

Single-Agent to Swarm: Proximal Policy Optimization (PPO)-Based Reinforcement Learning for Autonomous Tile Cleaning

Mohamed Elmohandes, Abdelrahman Ahmed, and Muhammed Noshin

Department of Computer Engineering, American University of Sharjah,
Sharjah, United Arab Emirates

{b00083108@alumni.aus.edu, b00087338@alumni.aus.edu, b00085566@alumni.aus.edu}

Abstract—This work presents the first fully decentralized, communication-free swarm reinforcement learning framework for robotic cleaning, where each e-puck robot independently learns a Proximal Policy Optimization (PPO) policy using only local sensing and individual rewards. Building upon the single-agent PPO framework of Pathmakumar et al. [1], we extend it to a multi-agent system with no shared critic, no centralized controller, and no inter-agent communication. A rigorous comparison of four training schedules in a 10×10 Webots environment shows that episodic updates with 1000 steps yield the most stable convergence, while longer or per-step updates introduce instability or overfitting. Scaling the system to two and six agents, we observe that the 2-agent swarm achieves consistent coverage gains and reward improvements, while the 6-agent swarm, deployed in a 20×20 grid, faces congestion yet retains a positive learning trajectory. Incorporating proximity-based penalties further stabilizes behavior without degrading performance. These results establish a reproducible and scalable baseline for multi-agent PPO in cleaning tasks, demonstrating that effective coordination can emerge purely from environmental interaction in the absence of communication or shared policy structures.

Index Terms—Reinforcement Learning (RL), Deep Reinforcement Learning, Single-Agent Systems, Multi-Agent Systems, Swarm Robots, Cleaning Robots, Autonomous Robots, Artificial Intelligence,

I. INTRODUCTION

Autonomous robotic cleaning has become an essential application in service robotics, particularly in indoor environments where repetitive tasks such as vacuuming, sweeping, and disinfection must be executed efficiently and with minimal human intervention. These systems are now routinely deployed in residential, industrial, and commercial settings, driven by increasing demands for hygiene, labor efficiency, and autonomous maintenance solutions [2], [3]. While single-agent cleaning robots remain prevalent, they are inherently constrained by sequential task execution, limited field of view, and poor adaptability to complex, dynamic environments. These limitations often result in redundant coverage, increased task duration, and incomplete cleaning, particularly in large or cluttered spaces.

To address these challenges, researchers have explored multi-agent robotic systems wherein multiple robots operate simultaneously within a shared environment. This architectural paradigm offers substantial advantages, such as parallelized

coverage, improved fault tolerance, and distributed workload balancing. However, it also introduces non-trivial challenges in coordination, collision avoidance, task allocation, and real-time communication, especially under conditions of partial observability and decentralized control [4], [5]. Existing rule-based or heuristic coordination schemes often fail to generalize across diverse environmental layouts or varying agent configurations, making them brittle and difficult to scale.

Reinforcement learning (RL) provides a promising alternative by enabling agents to learn optimal policies through interaction with their environment. In robotic cleaning, RL has been successfully applied to tasks such as obstacle avoidance, area exploration, and coverage optimization. When extended to the multi-agent domain, RL, specifically multi-agent reinforcement learning (MARL), enables the emergence of cooperative behavior through decentralized policy learning. Empirical studies have demonstrated that MARL frameworks can outperform static rule-based methods, particularly in dynamic or stochastic indoor settings where adaptability and distributed decision-making are critical [6], [7], [8]. Nevertheless, most MARL frameworks rely on shared policies, centralized critics, or explicit inter-agent communication, which can become impractical or inefficient in large-scale or communication-constrained deployments [9].

In this paper, we introduce a swarm-based reinforcement learning framework for robotic cleaning, where each agent learns its own control policy using Proximal Policy Optimization (PPO). Unlike conventional multi-agent systems that rely on shared policies or explicit communication, our approach enables fully independent learning under partial observability. Each robot operates using local observations and optimizes its behavior without access to global state information or peer coordination. This formulation scales naturally to large teams and allows complex behaviors to emerge from simple, decentralized policies. We build directly on the single-agent PPO model of Pathmakumar et al. [1], extending it to a robust swarm architecture capable of handling diverse layouts, agent densities, and coverage demands.

The main contributions of this work are as follows:

- We introduce the first RL framework for swarm-based robotic cleaning, where each agent independently learns its own PPO policy under partial observability and without communication.

- We extend the single-agent PPO formulation by Pathmakumar et al. [1] to a fully decentralized swarm setting, enabling scalable and emergent behavior without policy sharing or centralized control.
- We develop a simulation environment for evaluating swarm learning across diverse indoor layouts, varying agent densities, and limited sensing conditions.
- We expanded upon previous works [1] by adding multiple agents, to make a swarm of RL-trained bots.

The remainder of this paper is structured as follows: Section II introduces key concepts in robotic cleaning architectures and reinforcement learning. Then, Section III covers related works. Section IV describes the proposed decentralized multi-agent framework and training setup. Section V presents the experimental results followed by a discussion in Section VI. Then, Section VII discusses the open challenges and future work to be done in this field of research. Finally, Section VIII summarizes our findings and outlines future directions for MARL-based robotic cleaning.

II. BACKGROUND

A. Coordination Paradigms

Robotic systems are generally categorized into three architectural paradigms: single-agent, multi-agent, and swarm-based systems. These paradigms differ significantly in terms of scalability, autonomy, and coordination mechanisms. In single-agent robotic systems (SARS), an autonomous robot operates independently, relying solely on its own sensors and decision-making processes [10]. Multi-agent robotic systems (MARS) involve multiple agents, either homogeneous or heterogeneous, that coordinate through explicit communication or task allocation to improve performance [11]. In contrast, swarm robotics (SR) draws inspiration from natural systems such as insect colonies and bird flocks. Swarm systems consist of a large number of simple, homogeneous agents that coordinate indirectly through local interactions and environmental cues, without relying on centralized control or global communication [12].

While all three paradigms have been explored in robotic cleaning, swarm-based systems offer unique advantages in scalability, robustness, and fault tolerance [13]. However, most existing swarm implementations depend on hand-crafted heuristics or rule-based behaviors, which limit adaptability and performance in dynamic environments. Although Reinforcement Learning (RL) has shown promise in both SARS and MARS by enabling adaptive behavior through trial and error interaction, its application to swarm systems remains underexplored [14].

B. Reinforcement Learning (RL)

RL is a computational framework for sequential decision-making. It allows an agent to learn optimal behavior by interacting with an environment and receiving feedback in the form of rewards. An RL problem is typically formulated as a Markov Decision Process (MDP), defined by the tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$. Here, \mathcal{S} represents the set of states, \mathcal{A} denotes

the set of actions, P is the state transition probability distribution, R is the reward function, and $\gamma \in [0, 1)$ is the discount factor that balances immediate and future rewards [15].

The objective is to learn a policy $\pi(a | s)$ that maximizes the expected return:

$$G_t = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \right].$$

This return represents the discounted sum of rewards the agent receives over time.

RL has proven effective in domains where explicit programming is infeasible due to environmental uncertainty, high dimensionality, or continuous state spaces [16]. In robotic systems, it supports real-time learning, closed-loop control, and task adaptation, making it especially useful for applications like navigation, exploration, and cleaning [17].

1) *Model-Free and Model-Based RL*: Reinforcement learning methods can be broadly grouped based on whether they explicitly model the environment. This distinction leads to two main categories: model-free and model-based algorithms [15].

Model-free approaches, such as Q-learning, SARSA, and actor-critic variants, estimate value functions or policies directly from experience without attempting to learn the environment's transition dynamics. These methods are widely adopted in robotics due to their conceptual simplicity, robustness to environmental uncertainty, and ability to scale to high-dimensional state and action spaces [18].

In contrast, model-based methods attempt to learn or approximate the environment's dynamics, denoted $\hat{P}(s' | s, a)$, and optionally the reward function $\hat{R}(s, a)$. This enables the agent to plan ahead by simulating interactions internally before taking physical actions. While model-based RL can offer higher sample efficiency, it is often sensitive to modeling errors. These limitations are especially pronounced in real-time, stochastic, or unstructured settings, which hinders its practical deployment in complex robotic applications [19].

C. Proximal Policy Optimization (PPO)

PPO [20] is a first-order, model-free reinforcement learning algorithm that belongs to the policy-gradient family. It is designed to improve both the stability and reliability of training by penalizing large policy updates. Unlike Trust Region Policy Optimization (TRPO) [21], which enforces a hard constraint on the KL divergence between successive policies, PPO simplifies the optimization by introducing a clipped surrogate objective.

Let the probability ratio between the new and old policy be defined as:

$$r_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}.$$

The clipped objective function is given by:

$$L^{\text{CLIP}}(\theta) = \mathbb{E}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right],$$

where \hat{A}_t is the estimated advantage function and ϵ is a small threshold, typically set to 0.1 or 0.2. This objective discourages excessively large policy updates, thereby ensuring stable learning even in high-variance environments.

PPO alternates between policy rollout using the current policy and multiple epochs of stochastic gradient ascent on the clipped objective using minibatch updates. This decoupling of sampling and optimization facilitates parallel execution, significantly reducing sample inefficiency relative to other policy-gradient methods [22].

In robotic control, PPO is especially favored due to its compatibility with continuous action spaces, resistance to noisy sensor inputs, and strong empirical performance in partially observable domains [23]. Its computational simplicity and ease of implementation have led to widespread adoption in both simulated and real-world robotic systems, including robotic manipulation [24], navigation [25], and coverage-based cleaning [1], [26].

III. RELATED WORKS

This section reviews prior literature on RL-based robotic cleaning, with an emphasis on coordination paradigms and architectural frameworks. We categorize related work into single-agent, swarm-based, and multi-agent systems, and highlight insights that inform our swarm learning framework.

Early studies in this domain primarily explored single-agent RL for spatial coverage and path planning. Moon et al. [26] developed a PPO-based agent enhanced with transfer learning and reward shaping to improve generalization across different indoor layouts. Moreover, Sun et al. [27] applied tabular Q-learning to agricultural cleaning robots, where agent behavior was influenced by livestock movement patterns. Their approach achieved a 67.6% reduction in collision rates while maintaining high cleaning coverage. While effective in structured environments, these systems were limited by their inability to scale or coordinate across multiple agents.

To address the limitations of isolated agents, researchers have explored decentralized architectures inspired by swarm intelligence. Swarm-based systems typically consist of large numbers of homogeneous agents coordinating through simple rules and local interactions, without explicit communication. Hu et al. [28] implemented a reactive obstacle avoidance strategy in Colias robots, allowing them to navigate cluttered environments autonomously. Maryasin et al. [29] proposed a scout-worker model based on bee foraging behavior to divide cleaning tasks spatially. These biologically inspired approaches offer robustness and scalability, but rely heavily on handcrafted heuristics and lack adaptability, making them brittle in dynamic environments.

Building upon the scalability of swarm architectures and the adaptability of learning-based methods, several studies have applied MARL to robotic cleaning and monitoring tasks. To enhance decentralized cooperation, Siddiqua et al. [30] introduced a MARL framework with periodic inter-agent state sharing, achieving a 45% improvement in warehouse cleaning efficiency. Addressing robustness under failure, Luis et al. [31] applied Double Deep Q-Learning (DDQL) to autonomous water-monitoring robots using fault-aware reward shaping. Barrionuevo et al. [32] extended this line of work by integrating dueling networks with prioritized replay to coordinate heterogeneous agents in garbage collection, achieving over 96% spatial coverage. For real-time adaptability,

Zhao et al. [33] proposed a collaborative Q-learning protocol using shared occupancy maps and dynamic goal reassignment, yielding 98% coverage with minimal overlap. To support generalization across layouts, Qureshi et al. [34] introduced a graph-based Q-network that abstracts spatial structure for scalable patrolling. Despite their merits, these systems often rely on centralized critics, synchronized training, or explicit communication-factors that limit scalability in swarm-based deployments [9].

To reduce reliance on centralized training while retaining coordination benefits, several hybrid approaches have been proposed that combine centralized planning with decentralized execution. Gupta et al. [35] introduced MARBLE, a coordination framework where a centralized planner manages routing and charging, while local policies are executed autonomously across a heterogeneous robot fleet. Building on decentralized coordination, Collette et al. [36] proposed a proximity-based task assignment strategy that enables implicit coordination without direct communication, albeit under the assumption of static task distributions. To improve learning stability, Kim et al. [37] implemented a hierarchical MARL architecture that separates global task allocation from low-level execution, thereby enhancing convergence and policy robustness. In parallel, Wang et al. [38] employed decentralized POMDPs in shape-adaptive mop robots, demonstrating improved task completion and coverage through hardware-level reconfiguration. While these strategies enhance adaptability and policy scalability, they frequently depend on partial centralization or predefined coordination hierarchies, limiting their applicability in fully distributed swarm scenarios.

Complementary to these coordination strategies, recent efforts have integrated environmental perception into learning frameworks to enable more context-aware cleaning behaviors. Din et al. [39] used a convolutional neural network (CNN) to detect soiled regions and guide a DDQL-based agent, achieving a 25% improvement in cleaning efficiency through targeted execution. Similarly, Caccavale et al. [40] utilized Wi-Fi signal strength to construct dynamic heatmaps for disinfection robots, which reduced system response times by 35% through spatial prioritization. These studies highlight the value of coupling sensory information with policy learning, yet they remain largely confined to single-agent or centrally coordinated frameworks, limiting their scalability in decentralized swarm deployments.

Among the reviewed literature, the work by Pathmakumar et al. [1] is most directly aligned with our methodology. Their PPO-based agent for single-robot cleaning used onboard sensing and a sparse-reward formulation to learn efficient dirt detection and coverage policies. Although effective in simulation, their framework was restricted to a single agent and did not address distributed coordination or policy scalability. In contrast, our work extends this formulation to a fully decentralized swarm setting, where each agent learns an independent PPO policy using only local observations and individual rewards. Unlike traditional MARL systems, our approach does not rely on shared critics, communication, or centralized supervision-making it a novel application of reinforcement learning to scalable, communication-free swarm

cleaning.

IV. METHODOLOGY

In this study, Webots, the robotics simulation platform, was used to develop and evaluate both single-agent and multi-agent robotic systems trained using reinforcement learning. The e-puck robot was used for the single and multi-agent approaches. Figure 1 shows the how the environment looks for a single agent implementation. Multi-agent implementation is very similar to this but with different dimensions and number of robots. Red tiles are obstacles with height and physics. The black tiles are dirt tiles while the white tiles are empty. We have created those tiles in the code using a supervisor controller. This code is used to manage the policy and the environment creation which includes the dirt locations, obstacle locations along with their dimensions. On the other hand, the robot controller is for moving the robot.

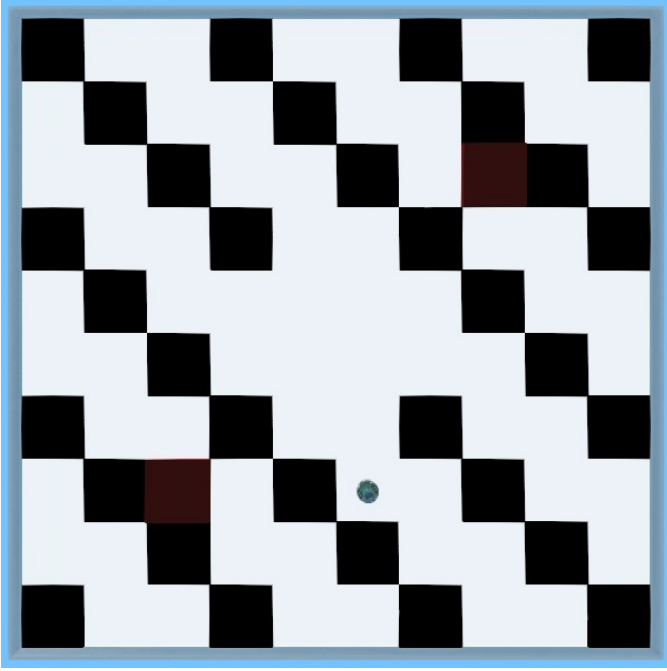


Fig. 1. Webots Environment

A. Hardware used

All experiments were conducted on consumer-grade hardware. For the single-agent setup, training was performed on a machine equipped with an NVIDIA RTX 3060 GPU, an AMD Ryzen 5 5600H CPU, and 16 GB of RAM. For the multi-agent experiments, a system with an NVIDIA RTX 3060 GPU, AMD Ryzen 7 4800HS CPU, and 16 GB of RAM was used. Both configurations provided sufficient computational resources to train the models efficiently, with GPU acceleration significantly speeding up the policy updates and inference during training. However the experiments take a long time to run, hence, we were limited on the number of experiments to perform due to time constraints. A more powerful machine

is required to run these simulations quicker as mentioned in [41] where they used an intel core-i7 CPU and an RTX-3080 GPU.

B. Single-Agent Implementation

For the single agent setup, an e-puck robot was deployed in a custom-designed 2x2 meter environment subdivided into a 10x10 grid, resulting in 100 individual tiles. Each tile could contain dirt, be empty, or represent an obstacle. A supervisor controller managed the overall environment, policy, and reward mechanism using the Proximal Policy Optimization (PPO) algorithm. The e-puck robot interacted with the environment by selecting one of three discrete actions: move forward, turn left, or turn right, based on its current 2D state. The robot had no global knowledge and operated solely based on the actions received from the supervisor. The supervisor in this case acted as a proxy to the robot's processor, to make the coding process easier, especially when expanding this to multiple agents. Communication between the robot and the supervisor was implemented using an emitter-receiver system on channel 1. The policy was iteratively updated based on the rewards gathered during episodes, and the robot learned to navigate and clean tiles effectively through exploration and exploitation over multiple episodes. Table I shows the parameters that we used in training the PPO algorithm. The reward system of our proposed approach is shown in table II. It is of note that we have forced the robot to stay within bounds due to several reasons. To discourage unnecessary wandering, a reward penalty was incorporated into the reward function. Although reinforcement learning agents can eventually learn these behavior patterns through extensive exploration, this would require significantly longer training durations. Due to time constraints, the wandering penalty and the obstacle collision penalty were predefined as part of the reward structure without actual effect on the outcome and they are never rewarded or deducted. The remaining rewards explained in the table are incorporated in the algorithm.

Our paper is inspired and builds up from the existing works of [1]. Table III show the parameters, variables and key information about the algorithm they used.

During training, the robot starts each episode at a randomly selected initial position within the environment and we make sure it starts in a safe location where there are no obstacles. The policy is updated at the end of every episode based on the accumulated rewards and transitions. Additionally, we experimented with an alternative setup for single agent system where the policy is updated at every timestep instead of after full episodes. A comparative analysis of both strategies is conducted to evaluate their impact on learning performance and convergence in section V.

C. Multi-Agent Implementation

As for the multi-agent approach, we simply expanded the single-agent code to accommodate multiple agents. This was done by adding a parameter that stores that number of robots being used. This was then used to loop through the robot names, and assigning each robot its own PPO policy. These

TABLE I
OUR PROPOSED PPO AGENT HYPERPARAMETERS AND CONFIGURATION

Component	Value	Explanation
State dimension	2	Robot's 2D position on the grid (x, y)
Action dimension	3	Three possible actions: forward, left, right
Hidden layers	2×128 neurons	Two dense layers of 128 neurons each
Learning rate	3e-4	Step size for gradient updates
Gamma (γ)	0.99	Discount factor for future rewards
Epsilon clip	0.2	Limits policy update deviation
Entropy coeff.	0.01	Adds randomness to encourage exploration
Epsilon decay	0.995	Decay rate for exploration-exploitation balance
Epsilon range	$0.9 \rightarrow 0.05$	Initial to final exploration rate

TABLE II
REWARD DESIGN FOR THE CLEANING ROBOT ENVIRONMENT

Condition	Reward	Explanation
Default step	-0.1	Penalizes movement to encourage efficiency
Out-of-bounds move	-5.0	Heavy penalty to discourage leaving grid limits
Collision with obstacle	-2.0	Penalty for hitting obstacles on the grid
Tile contains dirt	+5.0	Large positive reward for cleaning a dirty tile
First-time visit to tile	+0.2	Small reward for exploring new tiles
Agent stuck in same position for 5 steps	-0.5	Penalty for lack of movement or being stuck

policies learn independently of each other. There are no means of communication between the bots and although the supervisor is technically a "central controller" in this case, it is only due to the nature of our code and to simplify it instead of having pieces of code scattered across different controllers, one for each robot. This would be unsustainable once we decide to use an increasing number of robots. Therefore, this setup mimics a typical swarm set up. Each robot uses the same controller code, therefore they are all homogeneous. The initial states of each robot were determined in a similar manner to the single-agent robot, with the added condition that no two robots should start on the same tile. An additional punishment was added in the case that two robots found themselves in the same tile.

We tested firstly using two robots in the same environment as the single-agent and recorded the results. From there, we decided to experiment using six robots. Seeing as the environment would be somewhat overcrowded with six robots, we made the decision to increase its size to a 4x4 meter environment that is subdivided into 20x20 tiles. We also noticed that the simulation slowed down significantly with six robots, therefore we reduced the training to 150 episodes.

We then went back to the 10x10 tile environment with only two robots to see its performance improve on an episode by episode basis. We decided to increase the maximum steps per episode from 1000 to 50,000, and updated the policy every 1000 steps instead of every episode. This was after two failed attempts at the same task: one where we left the episode completely unbounded, and another where the episode was bounded to a maximum of 50,000 steps with policy

TABLE III
RL PARAMETERS FROM PRIOR WORK

Parameter	Value	Description
Environment	7x7 grid	Downsampled map
Agent	BELUGA robot	In-house cleaning robot
State	{pos, obstacle}	Position and obstacle profile
Action space	4 moves + 2 sensors	$\{p_0-p_3, S_0, S_1\}$
Algorithm	PPO	On-policy method
Hidden layers	3x512 (tanh)	Neural net config
Discount (γ)	0.99	PPO default
Clip param (ϵ)	0.2	PPO stability
Learning rate	0.001	Fixed rate
Training envs	4 layouts	Varied patterns
Max steps	35	Episode limit
Reward	$R_{dirt} + R_{unique} + R_{coll}$	Composite reward

updates occurring at the end of every episode. Both observed deteriorating performance that we were not able to record. The final variation we made was add sensors to the robots and communicated their values to the supervisor. Based on this, there was a punishment for when the sensors read a value that is closer than a given threshold to simulate collision avoidance.

V. RESULTS

A. Single-Agent

The grid in the single agent contains 32 tiles uncleaned. To evaluate the performance of the single agent learning strategy, we performed several tests that vary the learning rate. In the episode-based learning approach, where the policy is updated only at the end of each episode, we tested three configurations: (i) 1000 time steps over 300 episodes, (ii) 1000 time steps over 600 episodes, and (iii) 2000 time steps over 300 episodes. In addition, we tested a step-based learning approach in which the policy is updated after every individual action, enabling more frequent learning updates. We did it for 1000 time steps per episode for a total of 600 episodes. It is worth mentioning that the tests that had only a difference in number of episodes being trained were two independent tests, meaning that we initialized everything for each test. Figures 2–9 show the different plots for all the experiments we have performed on a single agent system.

We also performed an test where we wanted to plot the Episode versus the cumulative time steps taken to show that the more we train the less time it takes to clean all the tiles. We set a maximum cap of 50,000 time steps. This was motivated by the observation that, after completing most of the cleaning task, the agent tends to wander and consume excessive time, making cumulative time step analysis time-consuming. The results for this test are shown in Figure 10.

B. Multi-Agent

Starting with the results obtained from training the 2-bot swarm, we can see there is a general improvement throughout the episodes in 11. Going from -400 rewards to -100 in only 600 episodes. The lower graph is more insightful as it shows the tiles cleaned per episode per robot. It is notable that there is an overall upwards trend, with the major fluctuation occurring due to each robot beating the other to more tiles.

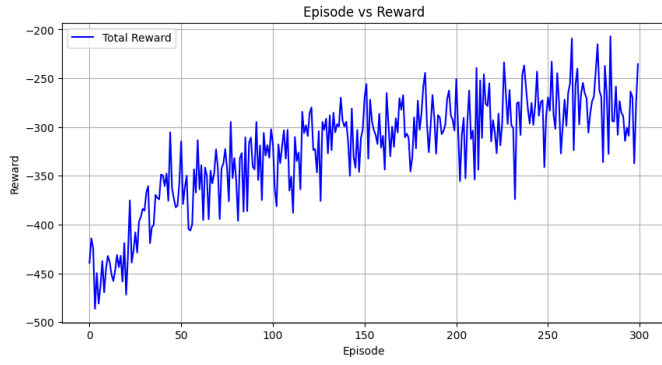


Fig. 2. Single Agent: Episode vs Total Reward (1000 Time steps per episode for 300 episodes)

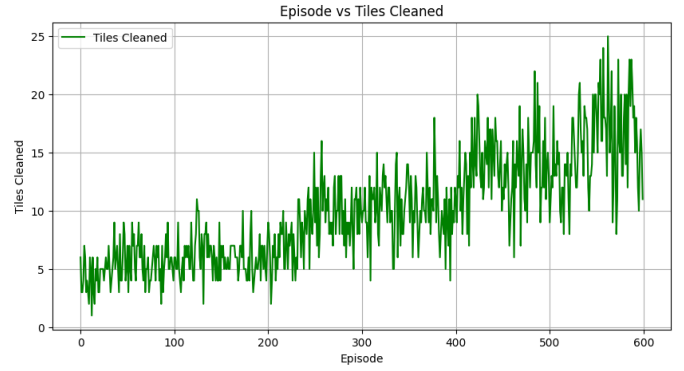


Fig. 5. Single Agent: Episode vs Total Tiles Cleaned (1000 Time steps per episode for 600 episodes)

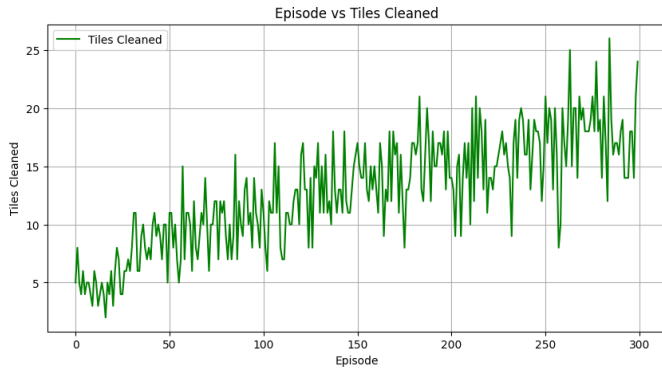


Fig. 3. Single Agent: Episode vs Total Tiles Cleaned (1000 Time steps per episode for 300 episodes)

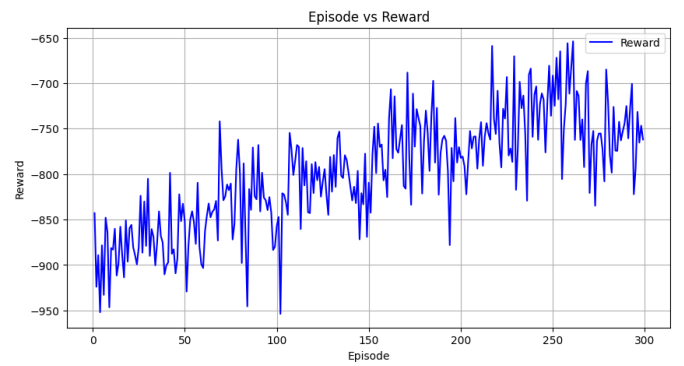


Fig. 6. Single Agent: Episode vs Total Reward (2000 Time steps per episode for 300 episodes)

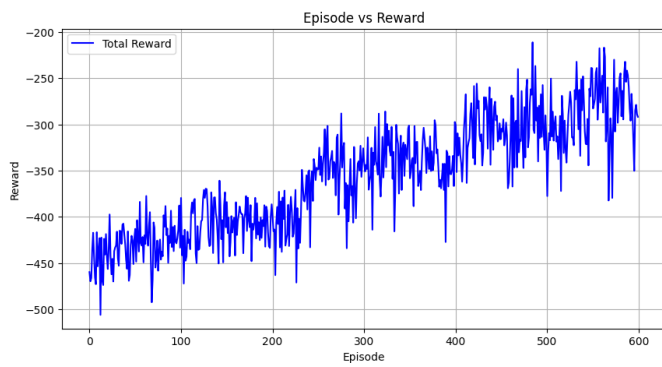


Fig. 4. Single Agent: Episode vs Total Reward (1000 Time steps per episode for 600 episodes)

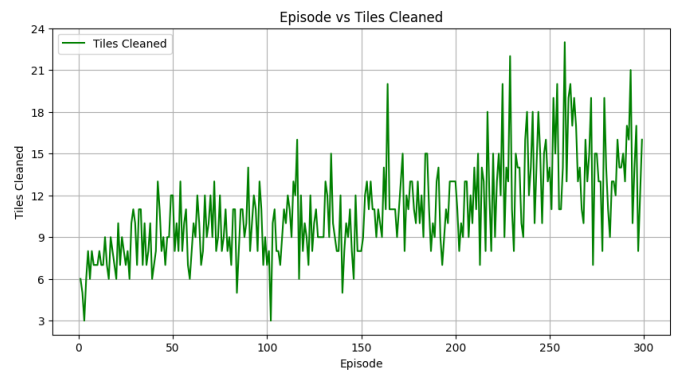


Fig. 7. Single Agent: Episode vs Total Tiles Cleaned (2000 Time steps per episode for 300 episodes)

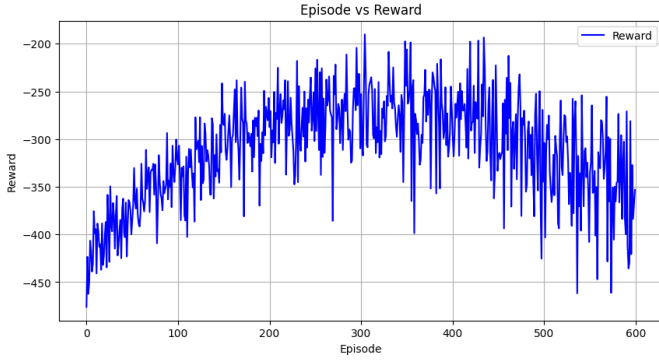


Fig. 8. Single Agent: Episode vs Total Reward (Update after 1 action taken with 1000 time steps per episode for 600 episodes)

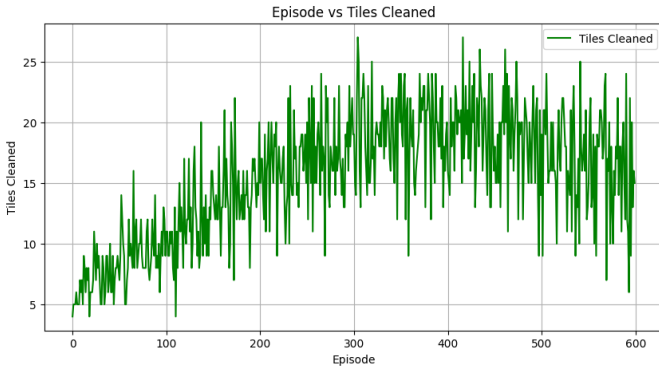


Fig. 9. Single Agent: Episode vs Total Tiles Cleaned (Update after 1 action taken with 1000 time steps per episode for 600 episodes)

Similar observations can be made for the 6-bot results as seen in 12.

13 shows the performance of the 2-bot swarm with sensor-based punishment for coming closer than a set threshold distance.

VI. DISCUSSION

Having presented our experimental results, we now conduct a detailed analysis of the performance characteristics observed

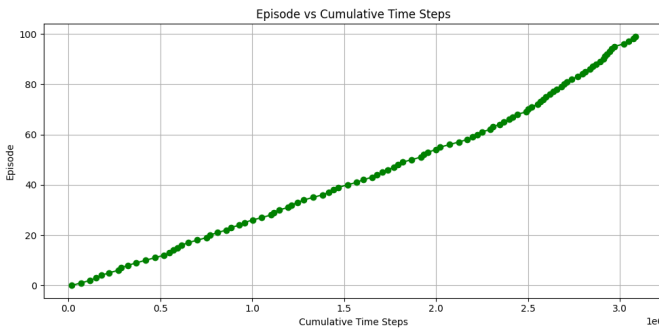


Fig. 10. Single Agent: Episodes vs Cumulative Timestep for cleaning the entire map or reaching a max timestep of 50k timesteps)

in both single-agent and multi-agent cleaning systems. We also used similar values as used by Pathmakumar et al. [1]. However, we had major differences such as the environment, reward system, PPO architecture, etc. We used the same algorithm they used and wanted to address their future work of having a multi-agent system.

A. Single-Agent

A common issue with all the single agent results is the amount of fluctuations. That can be due to several factors. For example, we made the robot initially start at random locations hence, some locations tend to be easier or difficult to start from compared to others. Another issue could be with the parameters such the learning rate as it causes learned values to fluctuate, harming long-term convergence.

1) *Episode-Based Learning*: All the tests in this subsection are episode-based learning. Our initial tests were done on a single agent with a 1000 time step per episode for 300 episodes. The results were performing well showing that the rewards over time increase as well as the tiles cleaned as shown in Figures 2–3. Following this test we wanted to know if the robot will perform even better after increasing the number of episodes and we can see that we get almost similar results which shows that we are approaching convergence as shown in Figures 4–5. We then wanted to test out increasing the number of time steps from 1000 to 2000 and we can see that it does not do as well as the 1000 time step approach as it will be penalized a lot for wandering in the extra time steps and will take longer to train as shown in Figures 6–7.

2) *Learning updated per action taken*: In this section we will analyze the approach where the policy is updated after every action taken. We can see from the Figure 10 that it reaches convergence level at around the same reward and tiles cleaned as in the previous tests, however, it performs worse after 300 episodes.

One potential reason for the observed decline in performance is overfitting to early episodes, where the agent may have learned a policy that performs well initially but fails to generalize as the training progresses. Additionally, an epsilon decay or exploration strategy failure could be limiting the agent's ability to discover better policies later in training, especially if exploration is reduced too quickly. Hence, hyperparameter tuning can lead to improving the performance. Another contributing factor could be sparse or misleading rewards, which make it difficult for the agent to accurately associate actions. Hence a revised reward system must be explored. Furthermore, the environment resets every 1000 time steps, and this long episode duration may lead to reward accumulation being diluted, making it harder for the agent to recognize the impact of beneficial actions amid the noise of many negative rewards.

3) *Episode vs Cumulative Time Steps*: In Figure 10 we see that the line seems to be a straight line with some curvature towards the end. Ideally this curve should be exponential which means that the time steps required for each episode will decrease as episode increases. The reason we do not have such a curve could be due to the difficulty of the problem, as well as the parameters we chose.

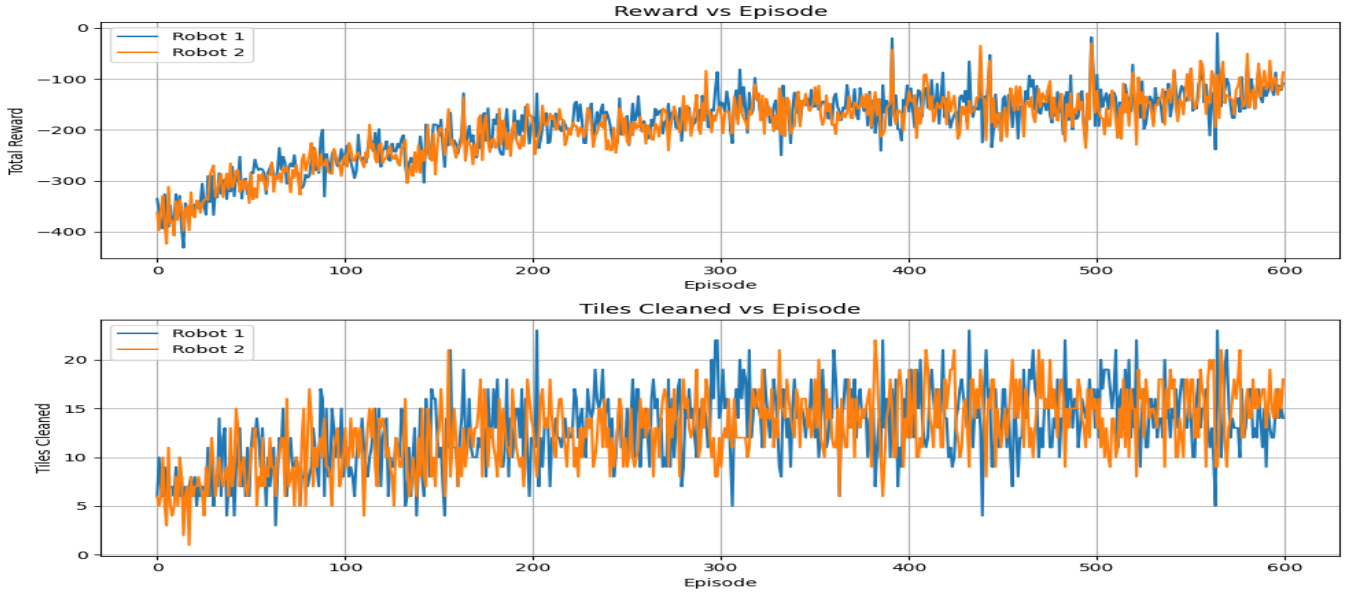


Fig. 11. 2-bot swarm results

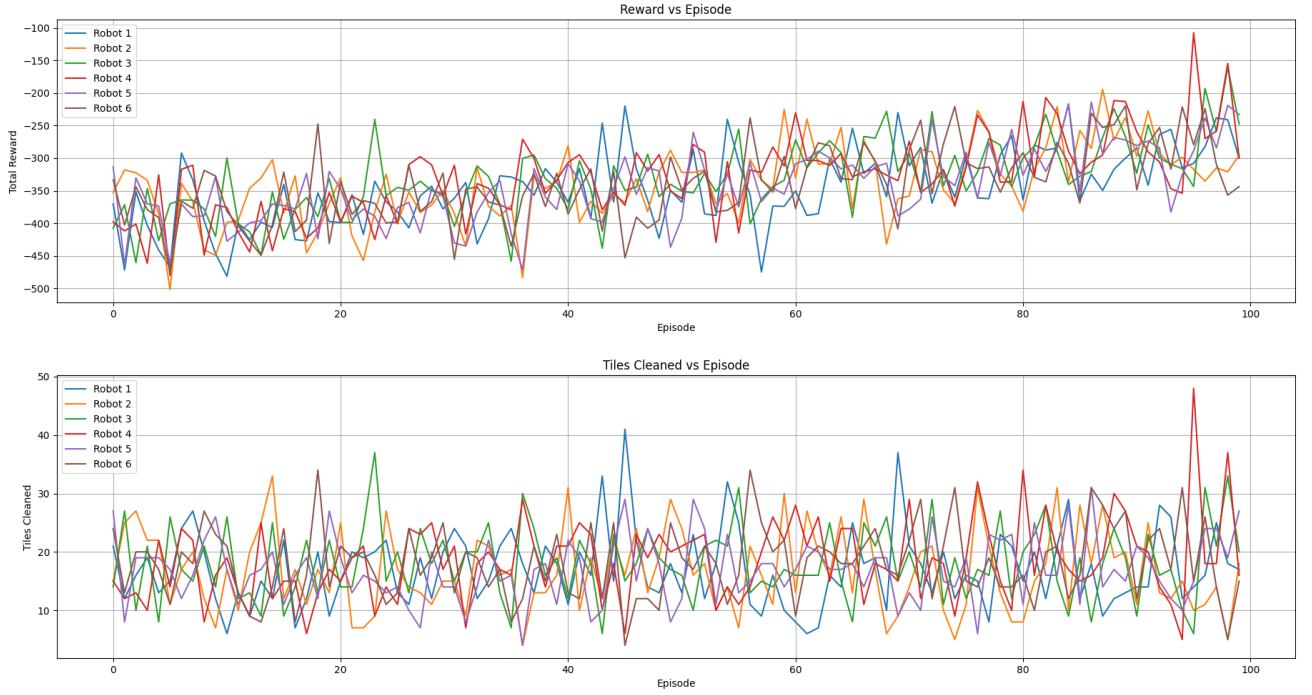


Fig. 12. 6-bot swarm results

B. Multi-Agent

One of the potential issues could be that we may have trained it for too long, at least in comparison to [1]’s 40 steps per episode. Although, the performance and overall improvement in training observed was decent, it could be improved by some tuning of the steps per episode, as well as increasing the number of episodes for the training.

It is notable that the swarm that used proximity sensors for collision avoidance 13 saw similar performance for the same duration of training as the swarm that did not 11. This is despite the fact that there are more factors that each robot will have to take into consideration. The improved collision avoidance could be causing improvements in the cleaning performance of the robots by avoiding movements that would not result in any benefit (moving towards a wall). Therefore

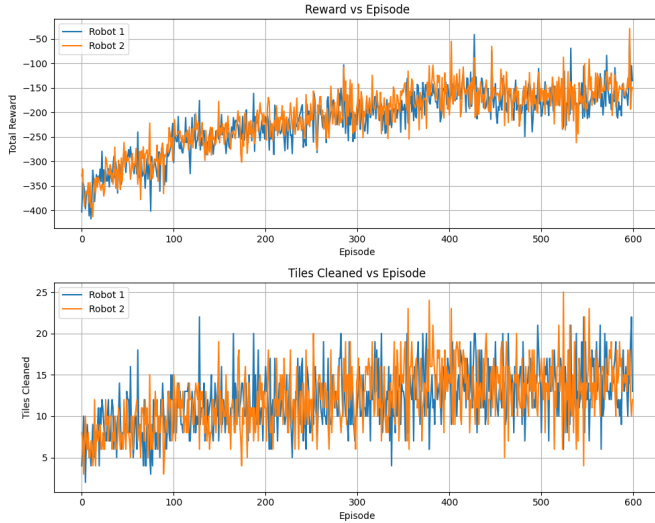


Fig. 13. Swarm With Proximity Sensors

the improvement allows the robots to overcome any setbacks that could have been caused by the increase in the number of factors affecting each individual robot's rewards.

Upon observing the results of the 6-bot swarm 12 and comparing it to the 2-bot swarm 11, we can see that the former's results are far noisier. This makes sense due to the presence of more robots on the map and therefore collisions are more likely. Keeping in mind that the environment was kept the same throughout the training, bots would likely have begun to "know" where to clean, but due to the number of robots present, tile "stealing" would have been likely an environment where they operate in such close proximity. It is of importance that the overall trend of rewards is increasing. However, looking at the cleaning per episode of the 6-bot swarm, it is harder to make out an upward trend. This is likely because the environment was too small to begin with, and the bots were able to clean it entirely seemingly from the first few episodes.

VII. OPEN CHALLENGES AND FUTURE DIRECTIONS

Implementing deep reinforcement learning for robotic cleaning tasks in a simulated environment presented several challenges due to time and limited hardware constraints. One major issue was the lack of diverse and dynamic testing environments, which limited the model's ability to generalize to new or changing scenarios. The environment remained mostly static, and as a result, the policy may not perform well when transferred to real-world conditions where obstacles and layouts vary. Additionally, due to time constraints and limited hardware resources, it was not feasible to train over extended episodes. A better designed reward system to improve the learning as well as hyperparameter tuning is required to improve the overall performance of the system.

A limitation of our is that we assume that there is a tool attached to the robot meaning that when it passes over a clean tile it cleans it. Therefore a better, more realistic approach is

to have a cleaning action where the robot will wait for some time and perform different cleaning actions, like sweeping or mopping the dirty floor. This requires a modification to our approach.

Another potential limitation of our work could be due to the robot's lack of communication. This could possibly encourage the robots to compete for the same tiles. Although this may not necessarily lead to collisions due to the presence of a punishment for two robots being in the same tile, it will make the cleaning process less efficient in general. Communication may allow robots to communicate their routes and therefore give others a chance to reroute and search for different dirty tiles.

One of the issues with our implementation is that the collision avoidance for obstacles and the boundaries of the environment did not work as intended - it was designed to punish the robot only once it was physically inside the boundary or obstacle, which was impossible due to the physics engine. We attempted to fix this in one of our later tests of multi-agent systems where we included sensor-based punishment (13.) however we could not repeat all of the previous experiments due to time constraints.

Future directions could involve expanding the project to test in multiple, procedurally generated or dynamic environments to improve generalization and robustness. Incorporating real-time changes in the environment such as moving obstacles or variable dirt distributions which could help the model learn more adaptable policies. Another key step would be to explore different reinforcement learning algorithms. With more time and access to advanced hardware (e.g., GPUs or distributed training setups), deeper architectures or multi-agent reinforcement learning algorithms like MADDPG or DDPG could also be investigated for enhanced performance and coordination [41]. Additionally, we should consider expanding the environment to accommodate the training of larger swarms.

VIII. CONCLUSION

This paper explored the feasibility of decentralized RL for swarm-based robotic cleaning in environments where agents operate under partial observability and without communication. Through systematic experimentation, we found that stable and effective learning is achievable even when each robot independently optimizes its policy using only local feedback. The single-agent PPO experiments confirmed that episodic updates with well-tuned step counts are critical to stable convergence, while longer or per-step updates can hinder learning through over-penalization or overfitting. Scaling to a multi-agent setting, we observed that decentralized PPO agents can exhibit emergent coordination and improved task coverage despite environmental congestion and the absence of shared information. The use of proximity-based penalties proved valuable in stabilizing behavior in denser swarms. More broadly, these results provide compelling evidence that decentralized learning strategies can scale to real-world swarm robotics applications, where communication overhead or centralized control is impractical. However, the system remains limited by fixed action spaces, synthetic environments, and

homogeneous agents. Future extensions will explore heterogeneous robot capabilities, dynamic environments, real-time hardware deployment, and hierarchical learning to improve task specialization and efficiency. This study offers a reproducible foundation for advancing scalable, fault-tolerant multi-agent reinforcement learning in autonomous robotic cleaning systems.

IX. RESOURCES AND LINKS

A video containing a brief presentation and a demonstration of the robot in action is available at: <https://youtu.be/8lsbPWx0NwA>

The full source code and simulation files are uploaded on GitHub: https://github.com/elmohandes2002/AdvAI_Project

REFERENCES

- [1] T. Pathmakumar, R. E. Mohan, B. F. Gómez, and B. Ramalingam, "A reinforcement learning based dirt-exploration for cleaning-auditing robot," *Sensors*, vol. 21, no. 24, p. 8331, 2021.
- [2] M. Turan, M.-T. Lam, and T.-K. Ko, "A survey on autonomous cleaning robots," *Robotics and Autonomous Systems*, vol. 145, p. 103857, 2021.
- [3] H.-C. Yang, B.-C. Chen, H.-T. Hsiao, and S.-H. Lai, "Hierarchical coverage path planning for multi-robot cleaning in large environments," *IEEE Access*, vol. 8, pp. 91 852–91 864, 2020.
- [4] H. Cao, Z. Zhang, M. Chen, Z. Wu, and L. Tang, "An intelligent multi-robot cleaning system for large public environments," *IEEE Access*, vol. 8, pp. 47 694–47 707, 2020.
- [5] L. Zhao, Y. Wang, Q. Liu, and H. Luo, "Multi-robot task allocation with dynamic task adjustment in cleaning scenarios," *IEEE Transactions on Automation Science and Engineering*, vol. 18, no. 2, pp. 653–666, 2021.
- [6] M. A. Alsheikh, M. G. Rashed, and M. Alazab, "Multi-agent reinforcement learning: A survey," *IEEE Access*, vol. 10, pp. 10 308–10 327, 2022.
- [7] Y. Chen, W. Wang, M. Tan, R. Deng, and B. Luo, "Multi-robot coverage path planning for cleaning in unknown environments," *IEEE Transactions on Cognitive and Developmental Systems*, vol. 13, no. 4, pp. 917–926, 2020.
- [8] J. Qin, J. Ren, and M. Li, "Multi-agent reinforcement learning for multi-robot systems: A survey," *IEEE Access*, vol. 9, pp. 118 145–118 166, 2021.
- [9] J. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson, "Counterfactual multi-agent policy gradients," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [10] B. Yamauchi, "Frontier-based exploration using multiple robots," *Proceedings of the Second International Conference on Autonomous Agents*, pp. 47–53, 1998.
- [11] L. Panait and S. Luke, "Cooperative multi-agent learning: The state of the art," *Autonomous Agents and Multi-Agent Systems*, vol. 11, no. 3, pp. 387–434, 2005.
- [12] E. Şahin, "Swarm robotics: From sources of inspiration to domains of application," in *Swarm Robotics*. Springer, 2005, pp. 10–20.
- [13] M. Brambilla, E. Ferrante, M. Birattari, and M. Dorigo, "Swarm robotics: a review from the swarm engineering perspective," *Swarm Intelligence*, vol. 7, no. 1, pp. 1–41, 2013.
- [14] M. Hüttenrauch, A. Sosc, M. Neunert, and J. Buchli, "Deep reinforcement learning for swarm systems," in *Proceedings of the International Conference on Autonomous Agents and MultiAgent Systems (AAMAS)*, 2019, pp. 1253–1261.
- [15] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. MIT press, 2018.
- [16] L. P. Kaelbling, M. L. Littman, and A. W. Moore, "Reinforcement learning: A survey," *Journal of Artificial Intelligence Research*, vol. 4, pp. 237–285, 1996.
- [17] P. Kormushev, S. Calinon, and D. G. Caldwell, "Reinforcement learning for robotic manipulation: Autonomy, generalization, and adaptation," *Autonomous Robots*, vol. 35, no. 4, pp. 361–367, 2013.
- [18] J. Kober, J. A. Bagnell, and J. Peters, "Reinforcement learning in robotics: A survey," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1238–1274, 2013.
- [19] M. P. Deisenroth and C. E. Rasmussen, "Pilco: A model-based and data-efficient approach to policy search," *Proceedings of the 28th International Conference on Machine Learning (ICML)*, pp. 465–472, 2011.
- [20] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [21] J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel, "Trust region policy optimization," in *International conference on machine learning*. PMLR, 2015, pp. 1889–1897.
- [22] M. Andrychowicz, A. Raichuk, P. Stańczyk, O. Bachem, L. Espeholt, M. Michalski *et al.*, "What matters in on-policy reinforcement learning? a large-scale empirical study," in *International Conference on Learning Representations (ICLR)*, 2021.
- [23] Z. Huang, Z. Zhang, and M. Q.-H. Meng, "A survey on reinforcement learning for robotic manipulation," *IEEE Transactions on Neural Networks and Learning Systems*, 2022, early access.
- [24] I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron *et al.*, "Solving rubik's cube with a robot hand," in *arXiv preprint arXiv:1910.07113*, 2019.
- [25] Y. Chen, Z. Liu *et al.*, "Learning-based robotic navigation in complex environments: A survey," *IEEE Access*, vol. 8, pp. 104 468–104 486, 2020.
- [26] H. Moon, J. Kim, and S. Lee, "Path planning for cleaning robots using reinforcement learning," *IEEE Access*, vol. 10, pp. 12 345–12 356, 2022.
- [27] C. Sun, R. Van Der Tol, R. Melenhorst, L. A. P. Pacheco, and P. G. Koerkamp, "Path planning of manure-robot cleaners using grid-based reinforcement learning," *Computers and Electronics in Agriculture*, vol. 226, p. 109456, Nov. 2024. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/S0168169924008470>
- [28] C. Hu, Q. Fu, and S. Yue, "Colias IV: The Affordable Micro Robot Platform with Bio-inspired Vision," in *Towards Autonomous Robotic Systems*. Springer, Cham, 2018, pp. 197–208, iSSN: 1611-3349. [Online]. Available: https://link.springer.com/chapter/10.1007/978-3-319-96728-8_17
- [29] O. Y. Maryasin, "Bee-Inspired Algorithm for Groups of Cyber-Physical Robotic Cleaners with Swarm Intelligence," in *Cyber-Physical Systems: Modelling and Intelligent Control*. Springer, Cham, 2021, pp. 167–177, iSSN: 2198-4190. [Online]. Available: https://link.springer.com/chapter/10.1007/978-3-030-66077-2_13
- [30] S. Siddiqua and T. Rahman, "Information-driven coordination for multi-robot systems using deep marl," *Robotics and Autonomous Systems*, 2024, in press.
- [31] S. Y. Luis, D. G. Reina, and S. L. T. Marín, "A Multiagent Deep Reinforcement Learning Approach for Path Planning in Autonomous Surface Vehicles: The Ypacará Lake Patrolling Case," *IEEE Access*, vol. 9, pp. 17 084–17 099, 2021. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/9330612>
- [32] A. M. Barrionuevo, S. Y. Luis, D. G. Reina, and S. L. T. Marín, "Optimizing Plastic Waste Collection in Water Bodies Using Heterogeneous Autonomous Surface Vehicles with Deep Reinforcement Learning," Dec. 2024, arXiv:2412.02316 [cs]. [Online]. Available: <http://arxiv.org/abs/2412.02316>
- [33] Y. Zhao, S. Ahmed, and T. Lin, "Dynamic goal assignment in multi-robot cleaning using shared occupancy maps," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 7021–7028.
- [34] A. Qureshi, M. Kim, and T. Moon, "Graph-based reinforcement learning for layout-agnostic patrolling," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [35] A. Gupta, M. J. van der Schoor, J. Bräutigam, V. B. Justo, T. F. Umland, and D. Göhlich, "Autonomous Service Robots for Urban Waste Management - Multiagent Route Planning and Cooperative Operation," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 8972–8979, Oct. 2022, conference Name: IEEE Robotics and Automation Letters. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/9816107>
- [36] J. Collenette and B. Logan, "Multi-agent Control of Industrial Robot Vacuum Cleaners," in *Engineering Multi-Agent Systems*. Springer, Cham, 2020, pp. 87–99, iSSN: 1611-3349. [Online]. Available: https://link.springer.com/chapter/10.1007/978-3-030-66534-0_6
- [37] S. Kim, K. Cho, and Y. Park, "Hierarchical multi-agent reinforcement learning for scalable indoor cleaning," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 8322–8328.

- [38] Z. Wang, K. Tan, and Y. Li, "Reconfigurable mop robots with decentralized pomdp-based control for adaptive cleaning," *IEEE Robotics and Automation Letters*, vol. 6, no. 4, pp. 7276–7283, 2021.
- [39] A. Din, M. Y. Ismail, B. Shah, M. Babar, F. Ali, and S. U. Baig, "A deep reinforcement learning-based multi-agent area coverage control for smart agriculture," *Computers and Electrical Engineering*, vol. 101, p. 108089, Jul. 2022. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/S0045790622003445>
- [40] R. Caccavale, M. Ermini, E. Fedeli, A. Finzi, V. Lippiello, and F. Tavano, "A multi-robot deep Q-learning framework for priority-based sanitization of railway stations," *Applied Intelligence*, vol. 53, no. 17, pp. 20 595–20 613, Sep. 2023. [Online]. Available: <https://doi.org/10.1007/s10489-023-04529-0>
- [41] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments," Mar. 2020, arXiv:1706.02275 [cs] version: 4. [Online]. Available: <http://arxiv.org/abs/1706.02275>



Mohamed Elmohandes received his B.Sc. in Computer Engineering from the American University of Sharjah. He is currently pursuing his Master's degree in Computer Engineering. His research interests include deep learning, robotics, and multi-agent systems.



Abdelrahman Ahmed (to be edited) Abdelrahman Ahmed received his B.Sc. in Computer Engineering from the American University of Sharjah. He is currently pursuing his Master's degree in Computer Engineering. His research interests include deep learning, app-development, and multi-agent systems.



Muhammed Noshin is currently pursuing an M.Sc. in Machine Learning at the American University of Sharjah, where he also earned his B.Sc. in Computer Engineering in 2023. His research comes at the intersection of machine learning and optimization, where he focuses on the design and analysis of scalable, low-complexity inference algorithms and on enabling real-time processing of large-scale, high-dimensional datasets in complex systems.