

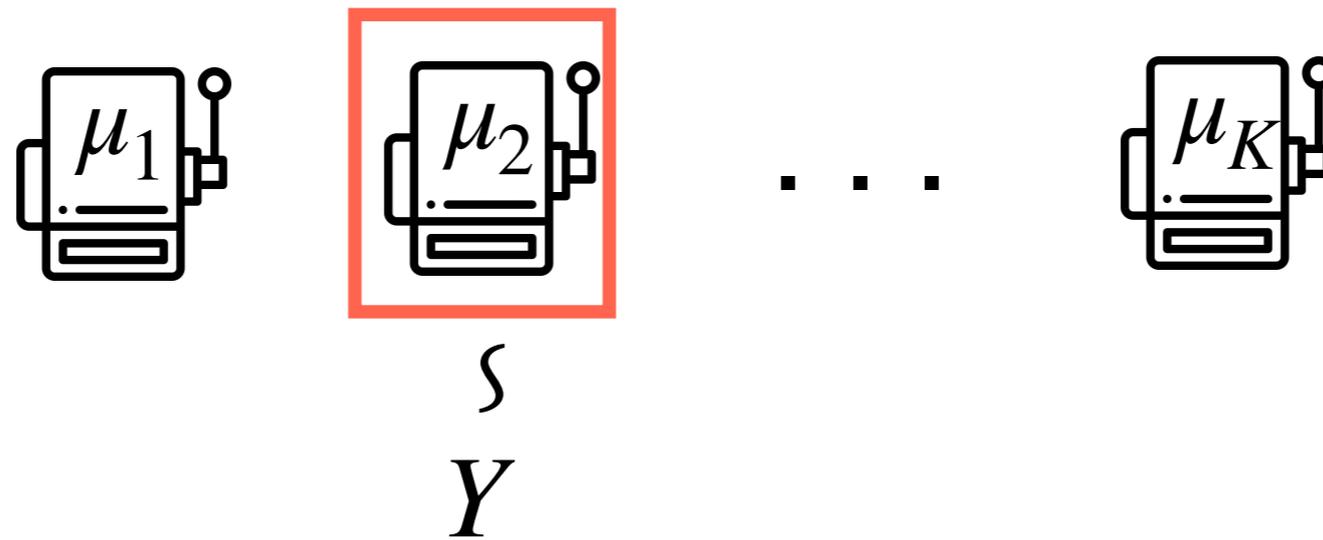
On conditional versus marginal bias in multi-armed bandits

Jaehyeok Shin¹, Aaditya Ramdas^{1,2} and Alessandro Rinaldo¹

Dept. of Statistics and Data Science¹,
Machine Learning Dept.², CMU



Stochastic Multi-armed bandits (MABs)



"Random reward"

Adaptive sampling scheme to maximize rewards / to identify the best arm

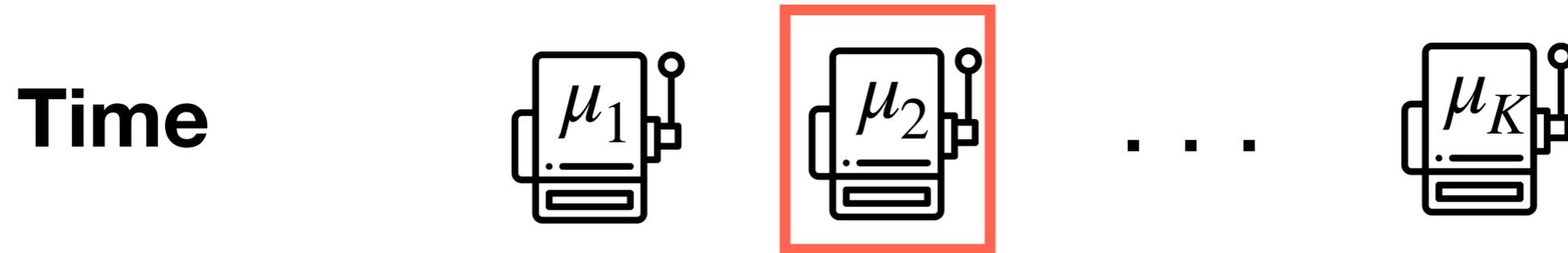


Adaptive sampling scheme to maximize rewards / to identify the best arm



$t = 1$

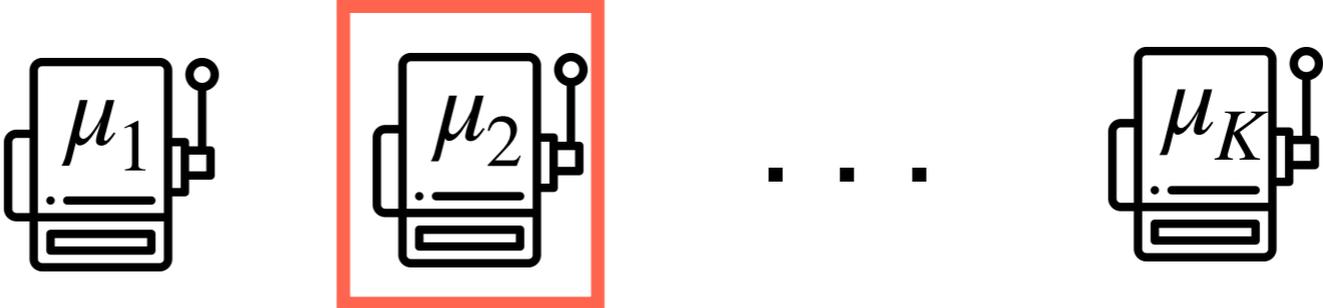
Adaptive sampling scheme to maximize rewards / to identify the best arm



$t = 1$

Adaptive sampling scheme to maximize rewards / to identify the best arm

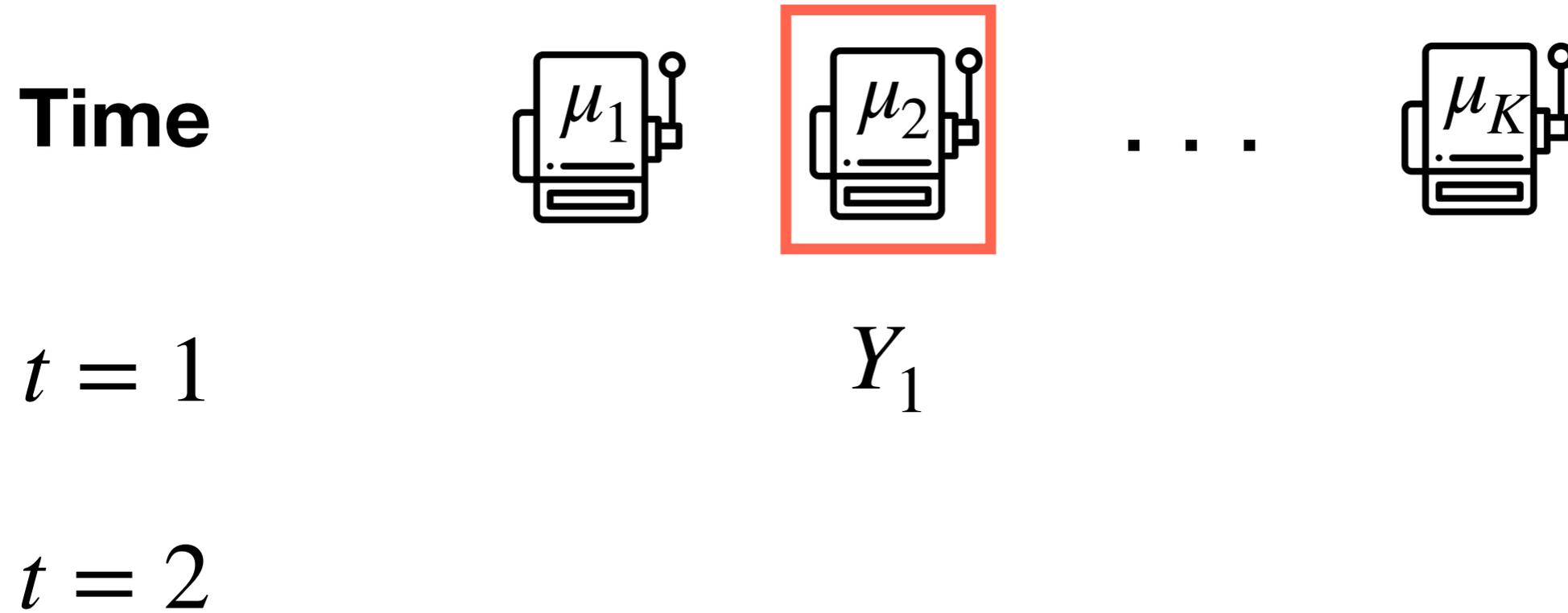
Time



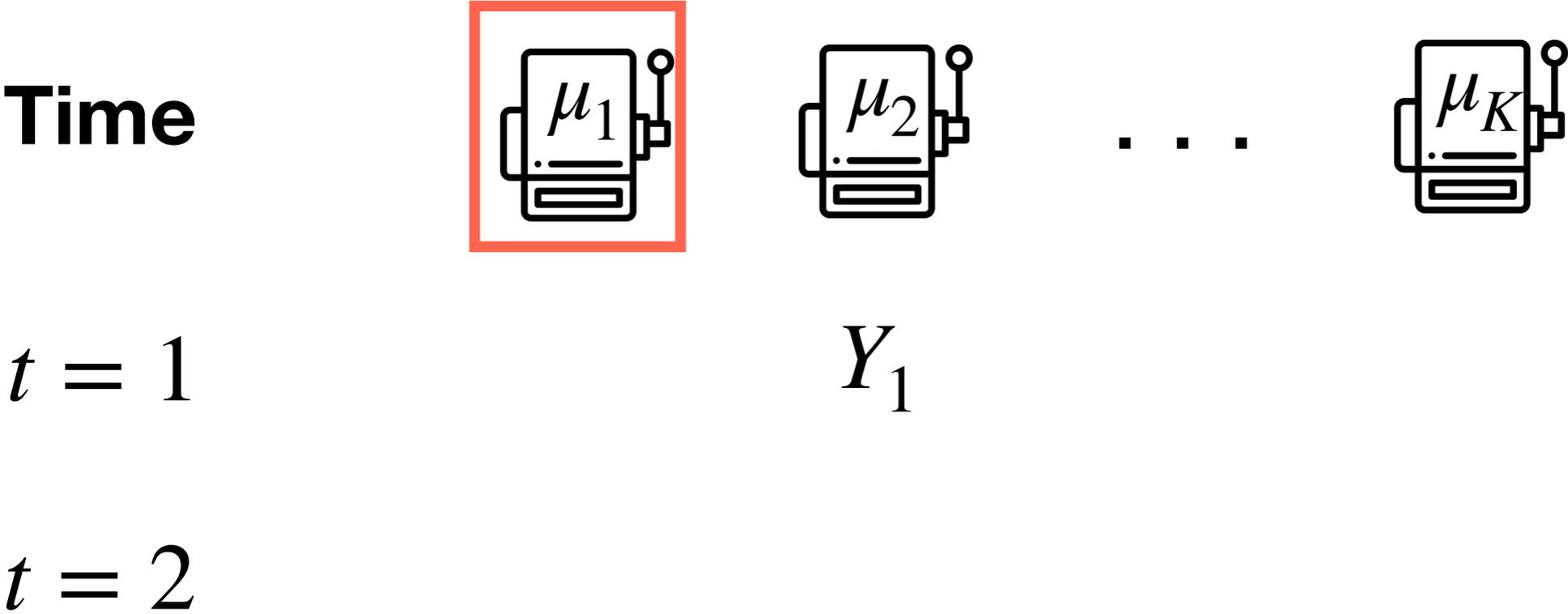
$t = 1$

Y_1

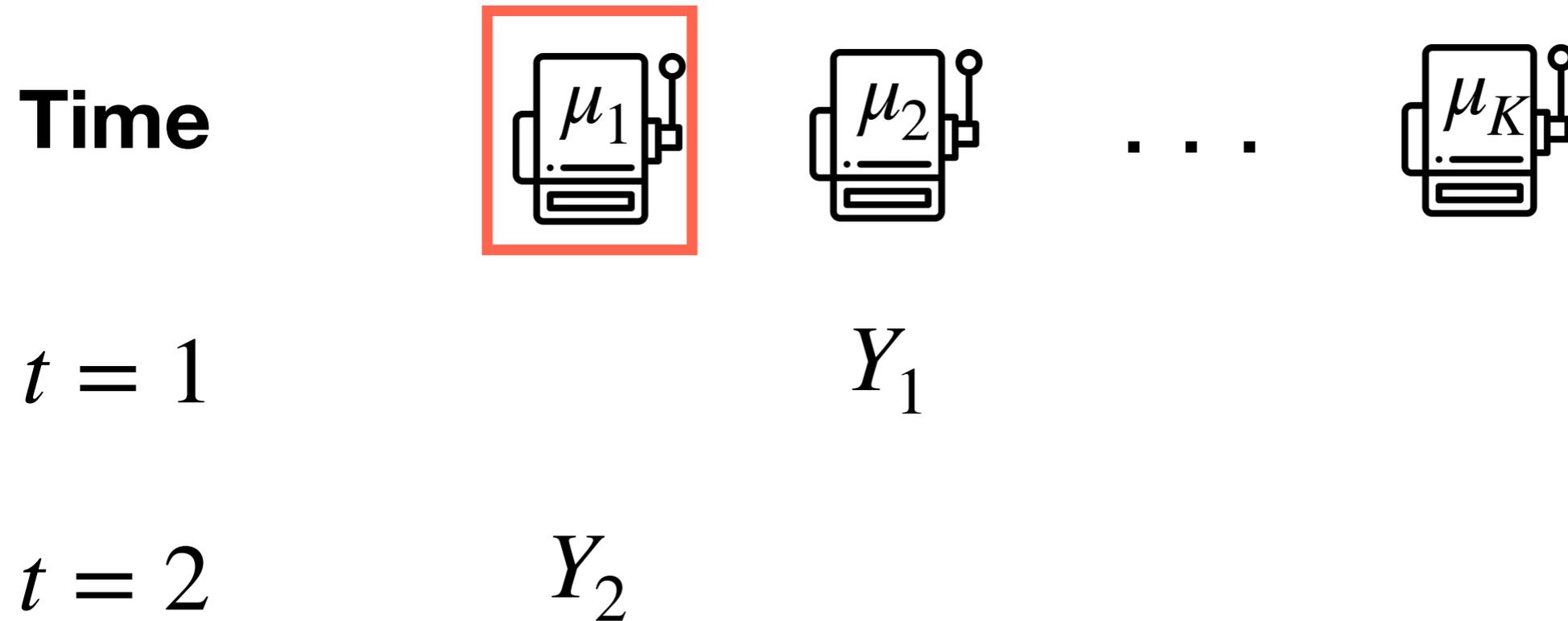
Adaptive sampling scheme to maximize rewards / to identify the best arm



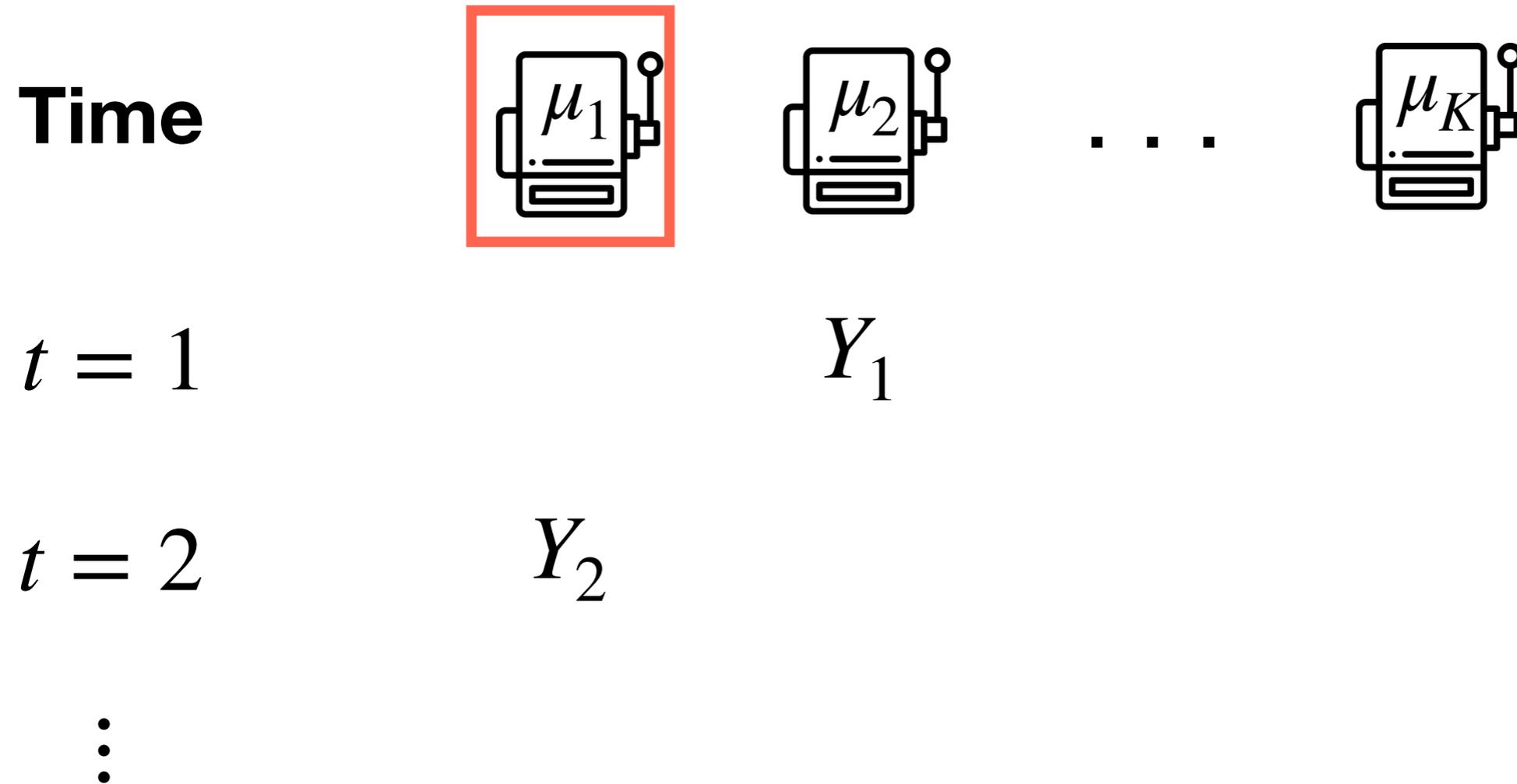
Adaptive sampling scheme to maximize rewards / to identify the best arm



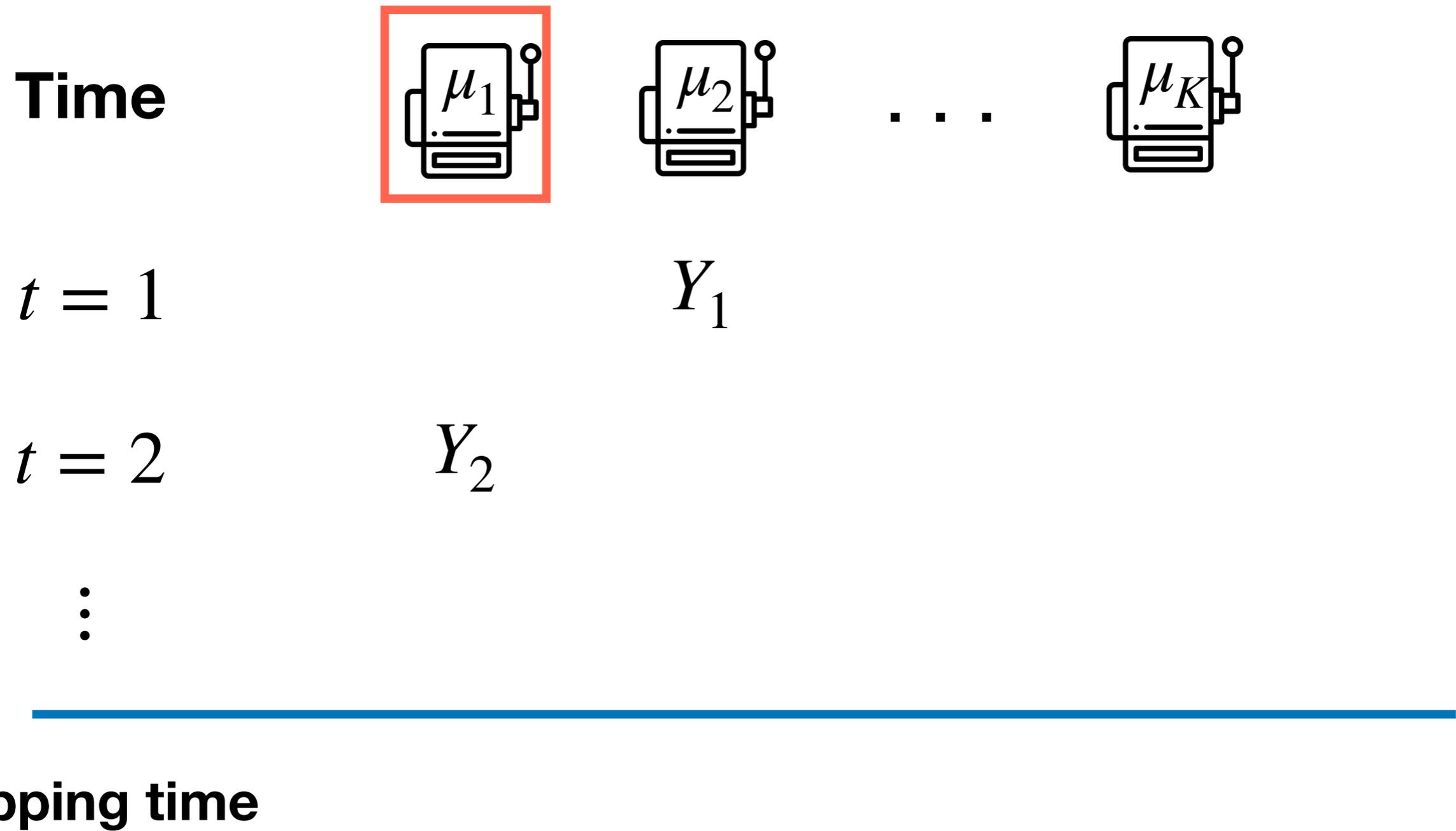
Adaptive sampling scheme to maximize rewards / to identify the best arm



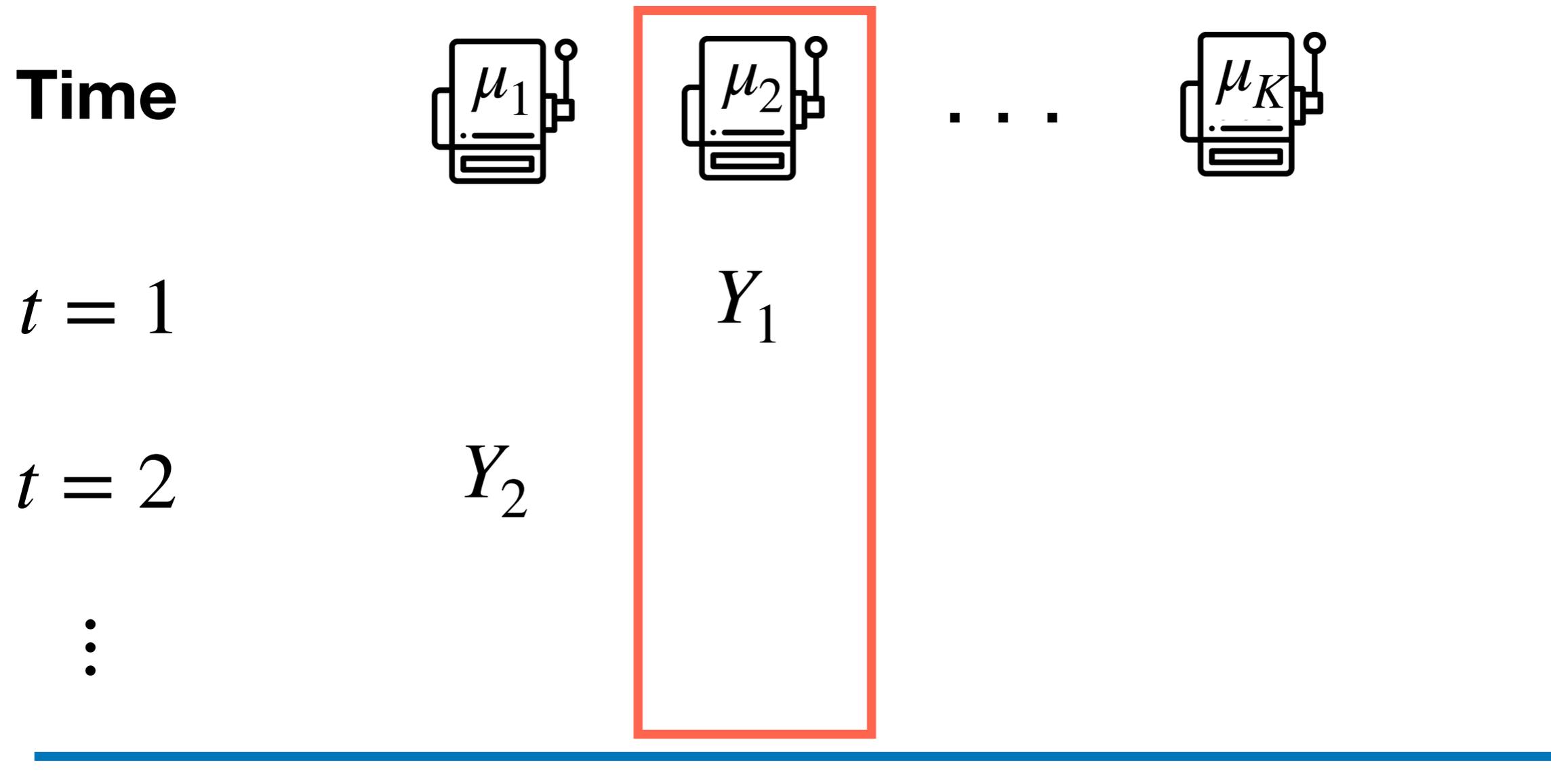
Adaptive sampling scheme to maximize rewards / to identify the best arm



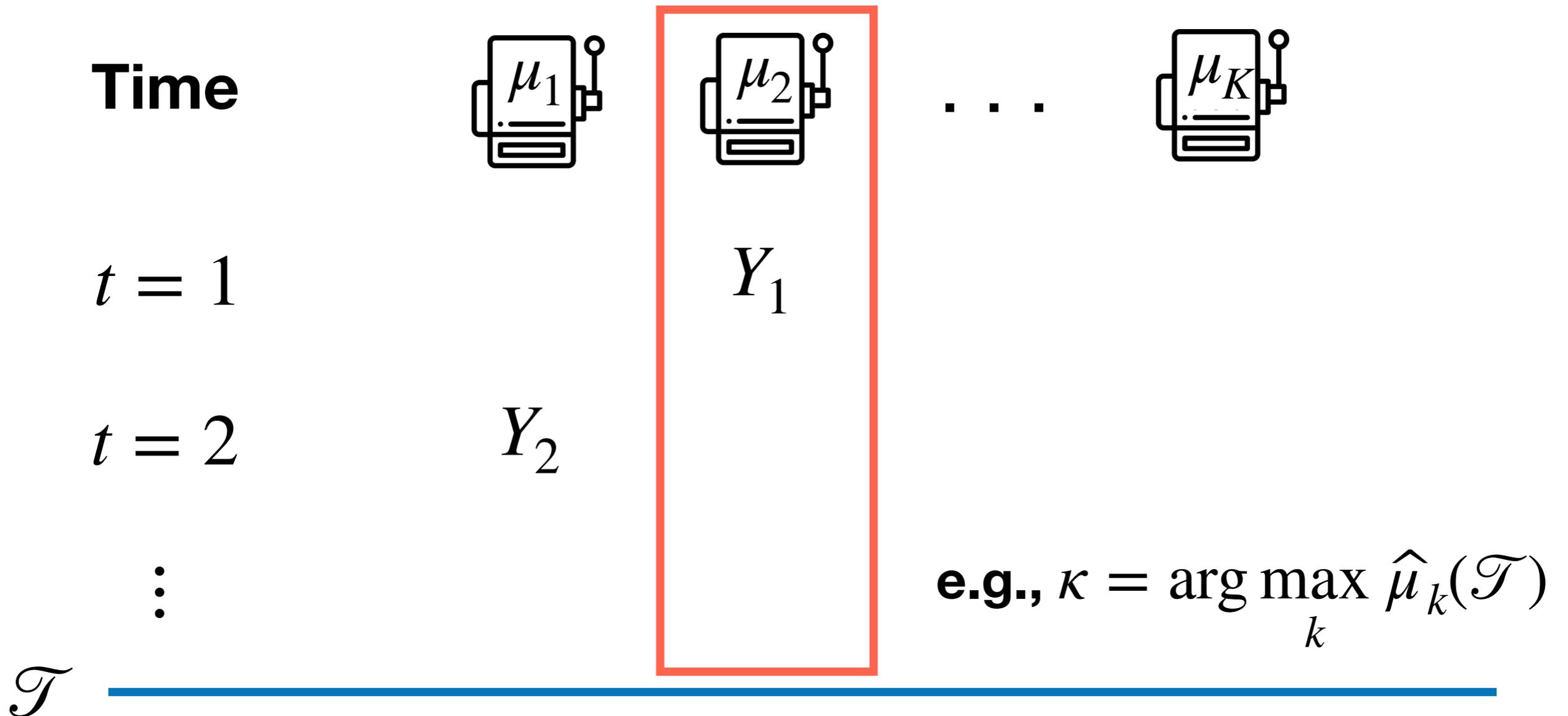
Adaptive sampling scheme to maximize rewards / to identify the best arm



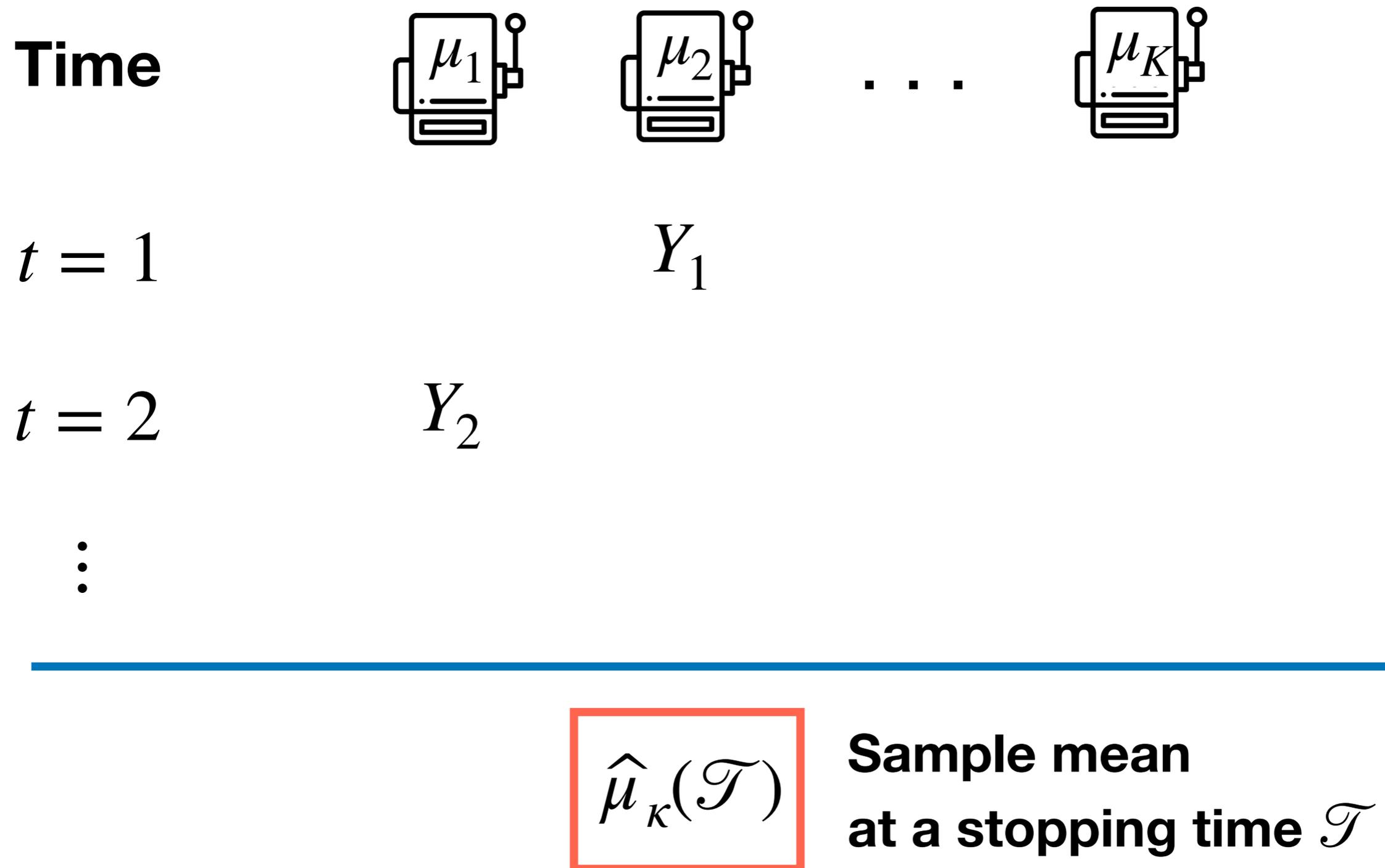
Collected data can be used to identify an interesting arm...



Collected data can be used to identify an interesting arm...



...and the data can be used to
conduct statistical inferences.



Q. Sign of the bias of sample mean?

$$\mathbb{E} \left[\hat{\mu}_\kappa(\mathcal{T}) - \mu_\kappa \right] \leq \text{or} \geq 0?$$

Xu et al. [2013] :

An informal argument why the sample mean is *negatively* biased for “optimistic” algorithms.

Villar et al. [2015] :

Demonstrate this *negative* bias in a simulation study motivated by using MAB for clinical trials.

Q. Sign of the bias of sample mean?

$$\mathbb{E} \left[\hat{\mu}_\kappa(\mathcal{T}) - \mu_\kappa \right] \leq \text{or} \geq 0?$$

Xu et al. [2013] :

An informal argument why the sample mean is *negatively* biased for “optimistic” algorithms.

Villar et al. [2015] :

Demonstrate this *negative* bias in a simulation study motivated by using MAB for clinical trials.

Nie et al. [2018]

Sample mean is **negatively** biased

$$\mathbb{E} \left[\hat{\mu}_k(t) - \mu_k \right] \leq 0$$

Fixed Arm

Fixed Time

for MABs designed to maximize cumulative reward.

Shin et al. [2019]

Introduced "monotonicity property" characterizing the bias of the sample mean for more general classes of MABs.

$$\mathbb{E} \left[\hat{\mu}_k(\mathcal{T}) - \mu_k \right]$$

Chosen Arm

Stopping Time

Nie et al. [2018]

Sample mean is **negatively** biased

$$\mathbb{E} \left[\hat{\mu}_k(t) - \mu_k \right] \leq 0$$

Fixed Arm Fixed Time

for MABs designed to maximize cumulative reward.

Shin et al. [2019]

Introduced "monotonicity property" characterizing the bias of the sample mean for more general classes of MABs.

$$\mathbb{E} \left[\hat{\mu}_k(\mathcal{T}) - \mu_k \right]$$

Chosen Arm Stopping Time

However, current understanding of bias is limited in two aspects.

1. Existing results concern the **bias of the sample mean** only.

However, current understanding of bias is limited in two aspects.

1. Existing results concern the **bias of the sample mean** only.

➔ We study the **bias of monotone functions** of the rewards.

However, current understanding of bias is limited in two aspects.

1. Existing results concern the **bias of the sample mean** only.

➔ We study the **bias of monotone functions** of the rewards.

2. Existing guarantees cover only the **marginal bias**.

However, current understanding of bias is limited in two aspects.

1. Existing results concern the **bias of the sample mean** only.

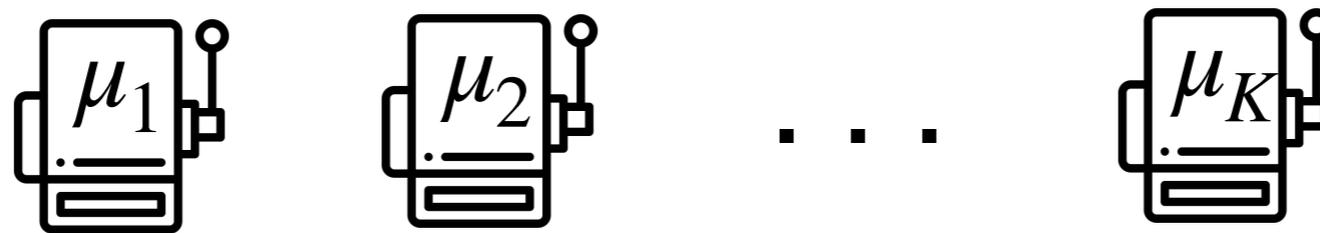
➔ We study the **bias of monotone functions** of the rewards.

2. Existing guarantees cover only the **marginal bias**.

➔ We extend previous results to cover the **conditional bias**.

Marginal vs Conditional bias

- **K prototype items**

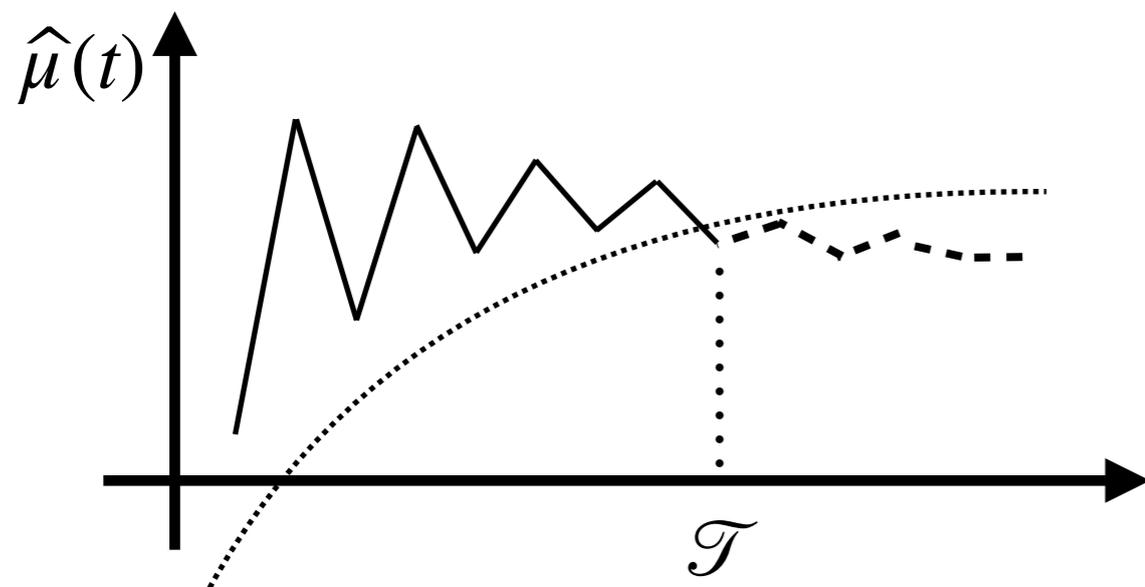


➔ Want to screen out some items by testing

$$H_{0k} : \mu_k \geq c \quad \text{vs} \quad H_{1k} : \mu_k < c \quad \text{for } k = 1, \dots, K.$$

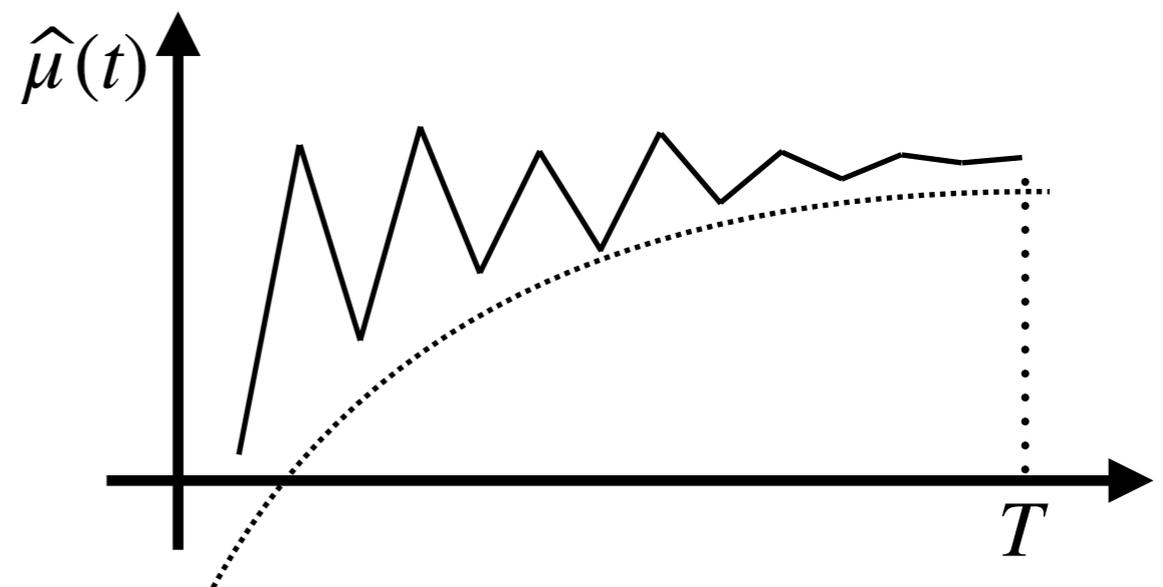
Marginal vs Conditional bias

$$H_0 : \mu \geq c \text{ vs } H_1 : \mu < c$$



"Screen out the item at \mathcal{T} "

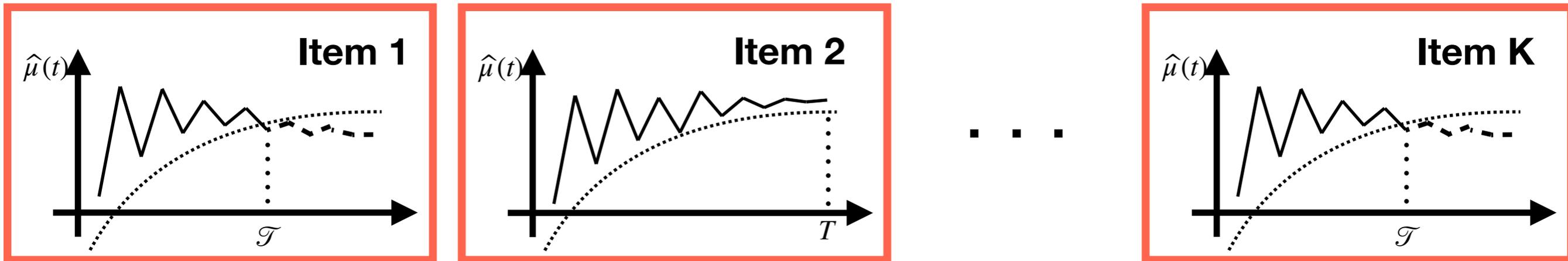
(Reject the null)



"Keep the item"

(Fail to reject the null)

Marginally,
the sample mean is **negatively** biased.

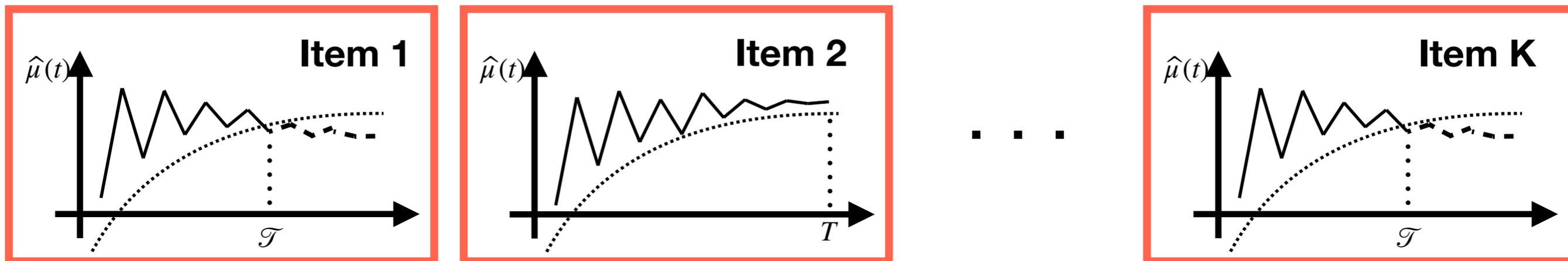


$$\mathbb{E} [\hat{\mu}_k - \mu_k] \leq 0, \quad k = 1, \dots, K$$

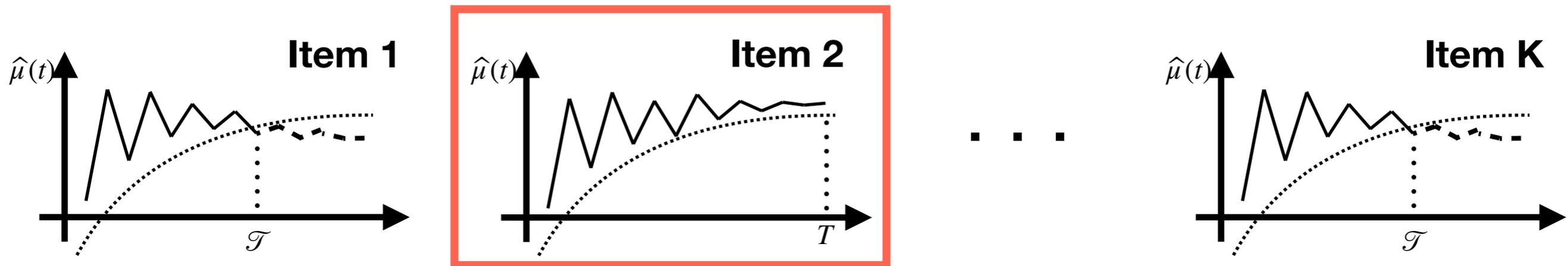
➔ **"Underestimating** the true mean revenue."

(e.g. Starr & Woodroffe [1968], Shin et al. [2019])

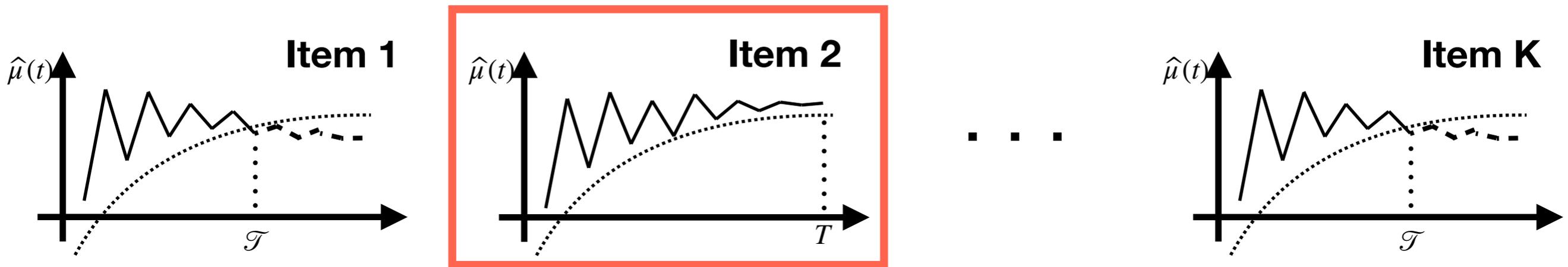
...however, we usually do not evaluate the sample mean every time.



...however, we usually do not evaluate the sample mean every time.



Conditioned on the "active" event, the sample mean is **positively** biased.



$$\mathbb{E} \left[\hat{\mu}_k - \mu_k \mid \text{item } k \text{ is active} \right] \geq 0, \quad k = 1, \dots, K$$

➡ "Overestimating the true mean revenue."

Conditional bias of the empirical cumulative distribution function (CDF)

For a fixed $y \in \mathbb{R}$,

$$\mathbb{E} \left[\widehat{F}_{k, \mathcal{T}}(y) - F_k(y) \mid C \right] \leq \mathbf{or} \geq 0?$$

e.g., $C = \{ \text{Reject the null} \}, \{ \text{Chosen as the best arm} \} \dots$

where

$\widehat{F}_{k, \mathcal{T}}$: **Empirical CDF of arm k at time \mathcal{T}**

F_k : **True CDF of arm k at time \mathcal{T}**

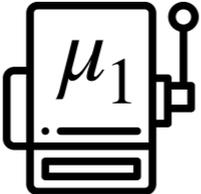
Tabular model of MABs

$$X_{\infty}^* \in \mathbb{R}^{N \times K} \quad \begin{cases} \begin{array}{cccc} \begin{array}{c} \mu_1 \\ \text{---} \\ \text{---} \\ \text{---} \end{array} & \begin{array}{c} \mu_2 \\ \text{---} \\ \text{---} \\ \text{---} \end{array} & \dots & \begin{array}{c} \mu_K \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \\ \wr \text{ i.i.d.} & \wr \text{ i.i.d.} & & \wr \text{ i.i.d.} \\ X_{1,1}^* & X_{1,2}^* & \dots & X_{1,K}^* \\ \\ X_{2,1}^* & X_{2,2}^* & \dots & X_{2,K}^* \\ \\ \vdots & \vdots & \vdots & \vdots \end{array} \end{cases} \quad \begin{array}{c} \vdots \\ \vdots \\ \vdots \\ \vdots \end{array}$$

"Hypothetical table"

Tabular model of MABs

Time

		...	
$X_{1,1}^*$	$X_{1,2}^*$...	$X_{1,K}^*$
$X_{2,1}^*$	$X_{2,2}^*$...	$X_{2,K}^*$
⋮	⋮	⋮	⋮

Tabular model of MABs

Time



$t = 1$

$$X_{1,1}^* \quad X_{1,2}^* \quad \dots \quad X_{1,K}^*$$

$$X_{2,1}^* \quad X_{2,2}^* \quad \dots \quad X_{2,K}^*$$

$\vdots \quad \vdots \quad \vdots \quad \vdots$

Tabular model of MABs

Time



...



$t = 1$

$X_{1,1}^*$

Y_1

...

$X_{1,K}^*$

$X_{2,1}^*$

$X_{2,2}^*$

...

$X_{2,K}^*$

⋮

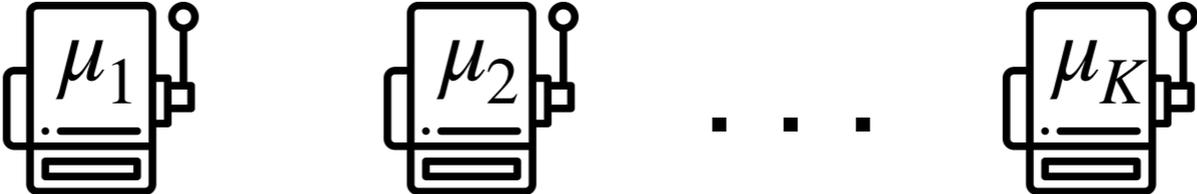
⋮

⋮

⋮

Tabular model of MABs

Time



$t = 1$



$t = 2$



Tabular model of MABs

Time			...	
$t = 1$	$X_{1,1}^*$	Y_1	...	$X_{1,K}^*$
$t = 2$	$X_{2,1}^*$	$X_{2,2}^*$...	$X_{2,K}^*$
	⋮	⋮	⋮	⋮

Tabular model of MABs

Time			...	
$t = 1$	$X_{1,1}^*$	Y_1	...	$X_{1,K}^*$
$t = 2$	Y_2	$X_{2,2}^*$...	$X_{2,K}^*$
	⋮	⋮	⋮	⋮

Hypothetical dataset

Hypothetical table

$$\mathcal{D}_{\infty}^* = X_{\infty}^* \cup \{W_t\}_{t=1}^{\infty}$$

Random seeds

Hypothetical dataset

Given $\mathcal{D}_\infty^* = X_\infty^* \cup \{W_t\}_{t=1}^\infty$

→ C , \mathcal{T} and $N_k(t)$ for each t and k can be expressed as some functions of \mathcal{D}_∞^* .

Monotone effect of a sample

Theorem

Suppose arm k has a finite mean. If $\frac{\mathbf{1}(C)}{N_k(\mathcal{T})}$ is an **increasing** function of each $X_{i,k}^*$ while keeping all other entries in \mathcal{D}_∞^* fixed then we have

$$\mathbb{E} \left[\widehat{F}_{k,\mathcal{T}}(y) - F_k(y) \mid C \right] \leq 0 \quad (\text{Negative conditional bias of the empirical CDF})$$

Monotone effect of a sample

Theorem

Suppose arm k has a finite mean. If $\frac{\mathbf{1}(C)}{N_k(\mathcal{T})}$ is an **increasing** function of each $X_{i,k}^*$ while keeping all other entries in \mathcal{D}_∞^* fixed then we have

$$\mathbb{E} \left[\widehat{F}_{k,\mathcal{T}}(y) - F_k(y) \mid C \right] \leq 0 \quad \text{(Negative conditional bias of the empirical CDF)}$$

$$\mathbb{E} \left[\widehat{\mu}_k(\mathcal{T}) - \mu_k \mid C \right] \geq 0 \quad \text{(Positive conditional bias of the sample mean)}$$

Monotone effect of a sample

Theorem

Suppose arm k has a finite mean. If $\frac{\mathbf{1}(C)}{N_k(\mathcal{T})}$ is a **decreasing** function of each $X_{i,k}^*$ while keeping all other entries in \mathcal{D}_∞^* fixed then we have

$$\mathbb{E} \left[\widehat{F}_{k,\mathcal{T}}(y) - F_k(y) \mid C \right] \geq 0 \quad (\text{Positive conditional bias of the empirical CDF})$$

Monotone effect of a sample

Theorem

Suppose arm k has a finite mean. If $\frac{\mathbf{1}(C)}{N_k(\mathcal{T})}$ is a **decreasing** function of each $X_{i,k}^*$ while keeping all other entries in \mathcal{D}_∞^* fixed then we have

$$\mathbb{E} \left[\widehat{F}_{k,\mathcal{T}}(y) - F_k(y) \mid C \right] \geq 0$$

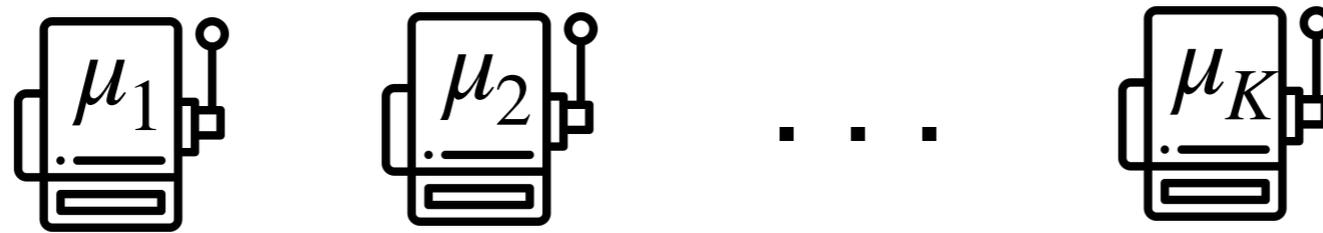
(Positive conditional bias of the empirical CDF)

$$\mathbb{E} \left[\widehat{\mu}_k(\mathcal{T}) - \mu_k \mid C \right] \leq 0$$

(Negative conditional bias of the sample mean)

E.g.: Best arm identification

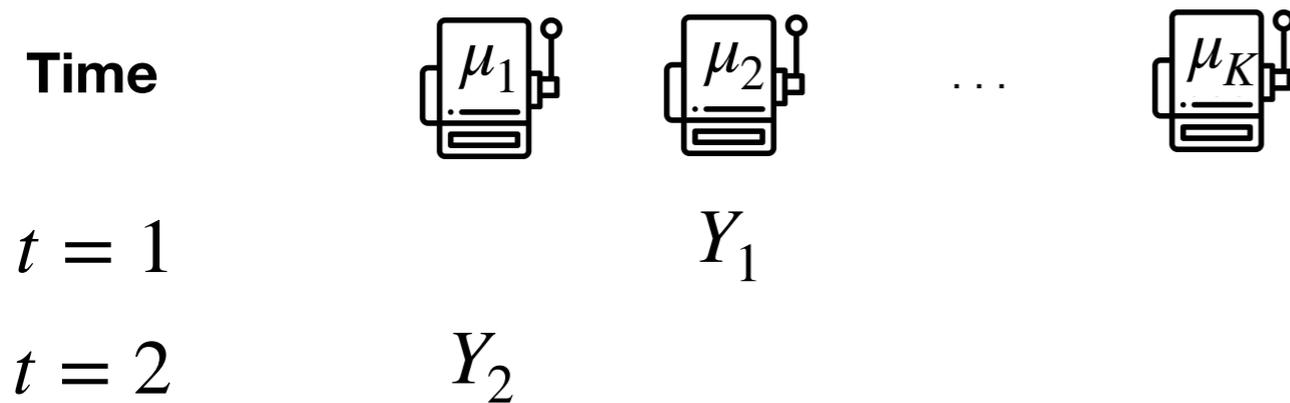
- K prototype items



➔ Want to figure out which one has the largest revenue.

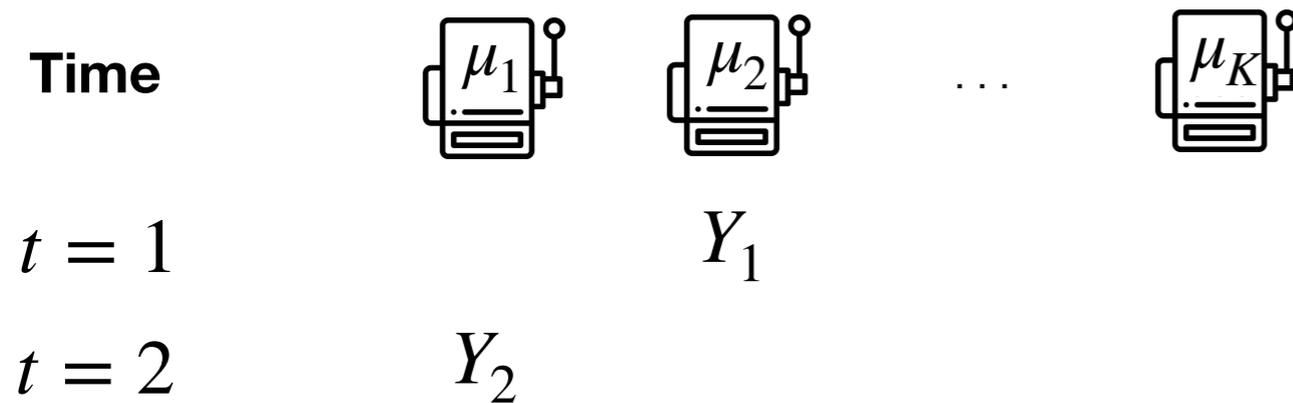
E.g.: Best arm identification

lil' UCB algorithm



E.g.: Best arm identification

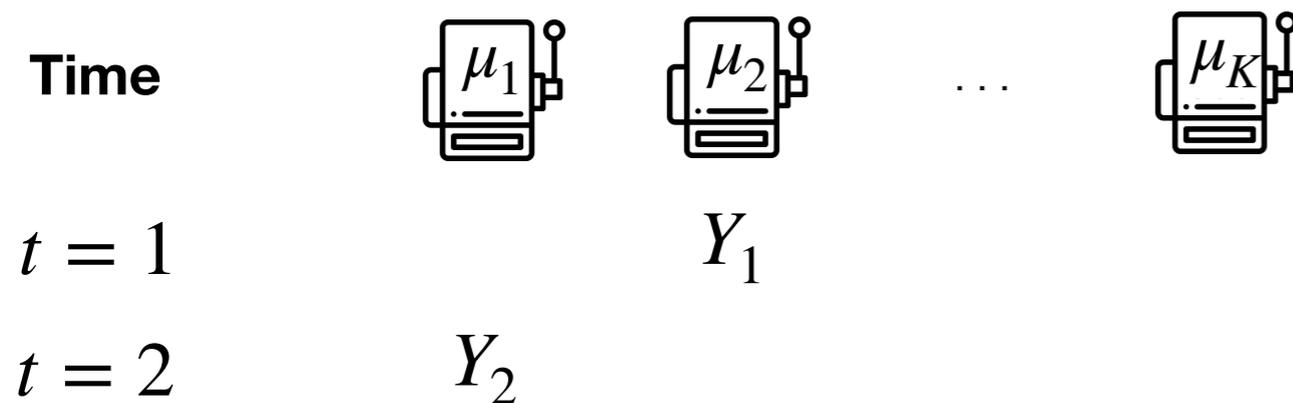
lil' UCB algorithm



$$A_t = \arg \max_k \hat{\mu}_k(t) + u(N_k(t)) \quad \text{(Upper confidence bound)}$$

E.g.: Best arm identification

lil' UCB algorithm

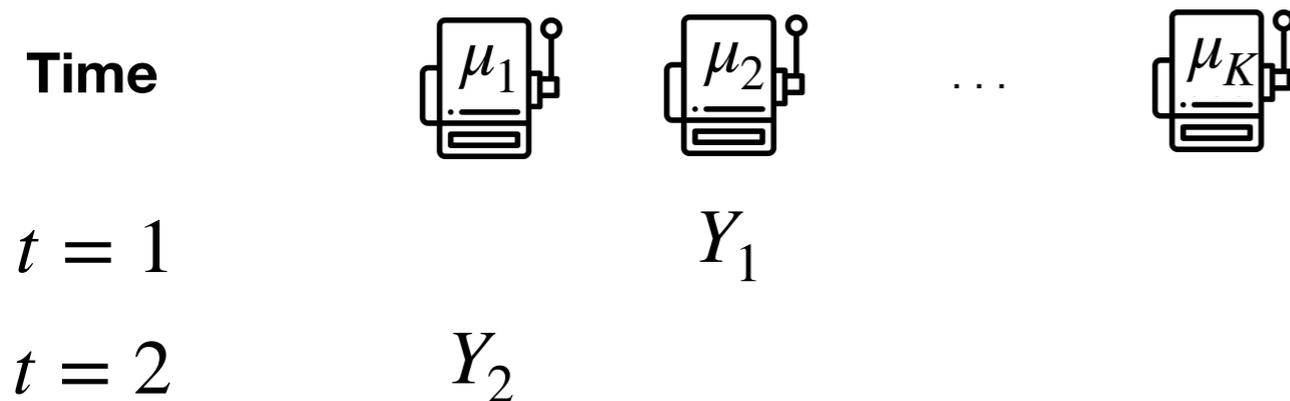


$$A_t = \arg \max_k \hat{\mu}_k(t) + u(N_k(t)) \quad \text{(Upper confidence bound)}$$

$$\mathcal{T} = \inf \left\{ t : \exists k, N_k(t) \geq 1 + \lambda \sum_{j \neq k} N_j(t) \right\}$$

E.g.: Best arm identification

lil' UCB algorithm



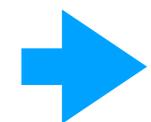
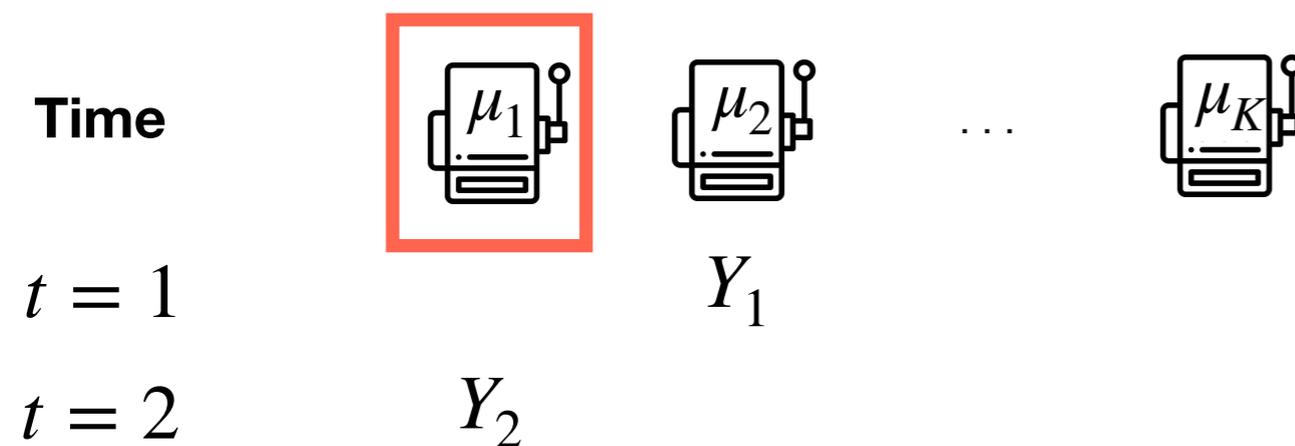
$$A_t = \arg \max_k \hat{\mu}_k(t) + u(N_k(t)) \quad \text{(Upper confidence bound)}$$

$$\mathcal{T} = \inf \left\{ t : \exists k, N_k(t) \geq 1 + \lambda \sum_{j \neq k} N_j(t) \right\}$$

$$\kappa = \arg \max_k N_k(\mathcal{T})$$

E.g.: Best arm identification

lil' UCB algorithm



a) Item 1 is chosen as the best.

b) Item 1 is NOT chosen as the best.

E.g.: Best arm identification

a) Item 1 is chosen as the best ($\kappa = 1$).



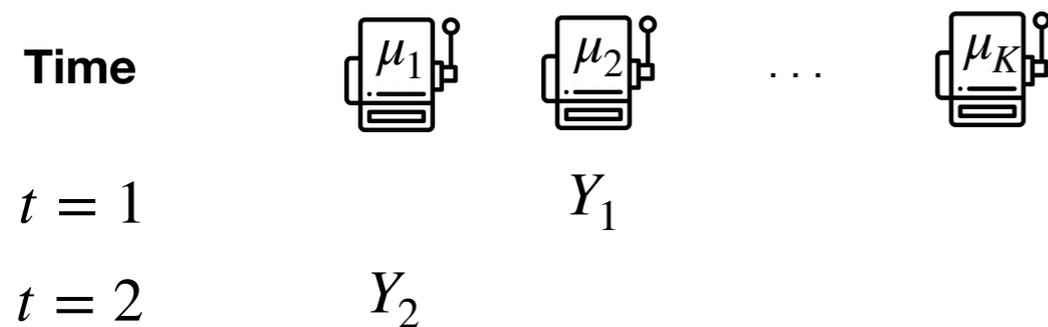
$$A_t = \arg \max_k \hat{\mu}_k(t) + u(N_k(t))$$

$$\mathcal{T} = \inf \left\{ t : \exists k, N_k(t) \geq 1 + \lambda \sum_{j \neq k} N_j(t) \right\}$$

$$\kappa = \arg \max_k N_k(\mathcal{T})$$

E.g.: Best arm identification

a) Item 1 is chosen as the best ($\kappa = 1$).



Sample from item 1 $\xrightarrow{\text{Increasing}}$ $\frac{\mathbf{1}(\kappa = 1)}{N_1(\mathcal{T})}$

$$A_t = \arg \max_k \hat{\mu}_k(t) + u(N_k(t))$$

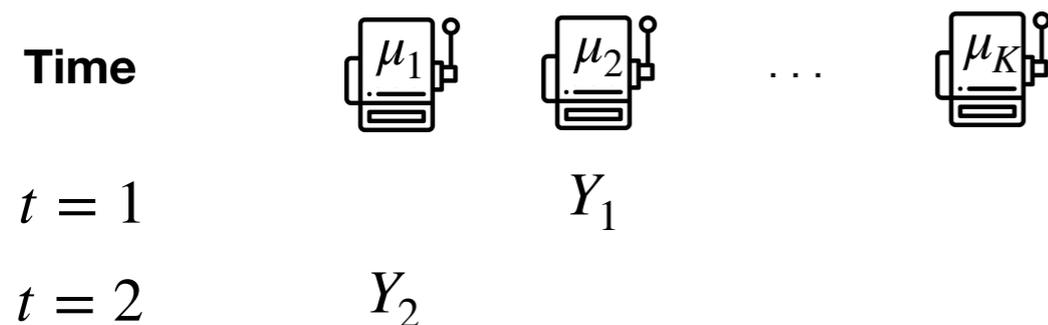
$$\mathcal{T} = \inf \left\{ t : \exists k, N_k(t) \geq 1 + \lambda \sum_{j \neq k} N_j(t) \right\}$$

$$\kappa = \arg \max_k N_k(\mathcal{T})$$

Negative conditional bias of the empirical CDF

E.g.: Best arm identification

a) Item 1 is chosen as the best ($\kappa = 1$).



Sample from item 1 $\xrightarrow{\text{Increasing}}$ $\frac{\mathbf{1}(\kappa = 1)}{N_1(\mathcal{T})}$

$$A_t = \arg \max_k \hat{\mu}_k(t) + u(N_k(t))$$

$$\mathcal{T} = \inf \left\{ t : \exists k, N_k(t) \geq 1 + \lambda \sum_{j \neq k} N_j(t) \right\}$$

$$\kappa = \arg \max_k N_k(\mathcal{T})$$

➔ **Negative** conditional bias of the empirical CDF

➔ **Positive** conditional bias of the sample mean

E.g.: Best arm identification

b) Item 1 is **NOT** chosen as the best ($\kappa \neq 1$).



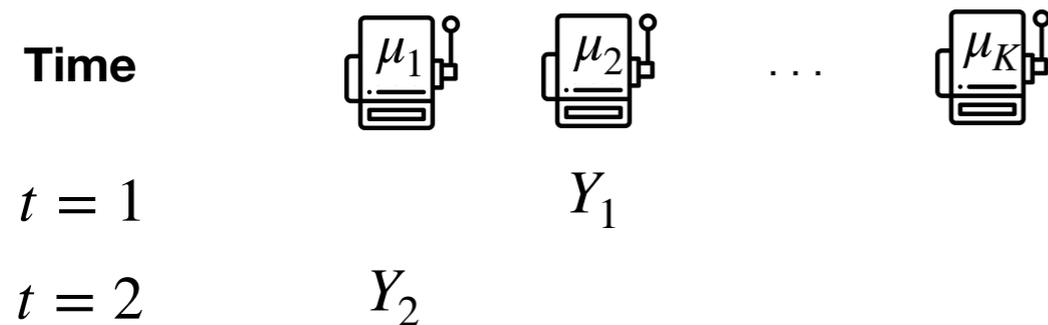
$$A_t = \arg \max_k \hat{\mu}_k(t) + u(N_k(t))$$

$$\mathcal{T} = \inf \left\{ t : \exists k, N_k(t) \geq 1 + \lambda \sum_{j \neq k} N_j(t) \right\}$$

$$\kappa = \arg \max_k N_k(\mathcal{T})$$

E.g.: Best arm identification

b) Item 1 is **NOT** chosen as the best ($\kappa \neq 1$).



Sample from item 1 $\xrightarrow{\text{Decreasing}}$ $\frac{\mathbf{1}(\kappa \neq 1)}{N_1(\mathcal{T})}$

$$A_t = \arg \max_k \hat{\mu}_k(t) + u(N_k(t))$$

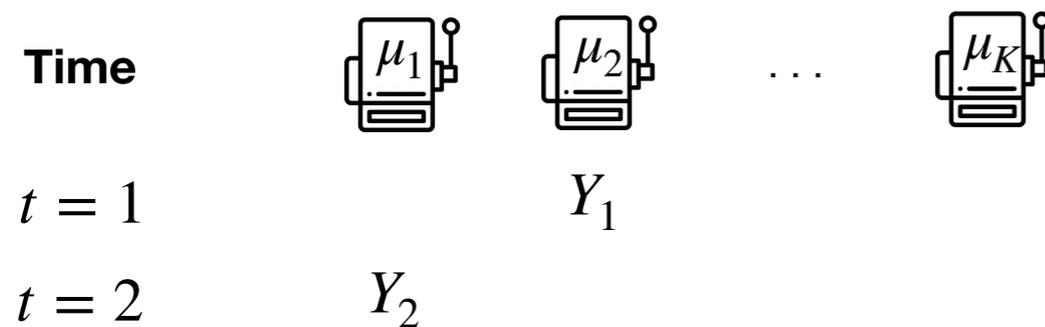
$$\mathcal{T} = \inf \left\{ t : \exists k, N_k(t) \geq 1 + \lambda \sum_{j \neq k} N_j(t) \right\}$$

$$\kappa = \arg \max_k N_k(\mathcal{T})$$

Positive conditional bias of the empirical CDF

E.g.: Best arm identification

b) Item 1 is **NOT** chosen as the best ($\kappa \neq 1$).



Sample from item 1 $\xrightarrow{\text{Decreasing}}$ $\frac{\mathbf{1}(\kappa \neq 1)}{N_1(\mathcal{T})}$

$$A_t = \arg \max_k \hat{\mu}_k(t) + u(N_k(t))$$

$$\mathcal{T} = \inf \left\{ t : \exists k, N_k(t) \geq 1 + \lambda \sum_{j \neq k} N_j(t) \right\}$$

$$\kappa = \arg \max_k N_k(\mathcal{T})$$

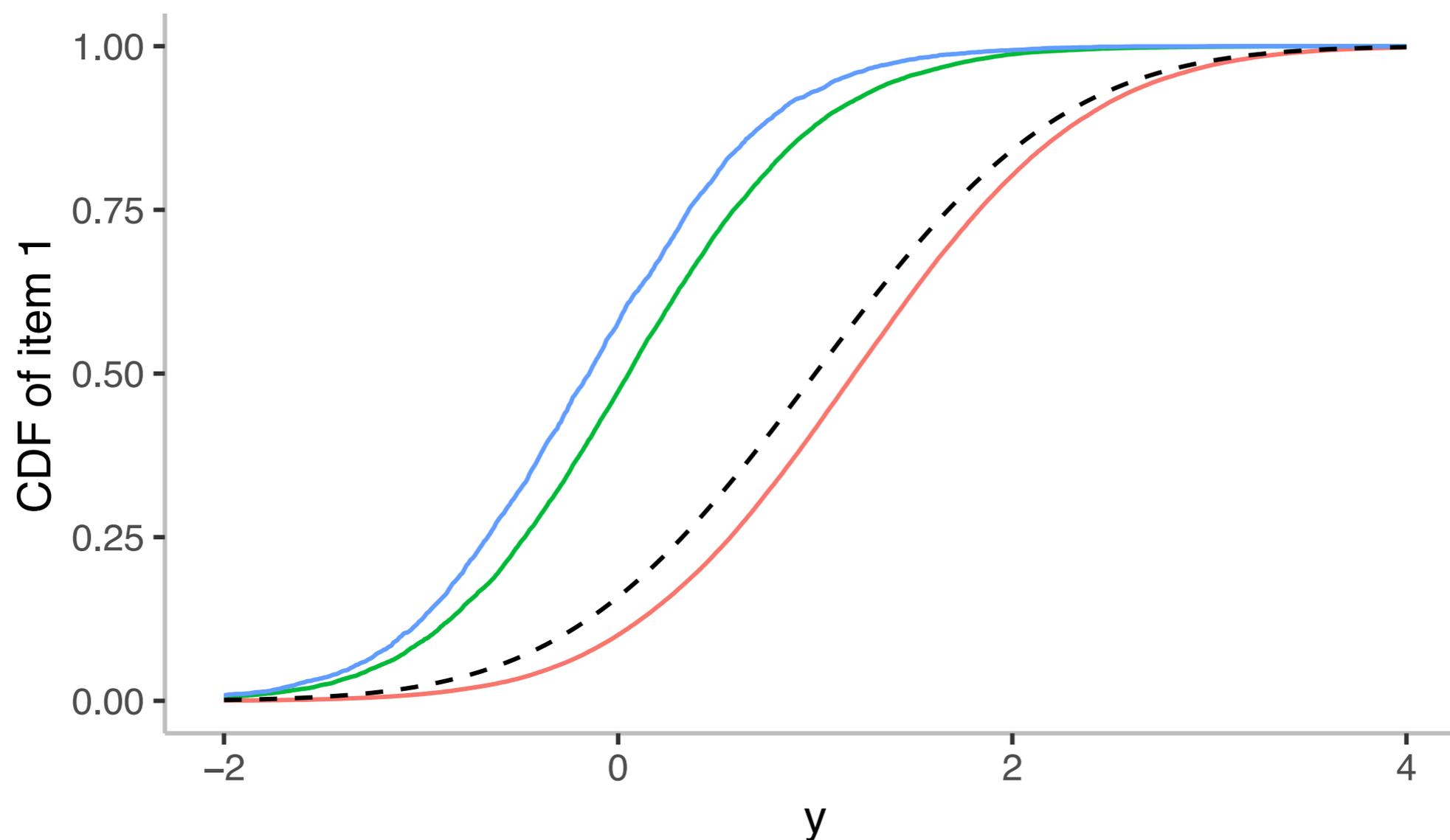
➔ **Positive** conditional bias of the empirical CDF

➔ **Negative** conditional bias of the sample mean

Average of the empirical CDF of item 1 conditioned on each event

lil'UCB on 3 items ($\mu_1 = 1$)

Mean bias = (0.2, -0.93, -1.14)



— Item 1 is chosen — Item 2 is chosen — Item 3 is chosen

Thank you!

On conditional versus marginal bias in multi-armed bandits

Jaehyeok Shin, Aaditya Ramdas and Alessandro Rinaldo



**Carnegie
Mellon
University**