Robust Planning over Restless Groups: Engagement Interventions for a Large-Scale Maternal Telehealth Program

Jackson A. Killian^{*},¹ Arpita Biswas^{*},¹ Lily Xu^{*},¹ Shresth Verma^{*},² Vineet Nair,² Aparna Taneja,² Aparna Hegde,³ Neha Madhiwalla,³ Paula Rodriguez Diaz,¹ Sonja Johnson-Yu,¹ Milind Tambe^{1,2}

¹Harvard University, ²Google Research, ³ARMMAN

{jkillian,arpitabiswas,lily_xu}@g.harvard.edu, {vermashresth, vineetn, aparnataneja}@google.com, {aparnahegde, neha}@armman.org, {prodriguezdiaz, sjohnsonyu}@g.harvard.edu, milindtambe@google.com

Abstract

In 2020, maternal mortality in India was estimated to be as high as 130 deaths per 100K live births, nearly twice the UN's target. To improve health outcomes, the non-profit AR-MMAN sends automated voice messages to expecting and new mothers across India. However, 38% of mothers stop listening to these calls, missing critical preventative care information. To improve engagement, ARMMAN employs health workers to intervene by making service calls, but workers can only call a fraction of the 100K enrolled mothers. Partnering with ARMMAN, we model the problem of allocating limited interventions across mothers as a restless multi-armed bandit (RMAB), where the realities of large scale and model uncertainty present key new technical challenges. We address these with GROUPS, a double oracle-based algorithm for robust planning in RMABs with scalable grouped arms. Robustness over grouped arms requires several methodological advances. First, to adversarially select stochastic group dynamics, we develop a new method to optimize Whittle indices over transition probability intervals. Second, to learn grouplevel RMAB policy best responses to these adversarial environments, we introduce a weighted index heuristic. Third, we prove a key theoretical result that planning over grouped arms achieves the same minimax regret-optimal strategy as planning over individual arms, under a technical condition. Finally, using real-world data from ARMMAN, we show that GROUPS produces robust policies that reduce minimax regret by up to 50%, halving the number of preventable missed voice messages to connect more mothers with life-saving maternal health information.

1 Introduction

Maternal mortality, the death of a mother¹ during pregnancy or within 42 days after childbirth, is an ongoing global health crisis. In India, the maternal mortality rate is particularly stark, estimated between 99 and 130 deaths per 100K births in 2020 (Meh et al. 2021; Gates Foundation 2021), significantly higher than Sustainable Development Goal 3.1



Figure 1: Mothers enrolled with ARMMAN receive life-saving preventative care information via voice messages throughout their pregnancy, childbirth, and neonatal period. Photo courtesy of ARMMAN.

target of 70 per 100K births (United Nations 2021). Tragically, most maternal deaths are preventable (HLPF Review of SDG3 2017), but lack of finances and awareness prevent mothers from seeking care, particularly in low-income communities (Carvalho, Salehi, and Goldie 2013).

To improve maternal health outcomes, we work with AR-MMAN, an India-based non-profit that provides free preventive care to millions of mothers by sending automated health voice messages, specifically targeted towards low-income communities (similar to MAMA (MomConnect 2021)). Mothers enrolled in the program receive weekly automated voice messages during pregnancy and up to one year after childbirth. Randomized control trials showed that ARM-MAN's messaging program significantly improves key indicators including treatment-seeking during complications, infant breastfeeding, and post-infancy weight (Murthy et al. 2019). However, ARMMAN found that nearly 38% of mothers disengage, missing critical health information. To improve engagement, ARMMAN employs health workers to provide service calls, but there are only tens of health workers compared to hundreds of thousands of mothers in a given service area — so interventions must be carefully targeted to maximize engagement.

Working with ARMMAN, we model this resource-limited intervention planning problem as a restless multi-armed bandit (RMAB), where each mother (arm) changes their weekly engagement (state) according to a stochastic Markov decision process. RMABs are PSPACE-hard to solve exactly (Papadimitriou and Tsitsiklis 1999) and even the more tractable, asymptotically optimal "Whittle index policy" (Whittle 1988) is challenging to compute at scale.

To improve the scalability of real-world RMAB planning, Mate et al. (2022) proposed to organize arms into a small number of *groups*, infer transition dynamics from

^{*}These authors contributed equally.

Copyright © 2023, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

¹We recognize that the term "mother" is imperfect, most notably by not reflecting transgender and non-binary identities. We highlight alternative language with discussion in Appendix A.

each group's data, then compute the Whittle index policy per group. While the scalability of their method is desirable for ARMMAN's problem setting, it ignores a key reality of *model uncertainty*: learning transition probabilities from historical data leads to imprecise and imperfect estimates (Sinha and Mahajan 2022) which must be accounted for in planning. Computing RMAB policies that are robust to model uncertainty has only recently been studied. Existing methods achieve robustness to *interval uncertainty* over model dynamics by planning against a model-controlling "nature" adversary to yield policies that minimize max regret (Killian et al. 2022; Xu et al. 2021). Robustness is desirable for ARMMAN's setting, but these methods require training deep reinforcement learning (RL) agents for each arm, so unfortunately do not scale past hundreds of arms.

To enable large-scale, robust intervention planning for ARMMAN, we bridge the gaps in previous works by introducing robust grouped RMAB. Our model achieves scalability by considering a grouped-arm paradigm and optimizing for minimax regret over the uncertain model dynamics per group. Unfortunately, the grouping abstraction breaks key assumptions used in previous robust RMAB work: that (1) policies improve by collecting samples of regret by evolving a joint state of all arms, and (2) the nature adversary controls the transitions of each arm individually. We overcome (1) by *decomposing regret per arm*, freeing the planner from relying on a cumbersome joint state to enable efficient group-abstracted planning. For (2), we prove that *restricting* the adversary to control dynamics only over groups does not change the equilibrium strategy, allowing us to leverage the scalable robust grouped model to find policies over hundreds of thousands of arms without sacrificing quality.

Our contributions are as follows. First, we introduce robust grouped RMABs with a minimax regret objective and propose a solution that employs the double oracle framework (McMahan, Gordon, and Blum 2003). The approach we propose is GROUPS: Group RMAB Oracles for Uncertainty-robust Planning at Scale. Second, we develop novel methods designed for robust grouped RMABs to implement the two oracles, the planner and adversary. Planning over groups of arms allows large scale-up but presents several new algorithmic challenges as we detail above. Third, we prove that the minimax regret-optimal strategy is the same whether the planner and adversary play at the individual or group level. Our proof enables massive scale-up as it is now sufficient to compute robust strategies over groups, instead of over individual arms. Finally, we demonstrate empirically on real data that GROUPS reduces worst-case regret up to 50% compared to baselines, representing potentially thousands of additional engagements with lifesaving information. We are working with ARMMAN to deploy GROUPS to positively impact maternal health.

2 Related Work

Mobile-based maternal health services are effective and affordable in low- and middle-income communities (Watterson, Walsh, and Madeka 2015; Tamrat and Kachnowski 2012). Successful programs include MatHealth in Uganda (Musiimenta et al. 2021), Aponjon in Bangladesh (Alam et al. 2017), ARMMAN in India (Murthy et al. 2019), and text4baby in the United States (Evans, Wallace, and Snider 2012). Our work is designed to support such programs.

Whittle (1988) introduced RMABs and proposed the *Whittle index policy*, which computes indices estimating each arm's "return on investment" then acts on arms with the top K. Weber and Weiss (1990) showed this policy is asymptotically optimal under a technical condition. Many RMAB studies assume known transition dynamics, although some recent works design methods to learn policies online (Wang, Huang, and Lui 2020; Nakhleh et al. 2021; Biswas et al. 2021; Killian et al. 2021; Wang et al. 2022). However, these online approaches require collecting a prohibitively large number of samples, limiting their real-world applicability in scenarios where the time horizon is short.

Most robust planning papers consider single-MDP (one arm) settings (Pinto et al. 2017; Lanctot et al. 2017; Li et al. 2019), rather than the budget-coupled N-MDP setting of RMAB. Even for single MDPs, optimizing criteria such as minimax regret (Braziunas and Boutilier 2007) requires searching massive strategy spaces; double oracle (McMahan, Gordon, and Blum 2003) is one approach to do so efficiently. Recent work combines double oracle with deep RL to solve for minimax regret-optimal robust policies for single MDPs (Xu et al. 2021). Killian et al. (2022) extended the idea to solve larger RMABs. Both Xu et al. (2021) and Killian et al. (2022) use deep RL which, if applied to a group setting, would need to explicitly account for the size of each group and state of each arm within each group, limiting their methods' ability to scale beyond hundreds of arms. For the large problem size that ARMMAN faces, our methods must scale to hundreds of thousands of arms.

Finally, robust planning for *stochastic* bandits is well studied (Maillard 2013; Huo and Fu 2017) However, stochastic bandits are stateless and lack passive rewards, and so are not expressive enough to model ARMMAN's setting.

3 Model

We consider grouped RMABs where N arms (enrolled mothers) comprise M groups. Each arm $n \in [N]$ follows an MDP $\langle S, A, P^n, r, \gamma \rangle$ where $s \in S := \{0, 1\}$ is the state space indicating whether a mother is engaging $(s_n = 1)$ or not engaging $(s_n = 0)$ with automated voice messages; r(s) = s is the reward function; $a \in \mathcal{A} := \{0, 1\}$ is the action space, i.e., {not intervene, intervene}; $P^n(s, a, s')$ is the probability that arm n transitions from state s to s' given action a; $\gamma \in [0, 1]$ is the discount factor. Let $s \in S^N$ and $a \in \mathcal{A}^N$ be the combined state and action vectors of all arms. At each timestep t, the task is to choose K mothers to intervene on (deliver service calls to) given the state s_t at time t.

Formally, we compute RMAB policies $\pi : S^N \to \mathcal{A}^N$ that respect a budget constraint $\|\pi(s_t)\|_1 = K$ for all t. For a given policy π and a fixed environment $P := \{P^n\}_{n \in [N]}$ representing a matrix of transition probabilities of all arms, the average discounted reward is $G(\pi, P) := \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_t) | \pi, P]$. Given P, the optimal policy which maximizes reward is $\pi_P^* := \max_{\pi} G(\pi, P)$. An asymptotically optimal RMAB policy is the Whittle index pol-

icy (WIP), which computes the Whittle index $W_P^n(s)$ for each arm n and state s, then intervenes on the arms with the greatest K indices. The Whittle index represents "return on investment," interpreted as a charge for acting that makes *no intervention* equally valuable as *intervention* in the long term. Let $Q_P^n(s, a, \lambda) = r(s) - \lambda a + \gamma \mathbb{E}_{s' \in S}[\max_{a' \in A} Q_P^n(s', a', \lambda)]$ be the long-term expected value of action a on arm n in state s. Then, for a given P, the Whittle index for arm n at state s is $W_P^n(s) = \min{\{\lambda : Q_P^n(s, 1, \lambda) = Q_P^n(s, 0, \lambda)\}}$.

Grouped RMAB For scalability, we organize arms into groups, extending the concept from Mate et al. (2022) to our more challenging *robust* setting, e.g., by clustering based on historical engagement patterns. We then estimate uncertainty intervals over transition probabilities per group. However, note that our robust policy computation steps in Section 4 are agnostic to the particular grouping and interval estimation methods. Let $\phi : [N] \rightarrow [M]$ be a surjective mapping of arms to groups and $\phi^{-1}(m)$ be the set of arms in group m. The uncertainty intervals are $\overline{P}_{s,a,s'}^m := [\underline{P}_{s,a,s'}^m, \overline{P}_{s,a,s'}^m]$ for all m, s, a, s'. Then let $\overline{\underline{P}}^m := \{\overline{\underline{P}}_{s,a,s'}^m\}_{s,a,s'}$ be the interval uncertainty matrix for group m across all states and actions. Importantly, though arms in the same group have the same uncertainty intervals, they may not have the same instantiated probabilities within those intervals.

Minimax regret We define regret for grouped RMAB as:

$$R(\pi, P) := G(\pi_P^{\star}, P) - G(\pi, P) , \qquad (1)$$

where P instantiates $P^m \in \overline{\underline{P}}^m$ for all groups $m \in [M]$. Our objective is to learn a policy π that minimizes max regret:

$$\min_{\pi} \max_{P} R(\pi, P) . \tag{2}$$

We choose minimax regret as our robust objective since it does not require probability distributions over the uncertainty intervals (Braziunas and Boutilier 2007). Such distributional information is scarce in our setting where $K \ll N$, giving us few samples of transitions for action a = 1.

4 Methodology

We introduce GROUPS (Group RMAB Oracles for Uncertainty-robust Planning at Scale), a four-step approach visualized end-to-end in Fig. 2. Step (3) is our key algorithmic contribution. In step (1), similar arms (mothers) are mapped into groups. In step (2), we combine data from arms in each group with historical engagement data, using bootstrapping to estimate uncertainty intervals $\overline{\underline{P}}^m$ for each group (Schomaker and Heumann 2018). In step (3), we compute a minimax regret-optimal policy over groups, where arms in a given group are treated as having the same transition probabilities, greatly improving computational efficiency. Critically, we show in Section 5 that this group-level planning is lossless - i.e., the policies we compute are the same minimax regret-optimal policies as would be computed if grouped arms were allowed different transition probabilities (within the same uncertainty intervals). In step (4), we map group-level policies back to individual-level policies by computing Whittle indices for each group $m \in [M]$, then assigning an index to each arm nwithin that group based on its current state s_n . Our policy is to intervene on mothers with the top K indices.

Double oracle In step (3), we adopt a double oracle (DO) framework (McMahan, Gordon, and Blum 2003), solving Eq. 2 by formulating the problem as a two-player zerosum game between the RMAB planner and nature adversary, where the players aim to minimize and maximize regret respectively. The planner's *pure strategy* space is the finite set of all feasible RMAB policies π ; the adversary has the continuous space of transition probabilities P within the uncertainty intervals $\overline{\underline{P}}^m$ for all $m \in [M]$. The algorithm maintains a finite pure strategy set for each player. For each iteration, we compute a mixed strategy Nash equilibrium (MSNE) on the game over the finite strategy sets. A *mixed* strategy is a probability distribution over pure strategies. In each iteration, the planner oracle computes a best response pure strategy π against the adversary's mixed strategy; π is added to the planner's finite strategy set. We follow a symmetric approach to compute a best response P for the adversary. Upon termination, we return the final planner mixed strategy, which is guaranteed, under mild conditions, to be an ϵ -optimal minimax solution (Xu et al. 2021). In practice, we terminate after T iterations (Lanctot et al. 2017). The key technical challenge of using the double oracle approach is designing *planner* and *adversary* oracles for *group* RMABs.

4.1 Planner Oracle: WI for Mixed Strategy

An adversary mixed strategy β contains tuples (P_i, β_i) where β_i is the probability of playing pure strategy P_i . Similarly, a planner mixed strategy α contains tuples (π_i, α_i) where α_i is the probability of playing pure strategy π_i .

The planner oracle must compute an intervention policy π that minimizes regret with respect to a given adversary mixed strategy β over environment settings P_i . Since β and thus all P_i are fixed, and only the second term of regret in Eq. 1 depends on π , minimizing regret is equivalent to maximizing reward, to ensure that mothers engage with as many voice messages as possible. However, existing rewardmaximizing RMAB algorithms assume a *single* environment P_i , versus a mixed strategy β over multiple P_i . To address this combinatorially hard problem, we develop a new heuristic approach that computes well-performing policies π based on strategically weighted combinations of Whittle indices.

Unfortunately, optimizing *exact* regret is at least PSPACE-hard (Papadimitriou and Tsitsiklis 1999). Previous work optimized regret of the Lagrange relaxation (Killian et al. 2022), but relied on joint arm states which does not scale. We introduce a decomposed notion of regret, allowing us to optimize regret of the full RMAB in a far more scalable way. We call this Whittle index regret: the sum of Whittle indices played by a policy π compared to the optimal WIP. The key is that the Whittle index is a measure of "reward if played" — so agents who play arms with low Whittle indexes in lieu of arms with high Whittle indexes will incur large regret. As a further advantage, this regret notion naturally extends to groups — since the Whittle index is a furch-



Figure 2: GROUPS pipeline for robust grouped RMABs. (1) Assign enrolled mothers (arms) to groups. (2) Estimate uncertainty intervals over transition probabilities. (3) Novelty of this work: Compute robust minimax regret–optimal policy via double oracle, where each oracle efficiently searches the large-scale strategy spaces by using the group abstraction. (4) To execute policies, translate group-level indices \tilde{I}_s^m to arm-level intervention policy.

tion only of transition probabilities and rewards, all of which are shared in a group under P_i — improving scaling.

Given states s, denote the set of arms pulled by policy π as $\Phi^{\pi}(s) = \{n \in [N] : \pi_n(s) = 1\}$ where $\pi_n(s)$ is the action on arm n. The planner's Whittle index regret $R_W^{planner}(s)$ is:

$$\sum_{\substack{(P_i,\beta_i)\\|\kappa|=K}} \beta_i \left[\max_{\substack{\kappa \subseteq [N]\\|\kappa|=K}} \left\{ \sum_{n \in \kappa} (W_{P_i}^n(\boldsymbol{s}^n)) \right\} - \sum_{n \in \Phi^{\pi}(\boldsymbol{s})} W_{P_i}^n(\boldsymbol{s}^n) \right].$$
(3)

The first term in Eq. 3 corresponds to a planner's optimal mixed strategy which plays the WIP corresponding to each setting of transition probabilities P_i in β . To minimize regret $R_W^{planner}$, we seek a policy π that plays Whittle indices as close as possible to the WIPs in the first term, which equivalently maximizes the second term. How to produce a *pure* strategy π that closely follows the *mixed* WIP policies of the first term is the key challenge. We start by making the first term more closely computable as a pure strategy with a relaxation that leads to relaxed regret, by moving the expectation over β_i inside the max over indices:

$$\max_{\substack{\kappa \subseteq [N] \\ |\kappa| = K}} \left\{ \sum_{n \in \kappa} \sum_{(P_i, \beta_i) \in \beta} \beta_i W_{P_i}^n(\boldsymbol{s}^n) \right\} .$$
(4)

We replace the first term of $R_W^{planner}(s)$ (from Eq. 3) with Eq. 4 to get $\tilde{R}_W^{planner}(s)$. This illuminates a heuristic for the planner oracle. Specifically, Eq. 4 can be computed exactly by a single policy π , meaning we can make $\tilde{R}_W^{planner}(s) = 0$ by finding a π equivalent to Eq. 4. To do so, we compute Whittle indices for each pure strategy P_i , compute the β_i weighted average index \tilde{I}_s^m for each group m and state s, then follow the greedy strategy of a WIP. Since the expectation over β_i is pushed through the max (Eq. 4) we have $\tilde{R}_W^{planner}(s) \leq R_W^{planner}$, but we show in appendix Fig. 4 that

Algorithm 1: WI4MS (Planner Oracle)

Input Adversary mixed strategy β

1: for $(P_i, \beta_i) \in \beta$ do // environment and probability *i* 2: for $\{m = 1 \text{ to } M\}$ and $\{s \in S\}$ do 3: $\tilde{I}[m, s] += \beta_i \times \text{COMPUTEWI}(m, s, P_i^m)$ 4: $\pi = \text{WIP}(\tilde{I})$ // implements Whittle index policy 5: return π // planner pure strategy

this weighted index policy performs well, despite this relaxation. We call this approach Whittle Index for Mixed Strategy (WI4MS), given in Alg. 1. Whittle indices are computed via COMPUTEWI described in Alg. 4 in the appendix.

4.2 Adversary Oracle: RegretMax Whittle Index

The adversary oracle must find one environment P that maximizes regret for the planner's current mixed strategy α over policies π_i to maximize the number of missed calls. To guide the search, we must address challenges both in maximizing regret of RMAB policies and in searching over a continuous strategy space \overline{P}^m . Our insight is to maximize regret by manipulating the optimal RMAB policy (a Whittle index policy) to simultaneously *minimize* the values of Whittle indices acted on by the planner and *maximize* indices that are not.

We utilize again the notion of Whittle index regret, redefined for the adversary oracle:

$$R_{W}^{adversary} = \mathbb{E}_{\boldsymbol{s}} \left[\sum_{n \in \Phi^{\pi_{P}^{\star}}(\boldsymbol{s})} W_{P}^{n}(\boldsymbol{s}^{n})) \mid \pi_{P}^{\star}, P \right]$$
$$- \sum_{(\pi_{i},\alpha_{i})\in\alpha} \alpha_{i} \left(\mathbb{E}_{\boldsymbol{s}} \left[\sum_{n \in \Phi^{\pi_{i}}(\boldsymbol{s})} W_{P}^{n}(\boldsymbol{s}^{n})) \mid \pi_{i}, P \right] \right). \quad (5)$$

Algorithm 2: RegretMaxWI (Adversary Oracle)

 \overline{P}^m , Mixed Input: strategies $(\alpha,\beta),$ intervals group-mean budget $K_M, P = []$

- 1: $\{L_s^m\}_{s\in\mathcal{S}}^{m\in[M]} = \text{MONTECARLO}(\alpha,\beta) \text{ // simulation}$ 2: $K_{\text{TH}} = \text{FINDTHRESH}(L,K_M) \text{ // returns action count}$ of $\lceil K_M \rceil^{\text{th}}$ group-state 3: for $\{m = 1 \text{ to } M\}$ and $\{s \in S\}$ do
- $obj[m,s] = \min \text{ if } (L_s^m \ge K_{\text{TH}}) \text{ else } \max$ 4:
- 5: for m = 1 to M do
- $P^m = MINMAXWHITTLEBQP(obj[m], \overline{P}^m)$ 6:
- 7: return P // Adversary pure strategy

Given an environment, P_i , Eq. 5 captures the difference in the Whittle indices collected by the optimal policy $\pi_{P_i}^{\star}$ versus the Whittle indices collected by the policies of the agent mixed strategies π_i . The WIP is a proxy for finding the most effective arms on which to intervene; intuitively, this means the adversary oracle should find P_i which maximizes the Whittle indices of arms played by the optimal policy but not played by the planner, and simultaneously minimizes the Whittle indices of arms played *only* by the planner policies.

The first challenge is to determine which arms the planner will act on in expectation. We propose a simple but effective solution which counts the number of times the armstate pairs are acted on during Monte Carlo simulation of the planner's mixed strategy. Since the adversary operates at the group level, we then aggregate arm-state counts into groupstate counts, denoted L_s^m for each group m and state s. The next question is which group-state indices to minimize or maximize. Intuitively, if we reduced all indices an equal amount, we would reduce reward but not regret since the optimal policy, i.e., the first term of Eq. 5, would reduce the same as the second. Thus, we need to strategically minimize some indices, but maximize others to induce an optimal policy that plays different arms. Specifically, we choose to minimize the indices of the top $K_M = \frac{K}{N/M}$ — i.e., the budget normalized by average group size — entries of L_s^m , approximating the top K choices of the agent mixed strategy in expectation. Then we maximize the Whittle indices of all group-state pairs below that threshold.

The second challenge is to find transition probabilities P that minimize or maximize the Whittle indices of a group over its transition probability intervals. This problem has general implications, e.g., for optimistic or pessimistic search over uncertainty sets in online learning. We derive a novel binary-quadratic program that, given a group and objective for each state (min, max, or null), computes a P^m that optimizes the indices for all states simultaneously, detailed in the appendix as MINMAXWHITTLEBQP (Eq. 18). We give the full adversary oracle algorithm, RE-GRETMAXWI, in Alg. 2 and empirically demonstrate its good performance in the appendix Fig. 5.

Theoretical Regret Guarantee 5

In Section 4, we proposed an approach to compute a minimax regret-optimal strategy against an adversary choosing the same transition probabilities for all arms in the same group from their corresponding intervals. However, arms within the same group may not have identical transition probabilities. Also, it is not intuitive that a minimax regretoptimal policy, when the adversary chooses the same transition probabilities for all the arms in a group, also minimizes max regret when the adversary chooses different transition probabilities for the arms in a group from their corresponding intervals. In this section, we show this is true under mild assumptions. In particular, the minimax regret-optimal strategy of the planner is the same against an adversary choosing transition probabilities at the group level as against an adversary choosing transition probabilities at the individual level.

Let Π be the planner's pure strategy space of all individual-level policies, i.e., all choices of subsets of arms with cardinality K. Then we define mixed strategy sets for the planner at *individual-level*, $\Delta_I(\Pi)$, and *group-level*, $\Delta_M(\Pi)$, where $\Delta_M(\Pi) \subseteq \Delta_I(\Pi)$ is a restricted set of mixed strategies in which the planner is indifferent between arms in the same group and state (see Appendix D.2 for definition). Next, let \mathscr{P} be the adversary's pure strategy space, containing all individual-level policies, i.e., choices of transition probabilities $\{P^n\}_{n \in [N]}$ respecting the given uncertainty intervals $\overline{P}^{\phi(n)}$. Similarly, we define mixed strategy sets for the adversary at individual-level, $\Delta_I(\mathscr{P})$, and group-level, $\Delta_M(\mathscr{P})$, where $\Delta_M(\mathscr{P}) \subseteq \Delta_I(\mathscr{P})$ is a restricted space that assigns same transition probabilities to all arms within a group.

For $X, Y \in \{I(individual), M(group)\}$, the regret game with X-level planner and Y-level adversary is noted as X/Y. The X/Y regret of a planner's mixed strategy $\alpha \in$ $\Delta_X(\Pi)$ against an adversary's mixed strategy $\beta \in \Delta_Y(\mathscr{P})$ is:

$$R(\alpha,\beta) := \sum_{i \in [|\Delta_X(\Pi)|]} \sum_{j \in [|\Delta_Y(\mathscr{P})|]} \alpha_i \beta_j R(\pi_i, P_j) ,$$

where α_i is the *i*th pure strategy of the X-level planner and β_i is the jth pure strategy of the Y-level adversary. Let α_{XY}^{\star} be the planner's mixed strategy of a X/Y game, defined:

$$\min_{\alpha \in \Delta_X(\Pi)} \max_{\beta \in \Delta_Y(\mathscr{P})} R(\alpha, \beta) = \max_{\beta \in \Delta_Y(\mathscr{P})} R(\alpha_{X,Y}^{\star}, \beta)$$

which holds since the regret game is a two-player zero sum game, making minimax regret equal to maximin reward. We call this the worst-case regret for $\alpha^{\star}_{X,Y}$.

We first show in Theorem 1^2 that, when all arms within the same group have the same transition intervals, the minimax I/I regret is equal to the minimax M/I regret.

Theorem 1. The worst-case regrets of $\alpha_{I,I}^{\star}$ and $\alpha_{M,I}^{\star}$ against an adversary operating at the individual level is equal:

$$\max_{\beta\in\Delta_{I}(\mathscr{P})}R(\alpha_{I,I}^{\star},\beta)=\max_{\beta\in\Delta_{I}(\mathscr{P})}R(\alpha_{M,I}^{\star},\beta)\;.$$

Similarly, in Theorem 2, we show that, when all arms within the same group have the same transition intervals, the minimax I/M regret is equal to the minimax M/M regret.

²Proofs of Theorem 1, 2, and 3 are given in Appendix E.

Theorem 2. The worst-case regrets of $\alpha_{I,M}^{\star}$ and $\alpha_{M,M}^{\star}$ against an adversary operating at the group level are equal:

$$\max_{\beta \in \Delta_M(\mathscr{P})} R(\alpha_{I,M}^{\star}, \beta) = \max_{\beta \in \Delta_M(\mathscr{P})} R(\alpha_{M,M}^{\star}, \beta)$$

Finally, we use these results to establish our main result in Theorem 3 that the worst-case regret of $\alpha_{M,M}^*$ is equal to the worst-case regret of $\alpha_{I,I}^*$ when (1) all arms in the same group have the same intervals and (2) there exists a surjective function ψ that maps $\Delta_I(\mathscr{P})$ to $\Delta_M(\mathscr{P})$ that preserves the regret ordering of planner and adversary strategies (formal definition and example ψ given in Appendix E.1).

Theorem 3. If there exists an order-preserving map, then the worst-case regret of $\alpha_{M,M}^{\star}$ is equal to that of $\alpha_{I,I}^{\star}$, against an individual-level adversary, that is,

$$\max_{\beta \in \Delta_I(\mathscr{P})} R(\alpha^{\star}_{M,M},\beta) = \max_{\beta \in \Delta_I(\mathscr{P})} R(\alpha^{\star}_{I,I},\beta) \ .$$

Theorems 1, 2, and 3 together establish that the minimax regret–optimal strategy is the same whether the planner and adversary play at individual or group level. In particular, this result ensures that, under some conditions, the minimax regret–optimal strategy obtained by our algorithm GROUPS, which implements group-level planner and adversary, is also minimax regret–optimal against an individual level adversary.

6 Experiments

6.1 Experiment Setup

ARMMAN maternal health domain Every week, ARM-MAN's automated system delivers prerecorded health messages to each enrolled mother with information tailored to the mother's gestational age. If mothers stop listening to the messages, healthcare workers can deliver interventions to try to improve mothers' engagement. We evaluate the increase in number of health messages mothers listen to using GROUPS to target interventions compared to existing baselines. To construct a simulation environment, we use a real anonymized dataset from ARMMAN's records of weekly program engagement data for 15,336 mothers (though we note that ARMMAN's larger service areas operate on the scale of hundreds of thousands). A mother is "engaged" if they listen to at least 30 seconds of a message that week. Thus, states are {not engaged, engaged} with rewards 0 and 1, respectively. To create an arm-group mapping, we run K-means clustering on the engagement data and compute uncertainty intervals via bootstrapping followed by multiple imputation to compute standard deviations of the means (Schomaker and Heumann 2018). Statistics on the uncertainty intervals and group sizes are shown in appendix Figs. 9 and 10. For details on the dataset and consent for collection, see appendix K.

In the experiments, the default parameters match the intervention setup used by ARMMAN, i.e., budget K = 100, N = 15,320 mothers, and M = 40 groups. For sensitivity analysis, we vary the budget, horizon, and number of mothers. Additional analysis varying uncertainty interval width, number of groups, and distribution of group sizes are included in appendix Fig. 7. Additional domains To demonstrate wider applicability, we include results from two additional domains. The TB domain is constructed from an anonymized dataset of daily adherence to tuberculosis medication (Killian et al. 2019). States, rewards, and groups were derived analogously to the maternal health setting; complete details are in appendix L, including group statistics in Figs. 11 and 12. In our experiments, the default setting has N = 8,350 arms, M = 60groups, budget K = N/10, and $A_{\sigma} = 3$, i.e., interval width of 3 standard deviations. We vary the budget, number of groups, and A_{σ} . Finally, we use the **Synthetic** benchmark domain from recent robust RMAB work (Killian et al. 2022). This domain considers three "arm types" [U, V, W]with different intervals, designed so that non-robust policies incur greater regret than robust ones. We augment the domain to allow homogeneous groups of each arm type, where the size and proportion of groups of each type may vary. In our experiments, the default setting has N = 18,000 arms, M = 36 groups, where 1/3 of groups are composed of each of the arm types, and budget K = 100. We run sensitivity analysis on K, the proportion of groups made up of each arm type, and a "block group" setting which joins all arms of a given type into a single group.

Evaluation To evaluate performance, we plan at the group level but simulate individuals within groups independently, where each individual undergoes state transitions based on their own state, action, and transition probabilities. All experiments use horizon H = 10 and report the average of 30 seeds. We measure total reward with discount factor $\gamma = 0.9$. In Fig. 3, we evaluate each approach in terms of regret (Eq. 1), computed by simulating each planner strategy against the full set of adversary pure strategies and selecting one that maximizes regret. Note, there is no actual deployment of the proposed algorithm; all results are simulated.

Baselines First, we compare against the state-of-the-art robust RMAB method, *DDLPO*, for small settings in which DDLPO can complete (Killian et al. 2022). For larger-scale experiments with tens of thousands of arms, no other robust methods are tractable, so we compare against several scalable non-robust baselines. Mate et al.'s non-robust baseline assumes all environment parameters take the median of their uncertainty intervals then computes a reward-maximizing WIP; this strategy was employed in a recent real-world pilot (Mate et al. 2022). We consider two additional non-robust variants which assume that all parameters take the lower bound of the uncertainty interval (*pessimist*) or the upper bound (*optimist*), then compute a WIP strategy. Finally, *random* plans a WIP strategy against an environment that is uniformly randomly sampled from the uncertainty intervals.

6.2 Results

Fig. 3 shows GROUPS outperforms baselines in terms of max regret across several settings. Fig. 3(a-c) shows results for the **maternal health** setting of ARMMAN. In particular, Fig. 3(c) shows that GROUPS scales past 300,000 arms, representing more than a $1000 \times$ increase over the robust state-of-the-art to meet a key need of real-world deployment settings. Moreover, across experiments, the max re-



Figure 3: Max regret (lower is better) incurred by GROUPS, our robust solution approach, compared to non-robust baselines across various settings. (a-c) Maternal health. For (c), the number of arms is increased by multiplying each group size by a constant factor, i.e., 1, 10, and 20, but M is constant. (d-f) TB. For (d), budgets are 5%, 10%, and 15% of N. (g-i) Synthetic. For (h), the x-axis is the fraction of groups of arm type U — the fraction of type V is always 0.33, and the remaining fraction are type W. For (i) the x-axis denotes the arm type that has been combined into a single group of 6000 arms, where the other two types are split across 12 groups each of size 500. In the maternal health and TB settings, regret can be interpreted, in real-world terms, as the maximum *preventable* missed health messages and doses, respectively, across the uncertainty space.

gret of GROUPS is nearly half that of the non-robust strategy used in Mate et al. (2022). In other words, our simulations demonstrate that compared to the best non-robust strategy **GROUPS could prevent mothers from missing thousands of pregnancy-related health messages, each containing potentially life-saving care information.**

On the **TB** domain (Fig. 3(d–f)), we see again that GROUPS performs well across various strategies for grouping and computing uncertainty, even with very imbalanced group sizes. On the **synthetic** domain (Fig. 3(g–i)), across various budgets and grouping strategies, the non-robust baselines vary in performance and are sometimes worse than random, demonstrating the need for reliable robust policies. Moreover, Table 1 shows that GROUPS even outperforms the state-of-the-art DDLPO in terms of regret on the synthetic benchmark dataset for problems sizes small enough for DDLPO to complete (i.e., N < 100). The superior performance of GROUPS is due to our Whittle-based policies which specialize to two-action settings, in contrast to the more general but highly stochastic deep learning–based policies of DDLPO.

Supported by Theorem 3, GROUPS scales significantly without incurring additional regret. In Appendix I, we demonstrate the significant runtime improvement of GROUPS as M decreases, holding N constant. The scalability of our approach is critical for robust RMAB solutions

Table 1: Regret of GROUPS vs. robust method DDLPO on **Synthetic**. We set M = N and K = 1 to match the evaluation in Killian et al. (2022). GROUPS incurs less regret.

	GROUPS	DDLPO
N = 6	0.64 ± 0.05	1.00 ± 0.06
N = 9	0.47 ± 0.06	0.98 ± 0.05
N = 12	0.45 ± 0.06	0.88 ± 0.05

to actually be deployed in real-world, low-resource settings.

7 Conclusion

The GROUPS algorithm we introduce presents several key advances to make RMABs more useful in practice, enabling simultaneous *scaleup* and *robustness to uncertainty*. We are working with ARMMAN to deploy GROUPS to positively impact maternal health, demonstrating the real-world capabilities this work enables. Most notably, **our simulation experiments demonstrate that our robust planning method could help ARMMAN prevent mothers from missing thousands of health messages**, a promising result that we hope to translate into practice to help deliver life-saving health information to otherwise under-served mothers.

8 Acknowledgments

J.A.K. was supported by an NSF Graduate Research Fellowship under grant DGE1745303. A.B. was supported by the Harvard Center for Research on Computation and Society. L.X. was supported by a Google PhD Fellowship, and was a Student Researcher at Google for part of the project.

References

Alam, M.; D'Este, C.; Banwell, C.; and Lokuge, K. 2017. The impact of mobile phone based messages on maternal and child healthcare behaviour: a retrospective cross-sectional survey in Bangladesh. *BMC Health Serv. Res.*, 17(1).

Biswas, A.; Aggarwal, G.; Varakantham, P.; and Tambe, M. 2021. Learn to Intervene: An Adaptive Learning Policy for Restless Bandits in Application to Preventive Healthcare. In *IJCAI*.

Braziunas, D.; and Boutilier, C. 2007. Minimax regret based elicitation of generalized additive utilities. In *UAI-07*.

Carvalho, N.; Salehi, A.; and Goldie, S. 2013. National and sub-national analysis of the health benefits and cost-effectiveness of strategies to reduce maternal mortality in Afghanistan. *Health Policy Plan*, 28(1).

Evans, W. D.; Wallace, J. L.; and Snider, J. 2012. Pilot evaluation of the text4baby mobile health program. *BMC public health*, 12(1).

Gates Foundation. 2021. Global Progress and Projections for Maternal Mortality. https://www.gatesfoundation. org/goalkeepers/report/2021-report/progress-indicators/ maternal-mortality/.

Glazebrook, K. D.; Hodge, D. J.; and Kirkbride, C. 2011. General notions of indexability for queueing control and asset management. *Ann Appl Probab*, 21(3): 876–907.

Green, H.; and Riddington, A. 2020. Gender inclusive language in perinatal services: Mission statement and rationale. *Brighton, England: Brighton and Sussex University Hospitals.*

Gribble, K. D.; Bewley, S.; Bartick, M. C.; Mathisen, R.; Walker, S.; Gamble, J.; Bergman, N. J.; Gupta, A.; Hocking, J. J.; and Dahlen, H. G. 2022. Effective communication about pregnancy, birth, lactation, breastfeeding and newborn care: the importance of sexed language. *Frontiers in global women's health*, 3.

Gurobi Optimization, L. 2021. Gurobi Optimizer Reference Manual.

HLPF Review of SDG3. 2017. High-Level Political Forum Thematic Review of Sustainable Goal 3.

Huo, X.; and Fu, F. 2017. Risk-aware multi-armed bandit problem with application to portfolio selection. *Royal Society open science*, 4(11): 171377.

Killian, J. A.; Biswas, A.; Shah, S.; and Tambe, M. 2021. Q-Learning Lagrange Policies for Multi-action Restless Bandits. In *KDD*.

Killian, J. A.; Wilder, B.; Sharma, A.; Shah, D.; Choudhary, V.; Dilkina, B.; and Tambe, M. 2019. Learning to Prescribe

Interventions for Tuberculosis Patients Using Digital Adherence Data. In *SIGKDD International Conference on Knowledge Discovery & Data Mining*.

Killian, J. A.; Xu, L.; Biswas, A.; and Tambe, M. 2022. Robust Restless Bandits: Tackling Interval Uncertainty with Deep Reinforcement Learning. *Conference on Uncertainty in Artificial Intelligence (UAI)*.

Lanctot, M.; Zambaldi, V.; Gruslys, A.; Lazaridou, A.; Tuyls, K.; Pérolat, J.; Silver, D.; and Graepel, T. 2017. A unified game-theoretic approach to multiagent reinforcement learning. *NeurIPS-17*, 30.

Li, S.; Wu, Y.; Cui, X.; et al. 2019. Robust multi-agent reinforcement learning via minimax deep deterministic policy gradient. In *AAAI*.

Maillard, O.-A. 2013. Robust risk-averse stochastic multiarmed bandits. In *ALT*. Springer.

Mate, A.; Madaan, L.; Taneja, A.; et al. 2022. Field Study in Deploying Restless Multi-Armed Bandits: Assisting Non-Profits in Improving Maternal and Child Health. *AAAI*.

McMahan, H. B.; Gordon, G. J.; and Blum, A. 2003. Planning in the presence of cost functions controlled by an adversary. In *ICML-03*.

Meh, C.; Sharma, A.; Ram, U.; Fadel, S.; Correa, N.; Snelgrove, J. W.; Shah, P.; Begum, R.; Shah, M.; Hana, T.; et al. 2021. Trends in maternal mortality in India over two decades in nationally representative surveys. *BJOG*.

MomConnect. 2021. Mobile Alliance for Maternal Action. https://www.jnj.com/our-giving/momconnectconnecting-women-to-care-one-text-at-a-time.

Murthy, N.; Chandrasekharan, S.; Prakash, M. P.; Kaonga, N. N.; Peter, J.; Ganju, A.; and Mechael, P. N. 2019. The impact of an mHealth voice message service (mMitra) on infant care knowledge, and practices among low-income women in India: findings from a Pseudo-Randomized controlled trial. *Maternal and child health journal*, 23(12): 1658–1669.

Musiimenta, A.; Tumuhimbise, W.; Pinkwart, N.; et al. 2021. A mobile phone-based multimedia intervention to support maternal health is acceptable and feasible among illiterate pregnant women in Uganda. *Digital Health*, 7.

Nakhleh, K.; Ganji, S.; Hsieh, P.-C.; Hou, I.; Shakkottai, S.; et al. 2021. NeurWIN: Neural Whittle Index Network For Restless Bandits Via Deep RL. *Advances in Neural Information Processing Systems*, 34.

Papadimitriou, C. H.; and Tsitsiklis, J. N. 1999. The complexity of optimal queuing network control. *Math. Oper. Res.*, 24(2).

Pinto, L.; Davidson, J.; Sukthankar, R.; and Gupta, A. 2017. Robust adversarial reinforcement learning. In *ICML*. PMLR.

Rioux, C.; Weedon, S.; London-Nadeau, K.; Paré, A.; Juster, R. P.; Roos, L. E.; Freeman, M.; and Tomfohr-Madsen, L. 2021. Gender-inclusive language in pregnancy-related research: why and how to improve current practices.

Schomaker, M.; and Heumann, C. 2018. Bootstrap inference when using multiple imputation. *Stat Med*, 37(14).

Sinha, A.; and Mahajan, A. 2022. Sensitivity of Whittle index policy to model approximation.

Tamrat, T.; and Kachnowski, S. 2012. Special delivery: an analysis of mHealth in maternal and newborn health programs and their outcomes around the world. *Matern Child Health J.*, 16(5).

United Nations. 2021. Sustainable Development Goal 3: Ensure healthy lives and promote well-being for all at all ages. https://sdgs.un.org/goals/goal3.

Wang, K.; Xu, L.; Taneja, A.; and Tambe, M. 2022. Optimistic Whittle Index Policy: Online Learning for Restless Bandits. *arXiv preprint arXiv:2205.15372*.

Wang, S.; Huang, L.; and Lui, J. 2020. Restless-UCB, an Efficient and Low-complexity Algorithm for Online Restless Bandits. *Advances in Neural Information Processing Systems*, 33: 11878–11889.

Watterson, J. L.; Walsh, J.; and Madeka, I. 2015. Using mHealth to improve usage of antenatal care, postnatal care, and immunization: a systematic review of the literature. *BioMed Res. Int.*

Weber, R. R.; and Weiss, G. 1990. On an index policy for restless bandits. *J. Appl. Probab.*, 27(3).

Whittle, P. 1988. Restless bandits: Activity allocation in a changing world. *J. Appl. Probab.*, 25(A).

Xu, L.; Perrault, A.; Fang, F.; Chen, H.; and Tambe, M. 2021. Robust Reinforcement Learning Under Minimax Regret for Green Security. In *UAI*.

A A Note on Language

Throughout this paper, we use the term "mother" to refer to pregnant women and people, birthing women and people, or postnatal women and people. We note that "mother" is gendered language that also typically refers to someone who has already given birth, which may not be the case for newly pregnant women and people who are enrolled in maternal health programs but have not yet given birth.

We stand by the need to provide compassionate medical care for trans and non-binary patients, and recognize that the word "mother" may not reflect the identity of all people. Some recent calls advocate the language "birthing person" or "birthing women and people" over "mother" (Rioux et al. 2021; Green and Riddington 2020). Others point out that language such as "birthing person" may be dehumanizing and go against best practices in health communication of making communication easy to understand for patients with low literacy or education or are communicating in their non-native language (Gribble et al. 2022).

For the above reasons and to keep our writing concise, in this paper we use "mother" while recognizing that the term is imperfect. We hope to move to healthcare language that is inclusive of all who are in need of those services. We also highlight the need to reevaluate other related terminology such as "maternal health" to move towards more inclusive language, and thus more inclusive care.

B Ethical Considerations

As our system aims to improve engagement with maternal health information, this deployment must be carefully trialed with guardrails developed before a large-scale deployment. We are collaborating with ARMMAN to carry out these trials and develop such guardrails. We recognize the system's recommendations are influenced by historical data, a potential source of bias, especially when data is scarce. A benefit of our method is to be robust to such scarcity-induced bias, a key step forward toward responsible deployment.

Thus, some of the key ethical considerations are in how GROUPS may be used to optimize resource allocation in real world settings. We discussed how we, specifically, will work toward responsible deployment with our partners, as well as how GROUPS represents a new capability in RMAB planning, by allowing one to encode uncertainty due to, e.g., data scarcity. However, we must also consider the impact of a system like GROUPS on society more broadly. For instance, with such a scalable and inherently black box optimization tool such as GROUPS, there may be a temptation to allow the scheduling and delivery of interventions to become fully automated. This could negatively impact, e.g., the availability of work for humans who may have previously scheduled or delivered the interventions, or negatively impact intervention recipients who perhaps could receive unwanted interventions with little option for humanmediated recourse. To avoid these negative impacts, we stress that our system should be seen as a supplemental tool on the toolbelt of intervention schedulers, to be considered among a range of existing criteria and expertise, rather than a replacement solution.

C Limitations

Our work takes a major step forward in scaling up RMAB solutions to perform robustly under uncertainty. However, it accomplishes this, in part, by taking advantage of the decomposable nature of the Whittle index and Whittle index regret, which both require the binary-action setting typically considered for RMABs. This could be limiting in domains where planners need to optimize over a suite of different *types* of interventions, rather than deciding between only {act, not act} for each arm. We note though, that our methods could be extended to the multi-action setting with, e.g., multi-action notions of the Whittle index (Glazebrook, Hodge, and Kirkbride 2011), and corresponding notions of multi-action index regret.

The scalability of our method also relies on the existence of a reasonable number of groups within data. We validate that our method runs hundreds of times faster than the state of the art when there are ~ 10 arms and groups (Fig. 8). We also show that our method runs in about 15 minutes for the real maternal health dataset which has ~ 15 K arms and 40 groups (Fig. 6). However, even our method may have difficulty scaling if there were, e.g., tens of thousands of groups in the data. Yet, in that case, the planner could create a smaller number of groups, each with more arms per group. Then uncertainty of each group might increase, but the planner could still use GROUPS to plan and achieve good performance. We demonstrate an example of such a tradeoff between number of groups and performance in Fig. 3(e). GROUPS performs better than the baselines across all cases.

D Additional Notation and Preliminaries

The number of mothers in group $m \in [M]$ is equal to γ_m . The total number of mothers is equal to $N = \sum_{m \in [M]} \gamma_m$. Further, $S^m = \times_{n \in [\gamma_m]} S_n^m$ (resp. $\mathcal{A}^m = \times_{n \in [\gamma_m]} \mathcal{A}_n^m$) denotes the set of different state-profiles (resp. action-profiles) of the mothers in group m. An element of S^m (resp. \mathcal{I}^m) is denoted as $s^m = (s_n^m)_{n \in [\gamma_j]}$ (resp. $\mathbf{a}^m = (a_n^m)_{n \in [\gamma_m]}$). Finally, $S = \times_{m \in [M]} S^m$ (resp. $\mathcal{A} = \times_{m \in [M]} \mathcal{A}^m$) denotes the different state-profile (action-profile) of all the N mothers, and s denotes an element of S.

A policy π is a map from S to A, such that for each $s \in S$, $\pi(s)$ has at most K actions as 1 (intervention) and the remaining are 0 (no-intervention). In particular, π maps a strategy profile of the N mothers to an action profile with at most k intervention actions. We use Π to denote the set of all such policies. We note that Π is equal to the set of pure strategies of the planner. Let $\alpha \in \Delta_I(\Pi)$ (respectively $\beta \in \Delta_I(\mathscr{P})$) be a mixed strategy. Then we use $P_{\alpha}(\pi)$ (resp. $P_{\beta}(P)$) to denote the probability of choosing $\pi \in \Pi_I$ (resp. $P \in \mathscr{P}_I$) under α (resp. β).

Finally, we often use $(\alpha_{X,Y}^*, \beta_{X,Y}^*)$ to denote the minimax maximin regret strategies of the planner and adversary respectively, when the planner plays mixed strategies from $\Delta_X(\Pi)$ and the adversary plays mixed strategies from $\Delta_Y(\mathscr{P})$ (see Section 5 for the definitions of $\Delta_X(\Pi)$ and $\Delta_Y(\mathscr{P})$), and $X, Y \in \{M, I\}$. Note that $\alpha_{X,Y}^*$ is the plan-

 Table 2: Notation table

Notation	Description
[N]	Set of N natural numbers $\{1, \ldots, N\}$
$[\underline{a}, \overline{a}]$	A real interval denoting all values of a such that $\underline{a} \leq a \leq \overline{a}$
N	Number of women (arms)
K	Number of interventions (service calls) that can be made each round
M	Number of groups
\mathcal{S}_n	Set of states (engagement status) of arm $n \in [N]$; $S_n = \{0, 1\}$ (non-engaging state and engaging state)
\mathcal{S}^N	Combinatorial set of states over all arms
s	Vector of states of all arms
\mathcal{A}_n	Set of actions (intervention decisions) on arm $n \in [N]$; $A_n = \{0, 1\}$ (no-intervention and intervention)
\mathcal{A}^N	Combinatorial set of actions over all arms
$P^n_{s,a,s'}$	Probability of transitioning from state s to s' on action a for an individual n
$\overline{\underline{P}}^{m}$	Uncertainty intervals for all transition probabilities for group m
π_i	Planner's <i>i</i> -th pure strategy (Whittle index policy)
P_i	Adversary's <i>i</i> -th pure strategy (instantiation of P^m for all m)
α	Planner's mixed strategy
β	Adversary's mixed strategy
$G(\pi, P)$	Expected reward when planner plays a pure strategy π and adversary plays a pure strategy P
$R(\pi, P)$	Expected regret when planner plays a pure strategy π and adversary plays a pure strategy P

ner's minimax strategy and $\beta^{\star}_{X,Y}$ is the adversary's maximin strategy.

D.1 Permutations on Π and \mathscr{P}

Let \mathfrak{G}_m be the set of all permutations of $[\gamma_m]$, where a permutation is a bijective map from $[\gamma_m]$ to $[\gamma_m]$. We use σ_m to denote the elements of \mathfrak{G}_m . Further, let $\mathfrak{G} = \mathfrak{G}_1 \times \cdots \times \mathfrak{G}_M$. We use $\sigma = (\sigma_1, \ldots, \sigma_M)$ to denote an element of \mathfrak{G} ; note that here $\sigma_m \in \mathfrak{G}_m$.

Let $\mathbf{a} = (a_n^m)_{n \in [\gamma_m], m \in [M]}$ be an action profile. Then for a $\sigma = (\sigma_1, \ldots, \sigma_M) \in \mathfrak{G}$, define $\sigma(\mathbf{a}) = (a_{\sigma_m(n)}^m)_{n \in [\gamma_m], m \in [M]}$. In particular σ permutes the actions corresponding to the mothers in group m according to σ_m . Now we show how a $\sigma \in \mathfrak{G}$ defines bijective maps on \mathscr{P} and Π . For every $\sigma \in \mathfrak{G}$ define the map $\phi_{\sigma} : \Pi \to \Pi$ (respectively $\psi_{\sigma} : \mathscr{P} \to \mathscr{P}$) as follows: denote $\phi_{\sigma}(\pi)$ as π' (resp. $\psi_{\sigma}(P)$ as P'), then $\pi'(\mathbf{s}) = \sigma(\pi(\mathbf{s}))$ (resp. $(P')_n^m = P_{\sigma_m(n)}^m)$). We make the following observation which will be helpful later on.

Observation 1. For all $\pi \in \Pi$, $P \in \mathscr{P}$, and $\sigma \in \mathfrak{G}$ the following holds $R(\pi, P) = R(\phi_{\sigma}(\pi), \psi_{\sigma}(P))$.

The following observation follows from the fact that for every $\sigma \in \mathfrak{G}$, ϕ_{σ} and ψ_{σ} define bijective maps on Π and Prespectively.

Observation 2. Let $\alpha \in \Delta_I(\Pi)$ and $\beta \in \Delta_I(\mathscr{P})$ be mixed strategies of the planner and adversary respectively. Then for any $\sigma \in \mathfrak{G}$ the following holds:

$$R(\alpha,\beta) = \sum_{\pi,\mathbf{p}} P_{\alpha}(\phi_{\sigma}(\pi)) \cdot P_{\beta}(\psi_{\sigma}(P)) \cdot R(\phi_{\sigma}(\pi),\psi_{\sigma}(P)) \,.$$

Next, we define permutations which are transpositions. We say $\sigma_m \in \mathfrak{G}_m$ is a transposition if there exists $n, n' \in$ $[\gamma_m]$ such that $\sigma_m(n) = n'$ and $\sigma_m(n') = n$, and for all $\ell \notin \{n, n'\}, \sigma_m(\ell) = \ell$. We say $\sigma = (\sigma_1, \ldots, \sigma_M) \in \mathfrak{G}$ is a transposition if for all $m \in [M], \sigma_m$ is a transposition. We note that every $\sigma \in \mathfrak{G}$ can be expressed as a composition of finitely many transpositions. We note the following observation.

Observation 3. Let $\sigma \in \mathfrak{G}$ be a transposition. Then for every $\pi \in \mathcal{A}$ and $P \in \mathscr{P}$, $\phi_{\sigma}(\phi_{\sigma}(\pi)) = \pi$ and $\psi_{\sigma}(\psi_{\sigma}(\mathscr{P})) = P$.

D.2 Mixed Strategies that Do Not Distinguish Between Mothers in the Same Group

We say a mixed strategy α of the planner is indifferent towards mothers from the same group if for all $\pi, \pi' \in \Pi$ such that there is a $\sigma \in \mathfrak{G}$ satisfying $\pi' = \phi_{\sigma}(\pi), P_{\alpha}(\pi) = P_{\alpha}(\pi')$. We use $\Delta_M(\Pi)$ to denote such mixed strategies.

If we define the probability of intervening on mother nunder a mixed strategy α as follows

$$P_{\alpha}(\text{intervene }n) = \sum_{\pi \in \Pi} P_{\alpha}(\pi) \sum_{\mathbf{s} \in \mathcal{S}} \mathbbm{1}\{a_{\mathbf{s},n}^{\pi} = 1\}$$

then it is easy to see that for an $\alpha \in \Delta_M(\Pi)$ and mothers n, n' in the same group $P_{\alpha}(\text{intervene } n) = P_{\alpha}(\text{intervene } n')$.

E Proofs of Theorem 1, 2, and 3

We refer the reader to Section D for additional notations and missing definitions used in the proof. Additionally, we now define order-preserving maps, that we will use in the proof of Theorem 3.

E.1 Order-Preserving Maps

To prove Theorem 3 from Section 3, we require the assumption that there is a map $\psi : \Delta_I(\mathscr{P}) \to \Delta_M(\mathscr{P})$ such that

- 1. For every $\alpha_1, \alpha_2 \in \Delta_M(\Pi)$ and $\beta \in \Delta_I(\mathscr{P})$, $R(\alpha_1, \beta) > R(\alpha_2, \beta)$ iff $R(\alpha_1, \psi(\beta)) > R(\alpha_2, \psi(\beta))$, and
- 2. For every $\alpha \in \Delta_M(\Pi)$ and $\beta_1, \beta_2 \in \Delta_I(\mathscr{P}) R(\alpha, \beta_1) > R(\alpha, \beta_2)$ iff $R(\alpha, \psi(\beta_1)) > R(\alpha_2, \psi(\beta_1))$

While (1) and (2) may not hold for general mixed strategies in $\Delta_I(\Pi)$, it is likely to be true for mixed strategies in $\Delta_M(\Pi)$, since these strategies do not distinguish between mothers from the same group. Next, we describe a possible candidate map.

First we define a map $\phi : \mathscr{P} \to \mathscr{P}$. Let $P \in \mathscr{P}$ be a pure strategy of the adversary which assigns different transition probabilities to the mothers in the same group, and for pure strategy \mathbf{p} let $P_{s,a,s'}^{m,n}$ be the probability of the mother n in group m transitioning from state s to s' under action a. We define $\phi(P) = \hat{P}$ as follows, \hat{P} assigns the same transition probability to all the mothers in a group by averaging the transition probabilities of the mothers in that group. In particular, $\hat{P}_{s,a,s'}^{m,n} = \hat{P}_{s,a,s'}^m = \sum_{m \in [\gamma_m]} p_{s,a,s'}^{m,n}$. Notice that the pure strategy $\phi(P)$ of the adversary is indifferent to mothers in the same group. Further, for a $P \in$ \mathscr{P} , let $\phi^{-1}(P) = \{P' \in P \mid \phi(P') = \mathbf{p}\}$. Now let $\beta \in \Delta_I(\mathscr{P})$. Then $\psi(\beta) = \beta'$ is defined as follows: $P_{\beta'}(P) = \sum_{P' \in \phi^{-1}(p)} P_{\beta}(P')$. Since, $\phi^{-1}(P) \neq \emptyset$ only if \mathbf{p} assigns the same transition probability to all the mothers in a group. Hence, we have $\psi(\beta) \in \Delta_M(\mathscr{P})$.

E.2 Proof of Theorem 1

We require the following proposition to prove the theorem.

Proposition 1. Suppose the minimax maximin regret strategies $(\alpha_{I,I}^*, \beta_{I,I}^*)$ is such that there exists a permutation $\sigma \in \mathfrak{G}$ satisfying $P_{\beta}(P) = P_{\beta}(\psi_{\sigma}(P))$ for every $P \in \mathscr{P}_{I}$. Then there exists a planner mixed strategy $\alpha' \in \Delta_{I}(\Pi)$ such that $P_{\alpha}(\phi_{\sigma}(\pi)) = P_{\alpha}(\pi)$ for every $\pi \in \Pi$, and $(\alpha', \beta_{I,I}^*)$ are minimax maximin regret strategies of the planner and adversary respectively at the individual level.

Proof. Suppose there exists $\pi_m, \pi_\ell \in \Pi_I$ such that $\phi_{\sigma}(\pi_m) = \pi_\ell$ but $P_{\alpha}(\pi_\ell) < P_{\alpha}(\pi_m)$. First, we show in this case that $R(\pi_m, \beta_{I,I}^*) = R(\pi_\ell, \beta_{I,I}^*)$.

$$R(\pi_{\ell}, \beta_{I,I}^{\star})_{(i)}^{\equiv} \sum_{\mathbf{p} \in P} P_{\beta_{I,I}^{\star}}(P) \cdot R(\pi_{\ell}, P)$$

$$\stackrel{\equiv}{\underset{(ii)}{\equiv}} \sum_{\mathbf{p} \in P} P_{\beta_{I,I}^{\star}}(\psi_{\sigma}(P)) \cdot R(\pi_{\ell}, \psi_{\sigma}(P))$$

$$\stackrel{\equiv}{\underset{(iii)}{\equiv}} \sum_{\mathbf{p} \in P} P_{\beta_{I,I}^{\star}}(P) \cdot R(\pi_{m}, P)$$

$$= R(\pi_{m}, \beta_{I,I}^{\star})$$

In the above equations: (i) follows from the definition of regret of a pure strategy π_{ℓ} on the adversary's mixed strategy $\beta_{I,I}^*$, (ii) follows as σ is a permutation (also see Observation 2) and from Observation 1 we have $R(\pi_m, P) =$ $R(\pi_{\ell}, \psi_{\sigma}(P))$ for all $P \in \mathscr{P}$, and (iii) follows by using $P_{\beta_{I,I}^*}(P) = P_{\beta_{I,I}^*}(\psi_{\sigma}(P))$ for all $P \in \mathscr{P}$. Now construct a mixed strategy α' such that $P_{\alpha'}(\pi_m) = P_{\alpha'}(\pi_{\ell}) =$ $\frac{P_{\alpha}(\pi_m)+P_{\alpha}(\pi_{\ell})}{2}, \text{ and for all } \pi \in \Pi_I \text{ such that } \pi \neq \pi_m \text{ and } \pi \neq \pi_{\ell} \text{ we have } P_{\alpha'}(\pi) = P_{\alpha}(\pi). \text{ Since } R(\pi_m, \beta_{I,I}^*) = R(\pi_{\ell}, \beta_{I,I}^*), \text{ it is easy to see that } (\alpha', \beta) \text{ is a minimax maximin regret strategies at the individual level, and from construction } P_{\alpha'}(\pi_m) = P_{\alpha'}(\pi_{\ell}). \square$

Proof of Theorem 1. If $\alpha_{I,I}^* \in \Delta_M(\Pi)$ then the Theorem follows. Hence, assume there are $\pi_m, \pi_\ell \in \Pi$ and a permutation $\sigma \in \mathfrak{G}$ such that $\phi_\sigma(\pi_m) = \pi_\ell$ but $P_\alpha(\pi_m) > P_\alpha(\pi_\ell)$. Here, we use the subscript m and ℓ to denote more and less. Observe that we may assume without loss of generality that σ is a transposition (see App. D). This is because any permutation in $\sigma \in \mathfrak{G}$ can be expressed as a composition of transpositions. Hence, assuming σ is a transposition, we have $\phi_\sigma(\pi_\ell) = \pi_m$.

Let $(\alpha_{I,I}^{\star}, \beta_{I,I}^{\star})$ be a minimax maximin regret strategies, that is, $\beta_{I,I}^{\star}$ is a regret maximizing mixed strategy of adversary against $\alpha_{I,I}^{\star}$, that is

$$\beta_{I,I}^{\star} \in \underset{\beta \in \Delta_{I}(\mathscr{P})}{\operatorname{argmax}} R(\alpha_{I,I}^{\star}, \beta) .$$

Construct β' such that $P_{\beta'}(P) = P_{\beta_{I,I}^*}(\psi_{\sigma}(P))$ for all $P \in \mathscr{P}$. Note that since σ is a transposition, we also have $P_{\beta'}(\psi_{\sigma}(P)) = P_{\beta_{I,I}^*}(P)$ for all $P \in \mathscr{P}$. Hence, as $\beta_{I,I}^*$ is regret maximizing for the adversary against the planner's mixed strategy $\alpha_{I,I}^*$, we have $R(\alpha_{I,I}^*, \beta') \leq R(\alpha_{I,I}^*, \beta_{I,I}^*)$.

First, we argue that $R(\alpha_{I,I}^{\star}, \beta') = R(\alpha_{I,I}^{\star}, \beta_{I,I}^{\star})$. Suppose $R(\alpha_{I,I}^{\star}, \beta') < R(\alpha_{I,I}^{\star}, \beta_{I,I}^{\star})$. Then writing the regret expressions for $R(\alpha_{I,I}^{\star}, \beta')$ and $R(\alpha_{I,I}^{\star}, \beta_{I,I}^{\star})$ we have

$$\sum_{\mathbf{p}} \sum_{\pi} P_{\beta'}(\psi_{\sigma}(P)) \cdot P_{\alpha_{I,I}^{\star}}(\pi) \cdot R(\pi, P)$$

$$< \sum_{\mathbf{p}} \sum_{\pi} P_{\beta_{I,I}^{\star}}(P) \cdot P_{\alpha_{I,I}^{\star}}(\pi) \cdot R(\pi, P)$$
(6)

Substituting $R(\pi, P) = R(\phi_{\sigma}(\pi), \psi_{\sigma}(P))$ for all π, P (see Observation 1) we have

$$\sum_{\mathbf{p}} \sum_{\pi} P_{\beta_{I,I}^{\star}}(\psi_{\sigma}(P)) \cdot P_{\alpha_{I,I}^{\star}}(\pi) R(\phi_{\sigma}(\pi), \psi_{\sigma}(P)) \quad (7)$$
$$< \sum_{\mathbf{p}} \sum_{\pi} P_{\beta_{I,I}^{\star}}(P) \cdot P_{\alpha_{I,I}^{\star}}(\pi) R(\pi, P)$$

Let α' be the mixed strategy of planner such that $P_{\alpha'}(\phi_{\sigma}(\pi)) = P_{\alpha_{I,I}^{\star}}(\pi)$. Substituting this in the above equation we have

$$\sum_{\mathbf{p}} \sum_{\pi} P_{\beta_{I,I}^{\star}}(\psi_{\sigma}(P)) \cdot P_{\alpha'}(\phi_{\sigma}(\pi)R(\phi_{\sigma}(\pi),\psi_{\sigma}(P)) \quad (8)$$
$$< \sum_{\mathbf{p}} \sum_{\pi} P_{\beta_{I,I}^{\star}}(P) \cdot P_{\alpha_{I,I}^{\star}}(\pi)R(\pi,P)$$

From Observation 2 we have,

$$\sum_{\mathbf{p}} \sum_{\pi} P_{\beta_{I,I}^{\star}}(\psi_{\sigma}(P)) \cdot P_{\alpha'}(\phi_{\sigma}(\pi)) \cdot R(\phi_{\sigma}(\pi), \psi_{\sigma}(P))$$
(9)

$$= R(\alpha', \beta_{I,I}^{\star})$$

Hence $R(\alpha', \beta_{I,I}^{\star}) < R(\alpha_{I,I}^{\star}, \beta_{I,I}^{\star})$. This contradicts the minimax theorem which states that $\alpha_{I,I}^{\star}$ is the regret minimizing mixed strategy of the planner against the adversary's mixed strategy $\beta_{I,I}^{\star}$. Hence, we have $R(\alpha_{I,I}^{\star}, \beta') = R(\alpha_{I,I}^{\star}, \beta_{I,I}^{\star})$. We note that the above equations also show that α' is the regret minimizing mixed strategy for the planner in response to adversary's mixed strategy β' .

Now construct $\tilde{\beta}$ such that

$$P_{\tilde{\beta}}(P) = \frac{P_{\beta_{I,I}^{\star}}(P) + P_{\beta'}(P)}{2} = \frac{P_{\beta_{I,I}^{\star}}(P) + P_{\beta_{I,I}^{\star}}(\psi_{\sigma}(P))}{2}$$

We now argue that $(\alpha_{I,I}^*, \hat{\beta})$ is a minimax maximin regret strategies at the individual level, that is, $\alpha_{I,I}^*$ is regret minimizing against $\tilde{\beta}$, and $\tilde{\beta}$ is regret maximizing against $\alpha_{I,I}^*$. Since σ is a transposition, this implies $P_{\tilde{\beta}}(P) = P_{\tilde{\beta}}(\psi_{\sigma}(P))$ for all $\mathbf{p} \in P$. Further, as $R(\alpha_{I,I}^*, \beta') = R(\alpha_{I,I}^*, \beta_{I,I}^*)$, we have $R(\alpha_{I,I}^*, \beta_{I,I}^*) = R(\alpha_{I,I}^*, \tilde{\beta})$, and hence, $\tilde{\beta}$ is a regret maximizing mixed strategy of the adversary against $\alpha_{I,I}^*$, that is,

$$\tilde{\beta} \in \operatorname*{argmax}_{\beta \in \Delta_{I}(\mathscr{P})} R(\alpha_{I,I}^{\star}, \beta)$$

Also, a similar argument, as from Equations 6 to 9, shows that $R(\alpha_{I,I}^*, \beta') = R(\alpha', \beta')$, and hence $\alpha_{I,I}^*$ is a regret minimizing mixed strategy for the planner in response to adversary's mixed strategy β' . This follows from the minimax theorem and that α' is the regret minimizing mixed strategy for the planner in response to adversary's mixed strategy β' . This together implies $\alpha_{I,I}^*$ is the regret minimizing mixed strategy for the planner in response to adversary's mixed strategy $\tilde{\beta}$. In particular, $(\alpha',)$

Now we use Proposition 1, which shows that there exists a $\tilde{\alpha}$ such that $P_{\tilde{\alpha}}(\pi_m) = P_{\tilde{\alpha}}(\pi_\ell)$,

$$\tilde{\beta} \in \operatorname*{argmax}_{\beta \in \Delta_{I}(\mathscr{P})} R(\tilde{\alpha}, \beta)$$

and $R(\tilde{\alpha}, \tilde{\beta}) = R(\alpha_{I,I}^{\star}, \tilde{\beta}) = R(\alpha_{I,I}^{\star}, \beta_{I,I}^{\star})$. We can repeat this process finitely many times to show to construct a mixed strategy α such that for any two policies $\pi, \pi' \in \mathcal{A}$, if there exists a σ such that $\phi_{\sigma}(\pi) = \pi'$ then $P_{\alpha}(\pi) = P_{\alpha}(\pi')$, and $\max_{\beta \in \Delta_I(\mathscr{P})}(R(\alpha, \beta)) = R(\alpha_{I,I}^{\star}, \beta_{I,I}^{\star})$. Since $\alpha \in \Delta_M(\Pi)$, we have, without loss of generality $\alpha = \alpha_{M,I}^{\star}$. \Box

E.3 Proof of Theorem 2

Let $(\alpha_{M,M}^{\star}, \beta_{M,M}^{\star})$ be a minmax-maximin regret strategies at the group level. Further, let $\pi, \pi' \in \Pi$ be such that $P_{\alpha_{M,M}^{\star}}(\pi) > P_{\alpha_{M,M}^{\star}}(\pi')$, and there exists a $\sigma \in \mathfrak{G}$ such that $\pi' = \phi_{\sigma}(\pi)$. Let α be an planner mixed strategy such that $P_{\alpha}(\pi) = P_{\alpha}(\pi') = \frac{P_{\alpha_{M,M}^{\star}}(\pi) + P_{\alpha_{M,M}^{\star}}(\pi')}{2}$. First observe that since $\beta \in \Delta_M(\mathscr{P})$, $R(\pi, \beta) = R(\pi', \beta)$. It follows from this that (α, β) is also a minmax regret solution at the group level. We can repeat this process finitely many times to show that for any two policies $\pi, \pi' \in \Pi$, if there exists a σ such that $\phi_{\sigma}(\pi) = \pi'$ then $P_{\alpha}(\pi) = P_{\alpha}(\pi')$.

E.4 Proof of Theorem 3

Let $(\alpha_{M,M}^{\star}, \beta_{M,M}^{\star})$ be a minimax-maximin regret strategies at the group level. Also, let β^{\star} be the regret maximizing strategy of the adversary at the individual level against the planner's mixed strategy $\alpha_{M,M}^{\star}$, that is

$$\beta^{\star} = \operatorname*{argmax}_{\beta \in \Delta_{I}(\mathscr{P})} R(\alpha^{\star}_{M,M},\beta)$$

Further, let $(\alpha_{I,I}^{\star}, \beta_{I,I}^{\star})$ be a minimax-maximin regret strategies at the individual level. In particular, from the minimax theorem

$$\beta_{I,I}^{\star} = \max_{\beta \in \Delta_{I}(\mathscr{P})} R(\alpha_{I,I}^{\star}, \beta)$$

Hence, we wish to show $R(\alpha_{M,M}^{\star}, \beta^{\star}) = R(\alpha_{I,I}^{\star}, \beta_{I,I}^{\star})$. Now let $(\alpha_{M,I}^{\star}, \beta_{M,I}^{\star})$ be a minimax-maximin regret strategies, when the planner plays at group level (from $\Delta_M(\Pi)$) and the adversary plays at the individual level (from $\Delta_I(\mathscr{P})$). Hence, again from the minimax theorem

$$\beta_{M,I}^{\star} = \max_{\beta \in \Delta_{I}(\mathscr{P})} R(\alpha_{M,I}^{\star}, \beta)$$

Recall from Theorem 1, we have

$$R(\alpha_{I,I}^{\star},\beta_{I,I}^{\star}) = R(\alpha_{M,I}^{\star},\beta_{M,I}^{\star})$$

Hence, to prove the theorem it is sufficient to show that $R(\alpha_{M,M}^{\star}, \beta^{\star}) = R(\alpha_{M,I}^{\star}, \beta_{M,I}^{\star})$. Since $(\alpha_{M,I}^{\star}, \beta_{M,I}^{\star})$) is a minimax-maximin strategy regret strategies, when the planner plays at group level (from $\Delta_M(\Pi)$) and the adversary plays at the individual level, we have

$$R(\alpha_{M,M}^{\star},\beta^{\star}) \ge R(\alpha_{M,I}^{\star},\beta_{M,I}^{\star})$$

Suppose for contradiction

$$R(\alpha_{M,M}^{\star},\beta^{\star}) > R(\alpha_{M,I}^{\star},\beta_{M,I}^{\star}) \tag{10}$$

Now we have the following two equations:

$$R(\alpha_{M,M}^{\star},\beta^{\star}) \ge R(\alpha_{M,M}^{\star},\beta_{M,I}^{\star}) \tag{11}$$

$$R(\alpha_{M,M}^{\star}, \beta_{M,I}^{\star}) \ge R(\alpha_{M,I}^{\star}, \beta_{M,I}^{\star})$$
(12)

Equation 11 follows from β^* being the regret maximizing strategy of the adversary at the individual level against the planner's mixed strategy $\alpha^*_{M,M}$, and Equation 12 follows from the minimax theorem and $\alpha^*_{M,I}$, $\beta^*_{M,I}$ being the minimax maximin regret strategies when the planner plays at the group level and the adversary plays at the individual level. Now corresponding to Equations 11 and 12 assuming there is an order preserving map (see App. E.1) $\psi : \Delta_I(\mathscr{P}) \to \Delta_M(\mathscr{P})$, we have

$$R(\alpha_{M,M}^{\star},\psi(\beta^{\star})) \ge R(\alpha_{M,M}^{\star},\psi(\beta_{M,I}^{\star}))$$
(13)

$$R(\alpha_{M,M}^{\star}, \psi(\beta_{M,I}^{\star})) \ge R(\alpha_{M,I}^{\star}, \psi(\beta_{M,I}^{\star}))$$
(14)

Equation 13 follows from property 1 of the order preserving map, and Equation 14 follows from property 2. From Equations 10, 13 and 14, we have

$$R(\alpha_{M,M}^{\star},\psi(\beta^{\star})) > R(\alpha_{M,I}^{\star},\psi(\beta_{M,I}^{\star}))$$
(15)

Finally we use property 2 of the order preserving map to claim the following two equations,

$$R(\alpha_{M,M}^{\star},\psi(\beta^{\star})) = R(\alpha_{M,M}^{\star},\beta_{M,M}^{\star})$$
(16)

$$R(\alpha_{M,I}^{\star},\psi(\beta_{M,I}^{\star})) = \max_{\beta \in \Delta_{M}(\mathscr{P})} R(\alpha_{M,I}^{\star},\beta)$$
(17)

Both equations require property 2 of the order-preserving map. Additionally, Equation 16, follows because β^* (resp. $\beta_{M,M}^{\star}$) is regret maximizing for adversary at the individual level (resp. at the group level) against planner's strategy $\alpha^{\star}_{M,M}$, and Equation 17 follows because $\beta^{\star}_{M,I}$ is regret maximizing for adversary at the individual level against planner's strategy α_{MI}^{\star} . Hence, from Equations 15, 16 and 17, we have

$$R(\alpha_{M,M}^{\star},\beta_{M,M}^{\star}) > \max_{\beta \in \Delta_{M}(\mathscr{P})} R(\alpha_{M,I}^{\star},\beta)$$

The above equation contradicts the worst-case minimality of $\alpha^{\star}_{M,M}$, when both players play at the group level. Hence, we have $R(\alpha_{M,M}^{\star}, \beta^{\star}) \geq R(\alpha_{M,I}^{\star}, \beta_{M,I}^{\star}).$

Minimizing/Maximizing Whittle Indices F

The binary quadratic program for simultaneously maximizing and/or minimizing the Whittle indices over one or more states of a group, given a set of interval parameter ranges on the transitions probabilities $[\underline{P}_{s,a,s'}, \overline{P}_{s,a,s'}]$, named MIN-MAXWHITTLEBQP is given as follows:

$$\min_{W_{s'}} \sum_{s' \in I(S)} \theta_{s'} W_{s'} \quad \text{Primary objective} \\
\min_{V^{s'}} \sum_{s' \in I(S)} V^{s'}(s', W_{s'}) \quad \text{Secondary objective} \\
\text{s.t.} \\
Q^{s'}(s, a, W_{s'}) = R(s) - W_{s'}C(a) + \gamma T(s, a, \cdot)^{\intercal} V^{s'}(\cdot, W_{s'}) \\
V^{s'}(s, W_{s'}) \ge Q^{s'}(s, a, W_{s'}) \\
V^{s'}(s, W_{s'}) \le Q^{s'}(s, a, W_{s'}) + b(s, a, s') M \\
\forall s \in S, a \in \mathcal{A}, s' \in I(S) \\
\sum_{a \in \mathcal{A}} b(s, a, s') = |\mathcal{A}| - 1 \\
\forall s \in S, s' \in I(S) \\
W_{s'} = \gamma \left[T(s', 1, \cdot)^{\intercal} V^{s'}(\cdot, W_{s'}) - T(s', 0, \cdot)^{\intercal} V^{s'}(\cdot, W_{s'}) \right] \\
\forall s' \in I(S) \\$$

$$T(s, a, s'') \in [\underline{P}_{s, a, s''}, \overline{P}_{s, a, s''}]$$

$$\forall s \in \mathcal{S}, a \in \mathcal{A}, s'' \in \mathcal{S}$$
(18)

Where I(S) is the set of all states for which the users wants to jointly optimize Whittle indices, $\theta_{s'} \in \{-1, 0, 1\}$ is the "sense" for the corresponding index to optimize, i.e., 1 to minimize, -1 to maximize, or 0 to not optimize the index for that state (note that θ corresponds to *obj* in Alg. 2), $W_{s'}$ is the Whittle index for state s', V and Q are the state and state-action value functions, respectively, C(a) = a is the

Algorithm 3: Double Oracle

Input: Grouped RMAB simulator and parameter uncertainty intervals $\overline{\underline{P}}^m$ for all groups.

Parameters: Number of iterations T**Output:** Agent mixed strategy α

- 1: $P_0 = \{P_0\}$, with P_0 selected at random
- 2: $\Pi_0 = \{\pi_{B_1}, \pi_{B_2}, \ldots\}$, where π_{B_i} are baseline and heuristic strategies
- 3: for epoch e = 1, 2, ..., T do
- Solve for (α_e, β_e) , mixed Nash equilibrium of regret 4: game with strategy sets P_{e-1} and Π_{e-1}
- 5:
- $\pi_e = \text{WI4MS}(\beta_e)$ $P_e = \text{Regret MAXWI}(\alpha_e)$ $P_e = P_e + \{P_e\} = P_e$ 6:

7:
$$P_e = P_{e-1} \cup \{P_e\}, \Pi_e = \Pi_{e-1} \cup \{\pi_e\}$$

8: return α_e

cost of an action, b is a binary variable that serves to enforce one of the Q constraints on V to be tight (ensuring the value function is solved, and thus Whittle index is valid), M is a large number, e.g., 10^4 , T are variables that hold the transition probabilities, R(s) = s is the reward, and γ is the discount factor. Note that all $T^{\intercal}V$ terms are quadratic, since both T and V are variables in this optimization.

G **Double Oracle and Whittle Index** Algorithms

The outer loop of the double oracle algorithm is given in Alg. 3. We also give COMPUTEWI in Alg. 4, our binarysearch based method for computing the Whittle index, given transition probabilities for a group P^m and a state s. Note also that COMPUTEWI could be implemented by using MAXWHITTLEBQP, and a small wrapper function to adjust the input appropriately. That is, MAXWHITTLEBQP expects intervals over transition parameters $\overline{\underline{P}}^m$, but computing the Whittle index only requires some choice of P^m in the intervals. Thus the wrapper needs to encode P^m as intervals, which can be accomplished by copying each transition probability of P^m to an interval with the same upper and lower bound, for each s, a, s'. Then, to get the Whitthe index for a certain state of the group m, the wrapper should pass in a sense that negative-one-hot encodes the desired state. E.g., if one wants to compute the index for state s = 1 for an arm with two states, the wrapper should pass in a sense $\theta = [0, -1]$. With the above described inputs, specifically since the intervals will have the same upper and lower bound, the quadratic terms in the MAXWHITTLE-BQP will become constant, effectively turning the binary quadratic program into a binary linear program that does not search over transition probabilities for the best Whittle index, but simply returns the Whittle index of the given transition probabilities. This represents a new way to compute Whittle indices that could also be of general interest.

Evaluating Each Oracle Η

Planner oracle We verify empirically that the planner oracle, described in Section 4.1, produces high quality best re-

Algorithm 4: Compute Whittle Index (ComputeWI)

Input: Group m, state s_I , transition probabilities P^m for group m, tolerance ϵ .

Output: Whittle index $W^m(s_I)$

- 1: $ub, lb = INITBSBOUNDS(P^m)$ // Return upper and lower bounds on $W^m(s_I)$ given P^m , e.g., 1, 0. // Now binary search for the Whittle index
- 2: while $ub lb > \epsilon$ do
- 3: $\lambda = \frac{ub+lb}{2}$
- 4: $a = VALUEITERATION(P^m, s_I, \lambda)$ // Run value iteration for the MDP defined by P^m with λ -adjusted reward function $r(s, a, \lambda) = s a\lambda$, and return corresponding $\pi^*(s_I)$
- 5: **if** a=0 **then**
- 6: $ub = \lambda$ // Charging too much, decrease
- 7: else if a=1 then
- 8: $lb = \lambda$ // Can charge more, increase
- 9: $W^m(s_I) = \frac{ub+lb}{2}$
- 10: return $W^m(s_I)$

sponses, i.e., reward-maximizing RMAB intervention policies π , across various problem sizes and intervals. As baselines, we compare against: No action (NA) which simulates the policy that takes action a = 0 on all arms at all time steps, representing a lower bound on reward; Ran**dom** which takes action a = 1 on K randomly chosen arms each round; and Brute Force which enumerates the entire feasible RMAB policy space, simulates the average reward of each policy, then returns the reward of the best-performing policy. Brute force can only be computed for small problem sizes, since its computation cost is exponential in N and K. We evaluate on test data generated by, for each seed, randomly sampling transition intervals $[\underline{P}^m_{s,a,s'}, \overline{P}^m_{s,a,s'}] \; \forall c, s, a, s',$ and randomly sampling a mixed nature strategy. Results, shown in Fig. 4, are reported as the average reward over ten random seeds. Our approach, WI4MS performs nearly as well as brute force for the small problem size, and outperforms all baselines as the problem size increases to the scale of the maternal health intervention problem.

Adversary oracle We verify empirically that the adversary oracle, described in Section 4.2, produces high quality best responses, i.e., regret-maximizing environments P, across various problem sizes and intervals. As baselines, we compare against: Random which selects an P by uniformly randomly sampling $P^m \in \overline{\underline{P}}^m \ \forall m \in [M]$ and **Brute** Force which (1) discretizes the adversary's pure strategy space into D = 3 uniformly spaced values for each interval $[\underline{P}_{s,a,s'}^m, \overline{P}_{s,a,s'}^m]$, (2) enumerates all possible combinations of the discrete environment setting, denoted P_d , (3) simulates the average regret induced by each P_d by simulating the optimal WIP against P_d and simulating some input planner mixed strategy α against P_d , then taking the difference, (4) then returns the regret of the best-performing P_d . Brute force can only be computed for small problem sizes, since its computation cost is exponential in N, K,

and *D*. For instance, even for N = 2, K = 1, and D = 3brute force enumerates ~60k choices of P_d . We evaluate on test data generated by randomly sample transition intervals $[\underline{P}_{s,a,s'}^m, \overline{P}_{s,a,s'}^m] \forall c, s, a, s'$. Also, since our adversary oracle implementation requires as input both a planner mixed strategy α and its corresponding mixed adversary β of some MSNE (see Alg. 2), we generate inputs to the oracle by first randomly generating agent and adversary pure strategy sets according to the sampled $[\underline{P}_{s,a,s'}^m, \overline{P}_{s,a,s'}^m]$, then running one iteration of the double oracle using these strategy sets to generate the required α and β . Results, shown in Fig. 5, are reported as the average reward over ten random seeds. Our approach, REGRETMAXWI performs nearly as well as brute force for the small problem size, and vastly outperforms the naive random baseline even as the problem size increases to the scale of the maternal health intervention problem.



Figure 4: Evaluating planner oracle best response quality (WI4MS). Objective is to maximize reward (higher is better). No Action simulates the policy that takes action a = 0 on all arms at all time steps, representing a lower bound on reward. Random takes action a = 1 on K randomly chosen arms each round. Brute Force enumerates the entire feasible RMAB policy space, simulates the average reward of each policy, then returns the reward of the best-performing policy. Brute force can only be shown for small problem sizes, since its computation cost is exponential in N and K. Our approach WI4MS performs well across all problem sizes.

I Runtime Scalability of GROUPS

In Fig. 6 we demonstrate the runtime improves achieved by our GROUPS robust planning method as the number of groups M decreases, with number of arms N held constant. Given Thm 3, this demonstrates that we can achieve large scaling up of our robust planning without losing performance, under some mild assumptions, e.g., similarity of groups of arms.

J Experiment Setup Details

All algorithms were implemented in Python 3.7.4 and mathematical programs were solved using Gurobi version 9.0.3 via the gurobipy interface (Gurobi Optimization 2021). Experiments were run on a cluster running CentOS with Intel(R) Xeon(R) CPU E5-2683 v4 @ 2.1 GHz with 8GB of RAM and four processors.



Figure 5: Evaluating adversary oracle best response quality. Objective is to maximize regret (higher is better). Our approach, **RegretMaxWI**, performs nearly as well as a discretized brute force algorithm for small problems, and continues to perform far better than naive strategies for larger problems where brute force is intractable.



Figure 6: Run time scalability of GROUPS. Holding the number of arms N constant, the runtime of GROUPS improves significantly as the number of groups M decreases.

To compute regret, we simulate each planner strategy against the full set of the adversary's pure strategies (i.e., environment parameter settings, including both those computed by the adversary oracle as well as baseline responses pessimist, median, optimist, and random) to determine that which maximizes regret.

K ARMMAN Consent for Data Collection and Analysis

In this section, we provide information about consent related to data collection, analyzing data, data usage and sharing. We highlight that this work is part of a long-standing research collaboration with ARMMAN, with continuous analysis of data performed in close consultation with ARMMAN researchers.

K.1 Secondary Analysis and Data Usage

This study falls into the category of secondary analysis. To evaluate the performance of the algorithms, we use unlinked, anonymized data generated during the course of implementation of the program, i.e., previously collected engagement trajectories of different beneficiaries participating in the service call program. The proposed algorithms are evaluated via a simulation-based method discussed in Section 6. This paper does not involve deployment of the proposed algorithm or any other baselines to the service call program.

K.2 Consent for Data Collection and Sharing

The consent for collecting data is obtained from each of the participants of the service call program. The data collection process is carefully explained to the participants to seek their consent before collecting the data. The data is anonymized before sharing with us to ensure anonymity. Data exchange and use was regulated through clearly defined exchange protocols including anonymization, read-access only to researchers and restricted use of the data for research purposes only.

K.3 Universal Accessibility of Health Information

To allay further concerns: this simulation study focuses on improving quality of service calls. Even in the intended future application, all participants will receive the same weekly health information by automated message regardless of whether they are scheduled to receive service calls or not. The service call program does not withhold any information from the participants nor conduct any experimentation on the health information. The health information is always available to all participants, and participants can always request service calls via a free missed call service. In the intended future application our algorithm may only help schedule **additional** service calls to help beneficiaries who are likely to drop out of the program.

L Domain Descriptions

Maternal health The data used in this paper are anonymized, aggregated summary statistics of engagement behavior of mothers enrolled in ARMMAN's mMtira program, and *does not* contain any demographic information.

From the full dataset of 23,008 mothers, we chose a subset of 15,336 mothers from this cohort who have a least one record of intervention, then computed summary statistics (i.e., frequentist transition probabilities) over that subset. To create an arm–group mapping, we run K-means clustering on those probabilities, and compute uncertainty intervals via bootstrapping followed by multiple imputation to compute standard deviations of the means (Schomaker and Heumann 2018).

We now describe how we set up the environment settings used for each simulation variant. As mentioned, the experiments in Figure 3(g)–(i) use the summary engagement statistics from ARMMAN. To vary the simulated number of mothers in Figure 3(i), we begin with the groupings from the summary data but scale up the number of mothers in each group by a factor of 10 and 20 to reach 153.2K and 306.4K mothers, respectively, with budget scaled accordingly to K = 1000 and K = 7000.

Statistics on the uncertainty intervals and group sizes for the summary ARMMAN dataset are displayed in Figures 9 and 10, respectively.

TB Derived from data obtained from (Killian et al. 2019), which contains anonymous records of daily adherence to tuberculosis (TB) medication. We used the 8,350 records with



Figure 7: More experimental results evaluating the regret (lower is better) incurred by GROUPS, our robust solution approach, compared to non-robust baselines across a variety of problem settings. Regret is interpreted, in real-world terms, as the maximum *preventable* missed messages across the uncertainty space. The problem setting used by ARMMAN, which we use as default values, are K = 100, H = 10, N = 15,320, *actual* group size, and M = 40 groups. Experiment (a) uses real data obtained from (Mate et al. 2022); (b) and (c) use randomly sampled data, as described in section L. Each result is averaged over 30 random seeds.



Figure 8: Run time scaling of DDLPO (Killian et al. 2022) vs GROUPS. GROUPS is hundreds of times faster than DDLPO.

at least 30 days of adherence data. Though not collected in a grouped RMAB setting, we augmented the data to have groups by running K-Means grouping on the passive transition probabilities, and simulated uncertainty intervals for passive and active transitions as A_{σ} standard deviations of the mean of group centers. We then simulated uncertainty intervals: (1) for the passive transition probabilities as A_{σ} standard deviations (default $A_{\sigma} = 3$) about each group center and (2) for the active transition probabilities by adding a random value $\eta \sim N(0.3, 0.3)$ to the corresponding passive transition (i.e., always preferable to act), and creating an uncertainty interval about the mean $1.5 \times$ width of the passive uncertainty (i.e., less knowledge of active transitions vs. passive). These specific values were chosen to calibrate uncertainties to roughly match those of the maternal health domain, for which uncertainty statistics were available.

Statistics on the uncertainty intervals and group sizes for the TB dataset are displayed in Figures 11 and 12, respectively. **Synthetic** A benchmark domain from Killian et al. (2022) comprised of three "arm types" [U, V, W], each with their own intervals, designed so that non-robust policies have higher regret than robust ones. Specifically,

$$T_{s=0}^{n} = \begin{bmatrix} 0.5 & 0.5\\ 0.5 & 0.5 \end{bmatrix}, \quad T_{s=1}^{n} = \begin{bmatrix} 1.0 & 0.0\\ 1 - p_{n} & p_{n} \end{bmatrix}$$
(19)
$$p_{U} \in [0.00, 1.00]$$

where
$$p_{V} \in [0.05, 0.90] .$$
$$p_{W} \in [0.10, 0.95]$$

We augment the domain to allow homogeneous groups of each arm type, where the size and proportion of groups of each type may vary.

Randomly sampled data Used in Figure 3(b)-(c), we randomly generate transition probabilities and group sizes, drawn from normal distributions. We ensure that these transition probabilities are valid, that is that the probability transitioning to (or remaining in) the good state (s' = 1) with intervention (a = 1) is always higher than the probability of not intervening (a = 0), and similarly that the probability when starting in the engaging state (s = 1) is always higher than the probability of starting in the not engaging state (s = 0).



Figure 9: Statistics on the uncertainty intervals from the AR-MMAN data, averaged over 40 groups. Left shows the distribution of interval widths over all 40 groups. Right shows the distribution of interval midpoints over all 40 groups. Some uncertainty intervals are wider than 0.8, but the majority have width below 0.4. Uncertainty intervals of most groups are in the range [0.3, 0.7].



Figure 10: Distribution of group sizes for the ARMMAN data.



Figure 11: Summary statistics on the uncertainty intervals from TB data obtained from (Killian et al. 2019), averaged over 60 groups. Left shows the distribution of interval widths over all groups. Right shows the distribution of interval midpoints over all groups. In general, transition probability medians are closer to 1 for this domain than the ARMMAN domain.



Figure 12: Distribution of group sizes for TB adherence data.