

---

# Fair Resource Allocation in Weakly Coupled Markov Decision Processes

---

Xiaohui Tu<sup>1,2</sup>

Yossiri Adulyasak<sup>1</sup>

Nima Akbarzadeh<sup>1,2</sup>

Erick Delage<sup>1,2</sup>

<sup>1</sup>GERAD & HEC Montréal, <sup>2</sup>Mila - Quebec AI Institute

## Abstract

We consider fair resource allocation in sequential decision-making environments modeled as weakly coupled Markov decision processes, where resource constraints couple the action spaces of  $N$  sub-Markov decision processes (sub-MDPs) that would otherwise operate independently. We adopt a fairness definition using the generalized Gini function instead of the traditional utilitarian (total-sum) objective. After introducing a general but computationally prohibitive solution scheme based on linear programming, we focus on the homogeneous case where all sub-MDPs are identical. For this case, we show for the first time that the problem reduces to optimizing the utilitarian objective over the class of “permutation invariant” policies. This result is particularly useful as we can exploit efficient algorithms that optimize the utilitarian objective such as Whittle index policies in restless bandits to solve the problem with this fairness objective. For more general settings, we introduce a count-proportion-based deep reinforcement learning approach. Finally, we validate our theoretical findings with comprehensive experiments, confirming the effectiveness of our proposed method in achieving fairness.

## 1 INTRODUCTION

Machine learning (ML) algorithms play a significant role in automated decision-making processes, influencing our daily lives. Mitigating biases within the ML pipeline is crucial to ensure fairness and gener-

ate reliable outcomes (Caton and Haas, 2024). Extensive research has been conducted to enhance fairness across various applications, such as providing job hiring services (van den Broek et al., 2020; Cimpean et al., 2024), assigning credit scores and loans (Kozdoi et al., 2022), and delivering healthcare services (Farnadi et al., 2021; Chen et al., 2023).

However, most real-world decision processes are sequential in nature and past decisions may have a long-term impact on equity (D’Amour et al., 2020). For example, if people are unfairly denied credit or job opportunities early in their careers, there would be long-term consequences on opportunities for advancement (Liu et al., 2018). Another motivating example is taxi dispatching. If certain areas are consistently prioritized over others, then there can be long-term disparities in service accessibility. This may lead to long waiting times for passengers in certain neighborhoods, while taxis run empty and seek passengers in other areas (Liu et al., 2021; Guo et al., 2023).

Fairness is a complex and multi-faceted concept, and there are many different ways in which it can be operationalized and measured. We resort to the generalized Gini social welfare function (GGF) (Weymark, 1981), which covers various fairness measures as special cases. The long-term impacts of fair decision dynamics have recently been approached using Markov decision processes (MDPs) (Wen et al., 2021; Puranik et al., 2022; Ghalme et al., 2022). Studying fairness in MDPs helps mitigate bias and inequality in decision-making processes and evaluate their broader societal and operational impacts across diverse applications.

To the best of our knowledge, we are the first to incorporate fairness considerations in the form of the GGF objective within weakly coupled Markov decision processes (WCMDPs) (Hawkins, 2003; Adelman and Mersereau, 2008), which can be considered as an extension of restless multi-arm bandit problems (RMABs) (Hawkins, 2003; Zhang, 2022) to multi-action and multi-resource settings. This model is particularly relevant to resource allocation problems, as it captures the complex interactions of coupled MDPs

(arms) over time restricted by limited resource availability, and allows the applicability of our work to various applications in scheduling (Saure et al., 2012; El Shar and Jiang, 2024), application screening (Gast et al., 2024), budget allocation (Boutilier and Lu, 2016), and inventory (El Shar and Jiang, 2024).

**Contributions** Our contributions are as follows. *Theoretically*, we reformulate the WCMDP problem with the GGF objective as a linear programming (LP) problem, and show that, under symmetry, it reduces to maximizing the average expected total discounted rewards, called the utilitarian approach. *Methodologically*, we propose a state count approach to further simplify the problem, and introduce a count proportion-based deep reinforcement learning (RL) method that can solve the reduced problem efficiently and can scale to larger cases by assigning resources proportionally to the number of stakeholders. *Experimentally*, we design various experiments to show the GGF-optimality, flexibility, scalability and efficiency of the proposed deep RL approach. We benchmark our approach against the Whittle index policy on machine replacement applications modeled as RMABs (Akbarzadeh and Mahajan, 2019), showing the effectiveness of our method in achieving fair outcomes under different settings.

There are two studies closely related to our work. The first work by Gast et al. (2024) considers symmetry simplification and count aggregation MDPs. They focus on solving an LP model repeatedly with a total-sum objective to obtain asymptotic optimal solutions when the number of coupled MDPs is very large, whereas we explicitly address the fairness aspect and exploit a state count representation to design scalable deep RL approaches. The second work by Siddique et al. (2020) integrates the fair Gini multi-objective RL to treat every user equitably. This fair optimization problem is later extended to the decentralized cooperative multi-agent RL by Zimmer et al. (2021), and further refined to incorporate preferential treatment with human feedback by Siddique et al. (2023) and Yu et al. (2023). In contrast, our work demonstrates that the WCMDP with the GGF objective and identical coupled MDPs reduces to a much simpler utilitarian problem, which allows us to exploit its structure to develop efficient and scalable algorithms. A more comprehensive literature review on fairness in resource allocation, MDPs, RL, and RMABs, is provided in Appendix A to clearly position our work.

## 2 BACKGROUND

We start by reviewing infinite-horizon WCMDPs and introducing the GGF for encoding fairness. We then

define the fair optimization problem and provide an exact solution scheme based on linear programming.

**Notation** Let  $[N] := \{1, \dots, N\}$  for any integer  $N$ . For any vector  $\mathbf{v} \in \mathbb{R}^N$ , the  $n$ -th element is denoted as  $v_n$  and the average value as  $\bar{v} = \frac{1}{N} \sum_{n=1}^N v_n$ . An indicator function  $\mathbb{I}\{x \in A\}$  equals 1 if  $x \in A$  and 0 otherwise. For any set  $X$ ,  $\Delta(X)$  represents the set of all probability distributions over  $X$ . We let  $\mathbb{S}^N$  be the set of all  $N!$  permutations of the indices in  $[N]$  and  $\mathcal{G}^N$  be the set of all permutation operators so that  $Q \in \mathcal{G}^N$  if and only if there exists a  $\sigma \in \mathbb{S}^N$  such that  $Q\mathbf{v}(n) = \mathbf{v}_{\sigma(n)}$  for all  $n \in [N]$  when  $\mathbf{v} \in \mathbb{R}^N$ .

### 2.1 The Weakly Coupled MDP

We consider  $N$  MDPs indexed by  $n \in \mathcal{N} := [N]$  interacting in discrete-time over an infinite horizon  $t \in \mathcal{T} := \{0, 1, \dots, \infty\}$ . The  $n$ -th MDP  $\mathcal{M}_n$ , also referred as sub-MDP, is defined by a tuple  $(\mathcal{S}_n, \mathcal{A}_n, p_n, r_n, \mu_n, \gamma)$ , where  $\mathcal{S}_n$  is a finite set of states with cardinality  $S$ , and  $\mathcal{A}_n$  is a finite set of actions with cardinality  $A$ . The transition probability function is defined as  $p_n(s'_n | s_n, a_n) = \mathbb{P}(s_{t+1,n} = s'_n | s_{t,n} = s_n, a_{t,n} = a_n)$ , which represents the probability of reaching state  $s'_n \in \mathcal{S}_n$  after performing action  $a_n \in \mathcal{A}_n$  in state  $s_n \in \mathcal{S}_n$  at time  $t$ . The reward function  $r_n(s_n, a_n)$  denotes the immediate real-valued reward obtained by executing action  $a_n$  in state  $s_n$ . Although the transition probabilities and the reward function may vary with the sub-MDP  $n$ , we assume that they are stationary across all time steps for simplicity. The initial state distribution is represented by  $\mu_n \in \Delta(\mathcal{S}_n)$ , and the discount factor, common to all sub-MDPs, is denoted by  $\gamma \in [0, 1)$ .

An infinite-horizon WCMDP  $\mathcal{M}^{(N)}$  consists of  $N$  sub-MDPs, where each sub-MDP is independent of the others in terms of state transitions and rewards. They are linked to each other solely through a set of  $K$  constraints on their actions at each time step. Formally, the WCMDP is defined by a tuple  $(\mathcal{S}^{(N)}, \mathcal{A}^{(N)}, p^{(N)}, \mathbf{r}, \mu^{(N)}, \gamma)$ , where the state space  $\mathcal{S}^{(N)}$  is the Cartesian product of individual state spaces, and the action space  $\mathcal{A}^{(N)}$  is a subset of the Cartesian product of action spaces, defined as  $\mathcal{A}^{(N)} := \{(a_1, \dots, a_N) \mid \sum_{n=1}^N d_{k,n}(a_n) \leq b_k, \forall k \in \mathcal{K}, a_n \in \mathcal{A}_n\}$ , where  $\mathcal{K} := [K]$  is the index set of constraints,  $d_{k,n}(a_n) \in \mathbb{R}_+$  represents the consumption of the  $k$ -th resource consumption by the  $n$ -th MDP when action  $a_n$  is taken, and  $b_k \in \mathbb{R}_+$  the available resource of type  $k$ .<sup>1</sup> We define an idle action that consumes no resources for any resource  $k$  to ensure that the feasible action space is non-empty.

<sup>1</sup>Actually,  $b_k \leq \sum_{n=1}^N \max_{a_n \in \mathcal{A}_n} d_{k,n}(a_n)$ , w.l.o.g.

The state transitions of the sub-MDPs are independent, so the system transits from state  $\mathbf{s}$  to state  $\mathbf{s}'$  for a given feasible action  $\mathbf{a}$  at time  $t$  with probability  $p^{(N)}(\mathbf{s}'|\mathbf{s}, \mathbf{a}) = \prod_{n=1}^N p_n(s'_n|s_n, a_n) = \prod_{n=1}^N \mathbb{P}(s_{t+1,n} = s'_n | s_{t,n} = s_n, a_{t,n} = a_n)$ . After choosing an action  $\mathbf{a} \in \mathcal{A}^{(N)}$  in state  $\mathbf{s} \in \mathcal{S}^{(N)}$ , the decision maker receives rewards defined as  $\mathbf{r}(\mathbf{s}, \mathbf{a}) = (r_1(s_1, a_1), \dots, r_N(s_N, a_N))$  with each component representing the reward associated with the respective sub-MDP  $\mathcal{M}_n$ . We employ a vector form for the rewards to offer the flexibility for formulating fairness objectives on individual expected total discounted rewards in later sections.

We consider stationary Markovian policy  $\pi : \mathcal{S}^{(N)} \times \mathcal{A}^{(N)} \rightarrow [0, 1]$ , with notation  $\pi(\mathbf{s}, \mathbf{a})$  capturing the probability of performing action  $\mathbf{a}$  in state  $\mathbf{s}$ . The initial state  $\mathbf{s}_0$  is sampled from the distribution  $\mu^{(N)}$ . Using the discounted-reward criteria, the state-value function  $V_n^\pi$  specific to the  $n$ -th sub-MDP  $\mathcal{M}_n$ , starting from an arbitrary initial state  $\mathbf{s}_0$  under policy  $\pi$ , is defined as  $V_n^\pi(\mathbf{s}_0) := \mathbb{E}_\pi [\sum_{t=0}^{\infty} \gamma^t r_n(s_{t,n}, a_{t,n}) | \mathbf{s}_0]$ , where  $\mathbf{a}_t \sim \pi(\mathbf{s}_t, \cdot)$ . The joint state-value vector-valued function  $\mathbf{V}^\pi(\mathbf{s}_0) : \mathcal{S}^{(N)} \rightarrow \mathbb{R}^N$  is defined as the column vector of expected total discounted rewards for all sub-MDPs under policy  $\pi$ , i.e.,  $\mathbf{V}^\pi(\mathbf{s}_0) := (V_1^\pi(\mathbf{s}_0), V_2^\pi(\mathbf{s}_0), \dots, V_N^\pi(\mathbf{s}_0))^\top$ . We define  $\mathbf{V}_0^\pi$  as the expected vectorial state-value under initial distribution  $\mu^{(N)}$ , i.e.,

$$\mathbf{V}_0^\pi := \mathbb{E}[\mathbf{V}^\pi(\mathbf{s}_0) | \mathbf{s}_0 \sim \mu^{(N)}]. \quad (1)$$

## 2.2 The Generalized Gini Function

The vector  $\mathbf{V}_0^\pi$  represents the expected utilities for sub-MDPs. A social welfare function aggregates these utilities into a scalar, measuring fairness in utility distribution with respect to a maximization objective.

Social welfare functions can vary depending on the values of a society, such as  $\alpha$ -fairness (Mo and Walrand, 2000), Nash social welfare (Fan et al., 2022; Mandal and Gan, 2022), or max-min fairness (Bistritz et al., 2020; Cousins et al., 2022). Following Siddique et al. (2020), we require a fair solution to meet three properties: efficiency, impartiality, and equity. This motivates the use of GGF from economics (Weymark, 1981), which satisfies these properties. For  $N$  sub-MDPs, GGF is defined as  $\text{GGF}_{\mathbf{w}}[\mathbf{v}] := \min_{\sigma \in \mathbb{S}^N} \sum_{n=1}^N w_n v_{\sigma(n)}$ , where  $\mathbf{v} \in \mathbb{R}^N$ ,  $\mathbf{w} \in \Delta(\mathcal{N})$  is non-increasing in  $n$ , i.e.,  $w_1 \geq w_2 \geq \dots \geq w_N$ . Intuitively, since  $\text{GGF}_{\mathbf{w}}[\mathbf{v}] = \sum_{n=1}^N w_n v_{\sigma^*(n)}$  with  $\sigma^*$  as the minimizer, which reorders the terms of  $\mathbf{v}$  from lowest to largest, it computes the weighted sum of  $\mathbf{v}$  assigning larger weights to its lowest components. When the order of sub-MDPs is fixed, we use the equivalent formulation  $\text{GGF}_{\mathbf{w}}[\mathbf{v}] = \min_{\sigma \in \mathbb{S}^N} \sum_n w_{\sigma(n)} v_n$  as per-

muting either vector results in the same outcome.

As discussed in Siddique et al. (2020), GGF can reduce to special cases by setting its weights to specific values, including the maxmin egalitarian approach ( $w_1 \rightarrow 1, w_2 \rightarrow 0, \dots, w_N \rightarrow 0$ ) (Rawls, 1971), regularized maxmin egalitarian ( $w_1 \rightarrow 1, w_2 \rightarrow \epsilon, \dots, w_N \rightarrow \epsilon$ ), leximin notion of fairness ( $w_k/w_{k+1} \rightarrow \infty$ ) (Rawls, 1971; Moulin, 1991), and the utilitarian approach formally defined below for the later use in reducing the GGF problem.

**Definition 2.1 (Utilitarian Approach)** *The utilitarian approach within the GGF framework is obtained by setting equal weights for all individuals, i.e.,  $\mathbf{w}_{1/N} := \mathbf{1}/N$  so that  $\text{GGF}_{\mathbf{w}_{1/N}}[\mathbf{v}] = \min_{\sigma \in \mathbb{S}^N} \sum_{n=1}^N \frac{1}{N} v_{\sigma(n)} = \frac{1}{N} \sum_{n=1}^N v_n = \bar{v}$ .*

The utilitarian approach maximizes average utilities over all individuals but does not guarantee fairness in utility distribution, as some sub-MDPs may be disadvantaged to increase overall utility. The use of GGF offers flexibility by encoding various fairness criteria in a structured way. Moreover,  $\text{GGF}_{\mathbf{w}}[\mathbf{v}]$  is concave in  $\mathbf{v}$ , which has nice properties for problem reformulation.

## 2.3 The GGF-WCMDP Problem

By combining GGF and the vectored values from the WCMDP in (1), the goal of the GGF-WCMDP problem (2) is defined as finding a stationary policy  $\pi$  that maximizes the GGF of the expected total discounted rewards, i.e.,  $\max_{\pi} \text{GGF}_{\mathbf{w}}[\mathbf{V}_0^\pi]$  that is equivalent to

$$\max_{\pi} \min_{\sigma \in \mathbb{S}^N} \mathbf{w}_{\sigma}^\top \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s_t, \mathbf{a}_t) \middle| \mathbf{s}_0 \sim \mu^{(N)} \right]. \quad (2)$$

We note that Lemma 3.1 in Siddique et al. (2020) establishes the optimality of stationary Markov policies for any multi-objective discounted infinite-horizon MDP under the GGF criterion. To obtain an optimal policy for the GGF-WCMDP problem (2), we introduce the following LP model with the GGF objective (GGF-LP):

$$\max_{\lambda, \nu, q} \sum_{i=1}^N \lambda_i + \sum_{j=1}^N \nu_j \quad (3a)$$

$$\text{s.t.} \quad \lambda_i + \nu_j \leq w_i \sum_{\mathbf{s} \in \mathcal{S}^{(N)}} \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} r_j(\mathbf{s}, \mathbf{a}) q(\mathbf{s}, \mathbf{a}), \quad \forall i, j \in \mathcal{N}$$

$$\sum_{\mathbf{a} \in \mathcal{A}^{(N)}} q(\mathbf{s}, \mathbf{a}) - \gamma \sum_{\mathbf{s}' \in \mathcal{S}^{(N)}} \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} q(\mathbf{s}', \mathbf{a}) p^{(N)}(\mathbf{s} | \mathbf{s}', \mathbf{a})$$

$$= \mu^{(N)}(\mathbf{s}), \quad \forall \mathbf{s} \in \mathcal{S}^{(N)}, \quad (3b)$$

$$q(\mathbf{s}, \mathbf{a}) \geq 0, \quad \forall \mathbf{s} \in \mathcal{S}^{(N)}, \quad \forall \mathbf{a} \in \mathcal{A}^{(N)}. \quad (3c)$$

See Appendix D.1 for details on obtaining model (3) that exploits the dual linear programming formula-

tion for solving discounted MDPs. Here,  $q(\mathbf{s}, \mathbf{a})$  represents the total discounted visitation frequency for state-action pair  $(\mathbf{s}, \mathbf{a})$ , starting from  $\mathbf{s}_0$ .

The dual form separates dynamics from rewards, with the expected discounted reward for sub-MDP  $n$  given by  $\sum_{\mathbf{s} \in \mathcal{S}^{(N)}} \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} r_n(\mathbf{s}, \mathbf{a}) q(\mathbf{s}, \mathbf{a})$ . The one-to-one mapping between the solution  $q(\mathbf{s}, \mathbf{a})$  and an optimal policy  $\pi(\mathbf{s}, \mathbf{a})$  is  $\pi(\mathbf{s}, \mathbf{a}) = q(\mathbf{s}, \mathbf{a}) / \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} q(\mathbf{s}, \mathbf{a})$ .

Scalability is a critical challenge in obtaining exact solutions as the state and action spaces grow exponentially with respect to the number of sub-MDPs, making the problem intractable. We thus explore approaches that exploit symmetric problem structures, apply count-based state aggregation, and use RL-based approximation methods, to address this scalability issue, which will be discussed next.

### 3 UTILITARIAN REDUCTION UNDER SYMMETRIC SUB-MDPS

In Section 3.1, we will formally define the concept of symmetric WCMDPs (definition 3.1) and prove that an optimal policy of the GGF-WCMDP problem can be obtained by solving the utilitarian WCMDP using “permutation invariant” policies. This enables the use of Whittle index policies in the RMAB setting while, for the more general setting, Section 3.2 proposes a count aggregation MDP reformulation that will be solved using deep RL in Section 4.

#### 3.1 GGF-WCMDP Problem Reduction

We start with formally defining the conditions for a WCMDP to be considered symmetric.

##### Definition 3.1 (Symmetric WCMDP)

A WCMDP is symmetric if

1. **(Identical Sub-MDPs)** Each sub-MDP is identical, i.e.,  $\mathcal{S}_n = \mathcal{S}$ ,  $\mathcal{A}_n = \mathcal{A}$ ,  $p_n = p$ ,  $r_n = r$ ,  $\mu_n = \mu$ , for all  $n \in \mathcal{N}$ , and for some  $(\mathcal{S}, \mathcal{A}, p, r, \mu, \gamma)$  tuple.
2. **(Symmetric resource consumption)** For any  $k \in \mathcal{K}$ , the number of resources consumed is the same for each sub-MDP, i.e.,  $d_{k,n}(a_n) = d_k(a_n)$  for all  $n \in \mathcal{N}$ , and for some  $d_k(\cdot)$ .
3. **(Permutation-Invariant Initial Distribution)** For any permutation operator  $Q \in \mathcal{G}^N$ , the probability of selecting the permuted initial state  $Q\bar{\mathbf{s}}_0$  is equal to that of selecting  $\bar{\mathbf{s}}_0$ , i.e.,  $\mu^{(N)}(\bar{\mathbf{s}}_0) = \mu^{(N)}(Q\bar{\mathbf{s}}_0)$ ,  $\forall \bar{\mathbf{s}}_0 \in \mathcal{S}^{(N)}$ ,  $\forall Q \in \mathcal{G}^N$ .

The conditions of symmetric WCMDP identify a class

of WCMDPs that is invariant under any choice of indexing for the sub-MDPs. This gives rise to the notion of “permutation invariant” policies (see definition 1 in Cai et al. (2021)) and the question of whether this class of policies is optimal for symmetric WCMDPs.

##### Definition 3.2 (Permutation Invariant Policy)

A Markov stationary policy  $\pi$  is said to be permutation invariant if the probability of selecting action  $\mathbf{a}$  in state  $\mathbf{s}$  is equal to that of selecting the permuted action  $Q\mathbf{a}$  in the permuted state  $Q\mathbf{s}$ , for all  $Q \in \mathcal{G}^N$ . Formally, this can be expressed as  $\pi(\mathbf{s}, \mathbf{a}) = \pi(Q\mathbf{s}, Q\mathbf{a})$ , for all  $Q \in \mathcal{G}^N$ ,  $\mathbf{s} \in \mathcal{S}^{(N)}$  and  $\mathbf{a} \in \mathcal{A}^{(N)}$ .

This symmetry ensures that the expected state-value function, when averaged over all trajectories, is identical for each sub-MDP, leading to a uniform state-value representation. From this observation, and applying Theorem 6.9.1 from (Puterman, 2005), we construct a permutation-invariant policy from any policy, resulting in uniform state-value representation (Lemma 3.3).

##### Lemma 3.3 (Uniform State-Value Representation)

If a WCMDP is symmetric, then for any policy  $\pi$ , there exists a corresponding permutation invariant policy  $\bar{\pi}$  such that the vector of expected total discounted rewards for all sub-MDPs under  $\bar{\pi}$  is equal to the average of the expected total discounted rewards for each sub-MDP, i.e.,  $\mathbf{V}_0^{\bar{\pi}} = \frac{1}{N} \sum_{n=1}^N \mathbf{V}_{0,n}^{\pi} \mathbf{1}$ .

The proof is detailed in Appendix B.3. Furthermore, one can use the above lemma to show that the optimal policy for the GGF-WCMDP problem (2) under symmetry can be recovered from solving the problem with equal weights, i.e., the utilitarian approach. Our main result is presented in the following theorem. See Appendix B.4 for a detailed proof.

**Theorem 3.4 (Utilitarian Reduction)** For a symmetric WCMDP, let  $\Pi_{\mathbf{1}/N, PI}^*$  be the set of optimal policies for the utilitarian approach that is permutation invariant, then  $\Pi_{\mathbf{1}/N, PI}^*$  is necessarily non-empty and all  $\pi_{\mathbf{1}/N, PI}^* \in \Pi_{\mathbf{1}/N, PI}^*$  satisfies

$$\text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\pi_{\mathbf{1}/N, PI}^*}] = \max_{\pi} \text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\pi}], \forall \mathbf{w} \in \Delta(N).$$

This theorem simplifies solving the GGF-WCMDP problem by reducing it to an equivalent utilitarian problem, showing that at least one permutation-invariant policy is optimal for the original GGF-WCMDP problem and the utilitarian reduction. Therefore, we can restrict the search for optimal policies to this specific class of permutation-invariant policies. The utilitarian approach does not compromise the GGF optimality and allows us to leverage more efficient and scalable techniques to solve the GGF-WCMDP problem (Eq. 2), such as the Whittle index

policies for RMABs, as demonstrated in the experimental section. We note that the utilitarian reduction theorem can be extended to a broader class of fairness measures, such as  $\alpha$ -fairness, as long as the fairness measure is concave, permutation invariant, and constant vector invariant (Corollary B.4.1). See Appendix B.5 for a detailed proof.

### 3.2 The Count Aggregation MDP

Assuming symmetry across all  $N$  sub-MDPs and using a permutation-invariant policy within a utilitarian framework allows us to simplify the global MDP by aggregating the sub-MDPs based on their state counts and tracking the number of actions taken in each state. Since each sub-MDP follows the same transition probabilities and reward structure, we can represent the entire system more compactly. This symmetry consideration is practical in many real-world applications where a large number of identical or interchangeable identities demand fair and efficient treatment, such as patients in healthcare or taxi drivers in public transportation services. By leveraging symmetry, we can reduce computational complexity for scalable fair solutions while inherently enforcing fairness as the policy treats all sub-MDPs equivalently.

Motivated by the symmetry simplification representation in Gast et al. (2024) for the utilitarian objective, we consider an aggregation  $\phi = (f, g_s)$ , where  $f : \mathcal{S}^{(N)} \rightarrow \mathbb{N}^S$  maps state  $\mathbf{s}$  to a count representation  $\mathbf{x}$  with  $x_s$  denoting the number of sub-MDPs in the  $s$ -th state. Similarly,  $g_s : \mathcal{A}^{(N)} \rightarrow \mathbb{N}^{S \times A}$  maps action  $\mathbf{a}$  to a count representation  $\mathbf{u}$ , where  $u_{s,a}$  indicates the number of MDPs at  $s$ -th state that performs  $a$ -th action. We can formulate the count aggregation MDP (definition 3.5). The details on obtaining the exact form are in Appendix C.

**Definition 3.5 (Count Aggregation MDP)** *The count aggregation MDP  $\mathcal{M}_\phi$  derived from a WCMDP  $(\mathcal{S}^{(N)}, \mathcal{A}^{(N)}, p^{(N)}, \mathbf{r}, \mu^{(N)}, \gamma)$  consists of the elements  $(\mathcal{S}_f^{(N)}, \mathcal{A}_{g_s}^{(N)}, p_\phi^{(N)}, \bar{r}_\phi, \mu_f^{(N)}, \gamma)$ .*

Both representations lead to the same optimization problem as established in Gast et al. (2024) when the objective is utilitarian. Using the count representation, the mean expected total discounted reward  $\bar{V}_0^{\pi_{1/N}}$  for a WCMDP  $\mathcal{M}^{(N)}$  with permutation invariant distribution  $\mu^{(N)}$  and equal weights  $\mathbf{w}_{1/N}$  (Theorem 3.4) is then equivalent to the expected total discounted mean reward  $\bar{V}_0^{\pi_\phi}$  for the count aggregation MDP  $\mathcal{M}_\phi$  given the policy  $\pi_\phi : \mathcal{S}_f^{(N)} \rightarrow \Delta(\mathcal{A}_{g_s}^{(N)})$  under aggregation mapping with initial distribution  $\mu_f^{(N)}$ , i.e.,  $\bar{V}_0^{\pi_{1/N}} = \frac{1}{N} \sum_{n=1}^N V_{0,n}^{\pi_{1/N}} = \frac{1}{S} \sum_{s=1}^S V_{0,s}^{\pi_\phi} = \bar{V}_0^{\pi_\phi}$ .

The objective in Eq. 2 is therefore reformulated as  $\max_{\pi_\phi} \bar{V}_0^{\pi_\phi}$ , i.e.,

$$\max_{\pi_\phi} \frac{1}{S} \mathbb{E}_{\pi_\phi} \left[ \sum_{t=0}^{\infty} \gamma^t \bar{r}_\phi(\mathbf{x}_t, \mathbf{u}_t) \mid \mathbf{x}_0 \sim \mu_f^{(N)} \right]. \quad (4)$$

An LP method is provided to solve the count aggregation MDP in Appendix D.2.

## 4 COUNT-PROPORTION-BASED DRL

We now consider the situation where the transition dynamics  $p_\phi^{(N)}$  are unknown and the learner computes the (sub-)optimal policy through trial-and-error interactions with the environment. In Section 4.1, we introduce a count-proportion-based deep RL (CP-DRL) approach. This method incorporates a stochastic policy neural network with fixed-sized inputs and outputs, designed for optimizing resource allocation among stakeholders under constraints with count representation. In Section 4.2, we detail the priority-based sampling procedure used to generate count actions.

### 4.1 Stochastic Policy Neural Network

One key property of the count aggregation MDP is that the dimensions of the state space  $\mathcal{S}_f^{(N)}$  and the action space  $\mathcal{A}_{g_s}^{(N)}$  are constant and irrespective of the number of sub-MDPs. To further simplify the analysis and eliminate the influence of  $N$ , we define the count state proportion as  $\bar{\mathbf{x}} = \mathbf{x}/N$  and the resource proportion constraint for each resource  $k$  as  $\bar{b}_k = b_k / (N \max_{a \in \mathcal{A}} d_k(a)) \in [0, 1]$ . This converts the states into a probability distribution, allowing generalization when dealing with a large number of agents. The stochastic policy network in Figure 1 is designed to handle the reduced count aggregation MDP problem (4) by transforming the tuple  $(\bar{\mathbf{x}}, \bar{\mathbf{b}})$  into a priority score matrix  $\mathbf{U}$  and a resource-to-use proportion vector  $\tilde{\mathbf{p}}$ , which are then used to generate count actions  $\mathbf{u}$  via a sampling procedure (discussed in Section 4.2).

The policy network features fixed-size inputs and outputs, enabling scalability in large-scale systems without requiring structural modifications when adjusting the number of resources or machines. The input consists of a fixed-size vector of size  $S + K$ , combining the count state proportion  $\bar{\mathbf{x}} \in [0, 1]^S$  and the resource proportion  $\bar{\mathbf{b}} \in [0, 1]^K$ . The policy network processes these inputs to produce outputs of size  $S \times A + K$ , which include a matrix  $\mathbf{U} \in (0, 1)^{S \times A}$  representing the priority scores for selecting count actions and a vector  $\tilde{\mathbf{p}} \in [0, 1]^K$  representing the proportion of resource

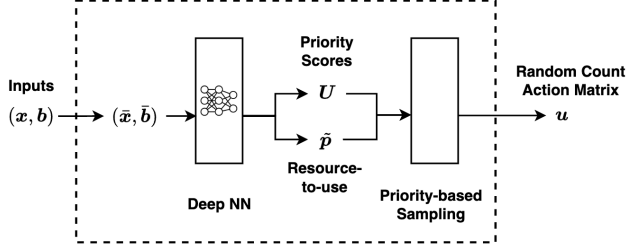


Figure 1: CP-based Stochastic Policy Neural Network

usage relative to the total available resources  $\bar{b}$ .

The advantages of adding additional resource proportion nodes  $\tilde{p}$  to the output layer are twofold. First, it reduces the computational effort required to ensure that the resource constraint is satisfied by restricting the resource-to-use proportions  $\tilde{p}$  to be element-wise proportional to  $\bar{b}$ . Second, since the optimal policy may not always use all available resources, we incorporate the additional nodes to capture the complex relationships between different states for more effective strategies to allocate resources.

## 4.2 Priority-based Sampling Procedure

Priority-based sampling presents a challenge since legal actions depend on state and resource constraints. To address this, a masking mechanism prevents the selection of invalid actions. Each element  $u_{s,a}^p \in U$  represents the priority score of taking the  $a$ -th action in the  $s$ -th state. When the state count  $x_s$  is zero, it implies the absence of sub-MDPs in this state, and the corresponding priority score is masked to zero. Legal priorities are thus defined for states with non-zero counts, i.e.,  $u_{s,a}^p = 0$  if  $x_s = 0$  for all  $s \in [S], a \in [A]$ .

Since the selected state-action pairs must also satisfy multi-resource constraints, we introduced a forbidden set  $\mathcal{F}$ , which specifies the state-action pairs that are excluded from the sampling process. The complete procedure is outlined in Algorithm 1. The advantage of this approach is that the number of steps does not grow exponentially with the number of sub-MDPs. In the experiments, after obtaining the count action  $u$ , a model simulator is used to generate rewards and the next state as described in Algorithm 2 in Appendix E. The simulated outcomes are used for executing policy gradient updates and estimating state values.

One critical advantage of using CP-DRL is its *scalability*. More specifically, the approach is designed to handle variable sizes of stakeholders  $N$  and resources  $K$  while preserving the number of aggregated count states constant for a given WCMDP. By normalizing

inputs to fixed-size proportions, the network can seamlessly adapt to different scales, making it highly adaptable. Moreover, the fixed-size inputs allow *flexibility* that the neural network is trained once and used in multiple tasks with various numbers of stakeholders and resource limitations.

---

### Algorithm 1 Count Action Sampling Based on Priority Scores

---

**Input:** Count state  $x$ , priority score matrix  $U$ , resource limitations  $b$ , resource-to-use proportion  $\tilde{p}$ , resource consumption function  $d(a)$

**Initialize:**  $\tilde{b} \leftarrow b \cdot \tilde{p}$ ,  $u \leftarrow \mathbf{0}_{S \times A}$ ,  $\mathcal{F} \leftarrow \emptyset$

Apply masking to  $U$  and update the forbidden set by  $\mathcal{F} \leftarrow \mathcal{F} \cup \{(s, a) \mid u_{s,a}^p = 0 \text{ for } s \in [S] \text{ and } a \in [A]\}$

**while**  $|\mathcal{F}| < S \times A$  **do**

    Sample a state-action index pair  $(s, a) \notin \mathcal{F}$  with the probability proportional to  $U$

**if**  $d_k(a) \leq \tilde{b}_k$  for all  $k$  **then**

        Update  $u_{s,a} \leftarrow u_{s,a} + 1$ ,  $x_s \leftarrow x_s - 1$

        Update  $\tilde{b}_k \leftarrow \tilde{b}_k - d_k(a)$  for all  $k$

**if**  $x_s = 0$  **then**

        Add all actions for the  $s$ -th state to forbidden set:  $\mathcal{F} \leftarrow \mathcal{F} \cup \{(s, a) \mid \forall a \in [A]\}$

**end if**

**else**

        Add  $(s, a)$  to forbidden set  $\mathcal{F} \leftarrow \mathcal{F} \cup \{(s, a)\}$

**end if**

**end while**

**Return:** Count action matrix  $u$

---

## 5 EXPERIMENTAL RESULTS

We apply our methods to the machine replacement problem (Delage and Mannor, 2010; Akbarzadeh and Mahajan, 2019), providing a scalable framework for evaluating the CP-DRL approach as problem size and complexity increase. We focus on a single resource ( $K = 1$ ) and binary action ( $A = 2$ ) for each machine, allowing validation against the Whittle index policy for RMABs (Whittle, 1988). We applied various DRL algorithms, including *Soft Actor Critic* (SAC), *Twin Delayed DDPG* (TD3), and *Proximal Policy Optimization* (PPO). Among these, PPO algorithm (Schulman et al., 2017) consistently delivers the most stable and high-quality performance. We thus choose PPO as the main algorithm for our CP-DRL approach (see Section 4). Our code is provided on GitHub.<sup>2</sup>

**Machine Replacement Problem** The problem consists of  $N$  identical machines following Markovian deterioration rules with  $S$  states representing aging stages. The state space  $\mathcal{S}^{(N)}$  is a Cartesian product. At each decision stage, actions  $a$  are applied to all machines under resource constraints, with action  $a_n$  representing operation (*passive* action) or replacement (*active* action). Resource consumption  $d_n(a_n)$  is

<sup>2</sup><https://github.com/x-tu/GGF-wcMDP>.

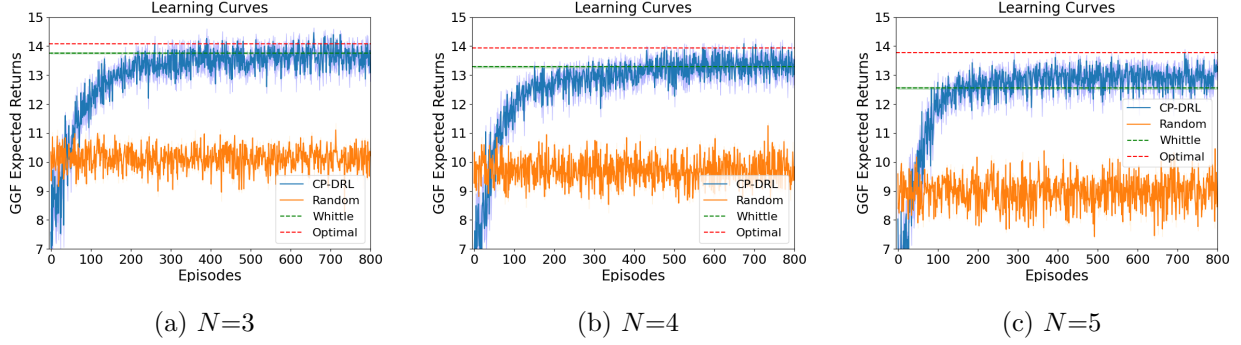


Figure 2: **(Colored) Learning Curves for Different Numbers of Machines ( $N$  from 3 to 5).** Experimental results for the *Exponential-RCCC* scenario are shown with y-axes starting at 7 for zoom-in. Red dashed lines represent the OPT values, green dashed lines show the WIP performance, blue lines depict CP-DRL learning curves over 800 episodes, and orange lines show the RDM performance. Shaded areas indicate the standard deviation across 5 runs.

1 for replacements and 0 for operations, with up to  $b$  replacements per time step. The costs range from 0 to 1, transformed to fit the reward representation by multiplying by -1 and adding 1. Machines degrade if not replaced and remain in state  $S$  until replaced. Refer to Appendix F.1 for cost structures and transition probabilities. We choose operational and replacement costs across two presets to capture different scenarios (see Appendix F.2 for details): *i) Exponential-RCCC* and *ii) Quadratic-RCCC*.

The goal is to find a fair policy that maximizes the GGF score over the expected total discounted mean rewards with count aggregation MDP. In cases like electricity or telecommunication networks (Nadarajah and Cire, 2024), where equipment is regionally distributed, a fair policy guarantees equitable operations and replacements, thereby preventing frequent failures in specific areas that lead to unsatisfactory and unfair results for certain customers.

**Experimental Setup** We designed a series of experiments to test the GGF-optimality, flexibility, scalability, and efficiency of our CP-DRL algorithm. We compare against seven benchmarks, including optimal solutions (OPT) from the GGF-LP model (3) for small instances solved with Gurobi 10.0.3, the Whittle index policy (WIP) for RMABs, and a random (RDM) agent that selects actions randomly at each time step and averages the results over 10 independent runs. Additionally, we implemented a simple DRL baseline, Vanilla-DRL (V-DRL), with a utilitarian objective. The stochastic policy network employs a fully connected neural network that maps the vector  $s$  to a  $N$ -dimensional probability vector. We also implemented two heuristics to complement the random agent approach. The oldest-first (OFT) approach selects the machine in the worst state, while the myopic (MYP) selects the machine that maximizes immediate reward.

We finally implemented an equal-resources (EQR) approach based on Li and Varakantham (2022a), which imposes that each machine be replaced once every  $N$  steps to ensure an equal distribution of resources.

GGF weights decay exponentially with a factor of 2, defined as  $w_n = 1/2^n$ , and normalized to sum to 1. We use a uniform distribution  $\mu^{(N)}$  over  $\mathcal{S}^{(N)}$  and set the discount factor  $\gamma = 0.95$ . We use Monte Carlo simulations to evaluate policies over  $M$  trajectories truncated at time length  $T$ . We choose  $M = 1,000$  and  $T = 300$  across all experiments. Hyperparameters for the CP-DRL algorithm are in Appendix F.4.

**Experiment 1 (GGF-Optimality)** We obtain optimal solutions using the OPT model for instances where  $N \in \{3, 4, 5\}$ , with each machine having  $S = 3$  states. We select indexable instances to apply the WIP method for comparison. Note that, the WIP method is particularly effective in this case as it solves the equivalent utilitarian problem (as demonstrated in the utilitarian reduction result in Section 3.1). In most scenarios with small instances, WIP performs near-GGF-optimal since resources are assigned impartially, making it challenging for CP-DRL to consistently outperform WIP. As shown in Figure 2, the CP-DRL algorithm converges toward or slightly below the OPT values across the scenarios for the *Exponential-RCCC* case. WIP performs better than the random agent but does not reach the OPT values, especially as the number of machines increases. CP-DRL either outperforms or has an equivalent performance as WIP but consistently outperforms the random policy.

**Experiment 2 (Flexibility)** The fixed-size input-output design allows CP-DRL to leverage multi-task training (MT) with varying machine numbers and resources. We refer to this multi-task extension as CP-DRL(MT). We trained the CP-DRL(MT) with  $N \in \{2, 3, 4, 5\}$ , randomly switching configurations at

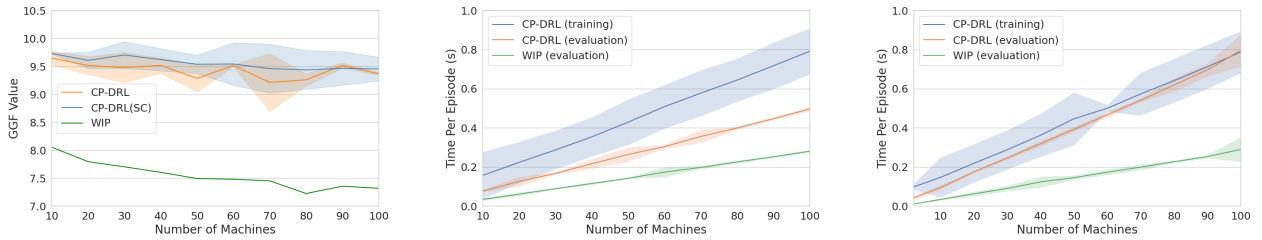


Table 1: GGF Scores (Exponential-RCCC)

$N$	OPT	WIP	CP-DRL	CP-DRL(MT)	V-DRL	OFT	MYP	EQR	RDM
2	14.19	14.07	<b>14.12</b> $\pm$ 0.01	14.11 $\pm$ 0.01	13.56 $\pm$ 0.00	5.84	12.59	10.05	9.67
3	14.08	13.75	<b>13.95</b> $\pm$ 0.02	13.89 $\pm$ 0.14	13.39 $\pm$ 0.00	7.92	12.32	11.67	10.13
4	13.94	13.27	<b>13.64</b> $\pm$ 0.05	13.59 $\pm$ 0.10	13.04 $\pm$ 0.01	9.02	12.86	12.03	9.74
5	13.77	12.47	12.96 $\pm$ 0.01	<b>13.28</b> $\pm$ 0.03	12.83 $\pm$ 0.00	10.01	12.08	11.87	8.95

Table 2: GGF Scores (Quadratic-RCCC)

$N$	OPT	WIP	CP-DRL	CP-DRL(MT)	V-DRL	OFT	MYP	EQR	RDM
2	16.17	<b>16.17</b>	16.14 $\pm$ 0.00	16.14 $\pm$ 0.00	15.36 $\pm$ 0.00	3.11	6.61	9.73	10.15
3	16.10	<b>16.09</b>	16.05 $\pm$ 0.00	16.05 $\pm$ 0.00	15.17 $\pm$ 0.01	6.16	6.63	12.73	11.83
4	16.01	<b>16.01</b>	15.94 $\pm$ 0.00	15.94 $\pm$ 0.00	15.01 $\pm$ 0.00	7.92	6.85	14.02	12.17
5	15.91	15.86	<b>15.87</b> $\pm$ 0.02	15.86 $\pm$ 0.02	14.73 $\pm$ 0.00	9.25	6.70	14.64	11.98



(a) GGF values for the number of machines  $N \in [10, 100]$  (b) Time per episode in seconds with a resource ratio  $b/N = 0.1$  (c) Time per episode in seconds with a resource ratio  $b/N = 0.5$

**Figure 3: (Colored) Scalability and Time Efficiency of CP-DRL.** Subfigures (a) and (b) show the scalability of CP-DRL with a fixed resource ratio of 0.1. Subfigure (a) presents GGF values across different machine counts, with intervals representing the standard deviation over 5 runs. Subfigure (b) and (c) depicts time per episode in seconds for a fixed resource ratio of 0.1 and 0.5, respectively. In all time plots, the green line represents WIP during MC evaluation, the blue line shows CP-DRL during training, and the orange line represents CP-DRL during MC evaluation.

the end of each episode over 2000 training episodes. CP-DRL(MT) was evaluated separately, and GGF values for WIP and RDM policies were obtained from 1000 Monte Carlo runs. The numbers following the plus-minus sign ( $\pm$ ) represent the variance across 5 experiments with different random seeds in Table 1 and 2. Variances for WIP and RDM are minimal and omitted, with bold font indicating the best GGF scores at each row excluding optimal values. As shown in Table 1, CP-DRL(MT) consistently achieves scores very close to the OPT values as the number of machines increases from 2 to 4. For the 5-machine case, CP-DRL(MT) shows slightly better performance than the single-task CP-DRL. In Table 2, the single- and multi-task CP-DRL agents show slight variations in performance across different machine numbers. For  $N = 5$ , CP-DRL achieves the best GGF score, slightly outperforming WIP. In both cases, the CP-DRL approach outperforms Vanilla-DRL, the three heuristic methods, and the random agent.

**Experiment 3 (Scalability)** We assess CP-DRL scalability by increasing the number of machines while keeping the resource proportion at 0.1 for the

*Exponential-RCCC* instances. We refer to this scaled extension as CP-DRL(SC). We vary the number of machines from 10 to 100 to evaluate CP-DRL performance as the problem size grows. We also use CP-DRL(SC), trained on 10 machines with 1 unit of resource, and scale it to tasks with 20 to 100 machines. Figure 3a shows CP-DRL and CP-DRL(SC) consistently achieve higher GGF values than WIP as machine numbers increase. CP-DRL(SC) delivers results comparable to separately trained CP-DRL, reducing training time while maintaining similar performance. Both WIP and CP-DRL show linear growth in time consumption per episode as machine numbers scale up.

**Experiment 4 (Efficiency)** In the GGF-LP model (3), the number of constraints grows exponentially with the number of machines  $N$  as  $N^2 + S^N$ , and the variables increase by  $2N + (N + 1) \cdot S^N$ . Using the count dual LP model (20) reduces the model size, but constraints still grow as  $\binom{N+S-1}{S-1}$  and variables increase by  $\binom{N+S-1}{S-1} \cdot A$ . These growth patterns create computational challenges as the problem size increases. A detailed time analysis for the GGF-LP and



count dual LP models with  $N$  from 2 to 7 is provided in Appendix F.5. In addition to the time per episode for a fixed ratio of 0.1 in Figure 3a, we analyze performance with a 0.5 ratio (Figure 3c) and varying machine proportions, keeping the number of machines fixed at 10. We evaluate CP-DRL over machine proportions from 0.1 to 0.9. The results show that the time per episode increases linearly with the number of machines, while training and evaluation times remain relatively stable. This indicates that the sampling procedure for legal actions is the primary bottleneck. Meanwhile, the resource ratio has minimal impact on computing times.

## 6 CONCLUSION

We incorporate the fairness consideration in terms of the generalized Gini function within the weakly coupled Markov decision processes, and define the GGF-WCMDP optimization problem. First, we present an exact method based on linear programming for solving it. We then derive an equivalent problem based on a utilitarian reduction when the WCMDP is symmetric, and show that the set of optimal permutation invariant policy for the utilitarian objective is also optimal for the original GGF-WCMDP problem. We further leverage this result by utilizing a count state representation and introduce a count-proportion-based deep RL approach to devise more efficient and scalable solutions. Our empirical results show that the proposed method using PPO as the RL algorithm consistently achieves high-quality GGF solutions. Moreover, the flexibility provided by the count-proportion approach offers possibilities for scaling up to more complex tasks and context where Whittle index policies are unavailable due to the violation of the indexability property by the sub-MDPs.

## Acknowledgements

Xiaohui Tu was partially funded by GERAD and IVADO. Yossiri Adulyasak was partially supported by the Canadian Natural Sciences and Engineering Research Council [Grant RGPIN-2021-03264] and by the Canada Research Chair program [CRC-2022-00087]. Nima Akbarzadeh was partially funded by GERAD, FRQNT, and IVADO. Erick Delage was partially supported by the Canadian Natural Sciences and Engineering Research Council [Grant RGPIN-2022-05261] and by the Canada Research Chair program [950-230057].

## References

Adelman, D. and Mersereau, A. J. (2008). Relaxations of weakly coupled stochastic dynamic programs. *Op-*

*erations Research*, 56(3):712–727.

- Akbarzadeh, N. and Mahajan, A. (2019). Restless bandits with controlled restarts: Indexability and computation of whittle index. In *2019 IEEE 58th conference on decision and control (CDC)*, pages 7294–7300. IEEE.
- Alandari, P. A., Klassen, T. Q., Creager, E., and McIlraith, S. A. (2023). Remembering to be fair: On non-markovian fairness in sequential decisionmaking (preliminary report). *arXiv preprint arXiv:2312.04772*.
- Argyris, N., Karsu, Ö., and Yavuz, M. (2022). Fair resource allocation: Using welfare-based dominance constraints. *European journal of operational research*, 297(2):560–578.
- Atwood, J., Srinivasan, H., Halpern, Y., and Sculley, D. (2019). Fair treatment allocations in social networks. *arXiv preprint arXiv:1911.05489*.
- Bertsimas, D., Farias, V. F., and Trichakis, N. (2012). On the efficiency-fairness trade-off. *Management Science*, 58(12):2234–2250.
- Bistriz, I., Baharav, T., Leshem, A., and Bambos, N. (2020). My fair bandit: Distributed learning of max-min fairness with multi-player bandits. In *International Conference on Machine Learning*, pages 930–940. PMLR.
- Boutilier, C. and Lu, T. (2016). Budget allocation using weakly coupled, constrained markov decision processes. In *UAI*.
- Cai, D., Lim, S. H., and Wynter, L. (2021). Efficient reinforcement learning in resource allocation problems through permutation invariant multi-task learning. In *2021 60th IEEE Conference on Decision and Control (CDC)*, pages 2270–2275.
- Caton, S. and Haas, C. (2024). Fairness in machine learning: A survey. *ACM Computing Surveys*.
- Chen, R. J., Wang, J. J., Williamson, D. F., Chen, T. Y., Lipkova, J., Lu, M. Y., Sahai, S., and Mahmood, F. (2023). Algorithmic fairness in artificial intelligence for medicine and healthcare. *Nature biomedical engineering*, 7(6):719–742.
- Chen, X. and Hooker, J. N. (2023). A guide to formulating fairness in an optimization model. *Annals of Operations Research*, 326(1):581–619.
- Cimpean, A., Jonker, C., Libin, P., and Nowé, A. (2024). A reinforcement learning framework for studying group and individual fairness. In *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems*, pages 2216–2218.
- Cousins, C., Asadi, K., and Littman, M. L. (2022). Fair e3: Efficient welfare-centric fair reinforcement

- learning. In *5th Multidisciplinary Conference on Reinforcement Learning and Decision Making. RLDM*.
- Creager, E., Madras, D., Pitassi, T., and Zemel, R. (2020). Causal Modeling for Fairness In Dynamical Systems. In *Proceedings of the 37th International Conference on Machine Learning*, pages 2185–2195. PMLR.
- D’Amour, A., Srinivasan, H., Atwood, J., Baljekar, P., Sculley, D., and Halpern, Y. (2020). Fairness is not static: Deeper understanding of long term fairness via simulation studies. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 525–534, Barcelona Spain. ACM.
- Delage, E. and Mannor, S. (2010). Percentile Optimization for Markov Decision Processes with Parameter Uncertainty. *Operations Research*.
- El Shar, I. and Jiang, D. (2024). Weakly coupled deep q-networks. *Advances in Neural Information Processing Systems*, 36.
- Elmalaki, S. (2021). Fair-iot: Fairness-aware human-in-the-loop reinforcement learning for harnessing human variability in personalized iot. In *Proceedings of the International Conference on Internet-of-Things Design and Implementation*, pages 119–132.
- Fan, Z., Peng, N., Tian, M., and Fain, B. (2022). Welfare and fairness in multi-objective reinforcement learning. *arXiv preprint arXiv:2212.01382*.
- Farnadi, G., St-Arnaud, W., Babaki, B., and Carvalho, M. (2021). Individual fairness in kidney exchange programs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 11496–11505.
- Gajane, P., Saxena, A., Tavakol, M., Fletcher, G., and Pechenizkiy, M. (2022). Survey on fair reinforcement learning: Theory and practice. *arXiv preprint arXiv:2205.10032*.
- Gast, N., Gaujal, B., and Yan, C. (2024). Reoptimization nearly solves weakly coupled markov decision processes.
- Ge, Y., Zhao, X., Yu, L., Paul, S., Hu, D., Hsieh, C.-C., and Zhang, Y. (2022). Toward pareto efficient fairness-utility trade-off in recommendation through reinforcement learning. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*, pages 316–324.
- Ghalme, G., Nair, V., Patil, V., and Zhou, Y. (2022). Long-term resource allocation fairness in average markov decision process (amdp) environment. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems*, pages 525–533.
- Guo, X., Xu, H., Zhuang, D., Zheng, Y., and Zhao, J. (2023). Fairness-enhancing vehicle rebalancing in the ride-hailing system. *arXiv preprint arXiv:2401.00093*.
- Hassanzadeh, P., Kreacic, E., Zeng, S., Xiao, Y., and Ganesh, S. (2023). Sequential fair resource allocation under a markov decision process framework. In *Proceedings of the Fourth ACM International Conference on AI in Finance*, pages 673–680.
- Hawkins, J. T. (2003). *A Lagrangian decomposition approach to weakly coupled dynamic optimization problems and its applications*. PhD thesis, Massachusetts Institute of Technology.
- Herlihy, C., Prins, A., Srinivasan, A., and Dickerson, J. P. (2023). Planning to fairly allocate: Probabilistic fairness in the restless bandit setting. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 732–740.
- Jabbari, S., Joseph, M., Kearns, M., Morgenstern, J., and Roth, A. (2017). Fairness in reinforcement learning. In *International conference on machine learning*, pages 1617–1626. PMLR.
- Jiang, J. and Lu, Z. (2019). Learning fairness in multi-agent systems. *Advances in Neural Information Processing Systems*, 32.
- Ju, P., Ghosh, A., and Shroff, N. B. (2023). Achieving fairness in multi-agent markov decision processes using reinforcement learning. *arXiv preprint arXiv:2306.00324*.
- Kozodoi, N., Jacob, J., and Lessmann, S. (2022). Fairness in credit scoring: Assessment, implementation and profit implications. *European Journal of Operational Research*, 297(3):1083–1094.
- Lan, T. and Chiang, M. (2011). An axiomatic theory of fairness in resource allocation. *George Washington University*, <http://www.seas.gwu.edu/tlan/papers/fairness.pdf>, Tech. Rep.
- Li, D. and Varakantham, P. (2022a). Efficient resource allocation with fairness constraints in restless multi-armed bandits. In *Uncertainty in Artificial Intelligence*, pages 1158–1167. PMLR.
- Li, D. and Varakantham, P. (2022b). Towards soft fairness in restless multi-armed bandits. *arXiv preprint arXiv:2207.13343*.
- Li, D. and Varakantham, P. (2023). Avoiding starvation of arms in restless multi-armed bandit. *International Foundation for Autonomous Agents and Multiagent Systems*.
- Liu, C., Chen, C.-X., and Chen, C. (2021). Meta: A city-wide taxi repositioning framework based on multi-agent reinforcement learning. *IEEE*

- Transactions on Intelligent Transportation Systems*, 23(8):13890–13895.
- Liu, L. T., Dean, S., Rolf, E., Simchowitz, M., and Hardt, M. (2018). Delayed Impact of Fair Machine Learning. In *Proceedings of the 35th International Conference on Machine Learning*, pages 3150–3158. PMLR.
- Liu, W., Liu, F., Tang, R., Liao, B., Chen, G., and Heng, P. A. (2020). Balancing between accuracy and fairness for interactive recommendation with reinforcement learning. In *Pacific-asia conference on knowledge discovery and data mining*, pages 155–167. Springer.
- Mandal, D. and Gan, J. (2022). Socially fair reinforcement learning. *arXiv preprint arXiv:2208.12584*.
- Mo, J. and Walrand, J. (2000). Fair end-to-end window-based congestion control. *IEEE/ACM Transactions on networking*, 8(5):556–567.
- Moulin, H. (1991). *Axioms of cooperative decision making*. Cambridge university press.
- Nadarajah, S. and Cire, A. A. (2024). Self-adapting network relaxations for weakly coupled markov decision processes. *Management Science*.
- Namazi, A. and Khodabakhshi, M. (2023). A novel game theoretic method on fair economic resource allocation with multiple criteria. *International Journal of Management Science and Engineering Management*, 18(3):170–176.
- Puranik, B., Madhow, U., and Pedarsani, R. (2022). Dynamic positive reinforcement for long-term fairness. In *ICML 2022 Workshop on Responsible Decision Making in Dynamic Environments*.
- Puterman, M. L. (2005). *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons.
- Rawls, J. (1971). A theory of justice. *Cambridge (Mass.)*.
- Reuel, A. and Ma, D. (2024). Fairness in reinforcement learning: A survey.
- Saure, A., Patrick, J., Tyldesley, S., and Puterman, M. L. (2012). Dynamic multi-appointment patient scheduling for radiation therapy. *European Journal of Operational Research*, 223(2):573–584.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Segal, M., George, A.-M., and Dimitrakakis, C. (2023). Policy fairness and unknown bias dynamics in sequential allocations. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*, EAAMO ’23, New York, NY, USA. Association for Computing Machinery.
- Siddique, U., Sinha, A., and Cao, Y. (2023). Fairness in preference-based reinforcement learning. *arXiv preprint arXiv:2306.09995*.
- Siddique, U., Weng, P., and Zimmer, M. (2020). Learning fair policies in multi-objective (deep) reinforcement learning with average and discounted rewards. In *International Conference on Machine Learning*, pages 8905–8915. PMLR.
- Sood, A., Jain, S., and Gujar, S. (2024). Fairness of exposure in online restless multi-armed bandits. *arXiv preprint arXiv:2402.06348*.
- van den Broek, E., Sergeeva, A., and Huysman, M. (2020). Hiring algorithms: An ethnography of fairness in practice. In *40th international conference on information systems, ICIS 2019*, pages 1–9. Association for Information Systems.
- Verma, S., Zhao, Y., Shah, S., Boehmer, N., Taneja, A., and Tambe, M. (2024). Group fairness in predict-then-optimize settings for restless bandits. In *The 40th Conference on Uncertainty in Artificial Intelligence*.
- Wen, M., Bastani, O., and Topcu, U. (2021). Algorithms for fairness in sequential decision making. In *International Conference on Artificial Intelligence and Statistics*, pages 1144–1152. PMLR.
- Weymark, J. A. (1981). Generalized gini inequality indices. *Mathematical Social Sciences*, 1(4):409–430.
- Whittle, P. (1988). Restless bandits: Activity allocation in a changing world. *Journal of applied probability*, 25(A):287–298.
- Yu, G., Siddique, U., and Weng, P. (2023). Fair deep reinforcement learning with generalized gini welfare functions. In *International Conference on Autonomous Agents and Multiagent Systems*, pages 3–29. Springer.
- Zhang, X. (2022). *Near-Optimality for Multi-action Multi-resource Restless Bandits with Many Arms*. PhD thesis, Cornell University.
- Zhang, X., Khaliligarekani, M., Tekin, C., et al. (2019). Group retention when using machine learning in sequential decision making: the interplay between user dynamics and fairness. *Advances in Neural Information Processing Systems*, 32.
- Zhang, X., Tu, R., Liu, Y., Liu, M., Kjellstrom, H., Zhang, K., and Zhang, C. (2020). How do fair decisions fare in long-term qualification? *Advances in Neural Information Processing Systems*, 33:18457–18469.

- Zhao, H. and Gordon, G. (2019). Inherent tradeoffs in learning fair representations. *Advances in neural information processing systems*, 32.
- Zimmer, M., Glanois, C., Siddique, U., and Weng, P. (2021). Learning fair policies in decentralized cooperative multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 12967–12978. PMLR.

## A RELATED WORK

Fairness-aware learning is increasingly integrated into the decision-making ecosystem to accommodate minority interests. However, naively imposing fairness constraints can actually exacerbate inequity (Wen et al., 2021) if the feedback effects of decisions are ignored. Many real-world fairness applications are not one-time static decisions (Zhao and Gordon, 2019) and can thus be better modeled with sequential decision problems, which still remain relatively understudied.

**Fairness in resource allocation** Fairness has been an important concern in resource allocation problems, where traditional approaches often build upon optimization frameworks, leveraging fairness-constrained optimization (Argyris et al., 2022; Chen and Hooker, 2023), game-theoretic concepts (Namazi and Khodabakhshi, 2023), or axiomatic principles (Lan and Chiang, 2011) to derive fair solutions. These models provide practical fair solutions, but do not account for the long-term impact of allocation decisions, which motivates fair-aware sequential decision-making frameworks based on Markov decision processes (MDPs) as discussed next.

**Fairness with dynamics** There are a few studies investigating fairness-aware sequential decision making without relying on MDPs. For instance, Liu et al. (2018) consider one-step delayed feedback effects, Creager et al. (2020) propose causal modeling of dynamical systems to address fairness, and Zhang et al. (2019) construct a user participation dynamics model where individuals respond to perceived decisions by leaving the system uniformly at random. These studies extend the fairness definition in temporally extended decision-making settings, but do not take feedback and learning into account that the system may fail to adapt to changing conditions. Alamdari et al. (2023) address this gap by introducing multi-stakeholder fairness as non-Markovian sequential decision making and developing a Q-learning based algorithm with counterfactual experiences to enhance sample-efficient fair policy learning.

**Fairness in Markov decision processes** Zhang et al. (2020) consider how algorithmic decisions impact the evolution of feature space of the underlying population modeled as MDPs but is limited to binary decisions. Ghalme et al. (2022) study a fair resource allocation problem in the average MDP setting and proposes an approximate algorithm to compute the policy with sample complexity bounds. However, their definition of fairness is restricted to the minimum visitation frequency across all states, potentially resulting in unbalanced rewards among sub-MDPs. Wen et al. (2021) develop fair decision-making policies in discounted MDPs, but the performance guarantees are achieved only under a loose condition. In contrast, our work takes into account a more comprehensive definition of fairness. Segal et al. (2023) investigate the impact of societal bias dynamics on long-term fairness and the interplay between utility and fairness under various optimization parameters. Additionally, Hassanzadeh et al. (2023) address a fair resource allocation problem similar to our work but in continuous state and action space. They define fairness to the agents considering all their allocations over the horizon under the Nash Social Welfare objective in hindsight.

**Fairness in reinforcement learning** Jabbari et al. (2017) initiate the *meritocratic fairness* notion from the multi-arm bandits setting to the reinforcement learning (RL) setting. Later, fairness consideration has been integrated in RL to achieve fair solutions in different domains, including a fair vaccine allocation policy that equalizes outcomes in the population (Atwood et al., 2019), balancing between fairness and accuracy for interactive user recommendation (Liu et al., 2020; Ge et al., 2022), and fair IoT that continuously monitors the human state and changes in the environment to adapt its behavior accordingly (Elmalaki, 2021). However, most work focuses on the impartiality aspect of fairness. Jiang and Lu (2019) investigate multi-agent RL where fairness is defined over agents and encoded with a different welfare function, but the focus is on learning decentralized policies in a distributed way. We refer readers to two literature review papers by Gajane et al. (2022) and Reuel and Ma (2024) on fairness considerations in RL, which provide comprehensive insights into current trends, challenges, and methodologies in the field.

**Fairness in restless multi-arm bandits** A line of work closely related to ours focuses on fairness in restless multi-arm bandits (RMABs). Li and Varakantham (2022a) first introduce the consideration of fairness in restless bandits by proposing an algorithm that ensures a minimum number of selections for each arm. Subsequent studies have explored similar individual fairness constraints, which aim to distribute resources equitably among arms but in a probabilistic manner. For instance, Herlihy et al. (2023) introduce a method that imposes a strictly positive lower bound on the probability of each arm being pulled at each timestep. Li and Varakantham (2022b, 2023) investigate fairness by always probabilistically favoring arms that yield higher long-term cumulative rewards. Additionally, Sood et al. (2024) propose an approach where each arm receives pulls in proportion to

its merit, which is determined by its stationary reward distribution. Our work differs from these approaches by explicitly aiming to prevent disparity and ensure a more balanced reward distribution among all arms through the generalized Gini welfare objective. The only work that considers the Gini index objective is by Verma et al. (2024), which develops a decision-focused learning pipeline to solve equitable RMABs. In contrast, our work applies to a more general setting on weakly coupled MDPs, and does not rely on the Whittle indexability of the coupled MDPs.

## B PROOFS OF SECTION 3

We start this section with some preliminary results regarding 1) the effect of replacing a policy with one that has permuted indices on the value function of a symmetric WCMDP (Section B.1); and 2) a well-known result from Puterman (2005) on the equivalency between stationary policies and occupancy measures (Section B.2). This is followed by the proof of Lemma 3.3, which helps establish our main result, Theorem 3.4 in Section B.4.

### B.1 Value Function under Permuted Policy for Symmetric WCMDP

**Lemma B.1** *If a weakly coupled WCMDP is symmetric (definition 3.1), then for any policy  $\pi$  and permutation operator  $Q$ , we have  $V_0^\pi = QV_0^{\pi^Q}$ , where the permuted policy  $\pi^Q(s, a) := \pi(Qs, Qa)$  for all  $(s, a)$  pairs.*

This lemma implies an important equivalency in symmetric weakly coupled MDPs with identical sub-MDPs. If we permute the states and actions of a policy, the permuted version of the resulting value function is equivalent to the original value function.

**Proof:** We can first show that for all  $t$ ,

$$\mathbb{P}^{\pi^Q}(s_t = s, a_t = a | s_0 = \bar{s}_0) = \mathbb{P}^\pi(s_t = Qs, a_t = Qa | s_0 = Q\bar{s}_0).$$

This can be done inductively. Starting at  $t = 0$ , we have that:

$$\begin{aligned} \mathbb{P}^{\pi^Q}(s_0 = s, a_0 = a | s_0 = \bar{s}_0) &= \pi^Q(s, a) \mathbb{I}\{s = \bar{s}_0\} \\ &= \pi(Qs, Qa) \mathbb{I}\{Qs = Q\bar{s}_0\} \\ &= \mathbb{P}^\pi(s_0 = Qs, a_0 = Qa | s_0 = Q\bar{s}_0). \end{aligned}$$

Next, assuming that  $\mathbb{P}^{\pi^Q}(s_t = s, a_t = a | s_0 = \bar{s}_0) = \mathbb{P}^\pi(s_t = Qs, a_t = Qa | s_0 = Q\bar{s}_0)$ , we can show that it is also the case for  $t + 1$ :

$$\begin{aligned} \mathbb{P}^{\pi^Q}(s_{t+1} = s', a_{t+1} = a' | s_0 = \bar{s}_0) &= \pi^Q(s', a') \sum_{s, a} \mathbb{P}(s_{t+1} = s' | s_t = s, a_t = a) \mathbb{P}^{\pi^Q}(s_t = s, a_t = a | s_0 = \bar{s}_0) \\ &= \pi(Qs', Qa') \sum_{s, a} p^{(N)}(s' | s, a) \mathbb{P}^\pi(s_t = Qs, a_t = Qa | s_0 = Q\bar{s}_0) \\ &= \pi(Qs', Qa') \sum_{s, a} p^{(N)}(Qs' | Qs, Qa) \mathbb{P}^\pi(s_t = Qs, a_t = Qa | s_0 = Q\bar{s}_0) \\ &= \mathbb{P}^\pi(s_{t+1} = Qs', a_{t+1} = Qa' | s_0 = Q\bar{s}_0), \end{aligned}$$

where we use the fact that the sub-MDPs are identical so that  $p^{(N)}(s' | s, a) = p^{(N)}(Qs' | Qs, Qa)$ .

We now have that,

$$\begin{aligned} V_0^{\pi^Q} &= \sum_{s, a} \sum_{\bar{s}_0} \mu^{(N)}(\bar{s}_0) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}^{\pi^Q}(s_t = s, a_t = a | s_0 = \bar{s}_0) r(s, a) \\ &= \sum_{s, a} \sum_{\bar{s}_0} \mu^{(N)}(\bar{s}_0) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}^\pi(s_t = Qs, a_t = Qa | s_0 = Q\bar{s}_0) r(s, a) \\ &= \sum_{s, a} \sum_{\bar{s}_0} \mu^{(N)}(Q\bar{s}_0) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}^\pi(s_t = Qs, a_t = Qa | s_0 = Q\bar{s}_0) r(s, a) \end{aligned}$$

$$\begin{aligned}
 &= \sum_{\mathbf{s}, \mathbf{a}} \sum_{\bar{\mathbf{s}}_0} \mu^{(N)}(Q\bar{\mathbf{s}}_0) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}^{\pi}(\mathbf{s}_t = Q\mathbf{s}, \mathbf{a}_t = Q\mathbf{a} | \mathbf{s}_0 = Q\bar{\mathbf{s}}_0) Q^{-1} \mathbf{r}(Q\mathbf{s}, Q\mathbf{a}) \\
 &= Q^{-1} \left( \sum_{\mathbf{s}, \mathbf{a}} \sum_{\bar{\mathbf{s}}_0} \mu^{(N)}(Q\bar{\mathbf{s}}_0) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}^{\pi}(\mathbf{s}_t = Q\mathbf{s}, \mathbf{a}_t = Q\mathbf{a} | \mathbf{s}_0 = Q\bar{\mathbf{s}}_0) \mathbf{r}(Q\mathbf{s}, Q\mathbf{a}) \right) \\
 &= Q^{-1} \left( \sum_{\mathbf{s}', \mathbf{a}'} \sum_{\bar{\mathbf{s}}'_0} \mu^{(N)}(\bar{\mathbf{s}}'_0) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}^{\pi}(\mathbf{s}_t = \mathbf{s}', \mathbf{a}_t = \mathbf{a}' | \mathbf{s}_0 = \bar{\mathbf{s}}'_0) \mathbf{r}(\mathbf{s}', \mathbf{a}') \right) = Q^{-1} \mathbf{V}_0^{\pi},
 \end{aligned}$$

where we first use the relation between  $P^{\pi^Q}$  and  $\pi$ , then exploit the permutation invariance of  $\mu^{(N)}$ . We then exploit the permutation invariance  $Q\mathbf{r}(\mathbf{s}, \mathbf{a}) = \mathbf{r}(Q\mathbf{s}, Q\mathbf{a})$ , and reindex the summations using  $\mathbf{s}' := Q\mathbf{s}$ ,  $\mathbf{a}' := Q\mathbf{a}$ , and  $\bar{\mathbf{s}}'_0 := Q\bar{\mathbf{s}}_0$ .  $\square$

## B.2 Mapping Between Stationary Policies and Occupancy Measures

We present results of Theorem 6.9.1 in Puterman (2005) to support Lemma 3.3. A detailed proof is provided in the book.

**Lemma B.2** (*Theorem 6.9.1 of Puterman (2005)*) *Let  $\Pi$  denote the set of stationary stochastic Markov policies and  $\mathcal{X}$  the set of occupancy measures. There exists a bijection  $h : \Pi \rightarrow \mathcal{X}$  such that for any policy  $\pi$ ,  $h(\pi)$  uniquely corresponds to its occupancy measure  $q_{\pi}$ . Specifically, there is a one-to-one mapping between policies and occupancy measures satisfying:*

1. For any policy  $\pi \in \Pi$ , the occupancy measure  $q_{\pi} : \mathcal{S}^{(N)} \times \mathcal{A}^{(N)} \rightarrow \mathbb{R}$  is defined as

$$q_{\pi}(\mathbf{s}, \mathbf{a}) := \sum_{\bar{\mathbf{s}}_0 \in \mathcal{S}^{(N)}} \mu^{(N)}(\bar{\mathbf{s}}_0) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}^{\pi}(\mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a} | \mathbf{s}_0 = \bar{\mathbf{s}}_0), \quad (5)$$

for all  $\mathbf{a} \in \mathcal{A}^{(N)}$  and  $\mathbf{s} \in \mathcal{S}^{(N)}$ .

2. For any occupancy measure  $q(\mathbf{s}, \mathbf{a}) : \mathcal{S}^{(N)} \times \mathcal{A}^{(N)} \rightarrow \mathbb{R}$ , the policy  $\pi_q$  is constructed as

$$\pi_q(\mathbf{s}, \mathbf{a}) := \frac{q(\mathbf{s}, \mathbf{a})}{\sum_{\mathbf{a}' \in \mathcal{A}^{(N)}} q(\mathbf{s}, \mathbf{a}')}, \quad (6)$$

for all  $\mathbf{a} \in \mathcal{A}^{(N)}$  and  $\mathbf{s} \in \mathcal{S}^{(N)}$ .

It follows that  $\pi = \pi_{q_{\pi}}$ .

Now, we show that the value function can be represented using occupancy measures.

**Lemma B.3** *For any policy  $\pi \in \Pi$ , and the occupancy measure  $q_{\pi}$  defined by (5), the expected total discounted rewards under the policy  $\pi$  can be expressed as:*

$$\mathbf{V}_0^{\pi} = \sum_{\mathbf{s} \in \mathcal{S}^{(N)}} \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} q_{\pi}(\mathbf{s}, \mathbf{a}) \mathbf{r}(\mathbf{s}, \mathbf{a}). \quad (7)$$

**Proof.** Expanding the expected total discounted rewards  $\mathbf{V}_0^{\pi}$  (as defined by equation 1), we have:

$$\mathbf{V}_0^{\pi} = \sum_{\bar{\mathbf{s}}_0 \in \mathcal{S}^{(N)}} \mu^{(N)}(\bar{\mathbf{s}}_0) \sum_{\mathbf{s} \in \mathcal{S}^{(N)}} \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} \sum_{t=0}^{\infty} \mathbb{P}^{\pi}(\mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a} | \mathbf{s}_0 = \bar{\mathbf{s}}_0) \gamma^t \mathbf{r}(\mathbf{s}, \mathbf{a}).$$

Rearranging the terms:

$$\mathbf{V}_0^{\pi} = \sum_{\mathbf{s} \in \mathcal{S}^{(N)}} \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} \left( \sum_{\bar{\mathbf{s}}_0 \in \mathcal{S}^{(N)}} \mu^{(N)}(\bar{\mathbf{s}}_0) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}^{\pi}(\mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a} | \mathbf{s}_0 = \bar{\mathbf{s}}_0) \right) \mathbf{r}(\mathbf{s}, \mathbf{a}).$$

Replacing the term in parentheses as the occupancy measure  $q_{\pi}(\mathbf{s}, \mathbf{a})$  in (5) leads directly to equation 7, which completes the proof.  $\square$



### B.3 Proof of Lemma 3.3

**Lemma 3.3 (Uniform State-Value Representation)** *If a WCMDP is symmetric (definition 3.1), then for any policy  $\pi$ , there exists a corresponding permutation invariant policy  $\bar{\pi}$  such that the vector of expected total discounted rewards for all sub-MDPs under  $\bar{\pi}$  is equal to the average of the expected total discounted rewards for each sub-MDP, i.e.,*

$$\mathbf{V}_0^{\bar{\pi}} = \frac{1}{N} \sum_{n=1}^N V_{0,n}^{\pi} \mathbf{1}.$$

**Proof by construction.** We first construct, for any fixed  $Q$ , the permuted policy  $\pi^Q(\mathbf{s}, \mathbf{a}) := \pi(Q\mathbf{s}, Q\mathbf{a})$  and characterize its occupancy measure  $q_{\pi^Q}$  as

$$q_{\pi^Q}(\mathbf{s}, \mathbf{a}) := \sum_{\bar{\mathbf{s}}_0 \in \mathcal{S}^{(N)}} \mu^{(N)}(\bar{\mathbf{s}}_0) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}^{\pi^Q}(\mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a} | \mathbf{s}_0 = \bar{\mathbf{s}}_0). \quad (8)$$

Next, we construct a new measure  $\bar{q}$  obtained by averaging all permuted occupancy measures  $q_{\pi^Q}$  for  $Q \in \mathcal{G}^N$  on all  $(\mathbf{s}, \mathbf{a})$  pairs as

$$\bar{q}(\mathbf{s}, \mathbf{a}) := \frac{1}{N!} \sum_Q q_{\pi^Q}(\mathbf{s}, \mathbf{a}). \quad (9)$$

One can confirm that  $\bar{q}$  is an occupancy measure, i.e.,  $\bar{q} \in \mathcal{X}$ , since each  $q_{\pi^Q} \in \mathcal{X}$  and  $\mathcal{X}$  is convex. Indeed, the convexity of  $\mathcal{X}$  easily follows from the fact that it contains any measure that it is the set of measures that satisfy constraints 3b and 3c.

From Lemma B.2, a stationary policy  $\bar{\pi}$  can be constructed such that its occupancy measure matches  $\bar{q}(\mathbf{s}, \mathbf{a})$ :

$$q_{\bar{\pi}}(\mathbf{s}, \mathbf{a}) := \sum_{\bar{\mathbf{s}}_0 \in \mathcal{S}^{(N)}} \mu(\bar{\mathbf{s}}_0) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}^{\bar{\pi}}(\mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a} | \mathbf{s}_0 = \bar{\mathbf{s}}_0) = \bar{q}(\mathbf{s}, \mathbf{a}), \forall \mathbf{s}, \mathbf{a}.$$

We can then derive the following steps:

$$\begin{aligned} \mathbf{V}_0^{\bar{\pi}} &= \sum_{\mathbf{s} \in \mathcal{S}^{(N)}} \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} q_{\bar{\pi}}(\mathbf{s}, \mathbf{a}) \mathbf{r}(\mathbf{s}, \mathbf{a}) \quad (\text{By Lemma B.3}) \\ &= \sum_{\mathbf{s} \in \mathcal{S}^{(N)}} \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} \bar{q}(\mathbf{s}, \mathbf{a}) \mathbf{r}(\mathbf{s}, \mathbf{a}) \quad (\text{By Lemma B.2}) \\ &= \frac{1}{N!} \sum_{Q \in \mathcal{G}^N} \sum_{\mathbf{s} \in \mathcal{S}^{(N)}} \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} q_{\pi^Q}(\mathbf{s}, \mathbf{a}) \mathbf{r}(\mathbf{s}, \mathbf{a}) \quad (\text{By construction in equation 9}) \\ &= \frac{1}{N!} \sum_{Q \in \mathcal{G}^N} \mathbf{V}_0^{\pi^Q} \quad (\text{By Lemma B.3}) \\ &= \frac{1}{N!} \sum_{Q \in \mathcal{G}^N} Q^{-1} \mathbf{V}_0^{\pi} \quad (\text{By Lemma B.1}) \\ &= \frac{1}{N!} \sum_{Q \in \mathcal{G}^N} Q^{-1} \begin{bmatrix} V_{0,1}^{\pi} \\ \vdots \\ V_{0,N}^{\pi} \end{bmatrix} \quad (\text{Vector form}) \\ &= \frac{1}{N!} \begin{bmatrix} (N-1)! \sum_{n=1}^N V_{0,n}^{\pi} \\ \vdots \\ (N-1)! \sum_{n=1}^N V_{0,n}^{\pi} \end{bmatrix} \quad (\text{Property of permutation group}) \\ &= \frac{1}{N!} (N-1)! \sum_{n=1}^N V_{0,n}^{\pi} \mathbf{1} \\ &= \frac{1}{N} \sum_{n=1}^N V_{0,n}^{\pi} \mathbf{1}. \end{aligned}$$

We complete this proof by demonstrating that  $\bar{\pi}$  is permutation invariant. Namely, for all  $Q \in \mathcal{G}^N$ , we can show that:

$$\begin{aligned}
 \bar{\pi}(Qs, Qa) &\propto \frac{1}{N!} \sum_{Q' \in \mathcal{G}^N} q_{\pi Q'}(Qs, Qa) \\
 &= \frac{1}{N!} \sum_{Q' \in \mathcal{G}^N} q_{\pi}(Q'Qs, Q'Qa) \\
 &= \frac{1}{N!} \sum_{Q'' \in \mathcal{G}^N} q_{\pi}(Q''s, Q''a) \\
 &= \frac{1}{N!} \sum_{Q'' \in \mathcal{G}^N} q_{\pi Q''}(s, a) \\
 &\propto \bar{\pi}(s, a).
 \end{aligned}$$

□

#### B.4 Proof of Theorem 3.4

We start with a simple lemma.

**Lemma B.4** *For any  $\mathbf{w}$  and any  $\mathbf{v} \in \mathbb{R}^N$ , we have that  $\text{GGF}_{\mathbf{1}/N}[\mathbf{v}] \geq \text{GGF}_{\mathbf{w}}[\mathbf{v}]$ .*

**Proof.** This simply follows from:

$$\text{GGF}_{\mathbf{1}/N}[\mathbf{v}] = \frac{1}{N} \mathbf{1}^\top \mathbf{v} = \left( \frac{1}{N!} \sum_{Q \in \mathcal{G}^N} Q\mathbf{w} \right)^\top \mathbf{v} \geq \min_{Q \in \mathcal{G}^N} (Q\mathbf{w})^\top \mathbf{v} = \text{GGF}_{\mathbf{w}}[\mathbf{v}].$$

□

**Theorem 3.4 (Utilitarian Reduction)** *For a symmetric WCMDP (definition 3.1), let  $\Pi_{\mathbf{1}/N, PI}^*$  be the set of optimal policies for the utilitarian approach (definition 2.1) that is permutation invariant, then  $\Pi_{\mathbf{1}/N, PI}^*$  is necessarily non-empty and all  $\pi_{\mathbf{1}/N, PI}^* \in \Pi_{\mathbf{1}/N, PI}^*$  satisfies*

$$\text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\pi_{\mathbf{1}/N, PI}^*}] = \max_{\pi} \text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\pi}], \forall \mathbf{w} \in \Delta(N).$$

**Proof.** Let us denote an optimal policy to the special case of the GGF-WCMDP problem (2) with equal weights as  $\pi_{\mathbf{1}/N}^*$ :

$$\pi_{\mathbf{1}/N}^* \in \arg \max_{\pi} \text{GGF}_{\mathbf{1}/N}[\mathbf{V}_0^{\pi}]. \quad (10)$$

Based on Lemma 3.3, we can construct a permutation invariant policy  $\bar{\pi}_{\mathbf{1}/N}^*$  satisfying

$$\bar{V}_0^{\pi_{\mathbf{1}/N}^*} \mathbf{1} = \mathbf{V}_0^{\bar{\pi}_{\mathbf{1}/N}^*}, \quad (11)$$

then with equation (13) and the fact that any weight vector  $\mathbf{w}$  must sum to 1, we have that

$$\text{GGF}_{\mathbf{1}/N}[\mathbf{V}_0^{\pi_{\mathbf{1}/N}^*}] = \bar{V}_0^{\pi_{\mathbf{1}/N}^*} = \text{GGF}_{\mathbf{w}} \left[ \frac{1}{N} \sum_{n=1}^N V_{0,n}^{\pi_{\mathbf{1}/N}^*} \mathbf{1} \right] = \text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\bar{\pi}_{\mathbf{1}/N}^*}], \forall \mathbf{w}.$$

Furthermore, given any  $\mathbf{w}$ , let us denote with  $\pi_{\mathbf{w}}^*$  any optimal policy to the GGF problem with  $\mathbf{w}$  weights. One can establish that:

$$\text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\pi_{\mathbf{w}}^*}] \geq \text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\bar{\pi}_{\mathbf{1}/N}^*}] = \text{GGF}_{\mathbf{1}/N}[\mathbf{V}_0^{\pi_{\mathbf{1}/N}^*}] \geq \text{GGF}_{\mathbf{1}/N}[\mathbf{V}_0^{\pi_{\mathbf{w}}^*}]. \quad (12)$$

Considering that the largest optimal value for the GGF problem is achieved when weights are equal (see Lemma B.4):

$$\text{GGF}_{\mathbf{1}/N}[\mathbf{V}_0^{\pi}] \geq \text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\pi}], \forall \pi, \forall \mathbf{w} \in \Delta(N).$$

The inequalities in (12) should therefore all reach equality:

$$\text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\pi_{\mathbf{w}}^*}] = \text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\bar{\pi}_{\mathbf{1}/N}^*}] = \text{GGF}_{\mathbf{1}/N}[\mathbf{V}_0^{\pi_{\mathbf{1}/N}^*}].$$

This implies that the bar optimal policy constructed from any optimal policy to the utilitarian approach remains optimal for any weights in the GGF optimization problem. Furthermore, it implies that there exists at least one permutation invariant policy that is optimal for the utilitarian approach.

Now, let us take any optimal permutation invariant policy  $\pi_{1/N, \text{PI}}^*$  to the utilitarian problem. The arguments above can straightforwardly be reused to get the conclusion that  $\bar{\pi}_{1/N, \text{PI}}^*$  have the same properties as the originally constructed  $\bar{\pi}$ . Namely, for all  $\mathbf{w}$ ,

$$\text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\pi_{1/N, \text{PI}}^*}] = \text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\bar{\pi}_{1/N, \text{PI}}^*}] = \text{GGF}_{1/N}[\mathbf{V}_0^{\pi_{1/N}^*}], \quad \forall \mathbf{w} \in \Delta(N).$$

Looking more closely at  $\bar{\pi}_{1/N, \text{PI}}^*$ , we observe that for any  $\mathbf{s}$  and  $\mathbf{a}$ :

$$\begin{aligned} \bar{\pi}_{1/N, \text{PI}}^*(\mathbf{s}, \mathbf{a}) &\propto \frac{1}{N!} \sum_{Q' \in \mathcal{G}^N} q_{\pi_{1/N, \text{PI}}^*, Q'}(\mathbf{s}, \mathbf{a}) \\ &= \frac{1}{N!} \sum_{Q' \in \mathcal{G}^N} q_{\pi_{1/N, \text{PI}}^*}(\mathbf{s}, \mathbf{a}) \\ &= q_{\pi_{1/N, \text{PI}}^*}(\mathbf{s}, \mathbf{a}) \\ &\propto \pi_{1/N, \text{PI}}^*(\mathbf{s}, \mathbf{a}). \end{aligned}$$

Hence, we have that  $\bar{\pi}_{1/N, \text{PI}}^* = \pi_{1/N, \text{PI}}^*$ . This thus implies that the permutation invariant  $\pi_{1/N, \text{PI}}^*$  already satisfied these properties, i.e.,

$$\text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\pi_{1/N, \text{PI}}^*}] = \text{GGF}_{\mathbf{w}}[\mathbf{V}_0^{\bar{\pi}_{1/N, \text{PI}}^*}] = \text{GGF}_{1/N}[\mathbf{V}_0^{\pi_{1/N}^*}], \quad \forall \mathbf{w} \in \Delta(N).$$

□

## B.5 Extension to Other Fairness Measures

The utilitarian reduction (Theorem 3.4) can be extended to a broader scope of fairness measures, where we replace the GGF measure in optimization problem (2) and define the  $\rho$ -WCMDP problem accordingly.

**Corollary B.4.1** *Let  $\rho : \mathcal{V} \rightarrow \mathbb{R}$ , with  $\mathcal{V} \subseteq \mathbb{R}^N$ , be a fairness measure that satisfies:*

- (Concavity): *The set  $\mathcal{V}$  is convex and  $\forall \mathbf{v}, \mathbf{w} \in \mathcal{V}$ , and  $\theta \in [0, 1]$ ,  $\rho[\theta \mathbf{v} + (1 - \theta)\mathbf{w}] \geq \theta \rho[\mathbf{v}] + (1 - \theta)\rho[\mathbf{w}]$*
- (Permutation invariance):  *$\forall \mathbf{v} \in \mathcal{V}$  and all  $Q \in \mathcal{G}^N$ , both  $Q\mathbf{v} \in \mathcal{V}$  and  $\rho[\mathbf{v}] = \rho[Q\mathbf{v}]$*
- (Constant vector invariant)  *$\forall \bar{\mathbf{v}} \in \mathbb{R} \cap \mathcal{V}$ ,  $\rho[\bar{\mathbf{v}}\mathbf{1}] = \bar{\mathbf{v}}$ .*

For a symmetric WCMDP such that  $\mathbf{V}_0^\pi \in \mathcal{V}$  for all  $\pi$ , let  $\Pi_{U, \text{PI}}^*$  be the set of optimal policies for the utilitarian approach that is permutation invariant, then  $\Pi_{U, \text{PI}}^*$  is necessarily non-empty and all  $\pi_{U, \text{PI}}^* \in \Pi_{U, \text{PI}}^*$  satisfies

$$\rho[\mathbf{V}_0^{\pi_{U, \text{PI}}^*}] = \max_{\pi} \rho[\mathbf{V}_0^\pi].$$

**Proof.** Similar to the proof of Theorem 3.4, we define the utilitarian fairness measure as

$$\rho_U[\mathbf{v}] = \frac{1}{N} \sum_{n=1}^N v_n.$$

Defining the optimal policy to the  $\rho_U$ -WCMDP problem with utilitarian objective as  $\pi_U^* \in \arg \max_{\pi} \rho_U[\mathbf{V}_0^\pi]$  and with Lemma 3.3, we can construct a permutation invariant policy  $\bar{\pi}_U^*$  satisfying

$$\bar{V}_0^{\bar{\pi}_U^*} \mathbf{1} = \mathbf{V}_0^{\pi_U^*} \mathbf{1} \in \mathcal{V}, \quad (13)$$

with  $\bar{V}_0^{\pi_U^*} := \frac{1}{N} \sum_{n=1}^N V_{0,n}^{\pi_U^*}$ . Then we have that

$$\rho_U[\mathbf{V}_0^{\pi_U^*}] = \bar{V}_0^{\pi_U^*} = \rho[\bar{V}_0^{\pi_U^*}] = \rho\left[\frac{1}{N} \sum_{n=1}^N V_{0,n}^{\pi_U^*} \mathbf{1}\right] = \rho[\mathbf{V}_0^{\bar{\pi}_U^*}],$$

where we exploit  $\rho[\bar{\mathbf{v}}\mathbf{1}] = \bar{\mathbf{v}}$  for all  $\bar{\mathbf{v}} \in \mathbb{R} \cap \mathcal{V}$ . Furthermore, let  $\pi^*$  be any optimal policy to the  $\rho$ -WCMDP

problem. One can establish that:

$$\rho[\mathbf{V}_0^{\pi^*}] \geq \rho[\mathbf{V}_0^{\bar{\pi}^*}] = \rho_U[\mathbf{V}_0^{\pi^*}] \geq \rho_U[\mathbf{V}_0^{\pi^*}]. \quad (14)$$

By Jensen's inequality and the fact that  $\rho[\cdot]$  is concave, it holds that  $\rho_U[\mathbf{V}_0^{\pi^*}] \geq \rho[\mathbf{V}_0^{\pi^*}]$  since for all  $\mathbf{v} \in \mathcal{V}$ , we have that

$$\rho[\mathbf{v}] = \frac{1}{N!} \sum_{Q \in \mathcal{G}^N} \rho[Q\mathbf{v}] \leq \rho[\frac{1}{N!} \sum_{Q \in \mathcal{G}^N} Q\mathbf{v}] = \rho[\frac{1}{N} \sum_{n=1}^N v_n \mathbf{1}] = \frac{1}{N} \sum_{n=1}^N v_n = \rho_U[\mathbf{v}],$$

where we use permutation invariance of  $\rho$ , followed with its concavity and its constant vector invariance. The inequalities in (14) should therefore all reach equality:

$$\rho[\mathbf{V}_0^{\pi^*}] = \rho[\mathbf{V}_0^{\bar{\pi}^*}] = \rho_U[\mathbf{V}_0^{\pi^*}].$$

The rest of the argument follows exactly as in the proof of Theorem 3.4 (see Appendix B.4).  $\square$

Now, we comment that the expected utility model  $\rho[\mathbf{v}] = u^{-1}\left(\frac{1}{N} \sum_{n=1}^N u(v_n)\right)$ , where  $u(\cdot)$  is a monotone and concave function, satisfies the three properties defined in Corollary B.4.1 and is a natural framework to measure fairness in resource allocation problems as discussed in Bertsimas et al. (2012). The concavity of  $u(\cdot)$  reflects a decreasing marginal utility on the allocated resource to an individual. This property promotes equitable distributions of resources by discouraging disparities in utility. A notable instance of this model is  $\alpha$ -fairness (Mo and Walrand, 2000; Ju et al., 2023), which is parameterized by  $\alpha > 0$  and takes the form

$$u_\alpha(v) := \begin{cases} \log(v) & \text{if } \alpha = 1, \\ \frac{v^{1-\alpha}}{1-\alpha} & \text{if } \alpha \neq 1. \end{cases}$$

The domain of  $\rho$  is restricted to non-negative if  $\alpha \neq 1$  and strictly positive otherwise. This function covers a range of fairness objectives, from the proportional fairness ( $\alpha = 1$ ) to the max-min fairness ( $\alpha \rightarrow \infty$ ).

## C COUNT AGGREGATION MDP

The exact form of the count aggregation MDP is obtained as follows.

**Feasible Action** The set of feasible actions in state  $\mathbf{x}$  is defined as

$$\mathcal{A}_{g_s}^{(N)}(\mathbf{x}) := \{\mathbf{u} \mid \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_k(a) u_{s,a} \leq b_k, \forall k \in \mathcal{K}; \sum_{a \in \mathcal{A}} u_{s,a} = x_s, \forall s \in \mathcal{S}\}.$$

**Reward Function** The average reward for all sub-MDPs is defined as

$$\bar{r}_\phi(\mathbf{x}, \mathbf{u}) = \frac{1}{N} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} u_{s,a} \cdot r(s, a).$$

**Transition Probability** The transition probability  $p_\phi^{(N)}(\mathbf{x}'|\mathbf{x}, \mathbf{u})$  is the probability that the number of sub-MDPs in each state passes from  $\mathbf{x}$  to  $\mathbf{x}'$  given the action counts  $\mathbf{u}$ . We define the pre-image  $f^{-1}(\mathbf{x}')$  as the set containing all elements  $\mathbf{s}' \in \mathcal{S}^{(N)}$  that map to  $\mathbf{x}'$ , then  $p_\phi^{(N)}(\mathbf{x}'|\mathbf{x}, \mathbf{u}) = \sum_{\mathbf{s}' \in f^{-1}(\mathbf{x}')} p^{(N)}(\mathbf{s}'|f^{-1}(\mathbf{x}), g_s^{-1}(\mathbf{u}))$ .

Given the equivalence of transitions within the pre-image set, for an arbitrary state-action pair  $(\mathbf{s}, \mathbf{a}) \in \phi^{-1}(\mathbf{x}, \mathbf{u})$ , the probability of transitioning from  $\mathbf{x}$  to  $\mathbf{x}'$  under action  $\mathbf{u}$  is the sum of the probabilities of all the individual transitions in the original space that correspond to this count aggregation transition. By using the transition probability in the product space, we obtain

$$p_\phi^{(N)}(\mathbf{x}'|\mathbf{x}, \mathbf{u}) = \sum_{\mathbf{s}' \in f^{-1}(\mathbf{x}')} p^{(N)}(\mathbf{s}'|\mathbf{s}, \mathbf{a}) = \sum_{\mathbf{s}' \in f^{-1}(\mathbf{x}')} \prod_{n=1}^N p_n(s'_n | s_n, a_n), \quad (15)$$

for any  $(\mathbf{s}, \mathbf{a})$  such that  $\mathbf{x} = f(\mathbf{s})$  and  $\mathbf{u} = g_s(\mathbf{a})$ .

**Initial distribution** By using a state count representation for symmetric weakly coupled MDPs, we know that  $\mathbf{1}^\top \mathbf{x} = N$ , so the cardinality of the set can be obtained through multinomial expansion that  $(s_1 + s_2 + \dots + s_N)^S = \sum_{x_1 + x_2 + \dots + x_S = N} \frac{N!}{x_1! x_2! \dots x_S!} s_1^{x_1} s_2^{x_2} \dots s_N^{x_S}$ . Intuitively, the term  $s_1^{x_1} s_2^{x_2} \dots s_N^{x_S}$  can represent distributing  $N$  identical objects (in this case, sub-MDPs) into  $S$  distinct categories (corresponding to different states). Thus, for each state count  $\mathbf{x}$ , the number of distinct ways to distribute  $N$  sub-MDPs into  $S$  states such that the counts

match  $\mathbf{x}$  is given by the multinomial coefficient  $|f^{-1}(\mathbf{x})| = \frac{N!}{x_1!x_2!\dots x_S!}$ . Given the initial distribution  $\mu^{(N)}$  is permutation invariant, the probability of starting from state  $\mathbf{x}$  in the initial distribution is

$$\mu_f^{(N)}(\mathbf{x}) = \sum_{\mathbf{s} \in f^{-1}(\mathbf{x})} \mu^{(N)}(\mathbf{s}) = |f^{-1}(\mathbf{x})| \cdot \mu^{(N)}(\bar{\mathbf{s}}) = \frac{N!}{x_1!x_2!\dots x_S!} \cdot \mu^{(N)}(\bar{\mathbf{s}}), \forall \mathbf{x}, \quad (16)$$

for any  $\bar{\mathbf{s}}$  such that  $f(\bar{\mathbf{s}}) = \mathbf{x}$ .

## D EXACT APPROACHES BASED ON LINEAR PROGRAMMING

### D.1 Optimal Solutions to the GGF-WCMDP Problem

First, we recall the dual linear programming (LP) methods to solve the MDP with discounted rewards when the transition and reward functions are known. The formulation is based on the Bellman equation for optimal policy, and is derived in Section 6.9.1 in detail by Puterman (2005).

The dual LP formulation for addressing the multi-objective joint MDP can be naturally extended to the context of vector optimization:

$$\begin{aligned} & \text{v-max} \quad \sum_{\mathbf{s} \in \mathcal{S}^{(N)}} \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} \mathbf{r}(\mathbf{s}, \mathbf{a}) q(\mathbf{s}, \mathbf{a}) \\ & \text{s.t.} \quad \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} q(\mathbf{s}, \mathbf{a}) - \gamma \sum_{\mathbf{s}' \in \mathcal{S}^{(N)}} \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} q(\mathbf{s}', \mathbf{a}) p^{(N)}(\mathbf{s} | \mathbf{s}', \mathbf{a}) = \mu^{(N)}(\mathbf{s}), \forall \mathbf{s} \in \mathcal{S}^{(N)}, \\ & \quad \quad \quad q(\mathbf{s}, \mathbf{a}) \geq 0 \quad \quad \quad \forall \mathbf{s} \in \mathcal{S}^{(N)}, \forall \mathbf{a} \in \mathcal{A}^{(N)} \end{aligned} \quad (17)$$

where any  $\mu^{(N)}(\mathbf{s}) > 0$  can be chosen, but we normalize the weights such that  $\sum_{\mathbf{s} \in \mathcal{S}^{(N)}} \mu^{(N)}(\mathbf{s}) = 1$  can be interpreted as the probability of starting in a given state  $\mathbf{s}$ .

We can now formally formulate the fair optimization problem by combining the GGF operator (Section 2.2) and the scalarizing function on the reward vector in (17):

$$\begin{aligned} & \text{max} \quad \text{GGF}_{\mathbf{w}}[\mathbf{v}] \\ & \text{s.t.} \quad \mathbf{v} = \sum_{\mathbf{s} \in \mathcal{S}^{(N)}} \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} \mathbf{r}(\mathbf{s}, \mathbf{a}) q(\mathbf{s}, \mathbf{a}) \\ & \quad \quad \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} q(\mathbf{s}, \mathbf{a}) - \gamma \sum_{\mathbf{s}' \in \mathcal{S}^{(N)}} \sum_{\mathbf{a} \in \mathcal{A}^{(N)}} q(\mathbf{s}', \mathbf{a}) p^{(N)}(\mathbf{s} | \mathbf{s}', \mathbf{a}) = \mu^{(N)}(\mathbf{s}), \forall \mathbf{s} \in \mathcal{S}^{(N)}. \\ & \quad \quad \quad q(\mathbf{s}, \mathbf{a}) \geq 0 \quad \quad \quad \forall \mathbf{s} \in \mathcal{S}^{(N)}, \forall \mathbf{a} \in \mathcal{A}^{(N)} \end{aligned} \quad (18)$$

By adding a permutation matrix  $\mathcal{Q}$  to replace the permutation applied to the index set,  $\text{GGF}_{\mathbf{w}}[\mathbf{v}]$  is equivalently represented as

$$\text{GGF}_{\mathbf{w}}[\mathbf{v}] = \inf_{\mathcal{Q}: \mathcal{Q} \geq 0, \sum_i \mathcal{Q}_{ij} = 1, \forall j, \sum_j \mathcal{Q}_{ij} = 1, \forall i} \sum_{ij} w_i \mathcal{Q}_{ij} v_j. \quad (19)$$

This reformulation relies on  $w_1 \geq w_2 \geq \dots \geq w_N$  to confirm that at the infimum we have  $\min_{\sigma} \sum_{n=1}^N w_n v_{\sigma(n)}$ . Indeed, if  $w_1$  is not assigned to the lowest element of  $\mathbf{v}$ , then one can get a lower value by transferring the assignment mass from where it is assigned to that element to improve the solution. This form is obtained through LP duality on (19):

$$\sup_{\boldsymbol{\nu}, \boldsymbol{\lambda}: \lambda_i + \nu_j \leq w_i v_j, \forall i, j} \sum_{i=1}^N \lambda_i + \sum_{j=1}^N \nu_j,$$

which leads to

$$\begin{aligned} & \text{max}_{\boldsymbol{\nu}, \boldsymbol{\lambda}, \mathbf{q}} \quad \sum_{i=1}^N \lambda_i + \sum_{j=1}^N \nu_j \\ & \text{s.t.} \quad \lambda_i + \nu_j \leq w_i v_j \quad \forall i, j = 1, \dots, N \end{aligned}.$$

Dual variable vectors are denoted by  $\boldsymbol{\lambda}$  and  $\boldsymbol{\nu}$ . Combining the constraints in (18), we can get the complete dual LP model with the GGF objective (GGF-LP) in (3).

## D.2 Solving Count Aggregation MDP by the Dual LP Model

Since the exact model for the count aggregation MDP  $\mathcal{M}_\phi$  is obtained (Appendix C), a dual LP model is formulated following Section 6.9.1 of Puterman (2005), but with count aggregation representation to solve (4):

$$\begin{aligned}
& \max \quad \sum_{\mathbf{x} \in \mathcal{S}_f^{(N)}} \sum_{\mathbf{u} \in \mathcal{A}_{g_s}^{(N)}} \bar{r}_\phi(\mathbf{x}, \mathbf{u}) q_\phi(\mathbf{x}, \mathbf{u}) \\
& \text{s.t.} \quad \sum_{\mathbf{u} \in \mathcal{A}_{g_s}^{(N)}} q_\phi(\mathbf{x}, \mathbf{u}) - \gamma \sum_{\mathbf{x}' \in \mathcal{S}_f^{(N)}} \sum_{\mathbf{u} \in \mathcal{A}_{g_s}^{(N)}} q_\phi(\mathbf{x}', \mathbf{u}) p_\phi^{(N)}(\mathbf{x}|\mathbf{x}', \mathbf{u}) = \mu_f^{(N)}(\mathbf{x}), \forall \mathbf{x} \in \mathcal{S}_f^{(N)} \quad (20) \\
& \quad q_\phi(\mathbf{x}, \mathbf{u}) \geq 0 \quad \forall \mathbf{x} \in \mathcal{S}_f^{(N)}, \forall \mathbf{u} \in \mathcal{A}_{g_s}^{(N)}
\end{aligned}$$

By choosing the initial distribution as  $\mu_f^{(N)}$ , the optimal solution  $q_\phi(\mathbf{x}, \mathbf{u}), \forall \mathbf{x}, \mathbf{u}$  is equivalent to the optimal solution to the corresponding weakly coupled MDP under the transformation  $\phi$ .

## E MODEL SIMULATOR

In the learning setting, the deep RL agent interacts with a simulated environment through a model simulator (Algorithm 2), which leads to the next state  $\mathbf{x}'$  and the average reward  $\bar{r}_\phi$  across all coupled MDPs.

---

### Algorithm 2 Simulation of Transition Dynamics

---

**Input:** count state  $\mathbf{x}$ , count action  $\mathbf{u}$ , transition probability  $p(s'|s, a)$  and reward function  $r(s, a)$ .

**Initialize:** next state  $\mathbf{x}' \leftarrow \mathbf{0}$ , average reward  $\bar{r}_\phi \leftarrow 0$

**for**  $s = 1, \dots, S$  **do**

**for**  $a = 1, \dots, A$  **do**

**while**  $u_{s,a} > 0$  **do**

            Sample the next state index  $s' \in [S]$  according to the probability distribution  $p(\cdot|s, a)$

$x'_{s'} \leftarrow x'_{s'} + 1$

$\bar{r}_\phi \leftarrow \bar{r}_\phi + \frac{1}{N} \cdot r(s, a)$

$u_{s,a} \leftarrow u_{s,a} - 1$

**end while**

**end for**

**end for**

**Return:** next state  $\mathbf{x}'$ , average reward  $\bar{r}_\phi$

---

## F EXPERIMENTAL DESIGN

### F.1 Parameter Setting

This section details the construction of the components used to generate the test instances based on Akbarzadeh and Mahajan (2019), including the cost function, the transition matrix, and the reset probability. This experiment uses a synthetic data generator implemented on our own that considers a system with  $S$  states and binary actions ( $A = 2$ ), where the two possible actions are to operate or to replace. After generating the cost matrix of the size  $\mathcal{S}^{(N)} \cdot \mathcal{A}^{(N)} \cdot N$ , we normalize the costs to the range  $[0, 1]$  by dividing each entry by the maximum cost over all state action pairs. This ensures that the discounted return always falls within the range  $[0, \frac{1}{1-\gamma}]$ .

**Cost function** The cost function  $c(s)$  for  $s \in [S]$  can be defined in five ways: 1) *Linear*:  $c(s) = s - 1$ , where the cost increases linearly with the state index; 2) *Quadratic*:  $c(s) = (s - 1)^2$  with a more severe penalty for higher states compared to the linear case; 3) *Exponential*:  $c(s) = e^{s-1}$ , which leads to exponentially increasing costs; 4) *Replacement Cost Constant Coefficient (RCCC)*:  $c(s) = 1.5(S - 1)^2$ , which is based on a constant ratio of 1.5 to the maximum quadratic cost; 5) *Random*:  $c(s)$  is randomly generated within the range  $[0, 1]$ .

**Transition function** The transition matrix for the deterioration action is constructed as follows. Once the machine reaches the  $S$ -th state, it remains in that worst state indefinitely until being successfully reset by a replacement action. For the  $s$ -th state  $s \in [S - 1]$ , the probability of remaining in the same state at the next step is given by a model parameter  $p_m \in [0, 1]$ , and the probability of transitioning to the  $(s + 1)$ -th state is  $1 - p_m$ .

**Reset probability** When a replacement occurs, there is a probability  $p_s$  that the machine successfully resets to the first state, and a corresponding probability  $1 - p_s$  of failing to be repaired and following the deterioration rule. In our experiments, we only consider a pure reset to the first state with probability 1.

## F.2 Chosen Cost Structures

We consider two cost models to reflect real-world maintenance and operation dynamics:

- (i) *Exponential-RCCC*: In this scenario, operational costs increase exponentially with age, and exceed replacement costs in the worst state to encourage replacements. This scenario fits the operational dynamics of transportation fleets, such as drone batteries, where operational inefficiencies grow rapidly and can lead to significant damage to the drones.
- (ii) *Quadratic-RCCC*: In contrast to *scenario ii*), operational costs increase quadratically with machine age, while replacement costs remain constant and always higher than operational costs. This setup is typical for high-valued machinery, where replacement costs can be significant compared to operational expenses.

## F.3 Comparison of CP-DRL Algorithms

As presented in Tables 3 and 4, the PPO algorithm consistently shows high-quality performance with low variance compared to TD3 and SAC. This motivated us to choose PPO to implement the CP-DRL algorithm.

Table 3: GGF Scores (Exponential-RCCC)

	$N=2$	$N=3$	$N=4$	$N=5$
PPO	<b>14.11</b> $\pm$ 0.01	<b>13.89</b> $\pm$ 0.14	<b>13.59</b> $\pm$ 0.10	13.28 $\pm$ 0.03
TD3	11.83 $\pm$ 2.74	12.30 $\pm$ 1.14	12.63 $\pm$ 0.23	12.62 $\pm$ 0.12
SAC	14.00 $\pm$ 0.05	13.76 $\pm$ 0.03	13.50 $\pm$ 0.05	<b>13.30</b> $\pm$ 0.02

Table 4: GGF Scores (Quadratic-RCCC)

	$N=2$	$N=3$	$N=4$	$N=5$
PPO	<b>16.14</b> $\pm$ 0.00	<b>16.05</b> $\pm$ 0.00	<b>15.94</b> $\pm$ 0.00	<b>15.86</b> $\pm$ 0.02
TD3	15.98 $\pm$ 0.03	15.75 $\pm$ 0.12	15.56 $\pm$ 0.24	15.51 $\pm$ 0.34
SAC	16.06 $\pm$ 0.01	15.75 $\pm$ 0.08	15.69 $\pm$ 0.04	15.28 $\pm$ 0.05

## F.4 Hyperparameters

In our experimental setup, we chose PPO algorithm to implement the count-proportion based architecture. The hidden layers are fully connected and the Tanh activation function is used. There are two layers, with each layer consisting of 64 units. The learning rate for the actor is set to  $5 \times 10^{-4}$  and the critic is set to  $3 \times 10^{-4}$ . The Vanilla-DRL baseline uses the same network architecture as PPO, with two hidden layers of 64 neurons each, but with a softmax output layer for its stochastic policy.

## F.5 Additional Results on LP Solving Times

The results in Tables 5 and 6 provide details on solving the GGF-LP model and the count dual LP model on the *Quadratic-RCCC* instances as the number of machines  $N$  increases from 2 to 7. The state size is set to  $S = 3$ , the action size to  $A = 2$ , and the resource to  $b = 1$ . The first block of the tables shows the number of constraints and variables. The second block provides the model solving times, with the standard deviations listed in parentheses. The results are based on 5 runs. Notice that the LP solve time includes pre-solve, the wallclock time, and post-solve times. The wallclock time is listed separately to highlight the difference from the pure LP Solve time, but not included in the total time calculation.



Table 5: Statistics for GGF-LP Model (3)

	$N = 2$	$N = 3$	$N = 4$	$N = 5$	$N = 6$	$N = 7$
# Constraints	13	36	97	268	765	2236
# Variables	31	114	413	1468	5115	17510
Data Build (s)	0.0019 (0.00)	0.0085 (0.00)	0.1076 (0.01)	1.3698 (0.05)	17.7180 (0.52)	320.0141 (16.85)
LP Build (s)	0.0028 (0.00)	0.0185 (0.01)	0.1449 (0.01)	1.4673 (0.08)	20.7464 (4.00)	392.0150 (33.97)
LP Solve (s)	0.0187 (0.02)	0.0212 (0.00)	0.1846 (0.06)	1.2914 (0.09)	13.0377 (0.63)	138.1272 (2.33)
Wallclock Solve* (s)	0.0026 (0.00)	0.0018 (0.00)	0.0115 (0.00)	0.0493 (0.00)	0.7849 (0.16)	13.3167 (0.19)
LP Extract (s)	0.0022 (0.00)	0.0019 (0.00)	0.0044 (0.00)	0.0158 (0.00)	0.0711 (0.01)	0.2828 (0.03)
<b>Total Time (s)</b>	0.0256 (0.03)	0.0500 (0.01)	0.4416 (0.06)	4.1443 (0.07)	51.5732 (3.87)	864.4391 (56.75)

\*Wall clock solve time is included in the LP Solve time.

Table 6: Statistics for Count Dual LP Model (D.2)

	$N = 2$	$N = 3$	$N = 4$	$N = 5$	$N = 6$	$N = 7$
# Constraints	6	10	15	21	28	36
# Variables	24	40	60	84	112	144
Data Build (s)	0.0035 (0.00)	0.0134 (0.00)	0.1306 (0.01)	1.4401 (0.01)	17.4080 (0.40)	204.4550 (4.49)
LP Build (s)	0.0018 (0.00)	0.0031 (0.00)	0.0056 (0.01)	0.0081 (0.00)	0.0297 (0.01)	0.0501 (0.01)
LP Solve (s)	0.0375 (0.02)	0.0362 (0.02)	0.0466 (0.00)	0.0386 (0.01)	0.0745 (0.01)	0.1634 (0.06)
Wallclock Solve* (s)	0.0034 (0.00)	0.0053 (0.00)	0.0031 (0.00)	0.0026 (0.00)	0.0034 (0.00)	0.0105 (0.00)
LP Extract (s)	0.0053 (0.00)	0.0025 (0.00)	0.0687 (0.00)	0.0026 (0.00)	0.0051 (0.00)	0.0209 (0.00)
<b>Total Time (s)</b>	0.0481 (0.02)	0.0552 (0.02)	0.2515 (0.13)	1.4894 (0.01)	17.5173 (0.42)	204.6893 (4.44)

\*Wall clock solve time is included in the LP Solve time.