

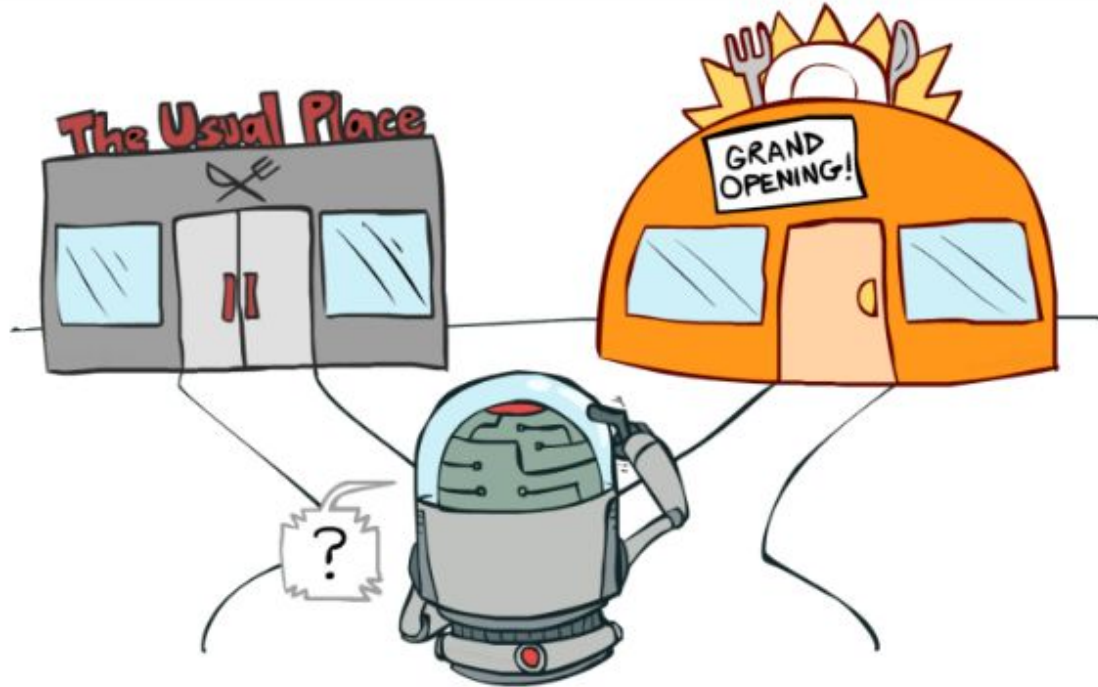
# Multi-Armed Bandits with Bounded Arm-Memory: Near-Optimal Guarantees for Best-Arm Identification and Regret Minimization

Arnab Maiti, Indian Institute of Technology, Kharagpur

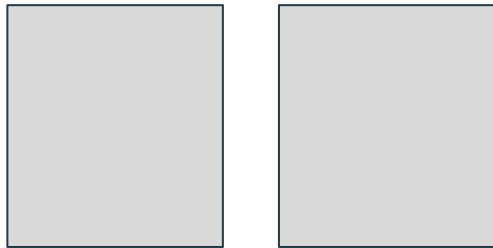
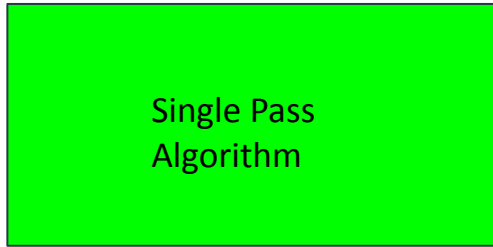
Vishakha Patil, Indian Institute of Science, Bangalore

Arindam Khan, Indian Institute of Science, Bangalore

# Multi-Armed Bandits

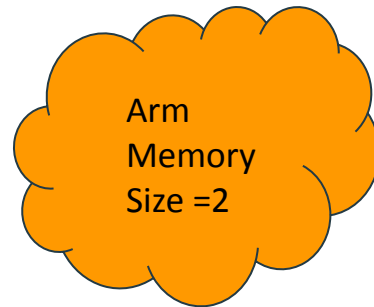
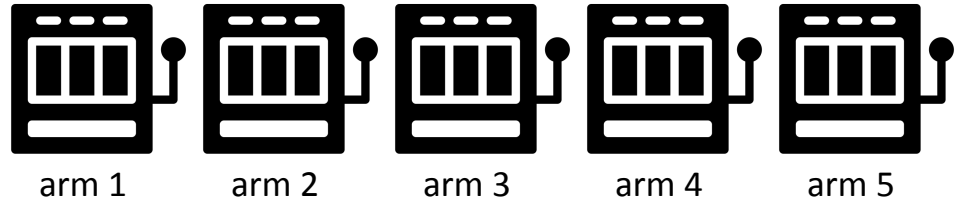


# Streaming Model

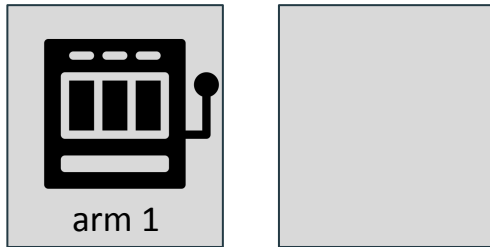
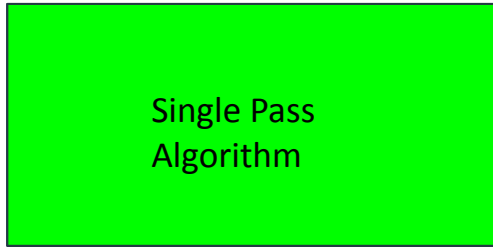


Arm Memory

Adversarial order arrival

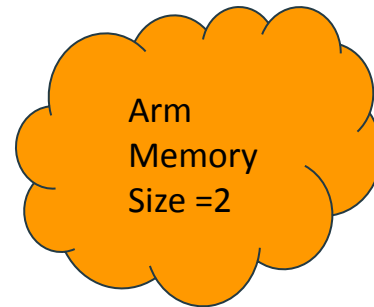
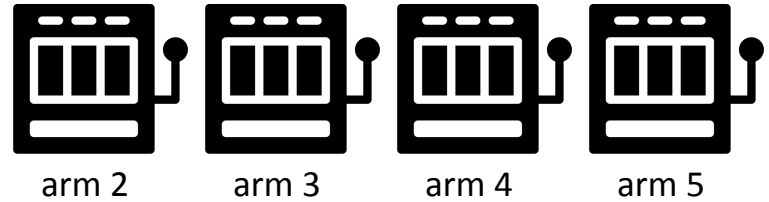


# Streaming Model

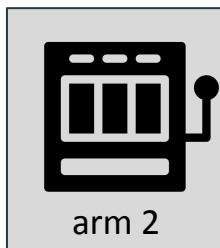
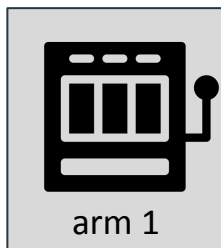
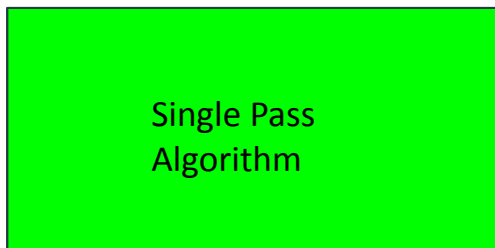


Arm Memory

Adversarial order arrival

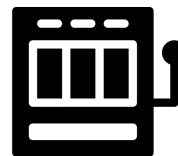


# Streaming Model

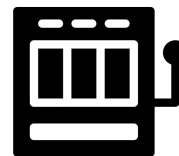


Arm Memory

Adversarial order arrival



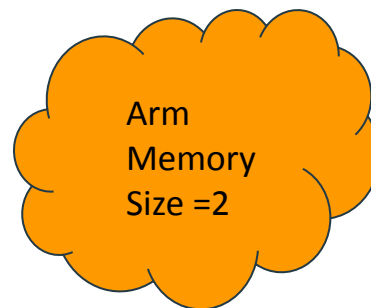
arm 3



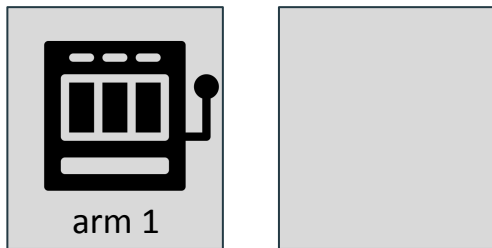
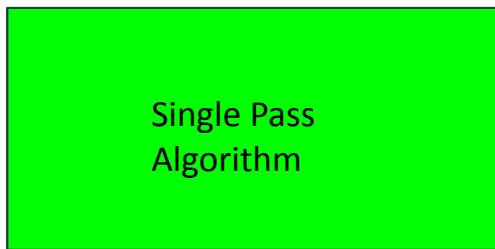
arm 4



arm 5

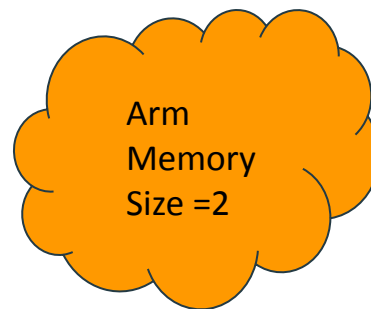
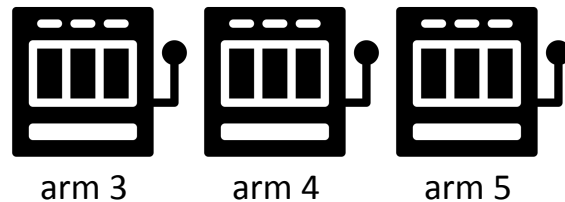


# Streaming Model



Arm Memory

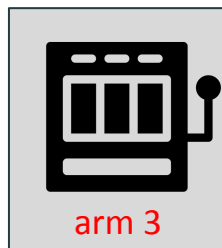
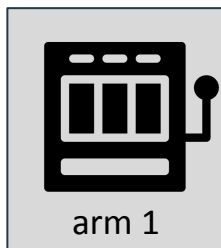
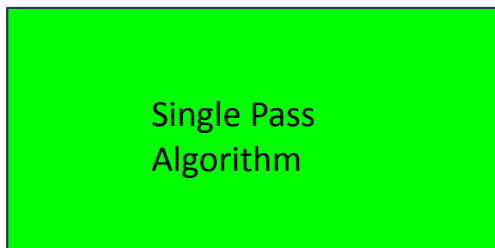
Adversarial order arrival



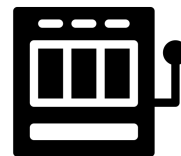
*Image source: slot machine by Deemak Daksina from the Noun Project*

# Streaming Model

Adversarial order arrival



Arm Memory



arm 4



arm 5

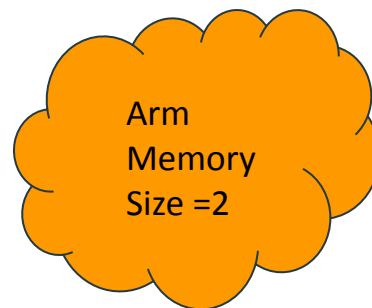
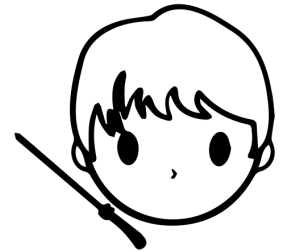
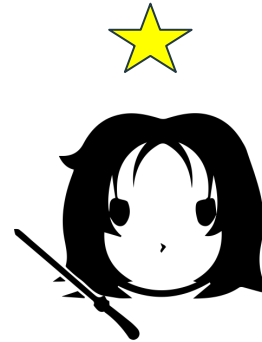
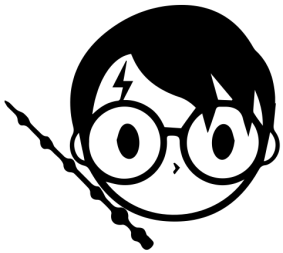
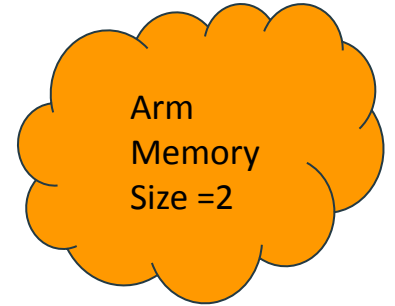


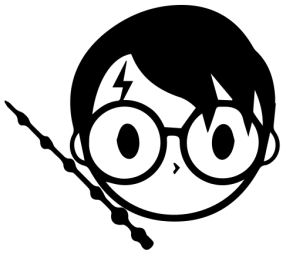
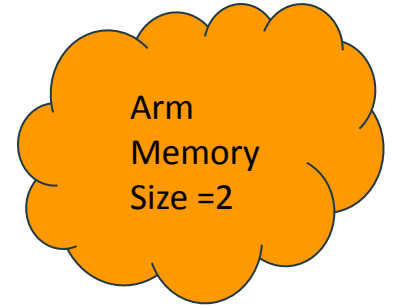
Image source: slot machine by Deemak Daksina from the Noun Project

# Application: Job Interviews

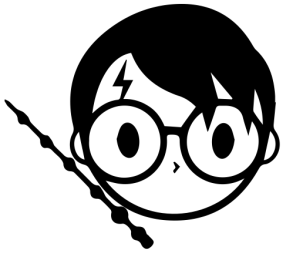
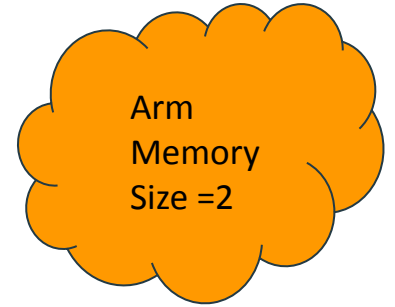




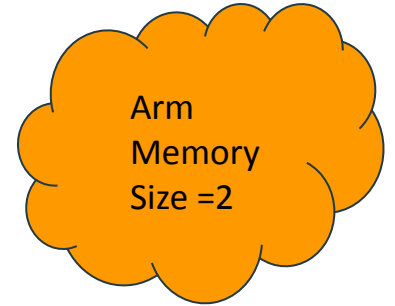
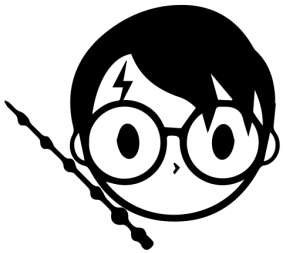
# Application: Job Interviews



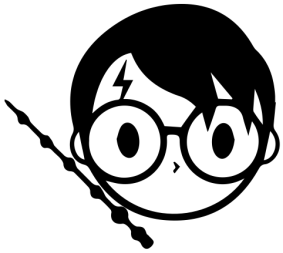
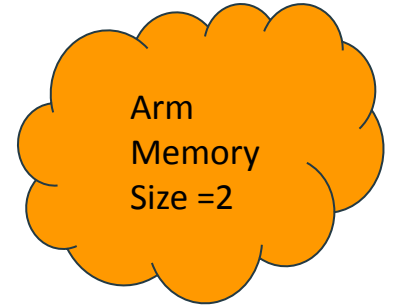
# Application: Job Interviews



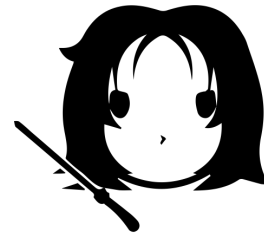
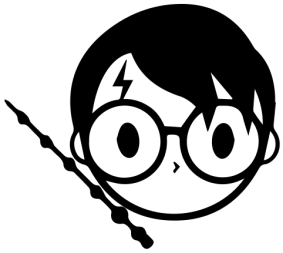
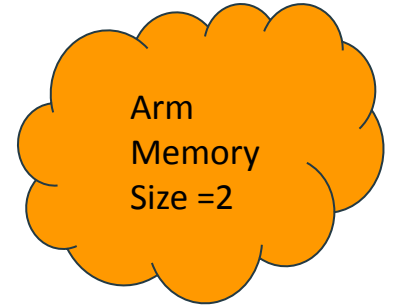
# Application: Job Interviews



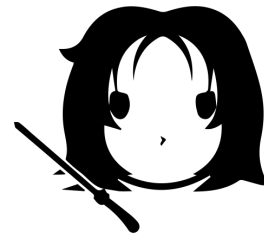
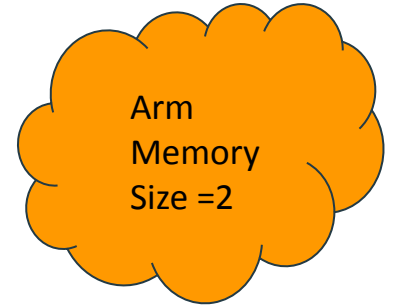
# Application: Job Interviews



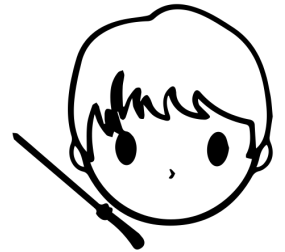
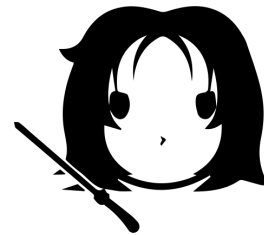
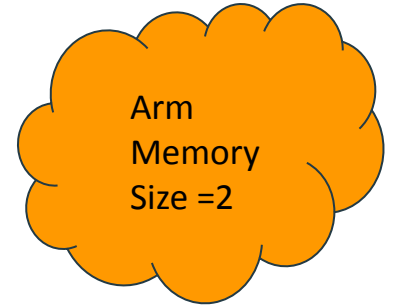
# Application: Job Interviews



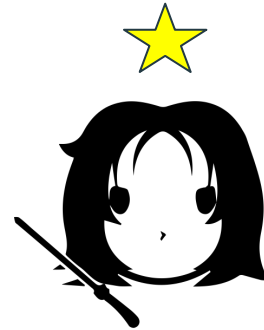
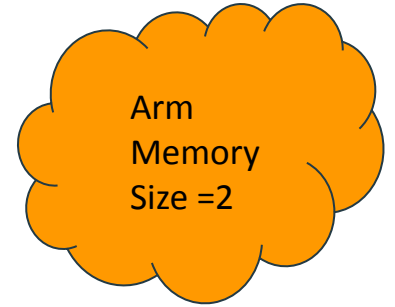
# Application: Job Interviews



# Application: Job Interviews



# Application: Job Interviews





# Regret Minimization

Stochastic MAB Instance with expected rewards

$$\langle \mu_1, \mu_2, \dots, \mu_n \rangle$$

Let  $\mu^* = \max(\mu_1, \dots, \mu_n)$

Let  $\mu_{i_t}$  be the mean of the arm sampled in the time step  $t$

Expected Cumulative Regret in  $T$  time steps

$$\mathbb{E}[R(T)] = \mu^* \cdot T - \sum_{t=1}^T \mathbb{E}[\mu_{i_t}]$$

Goal: Minimize  $\mathbb{E}[R(T)]$

# Best-Arm Identification

Stochastic MAB Instance with expected rewards

$$\langle \mu_1, \mu_2, \dots, \mu_n \rangle$$

$\mathcal{E}$ -best arm is an arm with mean at least  $\mu^* - \mathcal{E}$

$(\mathcal{E}, \delta)$ -PAC Algorithm: Finds an  $\mathcal{E}$ -best arm with probability at least  $1 - \delta$

Goal:  $(\mathcal{E}, \delta)$ -PAC algorithm such that total number of arm pulls is minimized.

# Results Overview

# Results Overview

**Main Result 1:** Any single-pass algorithm will incur at least

$$\Omega\left(\frac{n^{1/3}T^{2/3}}{m^{7/3}}\right)$$

expected cumulative regret when  $m < n$ .

This lower bound even holds for random order arrival.

# Results Overview

**Main Result 1:** Any single-pass algorithm will incur at least

$$\Omega\left(\frac{n^{1/3}T^{2/3}}{m^{7/3}}\right)$$

expected cumulative regret when  $m < n$ .

This lower bound even holds for random order arrival.

**Main Result 2:**  $(\epsilon, \delta)$  - PAC algorithm with

$O(r)$  arm memory

Optimal  $\mathcal{R}$ -round sample complexity

# Results Overview

**Main Result 1:** Any single-pass algorithm will incur at least

$$\Omega\left(\frac{n^{1/3}T^{2/3}}{m^{7/3}}\right)$$

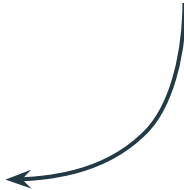
expected cumulative regret when  $m < n$ .

This lower bound even holds for random order arrival.

**Main Result 2:**  $(\varepsilon, \delta)$  - PAC algorithm with

$O(r)$  arm memory

Optimal  $\mathcal{R}$ -round sample complexity

$$O\left(\frac{n}{\varepsilon^2} (\mathbf{i}\log^{(r)}(n) + \log(\frac{1}{\delta}))\right)$$


# Results Overview

**Main Result 1:** Any single-pass algorithm will incur at least

$$\Omega\left(\frac{n^{1/3}T^{2/3}}{m^{7/3}}\right)$$

expected cumulative regret when  $m < n$ .

This lower bound even holds for random order arrival.

$$O\left(\frac{n}{\varepsilon^2} (\mathbf{i}\log^{(r)}(n) + \log(\frac{1}{\delta}))\right)$$

**Main Result 2:**  $(\varepsilon, \delta)$  - PAC algorithm with

$O(r)$  arm memory

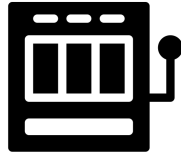
Optimal  $\mathcal{R}$ -round sample complexity

**Corollary:**  $(\varepsilon, \delta)$ - PAC algorithm with  $O(\log^* n)$  arm memory with optimal

worst-case sample complexity

# Regret Minimization Lower Bound

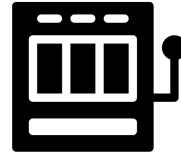
Input Instance 1



$1/2 + \epsilon$

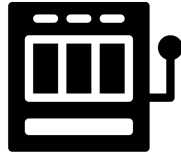


$1/2$

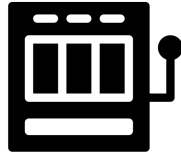


$1/2$

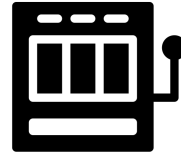
Input Instance 2



$1/2$



$1/2 + \epsilon$



$1/2$

Input Instance 3



$1/2$



$1/2$



$1$

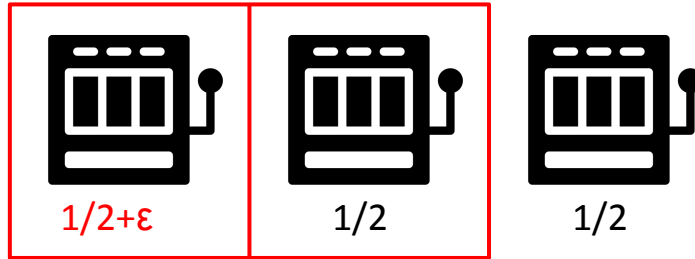
Means of the  
arms are  
specified below it

Arm  
Memory  
Size = 2



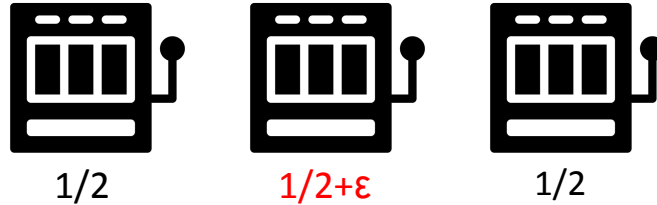
# Regret Minimization Lower Bound

Input Instance 1

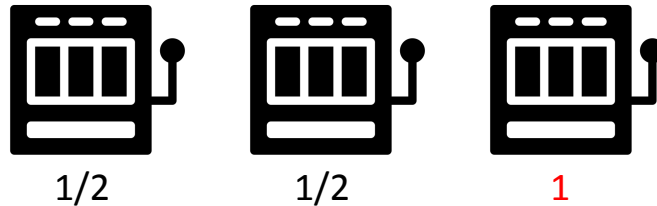


Means of the arms are specified below it

Input Instance 2



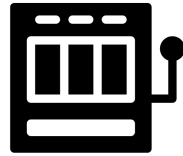
Input Instance 3



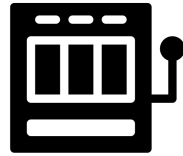
Arm Memory Size = 2

# Regret Minimization Lower Bound

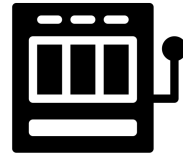
Input Instance 1



$1/2 + \epsilon$



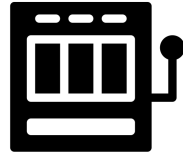
$1/2$



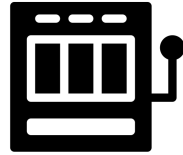
$1/2$

Means of the arms are specified below it

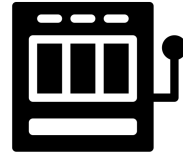
Input Instance 2



$1/2$



$1/2 + \epsilon$



$1/2$

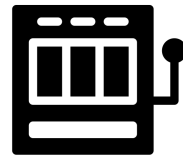
Input Instance 3



$1/2$



$1/2$

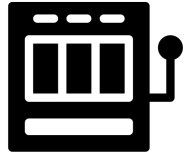


$1$

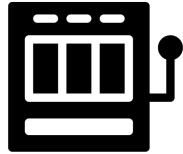
Arm  
Memory  
Size = 2

# Regret Minimization Lower Bound

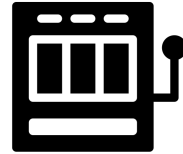
Input Instance 1



$1/2 + \epsilon$

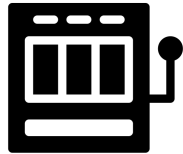


$1/2$

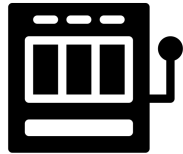


$1/2$

Input Instance 2



$1/2$



$1/2 + \epsilon$



$1/2$

Input Instance 3



$1/2$



$1/2$



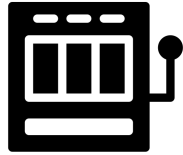
$1$

Means of the  
arms are  
specified below it

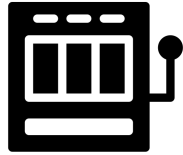
Arm  
Memory  
Size = 2

# Regret Minimization Lower Bound

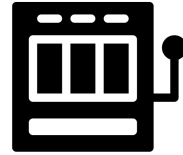
Input Instance 1



$1/2 + \epsilon$



$1/2$

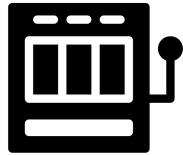


$1/2$

Input Instance 2



$1/2$

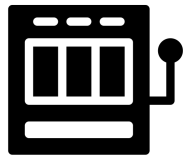


$1/2 + \epsilon$

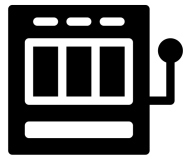


$1/2$

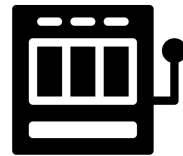
Input Instance 3



$1/2$



$1/2$



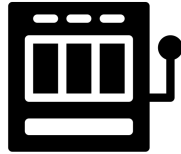
$1$

Means of the  
arms are  
specified below it

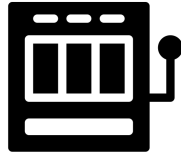
Arm  
Memory  
Size = 2

# Regret Minimization Lower Bound

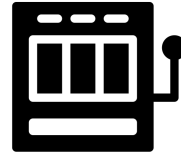
Input Instance 1



$1/2 + \epsilon$

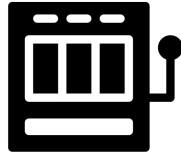


$1/2$

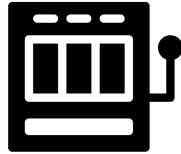


$1/2$

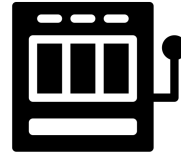
Input Instance 2



$1/2$



$1/2 + \epsilon$



$1/2$

Input Instance 3



$1/2$



$1/2$

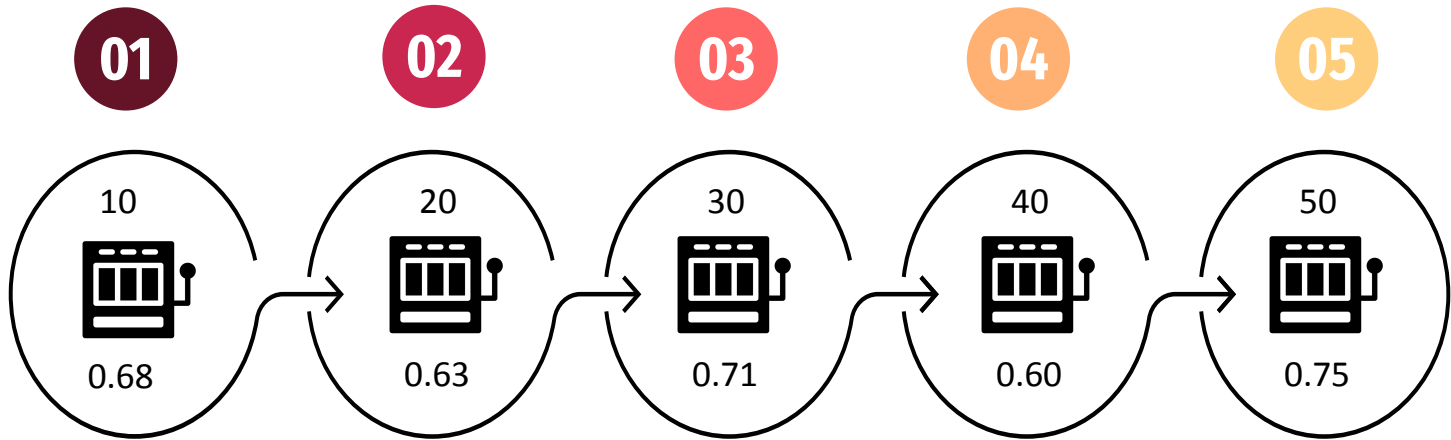


$1$

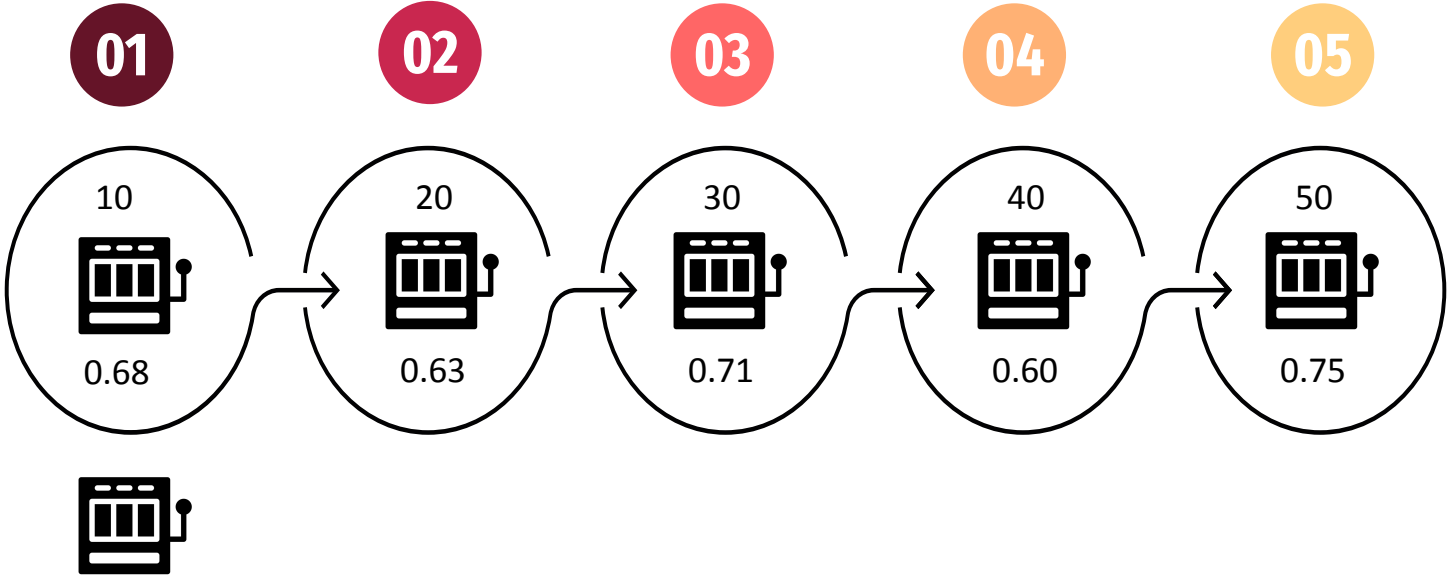
Means of the  
arms are  
specified below it

Arm  
Memory  
Size = 2

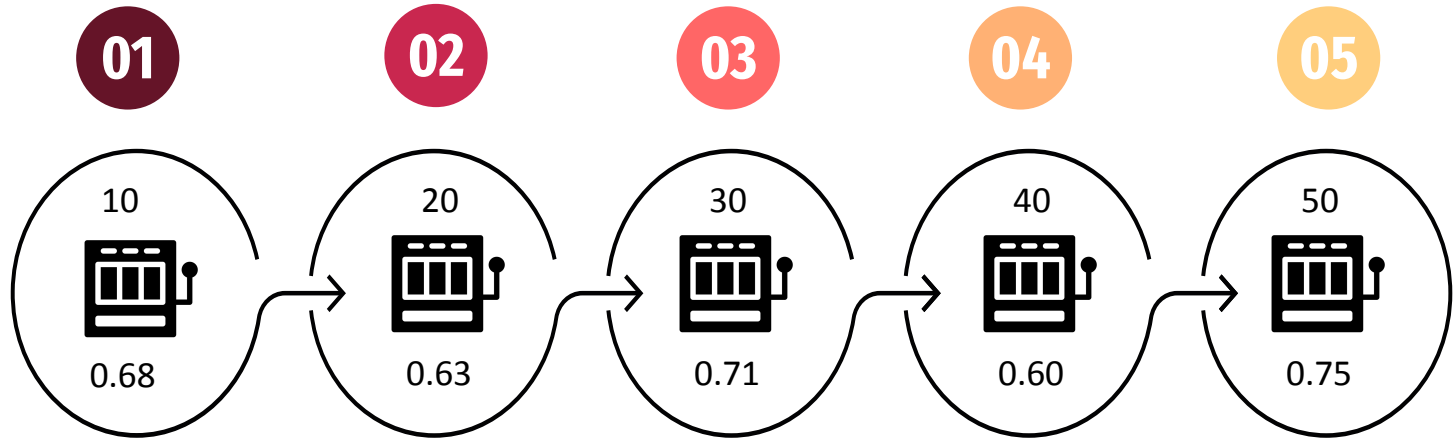
# Best-Arm Identification Algorithm



# Best-Arm Identification Algorithm



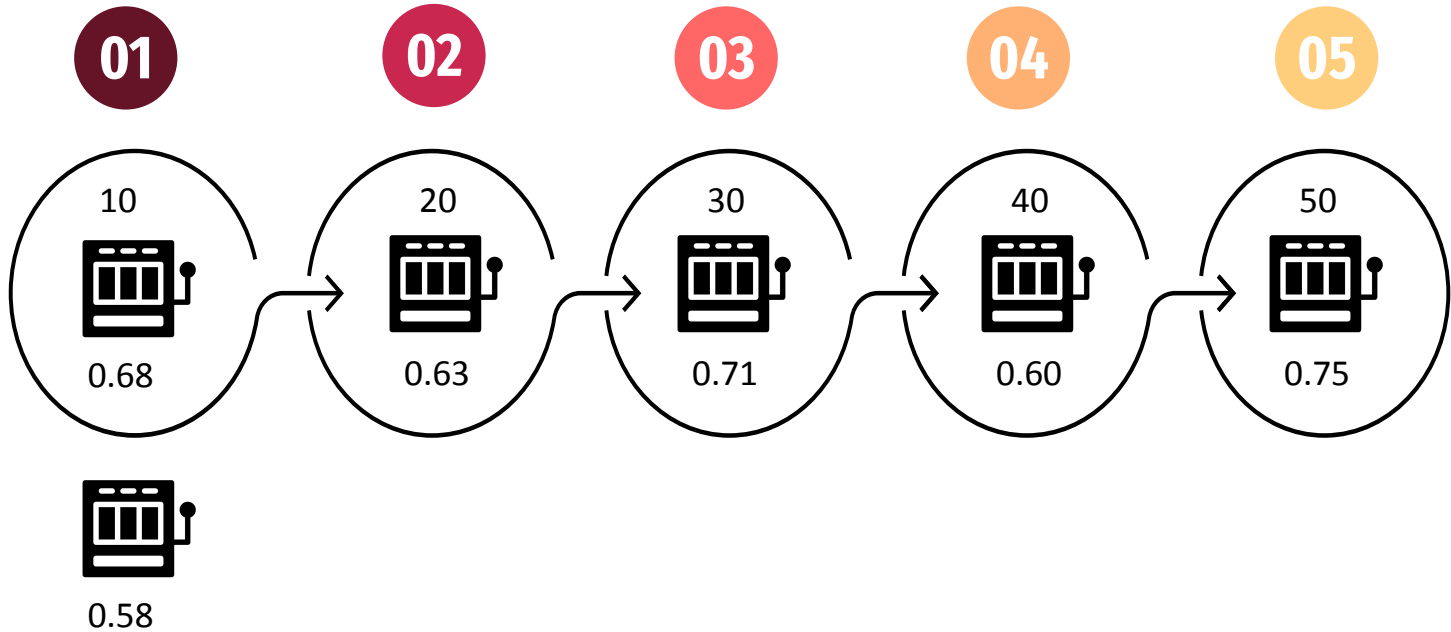
# Best-Arm Identification Algorithm



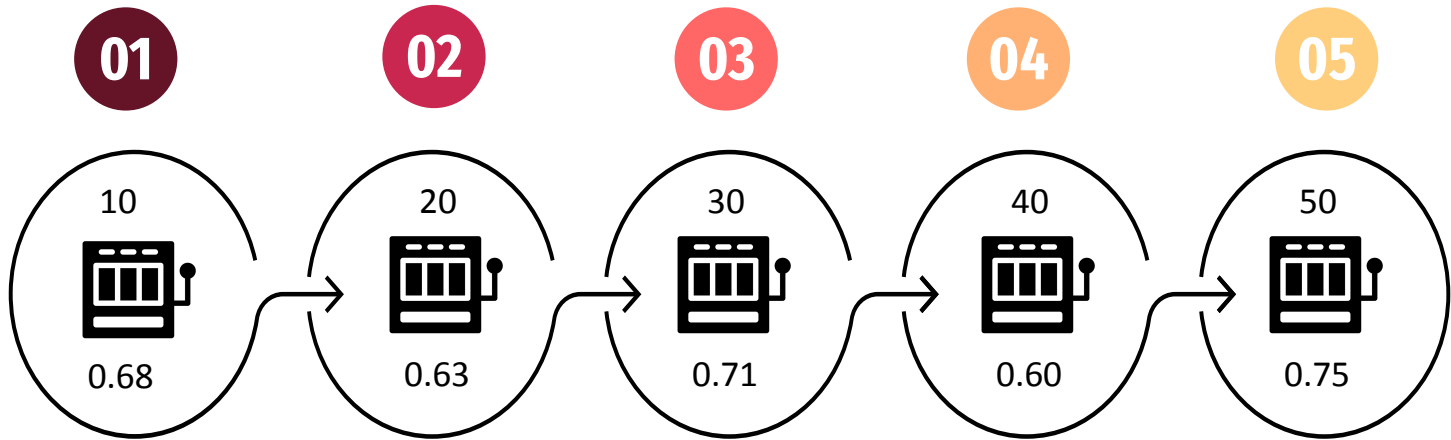
Sample for  $s_1$  times



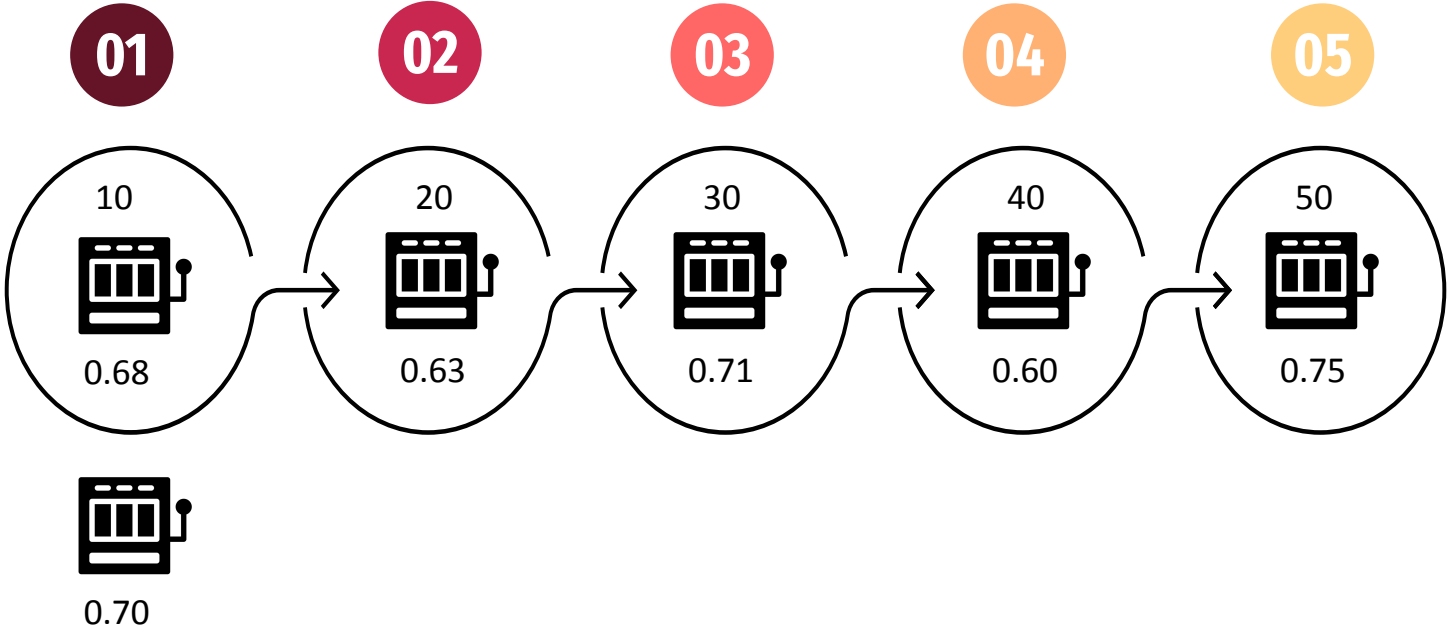
# Best-Arm Identification Algorithm



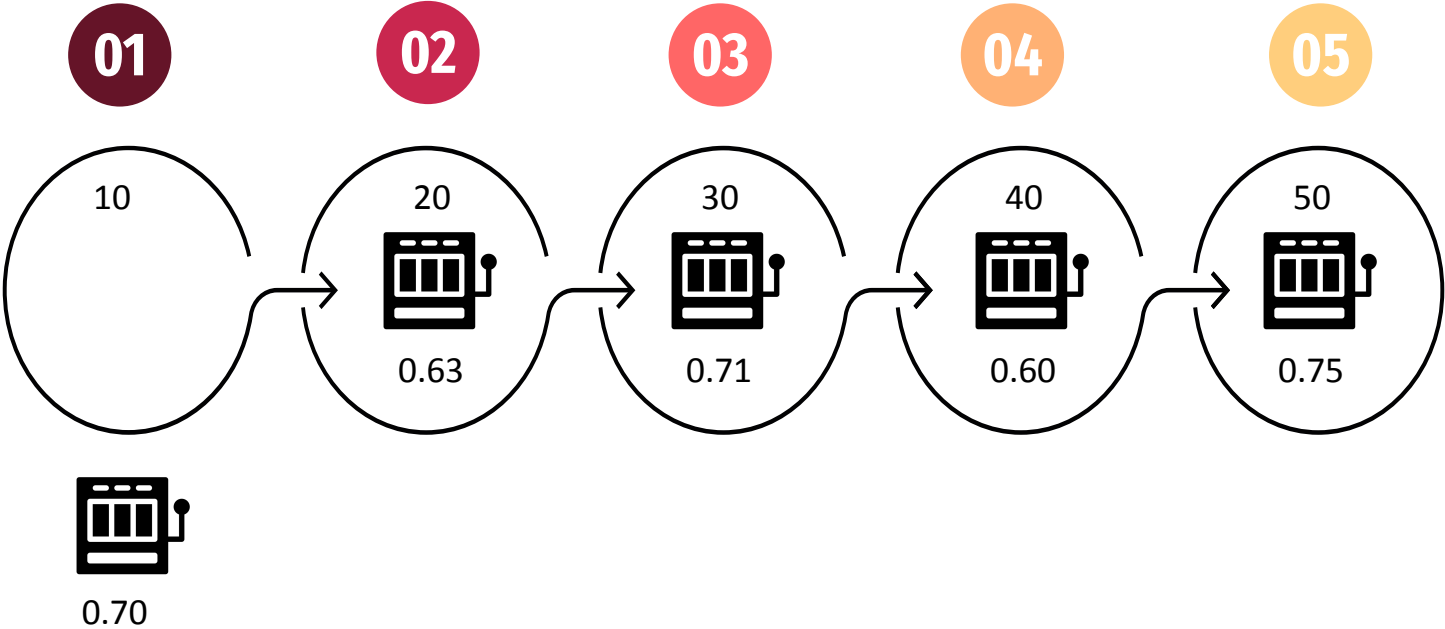
# Best-Arm Identification Algorithm



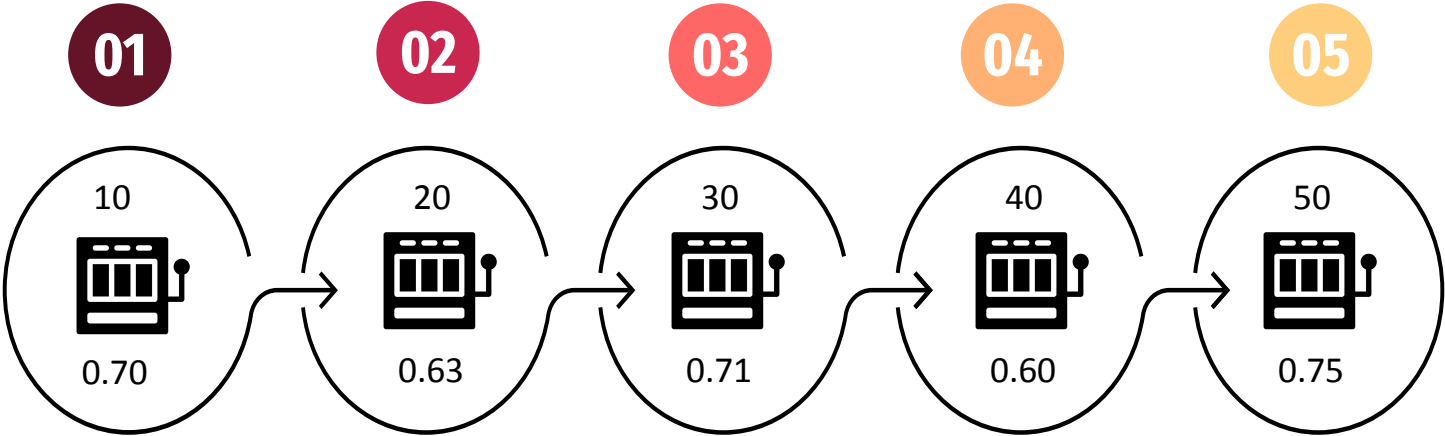
# Best-Arm Identification Algorithm



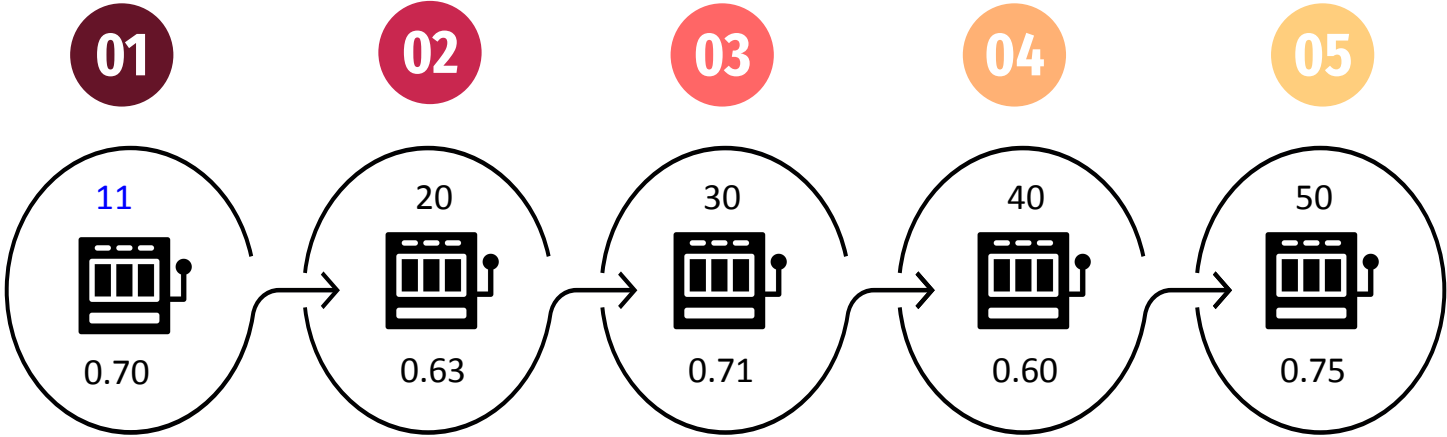
# Best-Arm Identification Algorithm



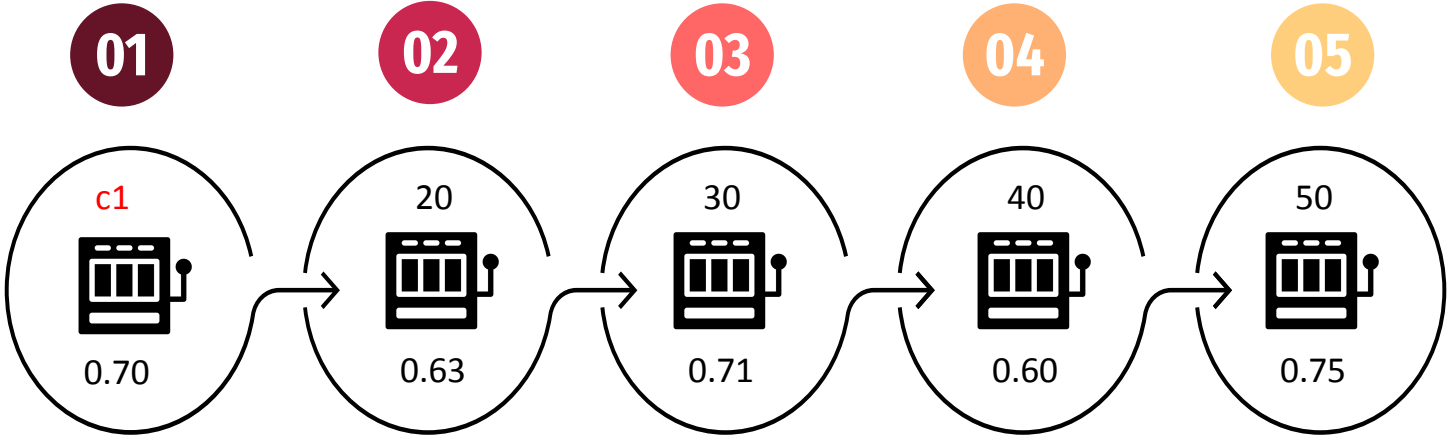
# Best-Arm Identification Algorithm



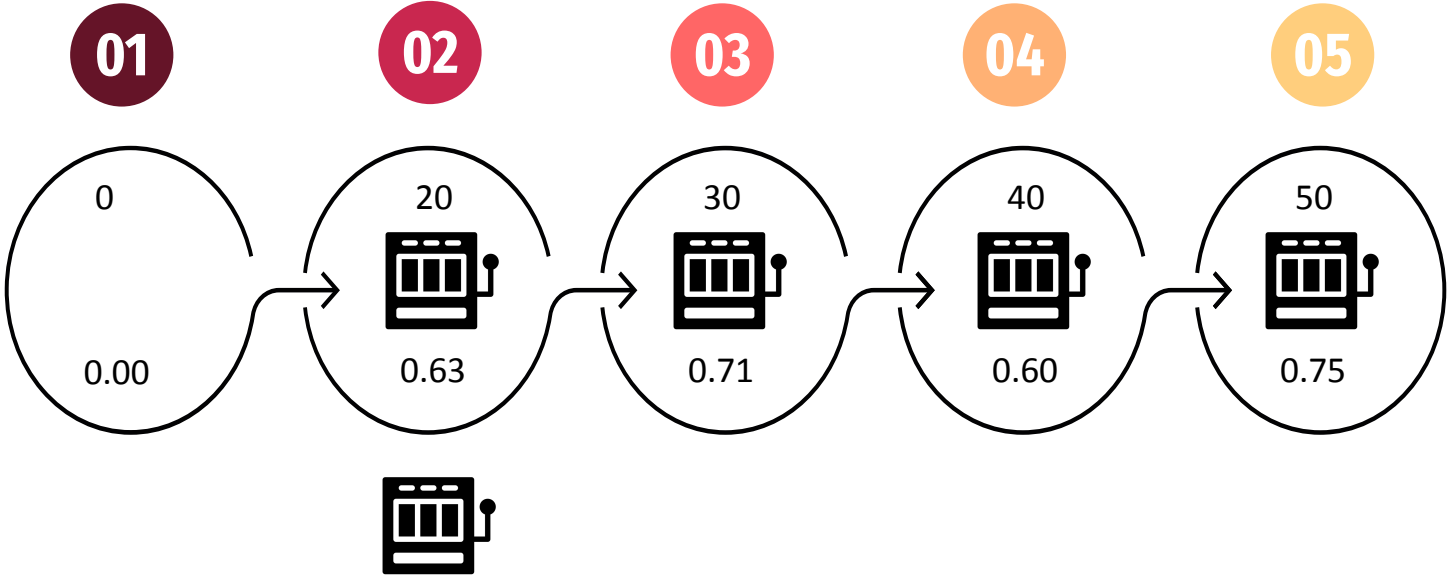
# Best-Arm Identification Algorithm



# Best-Arm Identification Algorithm

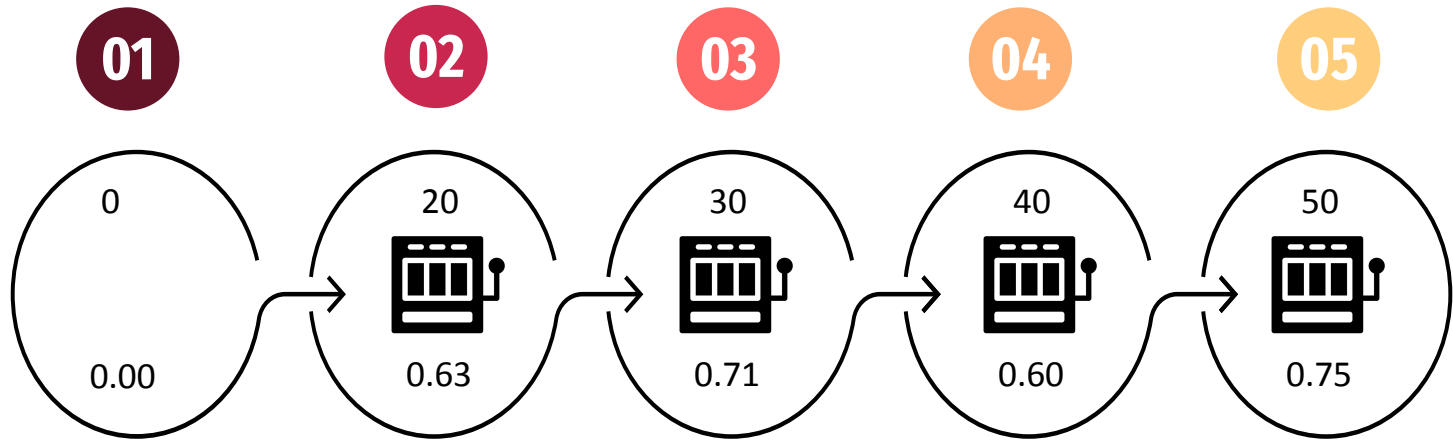


# Best-Arm Identification Algorithm



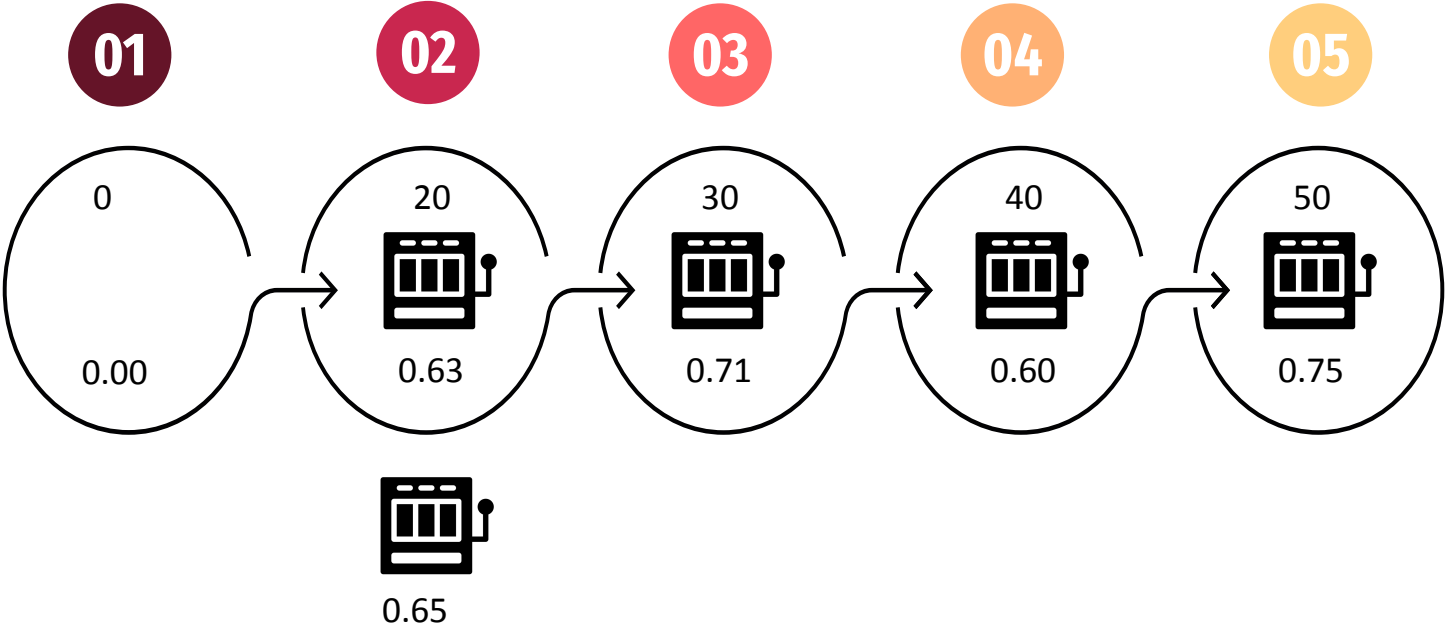


# Best-Arm Identification Algorithm

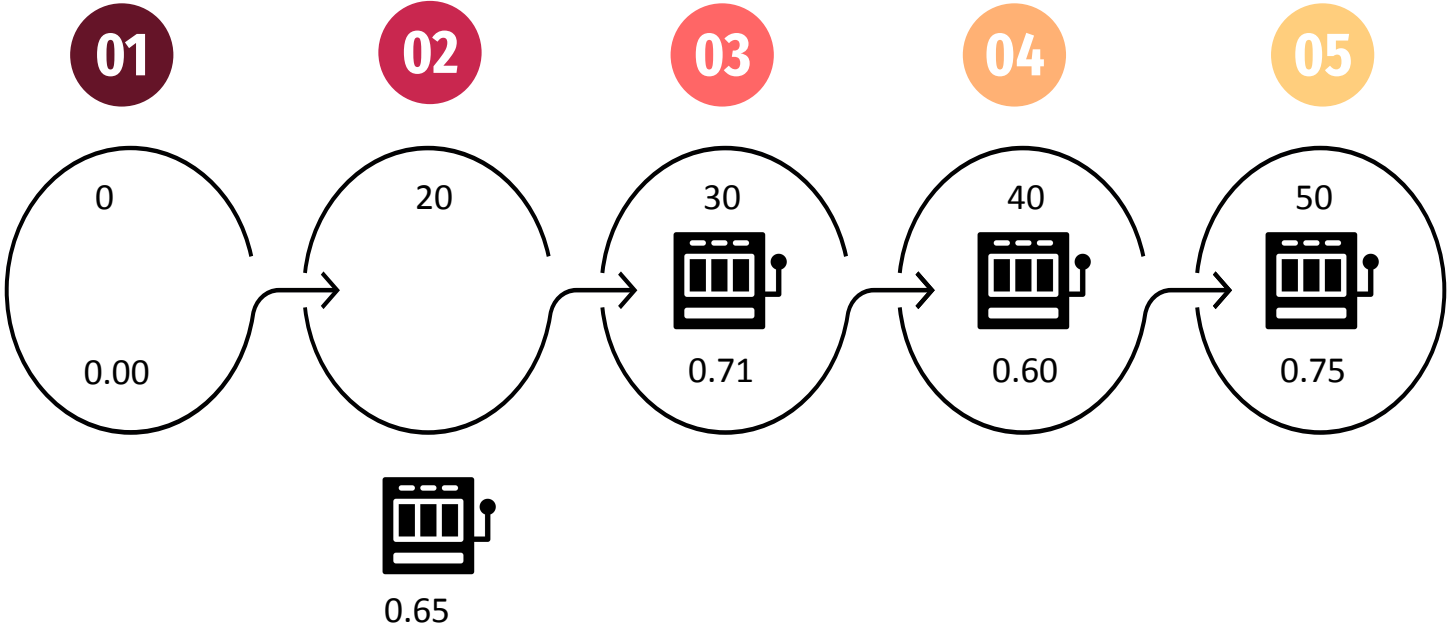


Sample for  $s_2$  times

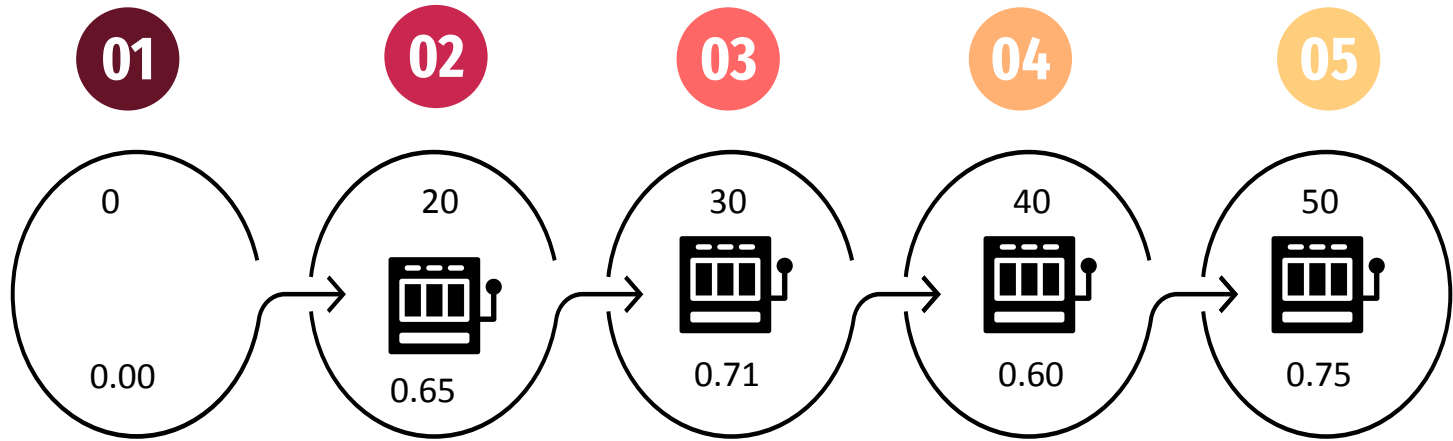
# Best-Arm Identification Algorithm



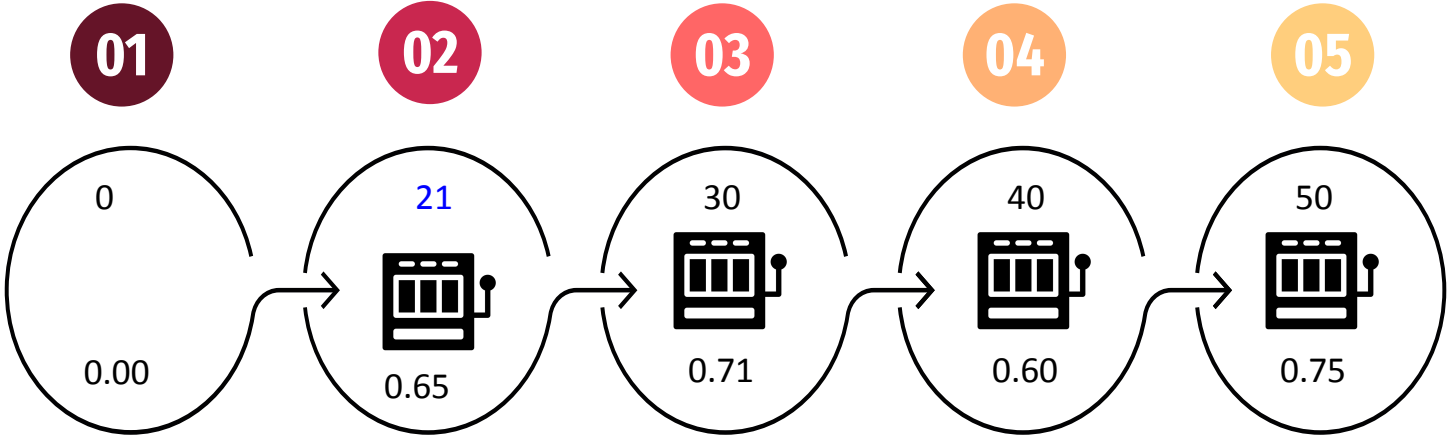
# Best-Arm Identification Algorithm



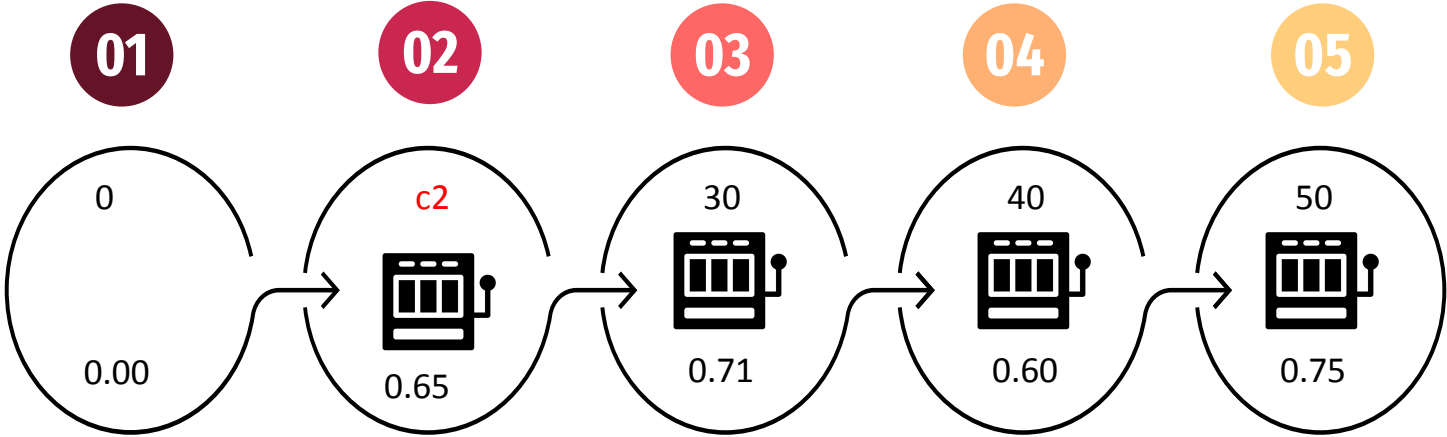
# Best-Arm Identification Algorithm



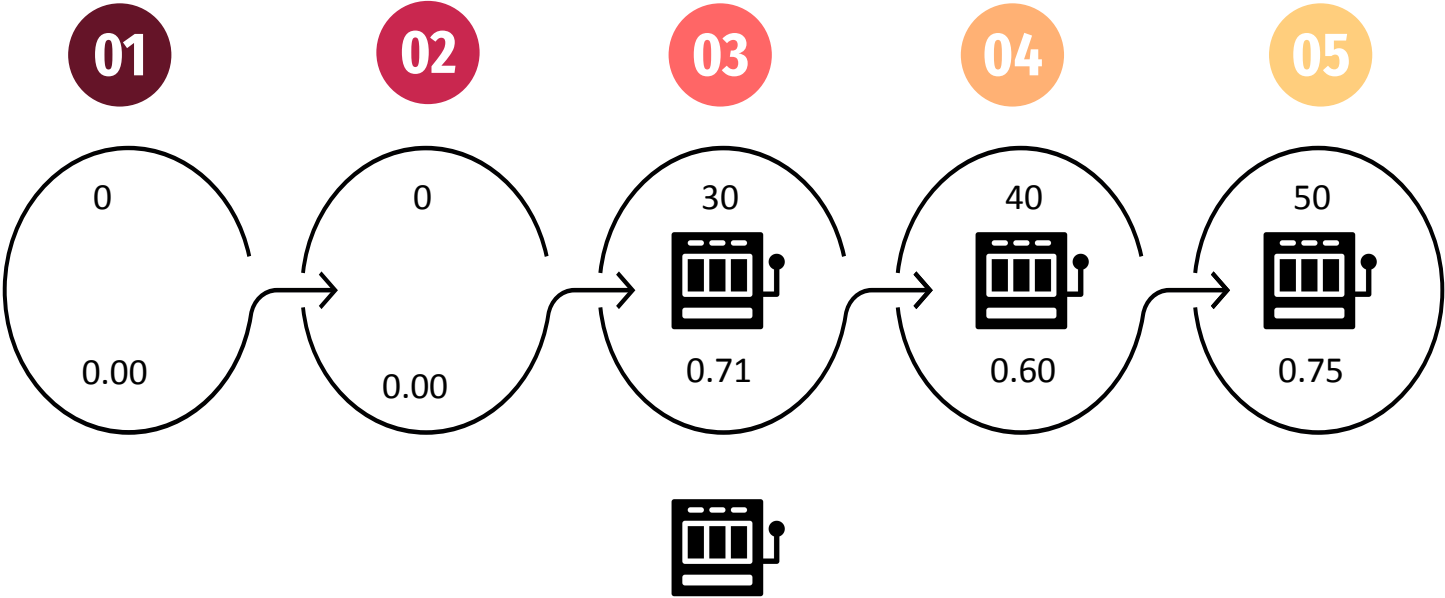
# Best-Arm Identification Algorithm



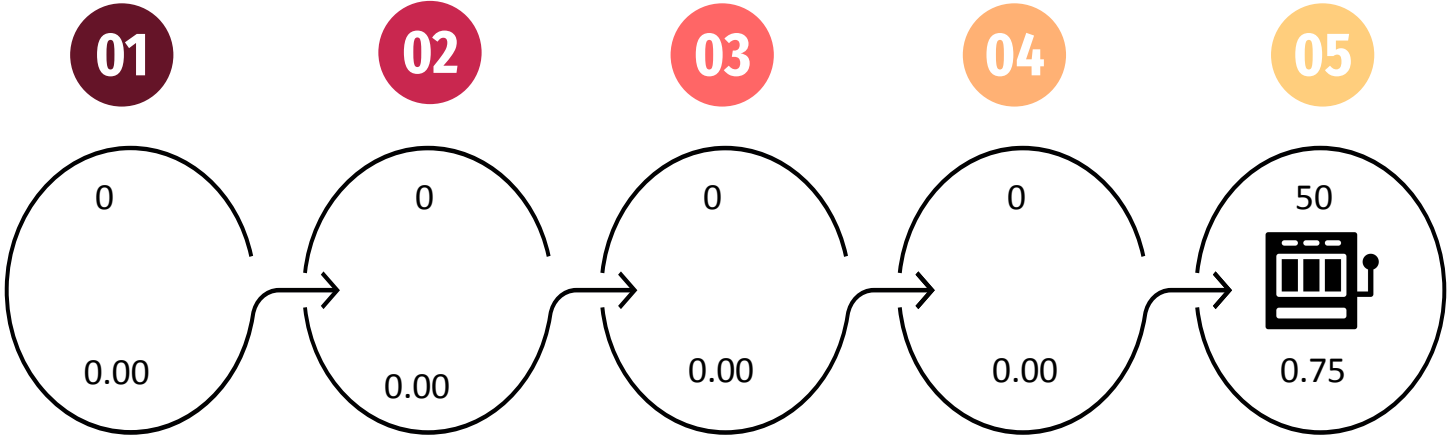
# Best-Arm Identification Algorithm



# Best-Arm Identification Algorithm



# Best-Arm Identification Algorithm





# Summary

**Main Result 1:** Any single-pass algorithm will incur at least

$$\Omega\left(\frac{n^{1/3}T^{2/3}}{m^{7/3}}\right)$$

expected cumulative regret when  $m < n$ .

This lower bound even holds for random order arrival.

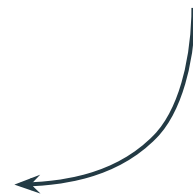
**Main Result 2:**  $(\varepsilon, \delta)$  - PAC algorithm with

$O(r)$  arm memory

Optimal  $\mathcal{R}$ -round sample complexity

**Corollary:**  $(\varepsilon, \delta)$ - PAC algorithm with  $O(\log^* n)$  arm memory with optimal

worst-case sample complexity

$$O\left(\frac{n}{\varepsilon^2} (\mathbf{i}\log^{(r)}(n) + \log\left(\frac{1}{\delta}\right))\right)$$


# Open Problems

- 1) Obtain instance-dependent lower bounds and upper bounds on the expected cumulative regret for single-pass MAB algorithms with bounded arm-memory.
- 2) Obtain lower bounds and upper bounds on the expected cumulative regret for  $k$ -pass MAB algorithms with bounded arm-memory where  $k > 1$ .
- 3) Obtain an  $(\epsilon, \delta)$ -PAC streaming algorithm with  $O(1)$  arm memory and optimal worst case sample complexity.

# Open Problems

- 1) Obtain instance-dependent lower bounds and upper bounds on the expected cumulative regret for single-pass MAB algorithms with bounded arm-memory.
- 2) Obtain lower bounds and upper bounds on the expected cumulative regret for  $k$ -pass MAB algorithms with bounded arm-memory where  $k > 1$ .
- 3) Obtain an  $(\epsilon, \delta)$ -PAC streaming algorithm with  $O(1)$  arm memory and optimal worst case sample complexity.

# Thank You