

# Towards a Pretrained Model for Restless Bandits via Multi-arm Generalization

Yunfan Zhao<sup>\*1</sup>, Nikhil Behari<sup>\*1</sup>, Edward Hughes<sup>2</sup>, Edwin Zhang<sup>1</sup>, Dheeraj Nagaraj<sup>2</sup>,  
Karl Tuyls<sup>2</sup>, Aparna Taneja<sup>2</sup> and Milind Tambe<sup>1,2</sup>

<sup>1</sup>Harvard University

<sup>2</sup>Google

## Abstract

Restless multi-arm bandits (RMABs) is a class of resource allocation problems with broad application in areas such as healthcare, online advertising, and anti-poaching. We explore several important question such as how to handle arms opting-in and opting-out over time without frequent retraining from scratch, how to deal with continuous state settings with nonlinear reward functions, which appear naturally in practical contexts. We address these questions by developing a pre-trained model (PreFeRMAB) based on a novel combination of three key ideas: (i) to enable fast generalization, we use train agents to learn from each other’s experience; (ii) to accommodate streaming RMABs, we derive a new update rule for a crucial  $\lambda$ -network; (iii) to handle more complex continuous state settings, we design the algorithm to automatically define an abstract state based on raw observation and reward data. PreFeRMAB allows general zero-shot ability on previously unseen RMABs, and can be fine-tuned on specific instances in a more sample-efficient way than retraining from scratch. We theoretically prove the benefits of multi-arm generalization and empirically demonstrate the advantages of our approach on several challenging, real-world inspired problems.

## 1 Introduction

Restless multi-arm bandits (RMABs), a class of resource allocation problems involving multiple agents with a global resource constraint, have found applications in various scenarios, including resource allocation in multi-channel communication, machine maintenance, and healthcare [Hodge and Glazebrook, 2015; Mate *et al.*, 2022]. RMABs have recently been studied from a multi-agent reinforcement learning perspective.

The usual RMAB setting considers a fixed number of arms, each associated with a known, fixed MDP with finite state and action spaces; the RMAB chooses  $K$  of  $N$  arms every round to optimize some long term objective. Even in this

<sup>\*</sup>Equal contribution.

setting, the problem has been shown to be PSPACE hard [Papadimitriou and Tsitsiklis, 1999]. Several approximation algorithms have been proposed in this setting [Whittle, 1988; Hawkins, 2003], particularly when MDP transition probabilities are fully specified, which are successful in practice. State-of-the-art approaches for binary action RMABs commonly provide policies based on the Whittle index [Whittle, 1988], an approach that has also been generalized to multi-action RMABs [Hawkins, 2003; Killian *et al.*, 2021b]. There are also linear programming-based approaches to both binary and multi-action RMABs [Zhang and Frazier, 2021]. Reinforcement learning (RL) based techniques have also been proposed as state-of-the-art solutions for general multi-action RMABs [Xiong and Li, 2023].

In this work, we focus on RL-based methods that provide general solutions to binary and multi-action RMABs, without requiring ground truth transition dynamics, or special properties such as indexability as required by other approaches [Wang *et al.*, 2023]. Unfortunately, several limitations exist in current RMAB solutions, especially for state of the art RL-based solutions, making them challenging or inefficient to deploy in real-world resource allocation problems.

The first limitation arises when dealing with arms that constantly opt-in (also known as streaming RMABs), which happens in public health programs where new patients (arms in RMABs) arrive asynchronously. Existing solutions either require ground truth transition probabilities [Mate *et al.*, 2021], which are often unknown in practice, or else require an entirely new model to be trained repeatedly, which can be extremely computationally costly and sample inefficient.

A second limitation occurs for new programs, or existing programs experiencing a slight change in the user base. In these situations, existing approaches do not provide a pretrained RMAB model that can be immediately deployed. In deep learning, pretrained models are the foundation for contemporary, large-scale image and text networks that generalize well across a variety of tasks [Bommasani *et al.*, 2021]. For real-world problems modeled with RMABs, establishing a similar pretrained model is essential to reduce the burden of training new RMAB policies from scratch.

The third limitation occurs in handling continuous state multi-action RMABs that have important applications [Sinha and Mahajan, 2022; Dusonchet and Hongler, 2003]. Naturally continuous domain state-spaces, such as patient adher-

40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83

ence, are often binned into manually crafted discrete state spaces to improve model tractability and scalability [Mate *et al.*, 2022], resulting in the loss of crucial information about raw observations.

In this work we present PreFeRMAB, a **Pretrained Flexible model for RMABs**. Using multi-arm generalization, PreFeRMAB enables *zero-shot* deployment for unseen arms as well as rapid fine-tuning for specific RMAB settings.

#### Our main contributions are:

- To the best of our knowledge, we are the first to develop a pretrained RMAB model with zero shot ability on entire sets of unseen arms.
- Whereas a general multiagent RL system could suffer from sample complexity exponential in the number of agents  $N$  [Gheshlaghi Azar *et al.*, 2013], we prove PreFeRMAB benefits from larger  $N$ , via multi-arm generalization and better estimation of the population distribution of arm features.
- Our pretrained model can be fine-tuned on specific instances in a more sample-efficient way than training from scratch, requiring less than 12.5% of samples needed for training a previous *multi-action* RMAB model in a healthcare setting [Verma *et al.*, 2023].
- We derive an update rule for a crucial  $\lambda$ -network, allowing changing numbers of arms without retraining. While streaming bandits received considerable attention [Liau *et al.*, 2018], we are the first to handle streaming *multi-action* RMABs with unknown transition dynamics.
- Our model accommodates both discrete and continuous states. To address the continuous state setting, where real-world problems often require nonlinear rewards [Riquelme *et al.*, 2018], we providing a StateShaping module to automatically define an abstract state.

## 2 Related Work

**RMABs with binary and multiple actions.** Solving an RMAB problem, even with known transition dynamics, is known to be PSPACE hard [Papadimitriou and Tsitsiklis, 1999]. For binary action RMABs, [Whittle, 1988] provides an approximate solution, using a Lagrangian relaxation to decouple arms and choose actions by computing so-called Whittle indices of each arm. It has been shown that the Whittle index policy is asymptotically optimal under the indexability condition [Weber and Weiss, 1990; Akbarzadeh and Mahajan, 2019]. The Whittle index was extended to a special class of *multi-action* RMABs with monotonic structure [Hodge and Glazebrook, 2015]. A method for more general *multi-action* RMABs based on Lagrangian relaxation was proposed by [Killian *et al.*, 2021b]. Weakly coupled Markov Decisions Processes (WCMDP), which generalizes multi-action RMABs to have multiple constraints, was studied by [Hawkins, 2003], who proposed a Langrangian decomposition approach. WCMDP was subsequently studied by [Adelman and Mersereau, 2008], who proposed improvements in solution quality at the expense of higher computational costs. While above methods developed for WCMDP require knowledge of ground truth transition dynamics, our

algorithm handles unknown transition dynamics, which is more common in practice [Wang *et al.*, 2023]. Additionally, the above works in *multi-action* settings do not provide algorithms for continuous state RMABs.

**Multi-agent RL and RL for RMABs.** RMABs are a specific instance of the powerful multi-agent RL framework used to model systems with multiple interacting agents in both competitive and co-operative settings [Shapley, 1953; Littman, 1994], for which significant strides have been made empirically [Jaques *et al.*, 2019; Yu *et al.*, 2022] and theoretically [Jin *et al.*, 2021; Xie *et al.*, 2020]. [Nakhleh *et al.*, 2021] proposed a deep RL method to estimate the Whittle index. [Fu *et al.*, 2019] provided an algorithm to learn a Q-function based on Whittle indices, states, and actions. [Avrachenkov and Borkar, 2022] and [Biswas *et al.*, 2021] developed Whittle index-based Q-learning methods with convergence guarantees. While the aforementioned works focus on binary action RMABs, [Killian *et al.*, 2021a] generalized this to multi-action RMABs using tabular Q-learning. A subsequent work [Killian *et al.*, 2022], which focussed on robustness against adversarial distributions, took a deep RL approach that was more scalable. However, existing works on multi-action RMABs do not consider streaming RMABs and require training from scratch when a new arm opts-in. Additionally, works built on tabular Q-learning [Fu *et al.*, 2019; Avrachenkov and Borkar, 2022; Biswas *et al.*, 2021; Killian *et al.*, 2021a] may not generalize to continuous state RMABs without significant modifications. Our pretrained model addresses these limitations, and enables zero-shot ability on a wide range of unseen RMABs.

**Streaming algorithms.** The streaming model, pioneered by [Alon *et al.*, 1996], considers a scenario where data arrives online and the amount of memory is limited. The model is adapted to multi-arm bandits (MAB), assuming that arms arrive in a stream and the number of arms that can be stored is limited [Liau *et al.*, 2018; Chaudhuri and Kalyanakrishnan, 2020]. The streaming model has recently been adapted to *binary action* RMABs with *known* transition probabilities [Mate *et al.*, 2021], but not studied in the more general and practical settings of *multi-action* RMABs with *unknown* transition dynamics. We aim to close this gap.

**Zero-shot generalization and fine-tuning.** Foundation models that have a strong ability to generalize to new tasks in zero shot and efficiently adapt to new tasks via fine-tuning have received great research attention [Bommasani *et al.*, 2021]. Such models are typically trained on vast data, such as internet-scale text [Devlin *et al.*, 2018] or images [Ramesh *et al.*, 2021]. RL has seen success in the direction of foundation models for decision making, using simulated [Team *et al.*, 2023] and real-world [Yu *et al.*, 2020] environments. A pretrained model for RMABs is needed [Zhao *et al.*, 2024]. To our knowledge, we are the first to realize zero-shot generalization and efficient fine-tuning in the setting of RMABs.

## 3 Problem Statement

We study multi-action RMABs with system capacity  $N$ , where existing arms have the option to opt-out (that is, the state-action-rewards corresponding to them are disregarded

197 by the model post opt-out), and new, unseen arms can request  
 198 to opt-in (that is, these arms are considered only post the opt-  
 199 in time). Such requests will be accepted if and only if the  
 200 system capacity permits. A vector  $\xi_t \in \{0, 1\}^N$  represents  
 201 the opt-in decisions:

$$\xi_{i,t} = \begin{cases} 1 & \text{if arm } i \text{ opts-in at round } t, \\ 0 & \text{otherwise.} \end{cases}$$

202 Notice that existing arms must opt-in in each round  $t$  to  
 203 remain in the system. For each arm  $i \in [N]$ , the state space  
 204  $\mathcal{S}_i$  can be either discrete or continuous, and the action space  
 205  $\mathcal{A}_i$  is a finite set of discrete actions. Each action  $a \in \mathcal{A}_i$   
 206 has an associated cost  $C_i(a)$ , with  $C_i(0)$  denoting a no-cost  
 207 passive action. The reward at a state is given by a function  
 208  $R_i : \mathcal{S}_i \rightarrow \mathbb{R}$ . We let  $\beta \in [0, 1)$  denote a discount factor. Each  
 209 arm has a unique feature vector  $\mathbf{z}_i \in \mathbb{R}^m$  that provides useful  
 210 information about the arm. Notice our model directly utilizes  
 211 feature information in its policy network, without requiring  
 212 intermediate steps to extract transition dynamics information  
 213 from features.

214 When the state space is discrete, each arm  $i \in [N]$  follows  
 215 a Markov Decision Process  $(\mathcal{S}_i, \mathcal{A}_i, C_i, T_i, R_i, \beta, \mathbf{z}_i)$ , where  
 216  $T_i : \mathcal{S}_i \times \mathcal{A}_i \times \mathcal{S}_i \rightarrow [0, 1]$  is a transition matrix representing  
 217 the probability of transitioning from the current state to the  
 218 next state given an action. In contrast, when the state space  
 219 is continuous, each arm  $i \in [N]$  follows a Markov Decision  
 220 Process  $(\mathcal{S}_i, \mathcal{A}_i, C_i, \Gamma_i, R_i, \beta, \mathbf{z}_i)$ , where  $\Gamma_i$  is a set of param-  
 221 eters encoding the transition dynamics. For example, in the  
 222 case that the next state moves according to a Gaussian distri-  
 223 bution,  $\Gamma_i$  may denote the mean and variance of the Gaussian.

224 For simplicity, we assume that  $\mathcal{S}_i, \mathcal{A}_i, C_i$ , and  $R_i$  are  
 225 the same for all arms  $i \in [N]$  and omit the subscript  $i$ . Note that  
 226 our algorithms can also be used in the general case where  
 227 rewards and action costs are different across arms. For ease  
 228 of notation, we let  $\mathbf{s} \in \mathbb{R}^N$  denote the state over all arms,  
 229 and we let  $\mathbf{A} \in \{0, 1\}^{N \times |\mathcal{A}|}$  denote one-hot-encoding of the  
 230 actions taken over all arms. The agent learns a policy  $\pi$  that  
 231 maps states  $\mathbf{s}$  and features  $\mathbf{z}$  to actions  $\mathbf{A}$ , while satisfying a  
 232 constraint that the sum cost of actions taken is no greater than  
 233 a given budget  $B$  in every timestep  $t \in [H]$ , where  $H$  is the  
 234 length of the horizon.

**Our goal is to learn an RMAB policy that maximizes the following Bellman equation** The key difficulty in learning such a policy is how to utilize features  $\mathbf{z}$  and address opt-in decisions  $\xi$ . These are important research questions not addressed in previous works [Killian *et al.*, 2022].

$$J(\mathbf{s}, \mathbf{z}, \xi) = \max_{\mathbf{A}} \left\{ \sum_{i=1}^N R(\mathbf{s}_i) + \beta \mathbb{E}[J(\mathbf{s}', \mathbf{z}, \xi) \mid \mathbf{s}, \mathbf{A}] \right\}, \quad (1)$$

$$\text{s.t. } \sum_{i=1}^N \sum_{j=1}^{|\mathcal{A}|} \mathbf{A}_{ij} c_j \leq B \quad \text{and} \quad \sum_{j=1}^{|\mathcal{A}|} \mathbf{A}_{ij} = 1 \quad \forall i \in [N],$$

235 where  $c_j \in \mathcal{C}$  is the cost of  $j^{\text{th}}$  action, and  $\mathbf{A}_{ij} = 1$  if action  
 236  $j$  is chosen on arm  $i$  and  $\mathbf{A}_{ij} = 0$  otherwise. Further, we  
 237 assume that the rewards  $R$  are uniformly bounded by  $R_{\max}$ .

## 4 Generalized Model for RMABs

We first provide an overview of key ideas and then discuss each of the ideas in more detail. (See Figure 3 in Appendix for an overview of the training procedure.)

### 4.1 Key Algorithmic Ideas

Several key algorithmic novelties are necessary for our model to address limitations of existing works:

**A pretrained model via multi-arm generalization:** We train agents to learn from each others' experience. Whereas a general multiagent RL system could suffer from sample complexity exponential in the number of arms  $N$  [Gheshlaghi Azar *et al.*, 2013], we prove PreFeRMAB benefits from a larger  $N$ , via generalization across arms.

**A novel  $\lambda$ -network updating rule for opt-in:** The opt-in and opt-out of arms induce a more complex form of the Lagrangian and add randomness to actions taken by agents. We provide a new  $\lambda$ -network update rule and train PreFeRMAB with opt-in and opt-out of arms, to enable zero-shot performance across various opt-in rates and accommodate streaming RMABs.

**Handling continuous states with StateShaping subrou-**  
**tine:** In the continuous state setting, real-world problems often require nonlinear rewards [Riquelme *et al.*, 2018], and naively using raw observations to train models may result in poor performance (see Figure 5). To tackle this challenge, we design the algorithm to automatically define an abstract state based on raw observation and reward data.

### 4.2 A Pretrained Model via Multi-arm Generalization

To enable multi-arm generalization, we introduce feature-based Q-values, together with a Lagrangian relaxation with features  $\mathbf{z}_i$  and opt-in decisions  $\xi_i$ :

$$\begin{aligned} & J(\mathbf{s}, \mathbf{z}, \xi, \lambda^*) \\ &= \min_{\lambda \geq 0} \left( \frac{\lambda B}{1 - \beta} + \sum_{i=1}^N \max_{a_i \in |\mathcal{A}|} \{Q(\mathbf{s}_i, a_i, \mathbf{z}_i, \xi_i, \lambda)\} \right), \quad (2) \\ & \text{s.t. } Q(\mathbf{s}_i, a_i, \mathbf{z}_i, \xi_i, \lambda) \\ &= \xi_i R(\mathbf{s}_i) - \xi_i \lambda c_{a_i} + \beta \mathbb{E}[Q(\mathbf{s}'_i, a_i, \mathbf{z}_i, \xi_i, \lambda) \mid \pi(\lambda)]. \end{aligned}$$

where  $Q$  is the Q-function,  $a_i$  is the action of arm  $i$ ,  $\mathbf{s}'_i$  is the state transitioned to from  $\mathbf{s}_i$  under action  $a_i$ , and  $\pi(\lambda)$  is the optimal policy under a given  $\lambda$ . Notice that this relaxation decouples the Q-functions of the arms, and therefore  $Q_i$  can be solved independently for a given  $\lambda$ .

Now we discuss how we use feature-based Q-values and how agents could learn from each other. During pretraining, having received arms' opt-in and out decisions (line 5), Algorithm 1 samples an action-charge  $\lambda$  based on updated opt-in decisions  $\xi$  and features  $\mathbf{z}_i$  (line 6). Next, from opt-in arms we collect trajectories (lines 7-14), which are later used to train a single pair of actor/critic networks for all arms, allowing the policy for one arm to benefit from other arms' data. After that, we update the policy network  $\theta$  and the critic network  $\phi$  (Line 16), using feature-based Q-values to compute advantage estimates for the actor in PPO update. Critically, feature-based Q-values updated with one arm's data,

284 improves the policy for other arms. In real-world problems  
 285 with missing feature entries or less informative features, it is  
 286 more important for agents to learn from each other (see Table  
 287 1 in Sec 5.2). Intuitively, if a model only learns from  
 288 homogeneous arms, then we should expect this model to per-  
 289 form poorly when used out-of-the-box on arms with com-  
 290 pletely different behaviors.

---

**Algorithm 1** PreFeRMAB (Training)

---

1: Input:  $n\_epochs$ ,  $n\_steps$ ,  $\lambda$ -update frequency  $K \in \mathbb{N}^+$ ,  
 and system capacity  $N$   
 2: Initialize actor  $\theta$ , critic  $\phi$ ,  $\lambda$ -network  $\Lambda$ , buffer = [], state  
 $\mathbf{s} \in \mathbb{R}^N$ , and features  $\mathbf{z}_i \in \mathbb{R}^m$   
 3: Initialize *StateShaping*, and set  $\bar{\mathbf{s}} \leftarrow \text{StateShaping}(\mathbf{s})$   
 4: **for** epoch = 1, 2, . . . ,  $n\_epochs$  **do**  
 5:     Receive opt-in/out requests and update  $\xi$  and  $\mathbf{z}_i$ .  
 6:     Compute  $\lambda = \Lambda(\bar{\mathbf{s}}, \{z_i\}_{i=1}^N, \xi)$   
 7:     **for** timestep  $t = 1, \dots, n\_steps$  **do**  
 8:         **for** Arm  $i = 1, \dots, N$  **do**  
 9:             **if** Arm  $i$  is opt-in (i.e.  $\xi_i = 1$ ) **then**  
 10:                 Sample an action  $a_i \sim \theta(\bar{s}_i, \lambda, \mathbf{z}_i)$   
 11:                  $s'_i, r_i = \text{Simulate}(s_i, a_i)$   
 12:                  $s'_i = \text{StateShaping}(s'_i)$   
 13:                 Add tuple  $(s_i, \bar{s}_i, a_i, r_i, \bar{s}'_i, \mathbf{z}_i)$  to buffer  
 14:                  $s_i \leftarrow s'_i, \bar{s}_i \leftarrow \bar{s}'_i$   
 15:             Add tuple  $(\lambda, \xi)$  to buffer  
 16:             Update the  $(\theta, \phi)$  pair using buffer.  
 17:             **if** epoch //  $K = 0$  **then**  
 18:                 Update  $\Lambda$  via Prop 2 using trajectories in buffer  
 19:                 Update  $\hat{r}(\cdot)$  in *StateShaping* using  $(s, r)$ -tuples in  
 buffer

---

291 We will now give a theoretical guarantee of multi-arm gen-  
 292 eralization by considering the following simplified setting  
 293 and assumptions where we do not consider opt-in and opt-  
 294 out. That is, in each epoch we draw a new set of  $N$  arms.  
 295 We let the distribution of arm features (i.e.  $\mathbf{z}_i$ ) to be denoted  
 296 by the probability measure  $\mu^*$ . Each ‘sample’ for our policy  
 297 network training consists of  $N$  features corresponding to  $N$   
 298 arms  $(\mathbf{z}_1, \dots, \mathbf{z}_N)$ , drawn i.i.d. from the distribution  $\mu^*$ . Call  
 299 the empirical distribution of  $(\mathbf{z}_i)$  to be  $\hat{\mu}$ . During training, we  
 300 receive  $n\_epochs$  i.i.d. draws of  $N$  arm features each, denoted  
 301 by  $\hat{\mu}_1, \dots, \hat{\mu}_{n\_epochs}$

302 Let  $\Theta$  denote the space of neural network weights of the  
 303 policy network (for clarity, we shorten the  $(\theta, \phi)$  in Algo-  
 304 rithm 1 to  $\theta$ ). The neural network inputs are Lagrangian mul-  
 305 tiplier  $\lambda$ , state of an arm  $s$ , its feature  $\mathbf{z}$  and the output is  
 306  $a \in \mathcal{A}$ . Let  $V(\mathbf{s}, \theta, \lambda, \hat{\mu})$  denote the discounted reward, av-  
 307 eraged over  $N$  arms with features  $\hat{\mu}$  obtained with the neural  
 308 network with parameter  $\theta$ , starting from the state  $\mathbf{s}$  (cumu-  
 309 lative state of all arms). The proposition below shows the  
 310 generalization properties of the output of Algorithm 1. The  
 311 proof and a detailed discussion of the assumptions and con-  
 312 sequences are given in Section D.

313 **Proposition 1.** *Suppose the following assumptions hold:*

314 1. *Algorithm 1 learns neural network weights  $\hat{\theta} \in \Theta$ ,*  
 315 *whose policy is optimal for each  $(\hat{\mu}_i, \lambda)$  for  $1 \leq i \leq$*

- $n\_epochs$  and  $\lambda \in [0, \lambda_{\max}]$  316
2. *There exists  $\theta^* \in \Theta$  which is optimal for every instance* 317  
 $(\hat{\mu}, \lambda)$ . 318
  3.  $\Theta = \mathcal{B}_2(D, \mathbb{R}^d)$ , *the  $\ell_2$  ball of radius  $D$  in  $\mathbb{R}^d$ .* 319
  4.  $|V(\mathbf{s}, \theta_1, \lambda, \hat{\mu}) - V(\mathbf{s}, \theta_2, \lambda, \hat{\mu})| \leq L\|\theta_1 - \theta_2\|$  *and* 320  
 $|V(\mathbf{s}, \theta, \lambda_1, \hat{\mu}) - V(\mathbf{s}, \theta, \lambda_2, \hat{\mu})| \leq L|\lambda_1 - \lambda_2|$  *for all* 321  
 $\theta_1, \theta_2, \theta \in \Theta$  *and  $\lambda_1, \lambda_2, \lambda \in [0, \lambda_{\max}]$ .* 322

*Then, the generalization error over unseen arms  $(\hat{\mu})$  sat-  
 isfies:*

$$\mathbb{E}_{\hat{\mu}, \hat{\theta}} \left[ \inf_{\lambda \in [0, \lambda_{\max}]} V(\mathbf{s}, \hat{\theta}, \lambda, \hat{\mu}) \right] \geq \mathbb{E}_{\hat{\mu}} \left[ \inf_{\lambda \in [0, \lambda_{\max}]} V(\mathbf{s}, \theta^*, \lambda, \hat{\mu}) \right] - \tilde{O} \left( \frac{1}{\sqrt{n\_epochs N}} \right) \quad (3)$$

*Here,  $\tilde{O}$  hides polylogarithmic factors in  $n\_epochs, N$  and con-  
 stants depending on  $d, D, L, \beta, \frac{B}{N}, c_j, R_{\max}$  and  $\lambda_{\max}$*  323 324

The assumption of existence of  $\theta^*$  is reasonable: This  
 means that there exists a neural network which gives the op-  
 timal policy for a family of single-arm MDPs indexed by  
 $(\mathbf{z}, \lambda)$ . Proposition 1 shows that when  $n\_epochs$  and  $N$  are  
 large, the Lagrangian relaxed value function of the learned  
 network is close to that of the optimal network. 325 326 327 328 329 330

*An important insight is that the generalization ability of the  
 PreFeRMAB network becomes better as the number of arms  
 per instance becomes larger.* This is counter intuitive since  
 a system with a larger number of agents are generally very  
 complex. Jointly, the arms form an MDP with  $|\mathcal{S}|^N$  states  
 and  $|\mathcal{A}|^N$  actions. General multi-agent RL problems with  $N$   
 arms thus can suffer from an exponential dependence on  $N$  in  
 their sample complexity for learning (see sample complexity  
 lower bounds in [Gheshlaghi Azar *et al.*, 2013]). However,  
 due to the structure of RMABs and the Lagrangian relaxation,  
 we achieve a better generalization with a larger  $N$ . Our proof  
 in the appendix shows that this is due to the fact that a larger  
 number of arms helps estimate the population distribution  $\mu^*$   
 of the arm features better. We show in Table 1 that indeed  
 having more number of arms helps the PreFeRMAB network  
 generalize better over unseen instances. 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346

### 4.3 A Novel $\lambda$ -network Updating Rule 347

In real-world health programs, we may observe new patients  
 constantly opt-in [Mate *et al.*, 2021]. The opt-in / opt-out  
 decisions render the updating rule in [Killian *et al.*, 2022] un-  
 usable and add additional randomness to actions taken by the  
 agent. To overcome this challenge and to stabilize training,  
 we develop a new  $\lambda$ -network updating rule. 348 349 350 351 352 353

**Proposition 2.** [ $\lambda$ -network updating rule] *The equation for  
 gradient descent for the objective (Eq 2) with respect to  $\lambda$ ,  
 with step size  $\alpha$  is:*

$$\Lambda_t = \Lambda_{t-1} - \alpha \left( \frac{B}{1 - \beta} \right) - \alpha \left( \sum_{i=1}^N \mathbb{E} \left[ \sum_{t=0}^H \xi_{i,t} \beta^t c_{i,t} + (1 - \xi_{i,t}) \beta^t c_{0,t} \right] \right),$$

*where  $c_{i,t}$  is the cost of the action taken by the optimal policy  
 on arm  $i$  in round  $t$ .* 354 355

Critically, this update rule allows PreFeRMAB to handle streaming RMABs, accommodating a changing number of arms without retraining and achieving strong zero-shot performance across various opt-in rates (see Table 3 and 14). Having established an updating rule, we provide a convergence guarantee. The proofs are relegated to Appendix E.

**Proposition 3** (Convergence of  $\lambda$ -network). *Suppose the arm policies converge to the optimal Q-function for a given  $\Lambda_t$ , then the update rule (in Prop 2) for the  $\lambda$ -network converges to the optimal as the number of training epochs and the number of actions collected in each epoch go to infinity.*

#### 4.4 Handling Continuous States with StateShaping

Real-world problems may require continuous states with nonlinear rewards [Riquelme *et al.*, 2018]. Existing RMAB algorithms either use a human-crafted discretization or fail to address challenging nonlinear rewards [Killian *et al.*, 2022]. Discretization may result in loss of information and fail to generalize to different population sizes. For example, the popular SIS epidemic model [Yaesoubi and Cohen, 2011] is expected to scale to a continuum limit as the population size increases to infinity, and a continuous state-space model can better handle scaling by using proportions instead of absolute numbers. Under nonlinear rewards, naively using raw observations in training may result in poor performance (see Figure 5). We provide a StateShaping module to improve model stability and performance.

---

#### Algorithm 2 StateShaping Subroutine

---

- 1: Input: estimator  $\in \{\text{Isotonic Regression, KNN}\}$ , states  $\mathbf{s} \in \mathbb{R}^N$ , data  $\mathcal{D}$  of  $(s, r)$  tuples
  - 2: Output  $\bar{\mathbf{s}} = \mathbf{s}$  if no normalization is desired
  - 3: Compute
 
$$r_{\min} = \min_{s': s' \in \mathcal{D}} r(s'), \quad r_{\max} = \max_{s': s' \in \mathcal{D}} r(s')$$

$$s_{\min} = \min_{s' \in \mathcal{D}} s', \quad s_{\max} = \max_{s' \in \mathcal{D}} s'$$
  - 4: Compute  $\hat{r}(s_i)$  using the choice of Estimator.
  - 5: Output  $\bar{\mathbf{s}}$ , where  $\bar{s}_i = \frac{\hat{r}(s_i) - r_{\min}}{r_{\max} - r_{\min}} (s_{\max} - s_{\min}), \forall i$
- 

In Algorithm 2, users can choose whether to obtain abstract state [Abel *et al.*, 2018] from raw observations (lines 2). We compute ranges of reward and raw observations, and obtain an reward estimate (lines 3-4). After that, we automatically refine the raw observation such that reward is a linear function of the abstract state (line 4), improving model stability for challenging reward functions. Here a key assumption is that reward is an increasing function of the raw observation, which is common in RMABs [Killian *et al.*, 2022]. Notice as we collect more observations, the accuracy of the reward estimator  $\hat{r}(\cdot)$  will improve (it is updated in line 24 of Algorithm 1).

StateShaping instantiates the idea of state abstraction, which is shown to improve generalizability and robustness [Li *et al.*, 2006], in the RMAB context for continuous states. Applying Theorem 3 in [Li *et al.*, 2006] to the Lagrangian relaxation (Eq. 2), we have that an optimal policy learnt using

the abstract state space is guaranteed to be also optimal in the ground MDP (defined by raw observations).

#### 4.5 Inference using Pretrained Model

An important difference between training and inference is that during inference time, we strictly enforce the budget constraint on the trained model, by greedily selecting highest probability actions until the budget is reached. The rest of the inference components are similar to the training component.

---

#### Algorithm 3 PreFeRMAB (Inference)

---

- 1: Input: : States  $\mathbf{s}$ , costs  $C$ , budget  $B$ , features  $\mathbf{z}_i \in \mathbb{R}^m$ , opt-in decisions  $\boldsymbol{\xi}$ , agent actor  $\theta$ ,  $\lambda$ -network, *StateShaping* routine with trained estimator  $\hat{r}(\cdot)$
  - 2: Compute  $\lambda = \Lambda(\bar{\mathbf{s}}, \{\mathbf{z}_i\}_{i=1}^N, \boldsymbol{\xi})$
  - 3: **for** Arm  $i = 1, \dots, N$  **do**
  - 4:   **if** Arm  $i$  is opt-in (i.e.  $\xi_i = 1$ ) **then**
  - 5:      $\bar{s}_i = \text{StateShaping}(s_i)$
  - 6:     Compute  $p_i \sim \theta(\bar{s}_i, \lambda, \mathbf{z}_i)$
  - 7:  $\mathbf{a} = \text{GreedyProba}(p, C, B)$    ▷ Greedily select highest probability actions until budget B is reached
- 

### 5 Experimental Evaluation

We provide experimental evaluations of our model in *three separate domains*, including a synthetic setting, an epidemic modeling setting, as well as a maternal healthcare intervention setting. We first describe these three experimental domains. Then, we provide results for PreFeRMAB in a *zero-shot evaluation setting*, demonstrating the performance of our model on *new, unseen test arms* drawn from distributions distinct from those in training. Here, we demonstrate the flexibility of PreFeRMAB, including strong performance across domains, state representations (discrete vs. continuous), and over various challenging reward functions. Finally, we demonstrate the strength of using PreFeRMAB as a *pre-trained model*, enabling faster convergence for fine-tuning on a specific set of evaluation arms.

In Appendix B, we provide **ablation studies** over (1) a wider range of opt-in rates (2) different feature mappings (3) DDLPO topline with and without features (4) more problem settings. For hyperparameter details, we refer to Appendix A.

#### 5.1 Experimental Settings

**Features:** In all experiments, we generate features by projecting parameters that describe the ground truth transition dynamics into features using randomly generated projection matrices. The dimension of feature equals the number of parameters required to describe the transition dynamics. In Appendix B, we provide results on different feature mappings.

**Synthetic:** Following [Killian *et al.*, 2022], we consider a synthetic dataset with binary states and binary actions. The transition probabilities for each arm  $i$  are represented by matrices  $T_{s=0}^{(i)}$  and  $T_{s=1}^{(i)}$  for arm  $i$  at states 0 and 1 respectively:

$$T_{s=0}^{(i)} = \begin{bmatrix} p_{00} & 1 - p_{00} \\ p_{01} & 1 - p_{01} \end{bmatrix}, \quad T_{s=1}^{(i)} = \begin{bmatrix} p_{10} & 1 - p_{10} \\ p_{11} & 1 - p_{11} \end{bmatrix}$$

Each  $p_{jk}$  corresponds to the probability of transitioning from state  $j$  to state 0 when action  $k$  is taken. These values are sampled uniformly from the intervals:

$$p_{00} \in [0.4, 0.6], p_{01} \in [0.4, 0.6], p_{10} \in [0.8, 1], p_{11} \in [0.0, 1]$$

**SIS Epidemic Model:** Inspired by the vast literature on agent-based epidemic modeling, we adapt the SIS model given in [Yaesoubi and Cohen, 2011], following a similar experiment setup as described in [Killian *et al.*, 2022]. Arms  $p$  represent a subpopulation in distinct geographic regions; states  $s$  are the number of uninfected people within each arm’s total population  $N_p$ ; the number of possible states is  $S$ . Transmission within each arm is guided by parameters:  $\kappa$ , the average number of contacts within the arm’s subpopulation in each round, and  $r_{infect}$ , the probability of becoming infected after contact with an infected person.

In this setting, there is a budget constraint over interventions. There are three available intervention actions  $a_0, a_1, a_2$  that affect the transmission parameters:  $a_0$  represents no action;  $a_1$  represents messaging about physical distancing;  $a_2$  represents the distribution of face masks. We discuss additional details in Appendix A.

**ARMMAN:** Similar to the set up in [Biswas *et al.*, 2021; Killian *et al.*, 2022], we model the real world maternal health problem as a discrete state RMAB. We aim to encourage engagement with automated health information messaging. There are three possible states, presenting self-motivated, persuadable, and lost cause. The actions are binary. There are 6 uncertain parameters per arm, sampled from uncertainty intervals of 0.5 centered around the transition parameters that align with summary statistics given in [Biswas *et al.*, 2021].

**Continuous State Modeling:** Continuous state restless bandits have important applications [Lefèvre, 1981; Sinha and Mahajan, 2022; Dusonchet and Hongler, 2003]. By not explicitly having a switch in the model (switching between discrete and continuous state space), we enable greater model flexibility. To demonstrate this, we consider both a Continuous Synthetic and a Continuous SIS modeling setting. We provide details of these settings in Appendix A.3.

We present **additional details**, including hyperparameters and StateShaping illustration in Appendix A.

## 5.2 PreFeRMAB Zero-Shot Learning

We first consider three challenging datasets in the discrete state space. After that, we present results on datasets with continuous state spaces with more complex reward functions and transition dynamics.

**Pretraining.** For each pretraining iteration, we sample from a binomial with mean 0.8 to determine which arms will be opted-in given system capacity  $N$ . For new arms, we sample new transition dynamics to allow the model to see a wider range of arm features.

**Evaluation.** We compare PreFeRMAB to *Random Action* and *No Action* baselines. In every table in this subsection, we present the *reward per arm* averaged over 50 trials, on *new, unseen arms* arm sampled from the testing distribution.

**Multi-arm Generalization:** Table 1 on Synthetic illustrates that PreFeRMAB, learning from *multi-arm generalization*, achieves stronger performance when the number of

System capacity $N = 21$ . Budget $B = 7$ .			
# Unique training arms	45	33	21
No Action	2.88 $\pm$ 0.17	2.88 $\pm$ 0.17	2.88 $\pm$ 0.17
Random Action	3.25 $\pm$ 0.22	3.25 $\pm$ 0.22	3.25 $\pm$ 0.22
PreFeRMAB (2/4 Feats. Masked)	3.81 $\pm$ 0.23	3.79 $\pm$ 0.22	3.59 $\pm$ 0.21
PreFeRMAB (1/4 Feats. Masked)	3.92 $\pm$ 0.24	3.70 $\pm$ 0.21	3.58 $\pm$ 0.20
<b>PreFeRMAB (0/4 Feats. Masked)</b>	4.02 $\pm$ 0.26	3.80 $\pm$ 0.22	3.78 $\pm$ 0.21

Table 1: Multi-arm generalization results on Synthetic (opt-in 100%). With the same total amount of data, PreFeRMAB achieves stronger performance when pretrained on more unique arms, especially when input arm features are masked.

unique arms (i.e. arms with unique features) seen during pre-training increases. Additionally, in practice arm features may be missing or not always reliable, such as in real-world ARMMAN data [Mate *et al.*, 2022]. Our results demonstrate that when features are masked, arms could learn from similar arms’ experience.

Wasserstein Distance	0.05	0.10	0.15	0.20	0.25
System capacity $N = 48$ . Budget $B = 16$ .					
No Action	3.07 $\pm$ 0.10	2.89 $\pm$ 0.08	2.68 $\pm$ 0.07	2.49 $\pm$ 0.09	2.35 $\pm$ 0.07
Random Action	3.49 $\pm$ 0.09	3.25 $\pm$ 0.09	2.99 $\pm$ 0.16	2.80 $\pm$ 0.17	2.57 $\pm$ 0.17
<b>PreFeRMAB</b>	4.50 $\pm$ 0.09	4.30 $\pm$ 0.10	3.81 $\pm$ 0.17	3.79 $\pm$ 0.18	3.46 $\pm$ 0.12
System capacity $N = 96$ . Budget $B = 32$ .					
No Action	3.09 $\pm$ 0.08	2.88 $\pm$ 0.04	2.74 $\pm$ 0.05	2.62 $\pm$ 0.06	2.49 $\pm$ 0.06
Random Action	3.44 $\pm$ 0.14	3.23 $\pm$ 0.09	3.05 $\pm$ 0.09	2.90 $\pm$ 0.10	2.70 $\pm$ 0.11
<b>PreFeRMAB</b>	4.44 $\pm$ 0.13	4.26 $\pm$ 0.13	4.12 $\pm$ 0.13	3.97 $\pm$ 0.16	3.75 $\pm$ 0.12

Table 2: Results on Synthetic (opt-in 100%). For each system capacity, we pretrain a model and present zero-shot results under different amounts of distributional shift.

Number of arms System capacity	80%	85%	90%	95%	100%
Parameters $a_1^{eff}, a_1^{eff}$ are uniformly sampled from [2, 8].					
No Action	5.23 $\pm$ 0.17	5.27 $\pm$ 0.16	5.28 $\pm$ 0.16	5.26 $\pm$ 0.14	5.28 $\pm$ 0.13
Random Action	6.94 $\pm$ 0.15	7.00 $\pm$ 0.16	7.03 $\pm$ 0.15	6.97 $\pm$ 0.14	6.99 $\pm$ 0.12
<b>PreFeRMAB</b>	7.64 $\pm$ 0.27	7.75 $\pm$ 0.25	7.96 $\pm$ 0.18	7.80 $\pm$ 0.16	7.82 $\pm$ 0.11
Parameters $a_1^{eff}, a_1^{eff}$ are uniformly sampled from [3, 9].					
No Action	5.29 $\pm$ 0.16	5.30 $\pm$ 0.17	5.29 $\pm$ 0.15	5.26 $\pm$ 0.14	5.28 $\pm$ 0.13
Random Action	7.21 $\pm$ 0.15	7.28 $\pm$ 0.18	7.26 $\pm$ 0.15	7.22 $\pm$ 0.13	7.22 $\pm$ 0.12
<b>PreFeRMAB</b>	7.77 $\pm$ 0.29	7.87 $\pm$ 0.28	7.90 $\pm$ 0.22	7.95 $\pm$ 0.16	7.95 $\pm$ 0.11

Table 3: Results on SIS ( $N = 20, B = 16, S = 150$ ). We pretrain a model and present zero-shot results on various distributions. During training,  $a_1^{eff}, a_1^{eff}$  are uniformly sampled from [1, 7].

**Discrete State Settings with Different Distributional Shifts:** Results on Synthetic (Table 2) shows PreFeRMAB consistently outperforms under varying amounts of distributional shift, measured in Wasserstein distance. Results on SIS (Table 3) shows PreFeRMAB performs well in settings with large state space  $S = 150$  and multiple actions, under various testing distributions and opt-in rates. Results on ARMMAN (Table 4) shows PreFeRMAB could handle more *challenging* settings that mimics the scenario of a real-world non-profit organization using RMABs to allocate resources.

**Continuous State Settings:** Our results (Figure 1) show

Number of arms System capacity	80%	85%	90%	95%	100%
40% motivated, 20% persuadable, and 40% lost cause.					
No Action	2.39±0.30	2.32±0.28	2.16±0.25	2.26±0.28	2.25±0.30
Random Action	3.04±0.40	3.06±0.38	3.00±0.36	3.14±0.36	3.24±0.32
<b>PreFeRMAB</b>	5.47±0.41	5.00±0.37	4.95±0.29	5.34±0.27	5.03±0.37
40% motivated, 40% persuadable, and 20% lost cause.					
No Action	2.07±0.29	2.19±0.28	2.19±0.30	2.05±0.24	2.17±0.28
Random Action	3.04±0.37	3.02±0.33	2.99±0.30	3.12±0.31	3.15±0.29
<b>PreFeRMAB</b>	5.06±0.36	4.81±0.35	5.13±0.34	5.01±0.26	5.00±0.28

Table 4: Results on ARMMAN ( $N = 25, B = 7, S = 3$ ). We pretrain a model and present zero-shot results on various testing distributions. During training, the proportion of self-motivated, persuadable, and lost cause arms are 20%, 20%, and 60% respectively.

510 that StateShaping is crucial in handling continuous states,  
 511 where the reward function can be more challenging. We pro-  
 512 vide additional evaluations in Table 5, showing PreFeRMAB  
 513 outperforms in complex transition dynamics. More details  
 514 are provided in Appendix A.

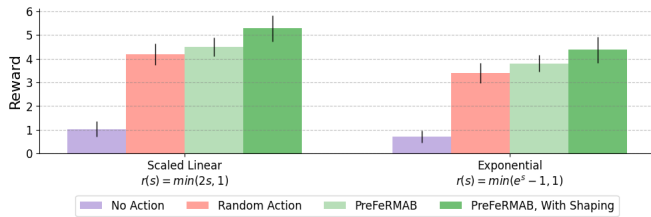


Figure 1: Results for *Continuous Synthetic* domain ( $N=21, B=7.0$ ) with challenging rewards  $r(s)$ .

Number of arms System capacity	80%	85%	90%	95%	100%
Continuous Synthetic ( $N=21, B=7.0, S=2$ )					
No Action	0.70±0.25	0.71±0.25	0.70±0.23	0.70±0.23	0.66±0.21
Random Action	3.44±0.48	3.43±0.45	3.45±0.41	3.37±0.41	3.20±0.38
<b>PreFeRMAB</b>	3.94±0.31	3.76±0.33	4.01±0.31	4.02±0.29	3.67±0.27
Continuous SIS Model ( $N=20, B=16$ )					
No Action	5.64±0.25	5.68±0.19	5.57±0.21	5.48±0.18	5.62±0.17
Random Action	7.11±0.22	7.31±0.22	7.23±0.22	7.24±0.21	7.18±0.17
<b>PreFeRMAB</b>	7.91±0.17	8.08±0.11	7.95±0.14	7.98±0.13	7.82±0.12

Table 5: Results on continuous states. For each problem instance, we pretrain a model.

515 **Comparison with an Additional Baseline:** DDLPO  
 516 [Killian *et al.*, 2022] could not handle distributional shifts  
 517 or various opt-in rates, the more challenging settings that  
 518 PreFeRMAB is designed for. Nevertheless, we provide  
 519 comparisons with DDLPO in settings with no distributional  
 520 shift (Table 6, see also Appendix B.4). Notably, PreFe-  
 521 MAB zero-shot performance on *unseen* arms is near that of  
 522 DDLPO, which is trained and tested on the same set of arms.

### 5.3 PreFeRMAB Fast Fine-Tuning

524 Having shown the zero-shot results of PreFeRMAB, we now  
 525 demonstrate finetuning capabilities of the pretrained model.  
 526 In Figure 2, we compare the number of samples required

Number of arms System capacity	80%	85%	90%	95%	100%
Synthetic with $N = 96, B = 32, S = 2$ .					
No Action	3.22±0.12	3.24±0.12	3.19±0.11	3.18±0.11	3.18±0.11
Random Action	3.62±0.13	3.66±0.13	3.58±0.13	3.60±0.12	3.60±0.12
<b>PreFeRMAB</b>	4.63±0.12	4.71±0.12	4.53±0.13	4.47±0.12	4.61±0.10
DDLPO (topline)	n/a	n/a	n/a	n/a	4.58±0.13
SIS with $N = 20, B = 16, S = 150$ .					
No Action	5.33±0.16	5.30±0.15	5.31±0.14	5.29±0.13	5.28±0.13
Random Action	7.03±0.17	7.13±0.16	7.02±0.14	7.11±0.13	7.06±0.13
<b>PreFeRMAB</b>	8.35±0.12	8.38±0.11	8.26±0.11	8.10±0.11	8.00±0.10
DDLPO (topline)	n/a	n/a	n/a	n/a	8.09±0.11
ARMMAN with $N = 25, B = 7, S = 3$ .					
No Action	2.12±0.26	2.30±0.29	2.29±0.27	2.19±0.23	2.26±0.25
Random Action	2.86±0.32	3.27±0.40	3.01±0.30	3.09±0.35	2.96±0.31
<b>PreFeRMAB</b>	5.06±0.34	5.26±0.33	4.68±0.33	4.75±0.35	4.61±0.27
DDLPO (topline)	n/a	n/a	n/a	n/a	4.68±0.09

Table 6: Comparison of PreFeRMAB zero-shot performance on unseen arms against that of DDLPO trained and tested on the same set of arms. For each problem instance, we pretrain a model.

527 to train DDLPO from scratch vs. the number of samples  
 528 for fine-tuning PreFeRMAB starting from a pre-trained model  
 529 (additional results in Appendix A.5). Results suggests the  
 530 cost of pretraining can be amortized over different down-  
 531 stream instances. A non-profit organization using RMAB  
 532 models may have new beneficiaries opting in every week, and  
 533 training a new model from scratch every week can be 3-20  
 534 times more expensive than fine-tuning our pretrained model.

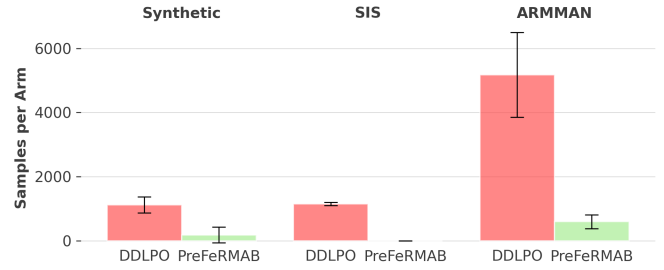


Figure 2: Comparison of samples per arm required by DDLPO and PreFeRMAB (fine-tuning using a pretrained model) to achieve maximum DDLPO reward across different environments. PreFeRMAB achieves the maximum topline reward with significantly fewer samples than DDLPO. Averages across training seeds are reported as interquartile means.

## 6 Conclusion

535 Our pretrained model (PreFeRMAB) leverages multi-arm  
 536 generalization, a novel update rule for a crucial  $\lambda$ -network,  
 537 and a StateShaping module for challenging reward functions.  
 538 PreFeRMAB demonstrates general zero-shot ability on un-  
 539 seen arms, and can be fine-tuned on specific instances in a  
 540 more sample-efficient way than training from scratch.  
 541

## Acknowledgments

542 The work is supported by Harvard HDSI funding.  
 543

## 544 Ethical Statement

545 The presented methods do not carry direct negative societal  
546 implications. However, training reinforcement learning mod-  
547 els should be done responsibly, especially given the safety  
548 concerns associated with agents engaging in extreme, unsafe,  
549 or uninformed exploration strategies. While the domains we  
550 considered such as ARMMAN do not have these concerns,  
551 the approach may be extended to extreme environments; in  
552 these cases, ensuring a robust approach to training reinforce-  
553 ment models is critical.

## 554 References

- 555 [Abel *et al.*, 2018] David Abel, Dilip Arumugam, Lucas  
556 Lehnert, and Michael Littman. State abstractions for life-  
557 long reinforcement learning. In *International Conference*  
558 *on Machine Learning*, pages 10–19. PMLR, 2018.
- 559 [Adelman and Mersereau, 2008] Daniel Adelman and  
560 Adam J Mersereau. Relaxations of weakly coupled  
561 stochastic dynamic programs. *Operations Research*,  
562 56(3):712–727, 2008.
- 563 [Akbarzadeh and Mahajan, 2019] N Akbarzadeh and A Ma-  
564 hajan. Restless bandits with controlled restarts: Indexabil-  
565 ity and computation of whittle index., 2019.
- 566 [Alon *et al.*, 1996] Noga Alon, Yossi Matias, and Mario  
567 Szegedy. The space complexity of approximating the fre-  
568 quency moments. In *Proceedings of the twenty-eighth an-  
569 nual ACM symposium on Theory of computing*, pages 20–  
570 29, 1996.
- 571 [Avrachenkov and Borkar, 2022] Konstantin E Avrachenkov  
572 and Vivek S Borkar. Whittle index based q-learning  
573 for restless bandits with average reward. *Automatica*,  
574 139:110186, 2022.
- 575 [Biswas *et al.*, 2021] Arpita Biswas, Gaurav Aggarwal,  
576 Pradeep Varakantham, and Milind Tambe. Learn to in-  
577 tervene: An adaptive learning policy for restless bandits  
578 in application to preventive healthcare. *arXiv preprint*  
579 *arXiv:2105.07965*, 2021.
- 580 [Bommasani *et al.*, 2021] Rishi Bommasani, Drew A Hud-  
581 son, Ehsan Adeli, Russ Altman, Simran Arora, Syd-  
582 ney von Arx, Michael S Bernstein, Jeannette Bohg, An-  
583 toine Bosselut, Emma Brunskill, et al. On the oppor-  
584 tunities and risks of foundation models. *arXiv preprint*  
585 *arXiv:2108.07258*, 2021.
- 586 [Chaudhuri and Kalyanakrishnan, 2020] Arghya Roy  
587 Chaudhuri and Shivaram Kalyanakrishnan. Regret  
588 minimisation in multi-armed bandits using bounded arm  
589 memory. In *Proceedings of the AAAI Conference on*  
590 *Artificial Intelligence*, volume 34, pages 10085–10092,  
591 2020.
- 592 [Devlin *et al.*, 2018] Jacob Devlin, Ming-Wei Chang, Ken-  
593 ton Lee, and Kristina Toutanova. Bert: Pre-training of  
594 deep bidirectional transformers for language understand-  
595 ing. *arXiv preprint arXiv:1810.04805*, 2018.
- 596 [Dusonchet and Hongler, 2003] Fabrice Dusonchet and M-O  
597 Hongler. Continuous-time restless bandit and dynamic

- scheduling for make-to-stock production. *IEEE Transac-* 598  
*tions on Robotics and Automation*, 19(6):977–990, 2003. 599
- [Fu *et al.*, 2019] Jing Fu, Yoni Nazarathy, Sarat Moka, and 600  
Peter G Taylor. Towards q-learning the whittle index for 601  
restless bandits. In *2019 Australian & New Zealand Con-* 602  
*trol Conference (ANZCC)*, pages 249–254. IEEE, 2019. 603
- [Gheshlaghi Azar *et al.*, 2013] Mohammad Ghesh- 604  
laghi Azar, Rémi Munos, and Hilbert J Kappen. Minimax 605  
pac bounds on the sample complexity of reinforcement 606  
learning with a generative model. *Machine learning*, 607  
91:325–349, 2013. 608
- [Hawkins, 2003] Jeffrey Thomas Hawkins. *A Lagrangian* 609  
*decomposition approach to weakly coupled dynamic opti-* 610  
*mization problems and its applications*. PhD thesis, Mas- 611  
sachusetts Institute of Technology, 2003. 612
- [Hodge and Glazebrook, 2015] David J Hodge and Kevin D 613  
Glazebrook. On the asymptotic optimality of greedy index 614  
heuristics for multi-action restless bandits. *Advances in* 615  
*Applied Probability*, 47(3):652–667, 2015. 616
- [Jaques *et al.*, 2019] Natasha Jaques, Angeliki Lazari- 617  
dou, Edward Hughes, Caglar Gulcehre, Pedro Ortega, 618  
DJ Strouse, Joel Z Leibo, and Nando De Freitas. Social 619  
influence as intrinsic motivation for multi-agent deep 620  
reinforcement learning. In *International conference on* 621  
*machine learning*, pages 3040–3049. PMLR, 2019. 622
- [Jin *et al.*, 2021] Chi Jin, Qinghua Liu, Yuanhao Wang, and 623  
Tiancheng Yu. V-learning—a simple, efficient, decen- 624  
tralized algorithm for multiagent rl. *arXiv preprint* 625  
*arXiv:2110.14555*, 2021. 626
- [Killian *et al.*, 2021a] Jackson A. Killian, Arpita Biswas, 627  
Sanket Shah, and Milind Tambe. Q-learning lagrange poli- 628  
cies for multi-action restless bandits. In *Proceedings of the* 629  
*27th ACM SIGKDD Conference on Knowledge Discovery* 630  
*& Data Mining*, KDD ’21, page 871–881, New York, NY, 631  
USA, 2021. Association for Computing Machinery. 632
- [Killian *et al.*, 2021b] Jackson A. Killian, Andrew Perrault, 633  
and Milind Tambe. Beyond ”to act or not to act”: Fast la- 634  
grangian approaches to general multi-action restless ban- 635  
dits. In *AAMAS*, pages 710–718, UK, 2021. AAMAS. 636
- [Killian *et al.*, 2022] Jackson A Killian, Lily Xu, Arpita 637  
Biswas, and Milind Tambe. Restless and uncertain: Ro- 638  
bust policies for restless bandits via deep multi-agent re- 639  
inforcement learning. In *Uncertainty in Artificial Intelli-* 640  
*gence*, pages 990–1000. PMLR, 2022. 641
- [Lefèvre, 1981] Claude Lefèvre. Optimal control of a 642  
birth and death epidemic process. *Operations Research*, 643  
29(5):971–982, 1981. 644
- [Li *et al.*, 2006] Lihong Li, Thomas J Walsh, and Michael L 645  
Littman. Towards a unified theory of state abstraction for 646  
mdps. In *AI&M*, 2006. 647
- [Liau *et al.*, 2018] David Liau, Zhao Song, Eric Price, and 648  
Ger Yang. Stochastic multi-armed bandits in constant 649  
space. In *International Conference on Artificial Intelli-* 650  
*gence and Statistics*, pages 386–394. PMLR, 2018. 651



- [Littman, 1994] Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In *Machine learning proceedings 1994*, pages 157–163. Elsevier, 1994.
- [Mate et al., 2021] Aditya Mate, Arpita Biswas, Christoph Siebenbrunner, Susobhan Ghosh, and Milind Tambe. Efficient algorithms for finite horizon and streaming restless multi-armed bandit problems. *arXiv preprint arXiv:2103.04730*, 2021.
- [Mate et al., 2022] Aditya Mate, Lovish Madaan, Aparna Taneja, Neha Madhiwalla, Shresth Verma, Gargi Singh, Aparna Hegde, Pradeep Varakantham, and Milind Tambe. Field study in deploying restless multi-armed bandits: Assisting non-profits in improving maternal and child health. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36:12017–12025, 06 2022.
- [Nakhleh et al., 2021] Khaled Nakhleh, Santosh Ganji, Ping-Chun Hsieh, I Hou, Srinivas Shakkottai, et al. Neurwin: Neural whittle index network for restless bandits via deep rl. *Advances in Neural Information Processing Systems*, 34:828–839, 2021.
- [Papadimitriou and Tsitsiklis, 1999] Christos H Papadimitriou and John N Tsitsiklis. The complexity of optimal queuing network control. *Mathematics of Operations Research*, 24(2):293–305, 1999.
- [Ramesh et al., 2021] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.
- [Riquelme et al., 2018] Carlos Riquelme, George Tucker, and Jasper Snoek. Deep bayesian bandits showdown: An empirical comparison of bayesian deep networks for thompson sampling. In *International Conference on Learning Representations*, 2018.
- [Shalev-Shwartz and Ben-David, 2014] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [Shapley, 1953] Lloyd S Shapley. Stochastic games. *Proceedings of the national academy of sciences*, 39(10):1095–1100, 1953.
- [Sinha and Mahajan, 2022] Amit Sinha and Aditya Mahajan. Robustness of whittle index policy to model approximation. *Available at SSRN 4064507*, 2022.
- [Team et al., 2023] Adaptive Agent Team, Jakob Bauer, Kate Baumli, Satinder Baveja, Feryal Behbahani, Avishkar Bhoopchand, Nathalie Bradley-Schmiege, Michael Chang, Natalie Clay, Adrian Collister, et al. Human-timescale adaptation in an open-ended task space. *arXiv preprint arXiv:2301.07608*, 2023.
- [Verma et al., 2023] Shresth Verma, Aditya Mate, Kai Wang, Neha Madhiwalla, Aparna Hegde, Aparna Taneja, and Milind Tambe. Restless multi-armed bandits for maternal and child health: Results from decision-focused learning. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, pages 1312–1320, 2023.
- [Vershynin, 2018] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [Wang et al., 2023] Kai Wang, Shresth Verma, Aditya Mate, Sanket Shah, Aparna Taneja, Neha Madhiwalla, Aparna Hegde, and Milind Tambe. Scalable decision-focused learning in restless multi-armed bandits with application to maternal and child health. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 12138–12146, 2023.
- [Weber and Weiss, 1990] Richard R Weber and Gideon Weiss. On an index policy for restless bandits. *Journal of applied probability*, 27(3):637–648, 1990.
- [Whittle, 1988] Peter Whittle. Restless bandits: Activity allocation in a changing world. *Journal of applied probability*, 25(A):287–298, 1988.
- [Xie et al., 2020] Qiaomin Xie, Yudong Chen, Zhaoran Wang, and Zhuoran Yang. Learning zero-sum simultaneous-move markov games using function approximation and correlated equilibrium. In *Conference on learning theory*, pages 3674–3682. PMLR, 2020.
- [Xiong and Li, 2023] Guojun Xiong and Jian Li. Finite-time analysis of whittle index based q-learning for restless multi-armed bandits with neural network function approximation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [Yaesoubi and Cohen, 2011] Reza Yaesoubi and Ted Cohen. Generalized markov models of infectious disease spread: A novel framework for developing dynamic health policies. *European journal of operational research*, 215(3):679–687, 2011.
- [Yu et al., 2020] Wenhao Yu, Jie Tan, Yunfei Bai, Erwin Coumans, and Sehoon Ha. Learning fast adaptation with meta strategy optimization. *IEEE Robotics and Automation Letters*, 5(2):2950–2957, 2020.
- [Yu et al., 2022] Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information Processing Systems*, 35:24611–24624, 2022.
- [Zhang and Frazier, 2021] Xiangyu Zhang and Peter I. Frazier. Restless bandits with many arms: Beating the central limit theorem, 2021.
- [Zhao et al., 2024] Yunfan Zhao, Nikhil Behari, Edward Hughes, Edwin Zhang, Dheeraj Nagaraj, Karl Tuyls, Aparna Taneja, and Milind Tambe. Towards zero shot learning in restless multi-armed bandits. *AAMAS (2-page extended abstract)*, 2024.

## 758 A Additional Experimental Details

### 759 A.1 Hyperparameters

760 In Table 7, we present hyperparameters used, with exceptions  
761 (1) for Continuous Synthetic, we use lambda scheduler dis-  
762 count rate = 0.95 (2) for Continuous SIS, we use training opt-  
763 in rate = 0.8.

764 In our experiments, all neural networks have 2 hidden lay-  
765 ers each with 16 neurons and tanh activation. The output layer  
766 has identity activation and its size is determined by the num-  
767 ber of actions (3 for SIS and Continuous SIS, and 2 for other  
768 environments).

769 The *lambda*-network training is similar to that in Killian *et*  
770 *al.*[2022]. After every  $n_{\text{subepochs}}$ , we update the  $\lambda$ -network  
771 and encourage the actor network to explore new parts of the  
772 state space immediately after the *lambda*-update (this explo-  
773 ration is controlled by the temperature parameter that weights  
774 the entropy term in the actor loss functions).

775 Different from Killian *et al.*[2022], we use a  $\lambda$ -network  
776 learning rate scheduler, which we found improves the perfor-  
777 mance and stability of the model.

Table 7: Hyperparameter values.

hyperparameter	value
training opt-in rate	0.8
agent clip ratio	2.0e+00
lambda freeze epochs	2.0e+01
start entropy coeff	5.0e-01
end entropy coeff	0.0e+00
actor learning rate	2.0e-03
critic learning rate	2.0e-03
lambda initial learning rate	2.0e-03
lambda scheduler discount rate	0.99
trains per epoch	2.0e+01
$n_{\text{subepochs}}$	4.0e+00

### 778 A.2 SIS Modeling (Discrete) Experimental Details

Recall that each arm  $p$  represents a subpopulation in distinct  
geographic regions. The state of each arm  $s$  is the number  
of uninfected people within the arm’s total population  $N_p$ .  
Transmission within each arm is guided by parameters:  $\kappa$ , the  
average number of contacts within the arm’s subpopulation in  
each round, and  $r_{\text{infect}}$ , the probability of becoming infected  
after contact with an infected person. The probability that a  
single uninfected person gets infected is then:

$$q = 1 - e^{-\kappa \cdot \frac{S-s}{S} \cdot r_{\text{infect}}},$$

779 where  $S$  is the number of possible states, and  $s \in [S]$  is the  
780 current state. Note  $\frac{S-s}{S}$  is the percentage of people who are  
781 currently infected. The number of infected people in the next  
782 timestep follows a binomial distribution  $B(S, q)$ .

783 Recall that there are three available intervention actions  
784  $a_0, a_1, a_2$  that affect the transmission parameters:  $a_0$  repre-  
785 sents no action;  $a_1$  represents messaging about physical dis-  
786 tancing, dividing  $\kappa$  by  $a_1^{\text{eff}}$ ;  $a_2$  represents the distribution of  
787 face masks, dividing  $r_{\text{infect}}$  by  $a_2^{\text{eff}}$ . The actions costs are  
788  $c = \{0, 1, 2\}$ . Following the implementation in [Killian *et*  
789 *al.*, 2022], these parameters are sampled within ranges:

$$\kappa \in [1, 10], r_{\text{infect}} \in [0.5, 0.99], a_1^{\text{eff}} \in [1, 10], a_2^{\text{eff}} \in [1, 10]$$

### A.3 Continuous States Experimental Details

We consider a synthetic dataset with continuous states and  
binary actions. For the current state  $s_i$  of arm  $i$ , and action  $a$ ,  
the next state  $s'_i$  is represented by the transition dynamic:

$$s'_i = \begin{cases} \text{clip}(s_i + \mathcal{N}(\mu_{i0}, \sigma_{i0}), 0, 1) & \text{if } a = 0 \\ \text{clip}(s_i + \mathcal{N}(\mu_{i1}, \sigma_{i1}), 0, 1) & \text{if } a = 1 \end{cases}$$

Where the transition dynamics are sampled uniformly from  
the intervals ( $\sigma_{i0} = \sigma_{i1} = 0.2$  is fixed):

$$\mu_{i0} \in [-0.5, -0.1], \quad \mu_{i1} \in [0.1, 0.5].$$

We also consider continuous state experiments in real-  
world settings. In the discrete state SIS Epidemic Model de-  
scribed above, each arm represents a subpopulation, and the  
state of that arm represents the number of uninfected people  
within the subpopulation. In real-world public health settings  
such as COVID-19 control, interventions like quarantine and  
mask mandates may be imposed on subpopulations of very  
large sizes such as an entire city. The SIS model is expected to  
scale to a continuum limit as the population size increases to  
infinity. Thus, a SIS model with population 1 million would  
behave roughly similar to that with population 1 billion in  
terms of the proportions. This notion is inherently captured  
by continuous models but not by those dealing with absolute  
numbers.

Following Killian *et al.*[2022], within an arm, any un-  
infected person will get infected with the same probabil-  
ity. Thus, the number of uninfected people in the next  
timestep follows a binomial distribution. It is well-known  
that a normal distribution  $\mathcal{N}(\mu, \sigma^2)$  well approximates a bi-  
nomial distribution  $B(n, p)$ , with the choice  $\mu = np$  and  
 $\sigma^2 = np(1 - p)$ , when  $n$  is sufficiently large.

### A.4 Distributional Shift Details

In real-world resources allocation problems, we may observe  
distribution shifts in arms, i.e., arms in testing are sampled  
from a distribution slightly different from that in training.  
In public health settings, a non-profit organization solving  
RMAB problems to allocate resources to beneficiaries may  
observe that beneficiaries’ behavior or feature information  
change over time [Wang *et al.*, 2023; Killian *et al.*, 2022].  
Additionally, a non-profit organization may have new bene-  
ficiaries joining who are in a different subpopulation. In Ta-  
ble 2, we provide ablation results illustrating that **PreFeR-  
MAB is robust to distribution shift in arms**. We mea-  
sure the shift in distribution using Wasserstein distance. The  
results demonstrate that even on arm samples from distri-  
butions that significantly deviate from that seen in training,  
PreFeRMAB still achieves strong performance and outper-  
forms baselines.

For each arm, the associated Markov Decision Process  
(MDP) has only two discrete states, the transition dynam-  
ics  $p(s'|s, a)$ , representing the probability of transitioning to  
state  $s'$  from state  $s$  given action  $a$ , can be described by four  
Bernoulli random variables, one for each combination of state  
and action. By introducing a uniform distribution shift, we  
can modify the transition probabilities of the Bernoulli ran-  
dom variable associated with each state-action pair by adding

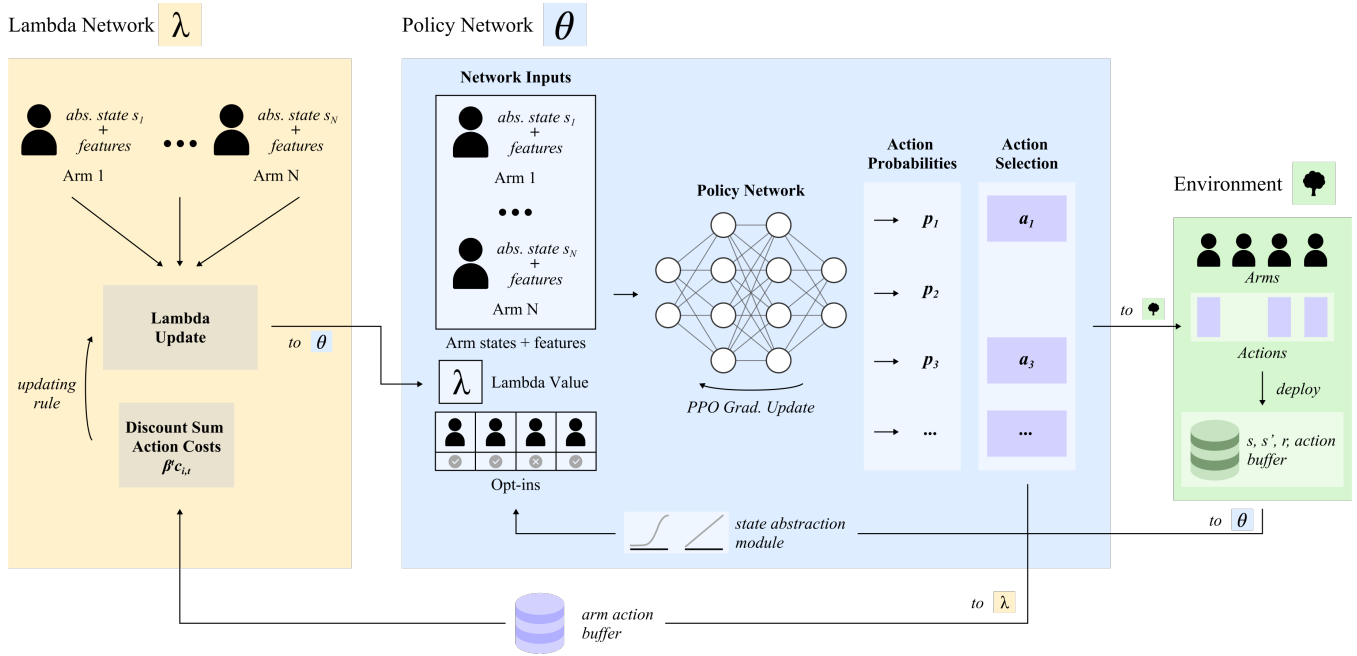


Figure 3: Overview of the PreFeRMAB training procedure. A trained model consists of a policy network, a critic network, a  $\lambda$ -network, and a StateShaping module. Arm states  $s_i$ , features  $z_i$ , and opt-in decisions  $\xi$  are passed through the policy network with an action-charge  $\lambda$ . The policy network independently predicts action probabilities for each arm, which are then greedily selected until the specified budget is reached. These selected actions are used with arm state, feature, and opt-in information to update the  $\lambda$ -network. Updated arm states  $s'$  and rewards  $r$  from the environment are then added to the buffer, and passed through the state abstraction module before being fed back through the policy network.

842 a constant  $\delta$  to the parameter of each Bernoulli distribution.  
 843 Consequently, this results in a consistent shift in the transition  
 844 probabilities across all states and actions.

845 Furthermore, this delta is exactly the Wasserstein distance  
 846 between the two distributions, which we show here. Suppose  
 847 we have two discrete probability distributions,  $P$  and  $Q$ , with  
 848 their respective probabilities associated with the outcomes  $x_i$   
 849 and  $y_j$ . The 1-Wasserstein distance, also known as the  
 850 Earth Mover's Distance  $W(P, Q)$ , can be calculated by solving  
 851 the following optimization problem:

$$W(P, Q) = \inf_{\gamma \in \Gamma(P, Q)} \sum_{i, j} |x_i - y_j| \gamma(x_i, y_j),$$

852 where  $\Gamma(P, Q)$  is the set of all joint distributions  $\gamma$  with  
 853 marginals  $P$  and  $Q$ , and  $\gamma(x_i, y_j)$  represents the amount of  
 854 "mass" moved from  $x_i$  to  $y_j$ .

855 Note  $W(P, Q)$  between two Bernoulli distributions with  
 856 parameters  $b_1$  and  $b_2$  can be succinctly determined as  $|b_2 - b_1|$ .  
 857 This is because each Bernoulli distribution has only two po-  
 858 tential outcomes, 0 and 1, and so moving mass from one out-  
 859 come to another across these distributions involves a shift of  
 860 probability mass  $|b_2 - b_1|$  across the one-unit distance be-  
 861 tween the two points. Therefore, without loss of generality  
 862 assuming  $b_2 \geq b_1$ , the Wasserstein distance simplifies to the  
 863 non-negative difference  $b_2 - b_1$ .

### 864 A.5 Fast Fine-Tuning

865 In subsection 5.3, we demonstrate that, in addition to strong  
 866 zero-shot performance, PreFeRMAB may also be used as

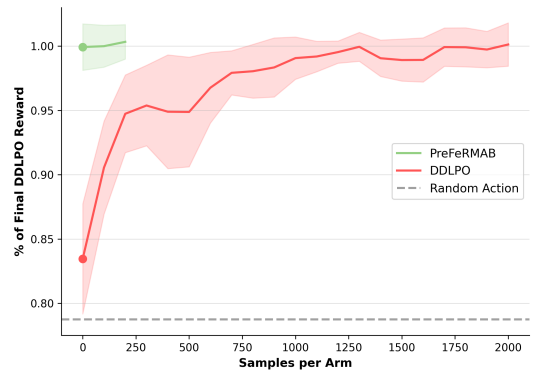


Figure 4: Comparison of the percentage of the final DDLPO (Kilian *et al.* [2022] topline) reward achieved by the number of samples per arm. In DDLPO, samples are used for training from scratch; in PreFeRMAB, samples are used to fine-tune a pretrained PreFeRMAB model. Results indicate that PreFeRMAB, from zero-shot results, achieves near-optimal performance, and requires a small fraction of the required DDLPO samples to achieve final DDLPO performance.

867 a *pretrained model* for fast fine-tuning in specific domains. 867  
 868 In particular, we demonstrate that we may start from a pre- 868  
 869 trained PreFeRMAB model, and train on additional samples 869  
 870 for a fixed environment (with fixed arm transition dynam- 870  
 871 ics). We showed that using this pre-trained model can help 871  
 872 achieve topline DDLPO performance in significantly fewer 872

873 fine-tuning samples than required by DDLPO to train from  
 874 scratch. In Figure 4, we further visualize these results in  
 875 training curves comparing DDLPO and PreFeRMAB. These  
 876 training curves, which plot the number of samples per arm  
 877 against the achieved percentage of final DDLPO reward, are  
 878 shown for the discrete-state synthetic environment setting for  
 879  $N=21, B=7.0$

880 The results in Figure 4 demonstrate that PreFeRMAB  
 881 shows both 1) strong zero-shot performance, achieving near-  
 882 topline reward with *no fine-tuning samples required*, as well  
 883 as 2) a significant reduction in the number of samples re-  
 884 quired to achieve final DDLPO performance. In particular,  
 885 we note that DDLPO, before training, achieves a reward only  
 886 marginally higher than the average Random Action reward.  
 887 Alternatively, PreFeRMAB begins, in a zero-shot setting,  
 888 with a much higher initial reward value. We also observe that  
 889 PreFeRMAB requires significantly fewer samples per arm to  
 890 achieve the final DDLPO reward. This is particularly critical  
 891 in high-stakes, real-world settings where continually sam-  
 892 pling arms from the environment may be prohibitively expen-  
 893 sive, especially for low-resource NGOs.

## 894 A.6 StateShaping

895 Figure 5 provides a simple example, illustrating how we adapt  
 896 states through the state abstraction procedure. In this particu-  
 897 lar example, the reward is an increasing function of the state,  
 898 and the reward plateaus at state 0.5, i.e.  $s \in [0.5, 1]$  achieve  
 899 the same reward. We map all raw observations in the range  
 900  $[0.5, 1]$  to abstract state 1. We note that this process is au-  
 901 tomatized, using data collected on arm states and reward from  
 902 prior (historical) samples to 1) estimate the reward of a cur-  
 903 rent arm, and 2) use this reward to normalize the arm state.  
 904 We demonstrate how these states are mapped to normalized  
 905 values in Table 5.

906 In Figure 5, we show results in two settings after 30 epochs  
 907 of training, evaluating on a separate set of test arms for zero-  
 908 shot evaluation. Continuous transition dynamics are used  
 909 directly as input features for training and evaluation. The  
 910 results illustrate that state abstraction can help achieve ad-  
 911 ditional performance gains for various challenging reward  
 912 functions. Drawing from prior literature on state abstraction,  
 913 this modular component of PreFeRMAB may also serve as a  
 914 placeholder for future automated state abstraction procedures  
 915 to improve generalizability and robustness of PreFeRMAB  
 916 across domains with challenging reward functions.

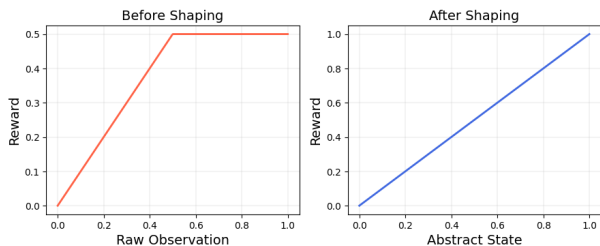


Figure 5: Illustration for StateShaping.

## B Ablation Studies

917 In this section, we provide ablation results over (1) a wider  
 918 range of opt-in rates than presented in the main paper (Ta-  
 919 ble 12) (2) different feature mappings, including linear and  
 920 non-linear feature transformations of the original transition  
 921 probabilities (3) DDLPO topline (Killian *et al.* [2022]) with  
 922 and without transition probability features as inputs (4) re-  
 923 sults in more problem settings. The ablation results showcase  
 924 that PreFeRMAB consistently achieve strong performance  
 925 and having access to feature information does not provide  
 926 PreFeRMAB an unfair advantage over DDLPO. 927

### B.1 Opt-in Rates

Number of arms System capacity	30%	40%	50%	60%	70%
System capacity $N = 21$ . Budget $B = 7$ .					
No Action	$3.09 \pm 0.31$	$3.10 \pm 0.32$	$3.12 \pm 0.29$	$3.14 \pm 0.25$	$3.16 \pm 0.22$
Random Action	$3.57 \pm 0.48$	$3.46 \pm 0.48$	$3.55 \pm 0.35$	$3.55 \pm 0.34$	$3.57 \pm 0.33$
<b>PreFeRMAB</b>	$3.78 \pm 0.72$	$3.75 \pm 0.70$	$4.16 \pm 0.57$	$4.52 \pm 0.54$	$4.45 \pm 0.42$
System capacity $N = 48$ . Budget $B = 16$ .					
No Action	$3.19 \pm 0.27$	$3.15 \pm 0.23$	$3.13 \pm 0.12$	$3.17 \pm 0.12$	$3.17 \pm 0.14$
Random Action	$3.44 \pm 0.26$	$3.43 \pm 0.23$	$3.46 \pm 0.20$	$3.47 \pm 0.17$	$3.44 \pm 0.17$
<b>PreFeRMAB</b>	$3.90 \pm 0.50$	$3.64 \pm 0.30$	$3.87 \pm 0.27$	$3.85 \pm 0.25$	$4.06 \pm 0.31$
System capacity $N = 96$ . Budget $B = 32$ .					
No Action	$3.21 \pm 0.17$	$3.17 \pm 0.20$	$3.17 \pm 0.15$	$3.18 \pm 0.14$	$3.17 \pm 0.13$
Random Action	$3.54 \pm 0.22$	$3.55 \pm 0.23$	$3.58 \pm 0.18$	$3.55 \pm 0.18$	$3.56 \pm 0.13$
<b>PreFeRMAB</b>	$3.93 \pm 0.33$	$3.72 \pm 0.23$	$3.79 \pm 0.18$	$4.02 \pm 0.21$	$4.16 \pm 0.25$

Table 8: Robustness to different opt-in rates with identity mapping. Evaluation follows Table 10: we run for 50 trials with 2 total number of states for each arm, and pretrain a model for each system capacity  $N$  and test generalization on different opt-in rates.

Number of arms System capacity	30%	40%	50%	60%	70%
System capacity $N = 21$ . Budget $B = 7$ .					
No Action	$3.13 \pm 0.36$	$3.19 \pm 0.34$	$3.17 \pm 0.27$	$3.17 \pm 0.23$	$3.15 \pm 0.24$
Random Action	$3.60 \pm 0.53$	$3.62 \pm 0.40$	$3.59 \pm 0.34$	$3.58 \pm 0.36$	$3.57 \pm 0.32$
<b>PreFeRMAB</b>	$3.79 \pm 0.49$	$3.76 \pm 0.49$	$3.79 \pm 0.39$	$3.89 \pm 0.39$	$3.86 \pm 0.46$
System capacity $N = 48$ . Budget $B = 16$ .					
No Action	$3.15 \pm 0.23$	$3.17 \pm 0.21$	$3.18 \pm 0.16$	$3.18 \pm 0.13$	$3.19 \pm 0.14$
Random Action	$3.41 \pm 0.29$	$3.52 \pm 0.27$	$3.50 \pm 0.20$	$3.49 \pm 0.22$	$3.48 \pm 0.17$
<b>PreFeRMAB</b>	$3.61 \pm 0.51$	$4.08 \pm 0.34$	$4.45 \pm 0.30$	$4.44 \pm 0.30$	$4.44 \pm 0.29$
System capacity $N = 96$ . Budget $B = 32$ .					
No Action	$3.18 \pm 0.18$	$3.17 \pm 0.13$	$3.18 \pm 0.14$	$3.19 \pm 0.14$	$3.16 \pm 0.15$
Random Action	$3.54 \pm 0.22$	$3.57 \pm 0.22$	$3.55 \pm 0.14$	$3.58 \pm 0.15$	$3.56 \pm 0.15$
<b>PreFeRMAB</b>	$3.74 \pm 0.21$	$3.68 \pm 0.23$	$3.76 \pm 0.21$	$3.98 \pm 0.19$	$4.06 \pm 0.22$

Table 9: Robustness to different opt-in rates with linear-mappings. Evaluation follows Table 10: we run for 50 trials with 2 total number of states for each arm, and pretrain a model for each system capacity  $N$  and test generalization on different opt-in rates.

929 Throughout the main paper, we provide results for evalu-  
 930 ation opt-in rates in the range 80%-100%. In Table 8 and  
 931 Table 9, we provide ablation results for **opt-in rates in a**  
 932 **wider range of 30%-70%**. During the training phase, we  
 933 maintain an expected opt-in rate of 80%, which may gen-  
 934 erally range from 70%-90% every training iteration. Given  
 935 this training configuration, we demonstrate strong results in  
 936 the main paper for evaluating on a similar range of test-time  
 937 opt-ins from 80% to 100%. However, we also further demon-  
 938 strate in Table 8 and Table 9 that our pretrained PreFeRMAB 938

939 model, despite a training opt-in rate around 80% in expect- 965  
 940 tation, achieves strong results on *testing* opt-in rates from a 966  
 941 substantially different range. These results highlight PreFeR- 967  
 942 MAB’s flexibility and ability to generalize to unseen opt-in 968  
 943 rates, which may be critical in real-world applications where 969  
 944 arms frequently exit and re-enter the environment. 970

## 945 B.2 Feature Mapping 971

Number of arms System capacity	80%	85%	90%	95%	100%
System capacity $N = 21$ . Budget $B = 7$ .					
No Action	3.46 ± 0.20	3.39 ± 0.19	3.40 ± 0.17	3.40 ± 0.18	3.22 ± 0.16
Random Action	3.80 ± 0.31	3.76 ± 0.30	3.79 ± 0.29	3.76 ± 0.31	3.58 ± 0.27
<b>PreFeRMAB</b>	4.57 ± 0.29	4.70 ± 0.33	4.70 ± 0.29	4.64 ± 0.29	4.37 ± 0.25
System capacity $N = 48$ . Budget $B = 16$ .					
No Action	3.22 ± 0.13	3.28 ± 0.13	3.22 ± 0.12	3.29 ± 0.12	3.21 ± 0.11
Random Action	3.56 ± 0.19	3.65 ± 0.18	3.56 ± 0.18	3.65 ± 0.17	3.57 ± 0.17
<b>PreFeRMAB</b>	3.94 ± 0.23	4.00 ± 0.23	3.86 ± 0.20	3.90 ± 0.16	3.73 ± 0.16
System capacity $N = 96$ . Budget $B = 32$ .					
No Action	3.24 ± 0.10	3.24 ± 0.10	3.21 ± 0.10	3.21 ± 0.10	3.20 ± 0.10
Random Action	3.61 ± 0.14	3.62 ± 0.15	3.57 ± 0.15	3.57 ± 0.13	3.57 ± 0.14
<b>PreFeRMAB</b>	4.35 ± 0.16	4.36 ± 0.15	4.29 ± 0.14	4.23 ± 0.13	4.20 ± 0.12

Table 10: Results on non-linearly transformed synthetic discrete states. We present final reward divided by the number of arms, averaged over 50 trials with each trial consisting on 10 rounds, for a total of 500 evaluations. The number of states  $S = 2$ . For each system capacity  $N$ , we pretrain a model

946 In our main paper, we use linear feature mapping, pro- 965  
 947 jecting true transition probabilities to features with randomly 966  
 948 generated projection matrices. This can be represented by 967  
 949  $\mathbf{y} = \mathbf{Ax}$ , where  $\mathbf{y}$  are the output features,  $\mathbf{A}$  is the transfor- 968  
 950 mation matrix, and  $\mathbf{x}$  denotes the ground truth arm transition 969  
 951 probabilities. To demonstrate the robustness of our approach 970  
 952 to various types of input features, we also consider **more 971**  
 953 **challenging, non-linear feature mappings**, which may intro- 972  
 954 duce higher representational complexity as compared to 973  
 955 linear feature mappings. For these ablation results, we use 974  
 956 a sigmoidal transformation, which can be expressed as  $\mathbf{y} =$  975  
 957  $\frac{1}{1 + \exp(-\mathbf{Ax})}$ . We demonstrate the results using these non- 976  
 958 linear feature mappings in Table 10. These results indicate 977  
 959 that PreFeRMAB consistently outperforms baselines under 978  
 960 various forms of feature mappings, and is robust to both linear 979  
 961 and non-linear input features. 980

## 962 B.3 DDLPO Topline with Features 984

Synthetic Experiment	N=21,B=7.0	N=48,B=16.0	N=96,B=32.0
<b>DDLPO, w/o Features</b>	4.63 ± 0.21	4.60 ± 0.18	4.35 ± 0.11
<b>DDLPO, w/ Features</b>	4.63 ± 0.23	4.59 ± 0.17	4.21 ± 0.11

Table 11: Performance comparison of Killian *et al.* [2022] DDLPO topline, with and without ground truth transition probabilities as input features. Results are shown for evaluation on a single, fixed training seed. The results suggest that transition probability features do not significantly improve the final performance of the topline DDLPO model—this implies that PreFeRMAB does not leverage these features for an unfair reward advantage.

963 In Table 11, we show that having access to **features** 985  
 964 **does not boost the performance of DDLPO**. Features help 986

965 PreFeRMAB generalize to unseen arms and achieve strong 966  
 967 zero-shot results, as demonstrated in the main paper. How- 968  
 969 ever, one may ask whether access to these features, as used by 970  
 971 PreFeRMAB, may provide an unfair reward advantage over 972  
 973 DDLPO, which in its original form [Killian *et al.*, 2022] does 974  
 975 not utilize feature information. That is, because input features 976  
 977 in our experiments are derived from the original arm transi- 978  
 979 tion probabilities, it may be the case that these are used to 979  
 980 achieve better performance. To determine whether there is an 981  
 982 advantage from utilizing these features, we modify the origi- 982  
 983 nal DDLPO model to accept ground truth transition proba- 983  
 984 bilities for each arm as feature inputs to the respective policy 984  
 985 networks. We present results for DDLPO with and without 985  
 986 input features, for a fixed seed, in Table Table 11. In this 986  
 987 table, we observe that across synthetic experiments for various 987  
 988 system capacities and budgets, DDLPO’s performance does 988  
 989 not improve given access to features. These results suggest 989  
 990 that PreFeRMAB is not leveraging the input features to gain 990  
 991 an unfair advantage in evaluation. 991

## 984 B.4 Different values of $N, B, S$ 984

985 We present results on a wider range of problem settings, 985  
 986 specifically different number of arms  $N$ , different budget  $B$ , 986  
 987 and (for SIS Epidemic Modeling only), different number of 987  
 988 possible states  $S$ . 988

Number of arms System capacity	80%	85%	90%	95%	100%
System capacity $N = 96$ . Budget $B = 32$ .					
No Action	3.22±0.12	3.24±0.12	3.19±0.11	3.18±0.11	3.18±0.11
Random Action	3.62±0.13	3.66±0.13	3.58±0.13	3.60±0.12	3.60±0.12
<b>PreFeRMAB</b>	4.63±0.12	4.71±0.12	4.53±0.13	4.47±0.12	4.61±0.10
DDLPO (topline)	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	4.58±0.13
System capacity $N = 48$ . Budget $B = 16$ .					
No Action	3.19±0.11	3.24±0.13	3.18±0.11	3.23±0.11	3.17±0.10
Random Action	3.61±0.17	3.70±0.21	3.56±0.18	3.67±0.17	3.58±0.16
<b>PreFeRMAB</b>	4.77±0.18	4.74±0.16	4.62±0.19	4.94±0.14	4.78±0.14
DDLPO (topline)	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	4.76±0.14
System capacity $N = 21$ . Budget $B = 7$ .					
No Action	3.44±0.21	3.43±0.19	3.41±0.20	3.38±0.17	3.22±0.16
Random Action	3.82±0.32	3.79±0.33	3.77±0.31	3.76±0.28	3.58±0.27
<b>PreFeRMAB</b>	4.20±0.27	4.46±0.23	4.48±0.23	4.74±0.26	4.56±0.23
DDLPO (topline)	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	4.81±0.14

Table 12: Results on Synthetic with discrete states. We present final reward divided by the number of arms, averaged over 50 trials. For each system capacity  $N$ , we pretrain a model. The DDLPO (topline) does not accomodate different opt-in rates and can only be used on 100% opt-in.

989 **Synthetic Evaluation:** We first evaluate the performance 989  
 990 of PreFeRMAB in the discrete state synthetic environment 990  
 991 setting described above. Table 12 illustrates these results. 991  
 992 In this synthetic setting, we find that PreFeRMAB is able 992  
 993 to consistently outperform Random Action and No Action 993  
 994 baselines, and achieve performance comparable to the topline 994  
 995 DDLPO approach. Critically, PreFeRMAB achieves good re- 995  
 996 ward outcomes across changing system capacity  $N$ , budgets 996  
 997  $B$ , as well as different opt-in rates. Additionally, we find 997  
 998 that PreFeRMAB achieves *near-topline results from zero-* 998  
 999 *shot learning* in the synthetic setting, compared to the topline 999  
 1000 DDLPO approach which is trained and evaluated on a fixed 1000

1001 set of arm transition dynamics for 100 epochs (we take the  
1002 best performance of DDLPO across the 100 epochs).

Number of arms System capacity	80%	85%	90%	95%	100%
Number of possible states per arm $S = 150$ .					
No Action	5.33 $\pm$ 0.16	5.30 $\pm$ 0.15	5.31 $\pm$ 0.14	5.29 $\pm$ 0.13	5.28 $\pm$ 0.13
Random Action	7.03 $\pm$ 0.17	7.13 $\pm$ 0.16	7.02 $\pm$ 0.14	7.11 $\pm$ 0.13	7.06 $\pm$ 0.13
<b>PreFeRMAB</b>	8.35 $\pm$ 0.12	8.38 $\pm$ 0.11	8.26 $\pm$ 0.11	8.10 $\pm$ 0.11	8.00 $\pm$ 0.10
DDLPO (topline)	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	8.09 $\pm$ 0.11
Number of possible states per arm $S = 100$ .					
No Action	5.28 $\pm$ 0.15	5.20 $\pm$ 0.13	5.30 $\pm$ 0.15	5.25 $\pm$ 0.14	5.27 $\pm$ 0.15
Random Action	6.95 $\pm$ 0.19	7.01 $\pm$ 0.16	7.11 $\pm$ 0.16	7.06 $\pm$ 0.15	7.07 $\pm$ 0.15
<b>PreFeRMAB</b>	7.88 $\pm$ 0.20	7.91 $\pm$ 0.19	7.99 $\pm$ 0.18	8.01 $\pm$ 0.17	8.02 $\pm$ 0.16
DDLPO (topline)	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	7.99 $\pm$ 0.08
Number of possible states per arm $S = 50$ .					
No Action	5.39 $\pm$ 0.15	5.47 $\pm$ 0.15	5.42 $\pm$ 0.13	5.44 $\pm$ 0.12	5.46 $\pm$ 0.12
Random Action	7.29 $\pm$ 0.17	7.33 $\pm$ 0.17	7.26 $\pm$ 0.14	7.38 $\pm$ 0.15	7.33 $\pm$ 0.12
<b>PreFeRMAB</b>	8.51 $\pm$ 0.08	8.37 $\pm$ 0.11	8.24 $\pm$ 0.07	8.10 $\pm$ 0.10	7.93 $\pm$ 0.09
DDLPO (topline)	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	8.04 $\pm$ 0.08

Table 13: Results on SIS Epidemic Model with discrete states. We present final reward divided by the number of arms, averaged over 50 trials. System capacity  $N = 20$  and budget  $B = 16$ . For each number of possible states per arm  $S$ , we pretrain a model. The DDLPO (topline) does not accommodate different opt-in rates and can only be used on 100% opt-in.

Number of arms System capacity	80%	85%	90%	95%	100%
System capacity $N = 25$ . Budget $B = 7$ .					
No Action	2.12 $\pm$ 0.26	2.30 $\pm$ 0.29	2.29 $\pm$ 0.27	2.19 $\pm$ 0.23	2.26 $\pm$ 0.25
Random Action	2.86 $\pm$ 0.32	3.27 $\pm$ 0.40	3.01 $\pm$ 0.30	3.09 $\pm$ 0.35	2.96 $\pm$ 0.31
<b>PreFeRMAB</b>	5.06 $\pm$ 0.34	5.26 $\pm$ 0.33	4.68 $\pm$ 0.33	4.75 $\pm$ 0.35	4.61 $\pm$ 0.27
DDLPO (topline)	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	4.68 $\pm$ 0.09
System capacity $N = 25$ . Budget $B = 5$ .					
No Action	2.14 $\pm$ 0.23	2.29 $\pm$ 0.26	2.24 $\pm$ 0.28	2.36 $\pm$ 0.24	2.19 $\pm$ 0.23
Random Action	2.68 $\pm$ 0.31	2.95 $\pm$ 0.36	2.75 $\pm$ 0.32	2.92 $\pm$ 0.26	2.69 $\pm$ 0.21
<b>PreFeRMAB</b>	4.10 $\pm$ 0.32	4.45 $\pm$ 0.40	4.39 $\pm$ 0.33	4.48 $\pm$ 0.34	3.95 $\pm$ 0.34
DDLPO (topline)	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	4.29 $\pm$ 0.25
System capacity $N = 50$ . Budget $B = 10$ .					
No Action	2.27 $\pm$ 0.24	2.31 $\pm$ 0.17	2.19 $\pm$ 0.22	2.21 $\pm$ 0.21	2.27 $\pm$ 0.23
Random Action	2.82 $\pm$ 0.26	2.91 $\pm$ 0.22	2.72 $\pm$ 0.23	2.69 $\pm$ 0.18	2.77 $\pm$ 0.23
<b>PreFeRMAB</b>	4.21 $\pm$ 0.30	3.98 $\pm$ 0.28	3.89 $\pm$ 0.28	3.68 $\pm$ 0.28	3.62 $\pm$ 0.26
DDLPO (topline)	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	<i>n/a</i>	4.08 $\pm$ 0.26

Table 14: Results on ARMMAN with discrete states. We present final reward divided by the number of arms, averaged over 50 trials. For each pair of  $(N, B)$ , we pretrain a model. The DDLPO (topline) does not accommodate different opt-in rates and can only be used on 100% opt-in.

1003 **SIS Evaluation:** Next, we evaluate the performance of  
1004 PreFeRMAB in the discrete-state SIS modelling setting. Table  
1005 13 illustrates these results. We evaluate PreFeRMAB  
1006 for  $N = 20, B = 16$  on three different number of possible  
1007 states per arm  $S = 50, 100, 150$ , representing the maximum  
1008 population of a region in the SIS setting. The results  
1009 shown demonstrate that PreFeRMAB performs well in *zero-*  
1010 *shot learning* in settings that model real-world planning prob-  
1011 lems, especially with larger state spaces and with multiple ac-  
1012 tions. We again find that PreFeRMAB achieves results compar-  
1013 able to the DDLPO topline with zero-shot testing, compar-  
1014 ed to DDLPO trained and evaluated on the same constant

set of arms.

**ARMMAN Evaluation:** We next evaluate the performance of PreFeRMAB in the discrete state ARMMAN modeling setting. Table 14 illustrates these results. In these experiments, we show performance for  $S = 3$  across 3 training configurations ( $(N = 25, B = 5)$ ,  $(N = 25, B = 7)$ ,  $(N = 50, B = 10)$ ) for 5 test-time opt-in rates. We observe that our approach again performs consistently well in a more *challenging* setting that models real-world planning problems across different system capacities, budgets, and opt-in rates. Specifically, we validate that PreFeRMAB can achieve higher average rewards for increased budgets given a fixed system capacity, which is expected as reward potential increases with higher budgets. Additionally, we see that PreFeRMAB again achieves *zero-shot* results comparable to the DDLPO topline reward, reaching  $\sim 90\%$  of the topline reward in zero-shot evaluation.

## C Multi-arm Generalization

In the main paper (Table 1), we presented results on Synthetic with  $N = 21, B = 7$ , demonstrating the benefit of multi-arm generalization. The results are obtained when the Wasserstein distance between training and testing distribution is 0.05 (see Sec A.4 for how we compute the Wasserstein distance). We provide additional results to further showcase the benefits of multi-arm generalization. Specifically, in Table 15, we present results for  $N = 12, B = 3$ .

# Unique arms	System capacity $N = 12$ . Budget $B = 3$ .				
	48	39	30	21	12
No Action	3.11 $\pm$ 0.31	3.11 $\pm$ 0.31	3.11 $\pm$ 0.31	3.11 $\pm$ 0.31	3.11 $\pm$ 0.31
Random Action	3.45 $\pm$ 0.31	3.45 $\pm$ 0.31	3.45 $\pm$ 0.31	3.45 $\pm$ 0.31	3.45 $\pm$ 0.31
<b>PreFeRMAB</b>	4.35 $\pm$ 0.28	4.28 $\pm$ 0.29	4.31 $\pm$ 0.27	4.04 $\pm$ 0.32	3.60 $\pm$ 0.30

Table 15: Multi-arm generalization results on Synthetic (opt-in 100%). With the same total amount of data, PreFeRMAB achieves stronger performance when pretrained on more unique arms.

## D Proof of Multi-Arm Generalization

In this section, we will shorten  $n_{\text{epochs}}$  to  $n$  for the sake of clarity. In this section, we let  $C_{\text{sys}}$  to denote a constant which depends on the parameters of the MDP such as budget per arm  $B/N$ , cost  $c_j$ , discount factor  $\beta$ ,  $\lambda_{\text{max}}$ ,  $R_{\text{max}}$ ,  $D$ ,  $d$  and  $L$ . It can denote a different constant in every appearance. We list the assumptions made in the statement of the proposition below for the sake of clarity.

**Assumption 1.** Suppose the learning algorithm learns neural network weights  $\hat{\theta}$ , whose policy is optimal for each  $(\hat{\mu}_i, \lambda)$  for  $i = 1, 2, \dots, n$  and  $\lambda \in [0, \lambda_{\text{max}}]$ . That is, it learns the optimal policy for every sample in the training data.

**Assumption 2.** There exists a choice of weights  $\theta^* \in \Theta$  which gives the optimal policy for every set of  $N$  features  $(\hat{\mu})$  drawn as the empirical distribution of i.i.d. samples from  $\mu^*$  and for every  $\lambda \in [0, \lambda_{\text{max}}]$ .

**Assumption 3.**  $\Theta = \mathcal{B}_2(D, \mathbb{R}^d)$ , the  $\ell_2$  ball of radius  $D$  in  $\mathbb{R}^d$ . We assume that

$$|V(\mathbf{s}, \theta_1, \lambda, \hat{\mu}) - V(\mathbf{s}, \theta_2, \lambda, \hat{\mu})| \leq L \|\theta_1 - \theta_2\|$$

$$|V(\mathbf{s}, \theta, \lambda_1, \hat{\mu}) - V(\mathbf{s}, \theta, \lambda_2, \hat{\mu})| \leq L|\lambda_1 - \lambda_2|$$

1057 Define the population average value function by  $\bar{V}(s, \theta) =$   
 1058  $\mathbb{E}_{\hat{\mu}} \inf_{\lambda \in [0, \lambda_{\max}]} V(\mathbf{s}, \theta, \lambda, \hat{\mu})$  and the sample average value  
 1059 function by  $\hat{V}(s, \theta) = \frac{1}{n} \sum_{j=1}^n \inf_{\lambda \in [0, \lambda_{\max}]} V(\mathbf{s}, \theta, \lambda, \hat{\mu}_j)$

1060 Now, consider:

$$\begin{aligned} \bar{V}(s, \hat{\theta}) - \bar{V}(s, \theta^*) &= \bar{V}(s, \hat{\theta}) - \hat{V}(s, \hat{\theta}) + \hat{V}(s, \hat{\theta}) - \hat{V}(s, \theta^*) \\ &\quad + \hat{V}(s, \theta^*) - \bar{V}(s, \theta^*) \\ &= \bar{V}(s, \hat{\theta}) - \hat{V}(s, \hat{\theta}) + \hat{V}(s, \theta^*) - \bar{V}(s, \theta^*) \\ &\geq -2 \sup_{\theta \in \Theta} |\bar{V}(s, \theta) - \hat{V}(s, \theta)| \end{aligned} \quad (4)$$

1061 The first step follows by adding and subtracting the same  
 1062 term. In the second step, we have used the fact that Assump-  
 1063 tions 1 and 2 imply that  $\hat{V}(s, \hat{\theta}) = \hat{V}(s, \theta^*)$ . In the third step,  
 1064 we have replaced the discrepancy between the sample average  
 1065 and the population average at specific points  $\hat{\theta}, \theta^*$  with the  
 1066 uniform bound over the parameter set  $\Theta$ .

1067 We use the Rademacher complexity bounds to bound this  
 1068 term. By [Shalev-Shwartz and Ben-David, 2014, Lemma  
 1069 26.2], we show the following:

1070 Let  $S$  denote the random training sample  $(\hat{\mu}_1, \dots, \hat{\mu}_n)$  and  
 1071  $P_0$  denote the uniform distribution  $\text{Unif}(\{-1, 1\}^n)$ . Then, for  
 1072 some numerical constant  $C$ , we have:

$$\mathbb{E}_S \sup_{\theta \in \Theta} |\bar{V}(s, \theta) - \hat{V}(s, \theta)| \leq C \mathbb{E}_S \mathcal{R}(\Theta \circ S)$$

1073 Where,  $\mathcal{R}(\Theta \circ S)$  is the Rademacher complexity:

$$\begin{aligned} \mathcal{R}(\Theta \circ S) &:= \\ &\frac{1}{n} \mathbb{E}_{\sigma \sim P_0} \sup_{\theta \in \Theta} \sum_{i=1}^n \sigma_i [\inf_{\lambda} V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i) - \mathbb{E}_{\hat{\mu}} \inf_{\lambda} V(\mathbf{s}, \theta, \lambda, \hat{\mu})] \end{aligned}$$

1074 Thus, to demonstrate the result, it is sufficient to show that:

$$\mathcal{R}(\Theta \circ S) \leq \frac{C_{\text{sys}} \text{polylog}(Nn)}{\sqrt{nN}} \quad (5)$$

1075 We will dedicate the rest of this section to demonstrate  
 1076 Equation (5). First we will state a useful Lemma which fol-  
 1077 lows from [Vershynin, 2018, Lemma 1.2.1]

1078 **Lemma 1.** Suppose a positive random variable  $X$  satisfies:  
 1079  $\mathbb{P}(X > t) \leq A \exp(-\frac{t^2}{2B})$  for some  $B > 0$ ,  $A > e$  and for  
 1080 every  $t \geq 0$  then for some numerical constant  $C$ , we have:

$$\mathbb{E}[X] \leq C \sqrt{B \log A}$$

*Proof.* From [Vershynin, 2018, Lemma 1.2.1], we have:

$\mathbb{E}X = \int_0^\infty \mathbb{P}(X > t) dt$ . Thus, we conclude:

$$\begin{aligned} \mathbb{E}X &\leq \int_0^\infty \min(1, A \exp(-\frac{t^2}{2B})) dt \\ &= \sqrt{2B \log A} + \int_{\sqrt{2B \log A}}^\infty A \exp(-\frac{t^2}{2B}) dt \\ &= \sqrt{2B \log A} + \int_0^\infty A \exp(-\frac{(t + \sqrt{2B \log A})^2}{2B}) dt \\ &\leq \sqrt{2B \log A} + \int_0^\infty \exp(-\frac{t^2}{2B}) dt \\ &\leq \sqrt{2B \log A} + \sqrt{2\pi B} \end{aligned} \quad (6)$$

In the fourth step we have used the fact that  $\exp(-(a + b)^2) \leq \exp(-a^2 - b^2)$  whenever  $a, b > 0$ .  $\square$

Define

$$v_i(\theta) := [\inf_{\lambda} V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i) - \mathbb{E}_{\hat{\mu}} \inf_{\lambda} V(\mathbf{s}, \theta, \lambda, \hat{\mu})].$$

We have the following lemma controlling how large  $v_i$  is for  
 any given  $\theta$ .

**Lemma 2.** For any  $\delta > 0$ , with probability at-least  $1 - \delta$ ,

$$\sup_i |v_i(\theta)| \leq \sqrt{\frac{C_{\text{sys}} \log(\frac{Nn}{\delta})}{N}}$$

Where  $C_{\text{sys}}$  depends on the system parameters.

*Proof.* First, we note that:

$$\begin{aligned} &|\inf_{\lambda} V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i) - \mathbb{E}_{\hat{\mu}} \inf_{\lambda} V(\mathbf{s}, \theta, \lambda, \hat{\mu})| \\ &\leq \mathbb{E}_{\hat{\mu}} |\inf_{\lambda} V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i) - \inf_{\lambda} V(\mathbf{s}, \theta, \lambda, \hat{\mu})| \\ &\leq \mathbb{E}_{\hat{\mu}} \sup_{\lambda \in [0, \lambda_{\max}]} |V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i) - V(\mathbf{s}, \theta, \lambda, \hat{\mu})| \\ &\leq \sup_{\lambda \in [0, \lambda_{\max}]} |V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i) - \mathbb{E}[V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i)]| \\ &\quad + \mathbb{E}_{\hat{\mu}} \sup_{\lambda \in [0, \lambda_{\max}]} |V(\mathbf{s}, \theta, \lambda, \hat{\mu}) - \mathbb{E}[V(\mathbf{s}, \theta, \lambda, \hat{\mu})]| \end{aligned} \quad (7)$$

In the last step, we have used the fact that  $\hat{\mu}$  and  $\hat{\mu}_i$   
 are identically distributed and hence  $\mathbb{E}[V(\mathbf{s}, \theta, \lambda, \hat{\mu})] =$   
 $\mathbb{E}[V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i)]$ . Note that by definition, the value func-  
 tion  $V(s, \theta, \lambda, \hat{\mu}_i) = \frac{1}{N} \sum_{j=1}^N V(s_j, \theta, \lambda, \mathbf{z}_j)$ . Thus, it is  
 clear that  $|V(s, \theta, \lambda, \hat{\mu}_i)| \leq A \frac{1 + \lambda_{\max}}{(1 - \beta)} =: V_{\max}$  where  $A$   
 is a constant which depends on the cost parameters  $c_j, \frac{B}{N}$   
 and the maximum reward. Take  $\hat{\mu}_i := (\mathbf{z}_1^{(i)}, \dots, \mathbf{z}_N^{(i)})$   
 and  $\hat{\mu} := (\mathbf{z}_1, \dots, \mathbf{z}_N)$ .

Thus, for a given  $\lambda$ , we have:  $V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i) -$   
 $\mathbb{E}V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i)$  has zero mean and

$$\begin{aligned} &V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i) - \mathbb{E}V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i) \\ &= \frac{1}{N} \sum_{j=1}^N [V(s_j, \theta, \lambda, \mathbf{z}_j^{(i)}) - \mathbb{E}V(s_j, \theta, \lambda, \mathbf{z}_j^{(i)})] \end{aligned} \quad (8)$$

It is an average of  $N$  i.i.d. zero mean random variables, bounded almost surely by  $2V_{\max}$ . Therefore, using the Azuma-Hoeffding inequality ([Vershynin, 2018]), we have:

$$\mathbb{P}(|V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i) - \mathbb{E}V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i)| > t) \leq C \exp(-\frac{c_1 N t^2}{V_{\max}^2})$$

1094 Only in this proof, let  $|V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i) - \mathbb{E}V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i)| =:$   
1095  $H(\lambda)$  for the sake of clarity. Let  $B \subseteq [0, \lambda_{\max}]$  be any finite  
1096 subset. Then, by union bound, we have:

$$\mathbb{P}\left(\sup_{\lambda \in B} H(\lambda) > t\right) \leq C|B| \exp(-\frac{c_1 N t^2}{V_{\max}^2}) \quad (9)$$

1097 Suppose  $B$  is an  $\epsilon$ -net for the set  $[0, \lambda_{\max}]$  for some  $\epsilon > 0$ .  
1098 This can be achieved with  $|B| = \frac{\lambda_{\max}}{\epsilon}$ . Let  $f : [0, \lambda] \rightarrow B$   
1099 map  $\lambda$  to the closest element in  $B$

$$\begin{aligned} \sup_{\lambda \in [0, \lambda_{\max}]} H(\lambda) &= \sup_{\lambda \in [0, \lambda_{\max}]} H(f(\lambda)) + H(\lambda) - H(f(\lambda)) \\ &\leq \sup_{\lambda \in [0, \lambda_{\max}]} H(f(\lambda)) + 2L\epsilon \\ &\leq \sup_{\lambda \in B} H(\lambda) + 2L\epsilon \end{aligned} \quad (10)$$

1100 Taking  $\epsilon = \frac{1}{\sqrt{N}}$ , we conclude from Equation (9) that with  
1101 probability at-least  $1 - \delta$ :

$$\sup_{\lambda \in [0, \lambda_{\max}]} H(\lambda) \leq C_{\text{sys}} \sqrt{\frac{\log(\frac{N}{\delta})}{N}}$$

The same concentration bounds hold for  $\sup_{\lambda} |V(\mathbf{s}, \theta, \lambda, \hat{\mu}) - \mathbb{E}V(\mathbf{s}, \theta, \lambda, \hat{\mu})|$  and integrating the tails (Lemma 1), we bound obtain the bound:

$$\sup_{\lambda \in [0, \lambda_{\max}]} |V(\mathbf{s}, \theta, \lambda, \hat{\mu}) - \mathbb{E}V(\mathbf{s}, \theta, \lambda, \hat{\mu})| \leq C_{\text{sys}} \sqrt{\frac{\log N}{N}}$$

1102 Applying a union bound over  $i = 1, \dots, n$ , conclude the  
1103 result.  $\square$

1104 We state the following folklore result regarding concentra-  
1105 tion of i.i.d. Rademacher random variables.

**Lemma 3.** *Given constants  $a_1, \dots, a_n \in \mathbb{R}$ , and  $\sigma_1, \dots, \sigma_n$  i.i.d Rademacher random variables, then for any  $\delta > 0$ , we have with probability at-least  $1 - \delta$ :*

$$\sum_{i=1}^n \sigma_i a_i \leq C \sqrt{\sum_i a_i^2} \sqrt{\log(\frac{1}{\delta})}$$

1106 Where  $C$  is a numerical constant

1107 We are now ready to prove Equation (5) and hence  
1108 complete the proof of Proposition 1. Given a data  
1109 set  $\hat{\mu}_1, \dots, \hat{\mu}_n$  and  $\theta \in \Theta$ , we let  $v_i(\theta) :=$   
1110  $[\inf_{\lambda} V(\mathbf{s}, \theta, \lambda, \hat{\mu}_i) - \mathbb{E}_{\hat{\mu}} \inf_{\lambda} V(\mathbf{s}, \theta, \lambda, \hat{\mu})]$ . Given a finite set  
1111  $\hat{\Theta} := \{\theta_1, \dots, \theta_H\} \subseteq \Theta$ , from Lemma 2, we have with prob-  
1112 ability  $1 - \delta$ ,

$$\sup_{\theta \in \hat{\Theta}} \sup_i |v_i(\theta)| \leq \sqrt{\frac{C_{\text{sys}} \log(\frac{nN|\hat{\Theta}|}{\delta})}{N}} =: R(\delta)$$

Therefore, with probability at-least  $1 - \delta$  over the random- 1113  
ness in  $\hat{\mu}_1, \dots, \hat{\mu}_n$ , we have: 1114

$$\begin{aligned} \mathbb{P}\left(\sup_{\theta \in \hat{\Theta}} \sum_{i=1}^n \sigma_i v_i(\theta) > t \mid \hat{\mu}_1, \dots, \hat{\mu}_n\right) \\ \leq C_1 |\hat{\Theta}| \exp\left(-\frac{C_2 t^2}{nR^2(\delta)}\right) \end{aligned} \quad (11)$$

We pick  $\hat{\Theta}$  to be an  $\epsilon$  net over  $\Theta$ . By [Vershynin, 2018, Corollary 4.2.13], we can take  $|\hat{\Theta}| \leq (\frac{3D}{\epsilon})^d$ . Let  $f : \Theta \rightarrow \hat{\Theta}$  be the map to its nearest element in  $\hat{\Theta}$ . Now, we have:

$$\begin{aligned} \sup_{\theta \in \Theta} \sum_i v_i(\theta) \sigma_i &= \sup_{\theta \in \Theta} \sum_i v_i(f(\theta)) \sigma_i + [v_i(\theta) - v_i(f(\theta))] \sigma_i \\ &\leq 2n\epsilon L + \sup_{\theta \in \hat{\Theta}} \sum_i v_i(\hat{\theta}) \sigma_i \end{aligned}$$

Combining this with Equation (11), we conclude that with 1115  
 $1 - \delta$  over the randomness in  $\hat{\mu}_1, \dots, \hat{\mu}_n$ , we have: 1116

$$\begin{aligned} \mathbb{P}\left(\sup_{\theta \in \Theta} \sum_{i=1}^n \sigma_i v_i(\theta) > t + 2n\epsilon L \mid \hat{\mu}_1, \dots, \hat{\mu}_n\right) \\ \leq C_1 |\hat{\Theta}| \exp\left(-\frac{C_2 t^2}{nR^2(\delta)}\right) \end{aligned} \quad (12)$$

Taking  $\epsilon = \frac{1}{n^{\frac{3}{2}} \sqrt{N}}$  and integrating the tails (Lemma 1), we 1117  
conclude that with probability at-least  $1 - \delta$  (with respect to 1118  
the randomness in  $\hat{\mu}_1, \dots, \hat{\mu}_N$ ). 1119

$$\mathbb{E}[\sup_{\theta \in \Theta} \sum_{i=1}^n \sigma_i v_i(\theta) \mid \hat{\mu}_1, \dots, \hat{\mu}_n] \leq C_{\text{sys}} \frac{R(\delta)}{\sqrt{n}} \text{polylog}(Nn)$$

Define the random variable

$$X := \mathbb{E}[\sup_{\theta \in \Theta} \sum_{i=1}^n \sigma_i v_i(\theta) \mid \hat{\mu}_1, \dots, \hat{\mu}_n]$$

Using the definition of  $R(\delta)$ , we have: 1120

$$\mathbb{P}(X > t) \leq C_1 \exp(-\frac{t^2 n N}{C_{\text{sys}} \text{polylog}(Nn)}).$$

We then apply Lemma 1 to the equation above to bound 1121  
 $\mathbb{E}X$  and conclude Equation (5). 1122

## E Proof for $\lambda$ -network Update Rule and Convergence 1123 1124

*Proof of Proposition 2.* We first consider a simple setting, where the opt-in and opt-out decisions of arms are fixed before training. Taking the derivative of the objective (Eq 2) with respect to  $\lambda$ , we obtain:

$$\frac{B}{1-\beta} - \sum_{i=1}^N \mathbb{E} \left[ \sum_{\substack{t \in [H] \\ \text{arm } i \text{ opts-in at } t}} \beta^t c_{i,t} + \sum_{\substack{t \in [H] \\ \text{arm } i \text{ opts-out at } t}} \beta^t c_{0,t} \right].$$



Now consider the general case that the opt-in and opt-out decisions are updated at each round during the training. We have

$$\Lambda_t = \Lambda_{t-1} - \alpha \left( \frac{B}{1-\beta} \right) + \alpha \left( \sum_{i=1}^N \mathbb{E} \left[ \sum_{t=0}^H \mathbb{I}\{\xi_{i,t} = 1\} \beta^t c_{i,t} + \mathbb{I}\{\xi_{i,t} = 0\} \beta^t c_{0,t} \right] \right),$$

1125 where the expectation is over the random variables  $\xi_{i,t}$  and  
 1126 the action chosen by the optimal policy. Rearranging and sim-  
 1127 plifying the right hand side terms, we obtain the  $\lambda$ -updating  
 1128 rule.  $\square$

1129 *Proof of Proposition 3.* The proof largely follows the proof  
 1130 of Proposition 2 in Killian *et al.*[2022].

1131 Since the max of piece-wise linear functions is a convex  
 1132 function, Equation 2 is convex in  $\lambda$ . Thus, it suffices to show  
 1133 (1) the gradient estimated using Proposition 2 is accurate and  
 1134 (2) all inputs (states, features, opt-in decisions) are seen in-  
 1135 finitely often in the limit. For (1), we note that training the  
 1136 policy network for a sufficient number of epochs under a fixed  
 1137 output of the  $\lambda$ -network ensures that Q-value estimates are ac-  
 1138 curate. With accurate Q-functions and corresponding optimal  
 1139 policies, the sampled cumulative sum of action costs is an un-  
 1140 biased estimator of expected cumulative sum of action costs.  
 1141 Critically, for the estimator to be unbiased, we do not strictly  
 1142 enforce the budget constraint during training, as in Killian  
 1143 *et al.*[2022]. In inference, we do strictly enforce the budget  
 1144 constraint. For (2), we note that during training, initial states  
 1145 are uniformly sampled, and opt-in decisions are also sampled  
 1146 from a fixed bernoulli distribution. For arms that newly opt-in,  
 1147 the features are uniformly sampled. Thus, both (1) and (2) are  
 1148 achieved.  $\square$