

# IMPERIAL

## **GINO-Q: Learning an Asymptotically Optimal Index Policy for Restless Multi-armed Bandits**

Gongpu Chen, Soung Chang Liew, Deniz Gündüz

Imperial College London & The Chinese University of Hong Kong  
AAAI2026, Jan 2026

# Restless Multi-armed Bandits (RMABs)

## Problem Statement

- M Markov arms, each corresponding to an MDP  $\mathcal{B}_i = (\mathcal{S}_i, \mathcal{A}_i, r_i, p_i)$ , where
  - state space  $\mathcal{S}_i$ , a finite set
  - action space  $\mathcal{A}_i = \{0, 1\}$
  - reward function  $r_i(\mathbf{s}, \mathbf{a})$
  - transition kernel  $p_i(\mathbf{s}_i^{t+1} | \mathbf{s}_i^t, \mathbf{a}_i^t)$
- all arms evolve independently based on their actions

$$P(\mathbf{s}_1^{t+1}, \dots, \mathbf{s}_M^{t+1} | \mathbf{s}_1^t, \dots, \mathbf{s}_M^t, \mathbf{a}_1^t, \dots, \mathbf{a}_M^t) = \prod_{i=1}^M p_i(\mathbf{s}_i^{t+1} | \mathbf{s}_i^t, \mathbf{a}_i^t)$$

- the arms are weakly coupled by a constraint:  $\sum_{i=1}^M a_i^t = N$  for all t  
→ only  $N < M$  arms can be activated at each time
- Objective:

$$\max_{\{a_i^t: i \in [M], t \geq 1\}} \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T \sum_{i=1}^M r_i(\mathbf{s}_i^t, \mathbf{a}_i^t) \right]$$

# Restless Multi-armed Bandits (RMABs)

## An MDP with large space and action spaces

An RMAB can be formulated as an MDP with

- joint state space  $\mathcal{S} = \mathcal{S}_1 \times \mathcal{S}_2 \times \cdots \times \mathcal{S}_M$
- joint action space  $\mathcal{A} = \{(a_1, a_1, \cdots, a_M) : \sum_{i=1}^M a_i = N\}$

# Restless Multi-armed Bandits (RMABs)

## An MDP with large space and action spaces

An RMAB can be formulated as an MDP with

- joint state space  $\mathcal{S} = \mathcal{S}_1 \times \mathcal{S}_2 \times \dots \times \mathcal{S}_M \Rightarrow$  grow exponentially with  $M$
- joint action space  $\mathcal{A} = \{(\mathbf{a}_1, \mathbf{a}_1, \dots, \mathbf{a}_M) : \sum_{i=1}^M \mathbf{a}_i = \mathbf{N}\} \Rightarrow$  a combinatorial space of size  $\binom{M}{N}$

Example: A moderate-scale RMAB with  $M = 100$ ,  $N = 25$ ,  $|\mathcal{S}_i| = 10$ . Then

$$|\mathcal{S}| = 10^{100}, |\mathcal{A}| = \binom{100}{25} = 2.4 \times 10^{23}$$

Hard to solve using standard RL methods!

SOTA solutions: Whittle-index-based RL methods  $\rightarrow$  Learn the **Whittle index policy** in model-free setting

# Whittle Index Policy

## Decouple the RMAB across arms

- Relaxation

$$\sum_{i=1}^M a_i^t = N, \forall t \quad \rightarrow \quad \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T \sum_{i=1}^M a_i^t \right] = N$$

- Lagrange Multiplier method

$$\max_{\{a_i^t\}} \inf_{\lambda} \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T \sum_{i=1}^M \left[ r_i(s_i^t, a_i^t) - \lambda a_i^t \right] \right] + N\lambda$$

- Decomposition (for any fixed  $\lambda$ )

$$J_i(\lambda) : \max_{\{a_i^t: t \geq 1\}} \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T \left[ r_i(s_i^t, a_i^t) - \lambda a_i^t \right] \right]$$

# Whittle Index Policy

## Decompose the RMAB into single-arm problems

Each single-arm problem  $J_i(\lambda)$  is an MDP associated with the arm  $\mathcal{B}_i$

$$J_i(\lambda) : \max_{\{a_i^t: t \geq 1\}} \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T \left[ r_i(s_i^t, a_i^t) - \lambda a_i^t \right] \right]$$

- state space  $\mathcal{S}_i$ , a finite set
- action space  $\mathcal{A}_i = \{0, 1\}$
- reward function  $r_i(\mathbf{s}, \mathbf{a}) - \lambda \mathbf{a}$   $\rightarrow$   $\lambda$  can be interpreted as the cost of action 1, called service charge
- transition kernel  $p_i(s_i^{t+1} | s_i^t, a_i^t)$

# Whittle Index Policy

## Definition and Indexability

### Whittle index policy

- Assign a Whittle index to each state of each arm:

$$W_i(\mathbf{s}) = \inf\{\lambda : Q_i(\mathbf{s}, 1, \lambda) = Q_i(\mathbf{s}, 0, \lambda)\}.$$

where  $Q_i(\mathbf{s}, a, \lambda)$  is the Q-function of  $J_i(\lambda)$

(Interpretation:  $W_i(\mathbf{s})$  is the infimum service charge that makes the two actions equally optimal)

- At each time, select the top N arms with the highest Whittle indices

### **The Whittle index policy is only applicable to indexable RMABs**

### Indexability

- Denote by  $\mathcal{E}_i(\lambda)$  the set of states in which action 0 is optimal for problem  $J_i(\lambda)$
- An arm  $\mathcal{B}_i$  is considered indexable if  $\mathcal{E}_i(\lambda)$  expands monotonically from the empty set to the entire state space  $\mathcal{S}_i$  as  $\lambda$  increases from  $-\infty$  to  $\infty$ .
- An RMAB is indexable if all its arms are indexable.

# Whittle-index-based learning methods

## Idea and limitation

**Basic idea:** Learn the Whittle index for each state of each arm in the model-free setting, and select the top  $N$  arms with highest indices at each time

- If an arm is indexable, there is a unique  $\lambda$  that satisfies

$$Q_i(\mathbf{s}, 1, \lambda) = Q_i(\mathbf{s}, 0, \lambda) \quad (1)$$

- learn a  $\lambda$  satisfying (1) for each state  $\mathbf{s}$  of each arm, use it as the Whittle index

### Limitations:

- The Whittle index policy is only applicable to indexable RMABs—NOT all RMABs are indexable
- Verifying indexability requires full knowledge of the system and considerable analytical effort

Existing Whittle-index-based learning methods assume the RMAB is always indexable, which is NOT true.

## What if using Whittle-index-based methods in non-indexable RMABs?

# Whittle-index-based learning methods

## Counter example

### A non-indexable arm

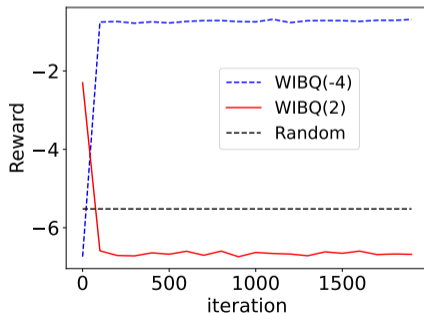
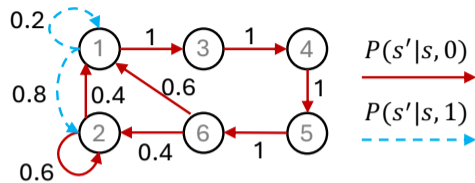
- Transition probabilities as shown in the right figure
- Reward function defined as  $r_i(1, 1) = -10$ ,  $r_i(1, 0) = -4$ ,  $r_i(2, 1) = r_i(2, 0) = 4$  and  $r_i(s, 1) = 0$ ,  $r_i(s, 0) = -2$  for  $s \in \{3, 4, 5, 6\}$ .

The arm is NOT indexable because

- state  $1 \in \mathcal{E}_i(\lambda)$  for  $-4 \leq \lambda \leq 2$  and
- state  $1 \notin \mathcal{E}_i(\lambda)$  for  $\lambda < -4$  and  $\lambda > 2$

When learning the Whittle index for state 1, the algorithm may randomly converge to either 2 or  $-4$

The performance can be arbitrarily poor!



# It's necessary to develop a method that does not require indexability

# Gain Index Policy

## Decouple the RMAB across arms

- Relaxation

$$\sum_{i=1}^M a_i^t = N, \forall t \quad \rightarrow \quad \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T \sum_{i=1}^M a_i^t \right] = N$$

- Lagrange Multiplier method

$$\max_{\{a_i^t\}} \inf_{\lambda} \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T \sum_{i=1}^M \left[ r_i(s_i^t, a_i^t) - \lambda a_i^t \right] \right] + N\lambda,$$

- Equivalent transform:

$$\text{Relaxed RMAB: } \inf_{\lambda} \max_{\{a_i^t\}} \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T \sum_{i=1}^M \left[ r_i(s_i^t, a_i^t) - \lambda a_i^t \right] \right] + N\lambda,$$

# Gain Index Policy

## Decompose the RMAB into single-arm problems

Each single-arm problem  $J_i(\lambda)$  is an MDP associated with the arm  $\mathcal{B}_i$

$$J_i(\lambda) : \max_{\{a_i^t: t \geq 1\}} \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T \left[ r_i(s_i^t, a_i^t) - \lambda a_i^t \right] \right]$$

- state space  $\mathcal{S}_i$ , a finite set
- action space  $\mathcal{A}_i = \{0, 1\}$
- reward function  $r_i(\mathbf{s}, \mathbf{a}) - \lambda a$   $\rightarrow$   $\lambda$  can be interpreted as the cost of action 1, called service charge
- transition kernel  $p_i(s_i^{t+1} | s_i^t, a_i^t)$

The Bellman optimality equation:

$$V_i(\mathbf{s}, \lambda) + \mathbf{g}_i(\lambda) = \max_{a \in \{0,1\}} \left\{ r_i(\mathbf{s}, a) - \lambda a + \sum_{s' \in \mathcal{S}_i} p_i(s' | \mathbf{s}, a) V_i(s', \lambda) \right\}, \mathbf{s} \in \mathcal{S}_i,$$

# Gain Index Policy

## Decompose the RMAB into Single-arm Problems

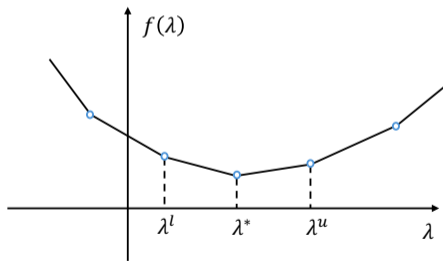
The relaxed RMAB problem reduces to

$$\inf_{\lambda} f(\lambda) \triangleq \sum_{i=1}^M g_i(\lambda) + N\lambda.$$

Denote by  $\lambda^*$  the optimal solution to the relaxed RMAB problem.

### Lemma

For any RMAB with bounded reward functions  $\{r_i\}_{i \in [M]}$ ,  $f(\lambda)$  is a piecewise linear and convex function. In addition, there always exists a bounded  $\lambda^*$  that achieves the minimum value of  $f(\lambda)$ .



# Gain Index Policy

## Definition and Properties

For any  $J_i(\lambda)$ , the Q-function (state-action value function) is:

$$Q_i(\mathbf{s}, \mathbf{a}, \lambda) \triangleq r_i(\mathbf{s}, \mathbf{a}) - \lambda \mathbf{a} + \sum_{\mathbf{s}' \in \mathcal{S}_i} p_i(\mathbf{s}' | \mathbf{s}, \mathbf{a}) V_i(\mathbf{s}', \lambda).$$

- We express  $V_i$  and  $Q_i$  as functions of  $\lambda$  to highlight their dependence on  $\lambda$ .

### Definition (Gain index policy)

For each arm  $i \in [M]$  and each state  $\mathbf{s} \in \mathcal{S}_i$ , a gain index is defined as:

$$W_i(\mathbf{s}) \triangleq Q_i(\mathbf{s}, 1, \lambda^*) - Q_i(\mathbf{s}, 0, \lambda^*). \quad (2)$$

Then the gain index of the  $i$ -th arm at time  $t$  is given by  $W_i(\mathbf{s}_i^t)$ . The gain index policy activates the  $N$  arms with the largest  $N$  gain indices, with ties broken arbitrarily.

# Gain Index Policy

## Definition and Properties

### Theorem (Asymptotical Optimality, informal)

If  $N/M$  is fixed, then the gain index policy is asymptotically optimal in the following sense:

$$\lim_{M \rightarrow \infty} \frac{1}{M} G_M^{\text{ind}} = \lim_{M \rightarrow \infty} \frac{1}{M} G_M^{\text{opt}},$$

where  $G_M^{\text{ind}}$  and  $G_M^{\text{opt}}$  denote the cumulative reward of the RMAB under the gain index policy and the optimal policy, respectively.

#### Remarks:

- The gain index policy tends to be optimal when the number of arms is large enough
- Applicable to all RMABs  $\rightarrow$  Does NOT require indexability
- Computing gain indices requires knowledge of the system

# How to learn gain indices in the model-free setting?

## GINO-Q

### Learn the gain indices via Q-learning

**Goal:** Learn the gain index for each state  $s$  of each arm:

$$W_i(s) \triangleq Q_i(s, 1, \lambda^*) - Q_i(s, 0, \lambda^*).$$

- need to learn  $\lambda^*$  and the associated Q-functions
- $\lambda^*$  is the minimizer of the convex function  $f(\lambda)$

## GINO-Q

### Estimate the gradient of $f(\lambda)$

#### Definition (Auxiliary MDP)

The auxiliary MDP associated with arm  $\mathcal{B}_i$  is defined as  $\mathcal{M}_i = (\mathcal{S}_i, \mathcal{A}_i, c, p_i)$ , where the state space  $\mathcal{S}_i$ , action space  $\mathcal{A}_i$ , and transition kernel  $p_i$  are identical to those of  $\mathcal{B}_i$ . The cost function is given by  $c(s, a) = a$  for all  $s \in \mathcal{S}_i$  and  $a \in \mathcal{A}_i$ .

For any policy  $\pi$ , the average cost  $h_i^\pi$  is determined by

$$D_i^\pi(s, a) + h_i^\pi = a + \sum_{s', a'} p_i(s'|s, a) \pi(a'|s') D_i^\pi(s', a').$$

$D_i^\pi(s, a)$  denotes the state-action value function of  $\mathcal{M}_i$  under policy  $\pi$ .

$$\frac{dg_i^\pi}{d\lambda} = -h_i^\pi, \quad \forall i \in [M]. \quad (3)$$

Let  $\pi_i^\lambda$  denote the optimal policy for problem  $J_i(\lambda)$ , then the derivative of function  $f(\lambda)$  is given by

$$f'(\lambda) = \sum_{i=1}^M \frac{dg_i(\lambda)}{d\lambda} + N = N - \sum_{i=1}^M h_i^{\pi_i^\lambda}. \quad (4)$$

# GINO-Q

## A Three-timescale stochastic approximation method

- Slow timescale: stochastic gradient-descent

$$\lambda^{t+1} = \lambda^t - \theta^t (\mathbf{N} - \sum_{i=1}^M \mathbf{h}_i^t).$$

- Medium timescale: Q learning

$$\mathbf{Q}_i(\mathbf{s}_i^t, \mathbf{a}_i^t) \leftarrow (1 - \beta_i^t) \mathbf{Q}_i(\mathbf{s}_i^t, \mathbf{a}_i^t) + \beta_i^t [r_i(\mathbf{s}_i^t, \mathbf{a}_i^t) - \lambda^t \mathbf{a}_i^t + \max_{\mathbf{a}} \mathbf{Q}_i(\mathbf{s}_i^{t+1}, \mathbf{a}) - \mathbf{g}_i^t].$$

- Fast timescale: SARSA learning to estimate  $f'(\lambda)$

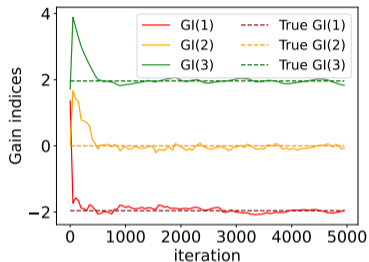
$$\mathbf{D}_i(\mathbf{s}_i^t, \mathbf{a}_i^t) \leftarrow (1 - \alpha_i^t) \mathbf{D}_i(\mathbf{s}_i^t, \mathbf{a}_i^t) + \alpha_i^t [\mathbf{a}_i^t + \mathbf{D}_i(\mathbf{s}_i^{t+1}, \mathbf{a}_i^{t+1}) - \mathbf{h}_i^t].$$

The timescales are controlled by the stepsizes:

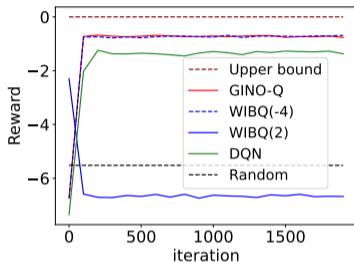
$$\alpha_i^t = \frac{C_1}{t}, \beta_i^t = \frac{C_2}{t\sqrt{\log t}}, \theta^t = \frac{C_3}{t \log t} \mathbf{1}\{t \pmod{C_4} = 0\}$$

# Experiments

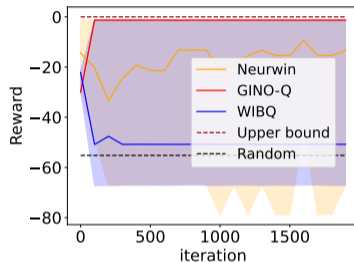
## Group 1: A Non-indexable RMAB



Gain index over training



Reward over training  
( $M = 10, N = 7$ )

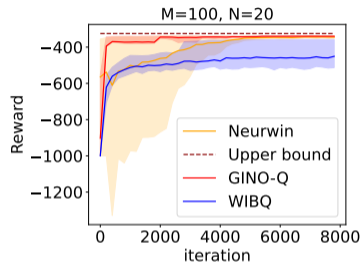
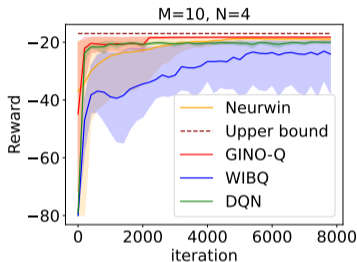
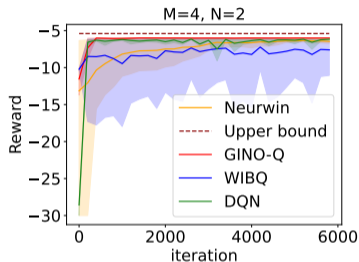


Reward over training  
( $M = 100, N = 70$ )

- GINO-Q indeed learns the gain indices
- While Whittle-index-based methods can perform poorly, GINO-Q always learn a near-optimal policy

# Experiments

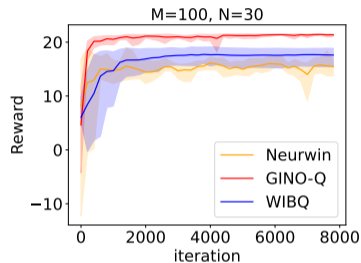
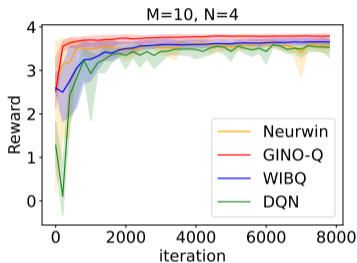
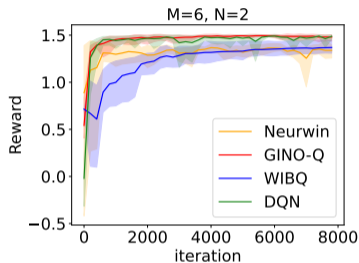
## Group 2: Channel Allocation



- Standard RL methods is infeasible when  $M$  is large
- GINO-Q achieves the best performance and converges fast

# Experiments

## Group 3: Patrol Scheduling



- Even in indexable RMABs, GINO-Q performs better than Whittle-index-based methods

## Conclusion

- RMAB is a widely used model for resource allocation, network scheduling, public health intervention, and more
- Standard RL methods are impractical due to exponential state spaces and combinatorial action spaces
- Most existing approaches rely on the Whittle index policy, which assumes indexability—a condition often unmet in practice
- In non-indexable RMABs, Whittle-index-based methods can perform arbitrarily poorly
- We propose **GINO-Q**, a learning algorithm based on the gain index policy, which does **not** require indexability
- GINO-Q is asymptotically optimal and achieves SOTA performance in non-asymptotical settings

# IMPERIAL

# Thank you

GINO-Q: Learning an Asymptotically Optimal Index Policy for Restless Multi-armed Bandits  
AAAI2026, Jan 2026