

13^{ème} Atelier en Évaluation des Performances

Deep Q-Learning with Whittle Index for Contextual Restless Multi-Armed Bandits

Presented by

PhD student Ibtihal EL MIMOUNI

PhD supervised by

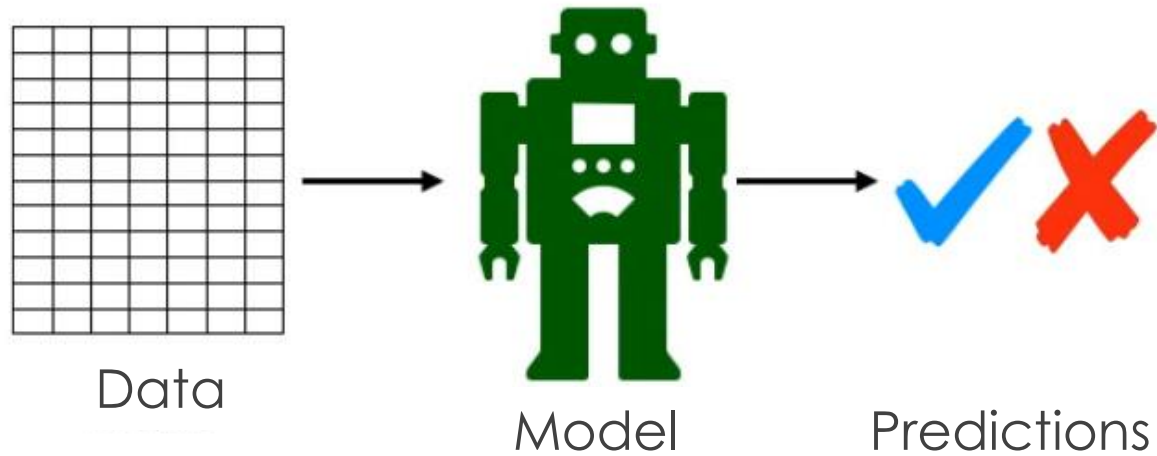
Dr Konstantin AVRACHENKOV & Mr Hervé BAILE

Tuesday, December 3rd 2024

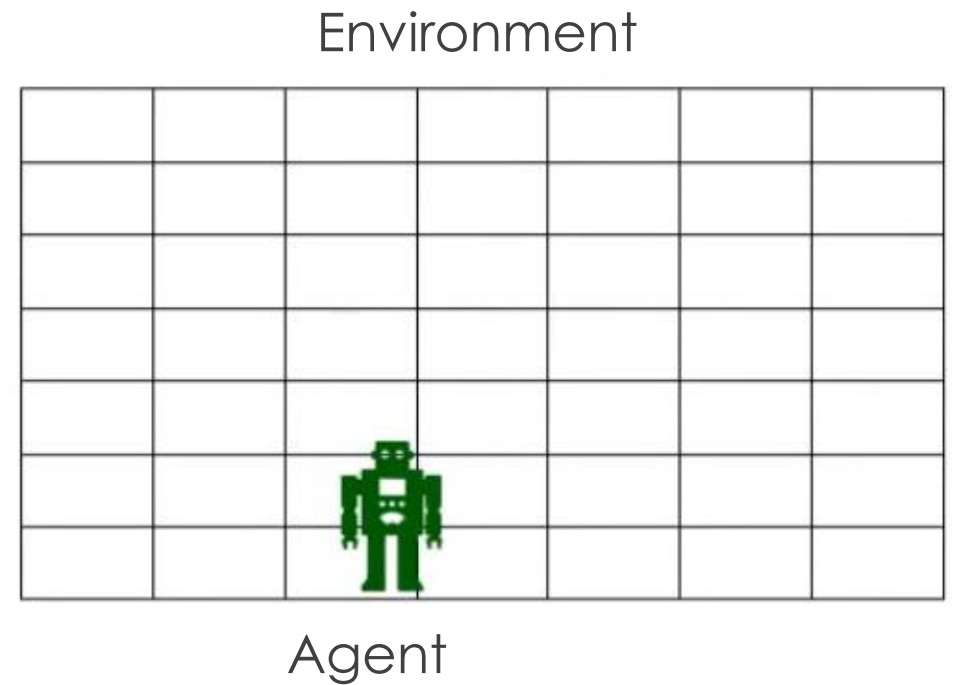
INTRODUCTION

MACHINE LEARNING

Predictive Machine Learning (ML)

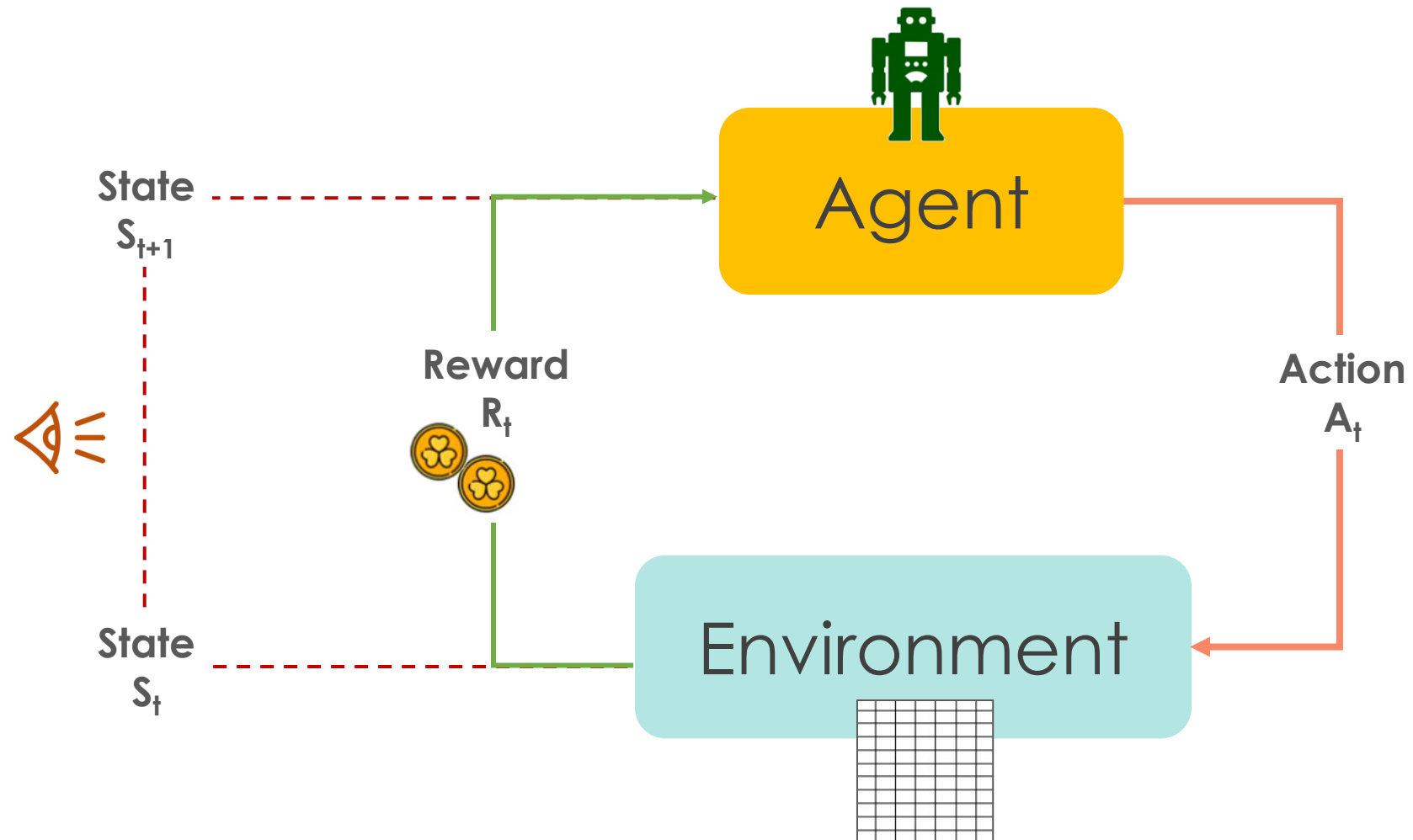


Reinforcement Learning (RL)



INTRODUCTION

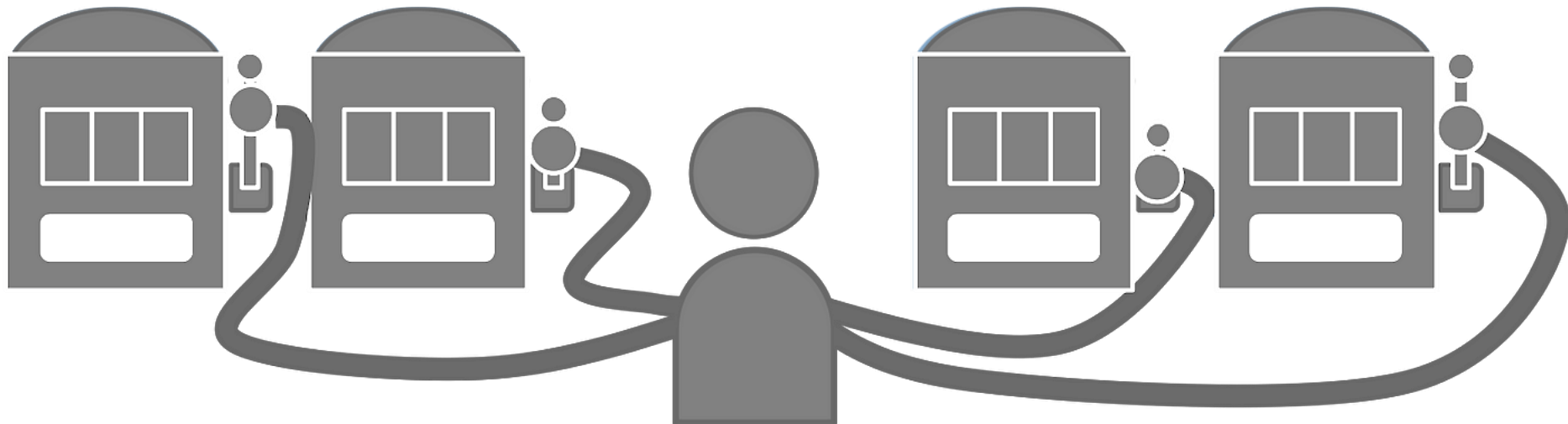
REINFORCEMENT LEARNING



INTRODUCTION

MULTI-ARMED BANDIT (MAB)

- **Multi-Armed:** Multiple options that an algorithm can choose from. Each option is like a **lever** or **arm** that can be pulled.
- **Bandit:** Analogy of casinos where there are multiple slot machines and the gambler has to decide which machine to play.

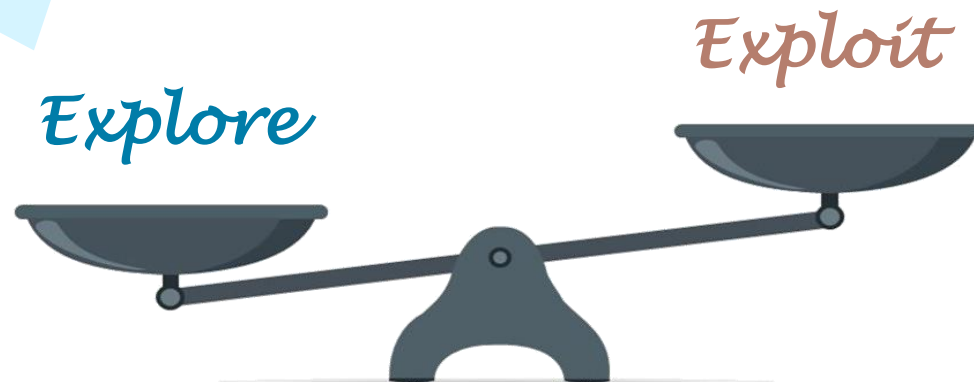


MULTI-ARMED BANDIT

EXPLORATION-EXPLOITATION TRADE-OFF

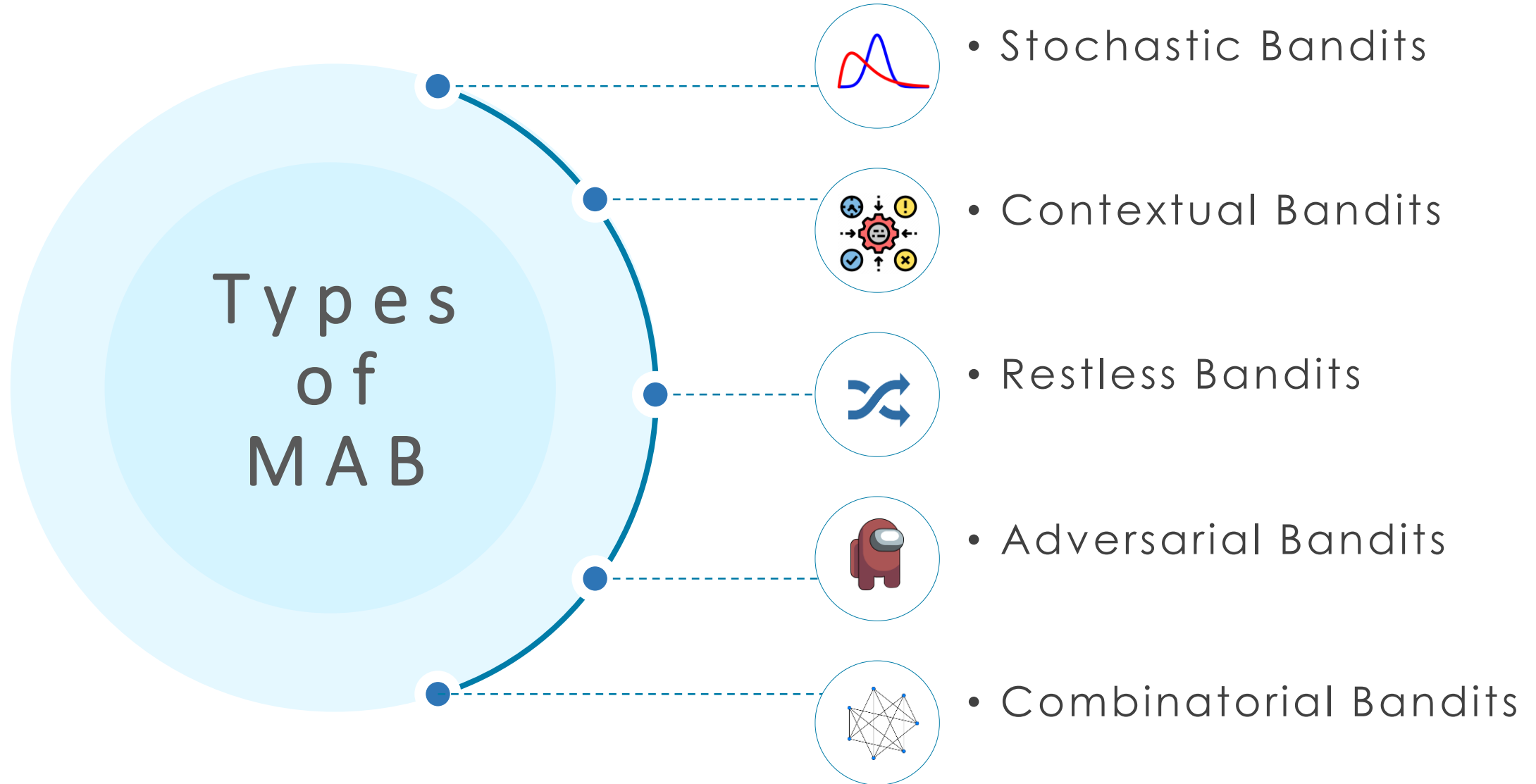
The agent **explores** different arms to see which ones yield the best reward.

The agent **exploits** the arms that have yielded the best rewards.



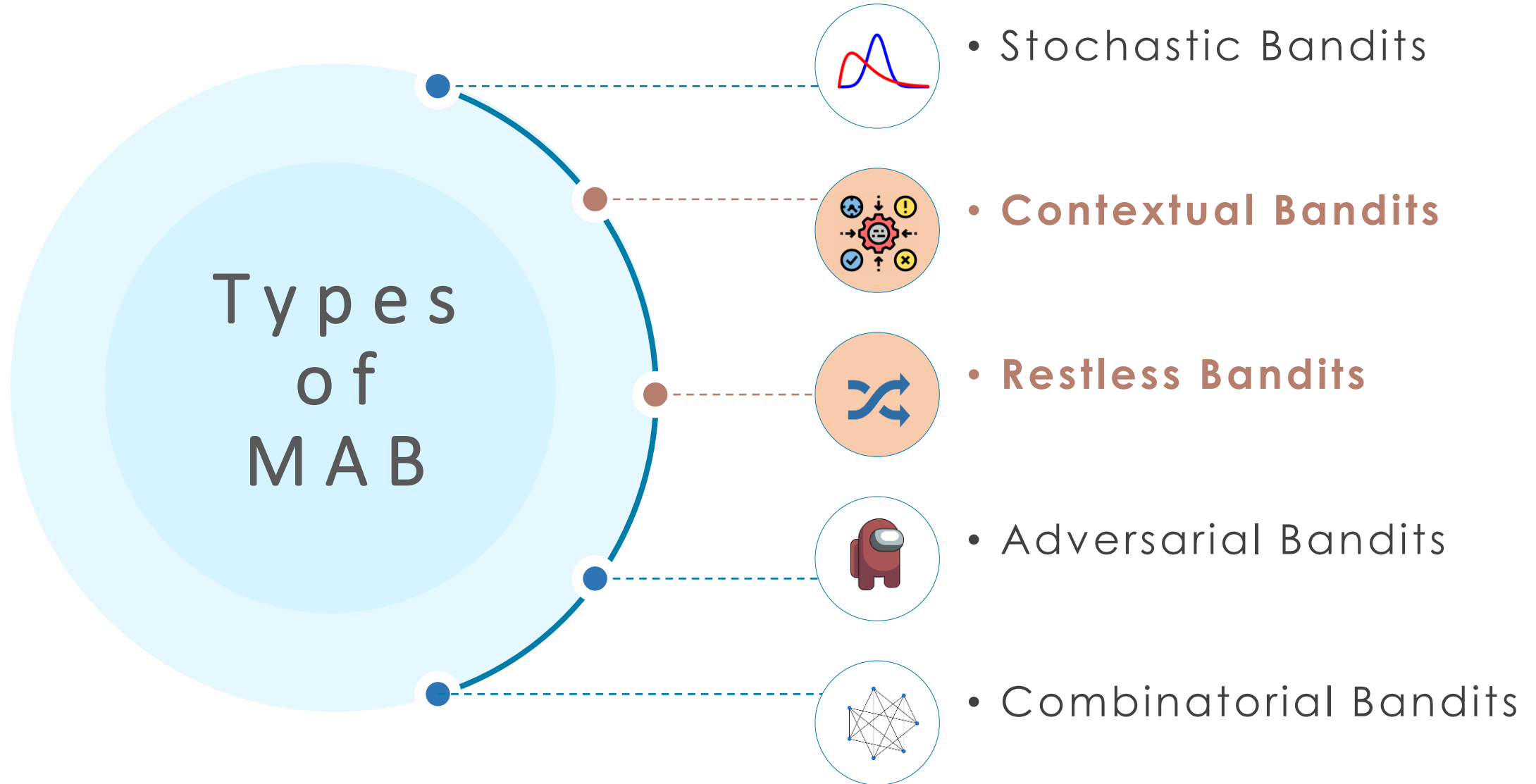
MULTI-ARMED BANDIT

OVERVIEW



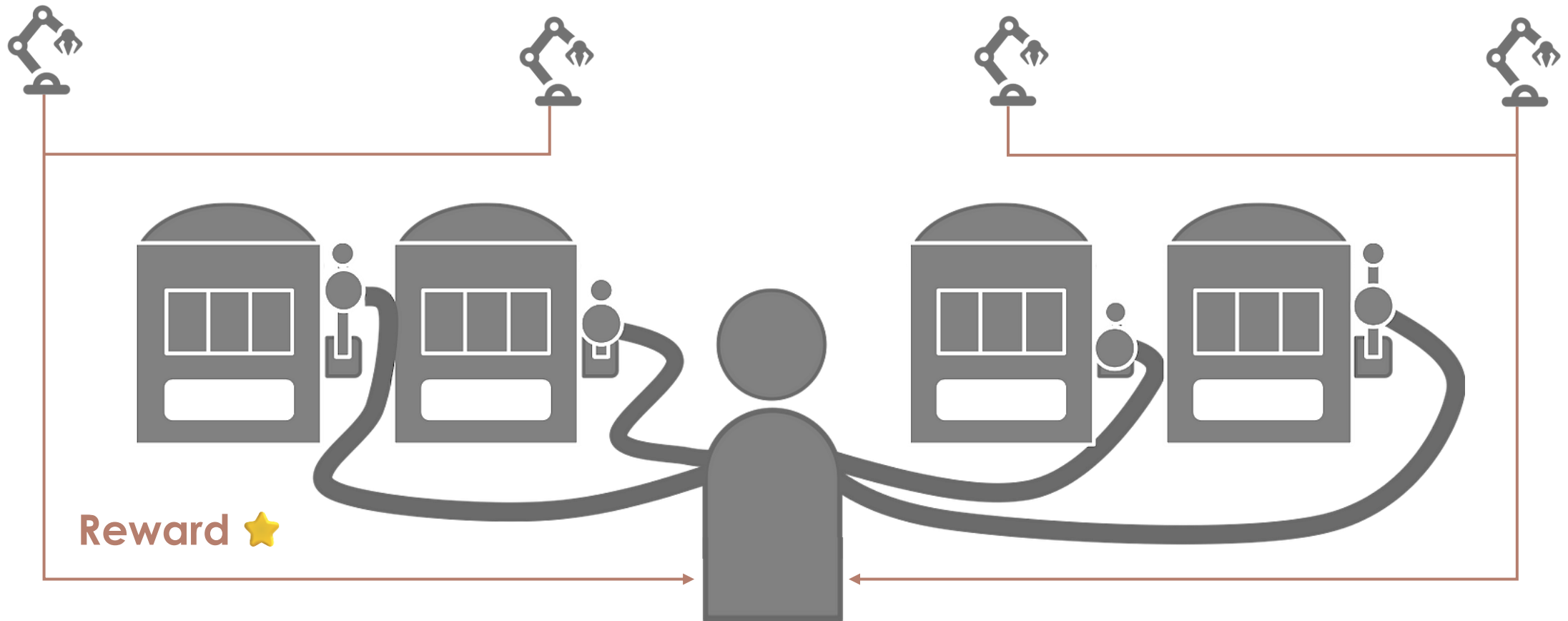
MULTI-ARMED BANDIT

OVERVIEW



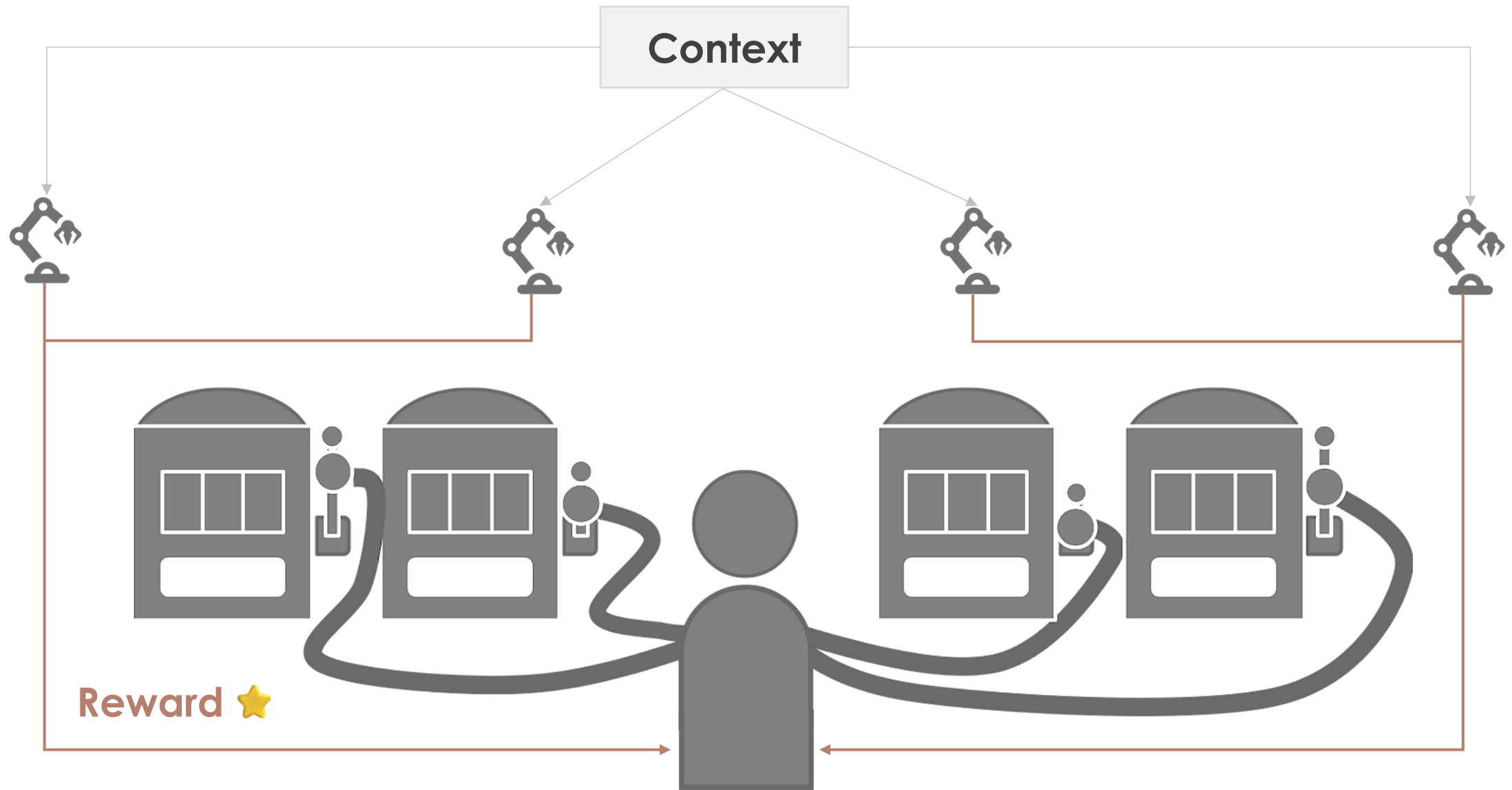
MULTI-ARMED BANDITS

OVERVIEW



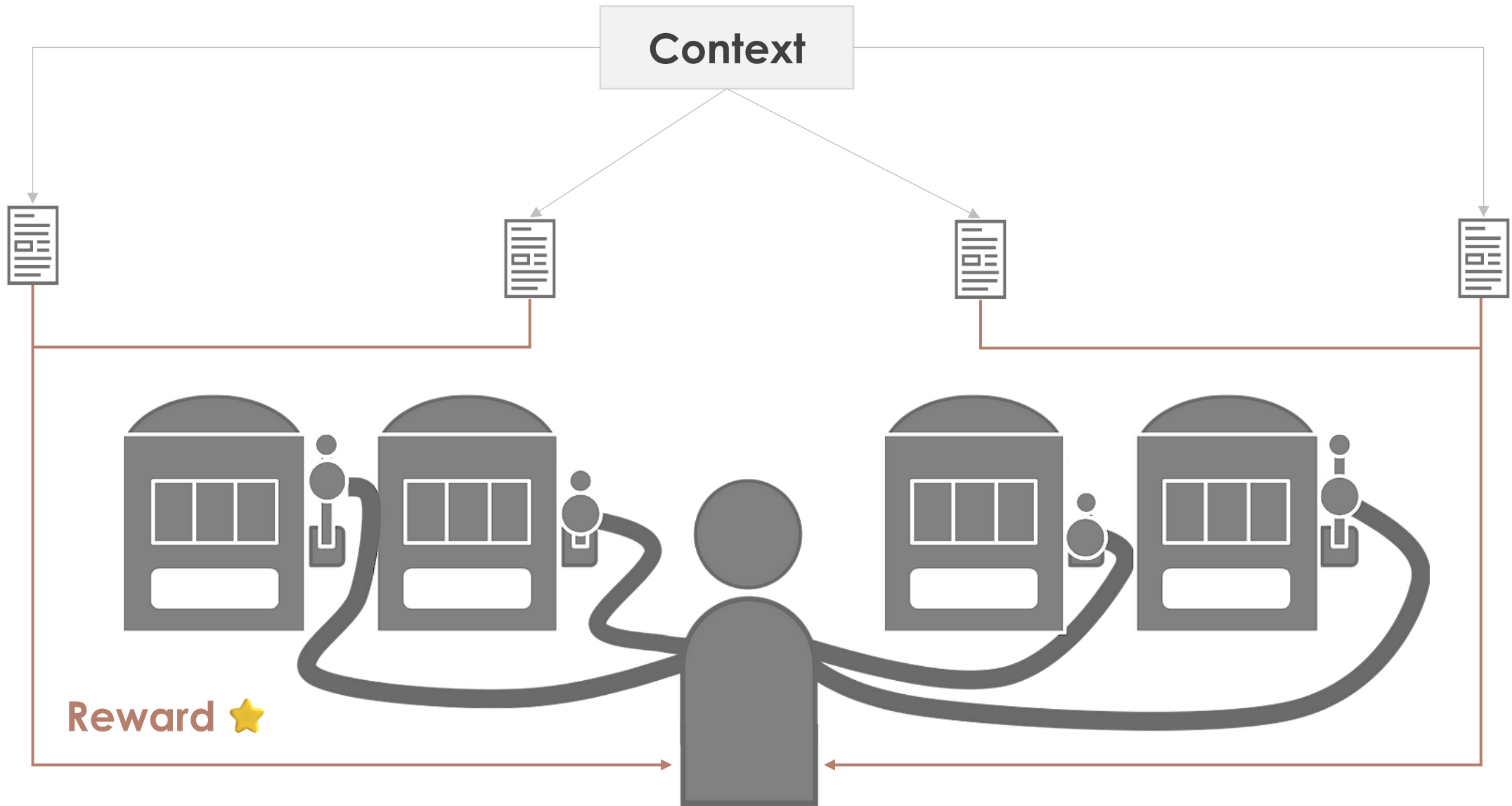
CONTEXTUAL MULTI-ARMED BANDITS

ADDING CONTEXT



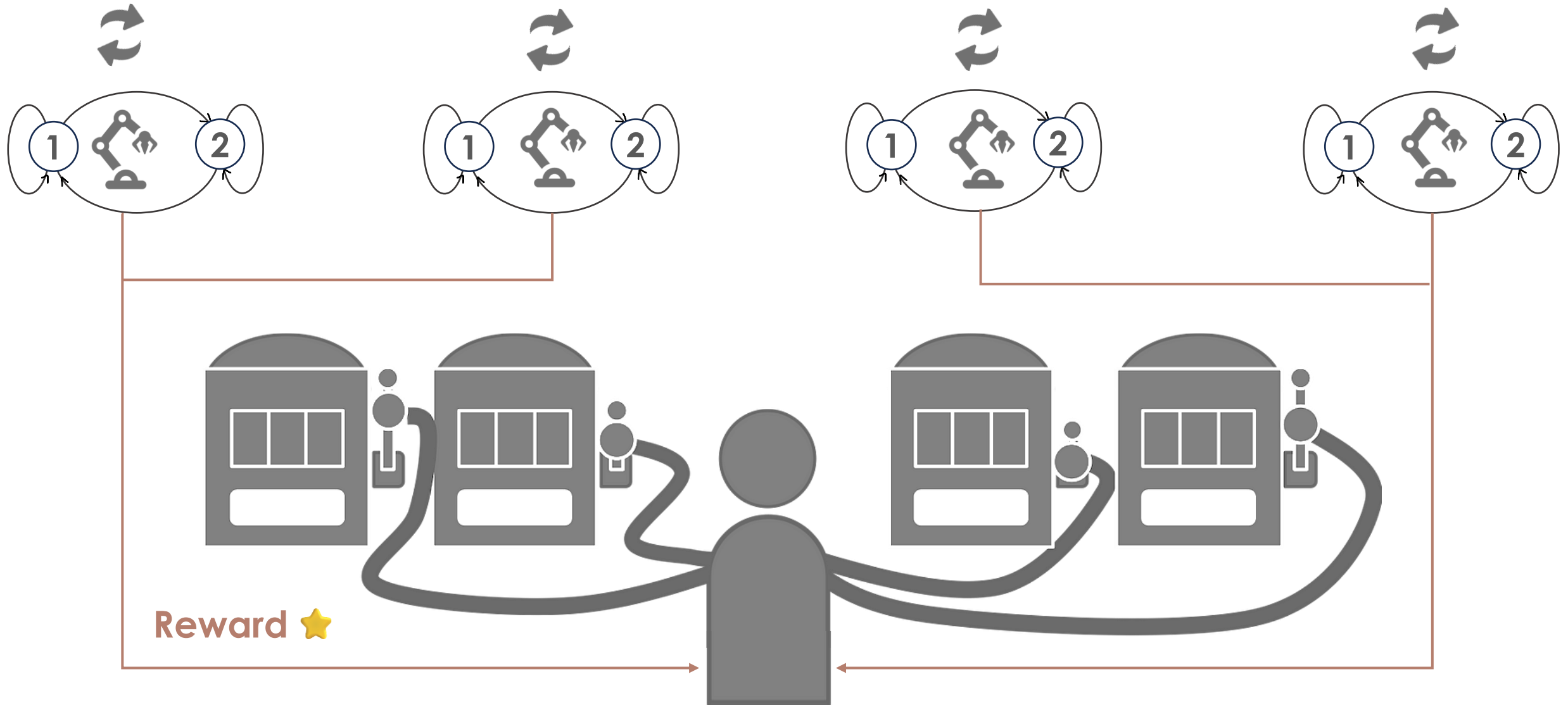
CONTEXTUAL MULTI-ARMED BANDITS

EXAMPLE: LinUCB



RESTLESS MULTI-ARMED BANDITS

RESTLESS ARMS



RESTLESS MULTI-ARMED BANDITS

MARKOVIAN BANDITS

- Restless Multi-Armed Bandits (or Restless Bandits) evolve independently in a controlled Markov manner, i.e. we have:

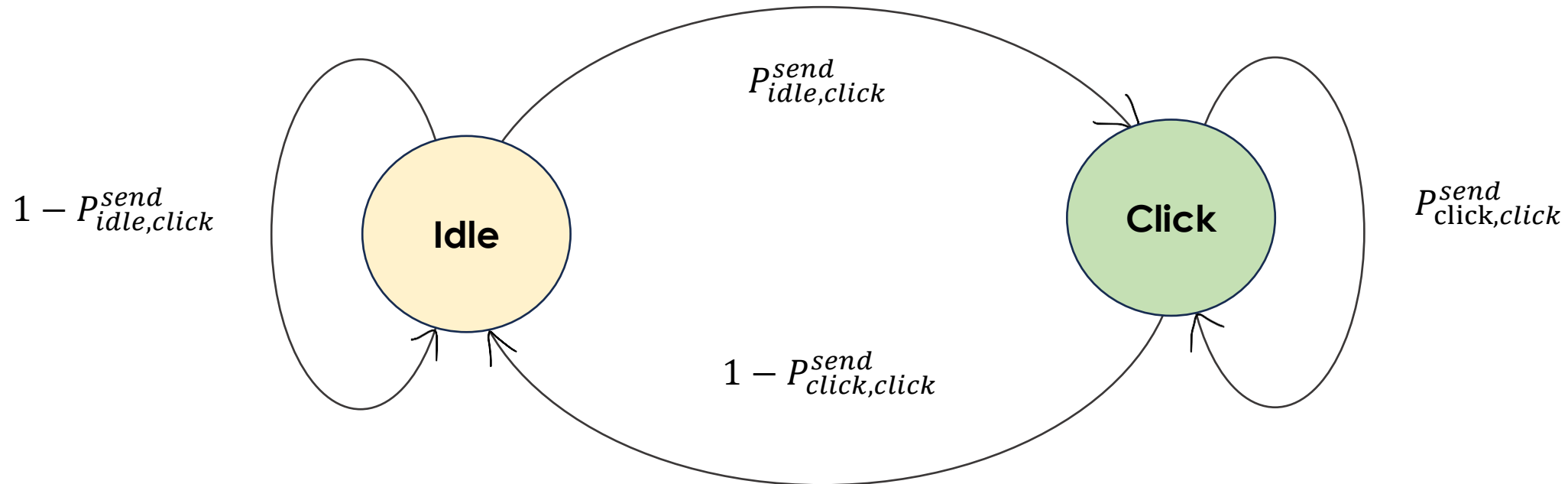
$$\begin{aligned} P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) &:= \mathbb{P}(\mathbf{S}_{t+1} = \mathbf{s}_{t+1} | \mathbf{S}_{0:t} = \mathbf{s}_{0:t}, \mathbf{A}_{0:t} = \mathbf{a}_{0:t}) \\ &= \prod_{i=1}^n \mathbb{P}(S_{t+1}^i = s_{t+1}^i | S_t^i = s_t^i, A_t^i = a_t^i) \end{aligned}$$

- $S_t := (S_t^1, \dots, S_t^n)$ and $A_t = (A_t^1, \dots, A_t^n)$ denote the joint state and actions at time t .

RESTLESS MULTI-ARMED BANDITS

A MARKOV DECISION PROCESS

- The state of **each arm** evolves according to a Markovian transition, even when the arm is not played.



RESTLESS BANDITS

FORMULATION

- A Restless Multi-Armed Bandit (RMAB) is defined by the controlled Markov processes $\mathbf{MDP} = (\mathbf{S}, \mathbf{A}, \mathbf{R}, \mathbf{P})$, such that:
 - $\mathbf{S} = \{1, 2, \dots, |\mathbf{S}|\}$ is the state space, where s_i is the state of arm i .
 - $\mathbf{A} = \{0, 1\}$ denotes the actions space with two actions: $a_i = 1$ to activate arm i , and $a_i = 0$ to keep arm i passive.
 - $r_i^t = R(s_i^t, a_i^t, s_i^t)$ is the reward function.
 - $P(s_i^{t+1} | s_i^t, a_i^t)$ denotes the transition probability from state s_i^t to state s_i^{t+1} .

RESTLESS BANDITS

OBJECTIVE

- Find a policy $\pi = (\pi_1, \dots, \pi_N)$ which maximizes the expected discounted reward i.e. π such that:

$$\text{maximize } \mathbb{E}_{s_0^{(i)} \sim \mu^{(i)}, a_t^{(i)} \sim \pi^{(i)}(s_t^{(i)}), s_{t+1} \sim p^{(i)}(\cdot | s_t^{(i)}, a_t^{(i)})} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t r^{(i)}(s_t^{(i)}, a_t^{(i)}) \right]$$

subject to the following hard constraint $\forall t$:

$$\sum_{i=1}^N a_t^{(i)} = M$$

RESTLESS BANDITS

OBJECTIVE

- Find a policy $\pi = (\pi_1, \dots, \pi_N)$ which maximizes the expected discounted reward i.e. π such that:

$$\text{maximize } \mathbb{E}_{s_0^{(i)} \sim \mu^{(i)}, a_t^{(i)} \sim \pi^{(i)}(s_t^{(i)}), s_{t+1} \sim p^{(i)}(\cdot | s_t^{(i)}, a_t^{(i)})} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t r^{(i)}(s_t^{(i)}, a_t^{(i)}) \right]$$

subject to the following relaxed constraint :

$$\mathbb{E}_{a_t^{(i)} \sim \pi^{(i)}(s_t^{(i)})} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t a_t^{(i)} \right] = \frac{M}{1 - \gamma}$$

WHITTLE INDEX

THREE OPTIMIZATION PROBLEMS

- **[Original]** Original problem :

$$\begin{aligned} & \text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t r^{(i)}(s_t^{(i)}, a_t^{(i)}) \right] \\ & \text{s.t. } \sum_{i=1}^N a_t^{(i)} = M, \quad \forall t \end{aligned}$$

- **[Relaxed]** Problem with Relaxed activation constraint : $\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t a_t^{(i)} = \frac{M}{1-\gamma}$

- **[Lagrange]** Introducing Lagrange Multiplier :

$$\text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t \left(r^{(i)}(s_t^{(i)}, a_t^{(i)}) + \tilde{\lambda} (1 - a_t^{(i)}) \right) \right]$$

WHITTLE INDEX

THREE OPTIMIZATION PROBLEMS

- **[Original]** Original problem :

$$\begin{aligned} & \text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t r^{(i)}(s_t^{(i)}, a_t^{(i)}) \right] \\ & \text{s.t. } \sum_{i=1}^N a_t^{(i)} = M, \quad \forall t \end{aligned}$$

- **[Relaxed]** Problem with Relaxed activation constraint : $\tilde{\lambda} \times \sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t a_t^{(i)} = \frac{M}{1-\gamma} \times \tilde{\lambda}$

- **[Lagrange]** Introducing Lagrange Multiplier :

$$\text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t \left(r^{(i)}(s_t^{(i)}, a_t^{(i)}) + \tilde{\lambda} (1 - a_t^{(i)}) \right) \right]$$

WHITTLE INDEX

DECOUPLED PROBLEM

- **[Lagrange]** The function is given by :

$$\text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t \left(r^{(i)}(s_t^{(i)}, a_t^{(i)}) - \tilde{\lambda} a_t^{(i)} \right) \right] + \tilde{\lambda}(M/(1 - \gamma)) \quad \text{s.t. } a_t^{(i)} \in \{0, 1\}, \forall t$$

- **[Decoupled Problem]** Solving the above for $\tilde{\lambda}$ is equivalent to solving the following for each $i \in [N]$. We can decouple this problem and neglect the last term.

$$\text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \left(r^{(i)}(s_t^{(i)}, a_t^{(i)}) - \tilde{\lambda} a_t^{(i)} \right) \right] \quad \text{s.t. } a_t^{(i)} \in \{0, 1\}, \forall t$$

WHITTLE INDEX

DECOUPLED PROBLEM

- **[Lagrange]** The Lagrange dual function is given by :

$$\text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t \left(r^{(i)}(s_t^{(i)}, a_t^{(i)}) - \tilde{\lambda} a_t^{(i)} \right) \right] + \tilde{\lambda} (M / (1 - \gamma)) \quad \text{s.t. } a_t^{(i)} \in \{0, 1\}, \forall t$$

cte

- **[Decoupled Problem]** Solving the above for $\tilde{\lambda}$ is equivalent to solving the following for each $i \in [N]$. We can decouple this problem and neglect the last term.

$$\text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \left(r^{(i)}(s_t^{(i)}, a_t^{(i)}) - \tilde{\lambda} a_t^{(i)} \right) \right] \quad \text{s.t. } a_t^{(i)} \in \{0, 1\}, \forall t$$

WHITTLE INDEX

THREE OPTIMIZATION PROBLEMS

- **[Original]** Original problem :

$$\begin{aligned} & \text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t r^{(i)}(s_t^{(i)}, a_t^{(i)}) \right] \\ & \text{s.t. } \sum_{i=1}^N a_t^{(i)} = M, \quad \forall t \end{aligned}$$

- **[Relaxed]** Problem with Relaxed activation constraint : $\tilde{\lambda} \times \sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t a_t^{(i)} = \frac{M}{1-\gamma} \times \tilde{\lambda}$

- **[Lagrange]** Introducing Lagrange Multiplier :

$$\text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t \left(r^{(i)}(s_t^{(i)}, a_t^{(i)}) + \tilde{\lambda} (1 - a_t^{(i)}) \right) \right]$$

WHITTLE INDEX

SOLVING THE DECOUPLED PROBLEM

$$\mathbb{E} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t \left(r^{(i)} \left(s_t^{(i)}, a_t^{(i)} \right) + \tilde{\lambda} \left(1 - a_t^{(i)} \right) \right) \right]$$

- The solution of the equation above is obtained by combining the solution to N independent problems.

$$V^i(s) = \max_{a \in \{0,1\}} \left[a \left(r^i(s, 1) + \gamma \cdot \sum_j p^i(j|s, 1) V^i(j) \right) + (1 - a) \left(r^i(s, 0) + \lambda + \gamma \cdot \sum_j p(j, |s, 0) V^i(j) \right) \right]$$

- We can rewrite the expression of the equation above as a function of the state-action $Q^i(s, a)$.

$$Q^i(s, a) = a \left(r^i(s, a) + \gamma \cdot \sum_j p^i(j|s, a) V^i(j) \right) + (1 - a) \left(r^i(s, a) + \lambda + \gamma \cdot \sum_j p(j|s, a) V^i(j) \right)$$

SOLVING THE DECOUPLED PROBLEM

BELLMAN OPTIMALITY EQUATION

- We have to find Q for each $i \in [N]$ such that the above is satisfied for some $\tilde{\lambda}$.

[Bellman Optimality Equation]

$$Q^*(s, a) = (1 - a)(r(s, a) + \tilde{\lambda}) + r(s, a) + \sum_{s' \in \mathcal{S}} p(s'|s, a) \max_{a' \in \mathcal{A}} Q^*(s', a') \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}$$

[Q-learning]

$$Q_{n+1}(s, a) = Q_n(s, a) + \alpha \mathbb{I}\{s_n = s, a_n = a\} \left((1 - a)(r(s, a) + \tilde{\lambda}) + r(s, a) + \max_{a' \in \mathcal{A}} Q_n(s', a') \right)$$

SOLVING THE DECOUPLED PROBLEM

MINIMUM SUBSIDY

- Rewriting the Bellman optimality equation for active actions 1 and 0 gives:

$$Q^*(s, 1) = r(s, 1) + \sum_{s' \in \mathcal{S}} p(s'|s, a) \max_{a' \in \mathcal{A}} Q^*(s', a')$$

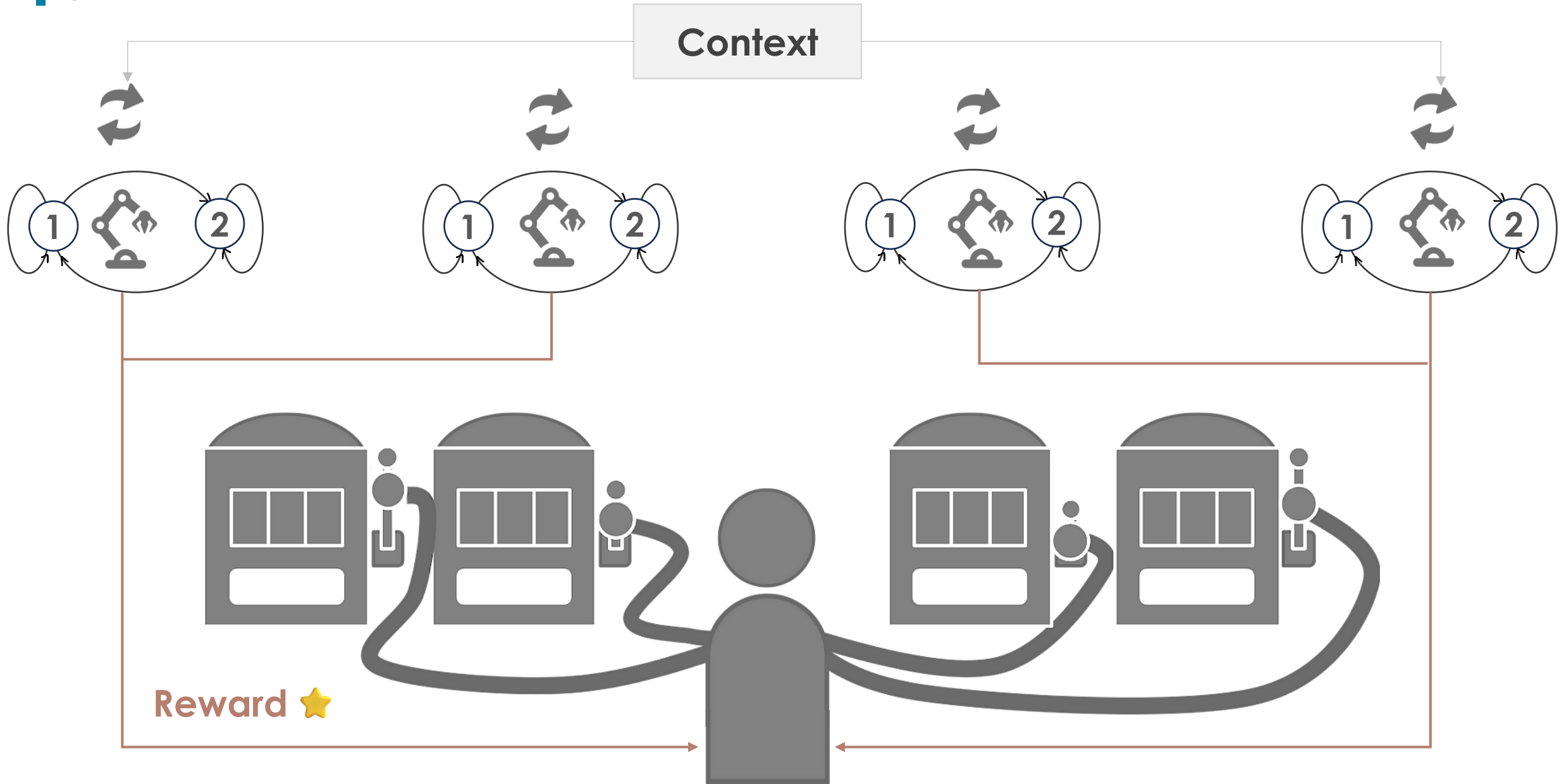
$$Q^*(s, 0) = r(s, 0) + \tilde{\lambda} + \sum_{s' \in \mathcal{S}} p(s'|s, a) \max_{a' \in \mathcal{A}} Q^*(s', a')$$

- Whittle defines the Whittle Index for state \hat{k} to be the value $\lambda(\hat{k})$ of $\tilde{\lambda}$ for which both active and passive actions are equally preferred in state \hat{k} i.e:

$$Q(\hat{k}, 1) = Q(\hat{k}, 0) \quad \text{for } \tilde{\lambda} = \lambda(\hat{k})$$

CONTEXTUAL RESTLESS BANDITS

OVERVIEW



RESTLESS BANDITS

RECALL

- A Restless Multi-Armed Bandit (RMAB) is defined by the controlled Markov processes $\mathbf{MDP} = (\mathbf{S}, \mathbf{A}, \mathbf{R}, \mathbf{P})$, such that:
 - $\mathbf{S} = \{1, 2, \dots, |\mathbf{S}|\}$ is the state space, where s_i is the state of arm i .
 - $\mathbf{A} = \{0, 1\}$ denotes the actions space with two actions: $a_i = 1$ to activate arm i , and $a_i = 0$ to keep arm i passive.
 - $r_i^t = R(s_i^t, a_i^t, s_i^t)$ is the reward function.
 - $P(s_i^{t+1} | s_i^t, a_i^t)$ denotes the transition probability from state s_i^t to state s_i^{t+1} .

CONTEXTUAL RESTLESS BANDITS

FORMULATION

- A **Contextual** Restless Multi-Armed Bandit (**CRMAB**) is defined by the controlled Markov processes $\mathbf{MDP}^c = (\mathbf{S}, \mathbf{C}, \mathbf{A}, \mathbf{R}, \mathbf{P})$, such that:
 - $\mathbf{S} = \{1, 2, \dots, |\mathbf{S}|\}$ is the state space, where s_i is the state of arm i .
 - $\mathbf{C} = \{1, 2, \dots, |\mathbf{C}|\}$, is the context space, where c_i is the context of arm i .
 - $\mathbf{A} = \{0, 1\}$ denotes the actions space with two actions: $a_i = 1$ to activate arm i , and $a_i = 0$ to keep arm i passive.
 - $r_i^t = R(s_i^t, c_i^t, a_i^t)$ is the reward function.
 - $P(s_i^{t+1} | s_i^t, a_i^t, c_i^t)$ denotes the transition probabilities.

CONTEXTUAL RESTLESS BANDITS

FORMULATION

- A Contextual Restless Multi-Armed Bandit (CRMAB) is defined by the controlled Markov processes $\mathbf{MDP}^c = (\mathbf{S}, \mathbf{C}, \mathbf{A}, \mathbf{R}, \mathbf{P})$, such that:
 - $\mathbf{S} = \{1, 2, \dots, |\mathbf{S}|\}$ is the state space, where s_i is the state of arm i .
 - $\mathbf{C} = \{1, 2, \dots, |\mathbf{C}|\}$, is the context space, where c_i is the context of arm i .
 - $\mathbf{A} = \{0, 1\}$ denotes the actions space with two actions: $a_i = 1$ to activate arm i , and $a_i = 0$ to keep arm i passive.
 - $r_i^t = R(s_i^t, c_i^t, a_i^t)$ is the reward function.
 - $P(s_i^{t+1} | s_i^t, a_i^t, c_i^t)$ denotes the transition probabilities

CONTEXTUAL RESTLESS BANDITS

OBJECTIVE

- Find a policy $\pi = (\pi_1, \dots, \pi_N)$ which maximizes the expected discounted reward i.e. π such that:

$$\text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} \sum_{i=1}^N \gamma^t r^{(i)}(s_t^{(i)}, c_t^{(i)}, a_t^{(i)}) \right]$$

- Emerging topic, only two works on CRMAB:

[1] X. Chen, I. Hou. “Contextual Restless Multi-Armed Bandits with Application to Demand Response Decision-Making”. In: arXiv preprint arXiv:2403.15640 **(2024)**

[2] B. Liang, L. Xu, A. Taneja, M. Tambe, and L. Janson. “A Bayesian Approach to Online Learning for Contextual Restless Bandits with Applications to Public Health”. In: arXiv preprint arXiv:2402.04933 **(2024)**

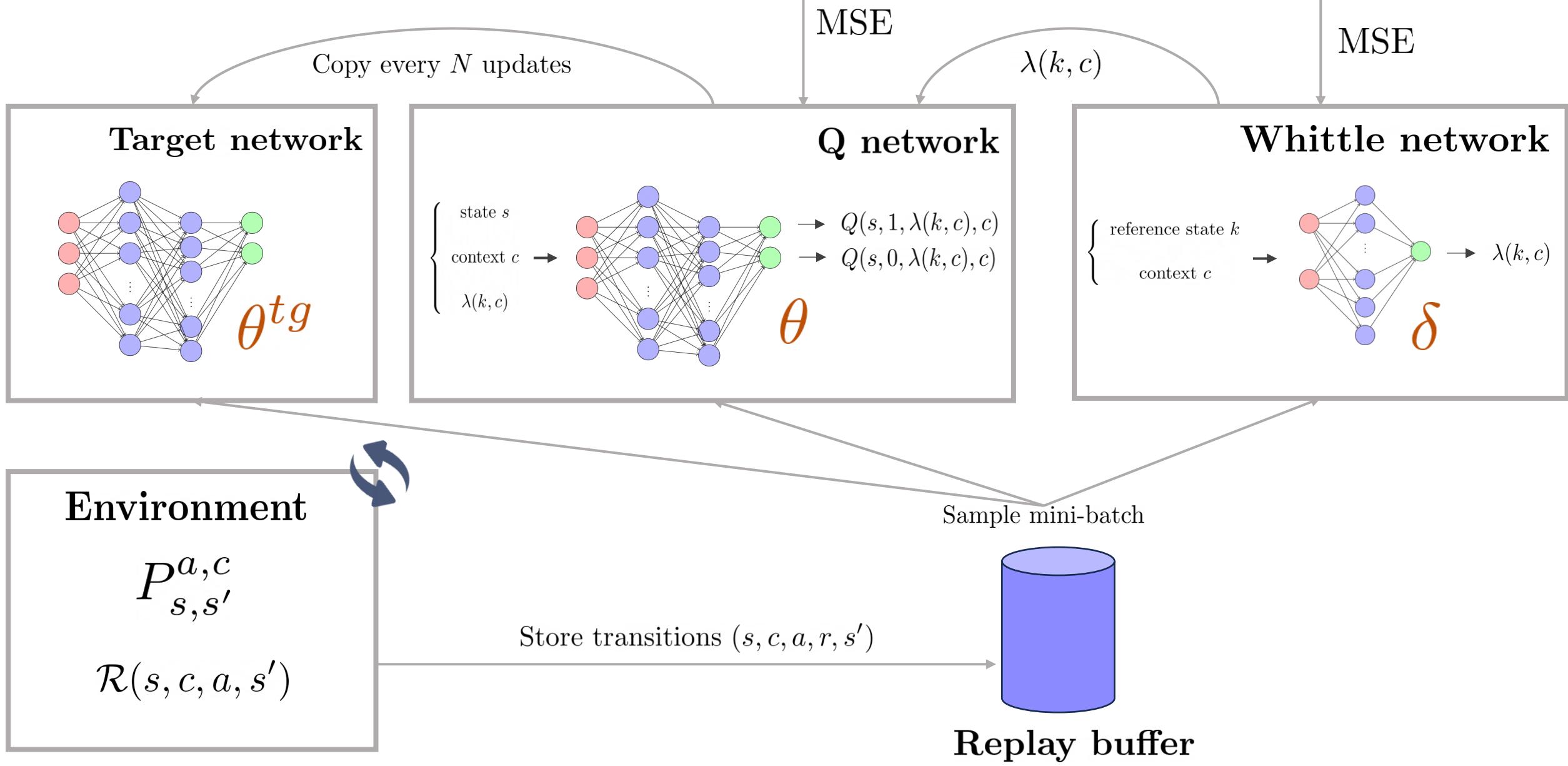
ALGORITHM

OVERVIEW

- Deep Q-learning the Whittle Index with Context (DQWIC).
- Combines **deep RL** and **Whittle index** theory within the CRMAB framework.
- Leverages **two Neural Networks (NN)**.
 - Q-network for approximating action-value functions,
 - Whittle-network for estimating Whittle indices.
- The learning process occurs through a **two time scale stochastic approximation**.

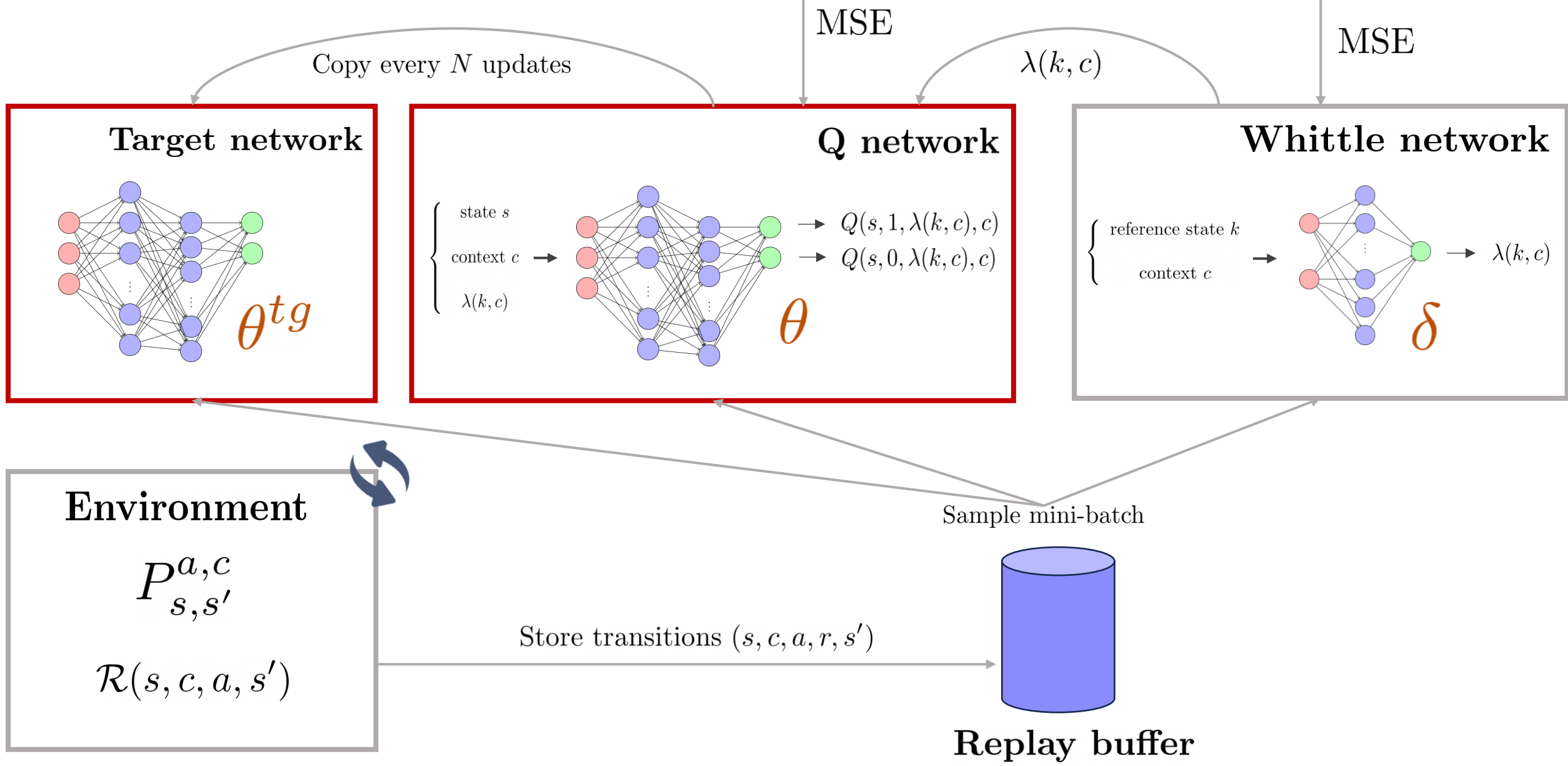
$$\mathcal{L}(\theta) := \mathbb{E} \left[\left\| Q_{\theta}(s, a, \lambda_{\delta}, c) - Q^{\text{target}}(s, a, \lambda_{\delta}, c) \right\|^2 \right]$$

$$\mathcal{E}(\delta) = \mathbb{E} \left[\left\| Q_{\theta}(k, 1, \lambda_{\delta}, c) - Q_{\theta}(k, 0, \lambda_{\delta}, c) \right\|^2 \right]$$



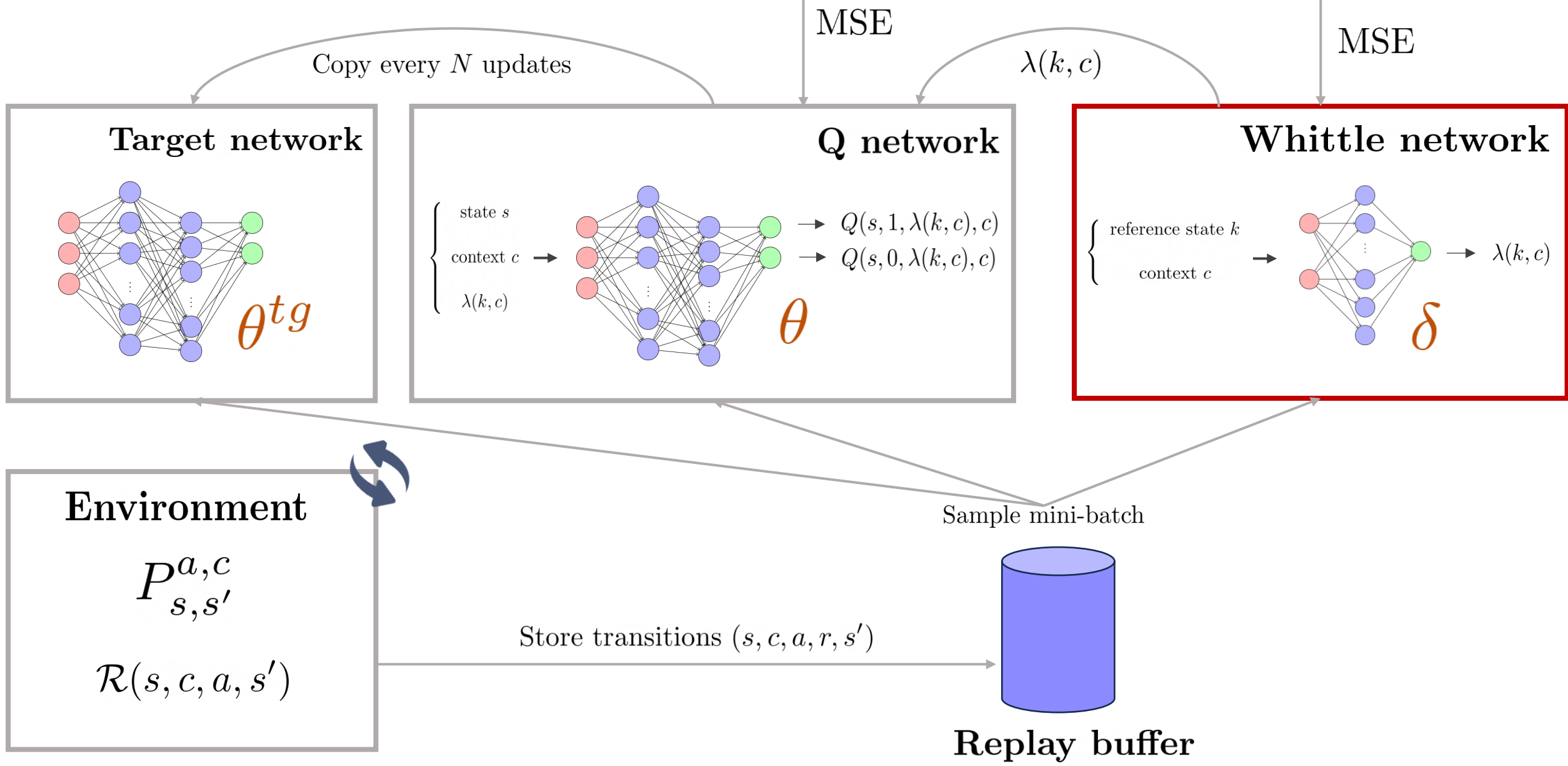
$$\mathcal{L}(\theta) := \mathbb{E} \left[\left\| Q_{\theta}(s, a, \lambda_{\delta}, c) - Q^{\text{target}}(s, a, \lambda_{\delta}, c) \right\|^2 \right]$$

$$\mathcal{E}(\delta) = \mathbb{E} \left[\left\| Q_{\theta}(k, 1, \lambda_{\delta}, c) - Q_{\theta}(k, 0, \lambda_{\delta}, c) \right\|^2 \right]$$



$$\mathcal{L}(\theta) := \mathbb{E} \left[\left\| Q_{\theta}(s, a, \lambda_{\delta}, c) - Q^{\text{target}}(s, a, \lambda_{\delta}, c) \right\|^2 \right]$$

$$\mathcal{E}(\delta) = \mathbb{E} \left[\left\| Q_{\theta}(k, 1, \lambda_{\delta}, c) - Q_{\theta}(k, 0, \lambda_{\delta}, c) \right\|^2 \right]$$



$$\mathcal{L}(\theta) := \mathbb{E} \left[\left\| Q_{\theta}(s, a, \lambda_{\delta}, c) - Q^{\text{target}}(s, a, \lambda_{\delta}, c) \right\|^2 \right]$$

$$\mathcal{E}(\delta) = \mathbb{E} \left[\left\| Q_{\theta}(k, 1, \lambda_{\delta}, c) - Q_{\theta}(k, 0, \lambda_{\delta}, c) \right\|^2 \right]$$

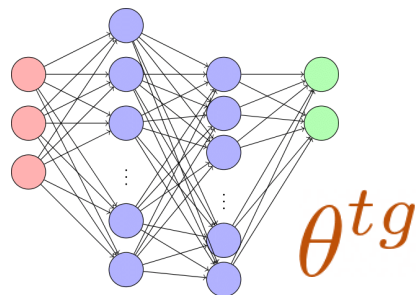
MSE

MSE

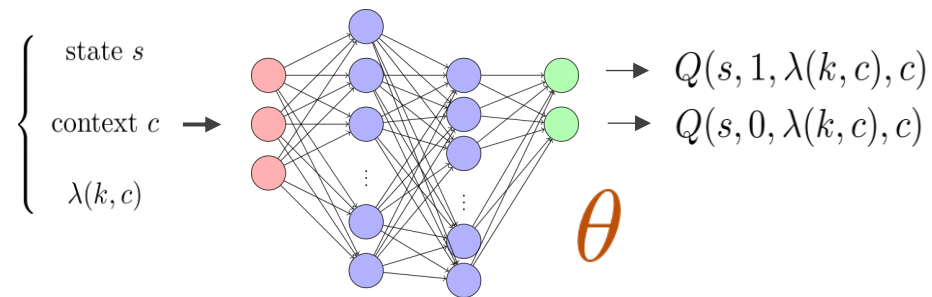
Copy every N updates

$\lambda(k, c)$

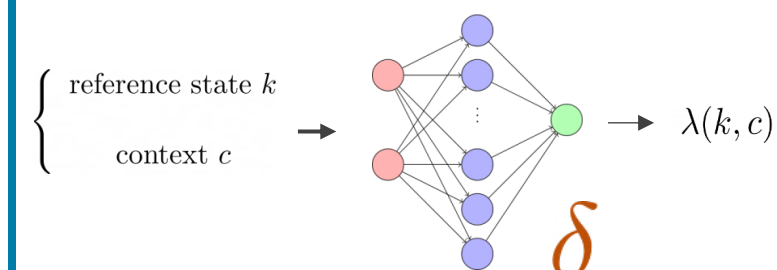
Target network



Q network



Whittle network



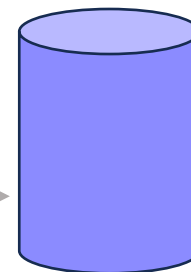
Environment

$$P_{s, s'}^{a, c}$$

$$\mathcal{R}(s, c, a, s')$$

Sample mini-batch

Store transitions (s, c, a, r, s')



Replay buffer

DQWIC

PSEUDOCODE

Algorithm: Deep Q-learning Whittle Index with Context (DQWIC)

Initialize: Q-network parameters θ , Whittle network parameters δ , target network parameters θ^{tg} , replay memory \mathcal{D} , and hyperparameters ϵ, γ .

```
1: Get initial state  $s$  and context  $c$ 
2: for each time step  $t$  do
3:   if Uniform[0, 1] <  $\epsilon$  then
4:     Explore by selecting  $M$  random arms
5:   else
6:     Exploit by selecting top  $M$  arms with the highest Whittle indices
7:   end if
8:   Execute action  $a$ , observe context  $c$ , next state  $s'$  and reward  $r$ 
9:   Store transitions  $(s, c, a, r, s')$  in replay memory  $\mathcal{D}$ 
10:  if  $|\mathcal{D}| >$  batch size then
11:    /* On a faster time scale */
12:    Sample mini-batch of transitions from  $\mathcal{D}$ 
13:    Compute  $Q^{\text{target}}$ 
14:    Compute loss:  $\mathcal{L}(\theta)$ 
15:    Update Q-network parameters  $\theta$ 
16:    /* On a slower time scale */
17:    Sample batch of  $(k, c)$  from  $\mathcal{D}$ 
18:    Compute loss:  $\mathcal{E}(\delta)$ 
19:    Update Whittle network parameters  $\delta$ 
20:  end if
21:  Update target network parameters  $\theta^{tg} \leftarrow \theta$ 
22: end for
```

DQWIC

PSEUDOCODE

1

Algorithm: Deep Q-learning Whittle Index with Context (DQWIC)

Initialize: Q-network parameters θ , Whittle network parameters δ , target network parameters θ^{tg} , replay memory \mathcal{D} , and hyperparameters ϵ, γ .

1: Get initial state s and context c

2: **for** each time step t **do**

3: **if** Uniform[0, 1] < ϵ **then**

4: Explore by selecting M random arms

5: **else**

6: Exploit by selecting top M arms with the highest Whittle indices

7: **end if**

8: Execute action a , observe context c , next state s' and reward r

9: Store transitions (s, c, a, r, s') in replay memory \mathcal{D}

10: **if** $|\mathcal{D}| >$ batch size **then**

11: */* On a faster time scale */*

12: Sample mini-batch of transitions from \mathcal{D}

13: Compute Q^{target}

14: Compute loss: $\mathcal{L}(\theta)$

15: Update Q-network parameters θ

16: */* On a slower time scale */*

17: Sample batch of (k, c) from \mathcal{D}

18: Compute loss: $\mathcal{E}(\delta)$

19: Update Whittle network parameters δ

20: **end if**

21: Update target network parameters $\theta^{tg} \leftarrow \theta$

22: **end for**

2

3

4

5

1 Parameter initialization

- Q-network params θ
- Whittle network params δ
- Replay memory D
- Exploration param ϵ
- Discount factor γ
- State s and context c

Algorithm: Deep Q-learning Whittle Index with Context (DQWIC)

Initialize: Q-network parameters θ , Whittle network parameters δ , target network parameters θ^{tg} , replay memory \mathcal{D} , and hyperparameters ϵ, γ .

1: Get initial state s and context c

2: **for** each time step t **do**

3: **if** Uniform[0, 1] < ϵ **then**

4: Explore by selecting M random arms

5: **else**

6: Exploit by selecting top M arms with the highest Whittle indices

7: **end if**

8: Execute action a , observe context c , next state s' and reward r

9: Store transitions (s, c, a, r, s') in replay memory \mathcal{D}

10: **if** $|\mathcal{D}| >$ batch size **then**

11: */* On a faster time scale */*

12: Sample mini-batch of transitions from \mathcal{D}

13: Compute Q^{target}

14: Compute loss: $\mathcal{L}(\theta)$

15: Update Q-network parameters θ

16: */* On a slower time scale */*

17: Sample batch of (k, c) from \mathcal{D}

18: Compute loss: $\mathcal{E}(\delta)$

19: Update Whittle network parameters δ

20: **end if**

21: Update target network parameters $\theta^{tg} \leftarrow \theta$

22: **end for**

Algorithm: Deep Q-learning Whittle Index with Context (DQWIC)

Initialize: Q-network parameters θ , Whittle network parameters δ , target network parameters θ^{tg} , replay memory \mathcal{D} , and hyperparameters ϵ, γ .

1: Get initial state s and context c

2: **for** each time step t **do**

3: **if** Uniform[0, 1] < ϵ **then**

4: Explore by selecting M random arms

5: **else**

6: Exploit by selecting top M arms with the highest Whittle indices

7: **end if**

8: Execute action a , observe context c , next state s' and reward r

9: Store transitions (s, c, a, r, s') in replay memory \mathcal{D}

10: **if** $|\mathcal{D}| >$ batch size **then**

11: */* On a faster time scale */*

12: Sample mini-batch of transitions from \mathcal{D}

13: Compute Q^{target}

14: Compute loss: $\mathcal{L}(\theta)$

15: Update Q-network parameters θ

16: */* On a slower time scale */*

17: Sample batch of (k, c) from \mathcal{D}

18: Compute loss: $\mathcal{E}(\delta)$

19: Update Whittle network parameters δ

20: **end if**

21: Update target network parameters $\theta^{tg} \leftarrow \theta$

22: **end for**

1 Parameter initialization

2 Selection of M arms

- Epsilon-greedy:
 - Explore with ϵ
 - Exploit with $1 - \epsilon$
- Observe transitions and store them in replay memory \mathcal{D}

Algorithm: Deep Q-learning Whittle Index with Context (DQWIC)

Initialize: Q-network parameters θ , Whittle network parameters δ , target network parameters θ^{tg} , replay memory \mathcal{D} , and hyperparameters ϵ, γ .

1: Get initial state s and context c

2: **for** each time step t **do**

3: **if** Uniform[0, 1] < ϵ **then**

4: Explore by selecting M random arms

5: **else**

6: Exploit by selecting top M arms with the highest Whittle indices

7: **end if**

8: Execute action a , observe context c , next state s' and reward r

9: Store transitions (s, c, a, r, s') in replay memory \mathcal{D}

10: **if** $|\mathcal{D}| >$ batch size **then**

11: */* On a faster time scale */*

12: Sample mini-batch of transitions from \mathcal{D}

13: Compute Q^{target}

14: Compute loss: $\mathcal{L}(\theta)$

15: Update Q-network parameters θ

16: */* On a slower time scale */*

17: Sample batch of (k, c) from \mathcal{D}

18: Compute loss: $\mathcal{E}(\delta)$

19: Update Whittle network parameters δ

20: **end if**

21: Update target network parameters $\theta^{tg} \leftarrow \theta$

22: **end for**

1 Parameter initialization

2 Selection of M arms

3 Run on fast timescale

Algorithm: Deep Q-learning Whittle Index with Context (DQWIC)

Initialize: Q-network parameters θ , Whittle network parameters δ , target network parameters θ^{tg} , replay memory \mathcal{D} , and hyperparameters ϵ, γ .

1: Get initial state s and context c

2: **for** each time step t **do**

3: **if** Uniform[0, 1] < ϵ **then**

4: Explore by selecting M random arms

5: **else**

6: Exploit by selecting top M arms with the highest Whittle indices

7: **end if**

8: Execute action a , observe context c , next state s' and reward r

9: Store transitions (s, c, a, r, s') in replay memory \mathcal{D}

10: **if** $|\mathcal{D}| >$ batch size **then**

11: */* On a faster time scale */*

12: Sample mini-batch of transitions from \mathcal{D}

13: Compute Q^{target}

14: Compute loss: $\mathcal{L}(\theta)$

15: Update Q-network parameters θ

16: */* On a slower time scale */*

17: Sample batch of (k, c) from \mathcal{D}

18: Compute loss: $\mathcal{E}(\delta)$

19: Update Whittle network parameters δ

20: **end if**

21: Update target network parameters $\theta^{tg} \leftarrow \theta$

22: **end for**

1 Parameter initialization

2 Selection of M arms

3 Run on fast timescale

4 Run on slow timescale

Algorithm: Deep Q-learning Whittle Index with Context (DQWIC)

Initialize: Q-network parameters θ , Whittle network parameters δ , target network parameters θ^{tg} , replay memory \mathcal{D} , and hyperparameters ϵ, γ .

1: Get initial state s and context c

2: **for** each time step t **do**

3: **if** Uniform[0, 1] < ϵ **then**

4: Explore by selecting M random arms

5: **else**

6: Exploit by selecting top M arms with the highest Whittle indices

7: **end if**

8: Execute action a , observe context c , next state s' and reward r

9: Store transitions (s, c, a, r, s') in replay memory \mathcal{D}

10: **if** $|\mathcal{D}| >$ batch size **then**

11: */* On a faster time scale */*

12: Sample mini-batch of transitions from \mathcal{D}

13: Compute Q^{target}

14: Compute loss: $\mathcal{L}(\theta)$

15: Update Q-network parameters θ

16: */* On a slower time scale */*

17: Sample batch of (k, c) from \mathcal{D}

18: Compute loss: $\mathcal{E}(\delta)$

19: Update Whittle network parameters δ

20: **end if**

21: Update target network parameters $\theta^{tg} \leftarrow \theta$

22: **end for**

1 Parameter initialization

2 Selection of M arms

3 Run on fast timescale

4 Run on slow timescale

5 NN parameter update

Algorithm: Deep Q-learning Whittle Index with Context (DQWIC)

Initialize: Q-network parameters θ , Whittle network parameters δ , target network parameters θ^{tg} , replay memory \mathcal{D} , and hyperparameters ϵ, γ .

1: Get initial state s and context c

2: **for** each time step t **do**

3: **if** Uniform[0, 1] < ϵ **then**

4: Explore by selecting M random arms

5: **else**

6: Exploit by selecting top M arms with the highest Whittle indices

7: **end if**

8: Execute action a , observe context c , next state s' and reward r

9: Store transitions (s, c, a, r, s') in replay memory \mathcal{D}

10: **if** $|\mathcal{D}| >$ batch size **then**

11: */* On a faster time scale */*

12: Sample mini-batch of transitions from \mathcal{D}

13: Compute Q^{target}

14: Compute loss: $\mathcal{L}(\theta)$

15: Update Q-network parameters θ

16: */* On a slower time scale */*

17: Sample batch of (k, c) from \mathcal{D}

18: Compute loss: $\mathcal{E}(\delta)$

19: Update Whittle network parameters δ

20: **end if**

21: Update target network parameters $\theta^{tg} \leftarrow \theta$

22: **end for**

1 Parameter initialization

2 Selection of M arms

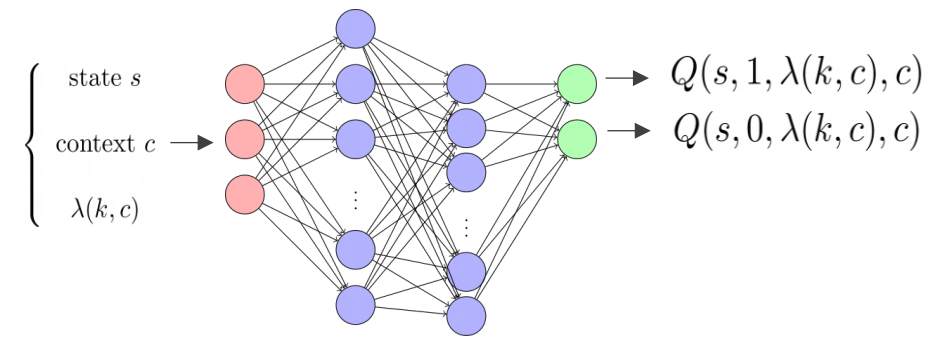
3 Run on fast timescale

4 Run on slow timescale

5 NN parameter update

DQWIC

Q-NETWORK



- On a fast time scale, sample a batch (which contains $[s, c, a, r, s']$) of size B from the memory D .
- The Q-network computes the predicted Q-values for all possible actions as follows:

$$Q^{\text{target}}(s_t, a_t, \lambda_\delta, c_t) = (1 - a_t)(r_0(s_t) + \lambda_\delta) + a_t r_1(s_t) + \gamma \max_{a \in \{0,1\}} Q^{\theta_{tg}}(s'_t, a, \lambda_\delta, c_t)$$

- Target network is a copy of the Q-network.
- $\max_{a \in \{0,1\}} Q^{\theta_{tg}}(s'_t, a, \lambda_\delta, c_t)$ is the maximum Q-value for the next state s'_t over all possible actions, predicted by the target network parameters θ^{tg} .

DQWIC

PARAMETER UPDATE

- The Q-network parameters θ are updated by minimizing the Mean Squared Error (MSE) between the predicted Q-values and the target Q-values.

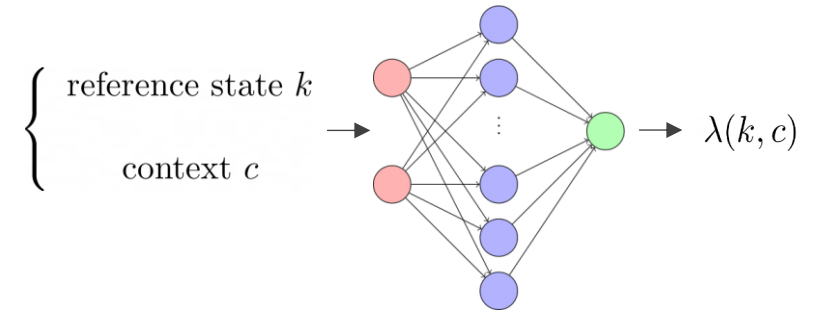
- The loss function is defined as:

$$\mathcal{L}(\theta) := \mathbb{E} \left[\left\| Q_{\theta}(s, a, \lambda_s, c) - Q^{\text{target}}(s, a, \lambda_s, c) \right\|^2 \right]$$

- The parameters θ of the Q-network are updated using backpropagation to minimize this loss.
- The parameters of the target network θ^{tg} are periodically synchronized with those of the Q-network θ , for instance every 100 iterations, to stabilize training.

DQWIC

WHITTLE NETWORK



- On a slow time scale, a batch of reference states k and contexts c are sampled from the replay memory D . These samples are used to compute the MSE loss function $E(\delta)$ based on the difference between the Q-values for the active and passive actions.

$$\mathcal{E}(\delta) = \mathbb{E} \left[\|Q_{\theta}(k, 1, \lambda_s, c) - Q_{\theta}(k, 0, \lambda_s, c)\|^2 \right]$$

- The parameters δ of the Whittle network are updated using gradient descent to minimize the loss $E(\delta)$.

APPLICATION

RECOMMENDER SYSTEMS

- A Recommender System (RS) is:
 - an information filtering system
 - suggests items to users based on their past preferences or behavior.
 - used in many domains such as in: e-commerce, entertainment, and social media, etc.
- Let $R(u, j)$ represent the interaction between user u and item j . This could be:
 - **Explicit:** Ratings (e.g., $R(u, j) \in \{1, 2, \dots, 5\}$).
 - **Implicit:** Binary clicks or preferences (e.g., $R(u, j) \in \{0, 1\}$).

RECOMMENDER SYSTEM

EXAMPLES

Because you watched this show ...

NETFLIX

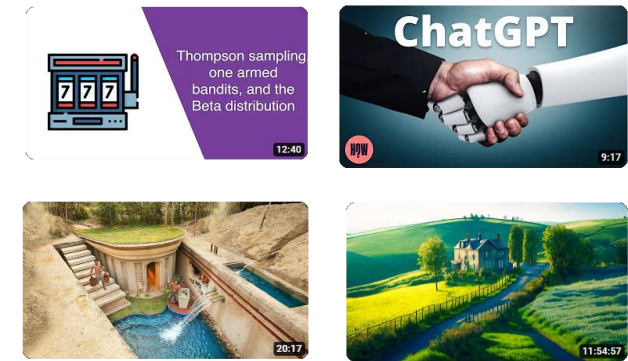


Customers like you also bought ...

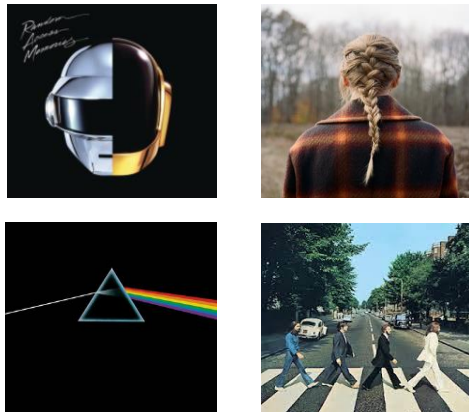
amazon



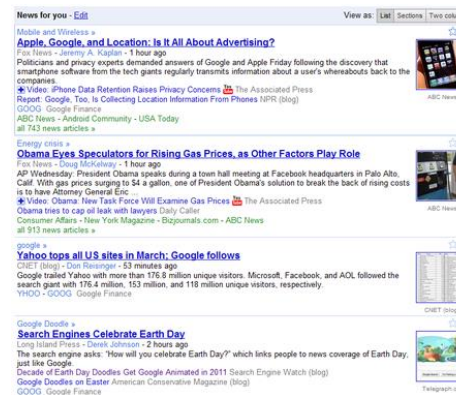
Videos recommended for you ...



See what's popular now ...



Recommended articles for you ...



Items similar to your purchase ...

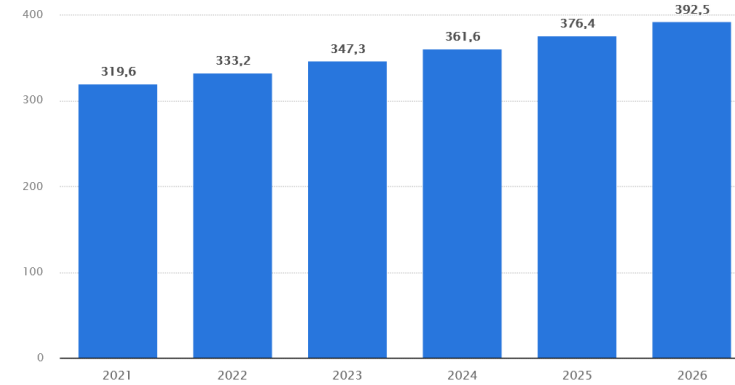
DECATHLON



EMAIL RECOMMENDER SYSTEM

MOTIVATION

- Sending mass emails results in suboptimal performance:
 - bad domain reputation with potential **spamming**,
 - overwhelming recipients with irrelevant content,
 - decreasing user engagement,
 - negative **user experience**.



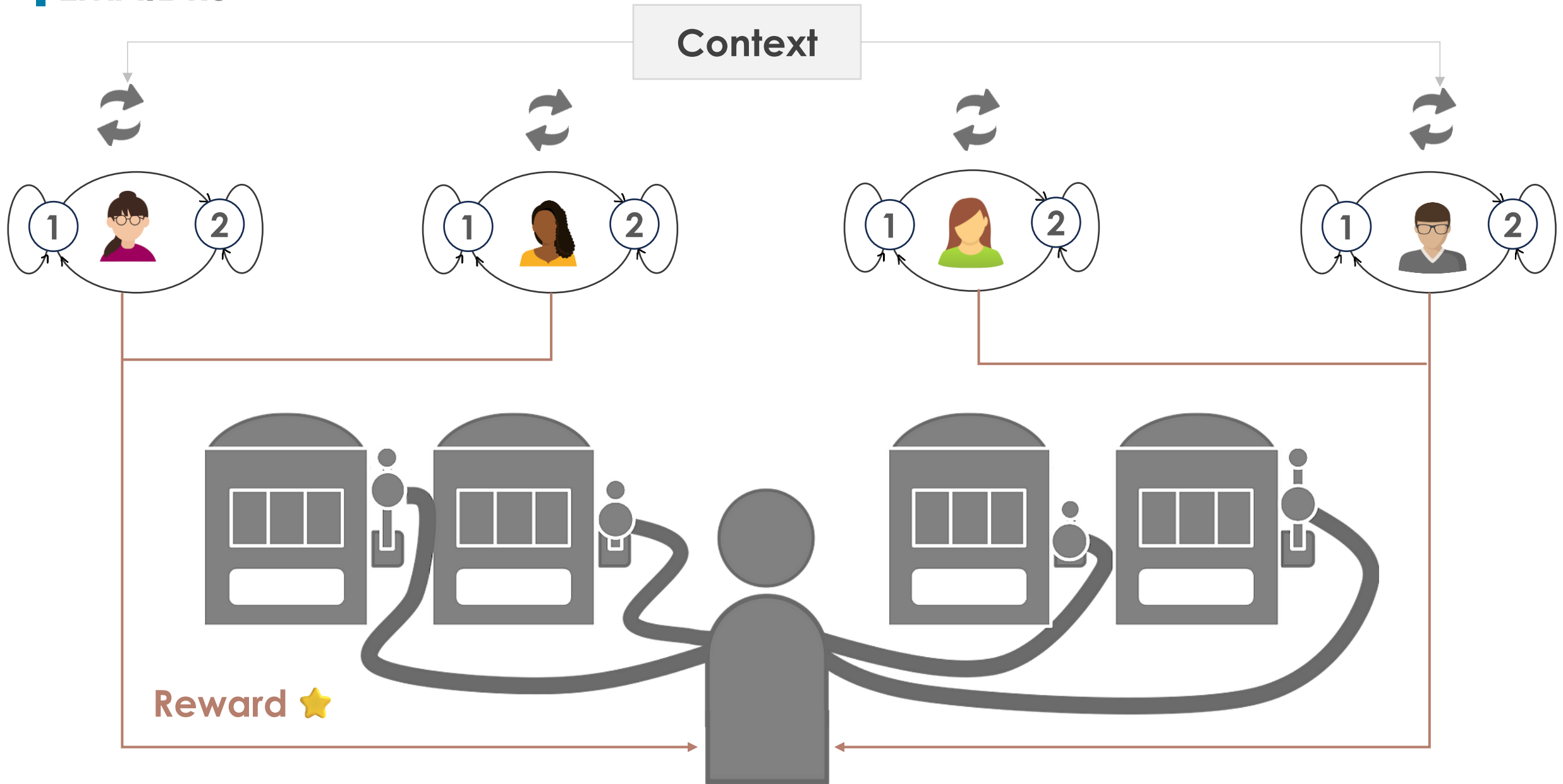
~ 347 billion
daily emails
(2023)



~ 46 % of
emails are
spams
(2023)

CONTEXTUAL RESTLESS BANDITS

EMAIL RS



EMAIL RECOMMENDER SYSTEM

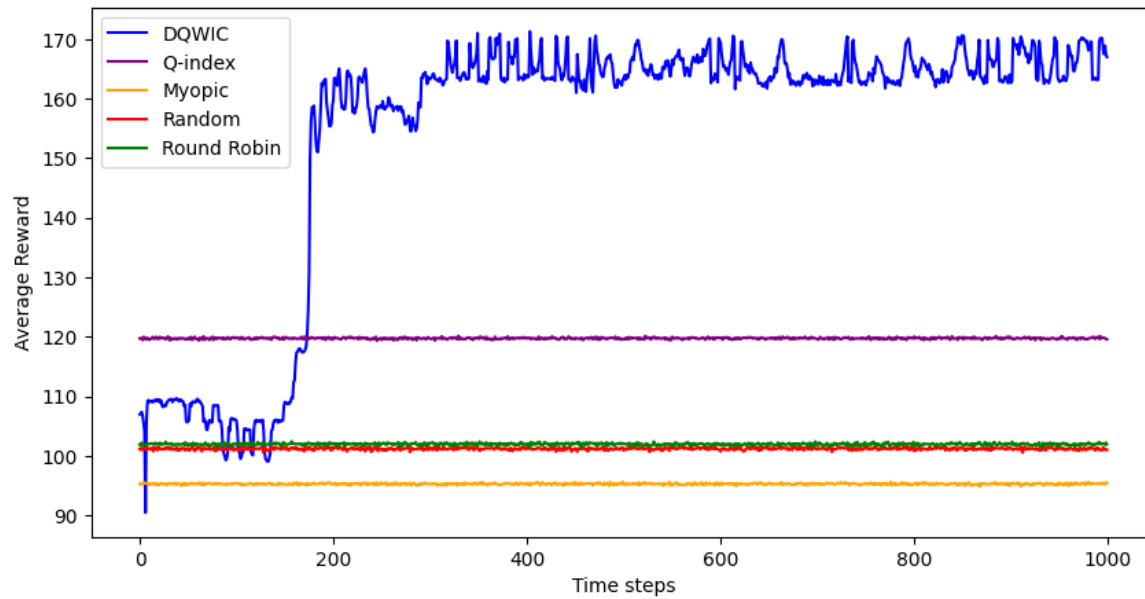
BASELINES

Policy	Definition
Random	Selects users randomly without considering engagement information.
Myopic	Selects users who are most likely to lead to immediate conversions, such as those who have recently opened emails.
Round-robin	Selects to contact users sends emails evenly in a cyclic order.
Q-index	The users are ranked according to the active Q-values approximated by a deep neural network.

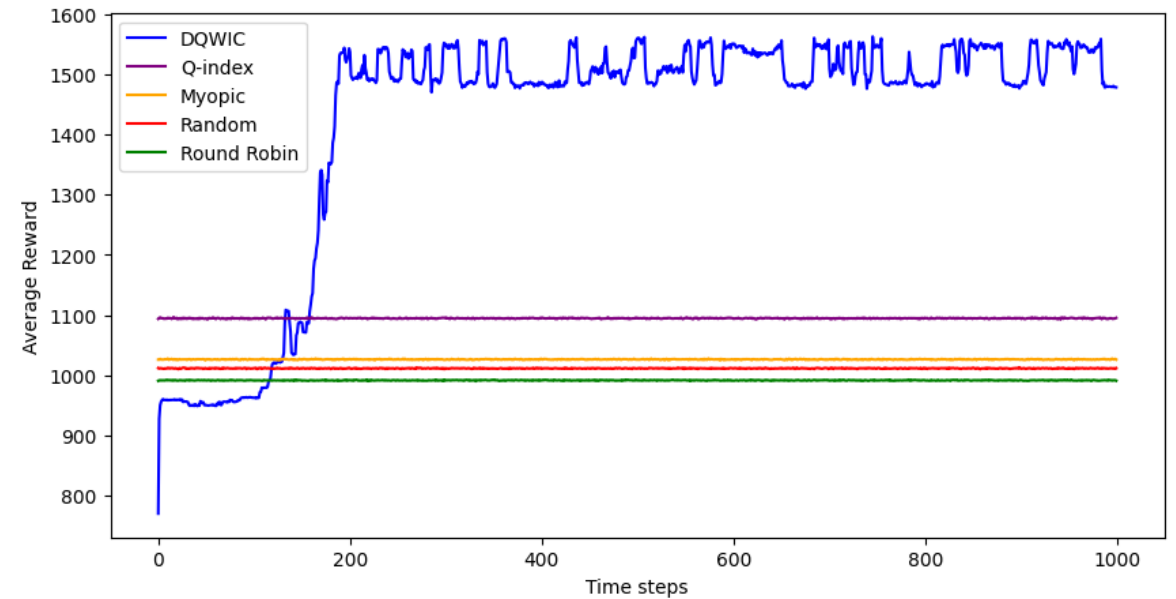
EMAIL RECOMMENDER SYSTEM

RESULTS

- **Q-network:** 2 hidden layers with 130 and 50 neurons, respectively, connected by ReLU.
- **Whittle network:** 1 hidden layer with 250 neurons, connected by ReLU.



Real Data | N = 100, M = 10



Synthetic Data | N = 1000, M = 100

CONCLUSION

- We proposed deep Q-learning with Whittle index for **CRMAB**, and applied it to an **email recommender system**.
- One important feature of our algorithm is the usage of only three neural networks for many (potentially thousands) heterogeneous arms. This is possible thanks to the **addition of context to the RMAB model**.
- **Experiments** on both **synthetic and real-world data** show that DQWIC outperforms existing baselines.
- **Future work** will focus on incorporating multiple actions, and enhancing fairness.

THANK YOU

REFERENCES

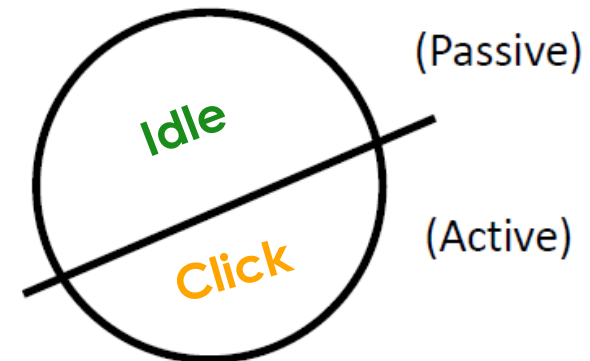
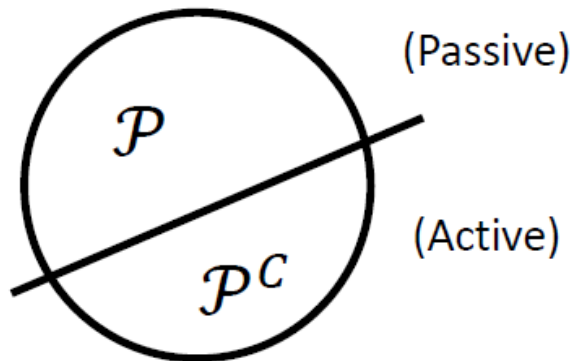
- [1]** X. Chen, I. Hou, et al. “Contextual Restless Multi-Armed Bandits with Application to Demand Response Decision-Making”. In: arXiv preprint arXiv:2403.15640 (2024)
- [2]** B. Liang, L. Xu, A. Taneja, M. Tambe, and L. Janson. “A Bayesian Approach to Online Learning for Contextual Restless Bandits with Applications to Public Health”. In: arXiv preprint arXiv:2402.04933 (2024).

WHITTLE INDEX

SOLUTION TO THE DECOUPLED PROBLEM

- In general, the optimal policy divides the state space into two subsets:
 - Let $\mathcal{P}(\tilde{\lambda})$ be the set of ALL states for which it is **optimal to be idle** when the playing charge is $\tilde{\lambda}$.
 - Let $\mathcal{P}^c(\tilde{\lambda})$ be the set of ALL states for which it is **optimal to be in click** when the playing charge is $\tilde{\lambda}$.
- **Optimal policy:** play, if $s_t^{(i)} \in \mathcal{P}^c(\tilde{\lambda})$; stop, otherwise.

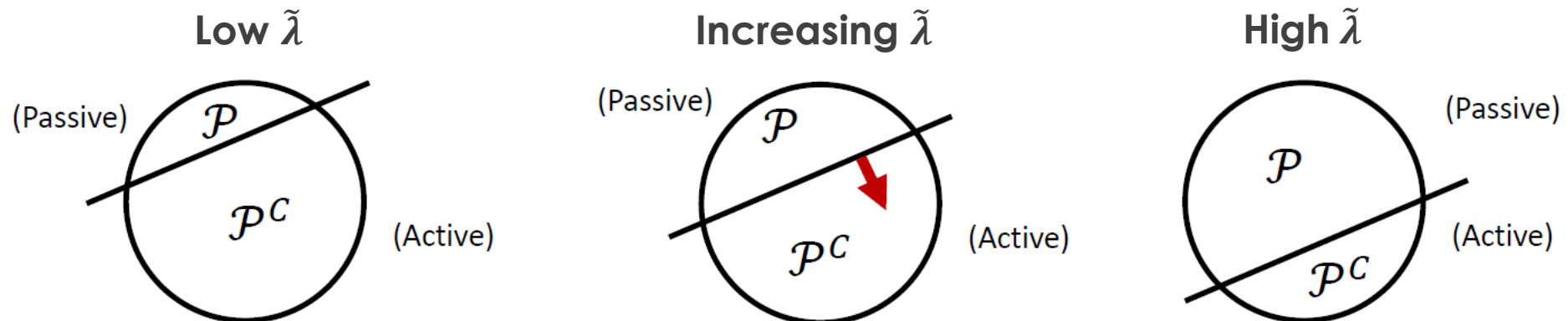
State space with $\tilde{\lambda}$



WHITTLE INDEX

INDEXABILITY

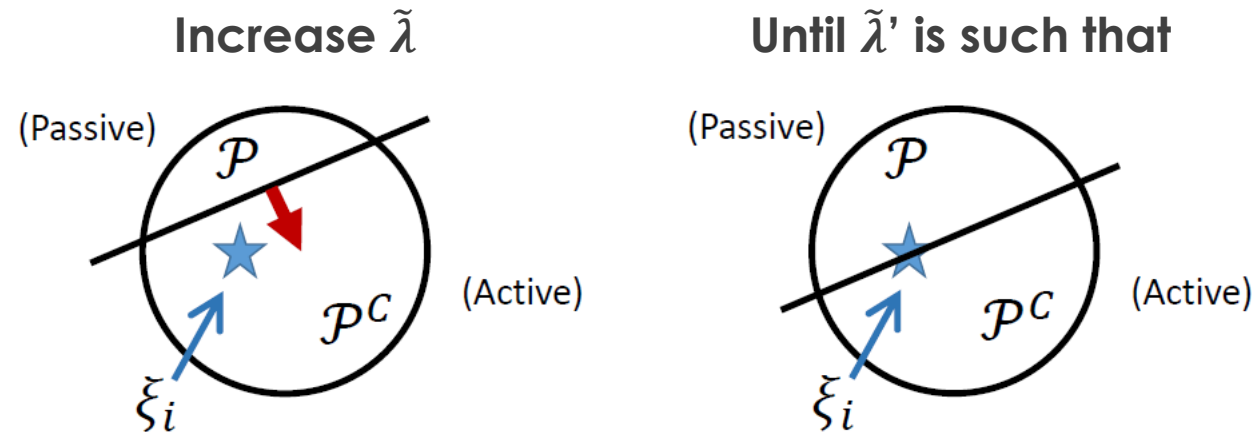
- An arm α is considered to be indexable if the set P increases monotonically from Φ to the entire state space as we increase $\tilde{\lambda}^\alpha$ (the subsidy for passive action).
- The RMAB problem is indexable if the Decoupled Problem is indexable for all bandits.



WHITTLE INDEX

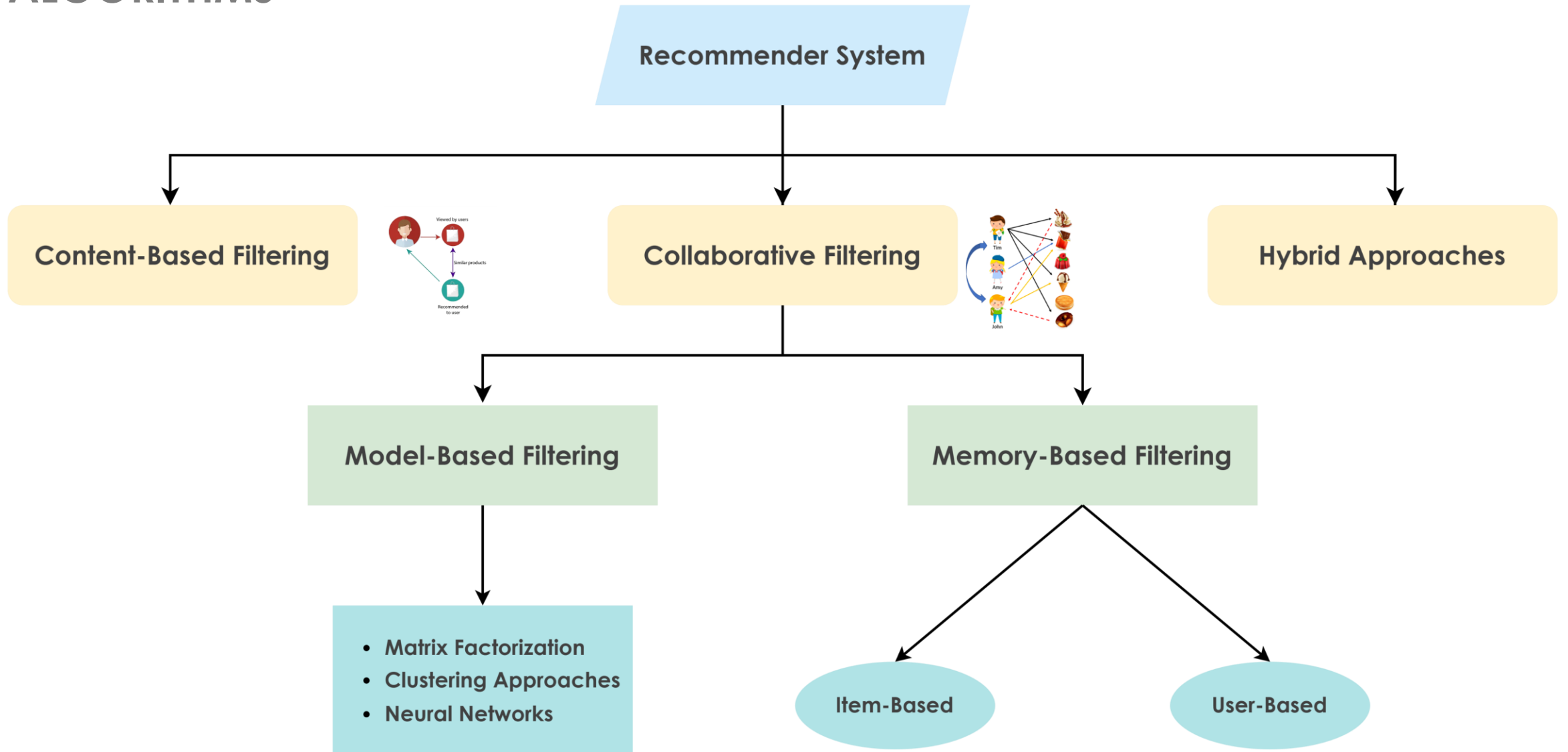
DEFINITION

- For an arm α , Whittle index is defined as the $\lambda^\alpha = \inf \{\lambda^\alpha : \alpha \in P\}$.
- It is the infimum subsidy λ^α one has to pay so that it is equally desirable to give an arm α active and passive action.



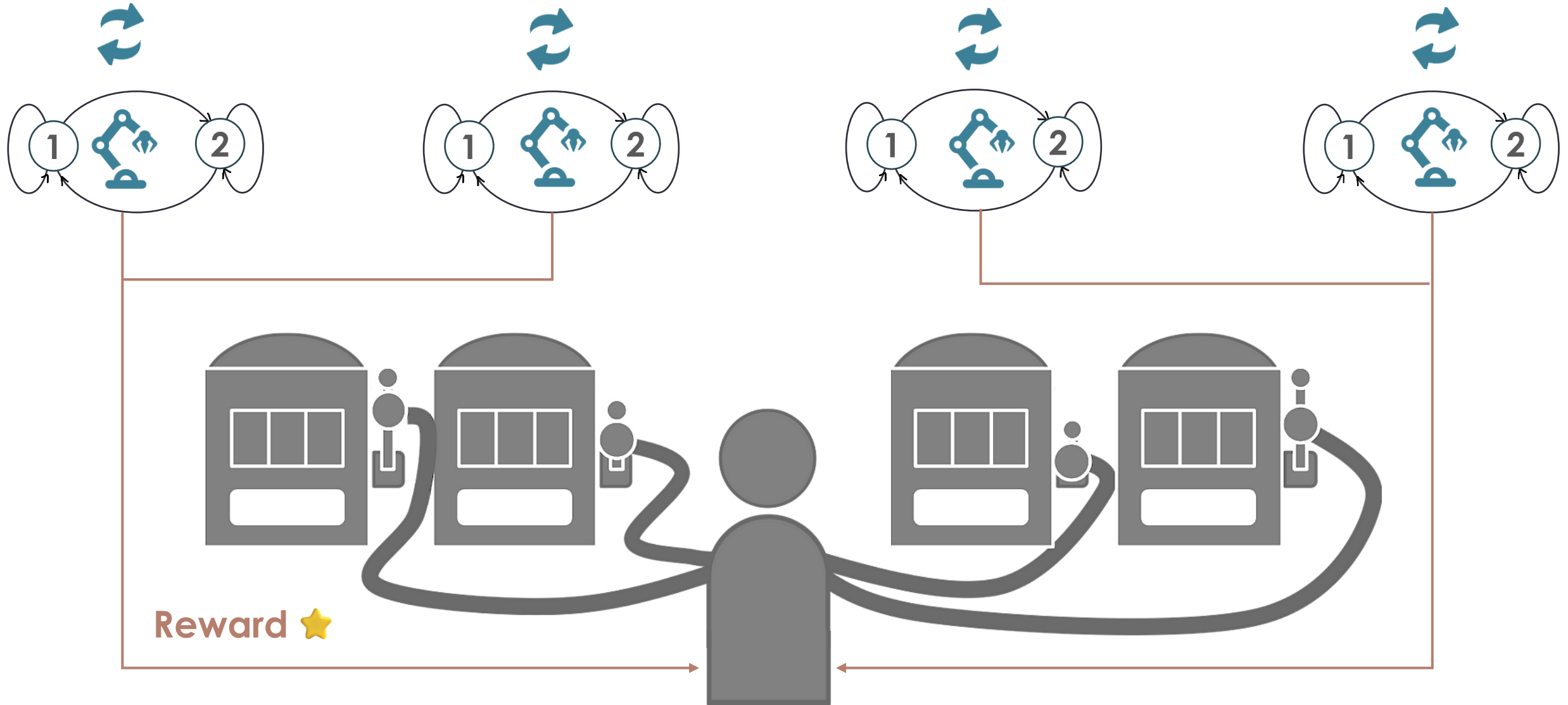
RECOMMENDER SYSTEM

ALGORITHMS



RESTLESS BANDITS

RECALL



APPLICATION

RECOMMENDER SYSTEMS

- A Recommender System (RS) is:
 - an information filtering system
 - suggests items to users based on their past preferences or behavior.
 - used in many domains such as in: e-commerce, entertainment, and social media, etc.
- Let $R(u, j)$ represent the interaction between user u and item j . This could be:
 - **Explicit:** Ratings (e.g., $R(u, j) \in \{1, 2, \dots, 5\}$).
 - **Implicit: Binary clicks or preferences (e.g., $R(u, j) \in \{0, 1\}$).**