# Safe Multi-Agent Reinforcement Learning
# for Behavior-Based Cooperative Navigation

Murad Dawood    Sicong Pan    Nils Dengler    Siqi Zhou    Angela P. Schoellig    Maren Bennewitz

*Abstract*— In this paper, we address the problem of behavior-based cooperative navigation of mobile robots using safe multi-agent reinforcement learning (MARL). Our work is the first to focus on cooperative navigation without individual reference targets for the robots, using a single target for the formation's centroid. This eliminates the complexities involved in having several path planners to control a team of robots. To ensure safety, our MARL framework uses model predictive control (MPC) to prevent actions that could lead to collisions during training and execution. We demonstrate the effectiveness of our method in simulation and on real robots, achieving safe behavior-based cooperative navigation without using individual reference targets, with zero collisions, and faster target reaching compared to baselines. Finally, we study the impact of MPC safety filters on the learning process, revealing that we achieve faster convergence during training and we show that our approach can be safely deployed on real robots, even during early stages of the training.

## I. INTRODUCTION

Cooperative navigation of unmanned vehicles has gained considerable attention due to its applications in missions such as search and rescue [1], surveillance [2], and payload transfers [3]. As outlined in the literature [4], cooperative navigation can be achieved through various strategies, including leader-follower, virtual-structures, and behavior-based approaches. Behavior-based methods [5], [6] offer more flexibility, as the robots' actions adapt to their own observations. The goal in behavior-based cooperative navigation is to maintain relative distances, avoid collisions, and jointly reach target locations.

In this work, we focus on safe learning for behavior-based cooperative navigation. Safety in this context means eliminating all collisions during both training and execution, achieved through safe reinforcement learning (RL). RL has been successful in cooperative navigation [7]–[9], but challenges remain due to the sim-to-real gap. This gap arises from real-world sensor noise and unmodeled nonlinearities in simulations. Unlike approaches that focus on improving simulation realism [10], [11], we emphasize the need to ensure RL is safe for real-world deployment. Although safety in RL
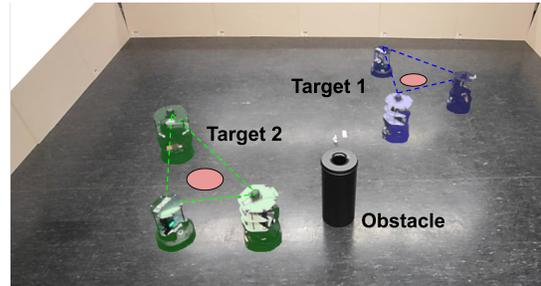
**Fig. 1:** Real-world example for the behavior-based cooperative navigation control. The robots start from random locations and navigate cooperatively to reach the targets for the centroid of the formation (shown in red) while aiming to maintain the predefined distances with respect to each other. The blue and green shades show the robots team at the first and second goals, respectively.

has advanced [12]–[14], it remains underexplored in behavior-based navigation. Furthermore, the impact of safety filters on RL training has not been fully studied. In this work, we investigate the effects of applying a model predictive control (MPC)-based safety filter to RL in the context of behavior-based cooperative navigation.

To facilitate the deployment of behavior-based cooperative navigation and enable the scalability of our approach, i.e., the application of the learned policy to a team with a higher number of robots without retraining, we use a reference target for the centroid of the formation. This modification is promising for coordinating teams of robots as in the StarCraft multi-agent challenge [15]. To maintain the distances between the robots, each robot considers the distances to its two neighbors. Optimization-based controllers struggle in solving this task, since these controllers require a reference target location for each individual robot and might even assume that each robot has access to all its neighbors positions [5], [6], [16]–[21]. Therefore, we use multi-agent reinforcement learning (MARL) to achieve the desired behavior. Figure 1 shows a real-world example of our approach, where the robots get into formation to reach the target locations for the centroid while avoiding unsafe actions.

The primary contributions of this paper can be summarized as follows: **(i) Safe learning of behavior-based cooperative navigation.** We introduce the application of safe RL to behavior-based navigation. This achievement is notable because it ensures the agents' safety during the training phase, a previously unaddressed challenge in behavior-based cooperative navigation, as well as during the execution phase. Additionally, we studied the effect of the safety layer on training efficiency, showing improved convergence compared to baselines. **(ii) Eliminating the need for individual path**

**planners for each robot**. We study the behavior of robots in the cooperative navigation task to use less information, resulting in a policy that is adaptable to more robots without the need of retraining. Our approach requires a single reference target for the centroid of the formation to navigate. This setup has not been addressed before in the field of cooperative navigation using MARL so far. **(iii) Behavior-based navigation of real robots.** We demonstrate the performance of the learned policy on real robots. Employing the distributed MPC safety filters ensures zero collisions throughout the experiments. Additionally, we show that it is safe to train the policy on the real robots, as our approach eliminates all collisions, even during the early stages of training.



**Fig. 2:** The observation space per robot includes lidar readings (red lines), distances and headings to the goal ($d_i^g$, $\theta_i^g$), two neighbors ($d_{ij}$, $\theta_{ij}$), and the closest obstacle ($d_i^{obs}$, $\theta_i^{obs}$). Additionally, robots have information about the centroid's distance to the goal ($d_c^g$) $d_c^g$.

## II. RELATED WORK

**Cooperative Navigation Using RL** Several studies applied RL to achieve cooperative navigation using different formulations. In [7] the authors used multi-agents proximal policy optimization [22] along with curriculum learning, reference targets, and relative positions of the robots with respect to each other. Additionally, [19] applied PPO for formation control of quadrotors, incorporating individual goals for each unit. The authors in [8] used imitation learning along with RL to formulate a distributed formation controller based on a leader-follower scheme. In [9] the authors used policy gradients along with a graph convolutional network (GCN) and reference targets for each robot to achieve static formations of drones. [23] and [24] achieved formation control of maritime unmanned surface vehicles using deep deterministic policy gradient (DDPG) and deep Q-networks, respectively, based on leader-follower approaches. **Different from the previous works**, we use a multi-agent extension of the soft-actor-critic (SAC), since the SAC showed improved sample efficiency when compared to other approaches, including on-policy methods [25]. Additionally, we only use reference targets for the centroid of the formation and information about only two neighbors for each robot. More importantly, we ensure collision-free training by using MPC-based safety filters, and we use the attention mechanism to account for the interaction between the robots.

**Safety in MARL:** Safety is a crucial aspect in MARL, with ongoing developments addressing this concern. [26] utilized control barrier functions (CBF) [27] along with multi-agent DDPG (MADDPG) for collision-free navigation of two agents. [12] explored the use of attention modules and decentralized safety shields based on CBF to reduce collisions in autonomous vehicle scenarios. Similarly, [13] implemented a centralized neural network as a safety layer with MADDPG in cooperative navigation. Moreover, [28] proposed using a predictive shielding approach, along with MADDPG, to ensure the safety of multi-agents in cooperative navigation scenarios. [29] presented RL-based model predictive control to achieve leader-follower formation control of mobile robots while ensuring their safety. **In contrast to the previous works**, we focus on behavior-based navigation where collisions between agents during navigation is more probable to occur, since the agents are required to navigate close to each other and not just have a single instance where their paths intersect. Additionally,
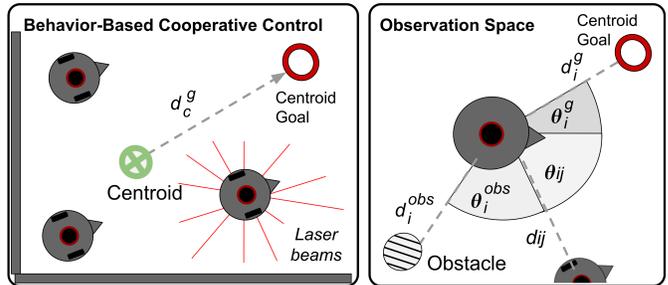
the previous studies focused on the task of having individual reference targets, which can be achieved by training a single RL policy and then deploy it to several robots as in [30]–[32]. Our task, on the other hand, necessitates using a MARL framework so that the robots interact with each other during the training to rely only on a reference location for the centroid. While the previous studies fell short of real-world experiments, we conducted extensive experiments with both trained and untrained policies on real robots to demonstrate the safety of our framework, and the performance of our behavior-based modification.

## III. PROBLEM STATEMENT

We consider the task of safely training a team of mobile robots to move the centroid of their formation to a desired target location while avoiding collisions with static obstacles and neighboring robots, and aiming to maintain predefined distances between robots. We assume that each robot is equipped with a lidar sensor to perceive its surroundings, each robot can localize itself within the environment, and that there is no pre-knowledge about the environment or existing obstacles. To ensure the safety of the robots at all times, each robot has its own safety layer to avoid collisions with its neighbors and the obstacles. The robots do not communicate with each other, but navigate cooperatively based on information between their centroid and the reference target, and the relative distances. Each robot receives relative distances and angles only about its two neighbors even if the team consists of more robots, as illustrated in Fig. 2. Finally, each robot gets the distance and relative angle to the closest obstacle and the target location of the centroid. We do not assume that the global information is available for each robot. This enables the scalability of our approach when applied to more robots, without retraining the policy.

## IV. OUR APPROACH

In this section, we formulate our work as a MARL problem, introduce the attention module for agent interaction, and describe the MPC-based distributed safety filters in detail.

### A. Multi-Agent Reinforcement Learning (MARL):

In this work, we adopt the Markov games' framework [33], which is formulated as a set of partial observable Markov decision processes (POMDPs). For the case of $N$ agents, we
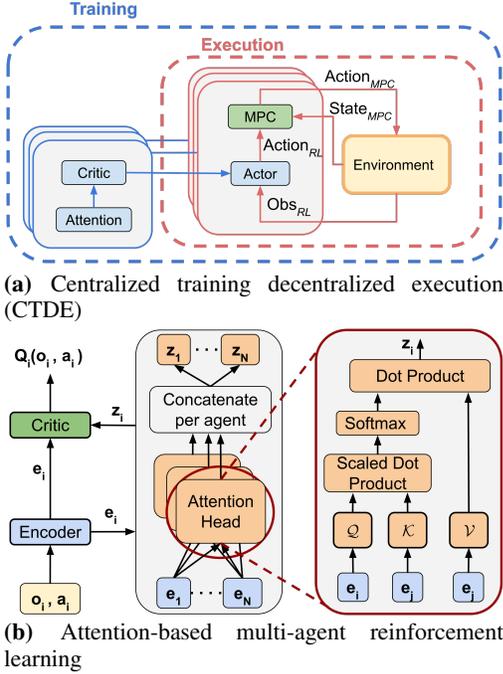
**(a)** Centralized training decentralized execution (CTDE)



**(b)** Attention-based multi-agent reinforcement learning

**Fig. 3:** Illustrations of the CTDE architecture and attention module. **(a)** At each time step, the agent interacts with the environment to receive the current observation $Obs_{RL}$, including relative information about the two neighbors and the closest obstacle, and outputs the $Action_{RL}$. The MPC controller receives the $State_{MPC}$, overrides any unsafe actions to prevent collisions, and sends $Action_{MPC}$ to the robot. During training (blue dashed), critics share all agents' observations, while during execution (red dashed), each actor accesses only its own observation. **(b)** In the attention-based critics, observations are encoded separately and fed into attention heads to calculate weights based on the query, and the key-value pairs. The attention output is then concatenated with the states and fed into the critics.

have a set $S$ that represents the augmented state space of all the agents, a set of actions $A_1, ..., A_N$ and a set of observations $O_1, ..., O_N$ for each agent. At each time step $t$, an agent $i$, receives an observation correlated with the state $o_i^t : S \rightarrow O_i$, which contains partial information from the global state $s^t \in S$. Where $s^t$ refers to $state_{RL}$ in Fig. 3a . The agent chooses an action $a_i^t$ according to a policy, $\pi_i : O_i \rightarrow A_i$, which maps each agent's observation to an action. The agent then receives a reward as a function of the state and the action taken, $r_i : S \times A_i \rightarrow \mathbb{R}$. For each agent, the tuple containing the state, the action, the reward, and the next state is added to the respective replay buffer $\mathcal{D}_i$.

Our approach is based on the soft-actor-critic (SAC) [25] algorithm. We extend SAC to the multi-agent domain and apply the centralized-training decentralized execution (CTDE) [34] scheme. During *training* all the agents share their observations with each other. However, during *execution*, each agent has access to its own observation only as shown in Fig. 3a.

For each agent we define the **observation space** as in Fig. 2. The observation includes 40 equiangular lidar readings, the relative distances $d_i^j$ and angles $\theta_i^j$ to its two neighbors, the relative distance $d_i^{obs}$ and angle $\theta_i^{obs}$ to the nearest obstacle, calculated from the lidar scan, the relative distance $d_i^{goal}$ and angle $\theta_i^{goal}$ to the goal of the centroid, the distance from the

centroid of the formation to the goal $d_{centroid}^{goal}$, and the actions of the robot for the previous time step $a_i^{t-1}$ which is the pair $(v_i^{t-1}, w_i^{t-1})$. For the **action space**, we control the linear and angular velocities $v$ and $w$. The **reward** function for each robot $i$ is defined as:

$$\mathbf{r_i^t} = \mathbf{r}_{goal} \cdot \mathbf{1}_{\text{goal reached}} + \mathbf{r}_i^{collision} \cdot \mathbf{1}_{\text{collision or stuck}}$$
$$+ (\mathbf{r}_i^{formation} + \mathbf{r}_i^{obs} + \mathbf{r}_{centroid}^{goal}) \cdot \mathbf{1}_{\text{otherwise}}$$

where:

- $r_{goal}$ and $r_i^{collision}$ are sparse rewards indicating if the centroid reached the target and if the robot has collided or hasn't been moving (stuck) for several consecutive steps, respectively.
- $r_i^{formation}$ is a continuous reward for maintaining the formation, calculated as the negative of the error between the reference distance among the agents and the actual distance.
- $r_i^{obs}$ is a continuous reward for keeping safe distance to the closest obstacle, calculated as the negative distance to the closest obstacle.
- $r_{centroid}^{goal}$ is a continuous reward for guiding the agents to move their centroid to the desired goal location, calculated as the negative distance between the centroid and the goal.

### B. Attention-Based Critics:

To effectively capture relative information among agents, we implement an attention module (Fig. 3b) following the framework of [35], which computes attention based on queries, keys, and values. Each agent's observation is encoded into embeddings $e_1 \ldots e_N$ by an encoder network. These embeddings are input to the $K$ and $V$ networks to generate keys and values, respectively, while the querying agent's embedding is processed by the $Q$ network. The computed attention weights are then used to enhance the encoded observations, which are concatenated with the agent-specific outputs $z_i$ from the attention module and input into the critics to compute Q-values.

### C. Model Predictive Safety Filter:

To ensure the safety of the agents, we implement distributed shields based on nonlinear model predictive control (NMPC). The MPC utilizes a mathematical model of the robot (1) to calculate the predicted states based on the current state and action. At each time step, the MPC solves an optimization problem (2) to find the optimal control sequence that satisfies a set of predefined constraints on the states and actions. Only the first action of the control sequence is applied while the rest of the sequence is discarded. Prior safety related studies have employed the MPC as a safety filter [28], [36], [37] due to its capability to handle systems with multiple states, controls, and constraints. To minimize the intervention of the MPC-safety filter and avoid disrupting the exploration of the RL agent, the mentioned studies align the MPC actions with those of the RL agent. In our work, we additionally, penalize the RL agent for deviating from the MPC safe actions, which in turn teaches the RL agent to not rely on the MPC-safety filter as we show in Sec.V-C.
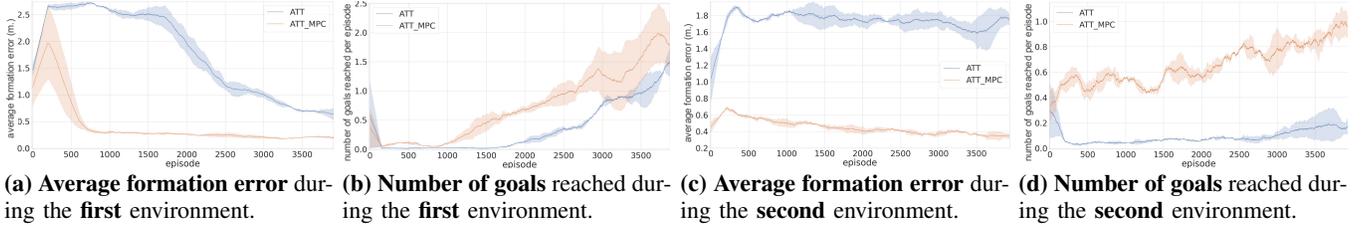
**(a) Average formation error** during the **first** environment. **(b) Number of goals** reached during the **first** environment. **(c) Average formation error** during the **second** environment. **(d) Number of goals** reached during the **second** environment.

**Fig. 4:** The figures show the impact of the safety filter on training, with bold lines representing the average and shaded areas indicating the standard deviation across three random seeds. **(a,c)** depict the average formation error in the first and second training environments, while **(b,d)** display the number of goals achieved per episode, in both environments. The agent with the MPC (**ATT_MPC**) consistently surpasses the pure learning agent (**ATT**) in reducing formation errors and increasing goal achievements. This demonstrates that incorporating a safety filter reduces the number of episodes required to achieve the desired performance compared to pure learning agents.

*1) Safety Filter Formulation::* The robot's state in NMPC ($State_{MPC}$ in Fig. 3a) is defined as $\mathbf{x} = [x, y, \theta]^T$, representing its position and heading. Controls are denoted by $\mathbf{a} = [v, \omega]^T$, matching the RL actions. A weighted Euclidean norm with a positive definite weighting matrix $R$ is denoted as $||\mathbf{x}||_R^2 = \mathbf{x}^T R \mathbf{x}$.

*2) Prediction Model::* The discrete-time model of the robot is:

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \begin{bmatrix} cos\theta_t & 0 \\ sin\theta_t & 0 \\ 0 & 1 \end{bmatrix} \mathbf{a}_t \Delta t \tag{1}$$

*3) Optimal Control Problem: :* To ensure the safety of the robot while minimizing the intervention of the safety filter, we formulate the optimization problem for each robot $i$ at time step $t$ as follows:

$$\min_{\substack{\mathbf{x}_{t:t+T|t}, \\ \mathbf{a}_{t:t+T-1|t}}} \|\mathbf{a}_{RL} - \mathbf{a}_{t|t}\|_{R_0}^2 + \sum_{k=1}^{T-1} \|\mathbf{a}_{t+k|t}\|_R^2 + \sum_{k=0}^{T} \|\frac{1}{e^{dist_{i,1}^{t+k|t}}}\|_D^2$$

$$+ \sum_{k=0}^{T} \frac{1}{e^{dist_{i,2}^{t+k|t}}}\|_D^2 + \sum_{k=0}^{T} \|\frac{1}{e^{dist_{obst}^{t+k|t}}}\|_D^2 \tag{2a}$$

$$\text{s.t.} \quad \mathbf{x}_{t|t} = \mathbf{x}_t, \tag{2b}$$

$$\mathbf{x}_{t+k+1|t} = f(\mathbf{x}_{t+k|t}, \mathbf{a}_{t+k|t}), \forall k = 0, 1, ..., T-1 \tag{2c}$$

$$\mathbf{a}_{t+k|t} \in \mathbf{A}, \qquad \forall k = 0, 1, ..., T-1 \tag{2d}$$

$$\mathbf{x}_{t+k|t} \in \mathbf{X}, \qquad \forall k = 0, 1, ..., T-1, \tag{2e}$$

where $T$ represents the prediction horizon. The notation $t + k|t$ indicates predictions at time $t+k$, assuming the current time is $t$. The optimization terms in (2a) are structured in three main components. The first term measures the deviation between the RL agent's proposed action $\mathbf{a}_{RL}$ ($Action_{RL}$ in Fig. 3a) and the initial MPC action $\mathbf{a}_0$ ($Action_{MPC}$), optimizing for minimal discrepancy. The Second term minimizes the magnitude of future control signals $\mathbf{a}_k$, promoting smoother transitions. Third, the distances $dist_{i,1}$ and $dist_{i,2}$ between the robot and its two neighbors, as well as $dist_{obst}$ to the nearest obstacle, are penalized to ensure safety. The exponential function applied to distances has proven effective in enhancing obstacle avoidance, as shown in previous research [38]. Weight matrices $R_0$, $R$, and $D$ are manually tuned to balance the contributions of the cost terms.

Constraints (2b–2e) ensure adherence to system dynamics and operational limits, defined as follows: (2b) sets the initial state, (2c) enforces the robot's dynamic model, (2d) and (2e) maintain the actions and states within predefined bounds,

where $\mathbf{X}$ and $\mathbf{A}$ denote the allowable sets of states and controls, respectively.

The safety filter, operating independently of the behavior-based navigation task, ensures robot safety by aligning with RL agent actions to maintain exploratory integrity during training. It is compatible with any MARL framework as it does not change the RL algorithm as we show in Sec.V-D. For implementation, we employed acados [39], a tool for solving nonlinear optimal control problems.

## V. EXPERIMENTAL EVALUATION

The main focus of this work is to achieve behavior-based navigation of mobile robots while ensuring the safety of the robots. We design the experiments so that we can: **(i)** study the effect of using the MPC on the training of the RL agents, showing that we can achieve collision free training to learn the task in fewer number of episodes (Sec. V-A), **(ii)** compare our approach against baselines in simulations, demonstrating that our approach outperforms state-of-the-art baselines, that our RL agent learns to not rely on the MPC-safety layer indefinitely, and that the MPC-safety layer can be integrated with other RL algorithms (Sec. V-B:V-D), **(iii)** show that our approach can be transferred on real robots (Sec. V-E), **(iv)** test the ability of the trained policy to generalize to more agents than used during training (Sec. V-F).

### A. Training With the MPC Filter

To evaluate the MPC's impact on agent training, we simulated two environments commonly used in the safe MARL literature [12], [13], [26], [28], [29]. The first environment is an empty walled setup to teach the robots to move the centroid to the goal. The second environment includes randomly placed obstacles for collision avoidance learning. We train the policies in both environments with and without the MPC. We utilized the same network parameters to isolate the MPC filter's effects. Additionally, we introduced a penalty in the RL policy for deviations from the MPC-safe actions. Comparative metrics include average formation error, number of goals reached per episode, and number of collisions. Each episode terminates once the maximum number of steps is reached, or one of the robots collides or is stuck.

The average formation error is defined as follows:

$$Formation\ Error = \frac{1}{N} \sum_{i=1}^{N} [dist_{i,j} - dist_{i,j}^{ref}]$$

| | Goals Reached W/o Obs. | | | Goals Reached W/ Obs. | | | S-Shaped Path | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Agent** | Success ↑ | Colls ↓ | Tout ↓ | Success ↑ | Colls ↓ | Tout ↓ | Goals Reached ↑ | Time (sec.) ↓ | Colls ↓ | Avg form err (m) ↓ |
| Ours | **99.5%** | **0** | **0.5%** | **96.2%** | **0** | **3.8%** | **100%** | **108.93 ± 14.8** | **0** | 0.392 ± 0.17 |
| MASAC [40] | 58% | 8.9% | 33.1% | 51.8% | 26.3% | 21.9% | 75.8% | 134.63 ± 21.3 | 2.8 ± 3.65 | 1.25 ± 0.15 |
| MADDPG [41] | 9.5% | 57.7% | 32.8% | 6.6% | 81.5% | 11.9% | 12.4% | 118.9 ± 28.39 | 4.16 ± 3.81 | 4.05 ± 4.06 |
| DMPC [38] | 13.9% | 0 | 86.1% | 6.9% | 0 | 93.1% | 13.4% | 162.7 ± 15.53 | 0 | **0.127 ± 0.184** |
| | | | | | | | **Can we execute our method without the MPC?** | | | |
| Ours without MPC | 98.6% | 1.4% | 0 | 96% | 2.5% | **1.5%** | 93.7% | 112.35 ± 16.7 | 0.45 ± 1.48 | 0.42 ± 0.23 |
| | | | | | | | **Integrating the MPC Safety Layer with the Baselines** | | | |
| MASAC [40] (With MPC) | 67.7% | 0 | 32.3% | 60.3% | 0 | 39.7% | 79.9% | 129.54 ± 38.7 | 0 | 0.98 ± 0.18 |
| MADDPG [41] (With MPC) | 10.4% | 0 | 89.6 % | 7.6% | 0 | 92.4% | 13.05% | 164.9 ± 41.02 | 0 | 3.56 ± 3.67 |

**TABLE I:** The table shows the performance of our approach against multi-agent SAC (MASAC), multi-agent DDPG (MADDPG), and decentralized model predictive control (DMPC) in simulation. Additionally, it demonstrates the results for testing our approach without the MPC-safety filter, and testing the baselines with the MPC-safety filter. The percentage reflects the proportion of achieved goals, collisions, and timeouts in relation to the 1000 trials in case of the goal reaching scenarios , and 120 tests for the S-shaped path test. Our approach makes zero collisions in all configurations, and takes less time to reach the desired locations. Disabling the MPC-filter results in fewer successful runs in all scenarios, as collisions tend to happen. Eliminating the collisions from the MASAC results in better performance with respect to reaching more goals. For the MADDPG, there is no noticeable improvement.
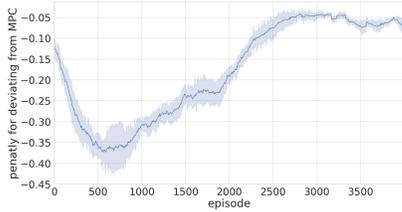


**Fig. 5:** The figure shows the penalty given to the agent for deviating from the MPC safe actions.The bold lines show the average while the shaded area show the standard deviation over three different random seeds. The penalty decreases through the training indicating that the agent learns to match the MPC safe actions.

where $dist_{i,j}^{ref}$ is the reference distance between robots $i$ and $j$, while $dist_{i,j}$ is the actual distance.

We denote the agents without the MPC filter as **ATT**, and those with the MPC as **ATT_MPC**. The results, displayed in Fig. 4, highlight the efficacy of the MPC filter in training enhancement. Notably, **ATT_MPC** reduces the average formation error to less than 0.5 meters within 1,000 episodes, a significant improvement over **ATT**, which requires about 4,000 episodes to achieve similar performance (Fig. 4a, 4c). This enhancement is largely due to the MPC filter's ability to eliminate collisions, facilitating more effective exploration by the robots.

Additionally, Fig. 4b, 4d illustrate that including the MPC-filter leads to reaching more goals per episode compared to the pure learning agents. These experiments show that we are able to decrease the number of episodes required to achieve the desired behavior, which is a crucial aspect for safe reinforcement learning, since fewer episodes means fewer resets of the environment in the real-world.

### B. Testing Against Baselines in Simulation

We evaluated our approach against three state-of-the-art baselines to validate its effectiveness in motion planning tasks. The baselines include:

- **MASAC:** A multi-agent variant of the Soft Actor-Critic (SAC) [40], as it outperformed several MARL baselines in the motion planning tasks.
- **MADDPG:** The Multi-Agent Deep Deterministic Policy Gradient [41], a standard benchmark in MARL studies.
- **DMPC:** A decentralized model predictive control approach [38], using the centroid's reference as the goal

for each robot, while maintaining inter-robot distances.

The RL agents were trained over two stages, each consisting of 4,000 episodes, with pre-trained agents loaded for further training in the second stage to adapt to obstacle avoiding scenarios. We compare between the agents in different scenarios.

**Reaching a Single Goal:** In the first scenario, we tested the ability of the robots to reach the desired goals starting from different configurations with and without obstacles in the environment. At the start of each episode, the locations of the robots, obstacles, and goal location are randomized. The results are summarized in Table I. Our approach outperforms the baselines in terms of all three metrics.

**S-Shaped Path Following:** In this scenario, we simulated S-shaped paths for the centroid of the robot formation to replicate a real-world application where robots must navigate cooperatively following a path, e.g. conducting surveillance across multiple locations. We defined eight distinct S-shaped paths and initiated the robots from 15 varied starting configurations for each path.

The performance was evaluated based on completion time, collisions, and formation error, with results summarized in Table I. Our approach demonstrated superior performance in completion time and collision avoidance, achieving the least time and zero collisions. While the DMPC showed the smallest formation error, it struggled with efficient positioning around the centroid, resulting in prolonged completion times. Among the learning-based methods, our approach exhibited the lowest formation errors, underscoring its effectiveness in maintaining formation while navigating.

### C. Can we execute our method without the MPC?

In this section, we investigated if the MPC-filter can be disabled after training. We mentioned in Sec.V-A that we added a penalty for deviating from the MPC-safe actions. In Fig. 5, we show that this penalty decreases over time, indicating the agents learn to adhere to the safe actions. To further test the agent without the MPC-safety filter, we used an agent previously trained with the MPC-filter, deactivated the MPC-safety filter and tested the agent in the same scenarios. Performance outcomes are detailed in Table I. While the agent without the MPC-filter achieved over 90% in all tests showing that it can be used without the MPC, it failed to
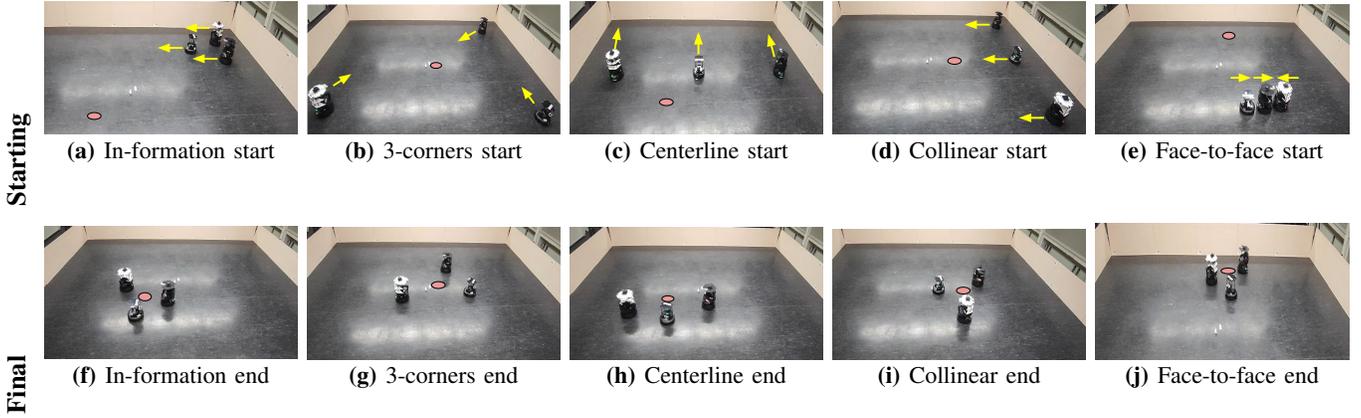
**(a)** In-formation start    **(b)** 3-corners start    **(c)** Centerline start    **(d)** Collinear start    **(e)** Face-to-face start

**(f)** In-formation end    **(g)** 3-corners end    **(h)** Centerline end    **(i)** Collinear end    **(j)** Face-to-face end

**Fig. 6:** Target Reaching Configurations test. The top row indicates the starting configurations, while the bottom row shows the final reached configurations by our approach. The yellow arrows indicate the starting orientation of the robots, while the red circles show the target location for the centroid of the formation. The configurations are arranged from left to right in terms of difficulty. The **Collinear** and **Facing each other** configurations were not experienced during the training, making them more difficult compared to the other configurations. However, our approach was able to successfully complete the tasks with zero collisions.

| | Reaching a Single Goal | | | | | |
|---|---|---|---|---|---|---|
| Agent | Ours | | | MASAC [40] | | |
| Configuration | Success ↑ | Time (sec.) ↓ | Colls ↓ | Success ↑ | Time (sec.)↓ | Colls ↓ |
| In-formation | 100% | **27.58 ± 7.82** | **0** | 100% | 44.7 ± 15.4 | 1.0 ± 0 |
| 3-corners | 100% | **22.3 ± 0.37** | **0** | 100% | 35.01 ± 16.6 | 1.0 ± 0 |
| Centerline | 100% | **17.71 ± 0.43** | **0** | 100% | 18.6 ± 0.21 | 1.5 ± 0.7 |
| Collinear | 100% | 15.94 ± 1.6 | **0** | 100% | **15.55 ± 3.5** | 0.5 ± 0.7 |
| Face-to-face | 100% | 45.7 ± 2.2 | **0** | 100% | **22 ± 6.5** | 2.5 ± 0.7 |

**TABLE II:** The table shows the performance of our approach against multi-agent SAC (MASAC). Each scenario is tested two times using each approach. Our approach makes zero collisions in all configurations. In terms of the time taken, our approach took less time to reach the desired locations for the centroids in the first three configurations. The last two configurations were previously unseen during the training which explains why our approach took more time to reach the target, nevertheless with zero collisions unlike the MASAC.

avoid collisions as effectively as it did with the MPC enabled, highlighting the MPC's critical role in ensuring safety.

### D. Integrating the MPC Safety Layer with the Baselines:

We evaluated the baseline models after being retrained with the MPC safety layer and evaluated their performance using the established tests. The outcomes are detailed in Table I. Incorporating the MPC with MASAC eliminated collisions, which increased the number of goals achieved but also led to more timeouts, indicating the MPC's intervention to prevent collisions. Adding the MPC to MADDPG did not yield improvements, potentially due to the inherent stability issues of DDPG [25]. The frequent timeouts suggest that the MPC often stops the robots to avoid collisions, resulting in timeouts due to the robots being stuck.

### E. Real-World Experiments

In the accompanying video, we demonstrated the behavior of initial (random) policies with and without the MPC filter on real robots, showcasing that our approach enables safe training on physical robots. However, training multiple robots in the real world is currently impractical due to time requirements. Our framework, although collision-free, requires 25 hours of simulation time, which translates to about 100 hours of real-world training.

Previous studies on safe MARL [12]–[14], [26], [28] have only evaluated final policies in simulations. In contrast, we validate our trained policies on real robots. The aim of

these experiments is to demonstrate the zero-shot transfer of behavior-based navigation and the safety of our approach.

For comparison, we selected MASAC as a baseline due to its superior simulation performance over MADDPG and DMPC (Sec. V-B), and trained both for the same number of episodes. The policies were then tested in real-world scenarios identical to those used in simulations.

**Reaching a Single Goal:** We tested five different starting configurations with a target location for the formation centroid, as shown in Fig. 6. Each configuration was tested twice for each policy, with results summarized in Table II. Our approach consistently achieved zero collisions across all configurations, including unseen scenarios, demonstrating superior generalization and safety. Additionally, our method reached the targets faster than MASAC, as it better coordinates robot movements.

In unseen scenarios, our approach required more maneuvers to avoid collisions, unlike MASAC, where robots reached the targets but failed to avoid collisions. For example, in the **face-to-face** configuration, our robots rotated and moved sequentially to avoid collisions, whereas MASAC robots collided and moved while touching, which explains the less time taken (these results are shown in the video). All tested configurations were reachable by both methods. We further tested more configurations that were not reachable by the baseline in the video. Both approaches were tested in obstacle avoidance scenarios. Fig. 7 show two of these scenarios, and the results are shown in Table III.

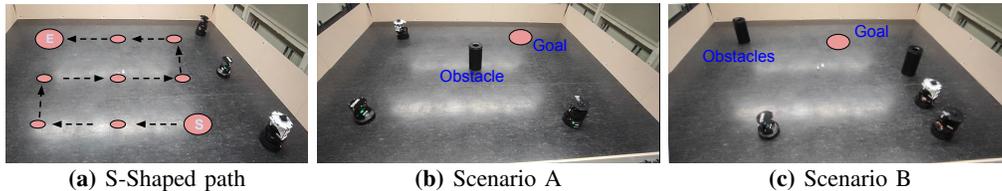**S-Shaped Path Following:** In the S-shaped path scenario,

**(a)** S-Shaped path      **(b)** Scenario A      **(c)** Scenario B

**Fig. 7:** The figures show the S-Shaped path tests and two scenarios of the obstacle avoidance tests. **(a)** shows the S-shaped path. The red circles indicate the target locations for the centroid of the formation, where the first and final targets are indicated by the letters **S**, **E**, respectively. **(b,c)** show two different scenarios for obstacle avoidance. The robots start from the shown configurations and navigate towards the goal while avoiding collisions with the obstacles shown in the figures.

| | **Goals Reaching With Obstacles** | | | | **S-Shaped Path** | | | |
|---|---|---|---|---|---|---|---|---|
| **Agent** | Success ↑ | Time (sec.) ↓ | Colls ↓ | Touts↓ | Goals reached ↑ | Completion time (sec.) ↓ | Collisions ↓ | Avg form err (m)↓ |
| Ours | **100%** | **31.43 ± 5.56** | **0** | **0** | **100%** | **143.2 ± 5.7** | **0** | **0.34 ± 0.016** |
| MASAC [40] | 0% | - | 100% | 0 | 96.2% | 212.2 ± 53.8 | 7.3 ± 1.52 | 0.58 ± 0.35 |

**TABLE III:** The aggregated results for the obstacle avoidance and the S-Shaped path tests. Our approach completes all four tests with zero collisions. The MASAC was not successful in reaching the goals in case of obstacle avoidance tests due to collisions with the obstacles in all tests, eventually resulting in timeouts. **The S-Shaped path test** was carried out three times for each approach. The MASAC reached 96.2% of the targets, since the robots failed to reach some targets within the allowed number of steps. During the tests, if the robots take more than 200 steps to reach a target, the next target in the path is given instead, and the previous target is counted as a timeout.
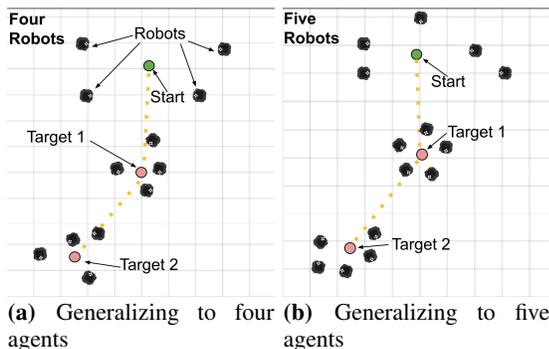


**(a)** Generalizing to four agents    **(b)** Generalizing to five agents

**Fig. 8:** Example scenarios **(a)** and **(b)** show the generalization of the trained policy to a larger team of robots. The robots successfully move their centroid from the starting location (green circle) to the target location (red circle). The yellow lines indicate the movement of the centroid through the different locations.

| **Goals Reaching Test With More Robots** | | | |
|---|---|---|---|
| **Number of Robots** | **Goals reached** | **Collisions** | **Timeouts** |
| Three | 99% | 0 | 1% |
| Four | 97% | 0 | 3% |
| Five | 91% | 0 | 9% |
| Six | 79% | 0 | 21% |
| Seven | 76% | 0 | 24% |
| Eight | 70% | 0 | 30% |

**TABLE IV:** The table shows the performance of our approach when generalizing to more robots, over 100 episodes. The percentage reflects the proportion of achieved goals, collisions, and timeouts in relation to the total number of trials. Our approach is able to generalize to more robots with zero collisions. However, as the number of robots increased, we observed an increase in timeouts, attributed to the conservative nature of the MPC filter.

Fig. 7a, we run the test three times for each approach and recorded the number of collisions, the average formation error over the whole path, number of timeouts, as well as the time taken to complete the path. The results are summarized in Table III. Our approach was able to successfully reach all targets without making any collisions, and in less times compared to the MASAC. Additionally, our approach had less average formation error compared to the MASAC over all the

runs.

### F. Generalizing to More Robots

Finally, we show that by relying only on information from two neighbors, our method can scale to more robots while maintaining a fixed observation size, regardless of the number of robots.

We evaluated the performance of the trained policy with up to eight robots, as shown in Fig. 8. For each robot, we defined two neighbors and tested the policy using the previously introduced goal reaching scenarios. Metrics included the number of goals reached, collisions, and timeouts, with results summarized in Table IV. The robots successfully cooperated to move their centroid to the target locations, with zero collisions. However, as the number of robots increased, we observed an increase in timeouts, attributed to the conservative nature of the MPC filter. These trade-offs are expected in safety-critical systems, where eliminating collisions is paramount.

### VI. CONCLUSIONS

In this study, we presented a safe multi-agent reinforcement learning (MARL) approach to achieve behavior-based cooperative navigation without individual reference targets in real-world scenarios. Our findings demonstrate that relying on the centroid of the formation is sufficient for robots to navigate cooperatively. By integrating MARL, and model predictive control (MPC)-based safety filters, we ensured zero collisions during training and achieved faster convergence. The inclusion of a safety layer not only eliminated collisions but also significantly improved training efficiency, leading to faster convergence of the MARL policies. This has crucial implications for addressing the sim-to-real gap, as it highlights the benefits of incorporating safety layers into RL training. Our trained policies, applied to real robots, outperformed baseline methods in various tests in terms of success rate, number of collisions and completion times, confirming the effectiveness of our approach. Our experiments show that training on real robots is safe, eliminating collisions even during early exploration stages.

# REFERENCES

[1] Y. Liu and G. Nejat, "Multirobot cooperative learning for semiautonomous control in urban search and rescue applications," *Journal of Field Robotics*, vol. 33, no. 4, pp. 512–536, 2016.

[2] R. Tallamraju, E. Price, R. Ludwig, K. Karlapalem, H. H. Bülthoff, M. J. Black, and A. Ahmad, "Active perception based formation control for multiple aerial vehicles," *IEEE Robotics and Automation Letters (RA-L)*, vol. 4, no. 4, pp. 4491–4498, 2019.

[3] R. T. Fawcett, L. Amanzadeh, J. Kim, A. D. Ames, and K. A. Hamed, "Distributed data-driven predictive control for multi-agent collaborative legged locomotion," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*. IEEE, 2023, pp. 9924–9930.

[4] G.-P. Liu and S. Zhang, "A survey on formation control of small satellites," *Proceedings of the IEEE*, vol. 106, no. 3, pp. 440–457, 2018.

[5] D. Xu, X. Zhang, Z. Zhu, C. Chen, P. Yang, *et al.*, "Behavior-based formation control of swarm robots," *mathematical Problems in Engineering*, vol. 2014, 2014.

[6] L. Quan, L. Yin, T. Zhang, M. Wang, R. Wang, S. Zhong, X. Zhou, Y. Cao, C. Xu, and F. Gao, "Robust and efficient trajectory planning for formation flight in dense environments," *IEEE Trans. on Robotics (TRO)*, 2023.

[7] Y. Yan, X. Li, X. Qiu, J. Qiu, J. Wang, Y. Wang, and Y. Shen, "Relative distributed formation and obstacle avoidance with multi-agent reinforcement learning," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*. IEEE, 2022, pp. 1661–1667.

[8] Z. Sui, Z. Pu, J. Yi, and S. Wu, "Formation control with collision avoidance through deep reinforcement learning using model-guided demonstration," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 6, pp. 2358–2372, 2020.

[9] A. Khan, E. Tolstaya, A. Ribeiro, and V. Kumar, "Graph policy gradients for large scale robot control," in *Conference on robot learning*. PMLR, 2020, pp. 823–834.

[10] N. Rudin, D. Hoeller, P. Reist, and M. Hutter, "Learning to walk in minutes using massively parallel deep reinforcement learning," in *Conference on Robot Learning*. PMLR, 2022, pp. 91–100.

[11] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 23–30.

[12] Z. Zhang, S. Han, J. Wang, and F. Miao, "Spatial-temporal-aware safe multi-agent reinforcement learning of connected autonomous vehicles in challenging scenarios," *arXiv preprint arXiv:2210.02300*, 2022.

[13] Z. Sheebaelhamd, K. Zisis, A. Nisioti, D. Gkouletsos, D. Pavllo, and J. Kohler, "Safe deep reinforcement learning for multi-agent systems with continuous action spaces," *arXiv preprint arXiv:2108.03952*, 2021.

[14] I. ElSayed-Aly, S. Bharadwaj, C. Amato, R. Ehlers, U. Topcu, and L. Feng, "Safe multi-agent reinforcement learning via shielding," *arXiv preprint arXiv:2101.11196*, 2021.

[15] M. Samvelyan, T. Rashid, C. S. De Witt, G. Farquhar, N. Nardelli, T. G. Rudner, C.-M. Hung, P. H. Torr, J. Foerster, and S. Whiteson, "The starcraft multi-agent challenge," *arXiv preprint arXiv:1902.04043*, 2019.

[16] P. Zhang, G. Chen, Y. Li, and W. Dong, "Agile formation control of drone flocking enhanced with active vision-based relative localization," *IEEE Robotics and Automation Letters (RA-L)*, vol. 7, no. 3, pp. 6359–6366, 2022.

[17] T. Xu, J. Liu, Z. Zhang, G. Chen, D. Cui, and H. Li, "Distributed MPC for trajectory tracking and formation control of multi-uavs with leader-follower structure," *IEEE Access*, pp. 1–1, 2023.

[18] S. Park and S.-M. Lee, "Formation reconfiguration control with collision avoidance of nonholonomic mobile robots," *IEEE Robotics and Automation Letters (RA-L)*, 2023.

[19] S. Batra, Z. Huang, A. Petrenko, T. Kumar, A. Molchanov, and G. S. Sukhatme, "Decentralized control of quadrotor swarms with end-to-end deep reinforcement learning," in *Conference on Robot Learning*. PMLR, 2022, pp. 576–586.

[20] V. K. Adajania, S. Zhou, A. K. Singh, and A. P. Schoellig, "Am-swarmx: Safe swarm coordination in complex environments via implicit non-convex decomposition of the obstacle-free space," *arXiv preprint arXiv:2310.09195*, 2023.

[21] L. Quan, L. Yin, C. Xu, and F. Gao, "Distributed swarm trajectory optimization for formation flight in dense environments," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 4979–4985.

[22] C. Yu, A. Velu, E. Vinitsky, Y. Wang, A. Bayen, and Y. Wu, "The surprising effectiveness of ppo in cooperative, multi-agent games," *arXiv preprint arXiv:2103.01955*, 2021.

[23] J. Xie, R. Zhou, Y. Liu, J. Luo, S. Xie, Y. Peng, and H. Pu, "Reinforcement-learning-based asynchronous formation control scheme for multiple unmanned surface vehicles," *Applied Sciences*, vol. 11, no. 2, p. 546, 2021.

[24] X. Zhou, P. Wu, H. Zhang, W. Guo, and Y. Liu, "Learn to navigate: cooperative path planning for unmanned surface vehicles using deep reinforcement learning," *IEEE Access*, vol. 7, pp. 165 262–165 278, 2019.

[25] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *International conference on machine learning*. PMLR, 2018, pp. 1861–1870.

[26] Z. Cai, H. Cao, W. Lu, L. Zhang, and H. Xiong, "Safe multi-agent reinforcement learning through decentralized multiple control barrier functions," *arXiv preprint arXiv:2103.12553*, 2021.

[27] A. D. Ames, S. Coogan, M. Egerstedt, G. Notomista, K. Sreenath, and P. Tabuada, "Control barrier functions: Theory and applications," in *2019 18th European control conference (ECC)*. IEEE, 2019, pp. 3420–3431.

[28] W. Zhang, O. Bastani, and V. Kumar, "Mamps: Safe multi-agent reinforcement learning via model predictive shielding," *arXiv preprint arXiv:1910.12639*, 2019.

[29] X. Zhang, Y. Peng, W. Pan, X. Xu, and H. Xie, "Barrier function-based safe reinforcement learning for formation control of mobile robots," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*. IEEE, 2022, pp. 5532–5538.

[30] B. Brito, M. Everett, J. P. How, and J. Alonso-Mora, "Where to go next: learning a subgoal recommendation policy for navigation in dynamic environments," *IEEE Robotics and Automation Letters (RA-L)*, 2021.

[31] A. P. Vinod, S. Safaoui, A. Chakrabarty, R. Quirynen, N. Yoshikawa, and S. Di Cairano, "Safe multi-agent motion planning via filtered reinforcement learning," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*. IEEE, 2022, pp. 7270–7276.

[32] M. Li, Y. Jie, Y. Kong, and H. Cheng, "Decentralized global connectivity maintenance for multi-robot navigation: A reinforcement learning approach," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 8801–8807.

[33] M. L. Littman, "Markov games as a framework for multi-agent reinforcement learning," in *Machine learning proceedings 1994*. Elsevier, 1994, pp. 157–163.

[34] J. Foerster, I. A. Assael, N. De Freitas, and S. Whiteson, "Learning to communicate with deep multi-agent reinforcement learning," *Advances in Neural Information Processing Systems*, vol. 29, 2016.

[35] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in Neural Information Processing Systems*, vol. 30, 2017.

[36] K. P. Wabersich, L. Hewing, A. Carron, and M. N. Zeilinger, "Probabilistic model predictive safety certification for learning-based control," *IEEE Transactions on Automatic Control*, vol. 67, no. 1, pp. 176–188, 2021.

[37] F. P. Bejarano, L. Brunke, and A. P. Schoellig, "Multi-step model predictive safety filters: Reducing chattering by increasing the prediction horizon," in *2023 62nd IEEE Conference on Decision and Control (CDC)*. IEEE, 2023, pp. 4723–4730.

[38] M. Dawood, N. Dengler, J. de Heuvel, and M. Bennewitz, "Handling sparse rewards in reinforcement learning using model predictive control," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*. IEEE, 2023, pp. 879–885.

[39] R. Verschueren, G. Frison, D. Kouzoupis, J. Frey, N. van Duijkeren, A. Zanelli, B. Novoselnik, T. Albin, R. Quirynen, and M. Diehl, "acados: a modular open-source framework for fast embedded optimal control," 2020.

[40] Z. He, L. Dong, C. Song, and C. Sun, "Multiagent soft actor-critic based hybrid motion planner for mobile robots," *IEEE transactions on neural networks and learning systems*, 2022.

[41] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. Pieter Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," *Advances in Neural Information Processing Systems*, vol. 30, 2017.