
SKILL-ALIGNED FAIRNESS IN MULTI-AGENT LEARNING FOR COLLABORATION IN HEALTHCARE *

Promise Osaine Ekpo, Brian La, Thomas Wiener, Saesha Agarwal,
 Arshia Agrawal, Gonzalo Gonzalez-Pumariaga, Angelique Taylor
 Cornell Tech, New York, NY, USA
 {poe6, byl8, tfw29, sa2388, aa2443, gg387, amt298}@cornell.edu

Lekan P. Molu
 Microsoft Research NYC
 patlekano@gmail.com

ABSTRACT

Fairness in multi-agent reinforcement learning (MARL) is often framed as a workload balance problem, overlooking agent expertise and the structured coordination required in real-world domains. In healthcare, equitable task allocation requires workload balance or expertise alignment to prevent burnout and overuse of highly skilled agents. Workload balance refers to distributing an approximately equal number of subtasks or equalised effort across healthcare workers, regardless of their expertise. We make two contributions to address this problem. First, we propose FairSkillMARL, a framework that defines fairness as the dual objective of workload balance and skill-task alignment. Second, we introduce MARLHospital, a customizable healthcare-inspired environment for modeling team compositions and energy-constrained scheduling impacts on fairness, as no existing simulators are well-suited for this problem. We conducted experiments to compare FairSkillMARL in conjunction with four standard MARL methods, and against two state-of-the-art fairness metrics. Our results suggest that fairness based solely on equal workload might lead to task-skill mismatches and highlight the need for more robust metrics that capture skill-task misalignment. Our work provides tools and a foundation for studying fairness in heterogeneous multi-agent systems where aligning effort with expertise is critical.

Keywords Multiagent reinforcement learning · Fairness · Healthcare simulator

1 Introduction

In multi-agent systems, agents must learn to cooperate while leveraging resources in their environment to achieve individual and collective objectives [1]. In this paper, we are interested in effective modeling methods in safety-critical healthcare workers (HCW) environments, such as emergency departments (ED). A distinguishing feature of these environments that necessitates new tools to model behaviour is the need for workers (henceforth referred to as agents) to perform overlapping tasks, and occasionally relieve other workers based on availability, energy levels, and skills. Common medical procedures in the ED are composed of subtasks that must be carried out sequentially to achieve a goal. This includes *individual tasks* such as the use of automated external defibrillator and administering medication, and *shared tasks* such as Cardiopulmonary resuscitation (CPR); CPR requires multiple agents to perform task-switching once they become fatigued from performing chest compressions. In real-world environments, many HCWs are overworked (i.e., burnout) when assigned tasks that do not align with their expertise (i.e., skill-task misalignment). When this happens, HCWs spend more time accomplishing tasks compared to a skilled HCW with relevant expertise who can perform the task more efficiently [2, 3, 4]. To prevent burnouts and to ensure fairness in EDs, in this work, we propose a framework for fairness in a multi-agent reinforcement learning (MARL) setting.

Prior work has quantified fairness as reward equality in social dilemmas [5], worst-case performance guarantees in traffic scheduling [6], balanced throughput in network control [7, 8], equitable effort distribution in cooperative navigation

**Citation:* Promise Osaine Ekpo, Brian La, Thomas Wiener, Saesha Agarwal, Arshia Agrawal, Gonzalo Gonzalez-Pumariaga, Lekan P. Molu, Angelique Taylor. Skill-Aligned Fairness in Multi-Agent Learning for Collaboration in Healthcare. Pages.... DOI:000000/11111.

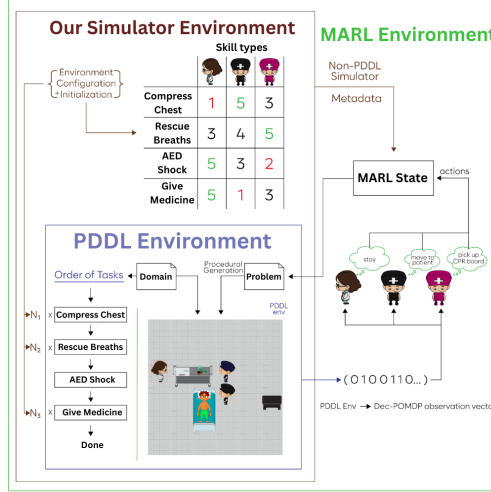


Figure 1: **MARLHospital**: Breakdown of the Planning Domain Definition Language (PDDL) environment wrapped by the MARL environment with 3 HCWs and a patient.

tasks [9], cumulative reward parity in hierarchical learning (such as Fair-Efficient Networks, [10] and demographic parity in resource allocation [11]. However, these definitions do not jointly account for the skill-task alignment and overutilization of agents who are repeatedly tasked with the most demanding actions. Furthermore, existing simulators lack realistic healthcare-themed environments that model agent energy levels, skill-task alignment, and workload in structured team configurations. More broadly, the intersection of MARL, healthcare, and fairness is understudied.

As a step towards a fairness framework for MARL for varying-skilled heterogeneous agents, such as in the ED, we require a multi-agent simulation environment that captures the challenges faced by HCWs during medical procedures. Specifically, various classes of MARL algorithms, including Independent learning (IL) [12, 13] and Centralized Training with Decentralized Execution (CTDE) algorithms [14, 15], have been benchmarked on cooperative applications where multiple agents must work together to achieve a goal. Other approaches, including Roles to DEcompose (RODE) [16], Role-oriented multi-agent reinforcement learning framework (ROMA) [17], learn dynamic subtask assignment LDSA [18], and role-transfer CTDE [19], explore role-switching. Simulators provide a controlled environment for modeling and analyzing cooperative tasks among multiple agents, thereby facilitating detailed optimization and evaluation of team dynamics and coordination [20, 21, 22, 23]. However, existing simulators lack realistic healthcare-themed environments that model agent energy levels, skill-task alignment, and workload.

To address these research gaps, we introduce MARLHospital, a simulation environment that captures HCWs’ interactions during medical procedures. MARLHospital is beneficial because it models structured teams with varying expertise, in tasks inspired by real-world medical procedures, and uniquely includes the execution of order-dependent and shared-task mode where agents alternate actions due to energy constraints, modelling agent fatigue across long time horizons, unlike prior work [20, 21, 22, 23]. Furthermore, we introduce FairSkillMARL, a framework that formulates fairness as a balance of workload and skill-task alignment during collaborative tasks, whilst factoring in the progressively dissipative energy levels of agents. This framing has the potential to improve the efficiency of healthcare teams by leveraging task-switching and agent expertise.

Our **contributions** are threefold: 1) We introduce a customizable hospital-themed game environment inspired by real-world settings. 2) We share insights about the performance of four standard MARL algorithms, including on-policy and off-policy methods, across varying team compositions, task difficulties, and energy levels in three healthcare benchmark tasks. 3) We introduce a FairSkillMARL framework, which redefines fairness as a composite disparity of agent skill-task alignment and workload, and our results demonstrate that FairSkillMARL shows competitive performance alongside state-of-the-art fairness metrics, particularly for simulated tasks and the need for more robust metrics to capture skill misalignment.

2 Related Work

In healthcare, fairness extends beyond efficiency to patient safety, requiring mechanisms that prevent HCWs fatigue from unfair workload and treatment delays caused by assigning tasks to workers outside their specialties. In social

dilemmas, inequity aversion reduces envy and guilt by penalizing outcome disparities, promoting cooperation [5]. Prior work has quantified fairness across multi-agent domains. In networked systems, fairness is often framed through worst-case guarantees, namely, maximizing the 5%-tile user data rate [6]. In traffic signal control, fairness is integrated into value functions to ensure low-traffic lanes are scheduled equally [7, 8]. These approaches focus on reward parity or throughput smoothing but assume homogeneous agents. Structural fairness has been explored in cooperative control by penalizing travel deviations to balance workload [9] and propose switching between fair and efficient subpolicies [10]. Social fairness constraints have also been encoded, such as demographic parity in Proximal Policy Optimization (PPO) [11]. While addressing various fairness goals, these works often overlook skill-agent compatibility and task-agent mismatch. In contrast, we define fairness for role differentiated, safety-critical domains along two axes: 1) workload balance and 2) skill-task alignment. This prevents overburdening skilled agents or skill-task misalignment in critical tasks, an underexplored challenge in decentralized MARL for healthcare settings.

To understand interactions between multiple agents, the majority of research efforts have captured popular settings such as games [24, 20] and cooking environments [25]. Despite their successes, many approaches did not take into account the need for fully customizable skill levels across agents or the evaluation of predefined team compositions inspired by human-human collaboration in safety-critical environments. Although simulators such as Overcooked-AI [20] and Robotouille [25] model tasks with temporal dependencies, they assume symmetric agent expertise and lack support for evaluating different team compositions. SMACv2 [23] supports heterogeneous agents but does not support structured task hierarchies and the shared task mode in MARLHospital. While Pommerman [26], CUISINEWORLD [27], and Melting Pot [28] focus on abstract coordination, they do not capture team heterogeneity or domain-specific expertise. VirtualHome [29] supports long-horizon tasks but lacks real-time multi-agent coordination, and AgentHospital [30] focuses on dialogue using Large Language Model rather than MARL benchmarking. In contrast, MARLHospital evaluates coordination in time-critical settings with configurable team compositions and shared task settings. A distinguishing feature of our work is we draw inspiration from modelling real-world safety-critical environments (i.e., the ED) where HCWs have varying skills, energy levels and workloads. Our work bridges gaps in prior work by considering agent skill and energy levels.

3 MARLHospital

We present MARLHospital, a framework for modeling multi-agent collaborative tasks in medical environments. We address the challenge of coordinating agents with varying expertise levels and in shared task mode, which is prevalent in multi-agent settings, particularly in ERs.

3.1 Notations and Preliminaries

We work in a decentralized Partially Observable Markov Decision Process (Dec-POMDP) [31], whereupon POMDPs in multi-agent settings under partial observations are formalized. Essentially, a Dec-POMDP is the tuple $(I, S, A, T, R, \Omega, O, \gamma)$, where $I = \{1, 2, \dots, n\}$ is a finite set of agents, S is a finite set of states, and $A = A_1 \times A_2 \times \dots \times A_n$ represents the joint action space, with A_i being the set of actions available to agent i . In our problem setting, each action $\mathbf{a}_t^i \in A_i$ is a vector representation of all the actions available to agent i at time t . The state at time t is $s_t \in S$, represented as 2-D coordinates x, y ; the state includes skill level information describing agent capabilities, and ongoing or recent actions $\mathbf{a}_t \in A$ (e.g., 'treatment' for ongoing treatment). The transition function $T : S \times A \times S \rightarrow [0, 1]$ specifies the probability $T(s_{t+1}|s_t, \mathbf{a}_t)$ of transitioning to state s_{t+1} given the joint actions \mathbf{a}_t taken in state s_t . Observations for each agent are drawn from $\Omega = \Omega_1 \times \Omega_2 \times \dots \times \Omega_n$, where Ω_i represents the observations available to agent i in the Planning Domain Definition Language (PDDL) environment. The observation function $O : S \times A \times \Omega \rightarrow [0, 1]$ defines the probability $O(\mathbf{o}_t|s_{t+1}, \mathbf{a}_t)$ of observing joint observation \mathbf{o}_t given joint action \mathbf{a}_t was taken and resulted in state s_{t+1} . Let $\mathbf{a}_t \in A$ denote the joint action at time step t . Discrete actions include movement, item manipulation (e.g., pick, place, stack, unstack), treatment (e.g., compress_chest and give_rescue_breaths) (see Supplemental Material for more details). The reward function $R : S \times A \rightarrow \mathbb{R}$ assigns rewards based on subgoal progress and task completion. The reward attainable in the simulator environment, based on the subgoal progress heuristic function H , is defined as $R(s_t, \mathbf{a}_t)$. For any given state s_t and joint action \mathbf{a}_t , the reward at timestep t is computed as: $R(s_t, \mathbf{a}_t) = (H(s_t) - H(s_{t-1}))$. Hence, the agents are rewarded to make progress towards task completion. A discount factor $\gamma \in [0, 1)$ balances immediate and future rewards.

The MARLHospital environment builds upon Robotouille [25] using the EPyMARL [32] setup for MARL algorithms.

3.2 Agent Energy Levels for Shared Tasks

Cardiopulmonary resuscitation (CPR) in the real world requires that HCWs alternate every 2 minutes [33]. In this sentiment, a feature of our simulator is the design of **shared tasks** with implementation that mandates alternation

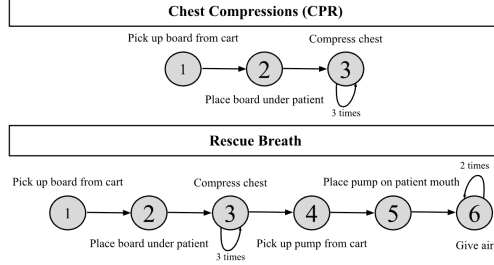


Figure 2: Task flow diagrams for the check compression (CPR) and rescue breath tasks in MARLHospital.

between agents. Specifically, we implemented this feature in the simulator only within the chest compression task and introduced a new energy cost constraint into the environment. We define three values, which are constant across the agents: 1) an energy cost incurred by agents for performing chest compressions, 2) a rate at which agents regain energy when not performing a chest compression, and 3) a maximum agent energy level. In more formal terms, at timestep t agent i has energy $e_{i,t}$, where $0 < e_{i,t} < e_{max}$. We use a constant energy cost to act c_a . In this setting, we define a “noop” action that models the neutral agent behavior when it is not equipped to execute an action. When not performing an energy-costly action, agents’ energy recharge at a constant rate; this is also described as c_a for a non-costly action a , but would have a negative value. For chest compressions, we would formally write the task $R(s_t, \mathbf{a}_t)$ for agent i at time t it as such:

$$\mathcal{R}_i(s_t) = \begin{cases} \text{compressions} & \text{if } c_{\text{compressions}} \leq e_{i,t} \\ \text{rest} & \text{else} \end{cases}$$

Consequently, a high cost-to-recharge ratio requires agents to alternate roles during CPR; otherwise, they take significantly more time to complete chest compressions.

3.3 Medical Tasks

We model the CPR and AED tasks using the standard procedural steps of “Adult Basic Life Support” from the American Red Cross code cards (see Figure 2). These cards provide visual flow charts guiding resuscitation procedures to capture real-world tasks based on evidence-based practices [?]. These procedural steps are modeled in the action space of MARLHospital.

4 FairSkillMARL

We modify the reward function from the preliminaries section to penalise both disparity in workload balance and skill-task misalignment in individual and shared tasks. Unlike prior definitions that primarily models fairness as the distribution of agent workloads (author?) [9], throughput parity [11]), we model the reality that task completion time and error rates depend on agent skill levels. A skilled agent completes complex tasks quickly with low error probability, while an unskilled agent requires more time and makes more mistakes on the same task. Our reward formulation prevents skilled agents from being overloaded with all complex tasks (causing fatigue) while also preventing unskilled agents from becoming bottlenecks on tasks beyond their capabilities. This captures the essential coordination challenge in ED, where performance depends on matching tasks to appropriate skill levels.

FairSkillMARL Objective Definition

The FAIRSkillMARL framework is a dual-pronged design that allows us to quantify fairness as both an agent-centric (burnout risk) and efficiency-centric (task success) objective, reflecting the unique operational demands of healthcare settings. Let N denote the number of agents and M the total number of subtasks completed in an episode. For each agent i , let $|\tau_i|$ be the count of subtasks assigned to agent i , and $\mathcal{E}_i(k) \in [0, 1]$ the skill level of agent i on subtask k . We define fairness as a linear combination of workload imbalance and skill-task misalignment **Workload imbalance** is measured using the **Gini Index**, L_1 as defined in (1), where $x_i = |\tau_i|$ denotes the number of subtasks assigned to agent i , n is the number of agents, and \bar{x} is the mean workload. The Gini index captures average pairwise workload disparity normalized by the mean, with $L_1 = 0$ indicating perfect balance and higher values denoting greater inequality in task allocation. This term is the proportion of subtasks completed relative to that agent’s skill levels in each subgoal; in other words, it measures the deviation from optimal task-to-skill assignments. A lower L_2 indicates better alignment of agent skills and the task they perform.

$$L_1 = \frac{\sum_{i=1}^n \sum_{j=1}^n |x_i - x_j|}{2n^2 \bar{x}} \quad (1)$$

$$L_2 = 1 - \frac{\sum_{t \in \mathcal{T}} \mathcal{S}_{i_t}(t)}{\sum_{t \in \mathcal{T}} \max_j \mathcal{S}_j(t)} \quad (2)$$

We combine these components into a *composite disparity* in Equ. 3, where $\alpha \in [0, 1]$ is a tunable parameter specifying the trade-off between workload balance and expertise alignment.

$$L_3 = \alpha L_1 + (1 - \alpha) L_2 \quad (3)$$

The shaped reward used to guide policy learning r_t is defined in Equ. 4, where $R(s_t, \mathbf{a}_t)$ is the original team reward at timestep t , and $\lambda > 0$ is a scaling factor controlling the strength of the fairness penalty.

$$r_t = R(s_t, \mathbf{a}_t) - \lambda L_3 \quad (4)$$

5 Experiment 1: MARLHospital

Our experiments were designed to answer the following questions: (1) How does varying task difficulty and different team compositions affect the performance of existing MARL algorithms and team efficiency? (2) How well do standard MARL algorithms perform in the shared task mode with varying energy levels, and how does this performance vary across team compositions?

5.1 Baseline Methods

Our experiments include standard MARL algorithms, including IL and CTDE, trainable on our custom simulation environment as building blocks for more sophisticated algorithms (i.e., EPyMARL library [32]² and PyMARL [22]³) [34, 35, 36, 37]. We investigate structured team compositions, including standard baselines that might facilitate different coordination strategies based on their decentralized and centralized training setups.

We summarize the IL and CTDE algorithms below.

- **IQL** [12] enables each IL agent to simultaneously learn its Q-function within the environment based on its local observation history and actions.
- **MAPPO** [38] is a CTDE policy gradient algorithm. Based on the PPO framework, thus enabling multiple updates on the same training batch to improve sample efficiency and stability, but it also uses a centralised state-value critic function conditioned on the state of the global environment.
- **VDN** [14] decomposes the CTDE team’s joint Q-value function into a sum of individual Q-value functions. The joint Q-value is trained using the Deep Q-Network (DQN) algorithm [39], and thus gradients are backpropagated to the individual agent networks.
- **QMIX** [15] extends CTDE VDN by using a mixing network to combine individual agent Q-values non-linearly, allowing for more complex value function factorization. The mixing network is constrained to maintain monotonicity in the relationship between the agent-specific and global Q-values, ensuring the optimal local actions and corresponding global joint actions are the same.

5.2 Training and Testing Procedure

We utilize the MARLHospital benchmark to evaluate the performance of the MARL algorithms with 3 agents to assess the success rate in achieving the goal. Experiments ran for 4 seeds with a maximum of 50 timesteps per episode. We perform task-specific hyperparameter tuning across all four algorithms (see Supplemental Materials). We train the baseline algorithms with four seeds for 2M timesteps for all off-policy algorithms and 20M timesteps for the on-policy algorithms, similarly done in EPyMARL [32]. Off-policy algorithms use an experience replay buffer for stabilized

²<https://github.com/uoel-agents/epymarl>

³<https://github.com/oxwhirl/pymarl>

Table 1: Success rates of MARL algorithms for Chest Compressions (P) at energy levels 0 and 3, and Rescue Breaths (C) at energy 0, with uniform (U), specialized (S), and forced cooperation (F) team compositions.

Method	Chest Compressions			Rescue Breaths			Chest Compressions		
	P-U	P-S	P-F	C-U	C-S	C-F	P-U	P-S	P-F
IQL [12]	0.64	0.21	0.51	0.46	0.00	0.32	0.62	0.80	0.58
MAPPO [38]	0.15	0.02	0.02	0.01	0.01	0.00	0.70	0.97	0.57
VDN [14]	0.85	0.79	0.70	0.80	0.74	0.61	0.84	0.78	0.85
QMIX [15]	0.84	0.68	0.54	0.78	0.60	0.40	0.79	0.60	0.63

policy learning. In terms of computational requirements, off-policy experiments were executed on 4 CPU cores on a compute system with two Intel Xeon Gold 6448Y processors (each with 32 cores, 2.10–4.10GHz, 60MB cache, PCIe 5.0), providing a total of 64 CPU cores. For all algorithms, we used 1 GPU on NVIDIA GPUs with x86 CPUs when available. We use **success rate** as a performance metric to assess how often agents reach the goal. This metric allows us to understand the stability and final performance of the trained policy across all evaluation episodes. Concretely, we report the total number of goals reached during test time over 4 seeds, with each evaluation consisting of 100 episodes.

5.3 Experimental Setup

We conducted two experiments to assess the performance of IL and CTDE algorithms. To understand the impact of short and long-term horizon tasks, we measure the **success rate** of CPR and rescue breaths. We then assess the effect of different team compositions on both tasks. Next, we proposed an energy function to model routine task-switching during CPR in the shared task mode, as inspired by the real-world situations. In doing so, we evaluate the impact of the energy level on the success rate of goal completion for CPR and rescue breath tasks with both skill levels in the observation space, with (energy cost of 3) and without (energy cost of 0) the energy level on the success rate of goal completion.

Task Difficulty We define the two goals as Partial (P) and Complete (C) with task difficulties based on the length of the time horizon. For the CPR goal, the HCWs must perform CPR, a short-time-horizon task that consists of picking up and placing a board under the patient and giving N chest compressions. For the rescue breaths, which is a longer horizon task, in addition to the CPR goal, the HCWs must perform the same actions in the CPR task, in addition to picking up the pump and placing it on the patient to give them air. Since we use procedural generation in the environment implementation, these goals can be modified in the environment configuration file.

Team Compositions Understanding team composition is essential for modeling clinical settings, as the structure and capabilities of a care team influence task performance, decision speed, and workload distribution. We define three team compositions that vary in expertise levels (EL) of agents: *uniform*, *specialized*, and *interdependent* teams. In *uniform Teams*, all agents possess identical capabilities and can perform every subtask with equal efficiency. This structure allows maximum flexibility as any agent can complete any part of the task independently, and coordination is optional. *Specialized Teams* allows agents to be more efficient in one particular subtask but still retains the ability to perform all others, although more slowly. This specific setting incentivizes but does not necessitate collaboration. In *Interdependent Teams (Forced Cooperation)*, agents are capable of performing only a subset of subtasks and cannot execute at least two others.

5.4 MARL Hospital Results

5.4.1 Impact of Task Difficulty on Performance.

The results in Table 1 with test success rates show that the CTDE algorithm (VDN) outperforms the DTDE algorithms in the CPR goal by a significant gap as seen in Figure 3. Additionally, in the rescue breaths goal, the CTDE (VDN, QMIX) algorithms outperform the DTDE algorithms. This indicates that CTDE approaches are more robust under increased task complexity. Also, we performed pairwise Welch’s t -tests on success rates for the chest compression task. **VDN significantly outperformed MAPPO** ($p = 0.0052$)

, highlighting the advantage of CTDE. Other comparisons, such as MAPPO vs. QMIX showed no significant difference.

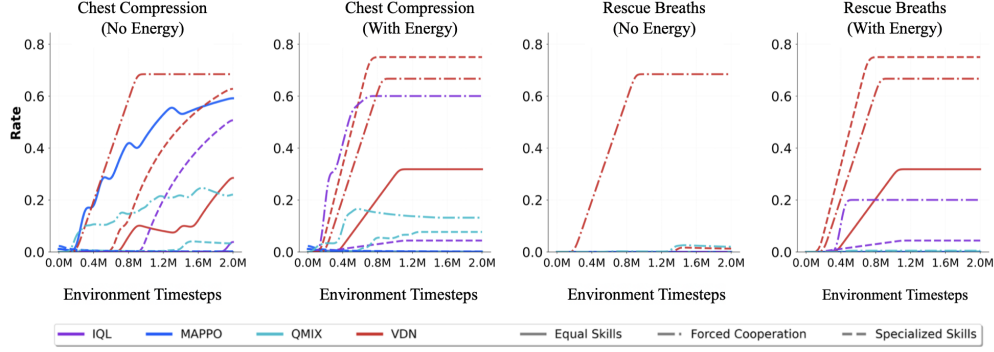


Figure 3: Test cumulative rates for the “chest compression” and “rescue breaths” tasks, with and without energy, showing the median.

5.4.2 Impact of Team Composition on Performance.

Using the same test success rates metric, the results show that CTDE specifically, VDN consistently outperformed DTDE methods across team compositions, indicating strong generalizability to different team dynamics (see Table 1). VDN has the highest success rate in the uniform skilled teams in both tasks. However, the forced cooperation teams had the lowest success rates across all algorithms, reflecting the coordination burden imposed by asymmetric skill levels of agents. This performance drop highlights the challenge of learning robust policies when agents are unable to act interchangeably.

5.4.3 Impact of Energy Function on Performance.

For the CPR goal, according to Table 1, the shared task mode (with an energy cost of “3” and recharge rate of “1”) impacts the MARL algorithm’s performance compared to the baseline scenario (energy cost of “0”). Under this setting, coordination becomes more challenging, as agents must alternate actions and manage limited energy budgets across a long-horizon task. Despite this, the CTDE algorithm VDN shows the best performance in both the uniform (PUE = 0.84) and forced cooperation (PFE = 0.85) teams, outperforming all DTDE baselines. This suggests that centralized training may help agents better learn coordination strategies under resource constraints. In contrast, DTDE methods such as IQL struggle perhaps due to a lack of shared state. MAPPO, despite struggling in the no energy team configurations achieves its best performance in the specialized setting with energy (PSE = 0.97), indicating that role specialization might mitigate coordination difficulty in this shared task setting. Surprisingly across both tasks as seen in Figure 3, agents tend to converge faster in the energy-constrained setting than in the no-energy baseline. This suggests that the structured turn-taking enforced by energy levels of agents in the shared task mode may simplify coordination rather than hinder performance, perhaps leading to more straightforward credit assignment and role specialization during learning.

6 Experiment 2: FairSkillMARL

Setup: We conduct three controlled experiments across 1M timesteps and 4 seeds to isolate the effects of different fairness components. **First**, to isolate the impact of L_2 , we compared skill-task alignment ($\alpha = 0.7$) against workload only fairness ($\alpha = 0$) across three fairness weights ($\lambda \in \{0.0, 1.0, 4.0\}$) using specialised teams due to their varying skill levels. **Second**, to verify that the improved fairness is a result of skill-task alignment, not the total team skill, we compare equal and specialised teams at fixed $\lambda = 1.0$ with $\alpha \in \{0.0, 0.7\}$. **Third**, we compare all on-policy and off-policy MARL algorithms from Experiment 1 to evaluate fair task distributions in terms of workload and skill-alignment using the efficiency-based reward function, $R(s_t, \mathbf{a}_t)$.

Baseline comparison: Using the success rate metric, we compare FairSkillMARL to two state-of-the-art fairness methods across 1M timesteps and 4 seeds to investigate how well our framework performs in terms of skill-task alignment, workload, and both. 1) Gini Index is similar to FairSkillMARL with $\alpha = 1$; thus, it measures agents’ contribution to sub-tasks in the reward function [40]. 2) Fair Efficient Network (FEN) measures agent resource utilization and penalizes agents when they deviate from the average utilization of all agents [41].

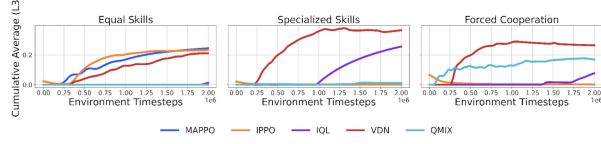


Figure 4: L_3 composite disparity over training timesteps for different algorithms and skill configurations with no energy. Lower values indicate fairer agent behavior.

6.1 FairSkillMARL Results

Interpreting the results, we consider the Gini index (workload imbalance), the percentage of contribution of the 3 agents (A %0, A1 %, A2 %), the success rate, range and the values of λ and α . The range is defined as the maximum contribution of an agent minus the minimum contribution of an agent.

6.1.1 Comparison across team compositions.

In equal teams under $\lambda = 1.0$, as seen in Table 3, FEN achieves a 0.194 success rate but poor utilisation for Agent 2 (4.33%, range = 46.07), indicating very unfair specialisation. In contrast, FairSkillMARL ($\alpha = 0.7$) achieves higher success (0.25) with a more balanced workload (Agent 2 = 30.92%, range = 31.45), confirming that the inclusion of skill levels impacts teams with equally skilled workers.

Table 2: Experiment 1: Impact of Skill-Task Alignment - Specialized Teams

Method	α	Success \uparrow	Gini \downarrow	Agent % (A0/A1/A2)	Range \downarrow
$\lambda = 0.0$ (No fairness incentive)					
Workload	0.0	0.11 \pm 0.10	0.27	33.25/36.61/30.14	6.47
FairSkill	0.7	0.10 \pm 0.06	0.32	39.64/39.66/20.69	18.97
FEN	—	0.09 \pm 0.05	0.30	28.49/38.17/33.34	9.68
$\lambda = 1.0$ (Moderate fairness)					
Workload	0.0	0.10 \pm 0.07	0.34	44.17/40.48/15.35	28.82
FairSkill	0.7	0.07 \pm 0.04	0.40	49.94/34.75/15.31	34.63
FEN	—	0.16 \pm 0.18	0.35	40.98/36.84/22.18	18.80
$\lambda = 4.0$ (Strong fairness)					
Workload	0.0	0.08 \pm 0.07	0.37	46.35/32.65/21.00	25.35
FairSkill	0.7	0.11 \pm 0.08	0.35	37.44/48.74/17.82	30.92
FEN	—	0.43 \pm 0.44	0.25	31.48/41.86/26.65	15.21

Table 3: Experiment 2: Team Composition Effects (Equal vs Specialized Teams, $\lambda = 1.0$)

Method	α	Success \uparrow	Gini \downarrow	Agent % (A0/A1/A2)	Range \downarrow
Specialized Teams					
Workload	0.0	0.10 \pm 0.07	0.34	44.17/40.48/15.35	28.82
FairSkill	0.7	0.07 \pm 0.04	0.40	49.94/34.75/15.31	34.63
FEN	—	0.16 \pm 0.18	0.35	40.98/36.84/22.18	18.80
Equal Skill Teams					
Workload	0.0	0.12 \pm 0.06	0.40	33.63/39.44/26.93	12.51
FairSkill	0.7	0.25 \pm 0.18	0.34	29.59/39.48/30.92	9.89
FEN	—	0.19 \pm 0.00	0.34	45.27/50.40/4.33	46.07

6.1.2 Comparison to FEN Baseline.

Under strong fairness enforcement ($\lambda = 4.0$), FEN achieves both higher success (0.429) and better workload balance (range = 15.21) than FairSkillMARL (0.113 success, range = 30.92). With no consideration of agent skill levels, FEN

distributes tasks more evenly (A0–A2: 31.5%–41.9%–26.7%). Meanwhile, FairSkillMARL over-relies on Agent 1 (48.7%) implying that in high-penalty settings, composite fairness objectives (workload + skill) can conflict as well as reducing performance and workload balance. Simpler shaping, as in FEN, may yield more stable coordination under strong constraints as seen in Table 2. In addition, we compare FairSkillMARL ($\alpha = 0.7$) to workload-only fairness ($\alpha = 0.0$) using paired two-tailed t -tests and Cohen’s d . At $\lambda = 0.0$, the difference is not significant ($t(3) = 1.21$, $p = 0.31$, $d = 0.69$). At $\lambda = 1.0$, FairSkillMARL shows a significant improvement ($t(3) = 3.42$, $p = 0.02$, $d = 1.91$), indicating that moderate fairness promotes effective coordination. At $\lambda = 4.0$, performance declines under skill alignment ($t(3) = -2.87$, $p = 0.041$, $d = -1.63$), suggesting that excessive fairness constraints may hinder learning.

6.1.3 Comparison Across Algorithms

We investigated the emergent fairness of standard MARL algorithms without fairness reward shaping. IQL and VDN exhibit severe agent imbalance, while MAPPO achieves balanced coordination. Under energy constraints, MAPPO remains robust, and value-based methods (IQL, VDN, QMIX) continue to show imbalance as seen in Figure 4.

7 Experiment 3: FairSkillMARL Ablation

We aim to investigate the impact of FairSkillMARL policies on success rates as we vary the α and λ hyperparameters to allow practitioners to control the learning process. When $\alpha = 1$, the objective prioritizes workload fairness exclusively, strongly penalizing uneven task distribution. On the other hand, $\alpha = 0$ prioritizes skill alignment, ensuring tasks are allocated to the most proficient agents regardless of workload balance. The penalty scaling factor λ modulates how much fairness impacts the shaped reward relative to efficiency. A small λ softly encourages fairness without significantly altering baseline policies, while a large λ can impose stricter fairness constraints at the cost of potential reductions in overall task performance. In our experiments, we set λ to 1 and vary α . We measure the success rate of agents in the rescue breathes tasks across 1M timesteps, and 4 seeds (see Table 4).

7.1 FairSkillMARL Ablation Results

Table 4: Experiment 3: FairSkillMARL Ablation - Optimal α Selection ($\lambda = 1.0$).

Method	α	Success \uparrow	Gini \downarrow	Agent % (A0/A1/A2)	Range \downarrow
Baseline Methods					
Workload	0.0	0.06 \pm 0.04	0.42	46.08/36.75/17.17	28.91
FEN	—	0.15 \pm 0.10	0.34	42.26/41.82/15.92	26.34
FairSkillMARL with Varying α					
FairSkill	0.5	0.09 \pm 0.02	0.39	45.10/36.63/18.27	26.83
FairSkill	0.7	0.13 \pm 0.06	0.33	44.73/35.28/19.99	24.74
FairSkill	1.0	0.05 \pm 0.02	0.32	47.10/39.39/13.71	33.39

The optimal skill-task alignment weight is $\alpha = 0.7$, which balances performance and fairness (36.7%). Pure workload balancing ($\alpha = 1.0$) achieves a marginally lower success rate with higher fairness disparity.

8 Conclusion

In this work, we introduced skill-task alignment as a fairness criterion for heterogeneous agents in MARL. We also introduced MARLHospital, a customizable simulation environment for benchmarking cooperative MARL algorithms on varying skill levels and energy levels across task difficulty levels. We evaluated 5 standard MARL algorithms across two tasks with three agents. We evaluated the effectiveness of IL and CTDE algorithms across different skill levels and energy cost constraints, finding that CTDE algorithms achieved the best performance in both the CPR task and the rescue breaths tasks. Our experiments with FairSkillMARL in MARLHospital are informative, achieving statistical improvements at $\alpha=0.7$ compared to pure workload balancing (0.13 vs 0.06 success rate, $p=0.02$), capturing skill misalignment in complex coordination tasks might require more robust metrics. By modeling agent capabilities explicitly, our method provides a foundation for fairness definitions that consider both workload and agent skills in heterogeneous teams. Future research can build upon this foundation to explore multi-objective optimization techniques, and applications to larger-scale heterogeneous multi-agent systems.

References

- [1] Dom Huh and Prasant Mohapatra. Multi-agent Reinforcement Learning: A Comprehensive Survey, July 2024.
- [2] Angelique Taylor, Tauhid Tanjim, Michael Joseph Sack, Maia Hirsch, Kexin Cheng, Kevin Ching, Jonathan St George, Thijs Roumen, Malte F. Jung, and Hee Rin Lee. Rapidly Built Medical Crash Cart! Lessons Learned and Impacts on High-Stakes Team Collaboration in the Emergency Room, February 2025. arXiv:2502.18688 [cs] version: 1.
- [3] Angelique Taylor, Tauhid Tanjim, Huajie Cao, and Hee Rin Lee. Towards Collaborative Crash Cart Robots that Support Clinical Teamwork. In *Proceedings of the 2024 ACM/IEEE International Conference on Human-Robot Interaction, HRI '24*, pages 715–724, New York, NY, USA, March 2024. Association for Computing Machinery.
- [4] Angelique Taylor, Hee Rin Lee, Alyssa Kubota, and Laurel D. Riek. Coordinating Clinical Teams: Using Robots to Empower Nurses to Stop the Line. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–30, November 2019.
- [5] Edward Hughes, Joel Z. Leibo, Matthew G. Phillips, Karl Tuyls, Edgar A. Duéñez-Guzmán, Antonio García Castañeda, Iain Dunning, Tina Zhu, Kevin R. McKee, Raphael Koster, Heather Roff, and Thore Graepel. Inequity aversion improves cooperation in intertemporal social dilemmas, September 2018.
- [6] Mingqi Yuan, Qi Cao, Man-On Pun, and Yi Chen. Multi-agent reinforcement learning-based fairness-aware scheduling for bursty traffic. In *2021 IEEE Global Communications Conference (GLOBECOM)*, pages 1–6. IEEE, 2021.
- [7] Wanqing Fang, Xintian Zhao, and Chengwei Zhang. Fairness-aware multi-agent reinforcement learning and visual perception for adaptive traffic signal control. *Optoelectronics Letters*, 20(12):764–768, December 2024.
- [8] Xingshuai Huang, Di Wu, and Benoit Boulet. Fairness-Aware Model-Based Multi-Agent Reinforcement Learning for Traffic Signal Control. September 2022.
- [9] Jasmine Jerry Aloor, Siddharth Nayak, Sydney Dolan, and Hamsa Balakrishnan. Cooperation and Fairness in Multi-Agent Reinforcement Learning, October 2024.
- [10] Jiechuan Jiang and Zongqing Lu. Learning Fairness in Multi-Agent Systems, October 2019. arXiv:1910.14472 [cs].
- [11] Gabriele Malfa, Jie Zhang, Michael Luck, and Elizabeth Black. *Fairness Aware Reinforcement Learning via Proximal Policy Optimization*. February 2025.
- [12] Ardi Tampuu, Tanel Matiisen, Dorian Kodelja, Ilya Kuzovkin, Kristjan Korjus, Juhan Aru, Jaan Aru, and Raul Vicente. Multiagent Cooperation and Competition with Deep Reinforcement Learning, November 2015.
- [13] Christian Schroeder de Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makovychuk, Philip H. S. Torr, Mingfei Sun, and Shimon Whiteson. Is Independent Learning All You Need in the StarCraft Multi-Agent Challenge?, November 2020.
- [14] Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel. Value-Decomposition Networks For Cooperative Multi-Agent Learning, June 2017.
- [15] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder de Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. QMIX: Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning, June 2018.

- [16] Tonghan Wang, Tarun Gupta, Anuj Mahajan, Bei Peng, Shimon Whiteson, and Chongjie Zhang. RODE: Learning Roles to Decompose Multi-Agent Tasks, October 2020.
- [17] Tonghan Wang, Heng Dong, Victor Lesser, and Chongjie Zhang. ROMA: Multi-Agent Reinforcement Learning with Emergent Roles, July 2020.
- [18] Mingyu Yang, Jian Zhao, Xunhan Hu, Wengang Zhou, Jiangcheng Zhu, and Houqiang Li. LDSA: Learning Dynamic Subtask Assignment in Cooperative Multi-Agent Reinforcement Learning, November 2022.
- [19] Dung Nguyen, Phuoc Nguyen, Svetha Venkatesh, and Truyen Tran. Learning to Transfer Role Assignment Across Team Sizes, April 2022.
- [20] Micah Carroll, Rohin Shah, Mark K. Ho, Thomas L. Griffiths, Sanjit A. Seshia, Pieter Abbeel, and Anca Dragan. On the Utility of Learning about Humans for Human-AI Coordination, January 2020.
- [21] Huaxiaoyue Wang, Gonzalo Gonzalez-Pumariega, Yash Sharma, and Sanjiban Choudhury. Demo2Code: From Summarizing Demonstrations to Synthesizing Code via Extended Chain-of-Thought, November 2023.
- [22] Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas Nardelli, Tim G. J. Rudner, Chia-Man Hung, Philip H. S. Torr, Jakob Foerster, and Shimon Whiteson. The StarCraft Multi-Agent Challenge, December 2019. arXiv:1902.04043 [cs].
- [23] Benjamin Ellis, Jonathan Cook, Skander Moalla, Mikayel Samvelyan, Mingfei Sun, Anuj Mahajan, Jakob N. Foerster, and Shimon Whiteson. SMACv2: An Improved Benchmark for Cooperative Multi-Agent Reinforcement Learning, October 2023.
- [24] William H. Guss, Brandon Houghton, Nicholay Topin, Phillip Wang, Cayden Codell, Manuela Veloso, and Ruslan Salakhutdinov. MineRL: A Large-Scale Dataset of Minecraft Demonstrations, July 2019.
- [25] Gonzalo Gonzalez-Pumariega, Leong Su Yean, Neha Sunkara, and Sanjiban Choudhury. Robotouille: An Asynchronous Planning Benchmark for LLM Agents, February 2025. arXiv:2502.05227 [cs].
- [26] Cinjon Resnick, Wes Eldridge, David Ha, Denny Britz, Jakob Foerster, Julian Togelius, Kyunghyun Cho, and Joan Bruna. Pommerman: A Multi-Agent Playground, September 2018. Publication Title: arXiv.org.
- [27] Ran Gong, Qiuyuan Huang, Xiaojian Ma, Hoi Vo, Zane Durante, Yusuke Noda, Zilong Zheng, Song-Chun Zhu, Demetri Terzopoulos, Li Fei-Fei, and Jianfeng Gao. MindAgent: Emergent Gaming Interaction, September 2023.
- [28] John P. Agapiou, Alexander Sasha Vezhnevets, Edgar A. Duéñez-Guzmán, Jayd Matyas, Yiran Mao, Peter Sunehag, Raphael Köster, Udari Madhushani, Kavya Kopparapu, Ramona Comanescu, D. J. Strouse, Michael B. Johanson, Sukhdeep Singh, Julia Haas, Igor Mordatch, Dean Mobbs, and Joel Z. Leibo. Melting Pot 2.0, October 2023.
- [29] Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Torralba. VirtualHome: Simulating Household Activities via Programs, June 2018.
- [30] Junkai Li, Siyu Wang, Meng Zhang, Weitao Li, Yunghwei Lai, Xinhui Kang, Weizhi Ma, and Yang Liu. Agent Hospital: A Simulacrum of Hospital with Evolvable Medical Agents, May 2024.
- [31] Frans Oliehoek and Christopher Amato. A Concise Introduction to Decentralized POMDPs. January 2016.
- [32] Georgios Papoudakis, Filippos Christianos, Lukas Schäfer, and Stefano V. Albrecht. Benchmarking Multi-Agent Deep Reinforcement Learning Algorithms in Cooperative Tasks, November 2021.
- [33] Sorravit Savatmongkornkul, Chaipaporn Yuksen, Sumalin Chumkot, Pongsakorn Atiksawedparit, Chetsadakon Jenpanitpong, Sorawich Watcharakitpaisan, Parama Kaninworapan, and Konwachira Maijan. Comparison of chest compression quality between 2-minute... : International Journal of Critical Illness and Injury Science.
- [34] Jakob Foerster, Ioannis Alexandros Assael, Nando de Freitas, and Shimon Whiteson. Learning to Communicate with Deep Multi-Agent Reinforcement Learning. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [35] He He, Jordan Boyd-Graber, Kevin Kwok, and Hal Daumé III. Opponent Modeling in Deep Reinforcement Learning, September 2016.
- [36] Yali Du, Lei Han, Meng Fang, Ji Liu, Tianhong Dai, and Dacheng Tao. LIIR: Learning Individual Intrinsic Reward in Multi-Agent Reinforcement Learning. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [37] Christopher Amato. A First Introduction to Cooperative Multi-Agent Reinforcement Learning, December 2024.
- [38] Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The Surprising Effectiveness of PPO in Cooperative, Multi-Agent Games, November 2022.

- [39] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing Atari with Deep Reinforcement Learning, December 2013.
- [40] Robert Busa-Fekete, Balazs Szorenyi, Paul Weng, and Shie Mannor. Multi-objective Bandits: Optimizing the Generalized Gini Index, June 2017. arXiv:1706.04933 [cs].
- [41] Jiechuan Jiang and Zongqing Lu. Learning Fairness in Multi-Agent Systems. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [42] Paul Maria Scheikl, Balazs Gyenes, Tornike Davitashvili, Rayan Younis, Andre Schulze, Beat P. Muller-Stich, Gerhard Neumann, Martin Wagner, and Franziska Mathis-Ullrich. Cooperative Assistance in Robotic Surgery through Multi-Agent Reinforcement Learning. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1859–1864, Prague, Czech Republic, September 2021. IEEE.
- [43] Xinzhi Zhang, Yeming Yang, Qingling Zhu, Qiuzhen Lin, Weineng Chen, Jianqiang Li, and Carlos A. Coello Coello. Multi-agent deep Q-network-based metaheuristic algorithm for Nurse Rostering Problem. *Swarm and Evolutionary Computation*, 87:101547, June 2024.
- [44] Qian Yue Hao, Fengli Xu, Lin Chen, Pan Hui, and Yong Li. Hierarchical Multi-agent Model for Reinforced Medical Resource Allocation with Imperfect Information. *ACM Transactions on Intelligent Systems and Technology*, 14(1):1–27, February 2023.
- [45] Dario Esposito, Davide Schaumann, Domenico Camarda, and Yehuda E. Kalay. Multi-Agent Modelling and Simulation of Hospital Acquired Infection Propagation Dynamics by Contact Transmission in Hospital Wards. In Yves Demazeau, Tom Holvoet, Juan M. Corchado, and Stefania Costantini, editors, *Advances in Practical Applications of Agents, Multi-Agent Systems, and Trustworthiness. The PAAMS Collection*, pages 118–133, Cham, 2020. Springer International Publishing.
- [46] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, February 2015.

A MARLHospital Environment

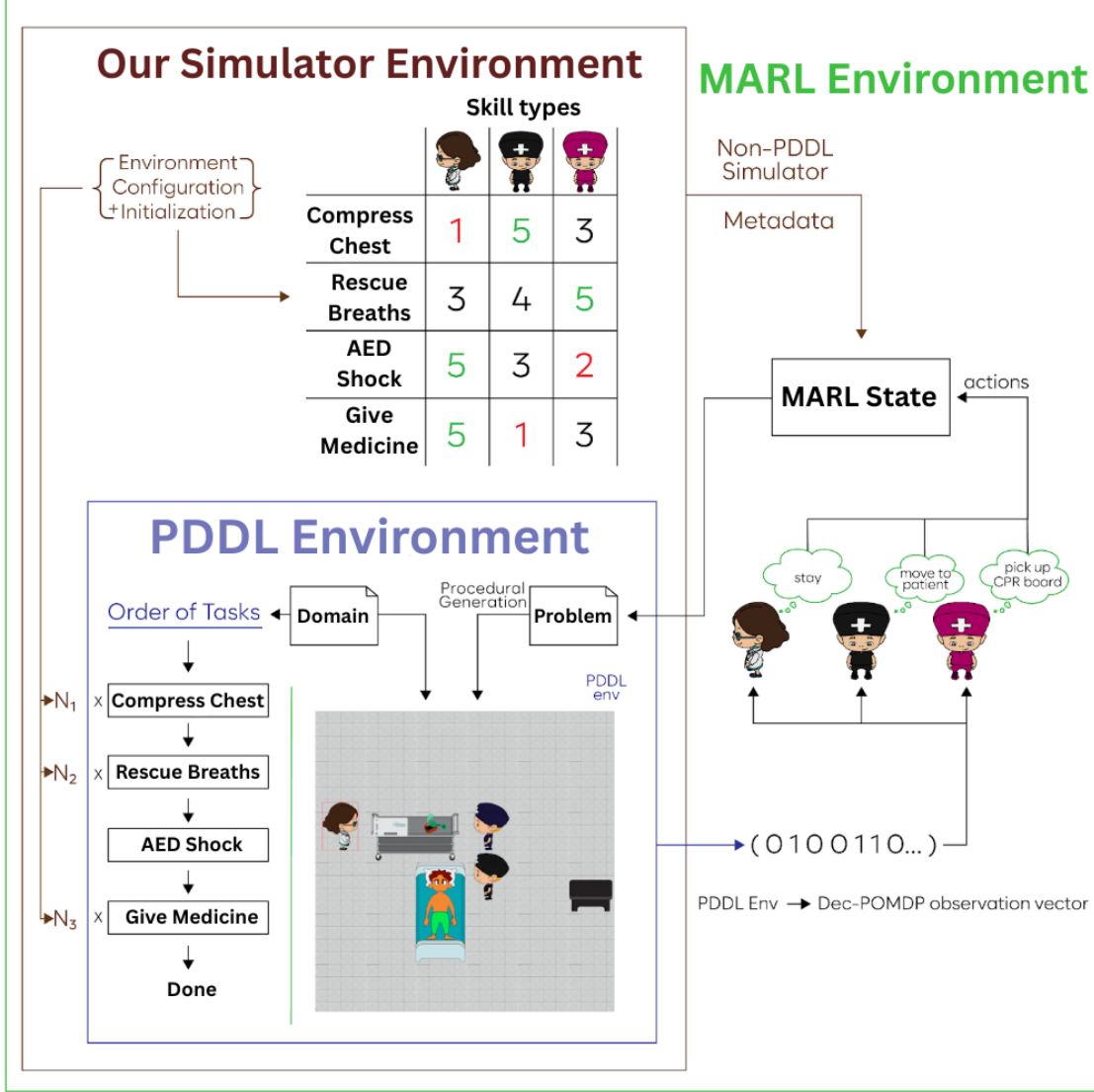


Figure 5: MARLHospital environment: Breakdown of the PDDL environment wrapped by the MARL environment with 3 HCW agents and a patient.

Figure 5 shows how the Planning Domain Definition Language (PDDL) and MARL environment are integrated within MARLHospital in the baseline configuration. A configuration file initializes 1) the agent, station, patient, and object positions in the PDDL environment, 2) the skill (and possibly initial energy) levels of the agents, and 3) the subtask requirements to reach the goal. In the default setting, agents receive a binary encoding of the PDDL environment as an observation; when skill information is included in the observation, that data is passed from the configuration as well. The MARL state converts agents' actions into pddl actions based on state factors such as skill level, energy, and subtask requirements marked by N_1 , N_2 , and N_3 in the diagram.

The remainder of this section provides an overview of the different components of the PDDL environment.

Players consist of the healthcare workers in the environment. Controlled by the RL agents, they take actions in the PDDL environment, and are allowed to move to locations on the grid adjacent to stations. See Figure 6 for player options.

Actions consist of the actions healthcare workers can perform within the environment. These include the following:



Figure 6: The PDDL players, controlled by RL agents

- **Move:** move the worker from station A to station B
- **Pick-up:** pick up an item at the worker's current station. Afterwards, the station is empty and the worker now holds the item.
- **Place:** place the item the worker is holding on the worker's current station. Before placing an item, the station must be empty. Afterwards player no longer holds an item.
- **Unstack:** unstack the item that is on top of another item at the worker's current station. The unstacked item is now held by the worker.
- **Stack:** stack an item on top of another item at the worker's current station. The worker no longer holds an item.
- **Stack under:** stack an item underneath another item at the worker's current station.
- **Chest compression:** perform a medical procedure in which the worker manually circulates the blood of the patient. Before compressing chest, the patient must have a CPR board underneath them.
- **Rescue breaths:** perform a medical procedure in which the worker directly provides oxygen to the patient. Before giving rescue breaths, the patient must have a breathing pump attached.
- **Give shock:** perform a medical procedure in which the worker uses a defibrillator to restart the patient's heart. The patient must have an AED device connected to them.
- **Give medicine:** perform a medical procedure in which the worker administers medicine to the patient through a syringe. The syringe must be on the patient.

moving between stations, picking up or placing items, stacking or unstacking objects (including stacking under), and performing medical procedures such as chest compressions and rescue breaths.

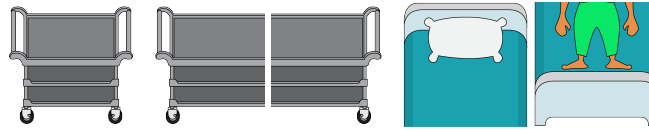


Figure 7: PDDL stations.

Stations are a class of immovable objects in the environment. Players can move to unoccupied cells adjacent to stations, as mentioned earlier, and items can be placed on stations. See available station assets in Figure 7.

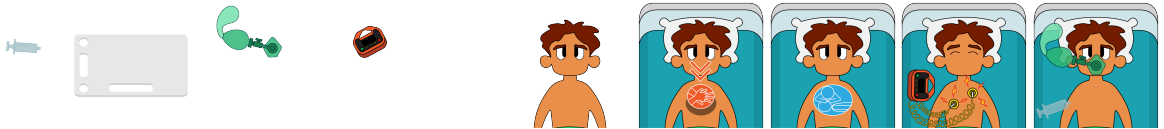


Figure 8: Item assets. *Left:* Items used for treatment. From left to right: syringe, CPR board, breathing pump, AED device. *Right:* Patient in various states. From left to right: default, after chest compressions, after rescue breaths, after AED shock, after given medicine

Items are interactable objects in the environment. Apart from the patient, the items in MARLHospital are tools necessary for the HCWs to complete tasks in treating the patient. They can be picked up and placed down by Players on stations. Items can be stacked on top of one another, though the player can only pick up the item at the top of a stack. To perform treatment actions, certain items need to be stacked on the patient based on that treatment action: a CPR board for `compress_chest`, a breathing pump for `give_rescue_breaths`, an AED device for `giving_shock`, and a syringe for `giving_medicine`.

There are two items with special properties in the environment. The first is the CPR board, which, unlike other items, can also be stacked under the patient. The second is the patient itself. Unlike other items, the patient cannot be moved

from the patient bed station, and it is the only possible predicate of treatment actions. The patient is rendered with different states based on the most recent task completed on them. See Figure 8 for item and patient assets.

Observation Space Size. The observation space consists of a total of 174 boolean values. Each of the three agents receives an individual observation vector of size 58. These 58 boolean features are structured as follows:

- **Patient State (4 booleans):** Indicates whether the patient has been chest compressed, rescue breathed, treated, or shocked.
- **Agent Locations (18 booleans):** Encodes the location of each of the three agents, where each agent can be at one of six discrete stations: `hospital_cart_right1`, `table1`, `hospital_cart1`, `hospital_cart_left1`, `patient_legs1`, or `patient_bed_station1`.
- **Held Items (9 booleans):** Indicates which item (if any) is held by each of the three agents. Possible items include: `pump1`, `cpr_board1`, and `patient1`.
- **Skill Levels (18 booleans):** Represents each agent’s skill level (unskilled, beginner, expert) for two skills: chest compression and rescue breathing.
- **Available Actions (9 booleans):** Encodes which actions are currently available to the agent, such as treating the patient, moving to any of the six stations, moving an item, or stacking an item under another.

B MARLHospital Environment

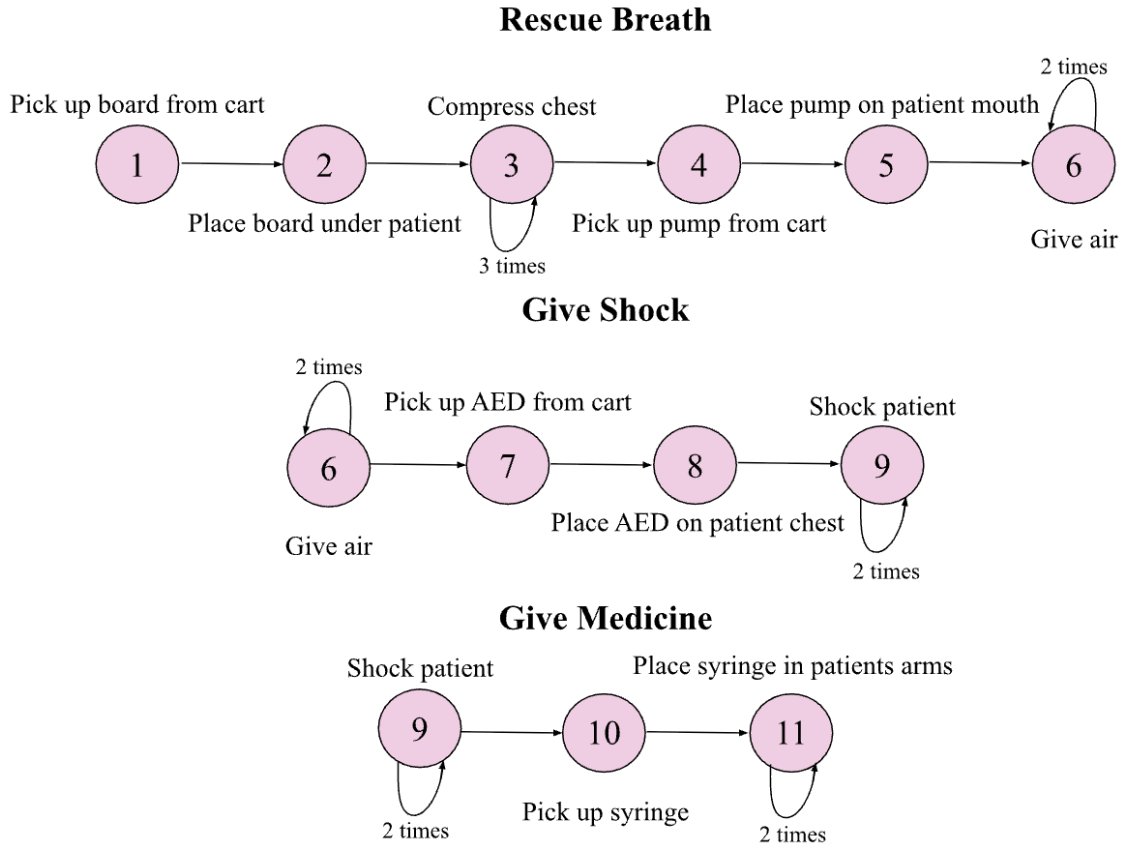


Figure 9: MARLHospital environment: Breakdown of the PDDL environment wrapped by the MARL environment with 3 HCW agents and a patient.

Table 5: Comparison between MARLHospital and other simulators

	Multi-Agent	MARL	Hospital Setting	Skill Specialization	Variable Skills	Energy Levels	Fairness Support
Overcooked-AI [20]	✓	✗	✗	✗	✗	✗	✗
Agent Hospital [30]	✓	✗	✓	✓ ¹	✗	✗	✗
MARL Robotic Surgery [42]	✓	✓	✓ ²	✓	✗	✗	✗
MARL for Nurse Rostering [43]	✓	✓	✓ ³	✗	✗	✗	✗
Cuisine World [27]	✓	✗	✗	✗	✗	✗	✗
VirtualHome [29]	✓	✗	✗	✗	✗	✗	✗
SMACv2 [23]	✓	✓	✗	✓	✗	✗	✗
Pommerman [26]	✓	✓	✗	✗	✗	✗	✗
Melting Pot 2.0 [28]	✓	✓	✗	✓	✗	✗	✗
Robotouille [25]	✓	✗	✗	✗	✗	✓	✗
HMARL for Medical Allocation [44]	✓	✓	✗	✗	✗	✗	✗
MA Hospital Infection Sim [45]	✓	✗	✓	✓	✗	✗	✗
<i>MARLHospital</i>	✓	✓	✓	✓	✓	✓	✓

C Selected Hyperparameters

For hyperparameter tuning in MARLHospital, we primarily conducted a grid search and did a sweep across several hyperparameters for all the algorithms, ensuring consistency with established MARL benchmarks. The hyperparameter selection process follows a structured approach that is compatible with prior algorithms’ work.

The hyperparameters for each algorithm, including IQL, MAPPO, VDN, and QMIX, are outlined in Tables 6 to 9. Key parameters include the hidden dimension size, learning rate, reward standardization, network type (e.g., GRU or fully connected networks), target update frequency, entropy coefficient (for policy gradient methods), and the number of steps used in bootstrapping (n-step returns).

For off-policy algorithms such as IQL, VDN, and QMIX, an experience replay buffer is employed to decorrelate samples and stabilize learning, following standard practices [46]. On-policy algorithms like MAPPO utilises parallel synchronous workers to mitigate sample correlation and improve stability during training.

The hyperparameter configurations for the IQL algorithm, with parameter sharing for MARLHospital is shown in Table 6.

Table 6: Hyperparameters for IQL with parameter sharing.

	MARLHospital
hidden dimension	64
learning rate	0.0005
reward standardisation	False
network type	GRU
evaluation epsilon	0.0
target update	20 (hard)

“The ‘target update 0.01 (soft)’ in our MAPPO implementations is based on Polyak averaging, where the target critic network is updated slowly using 1 ‘%’ new and 99 old weights for stability. Although we did not add it to the hyperparameter table, the PPO clipping coefficient is set to 0.2 and is used to constrain policy updates during training.

The hyperparameter configurations for the MAPPO algorithm, with parameter sharing for MARLHospital is shown in Table 7.

The hyperparameter configurations for the VDN algorithm, with parameter sharing for MARLHospital is shown in Table 8.

The hyperparameter configurations for the QMIX algorithm, with parameter sharing for MARLHospital, are shown in Table 9.

Table 7: Hyperparameters for MAPPO with parameter sharing.

MARLHospital	
hidden dimension	64
learning rate	0.002
reward standardisation	True
network type	GRU
entropy coefficient	0.01
target update	0.05 (soft)
clipping coef	0.2

Table 8: Hyperparameters for VDN with parameter sharing.

MARLHospital	
hidden dimension	64
learning rate	0.001
reward standardisation	False
network type	GRU
evaluation epsilon	0.0
target update	25 (hard)

Table 9: Hyperparameters for QMIX with parameter sharing.

MARLHospital	
hidden dimension	64
learning rate	0.001
reward standardisation	False
network type	GRU
evaluation epsilon	0.1
target update	25 (hard)

```

"config": {
  "num_compressions": {
    "patient": 3,
    "default": 3
  },
  "num_breaths": {
    "patient": 2,
    "default": 2
  },
  "num_shocks": {
    "patient": 2,
    "default": 2
  },
  "num_medicine_doses": {
    "patient": 2,
    "default": 2
  },
  "energy_levels": {
    "compresschest_cost": 0,
    "recharge_rate": 100,
    "max": 1000
  },
  "num_players": 3
}

```

Figure 10: JSON to configure the necessary steps to treat the patient

```

{
  "name": "robot",
  "x": 0,
  "y": 3,
  "direction": [
    1,
    0
  ],
  "actions": [
    "move",
    "moveitem",
    "compresschest",
    "giverescuebreaths",
    "giveshock",
    "givemedicine"
  ],
  "skill_info": {
    "compresschest": 1,
    "giverescuebreaths": 1,
    "giveshock": 1,
    "givemedicine": 1
  }
},

{
  "name": "robot",
  "x": 3,
  "y": 3,
  "direction": [
    -1,
    0
  ],
  "actions": [
    "move",
    "moveitem",
    "compresschest",
    "giverescuebreaths",
    "giveshock",
    "givemedicine"
  ],
  "skill_info": {
    "compresschest": 1,
    "giverescuebreaths": 1,
    "giveshock": 1,
    "givemedicine": 1
  }
},

{
  "name": "robot",
  "x": 4,
  "y": 3,
  "direction": [
    1,
    0
  ],
  "actions": [
    "move",
    "moveitem",
    "compresschest",
    "giverescuebreaths",
    "giveshock",
    "givemedicine"
  ],
  "skill_info": {
    "compresschest": 1,
    "giverescuebreaths": 1,
    "giveshock": 1,
    "givemedicine": 1
  }
}

```

Figure 11: JSON to configure the HCWs skill information for equal skill.

```

{
  "name": "robot",
  "x": 0,
  "y": 3,
  "direction": [
    1,
    0
  ],
  "actions": [
    "move",
    "moveitem",
    "compresschest",
    "giverescuebreaths",
    "giveshock",
    "givemedicine"
  ],
  "skill_info": {
    "compresschest": 3,
    "giverescuebreaths": 1,
    "giveshock": 1,
    "givemedicine": 1
  }
},

{
  "name": "robot",
  "x": 3,
  "y": 3,
  "direction": [
    -1,
    0
  ],
  "actions": [
    "move",
    "moveitem",
    "compresschest",
    "giverescuebreaths",
    "giveshock",
    "givemedicine"
  ],
  "skill_info": {
    "compresschest": 1,
    "giverescuebreaths": 2,
    "giveshock": 1,
    "givemedicine": 1
  }
},

{
  "name": "robot",
  "x": 4,
  "y": 3,
  "direction": [
    1,
    0
  ],
  "actions": [
    "move",
    "moveitem",
    "giverescuebreaths",
    "giveshock",
    "givemedicine"
  ],
  "skill_info": {
    "compresschest": 1,
    "giverescuebreaths": 1,
    "giveshock": 2,
    "givemedicine": 1
  }
}

```

Figure 12: JSON to configure the HCWs skill information for specialized roles.

```

{
  "name": "robot",
  "x": 0,
  "y": 3,
  "direction": [
    1,
    0
  ],
  "actions": [
    "move",
    "moveitem",
    "compresschest",
    "giverescuebreaths",
    "giveshock",
    "givemedicine"
  ],
  "skill_info": {
    "compresschest": 3,
    "giverescuebreaths": 0,
    "giveshock": 0,
    "givemedicine": 1
  }
},

{
  "name": "robot",
  "x": 3,
  "y": 3,
  "direction": [
    -1,
    0
  ],
  "actions": [
    "move",
    "moveitem",
    "compresschest",
    "giverescuebreaths",
    "giveshock",
    "givemedicine"
  ],
  "skill_info": {
    "compresschest": 0,
    "giverescuebreaths": 2,
    "giveshock": 0,
    "givemedicine": 1
  }
},

{
  "name": "robot",
  "x": 4,
  "y": 3,
  "direction": [
    1,
    0
  ],
  "actions": [
    "move",
    "moveitem",
    "compresschest",
    "giverescuebreaths",
    "giveshock",
    "givemedicine"
  ],
  "skill_info": {
    "compresschest": 0,
    "giverescuebreaths": 0,
    "giveshock": 2,
    "givemedicine": 1
  }
}

```

Figure 13: JSON to configure the HCWs skill information for required cooperation

Table 10: Comprehensive evaluation of FairSkillMARL and FEN using fairness metrics averaged over last 100 episodes. All fairness metrics (L1) and agent contributions are computed as moving averages to ensure robust measurement.

Experiment	λ	Method	α	Success Rate \uparrow	L1 \downarrow	Averaged Agent %			Range \downarrow
					L1	A0	A1	A2	
Experiment 1: Impact of Skill-Task Alignment (L2) - Specialized Teams									
$\lambda = 0.0$	0.0	Workload-only	0.0	0.110 ± 0.10	0.270	33.25	36.61	30.14	6.47
		FairSkillMARL	0.7	0.096 ± 0.06	0.323	39.64	39.66	20.69	18.97
		FEN	—	0.086 ± 0.05	0.300	28.49	38.17	33.34	9.68
$\lambda = 1.0$	1.0	Workload-only	0.0	0.104 ± 0.07	0.342	44.17	40.48	15.35	28.82
		FairSkillMARL	0.7	0.067 ± 0.04	0.398	49.94	34.75	15.31	34.63
		FEN	—	0.159 ± 0.18	0.351	40.98	36.84	22.18	18.80
$\lambda = 4.0$	4.0	Workload-only	0.0	0.081 ± 0.07	0.371	46.35	32.65	21.00	25.35
		FairSkillMARL	0.7	0.113 ± 0.08	0.351	37.44	48.74	17.82	30.92
		FEN	—	0.429 ± 0.44	0.254	31.48	41.86	26.65	15.21
Experiment 2: Team Composition Effects (Equal vs Specialized Teams, $\lambda = 1.0$)									
Specialized	1.0	Workload-only	0.0	0.104 ± 0.07	0.342	44.17	40.48	15.35	28.82
		FairSkillMARL	0.7	0.067 ± 0.04	0.398	49.94	34.75	15.31	34.63
		FEN	—	0.159 ± 0.18	0.351	40.98	36.84	22.18	18.80
Equal Skills	1.0	Workload-only	0.0	0.118 ± 0.06	0.401	33.63	39.44	26.93	12.51
		FairSkillMARL	0.7	0.245 ± 0.18	0.342	29.59	39.48	30.92	9.89
		FEN	—	$0.194 \pm \text{nan}$	0.338	45.27	50.40	4.33	46.07
Experiment 3: FairSkillMARL Ablation - Optimal α Selection ($\lambda = 1.0$)									
$\lambda = 1.0$	1.0	Workload-only ($\alpha = 0$)	0.0	0.060 ± 0.04	0.419	46.08	36.75	17.17	28.91
		FairSkillMARL	0.5	0.092 ± 0.02	0.388	45.10	36.63	18.27	26.83
		FairSkillMARL	0.7	0.133 ± 0.06	0.328	44.73	35.28	19.99	24.74
		FairSkillMARL ($\alpha = 1$)	1.0	0.050 ± 0.02	0.316	47.10	39.39	13.71	33.39
		FEN	—	0.146 ± 0.10	0.337	42.26	41.82	15.92	26.34

Note: Fairness metrics (L1) and agent contributions are averaged over the last 100 episodes. Success rates are calculated at convergence. Range = $\max(\text{Agent}\%) - \min(\text{Agent}\%)$.