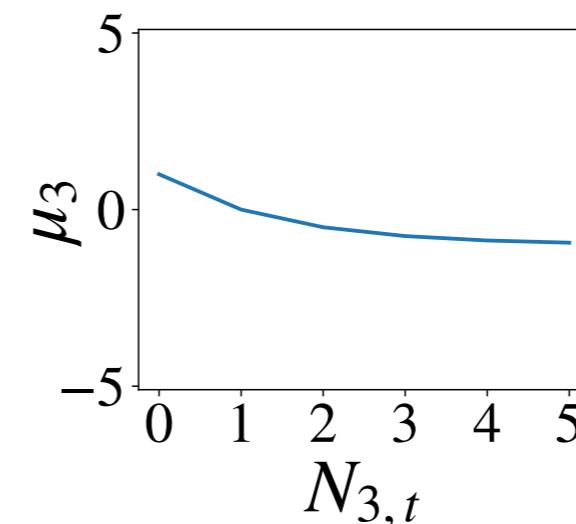
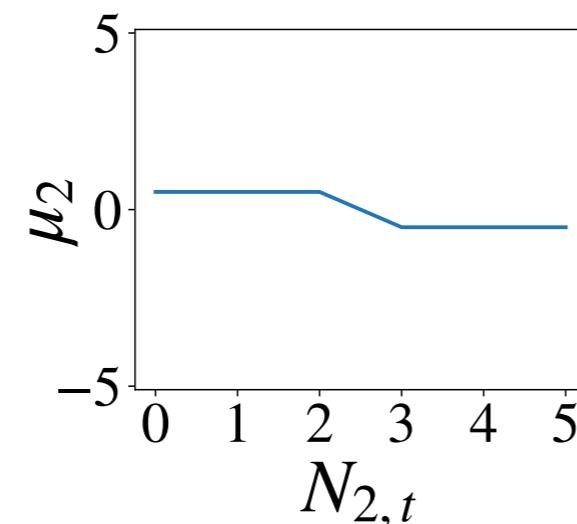
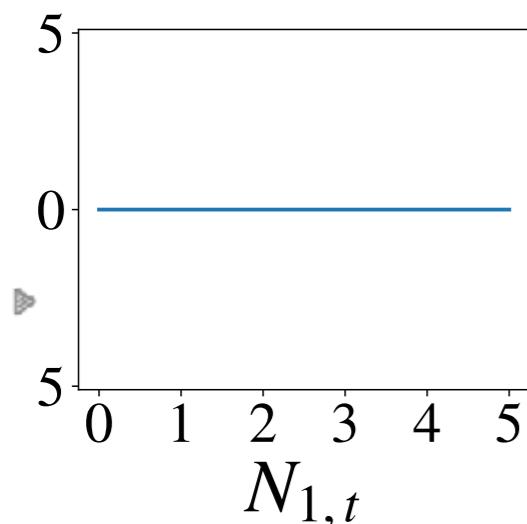
RESTED ROTTING BANDITS ARE ...

Stochastic bandits ...

- ▶ K arms
- ▶ At each round t , agent pulls arm i and receives a noisy reward $r_t \leftarrow \mu_i + \varepsilon_t$ (ε_t i.i.d. ; σ -sub-gaussian)
- ▶ Maximize cumulative reward: $J_T(\pi) = \mathbb{E} \left[\sum_{t \leq T} r_t \right]$

... with rotting arms

- ▶ $\{\mu_i\}$ are **non-increasing** functions of $N_{i,t}$ the **number of pulls of arm i** at time t
- ▶ $L \triangleq \max_{i \in K} \max_{n \leq T} \mu_i(n) - \mu_i(n+1)$

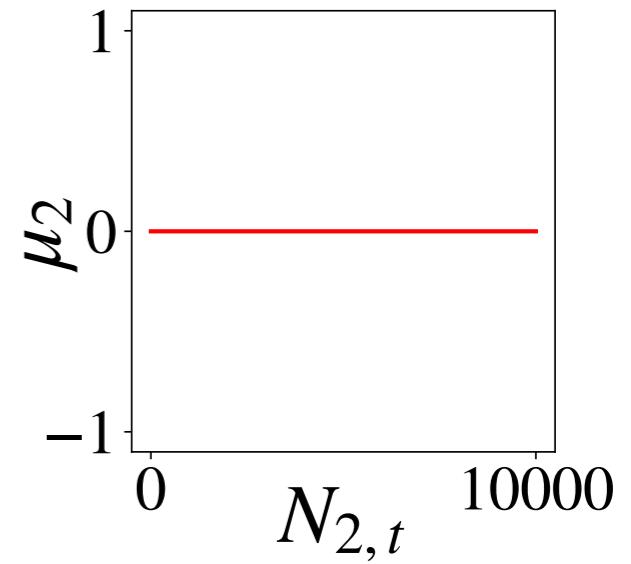
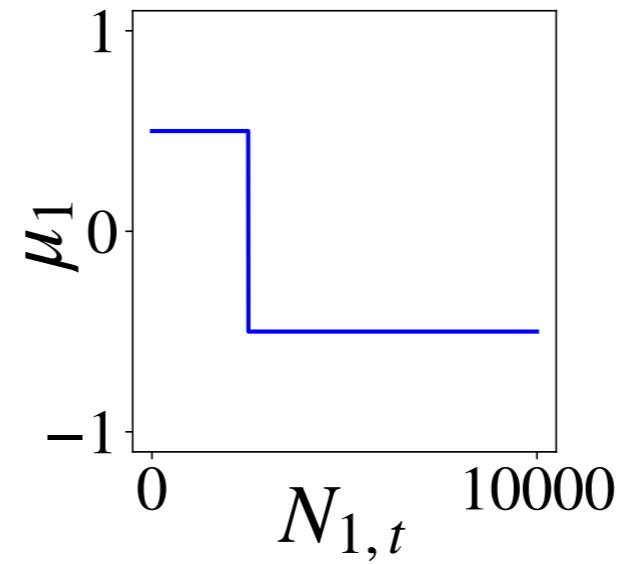


ORACLE & REGRET



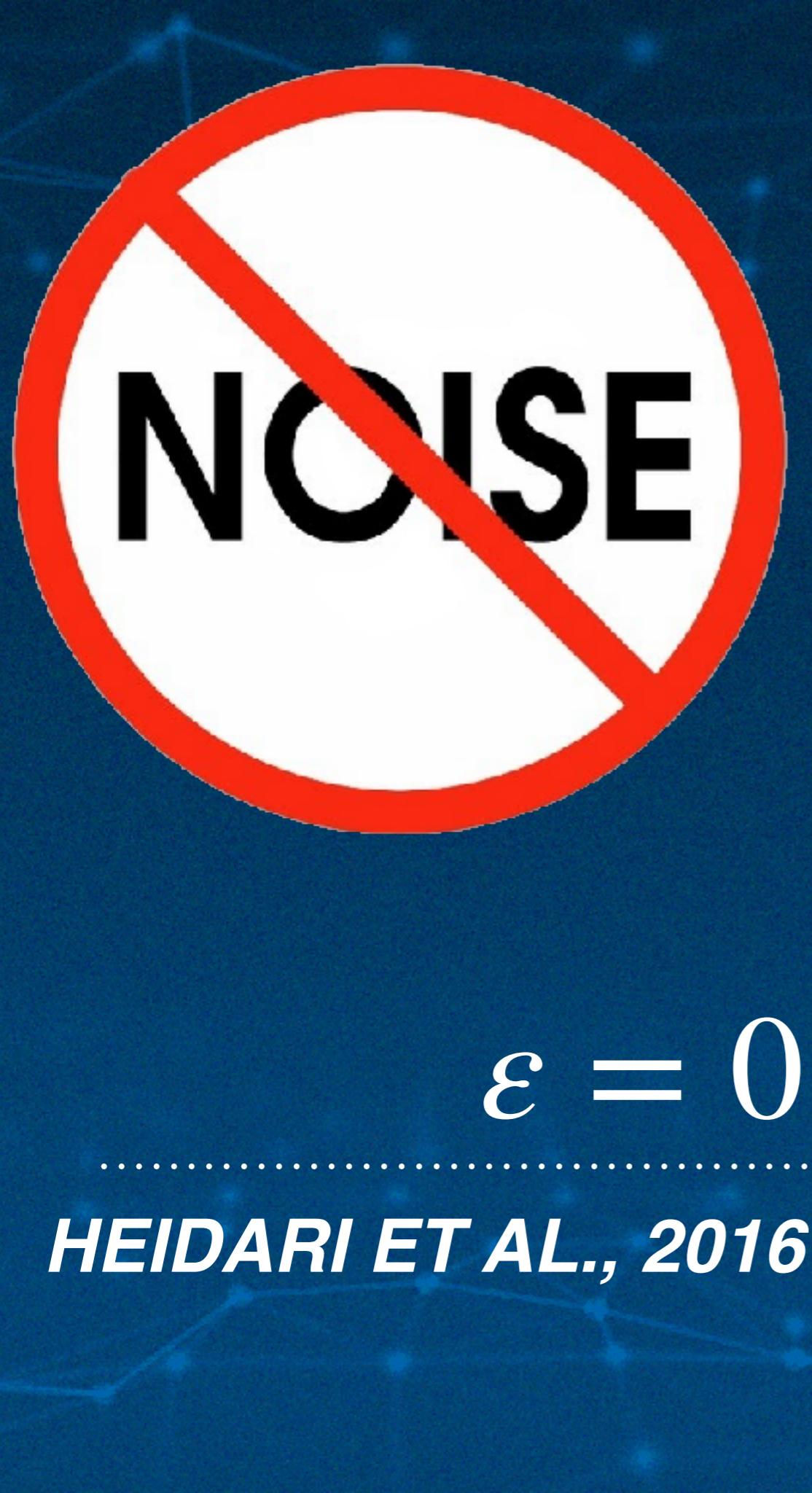
Optimal Oracle:

$$\pi^*(t) = \arg \max_i \mu_i(N_{i,t})$$



$$\begin{aligned} R_T(\pi) &= J_T(\pi^*) - J_T(\pi) \\ &= \sum_{i \in \text{UP}} \sum_{s=N_{i,T}^\pi + 1}^{N_{i,T}^*} \mu_i(s) - \sum_{i \in \text{OP}} \sum_{s=N_{i,T}^* + 1}^{N_{i,T}^\pi} \mu_i(s) \end{aligned}$$

$\left\{ \begin{array}{l} \text{UP : underpulled arms by } \pi \text{ w.r.t } \pi^* \\ \text{OP : overpulled arms by } \pi \text{ w.r.t } \pi^* \end{array} \right.$



Optimal **Oracle**:

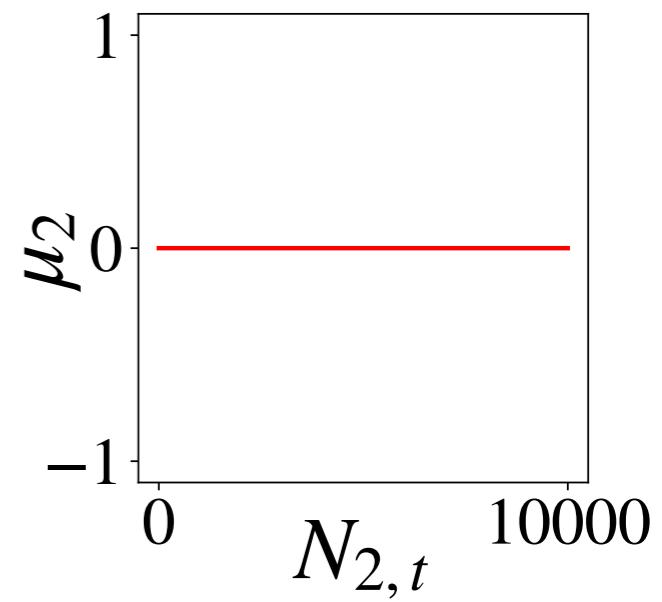
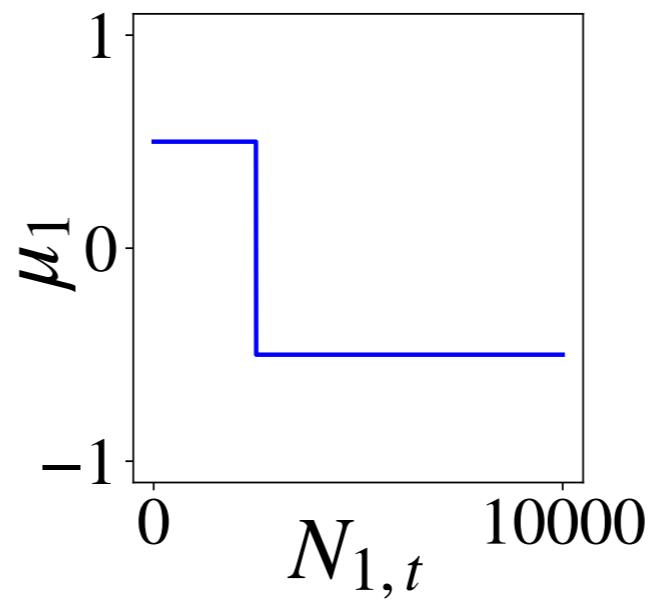
$$\pi_O(t) = \arg \max_i \mu_i(N_{i,t})$$

Online bandit policy:

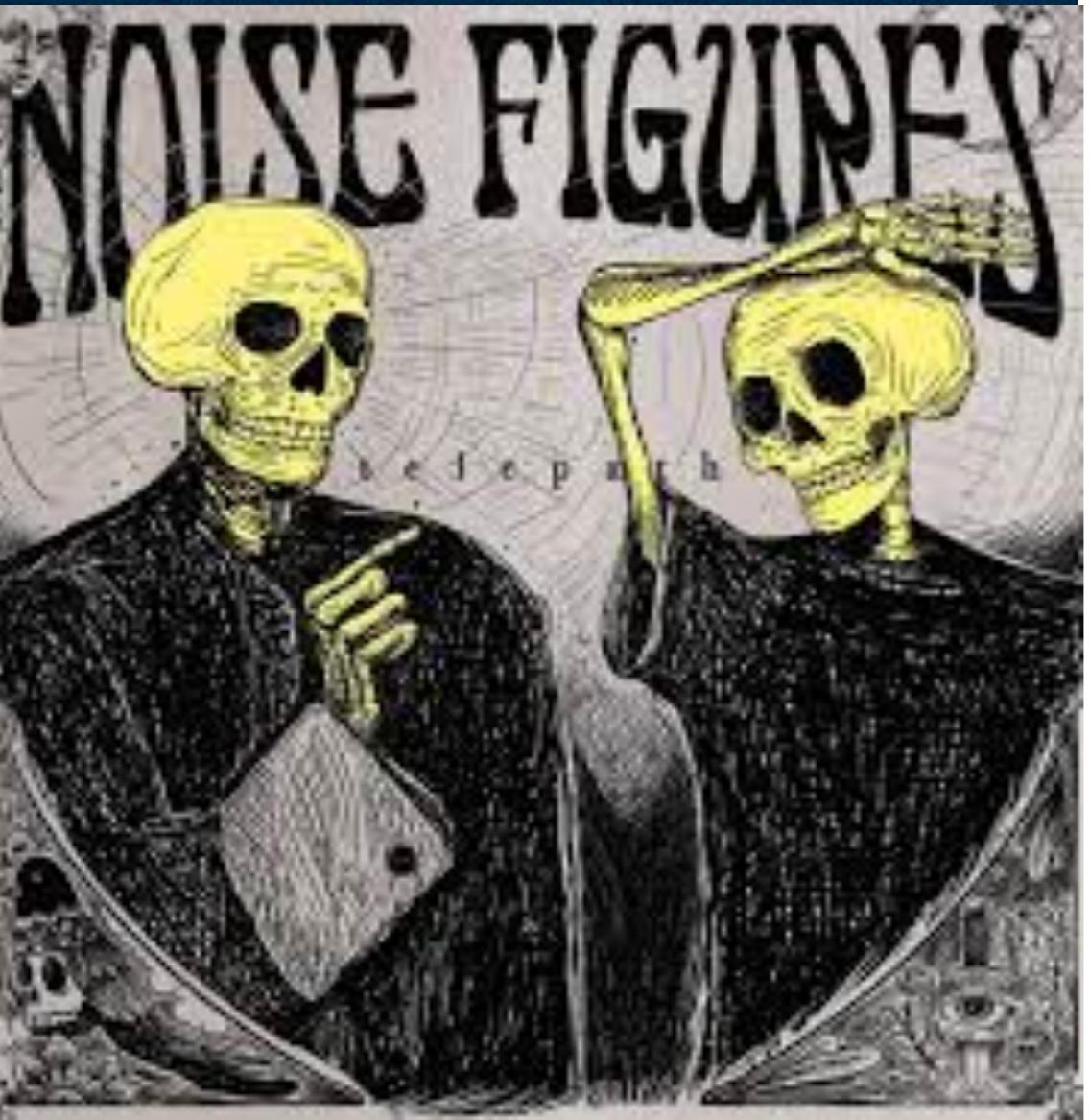
$$\pi_B(t) = \arg \max_i \mu_i(N_{i,t} - 1)$$

Minimax regret rate: $R_T(\pi_B) \leq (K - 1)L$

→ At most, 1 mistake per arm.



Rested : New mistake cancels the old one !



NOISE

LEVINE ET AL., 2017

NOISE = AVERAGE
(Which samples? How many?)

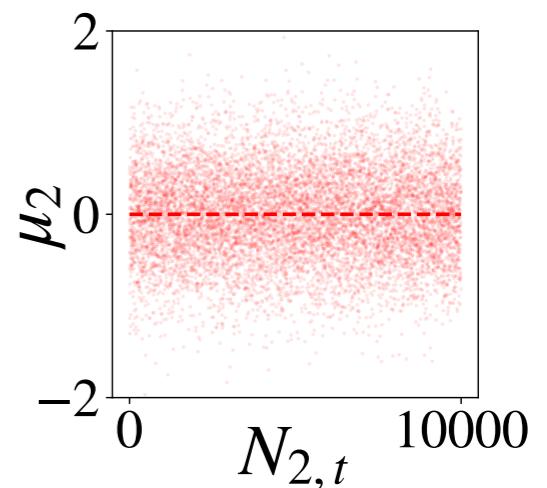
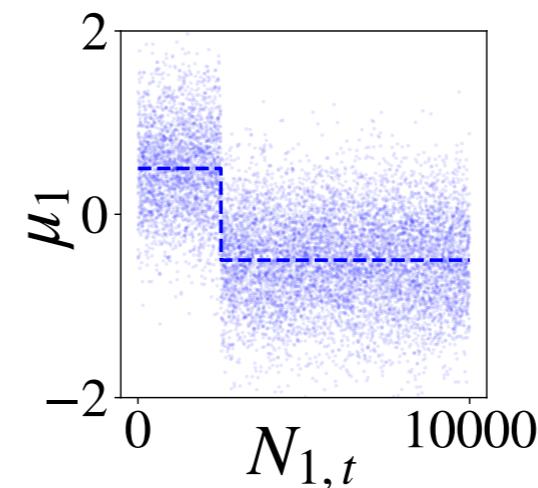
wSWA [LEVINE, 2017]

Algorithm 3 SWA (Levine et al., 2017)

Input: K, L, T, σ

- 1: $h \leftarrow \tilde{\mathcal{O}}\left(\left(\frac{\sigma T}{KL}\right)^{2/3}\right)$
- 2: **for** $t \leftarrow Kh + 1, Kh + 2, \dots$ **do**
- 3: SELECT : $\arg \max_{i \in \mathcal{K}} \hat{\mu}_i^h(N_{i,t})$
- 4: **end for**

Average of the h last samples



Regret due to bias: $\tilde{\mathcal{O}}(LKh)$

Regret due to variance : $\tilde{\mathcal{O}}\left(\sigma T \sqrt{h}^{-1}\right)$

Worst case regret : $\tilde{\mathcal{O}}(K^{1/3}T^{2/3})$

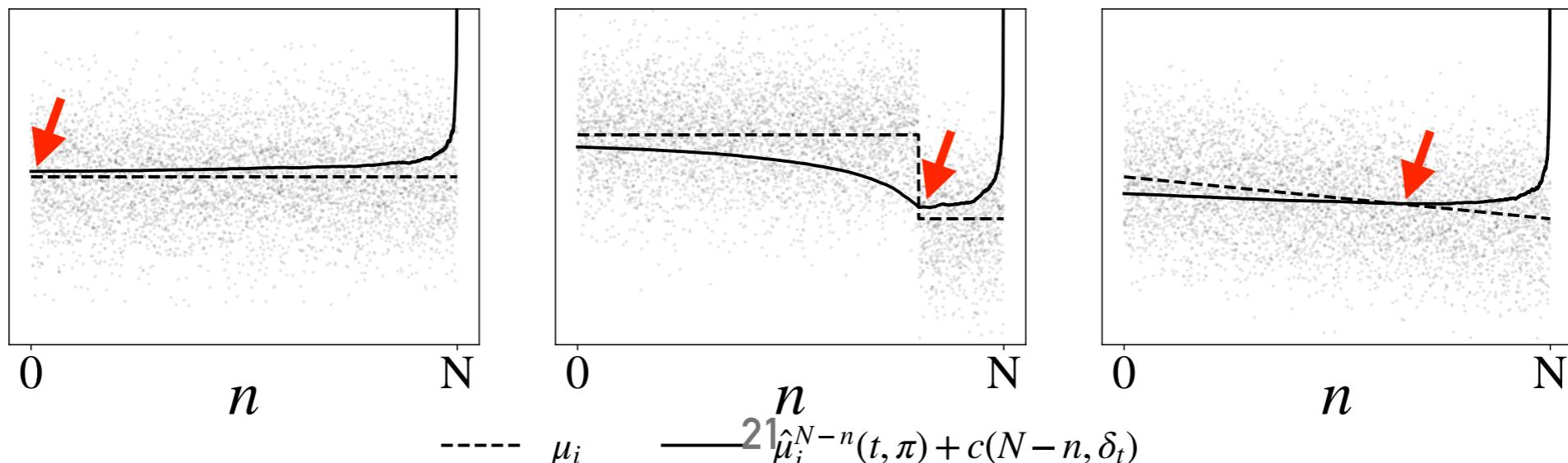
ROTTING ADAPTIVE WINDOW UCB

Algorithm 5 RAW-UCB

Input: K, σ, α

```
1: for  $t \leftarrow K + 1, K + 2, \dots$  do
2:    $\delta_t \leftarrow \frac{1}{t^\alpha}$ 
3:   SELECT :  $\arg \max_{i \in \mathcal{K}} \min_{h \leq N_{i,t}} (\hat{\mu}_i^h(N_{i,t}) + c(h, \delta_t))$ 
4: end for
```

- ▶ Rewards are decreasing: $\hat{\mu}_i^h(N_{i,t}) + c(h, \delta_t)$ is a UCB for the future value
- ▶ RAW-UCB selects the minimum (tightest) one (on h) as index of the arm



UPPER BOUNDS

Worst-case upper bound

$$\mathbb{E} [R_T(\pi_R)] \leq C\sigma\sqrt{KT \log(T)} + KL$$

Comparison w/ wSWA

$$\mathbb{E} [R_T(\pi_{\text{wSWA}})] = \tilde{O}(L^{1/3}\sigma^{2/3}K^{1/3}T^{2/3})$$

Problem-dependent upper bound

$$\mathbb{E} [R_T(\pi_R)] \leq \sum_{i \in \mathcal{K}} O\left(\frac{\log(T)}{\Delta_{i,h_{i,T}^+ - 1}}\right)$$

Comparison w/ wSWA

Pure worst-case strategy

$\Delta_{i,h}$ Difference between the average of the h first overpulls of arm i and the worst reward pulled by the optimal policy

$h_{i,T}^+$ High-probability upper bound on the number of overpulls for RAW-UCB

UPPER BOUNDS

Worst-case upper bound

$$\mathbb{E} \left[R_T(\pi_R) \right] \leq C\sigma\sqrt{\textcolor{red}{K}T \log(T)} + K\textcolor{red}{L}$$

Comparison w/ wSWA

$$\mathbb{E} \left[R_T(\pi_{\text{wSWA}}) \right] = \tilde{O} \left(\textcolor{red}{L}^{1/3} \sigma^{2/3} K^{1/3} T^{2/3} \right)$$

Problem-dependent upper bound

$$\mathbb{E} \left[R_T(\pi_R) \right] \leq \sum_{i \in \mathcal{K}} O \left(\frac{\log(T)}{\Delta_{i,h_{i,T}^+ - 1}} \right)$$

$\Delta_{i,h} = \Delta_i$ on a stationary bandits problem

$\Delta_{i,h_{i,T}^+ - 1}$ is a problem-dependent quantity

Comparison w/ wSWA

Pure worst-case strategy

PROOF SKETCH

1. Decompose the regret

$$\begin{aligned} R_T(\pi) &= \sum_{i \in \text{UP}} \sum_{s=N_{i,T}^\pi+1}^{N_{i,T}^\star} \mu_i(s) - \sum_{i \in \text{OP}} \sum_{s=N_{i,T}^\star+1}^{N_{i,T}^\pi} \mu_i(s) \\ &\leq \sum_{i \in \text{OP}} \sum_{s=N_{i,T}^\star+1}^{N_{i,T}^\pi} (\mu_T^+ - \mu_i(s)) \quad \mu_T^+ : \text{largest underpulled values} \\ &\sim \text{"number of mistakes"} \times \text{"average gap"} \end{aligned}$$

2. Use the adaptive window mechanism to get

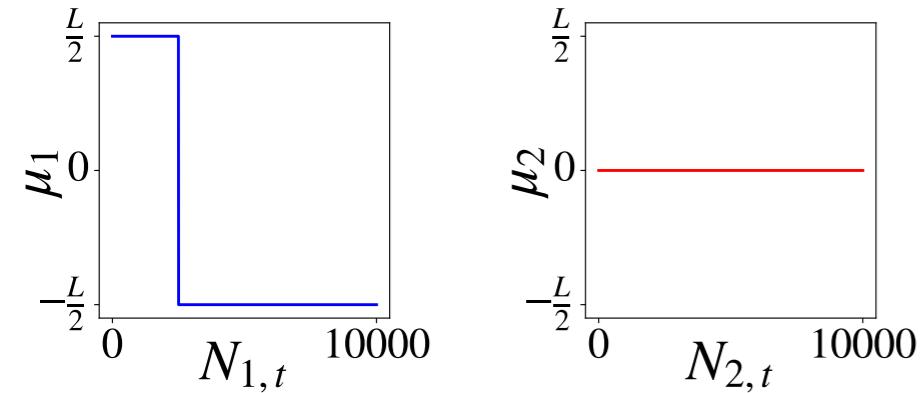
$$\text{"number of mistakes"} \sim \frac{\log T}{\text{"average gap"}^2}$$

3. Work on the gap to be problem-dependent

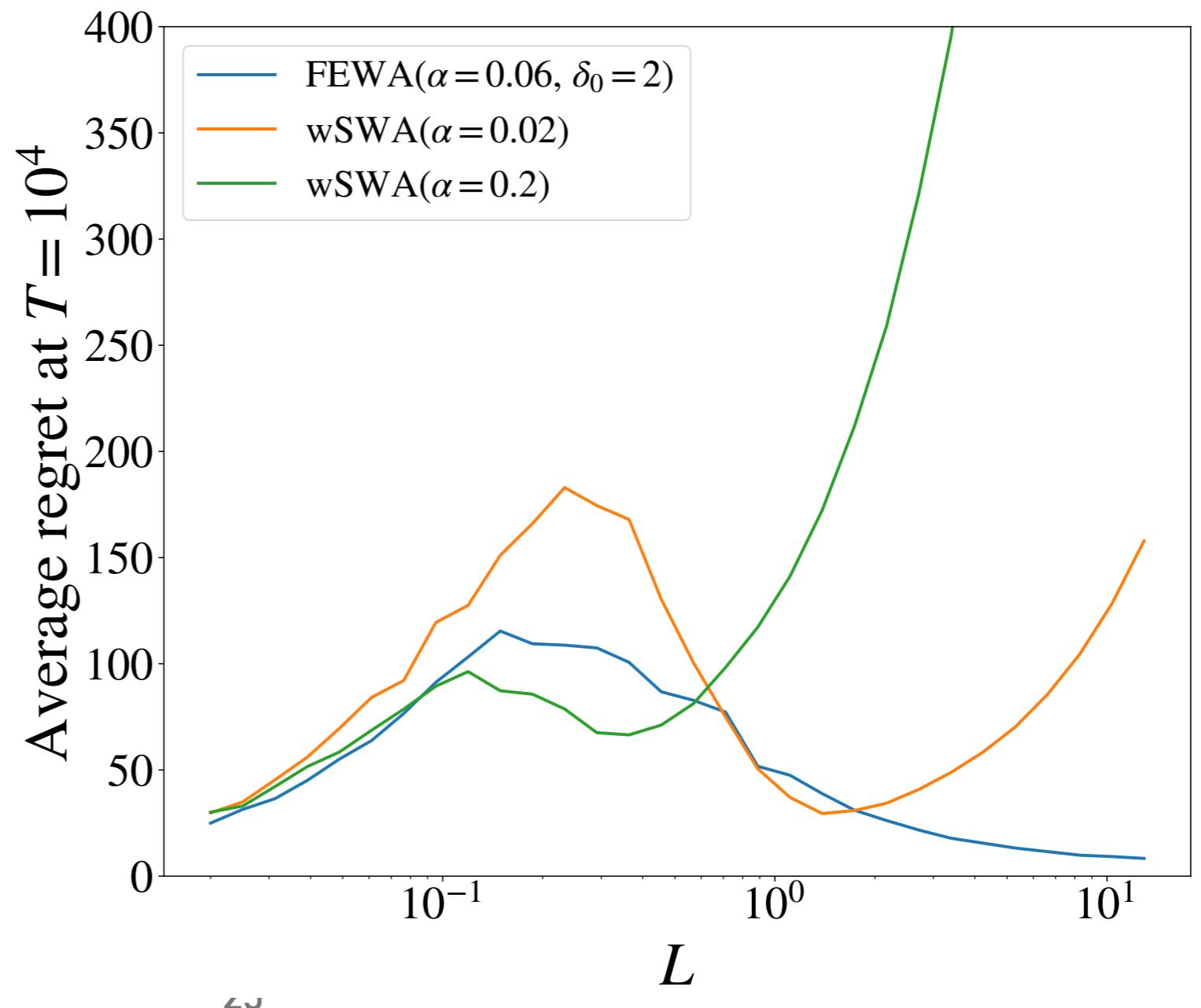
"number of mistakes" \rightarrow "maximal number of mistakes with high probability"

$$\mu_T^+ \rightarrow \min_i \min_{s \leq N_{i,T}^\star} \mu_i(s)$$

EXPERIMENTAL RESULTS



$\sigma = 1 ; L$



RESTED CONTRIBUTIONS

Algorithm	Sliding Window Average (Levine et al., 2017)	FEWA & RAW-UCB
Average(s)	Fixed single window	Multiple windows
Problem Independent Bound	$\tilde{\mathcal{O}}(L^{1/3}\sigma^{2/3}K^{1/3}T^{2/3})$	$\tilde{\mathcal{O}}(\sigma\sqrt{KT}) + \mathcal{O}(KL)$
Problem-Dependent Bound		$\mathcal{O}\left(\sum_{i \in \mathcal{K}} \frac{\log T}{\Delta_i^h}\right)$
Knowledge	$T \ L \ \sigma$	σ
Complexity per round	$\mathcal{O}(K^{1/3}T^{2/3})$	efficient algo $\mathcal{O}(K \log T)$

OUTLINE

1. Bandits for Adaptive Educational Systems
2. Rested Rotting Bandits
3. Restless Rotting Bandits

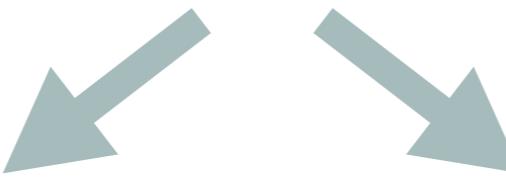
NON-STATIONARY RESTLESS BANDITS

- ▶ $\{\mu_i\}$ are functions of round t - restless
- ▶ Minimize cumulative regret w.r.t. the optimal strategy: $\sum_{t \leq T} \mu_{i_t^*}(t) - \mu_{i_t}(t)$
- ▶ Unlearnable if μ_i can change significantly at every round.
- ▶ Common settings
 - ◆ μ_i is piece-wise stationary with Υ_T pieces (*Garivier & Moulaines, 2011*)
 - ◆ μ_i has a permitted amount of change V_T (*Besbes et al., 2014*)
$$\sum_{t \leq T} \max_{i \in \mathcal{K}} |\mu_i(t) + 1 - \mu_i(t)| \leq V_T$$

Piece-wise stationary rate:

$$\Theta(\sqrt{K\Upsilon_T T})$$

(*Garivier & Moulaines, 2011*)



Variational budget rate:

$$\Theta(K^{1/3}V_T^{1/3}T^{2/3})$$

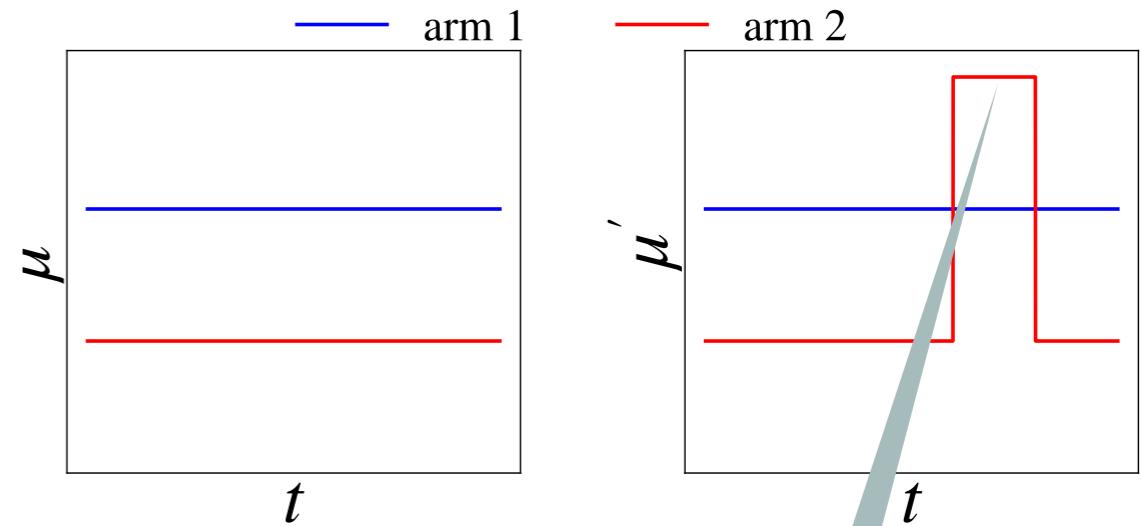
(*Besbes et al. , 2014*)

PROBLEM DEPENDENT GUARANTEES

EXP3.S (Auer et al. 2002b), an adversarial algorithm, matches the two minimax rates.
Can we get a problem-dependent bound?

Theorem (Garivier et al. 2011): Let π a policy suffering $R_T(\mu)$ on a 2-arms stationary bandits problem μ . Then, for T large enough, there exists a piece-wise stationary problem μ' such that π suffers:

$$R_T(\mu') \geq \frac{T}{22R_T(\mu)}$$



Quick increase

$$\tau = \frac{T}{R_T(\mu)/\Delta}$$

Corollary : No !

$\mathcal{O}(\sqrt{T})$ anywhere $\Leftrightarrow \mathcal{O}(\sqrt{T})$ everywhere



.....

RESTLESS ROTTING

Do we need to
spend $\mathcal{O}(\sqrt{T})$
exploration budget
when rewards
decay ?

RAW-UCB ON ROTTING RESTLESS BANDITS

RAW-UCB without knowing T , V_T nor Υ_T

Rotting Variational budget

$$\mathbb{E} [R_T(\pi_R)] \leq \tilde{\mathcal{O}} (\sigma^{2/3} K^{1/3} V_T^{1/3} T^{2/3})$$

minimax rate

Rotting Piecewise stationary

$$\mathbb{E} [R_T(\pi_R)] \leq \tilde{\mathcal{O}} (\sigma \sqrt{K \Upsilon_T T})$$

minimax rate

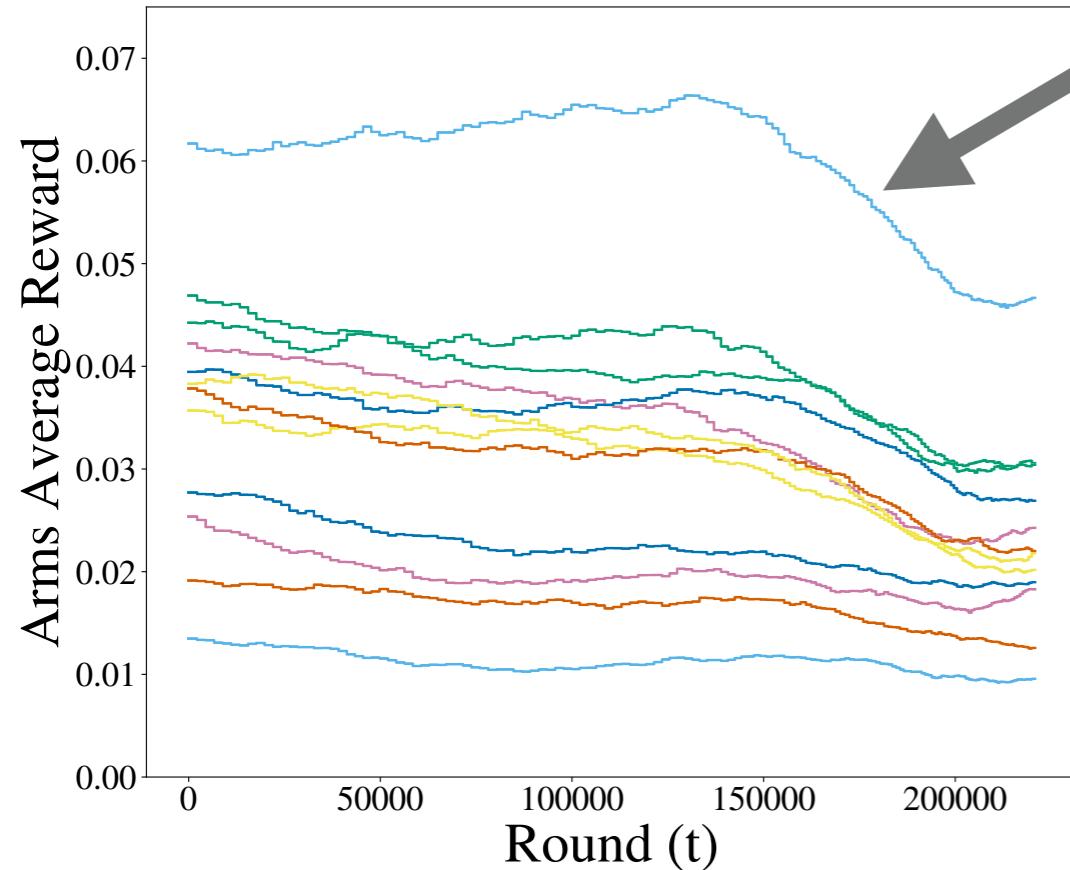
$$\mathbb{E} [R_T(\pi_R)] \leq \sum_{i \in \mathcal{K}} \sum_{k=1}^{\Upsilon_T} \frac{32\sigma^2 \log T}{\Delta_{i,k}} + \mathcal{O} (\sqrt{\log T})$$

problem-dependent bound!

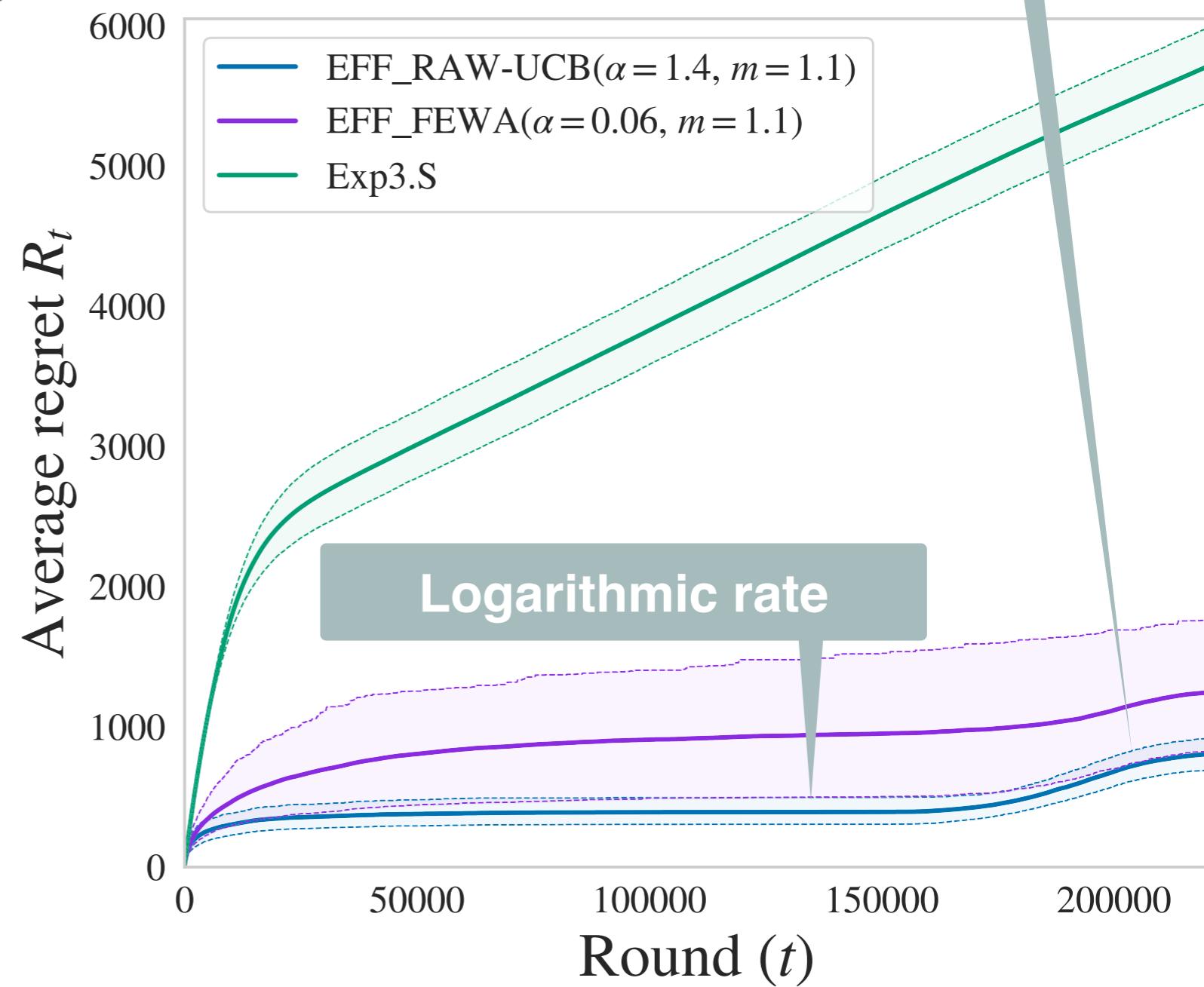
YAHOO! BENCHMARK

Significant decay :
exploration

Day 7 - $K = 12$



Decaying trend in the probability of click on an article from noon to 6am.



RESTLESS CONTRIBUTIONS

1

Rotting property makes restless bandits easier

- ✓ $\mathcal{O}(\log T)$ problem-dependent bound

2

UCB-like index is sufficient for rotting bandits

- ✓ No random exploration, no passive forgetting, no change-detection routine
- ✓ Near-optimal gap-dependent and minimax bounds
- ✓ Agnostic to Υ_T, V_T, T

3

RAW-UCB solves rested OR restless rotting bandits

- ✓ with the same tuning
- ✗ rested AND restless rotting bandits are **incompatible**

CONTRIBUTIONS TO EDUCATIONAL SYSTEMS

1. What is the objective ?

Target the most difficult topic

2. What is the « best » action ?

✓ Optimal policy selects different actions when the student progresses

3. How does action impact observation?

✓ Model the learning of the student

4. Can we learn something relevant with ~100 samples?

✓ Not harder than stationary bandits

CONTRIBUTIONS TO EDUCATIONAL SYSTEMS

1. What is the objective ?

Target the most difficult topic

✗ Target the easiest topic until it is mastered

2. What is the « best » action ?

✓ Optimal policy selects different actions when the student progresses

3. How does action impact observation?

✓ Model the learning of the student

4. Can we learn something relevant with ~100 samples?

✓ Not harder than stationary bandits

✗ RAW-UCB = Round Robin

THANK YOU !

Rotting Bandits are not Harder than Stochastic ones,
Julien Seznec, Andrea Locatelli, Alexandra Carpentier, Alessandro Lazaric, Michal Valko,
Artificial Intelligence and Statistics (2019), (with oral presentation)

A Single Algorithm for both Rested and Restless Rotting Bandits,
Julien Seznec, Pierre Ménard, Alessandro Lazaric, Michal Valko,
Artificial Intelligence and Statistics (2020)