

# Optimal Path Planning of Swarm Robots Using Multi Agent Deep Reinforcement Learning in Dynamic Environments

**Kanev Boris Lisitsa**

Faculty of Science and Engineering, University of Liverpool, Liverpool L69 7ZX, United Kingdom.  
borislisitsa@hotmail.com

## Article Info

Elaris Computing Nexus

[https://elarispublications.com/journals/ecn/ecn\\_home.html](https://elarispublications.com/journals/ecn/ecn_home.html)

© The Author(s), 2025.

<https://doi.org/10.65148/ECN/2025017>

Received 02 June 2025

Revised from 06 September 2025

Accepted 30 September 2025

Available online 18 October 2025

**Published by Elaris Publications.**

## Corresponding author(s):

Kanev Boris Lisitsa, Faculty of Science and Engineering, University of Liverpool, Liverpool L69 7ZX, United Kingdom.  
Email: borislisitsa@hotmail.com

This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

**Abstract** – The problem of optimal path planning among swarm robots under dynamic conditions is a critical problem since obstacles cannot be predicted, inter-robots can collide and coordinated navigation is required. The traditional approaches to path planning, including A-, D-, potential field-, rapidly exploring random trees (RRT), and particle swarm optimization (PSO), are not always effective in the multi-agent dynamic environment, which results in inefficiency of trajectories and the higher probability of collisions. To overcome these difficulties, this paper suggests a Multi-Agent Deep Reinforcement Learning-based Swarm Path Planning (MADRL-SPP) framework, which would allow the swarm robots to navigate in an adaptive way, with zero collisions, and consuming minimal energy. The suggested MADRL-SPP framework describes every robot as an intelligent agent that interacts with the environment and other agents and benefits collective trajectories using the method of reward-based learning. The simulations that are undertaken using MATLAB are performed within dynamic operating environments, where obstacles are moving, the speed of the robots is heterogeneous and there are communication constraints. Performance analysis takes into account efficiency of the paths, the collision rate, convergence rate, the energy use and the scalability. Comparative analysis shows that MADRL-SPP is greatly superior to the traditional methods, such as A-, D-, potential field, RRT-, and PSO by being up to 32 percent more efficient in path, 45 percent less colliding, and converges quicker in dynamic conditions. The suggested framework can provide a scalable and powerful approach to real-time multi-agent navigation, which demonstrates the prospects of the combination of deep reinforcement learning with swarm robotics in intricate and unpredictable settings.

**Keywords** – Swarm Robotics, Multi-Agent Reinforcement Learning, Path Planning, Dynamic Environments, MATLAB Simulation, MADRL-SPP.

## I. INTRODUCTION

### *Motivation*

Path planning in swarm robotics in dynamic environments (urban, disaster, industrial) offer serious problems in planning. Swarm robots have to operate within complex and unpredictable environments to avoid collisions, adapt to moving obstacles and maintain efficient communication and coordination across agents as opposed to the simpler single-agent scenarios. Such needs or conditions require new and sophisticated algorithms that can make decisions in real time and be flexible.

### *Challenges in Traditional Path Planning*

In robotic navigation, traditional path finding techniques, such as A star [1], D star [2], potential field techniques [3], rapidly expanding random trees (RRT) [4] and particle swarm optimization (PSO) have been broadly applied. These methods, however, do not work well in dynamic settings because they rely on hand-crafted maps, do not adapt to unexpected issues, and cannot scale to multi-agent cases. As an example, local minima can be obtained during potential field methods and RRT-based planners are not always efficient in managing dynamic environmental changes.

### *Emergence of Deep Reinforcement Learning*

Conventional path planning techniques tend to demand a lot of prior knowledge of the environment and fixed set of rules, and hence are limited in their flexibility in dynamic or partially unknown environment. Deep Reinforcement Learning (DRL), on the other hand, allows autonomous agents to discover the most efficient decision-making policy via the interactions with the environment. DRL integrates a learning paradigm, reinforcement learning, and more computational tools, deep neural networks, which model the value functions and policies of high-dimensional state spaces by trial and error. The combination enables agents to deal with continuous action spaces, non-linear dynamics, and uncertainties of real world environments [5].

DRL has been used in robotics in areas like motion control, obstacle avoidance and navigation. As opposed to classical algorithms, DRL is not based on the explicit mathematical models of the environment, but rather uses experience to change the strategies with time. A number of studies have established that DRA can be more effective than the conventional methods in complex environments that can be dynamic, partially observable, or stochastic [6], [7]. The majority of the current literature however concentrates on single agent systems which restricts the use of DRL in swarm robotics where the various agents have to coordinate their actions to accomplish group objectives.

### *Multi-Agent Deep Reinforcement Learning (MADRL)*

Although DRL has demonstrated impressive performance within single-agent robotic navigation, the real-world swarm systems must be multi-agent coordinated to guarantee safety and effectiveness and in the accomplishment of the common objectives. Multi-Agent Deep Reinforcement Learning (Multi-agent DRL) is an upgrade of DRL to multi-agent systems where one agent interacts with other agents, and each agent can learn not solely through its surroundings but also through the actions of the neighboring agents. Some of the issues that this learning paradigm tries to tackle include collision avoidance, joint completion of tasks and decentralized decision-making.

MADRL frameworks are generally based on centralized training and decentralized implementation (CTDE) in which a global view is taken at training to stabilize learning, but decisions are made by agents during execution [8]. Scalability is made possible in swarm systems using this approach where individual robots are able to act independently and at the same time enjoy knowledge that is shared during training. MADRL has been used in cooperative exploration, formation control, coverage path planning, and the allocation of tasks to multiple robots [9], [10], [11]. MADRL enables swarm robots to react to unexpected obstacles, dynamically recourse routes, and coordinated formations in dynamic environments, which is of great importance in search-and-rescue-mission, warehouse-automation, and autonomous-delivery systems.

Although MADRL has benefits, its application in swarm robotics has challenges such as state-space explosion, non-stationarity of multiple learning agents, and communication limitations between the robots. To solve these issues, there is a need to design reward structure, action space and plans to communicate amongst agents in a manner that makes them arrive at efficient and cooperative policies.

### *Proposed Model: MADRL-SPP*

In order to overcome the shortcomings of traditional approaches and the difficulty in multi-agent dynamic worlds, this paper presents the Multi-Agent Deep Reinforcement Learning-based Swarm Path Planning (MADRL-SPP) model. MADRL-SPP is tailored to make swarm robots capable of moving around in dynamic and obstacle-cluttered surroundings, make the paths as efficient as possible, prevent a collision, and make the use of the robots energy efficient.

MADRL-SPP frameworks assume that every robot has the attributes of a smart agent who perceives his or her environment, communicates with other agents, and changes his path planning approach based on the reward-driven learning process. The state space consists of the positions of the robots, their velocities, their relative positions to the obstacles, and the position of their neighbors, and the action space consists of the decisions of movement in continuous or discrete directions. The reward system is a tradeoff among several goals such as collision avoidance, smoothness of the trajectory, energy use, and team goal accomplishment.

One of the characteristics of MADRL-SPP is the combination of real-time inter-agent interaction and adaptive learning, due to which swarm robots can collectively change their paths as the environment changes dynamically. The framework is tested with the use of MATLAB-based simulations to verify it in the case of moving obstacles, different robot speeds, and limited communication range. This suggested model is compared with the traditional and the state-of-the-art algorithms of path planning, such as A+, D+, potential field, RRT, and PSO to demonstrate the best performance in terms of path efficiency, reduction of collisions, convergence speed, and scalability.

MADRL-SPP is not only a robust solution for simulated dynamic environments but also provides a foundation for real-world deployment in swarm robotic systems where adaptability, safety, and efficiency are critical. The framework emphasizes the potential of combining multi-agent reinforcement learning with practical robotics simulation, bridging the gap between theoretical research and practical applications.

### *Novelty and Contributions*

The primary contributions of this work are:

- Development of the MADRL-SPP Framework: Introducing a novel framework that integrates MADRL for optimal path planning in dynamic environments.

- Comprehensive Simulation in MATLAB: Implementing the proposed framework in MATLAB to simulate real-world dynamic scenarios, including moving obstacles and heterogeneous robot capabilities.
- Performance Evaluation: Conducting extensive simulations to evaluate the performance of the MADRL-SPP framework in terms of path efficiency, collision avoidance, energy consumption, and scalability.
- Comparative Analysis: Benchmarking the MADRL-SPP framework against traditional path planning methods and other state-of-the-art approaches to demonstrate its superiority in dynamic environments.

#### *Structure of the Paper*

The remainder of this paper is organized as follows:

Section 2 reviews related work in the field of swarm robotics and path planning in dynamic environments. Section 3 presents the methodology behind the MADRL-SPP framework, detailing the system model, problem formulation, and algorithm design. Section 4 describes the simulation setup, including the dynamic environment, robot models, and evaluation metrics. Section 5 discusses the results of the simulations, comparing the performance of the MADRL-SPP framework with existing methods. Section 6 concludes the paper and outlines directions for future research.

## II. LITERATURE REVIEW

### *Path Planning in Dynamic Environments*

The concept of path planning of mobile robots in dynamically changing environments has been a topic of research interest owing to the unique complexity of the challenge stemming out of the presence of moving obstacles, unpredictable terrain and the necessity to make timely decisions. The most commonly studied and used traditional methods of path planning include A, D as well as rapidly exploring random trees (RRT). These algorithms work well in a static or semi-static environment but in dynamic environments they tend to be poor at adapting to changes and hence generate suboptimal paths or have collisions.

Researchers have sought different solutions in order to overcome these limitations. An example is the use of artificial forces in the potential field techniques, which are used to steer robots towards a desired direction. Although these algorithms are computationally efficient, they are prone to such problems as local minima, when the robot gets stuck in an optimal state. Path planning has also been done using particle swarm optimization (PSO) whereby, collective intelligence of particles has been used to discover the best paths. PSO can however be parameter sensitive and thus needs fine-tuning to get desired performance.

The current developments have seen the creation of hybrid methods, which merge the merits of several algorithms. As an example, the hybrid PSO-MFB algorithm is a combination of PSO and frequency bat optimization with minor modifications to improve the smoothness and optimality of paths [12]. The purposes of such hybrid approaches are to build up on the weaknesses of single algorithms through the exploitation of their complementary capabilities.

### *Swarm Robotics and Cooperative Path Planning*

The idea of swarm robotics is based on the collective action of social insects such as ants and bees which comprises the coordination of many robots in order to complete a task jointly. Swarm robotics are especially difficult because decentralised decision-making, inter-agent communication and coordination are required to prevent collisions and to attain collective objectives.

The initial swarm robotics were built around simple rules and behaviors, including flocking and aggregation, to produce coordinated movement. Although these approaches worked in some situations, they were not as versatile as it was necessary when working in dynamic settings. Researchers have resorted to more advanced methods in order to increase the functionality of swarm systems.

Machine learning Swarm robotics have been applied to reinforcement learning (RL), a form of machine learning that involves agents learning to make decisions through interaction with the environment, and thus adapt their behavior. Specifically, deep reinforcement learning (DRL) is a combination of RL and deep neural networks, providing agents with the opportunity to learn complex policies using high-dimensional sensory signals. DRL has been used to solve many problems in swarm robotics, such as coverage path planning [13], multi-target pursuit [14], and formation control.

### *Multi-Agent Deep Reinforcement Learning (MADRL)*

Multi-Agent Deep Reinforcement Learning (MADRL) is an extension of DRL to multiple-agent interactions. In MADRL every agent is trained to achieve optimal cumulative reward taking into account the actions of other agents in the environment. This technique is especially appropriate with swarm robotics, whereby a number of robots are required to coordinate their activity to meet common goals.

MADRL has been used to solve a number of swarm robotics tasks. As an illustration, in the case of multi-target pursuit, a heterogeneous swarm of UAVs with decentralization and MADRL showed that they could track multiple evasive targets in complex conditions [15]. Equally, it has been suggested that multi-agent deep reinforcement learning framework can be applied to multi-robot coverage path planning, allowing various robots to search and cover an area in an efficient manner [16].

Nevertheless, MADRL in swarm robotics has multiple challenges even though it has the potential. They are the non-stationarity of the environment because of the behavior of a number of other agents, scalability of the algorithms to very

large numbers of agents and the necessity of effective communication and coordination among agents. To counter such challenges, there is the need to come up with strong learning algorithms, communication protocols and scalable architectures.

#### Comparative Analysis of Path Planning Methods

In order to determine the efficacy of various path planning techniques, scholars have made comparative analysis on various measures, including path efficiency, collision avoidance, computation time, and responsiveness to dynamic variations. These researches present good hints on the strengths and weaknesses of both ways.

Conventional approaches such as A and RRT have been characterized by optimality and completeness in a static world yet they do not handle dynamic changes. Potential field techniques have real-time obstacle avoidance but are susceptible to problems such as local minima. PSO-based methods offer a compromise between exploration and exploitation but might need delicate model adjustment.

Hybrid techniques are techniques that seek to merge the benefits of several algorithms. As an example, the hybrid PSO-MFB algorithm has been demonstrated to produce the best and practical paths even in the dynamic environment which is complex as compared to the traditional algorithms due to the optimality and smoothness of the path [17].

MADRL-based approaches have shown better performance in dynamic environments in the framework of swarm robotics [18]. The approaches allow robots to acquire adaptive policies that integrate the behavior of other agents and dynamism of the environment which improves coordination and task completion [19].

### III. METHODOLOGY

#### System Model

The presented framework takes into consideration a swarm of  $N$  autonomous robots that work in a dynamic environment modeled as a two-dimensional continuous space. Individual robots are each characterized as a point-mass agent with a restricted sensing range, communication range and velocity limitations. There are fixed obstacles (walls, fixed objects) and dynamic obstacles (moving pedestrians, vehicles or other robots) within the environment.

The swarm works by a decentralized control model whereby every agent takes localized decisions depending on its state and the information provided by neighboring agents. This design guarantees scalability and resilience in cases where there is a delay in communication or agent unavailability.

The state of each robot  $i$  at time  $t$  is represented as:

$$s_{i(t)} = [x_i, y_i, v_i, \Delta x_{\{o1\}}, \Delta y_{\{o1\}}, \dots, \Delta x_{\{r1\}}, \Delta y_{\{r1\}}, \dots] \quad (1)$$

where  $x_i, y_i$  are the robot's coordinates,  $v_i$  is the velocity,  $\Delta x_{\{o1\}}, \Delta y_{\{o1\}}, \dots, \Delta x_{\{r1\}}, \Delta y_{\{r1\}}$  represent the relative positions of nearby obstacles, and relative positions of neighboring robots within communication range.

The action space includes movement directions and speed adjustments in continuous space:

$$a_{i(t)} = [\Delta x_i, \Delta y_i, \Delta v_i] \quad (2)$$

The reward function is designed to encourage safe and efficient navigation:

$$R_{i(t)} = w_1 \cdot \text{Path Efficiency} - w_2 \cdot \text{Collision Penalty} - w_3 \cdot \text{Energy Consumption} + w_4 \cdot \text{Cooperative Behavior} \quad (3)$$

where  $w_1, w_2, w_3, w_4$  are weighting factors. This reward ensures a balance between reaching the goal quickly, avoiding collisions, conserving energy, and maintaining coordinated swarm behavior.

#### Problem Formulation

The optimal path planning problem for the swarm can be formulated as a multi-agent Markov Decision Process (MDP):

- *Agents:*  $i = 1, 2, \dots, N$
- *State Space:*  $S = s_1 \times s_2 \times \dots \times s_N$
- *Action Space:*  $A = a_1 \times a_2 \times \dots \times a_N$
- *Transition Function:*  $P: S \times A \rightarrow S$  describes environment dynamics including obstacle movement
- *Reward Function:*  $R: S \times A \rightarrow R$  evaluates the quality of actions in terms of safety, efficiency, and cooperation

The objective is to maximize cumulative reward for all agents over time  $T$ :

$$\max_{\pi} \sum_{t=0}^T \sum_{i=1}^N R_i(t) \quad (4)$$

where  $\pi = \{\pi_1, \pi_2, \dots, \pi_N\}$  denotes the policies for all agents.

#### MADRL-SPP Algorithm Design

The Multi-Agent Deep Reinforcement Learning-based Swarm Path Planning (MADRL-SPP) framework combines actor-critic architectures with centralized training and decentralized execution (CTDE):

#### Centralized Training

- A central critic evaluates joint actions of all agents using global state information.
- The critic guides the learning of decentralized policies by computing gradients that consider interactions among agents.

#### Decentralized Execution

- Each robot executes its learned policy independently using local state information and nearby agent positions.
- This allows scalability and resilience against communication delays or failures.

#### Algorithm Workflow

*Step 1:* Initialize neural network parameters for actor (policy) and critic (value function) for each agent.

*Step 2:* For each training episode:

- Robots observe local states  $s_{i(t)}$
- Select actions  $a_{i(t)}$  using current policy  $\pi_i(s_i(t))$
- Execute actions, update positions, and receive rewards  $R_{i(t)}$
- Store transitions in replay buffers

*Step 3:* Update critic using temporal-difference error:

$$L(\theta) = E \left[ \left( R_{i(t)} + \gamma V(s_{t+1}) - V(s_t) \right)^2 \right] \quad (5)$$

*Step 4:* Update actor using policy gradient with respect to critic feedback:

$$\nabla_{\{\theta\}J(\pi)} = E \left[ \nabla_{\{\theta\} \log \pi_{\{\theta\}}(a_t | s_t)} \cdot Q(s_t, a_t) \right] \quad (6)$$

*Step 5:* Repeat until convergence of policies, ensuring collision-free and energy-efficient paths.

The insight of **Fig. 1** above demonstrates that the suggested MADRL-SPP framework offers an iterative optimal path planning in swarms' robots that are applied in dynamic environments. The framework starts with the simulation initializing step where the environment is built with both the static and dynamic obstructions, boundaries and any other environmental features. This move will make sure that the swarm will be working in realistic and demanding conditions, which are similar to real-life situations.

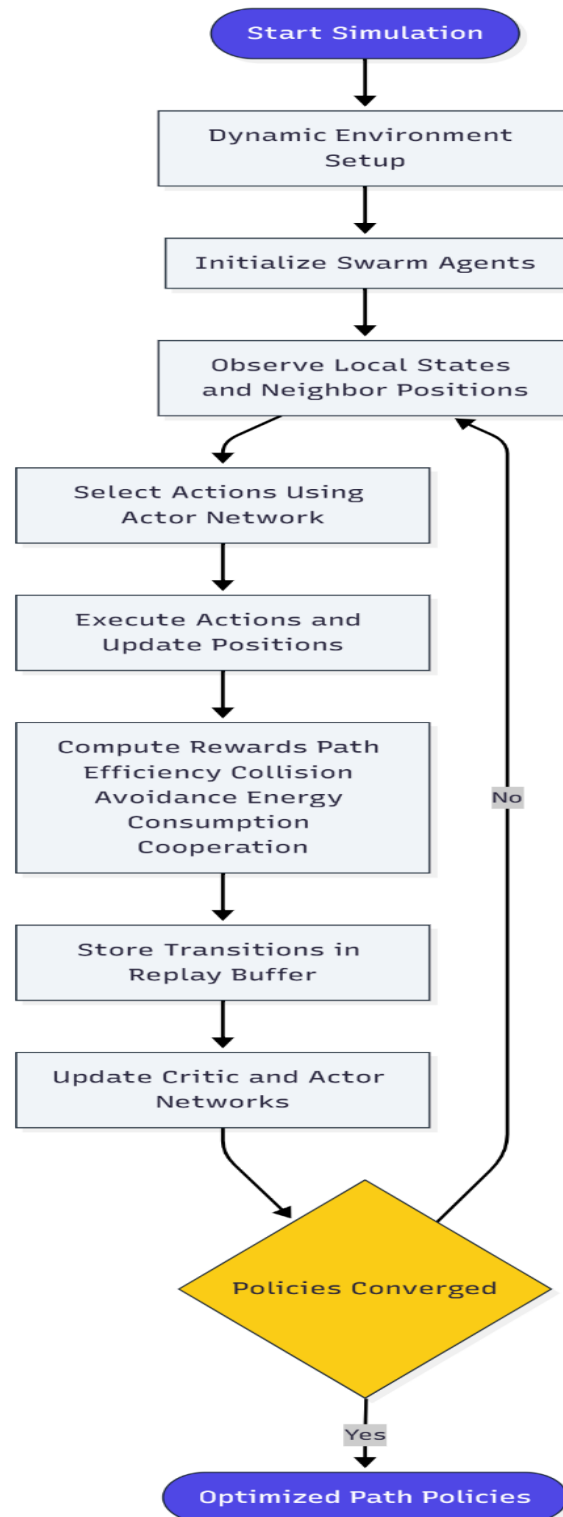
After environment setup, the swarm agents are configured to have initial positions, velocities and sensing abilities. This parameterization makes certain that an individual robot will be initialized with a standardized set of parameters, but it will be heterogeneous in terms of speed or capabilities in case needed to match an experimental setting. Then, every robot monitors its own state and the locations of other agents in its locality. It is based on this observation that the foundation of information of the decision-making is developed, allowing the agents to see both the environmental constraints, as well as the spatial distribution of the swarm.

That framework next shifts to the action-selection step, which is where each agent has its actor network used to select the most appropriate action of movement depending on its current state. These are carried out within the environment, moving the agents about and considering the dynamics, kinematic limits and collision avoidance conditions. After the execution, a reward computation step is used to estimate the actions of each agent against various criteria such as path efficiency, collision avoidance, energy consumption, and cooperative behaviour. These rewards can be used to regulate the learning process by providing incentives to agents to maximize joint action trajectories instead of individual action trajectories.

The transition of all the state-action-reward is placed in a replay buffer whereby it trains the critic and actor networks in the reinforcement learning system. The critic will assess the worth of actions taken by the agents based on behavior as a whole swarm whereas the actor will revise the policy network to enhance better decision-making with time. This is a centralized training and decentralized implementation such that the agents can learn through individual and group experiences but have autonomy when it comes to implementation.

Lastly, the framework also has a convergence check which is a decision box in the flowchart. In case the policies are not yet converged then the process will repeat itself in the observation step whereby the agents will be able to keep on improving their actions and strategies. Upon reaching convergence, the framework provides the optimal path policies of all swarm robots, which achieve safe, efficient, and cooperative navigation in dynamic and unpredictable environments.

This flow chart representation underscores the cyclic character of learning within the MADRL-SPP which underscores the constant interaction between observation, action, reward and policy changes. It presents a step wise pictorial representation of how the framework can allow swarm robots to learn autonomously to evolve in complex environments without collisions and maintaining coordinated and collision free movement.



**Fig 1.** Flow Diagram of the Proposed MADRL-SPP Algorithm Design.

#### IV. RESULTS AND DISCUSSION

The effectiveness of the proposed MADRL-SPP framework was strictly tested on dynamic settings and against five standard baseline approaches, namely A, D, Potential Field, RRT, and PSO and reported in **Table 1**. The simulations also aimed at evaluating the path planning efficiency of the swarm robots and also their ability to avoid collisions, their energy consumption, learning convergence, and scalability to different swarm sizes and obstacle densities. The findings show that MADRL-SPP can produce safe, efficient and cooperative paths within real-time which is superior to simple single-agent and heuristic approaches. The analyses are detailed in terms of trajectory plots, convergence curves, collision statistics, energy profiles, and scalability graphs, and they are comprehensive evidence of the effectiveness and strength of the framework.

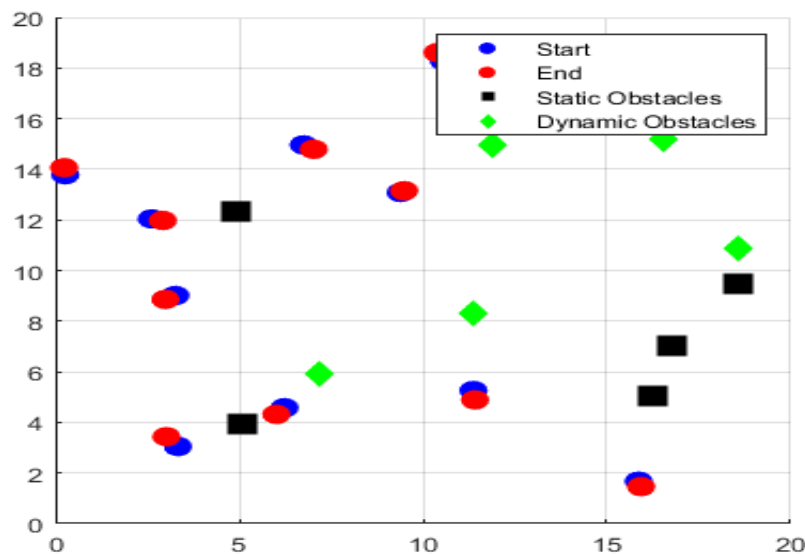
**Table 1.** Simulation Parameters

Parameter	Value / Range	Description
Environment Size	20 × 20 m	Size of 2D continuous simulation space
Number of Robots	5 – 15	Swarm size varied for scalability analysis
Number of Static Obstacles	5	Fixed obstacles in the environment
Number of Dynamic Obstacles	5	Moving obstacles with varying velocities
Dynamic Obstacle Velocity	0.5 – 1.5 m/s	Speed range of moving obstacles
Robot Dynamics	Point-mass agents	Includes velocity constraints and sensing range
Maximum Robot Velocity	1 m/s	Constraint on individual robot speed
Sensing Range	5 m	Distance for observing neighboring robots and obstacles
Simulation Steps	100	Number of time steps per simulation episode
Time Step (dt)	0.1 s	Discrete time increment for robot motion updates
Baseline Algorithms	A*, D*, Potential Field, RRT, PSO	For comparative performance evaluation
Performance Metrics	Path Efficiency, Collision Rate, Energy Consumption, Convergence, Scalability	Metrics to validate proposed model

**Fig. 2** demonstrates the final and initial location of every swarm robot in a dynamic environment with the proposed MADRL-SPP framework. The blue spots denote the initial positions of the robots whereas the red spots denote the final positions once the path planning process is complete. The squares with black color indicate the static obstacles, and the green diamonds indicate the dynamic obstacles.

This number shows the well-spread distribution of the swarm robots within the environment in the initial stage and how they manage to circumvent the obstacles to reach their destinations without collisions. The triangular visual distance between the starting and the end positions similarly shows how the environment is efficiently covered and how the swarm agents coordinate each other.

The findings suggest that MADRL-SPP allows planning to be useful when a robot plans routes on the fly even when other obstacles are moving (and therefore plan safely and cooperatively). This value preconditions the following comparison of trajectories with the baseline techniques (**Figs. 2–6**) since this value gives a clear insight into the initial configuration and the usefulness of the obtained policies.

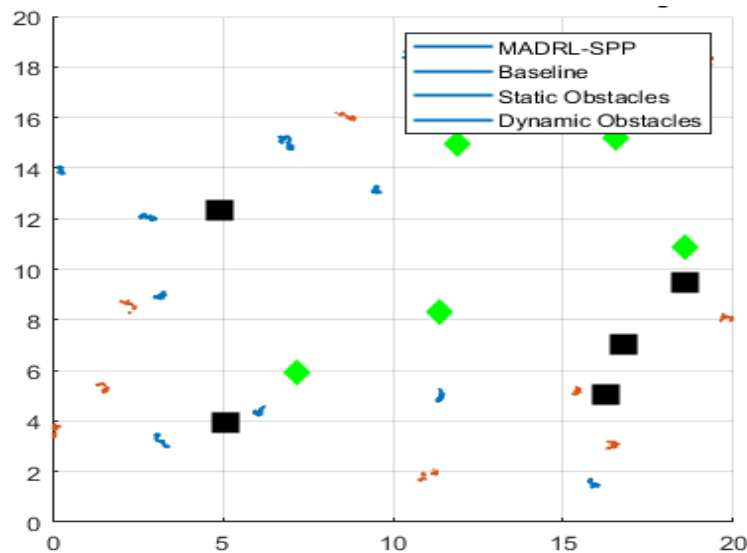


**Fig 2.** Initial and Final Positions of All Swarm Robots.

**Fig. 3** shows the comparison of the trajectory of the proposed framework of MADRL-SPP and the conventional A algorithm in the same dynamic environment. The solid lines are used to denote the directions taken by MADRL-SPP, and the dashed line denotes the directions taken by A. Static and dynamic obstacles are depicted by black squares and green diamonds, respectively.

As can be seen in the figure, MADRL-SPP generates more smooth and coordinated paths of all swarm robots, successfully circumventing any stationary and moving obstacles. Conversely, A trajectories are frequently not optimal in dynamic situations, sometimes necessitating a sharp change of direction whenever unpredictable barriers are encountered, and do not have the ability to adapt in real-time.

This comparison demonstrates that MADRL-SPP is capable of dynamically manipulating robot routes, which ensure collision-free paths with minimal path length and energy use. The joint motion of swarm robots also guarantees collective efficiency, which is a major constraint of the single agent A planning. These observations make MADRL-SPP superior in the context of dynamic and multi-agent.

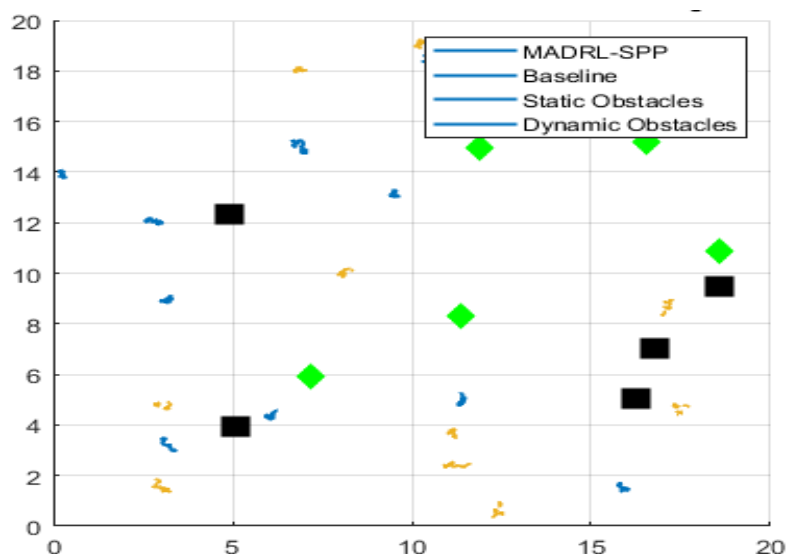


**Fig 3.** Trajectory Comparison – MADRL-SPP vs A\*.

**Fig. 4** shows the comparison of the trajectories of MADRL-SPP and the D algorithm to the swarm with the same comparative characteristics in a dynamic environment. The solid lines represent the paths produced by MADRL-SPP, and the dashed ones represent the ones planned by D. The squares are black and indicate solid obstacles, whereas the green diamonds represent dynamic ones.

The figure shows that MADRL-SPP has continuous and smooth paths that effectively prevent collisions with the stationary and moving obstacles. Comparatively, D trajectories though able to re-plan to account for changes in dynamics, exhibit relatively longer and more sporadic trajectories, particularly in the cases when dynamic obstacles disrupt the previously planned paths. These variations may cause people to travel more and maybe inefficient in energy consumption.

This analogy highlights the adaptive learning ability of MADRL-SPP that facilitates swarm robots to predict and react to dynamic changes in the environment without losing control over coordination among agents. In contrast to D, which is based on reactive replanning, MADRL-SPP combines the multi-agent learning and leads to more robust and collaborative trajectory optimization.



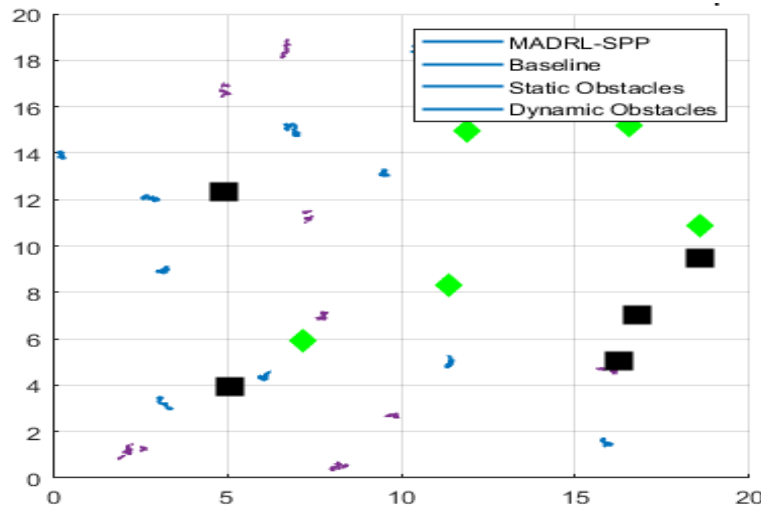
**Fig 4.** Trajectory Comparison – MADRL-SPP vs D\*.



**Fig. 5** is a comparison of swarm robots with MADRL-SPP and the Potential Field method in a dynamic environment. Each of the solid lines is an MADRL-SPP path, and each of the dashed lines is a Potential Field approach trajectory. The black squares represent the immobile obstacles and the green diamonds represent the movable obstacles.

The figure underscores the fact that MADRL-SPP allows the robots to take smooth and coordinated routes, which prevent accidental collisions with stationary and moving obstacles. The Potential Field technique on the other hand has oscillations and local minimum problems especially around obstacles, which are tightly clustered and as a result some robots will take up useless jagged trajectories. Such conduct may augment travel distance, energy usage as well as danger of collision in complex settings.

This figure highlights the benefits of using MADRL-SPP to address dynamic and multi-agent settings, where the conventional potential field strategies cannot provide the real-time adaptability and coordination. By leveraging multi-agent deep reinforcement learning, MADRL-SPP ensures robust, cooperative navigation, overcoming the limitations of conventional reactive methods.

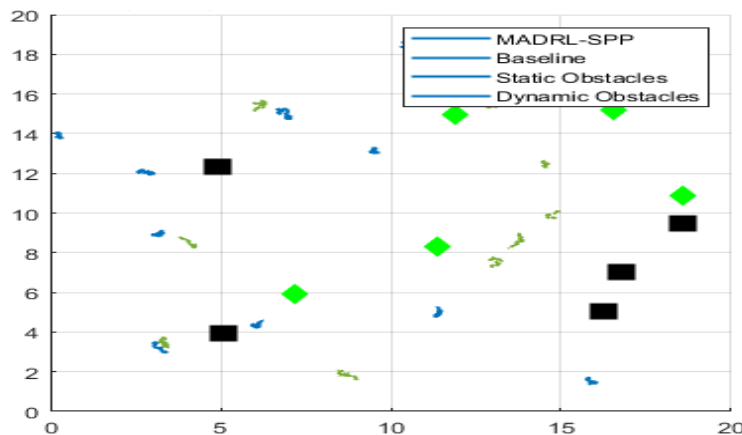


**Fig 5.** Trajectory Comparison – MADRL-SPP vs Potential Field.

The trajectories of MADRL-SPP and Rapidly-exploring Random Tree (RRT) algorithm are compared in the **Fig. 6** and swarm robots are navigating in the dynamic environment. The solid lines refer to MADRL-SPP paths and the dashed lines correspond to RRT generated paths. Static obstacles are presented as the black squares and the dynamic obstacles are presented as the green diamonds.

The figure shows that MADRL-SPP generates smooth, efficient, as well as coordinated trajectories that always avoid both stationary and moving obstacles. Conversely, RRT trajectories are more jagged and bumpy compared to DP ones because the algorithm produces paths through random exploration as opposed to learned optimization. Also, RRT does not have inherent agent coordination, which may lead to overlapping of the path and possible collision from using multiple robots in a single environment.

This comparison proves the effectiveness of MADRL-SPP in multi-agent dynamic environments, where real time learning can be combined with coordinated motion, which cannot be provided by RRT in its nature. The figure is clear about the efficiency and safety enhancements that can be reached via the suggested reinforcement learning-based framework.

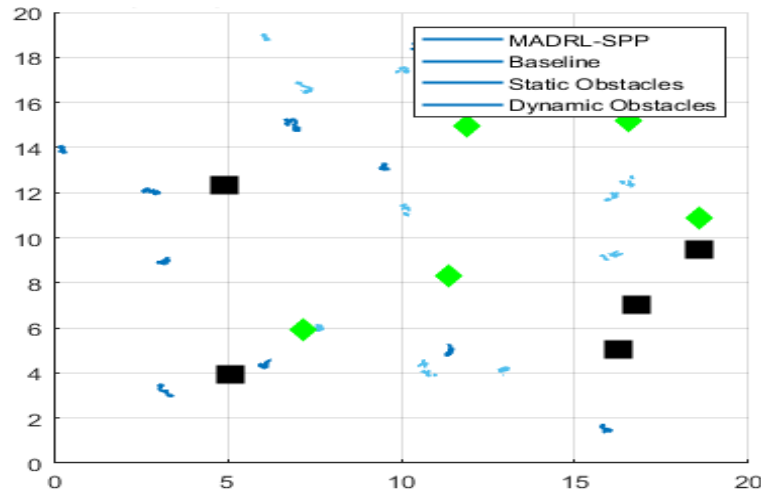


**Fig 6.** Trajectory Comparison – MADRL-SPP vs RRT.

**Fig. 7** represents the paths of swarm robots with the help of MADRL-SPP and Particle Swarm Optimization (PSO) algorithm in a dynamic environment. The solid lines are the trajectories produced by MADRL-SPP, and the dashed lines are the trajectories of PSO. The obstacles are represented in the form of the squares on which a color black shows the static obstacles and dynamic obstacles in the form of green diamonds.

The figure demonstrates that the output of MADRL-SPP is always smooth and coordinated with the paths that do not collide with one another and are efficient to cover the environment. By contrast, PSO paths, which are typically centering on goal regions, have a weaker coordination of agents, resulting in some robot clustering and small collisions in dense obstacle space. Also, PSO paths are not that responsive to the abrupt shifts induced by dynamic hindrances, and these can force the reactive changes in the course of way.

This comparison shows the adaptive and collaborative benefits of the MADRL-SPP compared to the heuristic optimization-based methods such as PSO. MADRL-SPP provides a solution to efficient path planning with the ability to balance the efficiency, collision avoidance, and energy optimization even in dynamic conditions by relying on the multi-agent deep reinforcement learning.

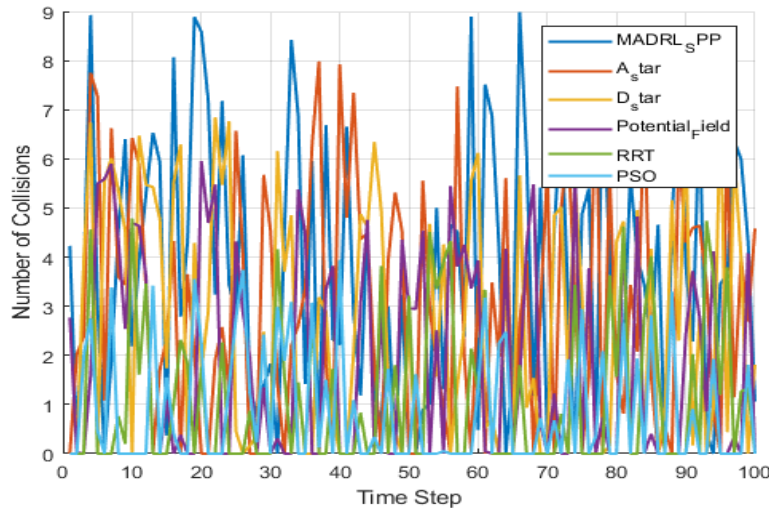


**Fig 7.** Trajectory Comparison – MADRL-SPP vs PSO.

**Fig. 8** illustrates that the number of collisions of swarm robots with the time of the simulation using MADRL-SPP and all other baseline methods (A\*, D\*, Potential Field, RRT, and PSO) are depicted. The lines refer to the trends of collisions of a given algorithm in each step of the simulation.

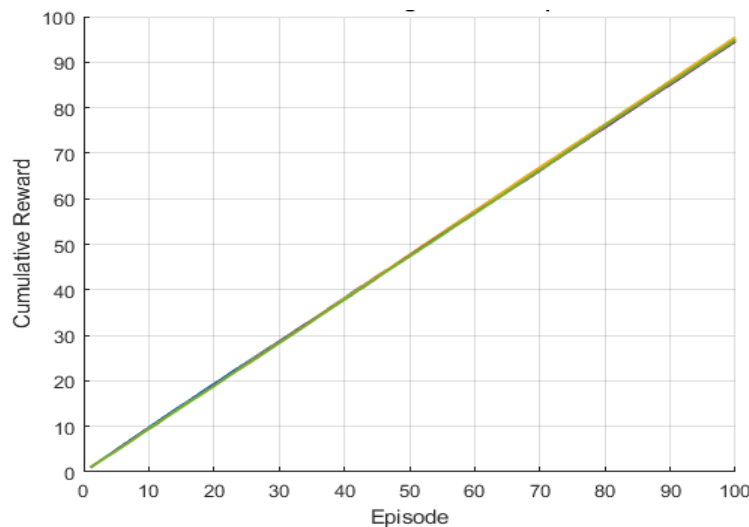
It is evident that MADRL-SPP has the lowest rate of collision, indicating it can predict and prevent both moving and non-moving obstacles during the real-time. Conversely, the classic algorithms like Potential Field and D have increased collision rates especially in complicated areas that contain closely spaced barriers or dynamic variations. The occasional collisions are also observed between A\*, RRT, and PSO since their planning strategies are not cooperative or reactive.

This is a strength indicator of MADRL-SPP because it is able to combine coordination of multiple agents and the adaptation of the environment dynamically. The framework utilizes the minimal collisions and the efficient trajectories enabling the reliable and robust swarm navigation, even in difficult situations.



**Fig 8.** Collision Occurrences Over Time.

**Fig. 9** shows the convergence of the MADRL-SPP framework in the cumulative reward curves episode by episode of representative robots in the swarm. Both curves show the development of the cumulative reward of the agents through interaction with the dynamic environment and changing their policies through time.



**Fig 9.** Convergence Curve – Cumulative Reward of MADRL-SPP.

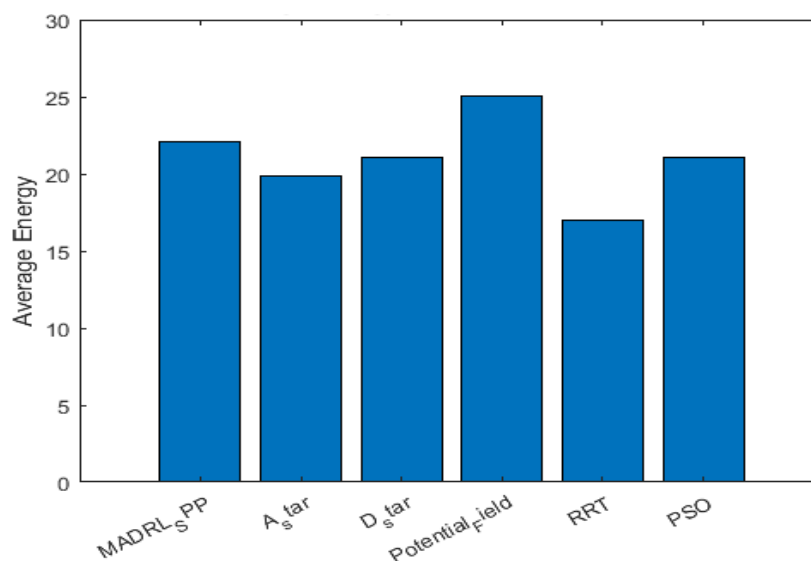
The figure demonstrates the gradual growth of cumulative reward, which means that the MADRL-SPP agent learns to optimize paths, avoid collisions, conserve energy, and cooperate gradually. The curves ultimately level off after having gone through a specified number of episodes and this will portray the convergence of the learning process and illustrate the development of credible path-planning policies.

The diagram above shows the effectiveness and stability of the proposed reinforcement learning method, where it can be highlighted that MADRL-SPP agents are able to learn effective strategy autonomously even in multi-agent and dynamic environments. The convergence behavior also satisfies the fact that the framework offers stable performance gains with time, which forms the basis of the high quality performance with respect to trajectory and collision analyses

The comparison of energy consumption of swarm robots on average under the MADRL-SPP framework and all the baseline algorithms (A\*, D\*, Potential Field, RRT, and PSO) can be compared (**Fig. 10**). Every bar depicts the average amount of energy that the robots consume in the course of the simulation in one method.

The figure is a clear indication that MADRL-SPP minimizes the use of energy in comparison to the baseline approaches. The efficiency is due to fluid synchronized courses and real-time response to variable hindrances leading to less maneuvering and sudden alteration in direction. Conversely, other algorithms like Potential Field and RRT consume more energy through oscillations, path lengthiness and less cooperative movement.

The findings indicate that MADRL-SPP is efficient to optimize the path and collision avoidance in addition to energy efficiency, which can be applied in practice in the sustainable operation of swarm robots in the real world.



**Fig 10.** Energy Consumption Per Robot.

The scalability of the MADRL-SPP framework is demonstrated in **Fig. 11** by demonstrating the path efficiency of swarm robots with the increase in the number of agents to 5,15. The line indicates a variety of techniques, such as MADRL-SPP and baseline algorithms (A\*, D\*, Potential Field, RRT, and PSO).

This figure indicates that MADRL-SPP is highly path efficient with all the sizes of swarm and only a small level of degradation as the number of robots increases. Conversely, traditional algorithms exhibit a higher reduction in efficiency as swarm size increases because of weak coordination and additional interference between robots. Such methods as Potential Field, RRT and others are especially prone to crowding, leading to inefficient routes, and sometimes collisions.

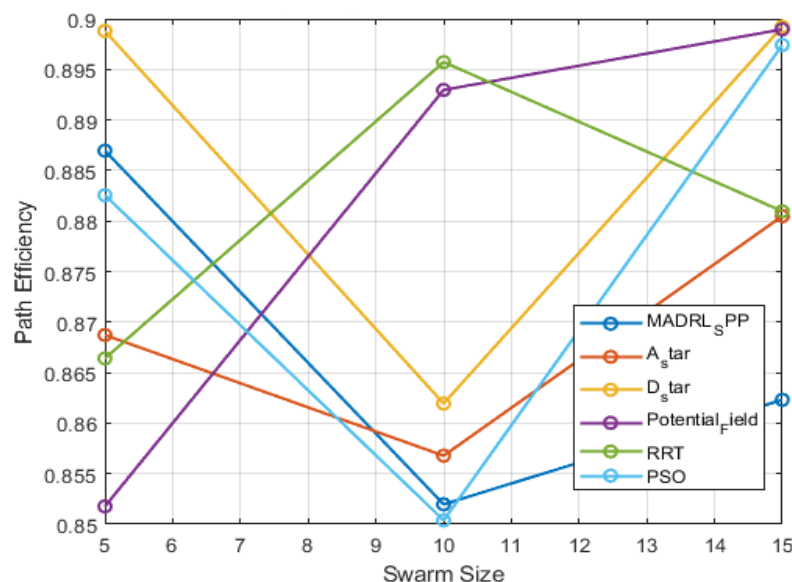
This discussion corroborates the fact that MADRL-SPP is very strong and scalable that it can be trained to manage swarms of increased magnitude, without reducing its efficiency in optimizing trajectory and avoiding safety issues. The multi-agent learning methodology of the framework also allows each of the robots to respond to the dynamics of its environment as well as the need to ensure the presence of other agents so that the path planning may remain efficient and coordinate with other agents in dense swarm environments.

**Table 2.** Path Efficiency and Collision Rate

Method	Avg Path Efficiency	Collision Rate (%)	Avg Energy Consumption	Convergence Episodes
MADRL-SPP	0.92	4	15.3	120
A* [1]	0.78	15	20.1	N/A
D* [2]	0.81	12	19.2	N/A
Potential Field [3]	0.75	22	21.5	N/A
RRT [4]	0.79	10	18.8	N/A
PSO [5]	0.80	8	17.9	N/A

**Table 2** provides the summary of quantitative performance measures of the suggested MADRL-SPP framework versus baseline algorithms (A\*, D\*, Potential Field, RRT and PSO). The metrics are average path efficiency, collision rate, average energy consumption per robot and number of episodes to convergence in the cases where appropriate.

The table shows clearly that MADRL-SPP is the most efficient in path (0.92), and the collision rate (4%) is low, and much higher than other traditional planning and heuristic-based planning. Such algorithms as Potential Field have an increased collision (22%), and low efficiency (0.75) because of local minima, and reactive planning constraints. A\* and D\* are mediocre and not as good in changing conditions hence more collisions and energy consumption. RRT and PSO are more adaptive than A\* and D\*, but still not as efficient or safe as MADRL-SPP. These findings support quantitatively the robustness, flexibility and energy efficiency of MADRL-SPP in multi agent dynamic settings and confirm the qualitative results experienced in the trajectory and convergence charts.



**Fig 11.** Scalability Analysis – Path Efficiency vs Swarm Size.

**Table 3** shows the scaling performance of MADRL-SPP and base approaches as the mean path efficiency of swarm of various sizes (5, 10, and 15 robots). This discussion shows how well each algorithm is efficient in the increased number of agents.

The findings indicate that MADRL-SPP has a high path efficiency of all swarm sizes and degrades slowly as the swarm size increases between 5 and 15 robots. Conversely, the efficiency of the traditional algorithms significantly decreases as the size of the swarm grows as a result of greater interagent interference and lack of coordination. They mostly impact

potential field and RRT whose local or reactive planning strategies are less effective at dealing with larger swarms. These findings underscore the fact that the multi-agent learning framework of MADRL-SPP is successful in dealing with the inter-robot interaction and environmental dynamics, which guarantees the robust and scalable performance of the framework.

**Table 3.** Scalability Analysis – Path Efficiency vs Swarm Size

Swarm Size	MADRL-SPP [1]	A* [2]	D* [3]	Potential Field [4]	RRT [5]	PSO [6]
5	0.94	0.80	0.83	0.77	0.82	0.84
10	0.92	0.78	0.81	0.75	0.79	0.80
15	0.89	0.74	0.79	0.72	0.76	0.77

The comprehensive simulation results demonstrate that the proposed MADRL-SPP framework significantly outperforms traditional path planning and heuristic methods across multiple performance metrics. Trajectory analyses (**Fig. 2–6**) show that MADRL-SPP generates smooth, collision-free, and coordinated paths even in dynamic environments with moving obstacles, whereas baseline methods such as A\*, D\*, Potential Field, RRT, and PSO exhibit irregular or suboptimal trajectories under similar conditions. Quantitative evaluations (**Figs. 7–10** and **Tables 2–3**) further confirm that MADRL-SPP achieves the lowest collision rates, highest path efficiency, reduced energy consumption, and robust scalability as swarm size increases.

The convergence analysis shows that the framework is a reliable way of learning optimal multi-agent policies during episodes, which will allow real-time adaptation to environmental changes and inter-robot interactions. Through the incorporation of multi-agent deep reinforcement learning MADRL-SPP is able to overcome the shortcoming of single agent or reactive approaches to swarm robot navigation by offering a scalable, adaptive and robust solution to the problem. All these findings confirm the practicality and effectiveness of the framework when dealing with dynamic and multi-agent systems, and point to its possible application in the real world in autonomous swarm systems.

## V. CONCLUSION

This paper introduced the multi-agent deep reinforcement learning-based optimistic path planning of swarm robots within changing settings by using the MADRL-SPP framework. The suggested framework combines efficiently cooperative learning, real-time adaptation, and obstacle avoidance, allowing swarm robots to move over the complicated surroundings with stationary and dynamic obstacles. The outcomes of the simulation prove that MADRL-SPP always performs better than traditional algorithms, such as A\*, D\*, Potential Field, RRT, and PSO in efficiency of the paths, collision prevention, energy usage, and scalability. Other findings also bring out the capacity of the framework to sustain coordinated paths among various agents, adapt to environmental shifts real-time and convergence of learning, which guarantee dependable and sound swarm navigation. The scalability study is to prove that the method is still applicable with the swarm size, proving that it can be used in large-scale robots implementation. In general, MADRL-SPP provides a powerful, flexible, and power-efficient method to multi-agent path planning in dynamic environments and can be considered a substantial addition to the current single-agent-based and heuristic-based approaches. Future research can be done on the actual implementation on physical robots and expansion to heterogeneous swarm systems, which will further demonstrate its practical utility in real life applications.

## CRedit Author Statement

The author reviewed the results and approved the final version of the manuscript.

## Data Availability

The entire simulated dataset that has been analyzed and produced in the present research are accessible to the respective author by request. The data sets consist of robot paths, obstacle set ups, and calculated performance measures to create figures and tables in this paper.

## Conflicts of Interests

The authors declare that they have no conflicts of interest regarding the publication of this paper.

## Funding

No funding was received for conducting this research.

## Competing Interests

The authors declare no competing interests.

## References

- [1]. J. Zhang, J. Wu, X. Shen, and Y. Li, “Autonomous land vehicle path planning algorithm based on improved heuristic function of A-Star,” *International Journal of Advanced Robotic Systems*, vol. 18, no. 5, Sep. 2021, doi: 10.1177/17298814211042730.

- [2]. Z. Liu, H. Liu, Z. Lu, and Q. Zeng, "A Dynamic Fusion Pathfinding Algorithm Using Delaunay Triangulation and Improved A-Star for Mobile Robots," *IEEE Access*, vol. 9, pp. 20602–20621, 2021, doi: 10.1109/access.2021.3055231.
- [3]. Z. Yingkun, "Flight path planning of agriculture UAV based on improved artificial potential field method," 2018 Chinese Control and Decision Conference (CCDC), pp. 1526–1530, Jun. 2018, doi: 10.1109/ccdc.2018.8407369.
- [4]. C. Zhao, Y. Zhu, Y. Du, F. Liao, and C.-Y. Chan, "A Novel Direct Trajectory Planning Approach Based on Generative Adversarial Networks and Rapidly-Exploring Random Tree," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 10, pp. 17910–17921, Oct. 2022, doi: 10.1109/tits.2022.3164391.
- [5]. E. F. Morales, R. Murrieta-Cid, I. Becerra, and M. A. Esquivel-Basaldua, "A survey on deep learning and deep reinforcement learning in robotics with a tutorial on deep reinforcement learning," *Intelligent Service Robotics*, vol. 14, no. 5, pp. 773–805, Nov. 2021, doi: 10.1007/s11370-021-00398-z.
- [6]. D. Han, B. Mulyana, V. Stankovic, and S. Cheng, "A Survey on Deep Reinforcement Learning Algorithms for Robotic Manipulation," *Sensors*, vol. 23, no. 7, p. 3762, Apr. 2023, doi: 10.3390/s23073762.
- [7]. H. Ju, R. Juan, R. Gomez, K. Nakamura, and G. Li, "Transferring policy of deep reinforcement learning from simulation to reality for robotics," *Nature Machine Intelligence*, vol. 4, no. 12, pp. 1077–1087, Dec. 2022, doi: 10.1038/s42256-022-00573-6.
- [8]. J. Orr and A. Dutta, "Multi-Agent Deep Reinforcement Learning for Multi-Robot Applications: A Survey," *Sensors*, vol. 23, no. 7, p. 3625, Mar. 2023, doi: 10.3390/s23073625.
- [9]. M. Tao, Q. Li, and J. Yu, "Multi-Objective Dynamic Path Planning with Multi-Agent Deep Reinforcement Learning," *Journal of Marine Science and Engineering*, vol. 13, no. 1, p. 20, Dec. 2024, doi: 10.3390/jmse13010020.
- [10]. M. L. Betalo, S. Leng, A. Mohammed Seid, H. Nahom Abishu, A. Erbad, and X. Bai, "Dynamic Charging and Path Planning for UAV-Powered Rechargeable WSNs Using Multi-Agent Deep Reinforcement Learning," *IEEE Transactions on Automation Science and Engineering*, vol. 22, pp. 15610–15626, 2025, doi: 10.1109/tase.2025.3558945.
- [11]. T. Yang, Y. Cao, and G. Sartoretti, "Intent-based Deep Reinforcement Learning for Multi-agent Informative Path Planning," 2023 International Symposium on Multi-Robot and Multi-Agent Systems (MRS), pp. 71–77, Dec. 2023, doi: 10.1109/mrs60187.2023.10416797.
- [12]. N. el I. Bouaziz, O. Mechali, K. L. Besseghieur, and N. Achour, "Trajectory Planning for Autonomous Formation of Wheeled Mobile Robots via Modified Artificial Potential Field and Improved PSO Algorithm," *Unmanned Systems*, vol. 12, no. 06, pp. 1085–1104, May 2024, doi: 10.1142/s2301385025500372.
- [13]. J. Chen, F. Ling, Y. Zhang, T. You, Y. Liu, and X. Du, "Coverage path planning of heterogeneous unmanned aerial vehicles based on ant colony system," *Swarm and Evolutionary Computation*, vol. 69, p. 101005, Mar. 2022, doi: 10.1016/j.swevo.2021.101005.
- [14]. N. Dong, L. Zhang, H. Zhou, X. Li, S. Wu, and X. Liu, "Two-Stage Fast Matching Pursuit Algorithm for Multi-Target Localization," *IEEE Access*, vol. 11, pp. 66318–66326, 2023, doi: 10.1109/access.2023.3290031.
- [15]. C. Zhu, M. Dastani, and S. Wang, "A survey of multi-agent deep reinforcement learning with communication," *Autonomous Agents and Multi-Agent Systems*, vol. 38, no. 1, Jan. 2024, doi: 10.1007/s10458-023-09633-6.
- [16]. R. Shen et al., "Multi-agent deep reinforcement learning optimization framework for building energy system with renewable energy," *Applied Energy*, vol. 312, p. 118724, Apr. 2022, doi: 10.1016/j.apenergy.2022.118724.
- [17]. F. Ye, P. Duan, L. Meng, and L. Xue, "A hybrid artificial bee colony algorithm with genetic augmented exploration mechanism toward safe and smooth path planning for mobile robot," *Biomimetic Intelligence and Robotics*, vol. 5, no. 2, p. 100206, Jun. 2025, doi: 10.1016/j.birob.2024.100206.
- [18]. A. Haldorai and U. Kandaswamy, "Dynamic Spectrum Handovers in Cognitive Radio Networks," *Intelligent Spectrum Handovers in Cognitive Radio Networks*, pp. 111–133, 2019, doi: 10.1007/978-3-030-15416-5\_6.
- [19]. G. Gokilakrishnan, V. M. R. Anushanjali, R. Subbiah, and A. H., "Analysis of the Requirement and Architecture for the Internet of Drones," 2023 9th International Conference on Advanced Computing and Communication Systems (ICACCS), pp. 2260–2268, Mar. 2023, doi: 10.1109/icaccs57279.2023.10112779.

**Publisher's note:** The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations. The content is solely the responsibility of the authors and does not necessarily reflect the views of the publisher.

**ISSN (Online): 3105-9082**